



The Hidden Complexity of Price Formation:

Exploring Microstructural Mechanisms in Financial Markets

Thèse de doctorat de l'Institut Polytechnique de Paris préparée à l'École polytechnique

École doctorale n°626 École doctorale de l'Institut Polytechnique de Paris (EDIPP)

Spécialité de doctorat: Physique

Thèse présentée et soutenue à Palaiseau, le 18 Septembre 2025, par

Guillaume Maitrier

Composition du Jury:

Emmanuel Bacry Président

Université Paris-Dauphine

Doyne Farmer Rapporteur

Oxford University

Fabrizio Lillo Rapporteur

University of Bologna

Cécile Monthus Examinatrice

Michael Benzaquen

Directeur de thèse CNRS & École Polytechnique

Jean-Philippe Bouchaud Co-directeur de thèse

CFM & Académice des Sciences

Grégoire Loeper Co-encadrant BNP Paribas CIB

Invité Mathieu Rosenbaum

Ecole Polytechnique

CNRS & CEA

 $Intentionally\ left\ blank.$

Remerciements

Je voudrais avant tout remercier chaleureusement mes trois directeurs de thèse pour ces trois années passées en un éclair. Merci pour votre grande bienveillance, votre patience et votre disponibilité, que ce soit sur les aspects techniques de la thèse ou sur d'autres plus généraux. J'ai une pensée particulière pour Jean-Philippe, à qui je dois l'ensemble de ce travail et bien plus. Merci pour ton extraordinaire implication, ce fut un privilège et un honneur de pouvoir travailler avec toi. Ces discussions m'accompagneront longtemps!

Je remercie aussi les équipes de CFM et de BNP avec qui j'ai eu la chance de collaborer, en particulier Julius Bonart, pour ses nombreux conseils.

I also want to deeply thank the jury members: Cecilia Monthus, Doyne Farmer, Fabrizio Lillo, and Emmanuel Bacry. I am especially grateful to Doyne Farmer and Fabrizio Lillo for accepting to be rapporteurs for this thesis, especially since I spent most of these three years trying to build upon the massive foundations they laid. I also want to warmly thank Mathieu Rosenbaum, for being part of my comité de suivi and present in this jury. This work would not be what it is without my stay in Japan, which was an incredible experience. I am sincerely redevable to Kiyoshi Kanazawa for welcoming me to his lab, and I deeply thank all the lab members—especially Yuki Sato—for those wonderful months we shared.

Then I would like to thank all the members of the EconophysiX lab —the atmosphere was truly amazing, and I learned so much thanks to all of you. A special 1/3+1/3 thanks for Max and Elia! Merci à Cecilia pour son soutien indéfectible, ainsi qu'à Salma, co-autrice idéale. Enfin, une pensée émue pour Antoine-Cyrus enfin, Bêche quoi - pour sa patience infinie et son humour douteux, je ne pouvais rêver meilleur "camarade" pour ces trois années. Merci!

Je profite aussi de cette occasion pour remercier les copaings, notamment de l'X et de Re-df qui sont de vrais piliers dans ma vie. Un merci particulier à ceux venus dropper et ceux de Faverges, qui ont beaucoup influencé cette thèse et mes choix futurs.

Je remercie profondément mes frères et soeurs (et leurs +1), qui sont pour moi de véritables exemples, mention spéciale pour Marianne, roc certifié depuis 26 ans, et Thibaut, pour la relecture de cette thèse. Un grand merci aussi à mes grands-parents qui m'accompagnent - pour certains à distance.

Enfin, je m'estime particulièrement chanceux d'avoir de tels parents, à qui ce manuscrit est dédié. Merci pour tout ce que vous m'apportez, chacun à votre manière, de façon si complémentaire.

Foreword

The story behind this work begins with the lectures on *Physics of Financial Markets* at École Polytechnique, a rather intriguing course, offered jointly to economics and physics majors, and held at the unfortunate hour of 8 a.m. on Fridays. Despite the timing, the course was a revelation for me. It showed how the elegance of theoretical results could be grounded in physical intuition and observation, and ultimately applied to a field as fascinating—and as central to our everyday lives—as economics.

To bridge the gap between engineering and economics, I decided to pursue an additional year focused on finance, as I was convinced that a PhD in Econophysics was one of the most exciting paths I could take. Therefore, I'm deeply grateful to Jean-Philippe Bouchaud and Michael Benzaquen for giving me the opportunity to join their lab, and to Grégoire Loeper for welcoming me at BNP Paribas CIB. Being part of these two teams—with their very different perspectives on markets, one rooted in physics and the other in a more mathematical and institutional approach—has been an incredibly enriching experience. These two experiences were further complemented by several exciting teaching roles as an Assistant Professor at ENSAE and Sorbonne University.

The initial research idea for this PhD was to study jumps and liquidity crises in the order book. And just like financial markets, life brings its own share of uncertainty: this first project, which I began in September 2022, ended up being the last one I completed in three years later... Perhaps a good illustration of what researchers often call "the long horizon of research". Indeed, it was interrupted by an amazing and unexpected opportunity to spend five months in Japan, starting in February 2024. I had the chance to join the Kanazawa lab—I am also deeply grateful to him—and work on a unique dataset containing trader identifiers. This experience opened the door to the study of price impact, which eventually became the main focus of my thesis and a subject I find truly fascinating.

Luckily, one of our projects, the VAR model, turned out to provide a natural transition from price impact to liquidity crisis. As a result, the structure of this thesis is as follows: it begins with a general introduction, which can be skipped by readers already familiar with financial markets. This is followed by two more technical introductions, one on price impact models and the other on market stability. Part II (Price Impact) and Part III (Market Stability) each present the technical studies we conducted. Finally, Part IV concludes the thesis and outlines several directions for future research.

List of publications

- [1] Elomari-Kessab S., **Maitrier G.**, Bonart J. & Bouchaud J.-P. (2024). *Microstructure Modes Disentangling the Joint Dynamics of Prices & Order Flow. Wilmott Magazine*, doi:10.54946/wilm.12074
- [2] Maitrier G., Loeper G., Kanazawa K. & Bouchaud J.-P. (2025). The "double" square-root law: Evidence for the mechanical origin of market impact using Tokyo Stock Exchange data. Under review, arXiv:2502.16246
- [3] Maitrier G., Loeper G. & Bouchaud J.-P. (2025). Generating realistic metaorders from public data. Under review, arXiv:2503.18199
- [4] Maitrier G. & Bouchaud J.-P. (2025). The Subtle Interplay between Square-root Impact, Order Imbalance & Volatility: A Unifying Framework. Under review, arXiv:2506.07711
- [5] Maitrier G., Loeper G. & Bouchaud J.-P. (2025). The Subtle Interplay between Square-root Impact, Order Imbalance & Volatility II: An Artificial Market Generator. Under review, arXiv:2509.05065
- [6] Maitrier G., Loeper G., Benzaquen M. & Bouchaud J.-P. (2025). A Generalized Santa-Fe-like Model to Understand Liquidity Crises. In preparation.

Contents

Ι	Mo	otivation and background 1					
1	Ger	eneral Introduction					
	1.1	Overv	iew of Modern Financial Markets	4			
		1.1.1	Their fundamental role	4			
		1.1.2	The limit order book	6			
		1.1.3	Market ecology	9			
		1.1.4	The liquidity game	11			
	1.2 Stylized Facts of the Order Flow						
		1.2.1	Empirical evidence of trade sign autocorrelation	14			
		1.2.2	Metaorders and the origins of long-memory in order flow .	15			
	1.3 Statistical Properties of Price Changes		16				
		1.3.1	The Louis Bachelier framework	16			
		1.3.2	Market efficiency and price diffusion	17			
		1.3.3	Multifractal dynamics of returns	18			
		1.3.4	Additional empirical regularities in price-volatility dynamics	18			
2	The	eoretic	al foundations for Price Impact	21			
	2.1	What	is price impact and how to measure it ?	22			
	2.2	The o	rigins of price impact theories —Information based models	23			
		2.2.1	The Kyle model	23			

Contents

4		pirical ot Law	Analysis of the Microscopic Foundations of the Square-	- 51
II	Pr	rice In	npact	49
	3.3	The U	nknowns of market stability	46
		3.2.4	Power-laws everywhere?	45
		3.2.3	Phase transition theory	44
		3.2.2	Agent-based models: The Santa Fe approach	43
		3.2.1	Hawkes processes	42
	3.2	Physic	eists' tools to model the Limit Order Book	42
		3.1.3	The excess volatility puzzle	41
		3.1.2	Price jumps in the Limit Order Book	41
		3.1.1	Flash crashes: a symptom of criticality ?	40
	3.1	Unstal	ble markets?	40
3	The	oretica	al foundations for Market stability	39
	2.6	The u	nknowns of price impact	37
	2.5		rt of systematic investing : Alpha versus Impact	34
		2.4.3	The puzzle of the Square-Root Law: A brief review of existing theories	31
		2.4.2	Data, the missing piece	31
		2.4.1	The Square-Root Law	29
	2.4	Metao	order impact —From empiric to models	29
		2.3.3	Aggregated impact	28
		2.3.2	The propagator model	26
		2.3.1	Permanent impact: The Lillo-Farmer model	26
	2.3	Marke	et orders impact and the propagator framework	25
		2.2.2	The Glosten–Milgrom framework	24
		2 2 2	TT1 C1 . 3.511 . 0	~ .

4.1	Introduction				
4.2	Data	description and preliminary observations	54		
	4.2.1	A unique dataset	54		
	4.2.2	General stylized facts about metaorders execution	56		
	4.2.3	Time scales & Market ecology	58		
4.3	Square	e-root impact: micro-scales & meso-scales	60		
	4.3.1	The "double" square-root impact of child orders	61		
	4.3.2	A non-linear propagator model	62		
4.4	From	single market orders to synthetic metaorders	64		
	4.4.1	The impact of single public market orders	64		
	4.4.2	Synthetic metaorders	65		
	4.4.3	Discussion	66		
4.5	The o	ther side of market orders: liquidity providers	67		
	4.5.1	Refill sequences	67		
	4.5.2	Strategic behaviour of liquidity providers	68		
4.6	Concl	usion	70		
		r Proxy: Examining the Puzzling Efficiency of Syntaorder Reconstruction	73		
5.1	Introd	luction	74		
5.2	The al	lgorithm	76		
5.3	Recov	ering metaorder stylized facts	79		
	5.3.1	Peak impact: the Square Root Law	79		
	5.3.2	Role of metaorder duration	81		
	5.3.3	Concave profile during metaorder execution	82		
	5.3.4	Metaorder decay post execution	83		
5.4	Concl	usion	85		

 $\mathbf{5}$

Contents

6			Framework for Market Microstructure: Reconciling e Root Law, Order Flow Dynamics & Price Dynamics 87	7
	6.1	Introd	luction	9
	6.2	A con	tinuous time description of order flow	1
		6.2.1	Model set-up	1
		6.2.2	Average number of metaorders	2
		6.2.3	Average trading activity and trading volume 92	2
	6.3	Order	flow imbalance	3
		6.3.1	Sign Imbalance	1
		6.3.2	Generalized Volume Imbalance	5
		6.3.3	The role of long-range correlations $between$ metaorders 98	3
		6.3.4	Empirical observations)
	6.4 The Impact-Diffusivity puzzle and a generalized propagator		6	
		6.4.1	Price diffusivity within the propagator model 106	3
		6.4.2	A generalized propagator model	3
		6.4.3	The role of metaorder autocorrelations	1
		6.4.4	The role of volume fluctuations	3
		6.4.5	The role of impact fluctuations	5
		6.4.6	Discussion	7
	6.5	Covar	iance between order flow imbalance and prices changes 117	7
		6.5.1	Without volume fluctuations	9
		6.5.2	With correlated metaorders	9
		6.5.3	With correlated metaorders and volume fluctuations 119	9
		6.5.4	With a random impact component	1
		6.5.5	With "informed" metaorders	1
		6.5.6	The correlation coefficient	2
		6.5.7	Empirical data	3
	6.6	Concl	usion	1

	On A	On Alpha Prediction and Permanent Impact				
7		om Abstraction to Animation: Artificial Market Simulator 137				
	7.1	Introd	uction	138		
	7.2	A brie	of reminder of the generalized propagator model	140		
	7.3	How t	o simulate our model?	142		
	7.4	Empir	rical stylized facts vs. simulations	145		
		7.4.1	The q -dependence of the autocorrelation of trades \dots	145		
		7.4.2	The scaling of the order flow imbalance	146		
		7.4.3	Recovering a diffusive price	149		
		7.4.4	Aggregated impact and anomalous rescaling	151		
		7.4.5	The covariance coefficient	152		
		7.4.6	The correlation coefficient	155		
	7.5 The puzzling effectiveness of proxy metaorders					
		7.5.1	How the acceptance window drives mapping accuracy	161		
	7.6	Conclusion				
II	I N	⁄Iarket	Stability	165		
8			cture modes in the Limit Order Book: of Marginal Instability	167		
	8.1	Introd	uction	168		
	8.2	Data 1	presentation	170		
		8.2.1	Variables of interest and intraday profile	171		
		8.2.2	A second coarse-graining	173		
		8.2.3	Box-Cox transformation	174		
	8.3	Micros	structure modes	175		
		8.3.1	PCA Analysis I: Raw data	175		

Contents

		8.3.2 PCA Analysis II: Binned data	7		
	8.4	A VAR model for flow dynamics	8		
		8.4.1 1-lag VAR model	9		
		8.4.2 Multi-lag VAR model	1		
	8.5	An attempt to model price impact	4		
	8.6	Conclusion and further discussions	8		
9	_	eneralized Santa Fe-like model to study liquidity crisis in the it Order Book 19	1		
	9.1	Introduction	2		
	9.2	Motivation : The initial Santa Fe model with feeedback 19	4		
	9.3	A 4 degree of freedom Santa Fe model	5		
	9.4	Investigating phase transitions in the 4DF Sante Fe model 19 $$	7		
		9.4.1 Stability maps	7		
		9.4.2 Finite-size scaling	8		
		9.4.3 Overview of exponents values	9		
		9.4.4 Study of spread explosion	2		
	9.5	Analytical Determination of the Stability Boundary 20	3		
	9.6	Conclusion	8		
IV	/ C	onclusion and Perspectives 21	1		
R	efere	nces 220	0		
		dices 23	5		
\mathbf{R}	Résumé substantiel en français 23				

			Contents
A	Chapter 4:	the Metaorder Proxy	239
В	Chapter 8:	the VAR model	243
\mathbf{C}	Chapter 9:	the Santa-Fe like model	245

Part I Motivation and background

Chapter 1

General Introduction

An efficient market is one in which price is within a factor 2 of value, ie, the price is more than half of value and less than twice value.

Fisher Black

Contents

Contents			
1.1	Ove	rview of Modern Financial Markets	4
	1.1.1	Their fundamental role	4
	1.1.2	The limit order book	6
	1.1.3	Market ecology	9
	1.1.4	The liquidity game	11
1.2	Styl	ized Facts of the Order Flow	13
	1.2.1	Empirical evidence of trade sign autocorrelation \dots	14
	1.2.2	Metaorders and the origins of long-memory in order flow	15
1.3	Stat	istical Properties of Price Changes	16
	1.3.1	The Louis Bachelier framework	16
	1.3.2	Market efficiency and price diffusion	17
	1.3.3	Multifractal dynamics of returns	18
	1.3.4	Additional empirical regularities in price-volatility dy-	
		namics	18

October 14th, 2013 —the Nobel Prize in Economic Sciences is awarded to three titans of financial thought: Eugène Fama, Lars Peter Hansen, and Robert Shiller, each honored for their pioneering contributions to our understanding of asset prices. And yet, beneath the surface of this shared recognition lies a striking paradox: their conclusions, though equally celebrated, stand in stark opposition.

Eugène Fama, following in the intellectual lineage of Louis Bachelier, is best known for formalizing the Efficient Market Hypothesis (EMH), according to which financial markets fully and instantaneously reflect all available information. In contrast, Robert Shiller, a leading figure in behavioral finance, argues that prices are often driven by irrational behavior, speculative bubbles, and psychological biases—forces that challenge the notion of market efficiency.

This enduring tension lies at the heart of modern financial economics and motivates the central question of this thesis: Why do prices move? To approach this question, we adopt a physicist-inspired methodology, in the spirit of Econophysics, seeking to infer macroscopic phenomena from underlying microscopic dynamics. In physics, this typically means studying the behavior of particles or spins to understand processes such as gas dynamics or magnetism. In finance, the comparable microscopic framework is market microstructure—the field that examines the mechanisms and frictions through which trading occurs and prices are formed.

When reviewing the literature and providing the state-of-the-art models in this area, it is difficult to surpass the seminal book *Trades*, *Quotes and Prices*: *Financial Markets under the microscope* [7]. It has served not only as a foundational reference for this thesis, but also as a true bedside book over these three years. Thus, I will aim to provide a concise introduction aligned with the work conducted over these three years, while interested readers are encouraged to consult this reference for more detailed content.

1.1 Overview of Modern Financial Markets

1.1.1 Their fundamental role

Financial markets play a central role in modern economies, even though their underlying principles date back centuries. Broadly speaking, their purpose can be summarized in three functions: facilitating the exchange of asset ownership between economic agents, determining asset values through transaction prices, and enabling companies to raise capital on a global scale. I will briefly mention the earliest form of trade —barter —which appears to be present even in the earliest human societies. The main drawback of barter is its extreme *illiquidity*: if one party isn't interested in what the other is offering, the transaction does not take

place —no matter how many goats you're willing to trade for a sack of grain. No mutual interest, no deal. That's why the use of an intermediary—such as gold or money —quickly established itself as the easiest way to trade. Indeed, it is much easier to agree on the value of a coin than to find someone who simultaneously wants your goat and happens to own a sack of grain. Let then to the earliest documented financial markets, which likely emerged in the early 17th century with the Initial Public Offering (IPO) of the Dutch East India Company on the Amsterdam Stock Exchange. This event marks the beginning of financial markets resembling those we know today, where agents can not only exchange goods but also invest in publicly traded companies. Again, one of the main challenges was finding liquidity—that is, a buyer willing to pay for the shares you wanted to sell. At the time, when stock ownership was mostly limited to the nobility, the usual method involved sending a valet to roam the city's pubs, loudly calling out in search of a willing counterpart.

Even if this form of market may seem far removed from the low-latency environments we are used to today (where transactions are typically effected at the microsecond level), it is still interesting to see that they weren't so different from modern markets, especially when it comes to bubble or crashes. For example, consider well-known events such as the Tulip Mania or the South Sea Bubble, during which even Isaac Newton famously suffered significant financial losses [8]. That been said, a quick look at the backgrounds of employees in today's major financial institutions suggests that physicists may have taken their revenge...

In this thesis, we will focus in particular on the second fundamental role of financial markets: assigning value to things through the disclosure of prices. To that end, let us first describe the well-known Efficient Market Hypothesis (EMH). The central question underlying the theory is: what is the best way to agree on the value of a given asset without relying on any central authority (such as legal or political power) to impose it? The answer is the following: bring everyone to the table and ask each participant to propose a price. Suppose that agents are independent, identically distributed (we will revisit this assumption in more technical terms in Section 1.3), rational and reasonably well-informed. Then, if we denote the fundamental value of the asset by p_0 , each agent submits a price $p_i = p_0 + \eta_i$, where η_i is white noise—that is, a random individual bias. By averaging, it indeed may allow on to estimate the fundamental value p_0 .

Building on this mechanism, and under the three standard assumptions of economic theory—agent rationality, perfect information, and utility maximization—the Efficient Market Hypothesis asserts that market prices are a reliable proxy for "fundamental value" of the underlying asset. Prices are assumed to instantaneously incorporate all available information at any given time. One consequence of this belief is the idea that markets should be in equilibrium: large price devi-

ations are mainly caused by exogenous events (i.e., external to the market), and markets should subsequently return to equilibrium. Another consequence—which may be questioned at first glance by observing the repartition of economic activity today—is that it should be impossible to systematically extract profits from the market.

Before going more deeply into these theories—which we will analyze quantitatively later—let us first outline the structures and principles governing the microscopic world of modern electronic financial markets.

1.1.2 The limit order book

Although many forms of asset exchange have emerged—especially over the past decade with the rise of blockchain and crypto-assets, which are typically traded in dark pools—two main conventional mechanisms still dominate today's markets. The first is Over-the-Counter (OTC) trading, typically used for specific asset classes such as large bond issues, complex derivative products, real estate etc. In this setting, transactions take place directly between buyer and seller, outside of a centralized exchange, and are often negotiated bilaterally. Information related to the transaction is usually not publicly disclosed, and liquidity in these markets is often provided by specialized institutions, commonly referred to as dealers. In this thesis, we focus exclusively on the second mechanism: assets traded on centralized exchanges via the limit order book. Those centralized orderbooks are usually used for liquid assets, such as stocks or futures. As we will study limit order books in detail later, let us begin by providing a thorough introduction.

Market venues: A given asset is usually traded on different stocks exchanges, and within a exchange it can also be traded on different limit orders books. All these possibilities are referred to as market venues, where liquidity is distributed. For instance, if a fund seeks to acquire a large quantity of a given stock, it is likely to execute portions of the order across multiple venues—buying part of it on the London Stock Exchange, another part on NASDAQ, and so on. Thanks to the digitization of financial markets, it has become increasingly easy to trade directly across multiple stock exchanges. Ensuring that prices remain consistent across these venues is typically the role of arbitrageurs (for example low latency market makers), as we will discuss later. Since arbitrageurs tend to operate with high efficiency, it is also a common practice to consider the aggregated limit order book—that is, the consolidation of all limit order books for a given asset across different venues. However, when relevant, we will also specify cases where data from a particular market venue has been used.

Types of orders: Even though some specific orders types could exist depending on the market venues (iceberg, block trades etc), one can list three basic universal kinds of order:

- Limit orders: These represent an intention to trade a specific quantity at a specified price. Once submitted to the exchange, a limit order is stored in the order book and made visible to all market participants. It constitutes the revealed liquidity of the book. If the trader wants to buy (resp. sell) the asset, the order is placed on the bid (resp. ask) side.
- Cancellation orders: These enable a trader to cancel an existing limit order in the order book. Naturally, a cancellation can only remove a limit order placed by the same trader.
- Execution (market) orders: These are used by a trader to actually exchange the asset by matching with revealed limit orders. Typically, the trader specifies a trade quantity, which is executed starting from the best available offers and continuing through less favorable ones if needed.

Type of auction: Different types of auction mechanisms are used across markets, depending on the exchange and the nature of the traded asset. We briefly describe the two most common formats: the single auction and the continuous double auction.

In a *single auction*, the limit order book evolves throughout the day—traders can post and cancel limit orders—but actual trades occur only at a predetermined time. For example, a trader may submit a market order (at 10 a.m. for example), yet retain the ability to cancel it before the auction is executed (at 12 p.m.). This type of mechanism, while rich in stylized facts, see [9] for reference, is not the focus of our study.

Instead, we concentrate on the *continuous double auction*, which is the standard mechanism in most modern stock exchanges. In this setting, both limit and market orders can be submitted at any time during trading hours. Transactions occur in continuous time: whenever a market order matches a limit order, the trade is executed immediately. As such, market orders are definitive—once submitted, they cannot be canceled. This continuous interaction between supply and demand shapes the dynamics of the limit order book in real time, and also define a real transaction price.

Finally, it is worth noting that some exchanges commonly combine both types of auction mechanisms. For example, the Tokyo Stock Exchange (which we will discuss in more detail in Section 4) operates with an opening and closing auction

at the start and end of the trading day, respectively, while the continuous double auction mechanism governs trading during the rest of the session.

Available information: Most of financial institutions working in the High Frequency field not only look at past transaction prices (it is called L2 data), but directly at the live limit order book (L3 data), which provide much more information than prices only, as it give one a hint of the live level of supply and demand. Indeed, it is well known for example that volume imbalance $\frac{Q^{\text{bid}}-Q^{\text{ask}}}{Q^{\text{bid}}+Q^{\text{ask}}}$ is a good predictor of the sign of the next price change, see [10]. While it is generally possible to record the price, quantity, and timestamp of each order in the order flow, the identity of the order's initiator remains unknown—except in a few highly specific market settings. Anonymity is deliberately maintained by the stock exchange as a core feature of the trading infrastructure. That being said, the availability of additional information can also open the door to new forms of market manipulation. One of the most well-known manipulative techniques is spoofing, extensively studied in [11]. The strategy works as follows: suppose you wish to sell at a favorable price. You place a genuine limit order on the ask side. To increase the chances of execution at that price, you simultaneously submit a very large buy limit order on the bid side, creating the illusion of strong buying interest. This artificial demand can move the market toward the ask, allowing your sell order to be executed. The bid-side order is then canceled before it can be matched. This tactic is, of course, illegal, as it misleads other market participants by generating false signals of demand or supply. Yet it turns out to be highly effective, as many traders heavily rely on the information visible in the limit order book.

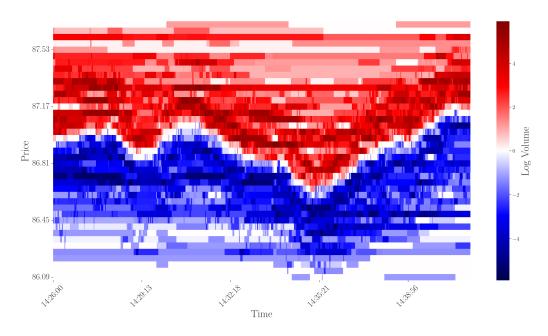


Figure 1.1: Snapshot of the order book for the U.S. stock AHAA-UQ, traded on the NYSE on 2022-11-02. Red indicates the volume available on the ask side, while blue shows the volume on the bid side. Prices are expressed in U.S. dollars.

Variables of Interest in the Limit Order Book The first key variable is the $mid\ price$, defined as the average of the best ask price p_a and the best bid price p_b . Closely related is the spread, which measures the difference of price between these best quotes. Introducing the $tick\ size\ \psi$, the smallest possible price increment in the limit order book (LOB), allows one to classify stocks as either $large\ tick$ or $small\ tick$.

Although somewhat simplified, this classification is based on the number of ticks—that is, discrete price levels—contained within the spread. For large tick stocks, this number is typically close to one, meaning the spread equals one tick: $p_a = p_b + \psi$. Conversely, for small tick stocks, the spread covers multiple ticks, usually two or three: $p_a = p_b + k\psi$ with $k \approx 2$ –3. Choosing an appropriate tick size is a subtle and important regulatory decision, see [12].

1.1.3 Market ecology

Now that we have a reasonable understanding of how the microstructure of financial markets works, let us introduce the key players—that is, the main institutions operating in these markets. This section has a twofold objective. First, it contributes to the broader understanding of modern financial markets. Second, we

will see that some specific participants either help stabilize the market or are often blamed for various stylized facts. For instance, in Chapter 4, we will use specific data to challenge part of a still very solid theory that attributes the so-called square-root law to the behavior of a particular group of market participants.

Participants:

The following list is by no means exhaustive, but it reflects the classification used in Chapter 4, where a more detailed and technical discussion is also provided. We divide market participants into four broad categories and simply list them here.

- Low-Frequency Traders. These are probably the most numerous participants in financial markets, although they account for only a small share of daily trading activity (including limit and cancellation orders). Typical examples include pension funds, asset managers, and some hedge funds that take directional positions and hold them over long periods. However, the label can be misleading: being a low-frequency trader (at least in our sense) does not mean that these participants never interact with markets at high frequency. Rather, it means their investment horizon is long, even though their execution is necessarily done in the high-frequency world—since acquiring a large position with a single order is virtually impossible, but would also be very costly (see next paragraph about the liquidity game). Therefore, low-frequency traders often either have an internal high-frequency execution team, or they delegate execution to a broker.
- Brokers. A broker is a trader who executes orders on behalf of clients—typically low-frequency ones. For example, an asset manager may want to buy 1,000 shares of a given company within two days. They contact a broker, whose job is to find counterparties in the market. The broker typically commits to an average price and takes care of the execution over time.
- Market Makers. These actors play a particularly interesting role in financial markets. Historically, market makers were hired by exchanges to reduce the spread between buyers and sellers. To illustrate, consider a market where one participant wants to sell an asset for 12€, while another is willing to buy it for 8€. The mid-price is 10€, but without anyone bridging the gap, no trades occur. This situation is harmful in two ways: first, the market becomes illiquid, making price discovery difficult and discouraging trading due to the risk of being unable to exit a position. Second, a wide spread is costly—whichever side initiates the trade must "cross the spread," potentially losing 4€ in this example.

The market maker solves this by quoting on both sides—say, buying at 9€ and selling at 11€. This provides better prices to both buyer and seller, while allowing the market maker to earn a 2€ spread if both sides are filled. Although this might sound like an easy "win-win" situation, the reality is quite the opposite. Today's market environment—with strong competition and significant price impact—makes profitability challenging for market makers. To summarize, market makers are usually paid by the exchange, to constantly quote, at a defined frequency, on both sides of the orderbook. They try to earn money from the spread while managing their inventory. we will come back later to the ongoing debate about whether market makers are beneficial or harmful to market stability, and try to offer some insights on the matter.

• **High-Frequency Traders.** This label can also be confusing, since brokers and market makers often operate at high frequency as well. Here, we define high-frequency traders as institutions that are neither market makers nor brokers, but that act on short-term signals and hold positions for very brief periods. These are typically systematic hedge funds operating at *intraday frequencies*. We use the term "systematic" because this type of trading must be automated at such time scales—unlike "discretionary" trading, where humans are still fast enough to make decisions themselves.

1.1.4 The liquidity game

Now that we have introduced the rules, the framework, and the main players, let us add a bit of strategy before diving into the technical details. As one might expect, the limit order book is the stage for fierce competition between agents, which we divide into two broad categories—even though the line between them is often blurred. On one side, we have the *liquidity providers*. These traders place limit orders in the book, aiming to trade at favorable prices since they don't have to cross the spread.

Think back to our simple example: it is clearly better for a seller to get executed at 12€ than to initiate a trade at 8€. However, placing limit orders isn't as easy or profitable as it might seem—it actually requires a great deal of skill, for three main reasons.

First, by revealing their intentions, liquidity providers expose themselves to adverse selection. In other words, they give away information about their own valuation of the asset, but when they get executed, they learn very little about what the counterpart knows. For instance, if you post a very large sell limit order at 12€ and it gets immediately filled, it might mean your price was too low—buyers

may have had better information than you, and you just got picked off.

The second issue is opportunity cost. In this ultra-fast environment, placing a limit order instead of executing immediately can turn out to be costly if the market moves away from you—especially once others have seen your order sitting in the book. In other words, it might have been better to trade at 8 right away than to quote 12 and not being executed, if the new best bid drops to 6 just after.

The final challenge we briefly highlight is the so-called *queue race*. In highly liquid markets, it is uncommon for a liquidity provider to be alone at a given price level. Consequently, there exists a persistent competition to secure priority in the order queue—that is, to be the first to be executed. This race, often done at microsecond timescales, is a key factor underlying of "need for speed" in modern financial markets [13]. It has driven substantial investments in low-latency infrastructure, such as dedicated fiber-optic or microwave transmission lines, particularly among high-frequency trading firms seeking even marginal timing advantages.

The second group consists of *liquidity takers*, who remove liquidity from the order book by executing outstanding limit orders. We have already discussed the spread cost associated with this type of execution, and the second part of this thesis will be devoted to analyzing its broader impact—an additional and often substantial component of transaction costs.

While the literature commonly distinguishes between these two groups—liquidity providers and liquidity takers—on the grounds that they are, broadly speaking, composed of different types of market participants, this separation is not always clear-cut. Indeed, liquidity provision is typically associated with market makers, whereas liquidity taking is more often carried out by brokers or lower-frequency institutional investors such as hedge funds. We will examine this paradigm more quantitatively in Chapter 4, but it is worth noting that in practice, most execution strategies rely on a combination of both roles. For instance, a fund seeking to acquire shares of a given asset may simultaneously execute market orders at the ask while placing limit orders at the bid. Some well-known hedge funds have even reported executing exclusively through limit orders, see [14], though being traditionally part of the buy side.

Latent Liquidity

One of the key consequences of the strategic nature of trading is that the majority of liquidity remains latent—that is, not yet revealed in the order book. This is a fundamental characteristic of limit order books: the visible liquidity at any given time typically represents only about $10^{-3} - 10^{-5}$ of the total daily traded volume. Liquidity is progressively revealed throughout the trading day as a result

of strategic interactions. This distinction is crucial when studying market impact, as relying solely on the observed (revealed) order book can be misleading. In reality, impact is more accurately understood in relation to the latent order book; see [15] and Chapter 2 for further discussion. This notion of latent liquidity also helps explain the splitting behavior that we will examine later. Suppose, for instance, that an investor wishes to acquire 1% of a company's shares. It is plausible that there exists a counterpart—or a set of counterparties—willing to sell such a quantity. However, because neither party wishes to reveal their intentions, the transaction must occur gradually. The buyer must acquire the shares incrementally, consuming the available revealed liquidity bit by bit rather than all at once, so as to avoid detection. As a result, completing the transaction may take several days. This illustrates a subtle interplay between the investor's predictive insight (i.e., the motivation behind the purchase) and the execution strategy used to implement the trade. As we will show in Chapter 2, the execution process itself plays a much more significant role than one might initially expect.

We are now concluding this first part of the introduction, which has provided a concise overview of the functioning of modern electronic markets. While not exhaustive, it highlights the key elements that will be central to the rest of this thesis.

At the heart of these markets lies the limit order book, the arena for a complex and captivating competition. On one side, long-term investors aim to build their positions discretely, seeking to avoid detection and minimize market impact. On the other hand, high-frequency institutions that actively manage their inventory—aiming to keep it close to zero—engage in rapid trading strategies to extract profit from the information revealed by these investors' actions.

While market microstructure offers a rich landscape of strategic interactions, it is also characterized by striking and well-documented stylized facts. We now turn to a brief overview of these empirical regularities.

1.2 Stylized Facts of the Order Flow

The order flow consists of the sequence of messages—limit orders, cancellations, and market orders—submitted to the exchange for a given asset. It is closely monitored by market participants, as one can view the market as a black box: it takes order flow as input and produces prices as output. The challenge, of course, lies in the fact that the internal mechanisms of this black box are highly complex. Understanding these mechanisms—at least partially—is one of the goals of this thesis.

For the sake of simplicity, we will focus in this introduction on the stylized facts of trade flow, which are among the most robust and well-documented regularities in market microstructure. We will examine limit and cancellation flows in more detail in Chapter 8, though these tend to be less universal and more dependent on market-specific behaviors. Indeed, limit orders and cancellations reflect trading intentions rather than actual transactions. As such, their statistical properties are often affected by high-frequency strategies—such as *jittering* ¹, where limit orders are placed and canceled in rapid succession—making their interpretation more delicate. In contrast, trades (executions) generally provide clearer and more interpretable signals of market activity.

1.2.1 Empirical evidence of trade sign autocorrelation

Trade-by-trade data is among the most accessible sources in market microstructure and is often the first dataset encountered when analyzing financial markets at high frequency. As a result, it has been extensively studied over the past 40 years. While many of the underlying mechanisms driving market behavior remain mysterious — a central motivation of this thesis — there is broad consensus on one key empirical fact: trade signs exhibit significant autocorrelation.

But what does this mean? Let ε_t denote the sign of a trade executed at time t (i.e., +1 for a buyer-initiated trade and -1 for a seller-initiated one). Autocorrelation implies that ε_t statistically influences the sign of future trades $\varepsilon_{t+\tau}$, even for large time lags τ . More intriguingly, the autocorrelation function of trade signs is observed to decay as a power law with respect to τ , indicating a long-memory process — past order flow continues to impact future order flow over surprisingly long timescales. Indeed, as evidenced by the data, see Chapter 6 for example, the autocorrelation function C(l) reads:

$$\mathbb{E}[\varepsilon_t \varepsilon_{t+\tau}] \sim \frac{c_0}{\tau^{\gamma}} \tag{1.1}$$

For most assets, the exponent γ typically lies in the range [0.4, 0.7], and $c_0 \approx 0.5$ for most of liquid assets. A detailed study of γ values for stocks listed on the Tokyo Stock Exchange, along with a technical discussion on how to compute an unbiased estimate of this exponent is provided in [16].

To illustrate the implications of this power-law behavior, let us consider the example introduced in [7]. Using realistic parameters $c_0 = 0.5$ and $\gamma = \frac{1}{2}$, the autocorrelation at a lag of 10,000 trades evaluates to approximately $C(10,000) \approx 0.005$. This indicates that when a buy market order is executed, the likelihood that another buy order will occur 10,000 trades later surpasses the probability of it being

¹Place and remove immediately orders in the spread, for strategic reasons

a sell order by more than 0.5%. Consequently, the trade flow exhibits a notable degree of predictability, even over extended timescales. This predictability seems to be an apparent contradiction with the concept of market efficiency, which asserts that price movements themselves are unpredictable. This phenomenon is often referred to as the *efficiency paradox* and will be a focal point of discussion in Chapter 2.

1.2.2 Metaorders and the origins of long-memory in order flow

Two primary mechanisms have been proposed to explain this phenomenon. The autocorrelation of trades may result from either *herding* or *splitting*. Herding refers to the behavior where traders follow trends: for instance, a buy trade can trigger a series of copy trades, as other traders want to imitate the initial trade and replicate it with further transactions. This explanation is realistic, as it is well-known that traders often do not act independently, and this behavior could potentially explain the formation of bubbles and crashes, as discussed in Chapter 3. Another aspect of herding is the distinction one can make between trend and value investors, see [17] for a more in-depth analysis.

Another perspective attributes this behavior to trade splitting. Due to the dynamics of the liquidity game and the latent liquidity it involves, traders are often forced to divide their initial orders — referred to as *metaorders* — into several smaller parts that are executed sequentially, known as *child orders*.

This assumption was initially introduced in the seminal paper by [18], and it is widely recognized in the literature as the LMF hypothesis, a term we shall also employ throughout this thesis. The authors propose a tractable quantitative framework that not only elucidates the power-law behavior of autocorrelation but also establishes a connection with the size of metaorders. Specifically, if we assume that metaorder sizes s are distributed according to a power-law distribution, $\Psi(s) \sim s^{-1-\mu}$, it can be demonstrated that such splitting results in an autocorrelated flow characterized by $C(\tau) \sim \tau^{-(\mu-1)}$. Why would the metaorders distributions be a power law ? Here again, this assumption could be nicely linked to the power-law distribution of company sizes, a phenomenon both predicted and empirically validated in numerous economic studies, see [19] for example.

Although this debate had persisted for some time—largely due to the inherent difficulty of observing the distribution of *all* metaorders in the market, as such information is typically proprietary—it was resolved in [16], thanks to a unique dataset that we also had the opportunity to work with (see Chapter 4). In that study, the authors gained access to the complete set of metaorders on the Tokyo Stock Exchange, allowing them to conclusively validate the splitting hypothesis.

While the order flow exhibits a much richer set of stylized facts—which we will explore in detail in Chapter 6—it is fair to say that, until now, most theoretical and empirical studies were primarily based on two key observations: trade signs are autocorrelated, and metaorders are embedded within the broader flow of orders. Naturally, inferring metaorders from the raw order flow is extremely challenging (despite many attempts, see [20, 21] for example, as such insights could lead to highly profitable strategies). In this sense, markets preserve a form of efficiency. Let us now turn to the central object of interest: the price process.

1.3 Statistical Properties of Price Changes

I will keep this section concise, but I strongly encourage interested readers to consult Jean-Philippe Bouchaud's lectures at College de France—both the manuscript and the video recordings—which are publicly available here: [22].

1.3.1 The Louis Bachelier framework

It is impossible to begin a rigorous study of price dynamics without acknowledging Louis Bachelier, a *French mathematician* whose 1900 doctoral thesis [23] laid the foundations of quantitative finance. In his work *Théorie de la Spéculation*, Bachelier introduced two core ideas:

- Asset price changes are fundamentally unpredictable, exhibiting random behavior.
- The accumulation of many small, independent price changes should, by the Central Limit Theorem, lead to a Gaussian distribution of returns.

This led Bachelier to model the price p(t) as a Brownian motion - 5 years before Albert Einstein :

$$p(t) = p(0) + \sigma W(t), \tag{1.2}$$

where W(t) is a standard Wiener process and σ the volatility. This implies that price increments $p(t) = p(t + \Delta t) - p(t)$ are independent and normally distributed:

$$\Delta p(t) \sim \mathcal{N}(0, \sigma^2 \Delta t).$$
 (1.3)

While elegant, this framework fails to capture key empirical features of financial time series, such as fat-tailed return distributions—price changes are not, in reality, Gaussian—and volatility clustering, whereby price changes exhibit temporal correlations. It is worth noting that this simplified model is nonetheless the basis the well-known Black-Scholes framework, widely used for option pricing (i.e., bets on the future value of an asset) and various risk modeling applications. Although

it is now widely acknowledged that the Black-Scholes model does not accurately describe real market behavior, it remains a standard reference and common language in the field of option pricing.

1.3.2 Market efficiency and price diffusion

Bachelier's insights anticipated the Efficient Market Hypothesis (EMH) formalized 70 years later. The EMH implies that asset prices incorporate all available information, and thus price changes are unpredictable. Mathematically, if \mathcal{F}_t is the information set at time t, then:

$$\mathbb{E}[p(t+\Delta t)|\mathcal{F}_t] = p(t). \tag{1.4}$$

Empirically, this translated in a very natural property of prices: they are diffusive.

$$Var[p(t + \Delta t) - p(t)] \propto \Delta t, \tag{1.5}$$

Although this property may appear simple at first glance, it is both theoretically challenging to reproduce—see Chapter 6—and fundamental to the functioning of financial markets. Indeed, if prices exhibited even slight deviations from diffusive behavior—being subdiffusive (where a positive price change is more likely to be followed by a negative one) or superdiffusive (where positive changes tend to be followed by further positive ones)—then the market would become arbitrageable mathematically speaking. In such cases, it would be possible to generate statistically systematic profits, violating one of the core principles of modern financial theory.

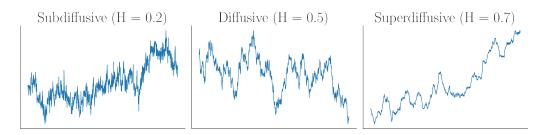


Figure 1.2: Example of three different processes, each characterized by their Hurst exponent. Financial markets are often considered efficient because their Hurst exponent is around the fragile equilibrium value of 0.5. Indeed, empirically, after a brief period of mean reversion, prices follow a nearly perfect diffusive process - see Chapter 7, rendering them essentially unpredictable. This property is very subtle to obtain, see Chapter 6

1.3.3 Multifractal dynamics of returns

Empirical studies have long shown that return distributions are not only heavytailed but also exhibit non-trivial scaling behaviors that are incompatible with monofractal. Notably, the moments of aggregated returns over varying time horizons scale anomalously as

$$\mathbb{E}[|r_t^{(\tau)}|^q] \propto \tau^{\zeta(q)},\tag{1.6}$$

where $\zeta(q)$ is a nonlinear function of q.

These observations are closely related to the phenomenon of volatility clustering, whereby large price movements tend to be followed by large movements (of either sign), and small movements by small ones—indicating persistent temporal dependence in volatility. While volatility clustering reflects the autocorrelation structure of volatility over time, multifractality captures a deeper, scale-invariant organization in the distribution of returns, extending the notion of temporal dependence to all statistical moments.

The idea of modeling financial time series with multifractal properties was notably pioneered by Benoît Mandelbrot, see [24], who proposed that price dynamics might follow a multifractal process rather than a Brownian motion. Building on this intuition, Bacry et al. introduced the Multifractal Random Walk (MRW), see [25], in which the return process is defined by a stochastic volatility model with long memory:

$$r_t = e^{\omega_t} \epsilon_t, \tag{1.7}$$

where ϵ_t is a standard white noise and ω_t is a Gaussian process with long memory.

Such models successfully replicate most of the statistical properties of prices, for single asset or even indices, see a recent study in [26]. They are also closely related with *rough volatility models* proposed in [27]. A microstructural interpretation of these multifractal properties will be partially explored in Chapter 6.

1.3.4 Additional empirical regularities in price-volatility dynamics

A key stylized fact in the dynamics of asset prices is the presence of temporal asymmetries linking returns and volatility, at least for stocks. The *leverage effect* refers to the empirical observation that negative returns tend to be followed by increases in volatility, whereas positive returns have a comparatively weaker effect. This phenomenon can be quantified by the leverage correlation function

$$L(\tau) = \mathbb{E}[r_t \sigma_{t+\tau}^2],\tag{1.8}$$

which typically takes negative values for small positive lags τ , indicating that past returns influence future volatility in an asymmetric way. This effect is often

interpreted through mechanisms such as volatility feedback, dynamic hedging, or investor risk aversion.

A closely related but subtly different phenomenon is the *Zumbach effect*, which captures the breaking of time-reversal symmetry in volatility dynamics. It refers to the empirical observation that past price trends —whether upward or downward —tend to increase future volatility, but the reverse effect is less important:

$$cov(\sigma_t^2, (R_{t,\tau})^2) < cov(\sigma_t^2, (R_{t,-\tau})^2), \tag{1.9}$$

Interested readers may refer to [28] for the theoretical framework, and to the excellent empirical analysis by [29], which will also serve as a foundation for Chapter 9.

Chapter 2

Theoretical foundations for Price Impact

Richard P. Feynman

Contents

Comen	5			
2	2.1	Wha	at is price impact and how to measure it?	22
2	2.2	The	origins of price impact theories —Information-	
		base	d models	23
		2.2.1	The Kyle model	23
		2.2.2	The Glosten–Milgrom framework	24
2	2.3	Mar	ket orders impact and the propagator framework	25
		2.3.1	Permanent impact: The Lillo-Farmer model	26
		2.3.2	The propagator model	26
		2.3.3	Aggregated impact	28
2	2.4	\mathbf{Met}	aorder impact —From empiric to models	29
		2.4.1	The Square-Root Law	29
		2.4.2	Data, the missing piece	31
		2.4.3	The puzzle of the Square-Root Law: A brief review of	
			existing theories	31
2	2.5	The	art of systematic investing: Alpha versus Impact	34
2	2.6	The	unknowns of price impact	37

2.1 What is price impact and how to measure it?

Price impact is arguably one of the most fascinating phenomena—at least from my perspective—and unquestionably one of the most critical for anyone engaged in trading, as it can transform an apparently profitable strategy into a money losing one: it is one of the underlying reasons behind the limitations — or the challenges — of backtesting. Moreover, price impact is one of those specific topics where academic research and industry practice closely intersect, since grasping it is essential both for practitioners and for researchers seeking to uncover the mechanisms behind price formation. A simple definition of price impact is the variation in price caused by someone interacting with the market. While it is well known that limit orders and cancellations can also impact prices (see [7], Chapter 13), most studies focus primarily on the impact of market orders, as they are directly associated with transactions. In line with this common approach, we will also concentrate in this introduction—and throughout the thesis—on the impact of market orders.

While the concept of price impact is easy to state, its rigorous measurement is more subtle. One must carefully specify: the impact of what, on which price, and measured when. Let us consider a simple example to illustrate this. Suppose a trade of sign $\varepsilon_t \in \{+1, -1\}$ is executed at time t, where $\varepsilon_t = +1$ denotes a buy order and $\varepsilon_t = -1$ a sell order. Let m_t denote the mid-price immediately before the trade.

Then, the price impact at a time lag $\tau > 0$ can be defined as the difference between the expected mid-price at time $t + \tau$ given that the trade occurred, and the hypothetical (counterfactual) mid-price that would have been observed at the same time had the trade not taken place:

$$\mathcal{I}(\tau) = \mathbb{E}[m_{t+\tau} \mid \varepsilon_t] - \mathbb{E}[m_{t+\tau} \mid \text{no trade at } t]$$

This formulation highlights a key difficulty: while the first term can be estimated from market data, the second term—what the price would have been in the absence of the trade—is unobservable and very difficult to obtain. Thus, when studying price impact, one has to be satisfied with an approximation of the real price impact. We then define for the rest of the thesis price impact as:

$$\mathcal{I}(\tau) = \mathbb{E}[m_{t+\tau} - m_t \mid \varepsilon_t]$$

Considerable effort has been devoted to developing frameworks for generating hypothetical market scenarios in order to answer a fundamental question: What

would the price have been if I had traded at a specific time? This lies at the heart of backtesting, which aims to simulate the outcome of a trading strategy based on historical data. Its reverse formulation—estimating the impact of a trade—relies on the same logic: inferring how the price would have evolved in the absence of that trade. We explore this challenge in Chapter 8, but as we will see, accurately reproducing the relevant stylized facts proves to be extremely difficult (see [21]). That said, although we have not tested it, we believe that the framework introduced in Chapters 6 and 7 holds promise for successfully reconstructing such a market.

2.2 The origins of price impact theories —Information-based models

The existence of price impact is, at once, both intuitive and paradoxical. From a basic supply-and-demand perspective, the logic is straightforward: scarcity drives value. If there is persistent buying pressure—ie one executed a trade at time t—the price should rise accordingly. However, this notion appears to contradict the Efficient Market Hypothesis (EMH), which asserts that prices instantly and fully reflect all available information. Under this view, price changes should come from variations in the asset's fundamental value, not from the act of trading itself. However, if a trader is perfectly informed, a change in the asset's fundamental value may be revealed through their trading activity. In this sense, the trade itself becomes the mechanism by which information enters the market. This idea lies at the heart of the following price impact model, that could be seen as the very first theory of market microstructure.

2.2.1 The Kyle model

The Kyle (1985) model provides a foundational framework in which price impact arises from asymmetric information. A single informed trader – knowing exactly the fundamental value v of an asset – trades continuously in a market alongside uninformed noise traders and a competitive market maker. The market maker sets prices based on aggregate order flow Y = x + u, where x is the insider's informed order and $u \sim \mathcal{N}(0, \sigma_u^2)$ a trading noise, due to uninformed traders executions.

Kyle shows that the equilibrium pricing rule is linear: to break even, the market maker should set: $p = \mathbb{E}[v \mid Y] = \lambda Y$, with impact parameter

$$\lambda = \frac{\sigma_v}{2\sigma_u},$$

where σ_v^2 is the variance of the fundamental asset value v. λ is often referred as

the Kyle-Lambda. The informed trader chooses x to maximize profits,

$$\max_{x} x(v - p) = x(v - \lambda(x + u)),$$

leading to $x = \beta(v - p)$. The model yields a linear and permanent price impact: $\Delta p = \lambda x$.

Kyle's framework provides an elegant link between price impact and the gradual revelation of private information, and it has long served as a foundational reference in market microstructure theory. It also open the way for the influential conceptual distinction between *informed traders* and *noise traders*, which has inspired a substantial body of subsequent research. While the model is theoretically appealing, two important concerns can be raised. First, the boundary between informed and uninformed trading is often ambiguous. A trader acting on a long-term signal—e.g., over a two-year horizon—may appear indistinguishable from a noise trader when viewed over intraday timescales, especially if they trade while prices are declining. Second, these models often assume that the aggregate impact of noise traders averages out, contributing no lasting price effect. However, in practice, trades devoid of fundamental information can exhibit herding behavior or follow misleading signals, potentially reinforcing price movements. This raises questions about the assumption that noise impact is neutral over time.

These issues point to a more complex interplay between information, behavior, and impact—topics we will explore further in Chapter 9.

2.2.2 The Glosten-Milgrom framework

The Glosten–Milgrom (GM) model reuse this dichotomy between informed and noises traders, to addresses a similar question through a different lens. In this setup, a monopolistic market maker faces uncertainty about traders' information. Let $q_t \in \{+1, -1\}$ denote the type of incoming order. The market maker assumes that informed trades occur with probability π ; otherwise the trade is uninformed.

After observing an order q_t , the market maker updates beliefs about the asset's value v via Bayes' theorem:

$$\mathbb{P}(\text{informed} \mid q_t) = \frac{\pi f(q_t \mid v)}{\pi f(q_t \mid v) + (1 - \pi)f(q_t)}.$$

She sets bid and ask prices such that expected post-trade value matches pre-trade price:

$$a = \mathbb{E}[v \mid q_t = +1], \qquad b = \mathbb{E}[v \mid q_t = -1].$$

This bid—ask spread compensates for adverse selection: the risk that informed traders extract profit.

Thus, the GM model is somewhat more sophisticated than the Kyle model and offers a rationale for how market makers should set the spread to break even. In continuous time, the GM framework leads to a linear permanent price impact when informed trades occur. However, a key limitation of the GM model is its assumption that the fundamental value of the asset is revealed at the end of the trade. For a more flexible framework, one can refer to the *Madhavan–Richardson–Roomans* (MRR) model; see [7], Chapter 16.2.1. Building on the same mechanism, another extended model is proposed in [30], where a Bayesian market maker reacts dynamically to the order flow. In this setting, the impact function becomes concave—more consistent with empirical observations—though the model also presents other limitations, see Chapter 4.

Limitations of Information-Based Models

Models based on informational asymmetries suffer from several key limitations, at least in my view:

- They typically predict a *linear* market impact, which is at odds with empirical findings—see Section 2.4 and Chapter 4.
- They often assume *fully rational behavior* from a specific class of agents, usually market makers. Yet, market impact appears to be a universal phenomenon, observed across various markets and participants, suggesting it should not depend on the strategic behavior of any one group. Furthermore, isn't it optimistic to base a theory on the rational behavior of agents?
- They assume that agents possess some form of information, typically
 understood as a predictive signal about future prices. This is a strong
 assumption, and even if such information exists, there is no reason to
 expect the associated impact to align precisely with the execution of
 their orders.
- They generally suppose that the impact of noise traders cancels out, neglecting the fact that uninformed trading can at least increase volatility and may even generate persistent price movements or trends, see Chapter 9

2.3 Market orders impact and the propagator framework

At the opposite end of information-based models are mechanical models, which tackle a fundamental question: how can a correlated order flow produce diffusive price behavior? Two main models address this issue—one argues that impact is permanent but depends on trading history, while the other proposes that impact is fixed in size but transient.

2.3.1 Permanent impact: The Lillo-Farmer model

Lillo and Farmer (LF) propose that price impact is **permanent**, but its magnitude varies depending on market conditions—especially liquidity. They aim to resolve the efficiency puzzle by showing that when order flow becomes predictable, the market adjusts impact accordingly to keep prices efficient.

Their model modifies the basic price impact equation as follows:

$$r_i = \frac{\epsilon_i f(v_i)}{\lambda_i} + \eta_i,$$

where:

- r_i is the return at trade i,
- $\epsilon_i \in \{-1, 1\}$ is the trade sign (buy/sell),
- $f(v_i)$ is a concave impact function of the trade volume v_i ,
- λ_i is a liquidity parameter
- η_i is a noise term.

Crucially, LF assume that liquidity adjusts dynamically: when buy trades are more likely (predictable), liquidity for buys increases (i.e., λ_i increases), thus reducing their impact. This ensures that the price remains diffusive despite persistent order flow. In other words, the only visibile impact is the impact of surprises, see [7] Chapter 13 for detail discussions.

The main limitation of this model is its reliance on the liquidity parameter λ_i , which must be dynamically - and optimally - adjusted by market participants.

Furthermore, we will present in the remainder of this thesis several arguments suggesting that the impact of orders cannot be permanent and ultimately decays to zero. We will discuss this in Chapter 6.

2.3.2 The propagator model

Another possible resolution to the memory-diffusivity conundrum is provided by the propagator model, first introduced in [31], which posits that each market order has a decaying impact on price. The core idea is that the price can be expressed as a linear superposition of past signed trades, each weighted by a time-decaying kernel:

$$p_t = \sum_{s < t} G(t - s)\varepsilon_s, \tag{2.1}$$

where $G(\cdot)$ is the *propagator function*, encoding how the impact of a single trade evolves over time.

One of the main strengths of the propagator model is its ability to reproduce price diffusivity despite the presence of correlated order flow. Indeed, it is well established that the order sign process ε_t exhibits long memory, with autocorrelation $\mathbb{E}[\varepsilon_t \varepsilon_{t+\tau}] \sim \tau^{-\gamma}$. If we set the propagator to decay as $G(\tau) \sim \tau^{-\beta}$, then the price variance evolves as²:

$$\mathbb{E}[(p_{t+T} - p_t)^2] \sim T^{2-2\beta+\gamma}.$$
 (2.2)

To ensure price diffusivity (i.e., variance growing linearly with time), the exponents must satisfy the condition $\beta = \frac{1-\gamma}{2}$.

This model yields two key insights:

- The impact of individual trades must eventually decay to zero.
- A precise balance—referred to as market efficiency—must exist between the autocorrelation of trade signs (γ) and the decay rate of the propagator (β) to preserve price diffusion.

Note that, although the LF model and the propagator model are based on phenomenologically different mechanisms, it is straightforward to show that they can be mathematically equivalent, see [7, 32].

Several extensions of the propagator framework have been proposed to account for additional empirical features, always with the aim of explaining returns solely through order flow. Without being exhaustive, I briefly mention two key contributions for interested readers:

• A detailed empirical evaluation of refined propagator models—particularly comparing Transient Impact Models (TIM) and History-Dependent Impact Models (HDIM)³—was conducted in [33]. This analysis led to the introduction of the Mixture Transient Impact Model in [34], which was later

$$r_t^{\mathrm{HDIM2}} = \sum_{\pi^{\prime\prime}} \delta_{\pi_t,\pi^{\prime\prime}} \sum_{\pi^\prime} \sum_{j > 0} \kappa^{\pi^\prime \pi^{\prime\prime}}(j) \, \delta_{\pi_{t-j},\pi^\prime} \, \varepsilon_{t-j},$$

where π_t, π_{t-j} indicate whether the mid-price changed following the (t-j)-th transaction.

²See [7], Chapter 13.2.1 for a detailed derivation

³HDIM2 introduces two distinct kernels $\kappa^{\pi'\pi''}(j)$ that depend both on the current event type and on the type of the most recent past event. The return is given by:

revisited and further simplified into the Constant Impact Model⁴ in [35]. These models enabled Patzelt et al. to reproduce mid-price dynamics from historical order flow with remarkable accuracy at high frequency—at least for small-tick stocks.

• An empirical study presented in [36], which aligns closely with the perspective of this thesis. By fitting the propagator kernel to both proprietary and public trade data, the authors demonstrate that short-term impact appears highly universal. This finding supports one of our central claims: trades initiated on private information (presumably unknown to the rest of the market, we hope) affect market prices in much the same way as typical, uninformed trades.

2.3.3 Aggregated impact

Aggregated impact is perhaps one of the first coarse-grained quantities to consider in market microstructure. While the impact of individual trades is highly noisy, the aggregated impact reveals much cleaner and more interpretable patterns, which we now describe.

We begin by defining a time window of length T, typically measured in number of trades, although it can also correspond to fixed real-time intervals. Within each window (setting the beginning of the window at t = 0), we compute the order flow imbalance I and the associated price change Δp :

$$I(T) = \sum_{T} \varepsilon_t, \quad \Delta p(T) = p_T - p_0$$

By averaging over many such intervals and repeating the procedure for various values of T, we find the following empirical relationship, for small imbalances:

$$\frac{\mathbb{E}[\Delta p \mid I, T]}{\sqrt{T}} = \frac{I}{T^{\chi}} \tag{2.3}$$

with $\chi \approx 0.75$, see [37]. For large imbalances, the impact saturates, see Chapter 6, 7. However, in general Eq. (2.3) means that the aggregated impact is linear in the imbalance, up to a rescaling factor that depends on the window size T.

Such a law has been extensively investigated in [7, 37] ⁵, as it appears to be remarkably universal across markets. Interestingly, this behavior is partly reproduced by

$$r_t^{\text{CIM2}} = \Delta_c \, \delta_{\pi_t, c} \, \varepsilon_t,$$

where Δ_c is a constant impact coefficient.

 $^{^4\}mathrm{CIM2}$ is a limiting case of HDIM2 where the impact is both instantaneous and constant. It takes the form:

⁵We will further confirm this behavior in Chapters 6 and 7 and demonstrate both empirically and theoretically that our proposed theory successfully reproduces this non trivial stylized fact.

the propagator framework (see [7], Chapter 13.4.3), although the predicted exponent differs from the empirical one.

Another major puzzle is how to reconcile the linear form of aggregated impact with the time-independent concave impact of metaorders, which we now introduce.

2.4 Metaorder impact —From empiric to models

2.4.1 The Square-Root Law

Understanding how prices react to the execution of metaorders —that is, sequences of orders initiated by the same trader due to liquidity scarcity —is a central challenge in market microstructure. Over the past two decades, a remarkably robust empirical finding has emerged across markets and asset classes 6 : the so-called *square-root law* of market impact. Formally, it states that the average price change I(Q) resulting from a metaorder of volume Q is of the form:

$$I(Q) = Y \sigma_D \left(\frac{Q}{V_D}\right)^{\delta}, \text{ with } \delta \approx \frac{1}{2}, Y \approx 0.5$$
 (2.4)

where σ_D is the daily volatility, and V_D the daily traded volume of the asset.

This result is both puzzling and profound.

From a practitioner's perspective, it offers a simple yet powerful rule for estimating execution costs —a key component in optimal execution, transaction cost analysis, and portfolio optimization. These costs are critical to the profitability of any strategy, as we will detail in Section 2.5.

For academics, it poses two fundamental question:

- How can such a *universal* and *simple* law emerge from the highly complex, noisy, and heterogeneous structure of financial markets?
- Why do trades affect asset prices, given that under the Efficient Market Hypothesis (EMH), price fluctuations should only reflect changes in the fundamental value of the asset?

By let's consider again Eq (2.4). The square-root law is surprising – and maybe counterintuitive – for several reasons.

• First, it contradicts naive expectations that impact should grow linearly with traded volume.

⁶The SQL has been widely validated on stocks [2, 38, 39] but also for futures [40] and even in OTC market, see [41].

Chapter 2.

- Second, it is largely independent of the execution schedule ⁷or the precise order-splitting strategy, suggesting a form of universality. Indeed, impact seems to depend only on one variable: Q.
- Finally, and most remarkably, it is concave: meaning that the market may overreact for small volumes, and under react for larges ones

Note that the prefactor, known as the Y-ratio, is also of significant interest, as it directly influences execution costs. While the concavity of the impact function has been extensively studied, the Y-ratio remains more elusive and difficult to analyze. In Chapter 6, we will derive an analytical expression for it, although thorough empirical investigations are still lacking.

As for the concavity exponent δ , its estimated value has varied across studies — typically between 0.3 and 0.7 —due to biases or limited data availability. However, a comprehensive analysis by Sato et al. [42] demonstrates that choosing $\delta = 0.5$ is the optimal and universal choice.

Eq (2.4), which describes the so-called peak impact, is accompanied by two other well-documented stylized facts regarding metaorders:

• Concave dynamic impact: During the execution of a metaorder, the impact evolves in a concave fashion. Specifically, for $\phi \in [0, 1]$:

$$I(\phi Q) \sim \sqrt{\phi}\sqrt{Q} \tag{2.5}$$

This suggests that the market tends to overreact at the beginning of the execution and underreact toward the end.⁸

• Post-execution relaxation: After the execution is complete, the impact begins to decay. Studying this relaxation is challenging, as price noise increasingly dominates the signal. A long-standing debate exists as to whether impact is truly permanent or ultimately decays to zero. While early studies, such as [43], reported that impact stabilizes around 2/3 of its peak value, more recent work suggests that this apparent plateau may result from restricting the measurement to intraday horizons. Given that relaxation unfolds over timescales much longer than the execution itself, [44] showed that, when measured over multiple days, impact indeed continues to decay and

⁷Of course, the SQL also has its limits. If a metaorder is executed too quickly or too aggressively, the SQL may no longer hold. However, since traders generally execute orders correctly, most metaorders observed in datasets tend to follow this law.

⁸Some argue that this effect could simply reflect typical intraday trading profiles. However, we have verified that this behavior persists even for metaorders executed entirely in the early morning or late afternoon, ruling out a purely intraday explanation.

eventually becomes very small (≈ 0.1 of the peak impact). We will revisit this question in Chapter 6.

Two technical facts that make the Square-Root Law even more puzzling

- The dynamics of price impact—both during and after execution—are particularly intriguing given that metaorders are indistinguishable from the background order flow. The market has no information about when a metaorder begins or ends.
- Most traders execute cautiously, consuming less than the available liquidity at the best quote (over 90% of child orders follow this rule). As a result, individual child orders typically have zero immediate impact.

2.4.2 Data, the missing piece

The SQL is both central to market microstructure and notoriously difficult to study for a given researcher. Observing it empirically requires access to executed metaorders—a type of proprietary and highly confidential data. Moreover, since such datasets usually originate from a single institution (as no one wants to share their trading portfolio), they often suffer from biases and limited scope.

Fortunately, in Chapter 4, we present an empirical analysis based on an unbiased dataset obtained directly from the exchange. Then, in Chapter 5, we introduce an algorithm that generates realistic metaorders from public trade data, making it possible for anyone to explore this phenomenon.

In fact, for those who can access real public trade data, we go one step further in Chapter 7 where we propose a method to simulate realistic market environments. Applying the metaorder proxy to these simulations yields results that are encouraging—though not yet perfect, and still the subject of ongoing research.

2.4.3 The puzzle of the Square-Root Law: A brief review of existing theories

The SQL has inspired a wide range of theoretical efforts aiming to explain its origin. However, to date, no model has provided a phenomenological explanation that fully aligns with empirical stylized facts. To illustrate the difficulty of this challenge, we first briefly present few theories that are known to be incompatible with SQL, before discussing the most prominent one—which we will unfortunately also refute in Chapter 4.

⁹Note however that such a law was deducted simply through a dimensional analysis in [45]

Chapter 2.

A striking example is the propagator model, which accurately describes price dynamics at the market order level but fails at the metaorder scale. In this framework, the predicted impact is given by 10 :

$$I^{\text{prop}}(Q) \sim T^{-\beta}Q$$
 (2.6)

This expression is incorrect in two ways: the peak impact is linear in Q, and it depends on the execution time T. In its standard form, the propagator model is therefore incompatible with the SQL. In Chapter 6, we propose a generalization of this model to address these shortcomings. Still, it's important to note that propagator-based approaches do not offer a *phenomenological* explanation of market impact.

A second theory, proposed by [43] and inspired by the LMF hypothesis, was appealing but ultimately refuted. The model predicted a relationship between the concavity exponent δ of the impact function and the exponent μ governing the distribution of metaorder sizes. However, using the extended TSE dataset, Sato et al. [42] showed that these two exponents are, in fact, empirically uncorrelated.

Another early line of reasoning relates the SQL to inventory management, as proposed in some of the first attempts to explain it [46, 47]. The idea is that a metaorder of size Q, fully absorbed by market-makers, is gradually unwound over a time horizon $T_{\rm off}$. The associated price risk scales as:

Risk
$$\sim \sigma \sqrt{T_{\text{off}}}$$
. (2.7)

Assuming $T_{\rm off} \propto Q$ and inversely proportional to the market trading rate V_T/T , we get:

$$\frac{T_{\text{off}}}{T} \sim \frac{Q}{V_T}.\tag{2.8}$$

Matching risk with compensation yields the impact:

$$I^{\rm inv}(Q) \sim Y \sigma \sqrt{\frac{Q}{V_T}},$$
 (2.9)

consistent with the empirical square-root law.

However, this argument neglects competition: inventory risk is diversifiable, and charging each metaorder for it would lead to excess profits—quickly arbitraged away by other liquidity providers. This would imply $Y \ll 1$, contrary to empirical findings where $Y = \mathcal{O}(1)$.

 $^{^{10}}$ see [7] Chapter 13.4.4 for exact derivation.

The Locally Linear Limit Order Book theory (LLOB) Let us pause briefly to introduce one of the most coherent models for the Square-Root Law - to my knowledge - which we will explore further using the Tokyo Stock Exchange (TSE) dataset. In [15], Donier et al. develop a minimal yet powerful theoretical framework to describe the impact of metaorders, based on a reaction—diffusion model of liquidity. The cornerstone of their approach is the concept of a *Locally Linear Order Book* (LLOB), which assumes that *latent* liquidity is linear around the mid price.

Let $\varphi(x,t) = \rho_B(x,t) - \rho_A(x,t)$ denote the difference between the latent bid and ask densities at price level x and time t. The dynamics of φ are governed by the following reaction–diffusion equation:

$$\frac{\partial \varphi}{\partial t} = D \frac{\partial^2 \varphi}{\partial x^2} - \nu \varphi + \lambda \operatorname{sign}(p_t - x), \qquad (2.10)$$

where:

- D is the diffusion coefficient, modeling the random reassessment of prices by agents,
- ν is the cancellation rate of latent orders,
- λ is the rate of new order arrival,
- p_t is the mid-price at time t, defined implicitly by $\varphi(p_t, t) = 0$.

In the limit of slow cancellations and slow order arrival (i.e., $\nu \to 0, \lambda \to 0$), the stationary solution to Eq. (2.10) becomes locally linear around the mid-price:

$$\varphi_{\rm st}(x) \approx -L(x-p_t),$$
(2.11)

where $L = \lambda/\sqrt{D\nu}$ is an effective liquidity parameter. To incorporate the effect of a metaorder, the model adds an external source term corresponding to the execution of a metaorder of size Q over a duration T at a continuous execution rate m = Q/T. The resulting price trajectory $y_t = p_t - p_0$ satisfies an integral equation of the form:

$$y_t = \frac{m}{L} \int_0^t \frac{ds}{\sqrt{4\pi D(t-s)}} \exp\left(-\frac{(y_t - y_s)^2}{4D(t-s)}\right).$$
 (2.12)

The behavior of such an integral depends crucially on the comparison between m and the market's intrinsic execution rate J = DL. Two limiting regimes emerge for the peak impact, ie I(Q) = y(T):

• Weak trading intensity $m \ll J$:

$$I(Q) \sim \sqrt{\frac{m}{J\pi}} \sqrt{\frac{Q}{L}}$$
 (2.13)

• Strong trading intensity $m \gg J$:

$$I(Q) \sim \sqrt{2\frac{Q}{L}} \tag{2.14}$$

Thus, the LLOB recover the SQL only in the regime of strong trading intensity. In plain word, to follow the SQL, the metaorder has to be executed fast enough to the the linear profile. However, in real market, m is typically much smaller than J, as market participant try to be not detected by the rest of the market.

A way to mend this issue was proposed in [48] with the introduction of slow (low frequency) and fast (high frequency) traders, and the decomposition of $J = J_{\text{fast}} + J_{\text{slow}}$.

Then by assuming that $J_{\text{fast}} \gg J_{\text{slow}}$ - thus $J \approx J_{\text{fast}}$ - one can obtain the SQL as the metaorder is executed only against slow traders after a quick transition phase. Thus, we don't need $m \gg J$ anymore, but just $m \gg J_{\text{slow}}$ which is much more realistic.

While this framework is certainly appealing—and, in my view, mathematically beautiful—the extended model faces two main issues. First, as we will show in Chapter 4, we find that $J \neq J_{\text{fast}}$. Second, high-frequency traders (HFTs) are always present during the execution of a metaorder. However, we will also argue that, in a certain sense, the actions of HFTs may cancel each other out, and that the resistance encountered during metaorder execution—leading to the SQL—is ultimately driven by slow traders. We will return to these points in Chapter 4.

Finally, note that the LLOB framework also makes non-trivial predictions about the post-execution decay: it predicts an infinite negative slope immediately after the end of the execution, followed by a decay in $t^{-1/2}$. This behavior is very close from what is observed empirically, see [44] for an empirical investigation of the decay.

2.5 The art of systematic investing: Alpha versus Impact

Although it is not the main focus of this thesis, let us briefly present the ecosystem of an investment strategy to shed light on why the SQL is so important. For interested readers, an extensive and insightful reference is K. Webster's *Handbook of Price Impact* [49] - I highly recommend.

Different teams for different skills: A hedge fund typically consists of multiple specialized teams. In this section, we focus on two key groups: the *alpha team*, which is responsible for designing predictive trading signals (or "alphas"), and the *execution team*, whose role is to carry out orders in a way that minimizes transaction costs and price impact. Of course, there are many types of alpha—ranging from high-frequency to low-frequency signals—derived through statistical arbitrage, machine learning, alternative data, and more. And naturally, execution strategies may also incorporate short-term alpha to optimize order placement and timing.

The execution paradox: One might think that the true art of investing lies entirely in alpha generation—that is, finding some way to predict future prices—while execution is merely a mechanical step to enter a position. However, in reality—at least to my modest knowledge—execution plays also a critical role.

Consider a portfolio manager receiving a prediction from the alpha team: a fancy machine learning model forecasts that a given stock will rise by 10% over the next month. Naturally, the manager wants to take a position and asks the execution team to buy Q shares (we will discuss how to determine Q in the next section). But in doing so, the very act of buying the stock creates market impact, pushing the price upward. If one believes that impact is inherently tied to information, this price move may appear as a realization of the alpha prediction—it looks like the market is validating the signal.

However, as we will argue throughout this thesis—and, I hope, convince the reader—most of the observed impact is mechanical in nature, not informational. Therefore, one should subtract this mechanical component from observed price changes. When doing so empirically, as we will comment in Section 2.6, impact often turns out to be much larger than the alpha prediction itself.

This leads to a fundamental tension: depending on whom you ask, there is either no such thing as impact—only skilled traders acting on alpha—or there is no meaningful alpha—just the mechanical footprint left by large trades.

What is the real alpha within the mechanical impact framework: A striking implication of the mechanical view of market impact is a shift in the way we think about alpha. Traditionally, alpha is associated with predicting economic fundamentals—such as anticipating macro announcements or firm-specific news before they are priced in. But in a market where trades – and even orders, see Chapter 8 – themselves move prices, another perspective emerges.

In such an environment, alpha increasingly reflects the ability to anticipate what

Chapter 2.

other market participants will do—because their actions will have an impact on prices. Rather than only forecasting the content of an economic release, it becomes more relevant to predict *how others will interpret and react to it*, and how this collective response will move the market. In this sense, alpha is not necessarily about knowing more, but about being ahead in understanding the flow of future orders.

We have a saying, which is that "The impact of others is our alpha, and vice versa." — Jean-Philippe Bouchaud

This helps explain certain seemingly paradoxical situations: it is not uncommon to see breaking news—such as the outbreak of a war or a major political shift—yet observe that stock prices remain flat, even for companies that are clearly economically exposed. While some strong believers in the EMH might argue that such news was already anticipated and priced in, a more plausible interpretation—especially when the event was unpredictable—is that prices didn't move simply because market participants did not react. In that sense, price reflects more the collective beliefs and reactions of traders than any objective economic reality.

Alpha, then, becomes less about discovering fundamental value and more about anticipating how others will trade. Understanding the timing, direction, and structure of future order flow can be more profitable than trying to determine what the "correct" price should be.

We will explore the relationship between impact and alpha in greater detail at the end of Chapter 6, in light of the new theories introduced in this study.

Optimal execution: How to size your trade The delicate balance between execution and alpha is at the heart of what is known as optimal execution—a problem that has sparked an extensive body of research, given the substantial financial implications involved (see [50] for reference). While a full description lies beyond the scope of this thesis, we briefly present a useful rule of thumb to estimate impact.

Suppose you have an alpha signal α , and you wish to trade a quantity Q in order to maximize expected profit. The key question becomes: knowing that your trade will move the market, how much of your alpha will be lost to impact? Or equivalently, how much should you trade?

Your expected PnL can be written as:

$$PnL = \alpha Q - Q \cdot I(Q), \qquad (2.15)$$

where I(Q) is the market impact function. Assuming permanent square-root impact, i.e. $I(Q) \sim \sqrt{Q}$, the optimal trade size is such that impact costs eat up to 2/3 of the alpha, representing a substantial cost. Note that this represents only the impact cost, under the permanent impact hypothesis. In reality, one must also account for impact decay, fees and spread costs, which further exacerbate the situation.

2.6 The unknowns of price impact

To conclude this brief introduction to the fascinating world of price impact, let us highlight some of the central open questions that continue to puzzle both academics and practitioners. Shedding light on them is one of the core ambitions of this thesis.

How can the propagator model and the square-root law be reconciled?

A fundamental challenge is to develop a unified framework that bridges microscopic models of individual market order impact—such as the propagator model—with the mesoscopic regularities observed in metaorder execution, ie the three stylized facts of the SQL. Since child orders are indistinguishable from the surrounding order flow, one would expect that a consistent theory could recover SQL-like behavior from first principles. Moreover, empirical studies show a striking asymmetry between the execution and relaxation phases of a metaorder: impact builds up rapidly during execution, while its decay is slower and resembles the long-memory relaxation patterns found in the order flow itself. A satisfactory model should naturally account for this asymmetry.

How to reconcile price diffusivity, volatility, and the square-root law?

Another pressing question concerns the statistical properties of prices resulting from trading activity. Models assuming a square-root impact often produce subdiffusive price dynamics, as shown for instance in the no-arbitrage framework of [51]. This suggests that metaorders alone may not be sufficient to explain the observed diffusive behavior of prices. It raises the possibility that price diffusivity originates from exogenous sources such as news, or from a component of volatility that is independent of trade execution. Establishing—or refuting—a robust link between metaorder execution and long-term volatility is therefore essential to our understanding of market dynamics.

Chapter 2.

• Is metaorder impact permanent or transient?

Perhaps the most controversial question concerns the nature of impact itself. While the square-root law is now widely accepted empirically, its interpretation remains debated. One school of thought argues that impact is (at least partly) permanent and reflects the informational content—the alpha—of the metaorder. Another perspective posits that impact must decay to zero, with any permanent component being negligible. So, is there truly a permanent component of impact? If so, how does it relate to alpha? And if impact fully reverts, then what exactly is alpha measuring? For interested readers, a lively ping-pong debate takes place between Gabaix et al. [52] and Bouchaud [53] on the *Inelastic Market Hypothesis*, related to those topics.

These fundamental questions are not merely theoretical curiosities: they directly impact how we understand, predict, and manage trading costs, and more broadly, how we conceive price formation itself. And remember, price formation is, after all, the mechanism through which value (or at least price) is assigned to goods, commodities, stocks etc... in modern economies. After a detailed empirical study in Chapter 4, leveraging a unique dataset, we will return to these core issues in Chapters 5 and 6, where we aim to provide new insights and possible answers through a novel and unified approach of market microstructure.

Chapter 3

Contents

Theoretical foundations for Market stability

Markets can remain irrational longer than you can remain solvent

John Maynard Keynes

3.1 Un	stable markets?	40
3.1.1	Flash crashes: a symptom of criticality ?	40
3.1.2	Price jumps in the Limit Order Book	41
3.1.3	The excess volatility puzzle	41
3.2 Ph	ysicists' tools to model the Limit Order Book	42
3.2.1	Hawkes processes	42
3.2.2	Agent-based models: The Santa Fe approach	43
3.2.3	Phase transition theory	44
3.2.4	Power-laws everywhere?	45

Market stability is a topic where physicists—particularly those studying complex systems—may feel very much at home. In nature, it is common for large, unexpected macroscopic events to emerge from the collective behavior of many interacting components, especially near critical points or during phase transitions.

3.1 Unstable markets?

3.1.1 Flash crashes: a symptom of criticality?

On May 6, 2010, U.S. equity markets experienced a dramatic collapse and recovery within minutes. The Dow Jones Industrial Average – and many others American indices, see Fig. 3.1 – dropped nearly 9% before rapidly rebounding without apparent economic reasons. And this kind of crisis is not a simple outlier, indeed those large financial breakdown have occurred many times in the last century (period for which there is easy access to reliable financial data.

Such "flash crashes" cannot be explained by standard models of market equilibrium, ie the EMH hypothesis or by the classic Brownian price theory we presented earlier one.

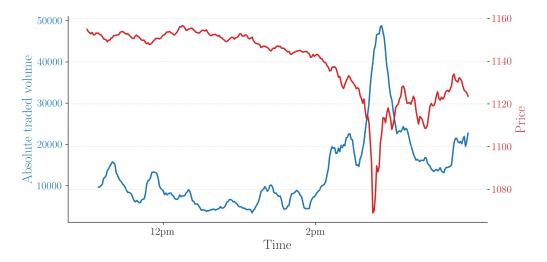


Figure 3.1: Intraday evolution of the S&P Mini on May 6, 2010. The red curve shows the mid-price, which presents a drop of nearly 10% without any identifiable economic trigger. The blue curve represents a moving average of the absolute traded volume.

Regarding the 2010 flash crash, extensive investigations have been conducted (see [54]). The event was triggered by a single, exceptionally large limit order executed by an algorithm developed by a day trader in the UK. That such a single trader could - temporarily - destabilize the entire U.S. market underscores how seemingly small fluctuations can cascade into large-scale market disruptions—much like avalanches in self-organized critical systems

3.1.2 Price jumps in the Limit Order Book

To not speak only of major liquidity crises, price jumps are a frequent and significant phenomenon observed in the limit order book. In particular, [55] analyzed 300 highly liquid NYSE stocks and found that, on average, there is roughly one anomalous price jump per day. Here, an anomalous price jump is defined as a one-minute binned return exceeding four standard deviations (4σ) , after proper volatility renormalization. If prices genuinely followed a Gaussian distribution—as posited by the EMH and assumed in the Black-Scholes framework—such extreme fluctuations would be extremely rare, with a daily probability of around 10^{-5} .

More interestingly, the authors differentiate between exogenous and endogenous jumps by cross-referencing jump times with news releases from Bloomberg¹¹. Two key findings emerge:

- The *temporal profile* of jumps differ substantially depending on their classification as endogenous or exogenous.
- Exogenous jumps, triggered by news, account for only about 1% of all detected jumps.

A notable extension of this work was provided by [56], who employed wavelet techniques to more accurately identify and characterize jumps in price series.

The frequent occurrence of these jumps and liquidity shocks highlights an underlying instability in financial markets. Crucially, most of these extreme events—extreme in size, though not necessarily in frequency as they happen once a day—appear to originate endogenously from the market's internal dynamics, underscoring again the complex, self-organized nature of market instability, that is yet to be proven.

3.1.3 The excess volatility puzzle

One of the most enduring critiques of traditional financial theory is the so-called excess volatility puzzle, first formally articulated by Robert Shiller in the early 1980s [57]. According to the Efficient Market Hypothesis (EMH), prices should reflect all available information, adjusting only in response to new, exogenous signals. Shiller's translates that by computing the volatility of a stock relative to the volatility of discounted future dividends. And by doing so, here revealed a stark discrepancy: the actual volatility of financial markets is an order of magnitude larger than what can be justified by fundamental news alone. This

¹¹Each jump was classified based on whether it coincided with a relevant news event for the corresponding stock near the jump time.

Chapter 3.

mismatch implies either that agents are systematically mispricing assets, or that prices are being driven by *endogenous factors* unrelated to fundamentals.

Numerous studies have sought to explain this longstanding enigma. More recently, Kurth et al. [17] revisited the issue through the lens of the Chiarella model, demonstrating that fundamental volatility could be up to four times smaller than observed market volatility—thereby reexamining the excess volatility puzzle. In particular, Guyon et al. [58] demonstrated that volatility is remarkably well captured by path-dependent models, giving yet another proof of the predominance of its endogenous component. Lastly, a very interesting line of research focuses on rough volatility, as discussed in [26, 27].

Intermediate Conclusion

The frequent occurrence of *liquidity crises* triggered by endogenous dynamics, the fact that most *price jumps* also originate endogenously, and the evidence that *volatility itself is largely endogenous*—all these observations point toward markets operating near a critical point. This stands - againin stark contrast with the EMH. We will dive deeper into these phenomena in Part III. But first, let us briefly outline the tools—and the hedge?—that physicists have when study this aspect of financial markets.

3.2 Physicists' tools to model the Limit Order Book

3.2.1 Hawkes processes

The concept of Hawkes processes was originally developed in the context of seismology, where it was observed that earthquakes tend to trigger subsequent events—aftershocks—resulting in a self-exciting temporal structure. In such a framework, the probability of an earthquake occurring increases if another one has recently taken place. This self-exciting behavior is also observed in financial markets: trading activity tends to cluster, with trades being more likely to occur shortly after other trades.

Hawkes processes constitute a particular class of inhomogeneous Poisson point processes in which the intensity function $\phi(t)$ encapsulates both an exogenous component and an endogenous component that depends on the past history of the process itself. Formally, a Hawkes process is defined as a point process whose intensity evolves according to:

$$\phi(t) = \phi_0(t) + \int_{-\infty}^t dN(u) \,\Phi(t - u), \tag{3.1}$$

where $\phi_0(t)$ is a deterministic baseline intensity, and $\Phi(t) \geq 0$ is a non-negative kernel function that quantifies the influence of past events on the current intensity. In most applications, the kernel $\Phi(t)$ is taken to be strictly decreasing to reflect the decaying influence of older events. Interested readers are referred again to [7], Chapter 9, for a more detailed presentation of Hawkes processes.

A key quantity of interest is the endogeneity ratio η , also known as the branching ratio, defined as:

$$\eta = \int_0^\infty \Phi(t) \, dt. \tag{3.2}$$

Empirical studies—see [59, 60]—suggest that in equity markets, η is often greater than 0.7 - when fitted on price changes. This indicates that financial markets operate close to a critical point, where small perturbations can propagate and amplify significantly ¹².

However, it is important to highlight that the application of Hawkes processes to financial data has received criticism. In particular, they may not be ideally suited for modeling order flow, which exhibits long memory. This long-range dependence can naturally drive fitted Hawkes models toward criticality, potentially obscuring the underlying dynamics. We will revisit this issue in Chapter 8.

That said, we will nonetheless make use of Hawkes processes in Chapter 9 to model feedback mechanisms and endogenous loops within market dynamics.

3.2.2 Agent-based models: The Santa Fe approach

Another powerful tool that has developed rapidly in recent years—particularly in economics—is the so-called agent-based model (ABM). In statistical physics, where systems are typically composed of a large number of identical elements (such as spins or atoms), it is common practice to simulate these systems to validate theoretical predictions. Over the past decade, this simulation-based approach has been increasingly adapted to economic problems. Instead of modeling atoms, one simulates a large number of economic agents, defines rules for their interactions, and observes the resulting macroscopic behavior. A compelling application of this methodology can be found in [61], where agent-based modeling is used to study the post-COVID economic recovery.

In the context of market microstructure, one of the first—if not the first—agentbased models was introduced during the Santa Fe conference, giving rise to what is now known as the *Santa Fe model* ¹³. This pioneering model explored the

 $^{^{12}\}eta = 1$ being the critical point...

¹³To be more precise, there are two distinct models referred to as the Santa Fe model: the first is a conventional agent-based model that represents various agent types, including informed

Chapter 3.

dynamics generated by a large number of traders placing orders randomly in a limit order book. Surprisingly, despite its minimal assumptions, the resulting price process more or less resembled real market prices, although it tended to be – strongly – mean-reverting. For a comprehensive overview of the stylized facts captured by this model, we refer the reader to Chapter 8 of [7].

Several extensions of the Santa Fe model have been proposed. For instance, in [62], the authors introduced a simple tuning of the aggressiveness of market orders—specifically, the proportion of the best quote consumed by a trade—and demonstrated that this modification alone suffices to recover a diffusive price behavior. Another significant reference is the work by Ravagnani et al. [63], in which the authors introduced an extended version of the Santa Fe model that incorporates the SQL. However, reconciling this model with price diffusion remains challenging.

Here, we are more interested in the development proposed in [64], where the authors added a single rule to the behavior of otherwise random traders: a feedback mechanism that makes agents responsive to past price trends. This modification led to the emergence of liquidity crises, and more precisely, revealed that the limit order book undergoes a second-order phase transition, that we will define in the next Section.

This result is especially intriguing, as liquidity crises—and extreme price jumps—are both frequent in financial markets and notoriously difficult to understand analytically. Agent-based models thus offer a promising and efficient framework to study such phenomena. In Chapter 9, we will develop a model in this spirit. But before doing so, we first introduce briefly the concept of phase transition.

3.2.3 Phase transition theory

Being exhaustive about phase transitions is clearly beyond the scope of this brief introduction ¹⁴. In a nutshell, it is well known in physics that a system—such as one liter of water—can exist in different phases (solid, liquid, or vapor). Transitions between these phases can be triggered by tuning just a few parameters, typically pressure and temperature. What is striking is that although the underlying constituents remain identical (water molecules, in this case), their organization leads to vastly different macroscopic properties: ice and vapor behave very differently, despite being composed of the same molecules.

traders and noise traders. The second model, which is used in this thesis, also acknowledges the presence of different agents but directly focus on aggregated observables, ie the different order flows.

¹⁴For interested readers, a great reference in my view for phase transition theory is: https://www.lpthe.jussieu.fr/~leticia/TEACHING/ICFP2021/PhaseTransitionsICFP-Chapter.pdf

A similar analogy can be drawn for limit order books. With just a few changes in the system—such as memory of market participant for example—the market can shift from a stable regime, where prices are well-defined and liquid, to a regime resembling a liquidity crisis. In this context, the transition is dynamic, since the number orders present in the orderbook may evolve over time.

What is particularly interesting is the *signature* of a phase transition—the point at which the system crosses from one phase to another. At this critical point, the system becomes scale-invariant or "fractal", and several observable quantities begin to follow power-law scaling. These include, depending on the system, the susceptibility, correlation length, specific heat and so on.

The exponents associated with these power laws serve to characterize the nature of the phase transition. Remarkably, most transitions fall into a limited number of *universality classes*, each defined by a specific set of critical exponents. This elegant framework from statistical physics offers deep insight into seemingly diverse systems undergoing transitions.

Returning to financial markets, the analogy with physical phase transitions proves surprisingly relevant. As we will demonstrate in Chapter 9, in a modified version of the Santa Fe model, several key quantities—such as the spread, susceptibility, and others—exhibit power-law scaling close to criticality. We will explore how to derive the corresponding critical exponents, thereby aiming to characterize the nature of these transitions.

3.2.4 Power-laws everywhere?

Let us use this brief discussion of phase transitions to address a common criticism of Econophysics—one that I have frequently encountered over the past three years—namely that "physicists see power laws everywhere." While this remark is not entirely unfounded, I would argue that physicists are, in fact, justified in doing so, for at least three key reasons:

• First, power laws indeed arise naturally in many physical systems and provide remarkably good fits to a wide range of empirical distributions observed in nature—such as the energy of earthquakes, solar flares, or avalanche sizes in sandpile models. As we discussed earlier, they often emerge from systems undergoing phase transitions. Power laws are also prevalent in economics, a classic example being the Pareto distribution of wealth: roughly 20% of individuals hold 80% of the wealth, and this pattern recursively holds within the wealthiest 20% - ie being self invariant. In short, once elements within a system begin interacting, the assumptions underpinning the Central Limit Theorem break down, and power-law signature can emerge, with extreme

Chapter 3.

events becoming more likely due to cascading effects and collective dynamics. For a more in-depth discussion, see the seminal work of Nassim Nicholas Taleb [65] and his *Incerto* series (*The Black Swan, Skin in the Game*, etc.).

- Second, it is well known that power laws can be efficiently approximated by a sum of exponentials, providing a practical mathematical tool to bridge exponential and power-law behaviors [66].
- Finally, in statistics, there is a well-known adage that every fancy model ultimately reduces to a form of linear regression. And indeed, it remains a powerful tool for data analysis. However, when describing phenomena that unfold over multiple timescales, it is often more appropriate to work with the logarithm of variables rather than the variables themselves. But then, performing a regression in log-log space is essentially equivalent to fitting a power law!

To conclude, in Econophysics—as I see it—researchers do not necessarily claim that underlying processes follow an exact power law. Rather, power laws often provide *robust and useful approximations* over a broad range of observed data, making them a valuable modeling tool. Crucially, many phenomena deviate significantly from Gaussian assumptions—widely held in economics or quantitative finance—and thus require fat-tailed distributions to be described accurately, a fact often overlooked by conventional wisdom!

3.3 The Unknowns of market stability

To summarize this introduction, financial markets exhibit several symptoms suggesting they operate near a critical point. At low frequencies, phenomena such as flash crashes and liquidity crises occur far more frequently than standard models would predict, often with terrible consequences for the global economy. At high frequencies, markets display an anomalously high rate of extreme price movements—for instance, price jumps exceeding 4σ occur on average once per day even among the most liquid stocks in the world. Notably, approximately 99% of these jumps are not associated with identifiable news events, indicating endogenous origins.

To investigate this apparent criticality, we will first introduce in Chapter 8 a Vector Autoregression (VAR) framework designed to capture the most probable direction of market evolution. By performing a principal component analysis (PCA) on the system's dynamics, we identify a dominant eigenvector whose associated eigenvalue approaches one. This suggests that the system enters a regime where it is almost "certain" to evolve in a specific direction, highlighting strong internal coordination—and this direction is exactly the one of a liquidity crisis.

Finally, in Chapter 9, we will propose an extended version of the Santa Fe model incorporating realistic feedback mechanisms. This will allow us to test whether, under plausible behavioral rules, the system still exhibits a second-order phase transition, thereby reinforcing the analogy between market dynamics and critical phenomena in physics.

Chapter 3.

Part II Price Impact

Chapter 4

Empirical Analysis of the Microscopic Foundations of the Square-Root Law

The scientist does not study nature because it is useful, he studies it because he takes pleasure in it, and he takes pleasure in it because it is beautiful.

Henri Poincaré

To better understand the Square Root Law, this thesis begins by analyzing price impact using a detailed dataset from the Japanese stock exchange, which contains trader IDs for all orders submitted between 2012 and 2018. Our analysis reveals that the square-root price impact law has microscopic roots, evident even at the level of individual child orders, provided there is enough time for the market to "digest" them. Additionally, we find that the mesoscopic impact of larger orders, or metaorders, results from a "double" square-root effect: a square-root dependence on individual trade volume coupled with an inverse square-root decay over time.

Chapter 4.

From: The "double" square-root law: Evidence for the mechanical origin of market impact from the Tokyo Stock Exchange

G. Maitrier, G. Loeper, K. Kanazawa, JP. Bouchaud

Contents

4.1	Intr	oduction	52
4.2	Data	a description and preliminary observations	54
	4.2.1	A unique dataset	54
	4.2.2	General stylized facts about metaorders execution $\ \ . \ \ .$	56
	4.2.3	Time scales & Market ecology	58
4.3	Squa	are-root impact: micro-scales & meso-scales	60
	4.3.1	The "double" square-root impact of child orders $\ \ .$	61
	4.3.2	A non-linear propagator model	62
4.4	Fron	n single market orders to synthetic metaorders	64
	4.4.1	The impact of single public market orders $\dots \dots$.	64
	4.4.2	Synthetic metaorders	65
	4.4.3	Discussion	66
4.5	The	other side of market orders: liquidity providers $$.	67
	4.5.1	Refill sequences	67
	4.5.2	Strategic behaviour of liquidity providers	68
4.6	Con	clusion	70

4.1 Introduction

As presented in Chapter 2, the Square-Root Law (SQL) for price impact is arguably one of the most robust empirical regularities discovered in the last 30 years – see [38, 40, 47, 67–69] and [7, 49] for reviews. It states that when executing a buy (resp. sell) meta-order of total size Q, sliced and diced into N child orders of size q = Q/N, the price on average moves up (down) by an amount proportional to \sqrt{Q} . Price impact is, quite remarkably, found to be approximately independent of both N and of the total time T needed to achieve full execution [70]. In other words, provided the participation rate is not too large, average price impact only depends – to a first approximation – on the total volume traded Q, but not on execution schedule [7].

Such a square-root dependence, and its apparent universality across a wide variety of markets [7], is surprising and non-intuitive.

Several theoretical ideas have been put forth in the literature to explain non-linear impact. Some models predict a concave price impact Q^{δ} with $\delta \leq 1$ related to the power-law tail exponent μ of the executed volume [71] or the power-law tail exponent γ of the time autocorrelation of the sign of market orders [18, 43]. However, as recently shown by Sato and Kanazawa [72] using ID-resolved data from the Tokyo Stock Exchange (TSE), the predicted relations between δ and μ or γ are not borne out by the data: whereas α and γ significantly differ between stocks, exponent δ remains stubbornly anchored around $\delta = 1/2$, i.e. the value corresponding to the square-root law.

Numerous other models exist, including those mentioned earlier. For more details, see [30, 68, 73, 74]. Despite its critical importance for both financial microstructure and an asset pricing [52, 53], the very origin of this central phenomenon remains a topic of debate. The universality of the phenomenon suggests a *purely mechanical*, rather than informational, origin, however this point is controversial and at odds with most of the economics literature on the subject, starting with the famous Kyle model [75]. The aim of the present Chapter is to give more credence to the "mechanical hypothesis" of the square-root market impact law using an ID-resolved data set from the TSE. Our main results are the following:

- The square-root impact law of metaorders is already valid for child orders, provided one waits long enough for the market to digest these orders.
- The square-root impact law of metaorders emerges from a "double" square-root behaviour: the dependence of the impact of child orders on their individual volume and the inverse square-root relaxation of this impact.
- There does not seem to be anything special about the impact of the child orders of a given metaorder in fact all market orders appear to impact prices in the same manner on average. Correspondingly, the impact of synthetic metaorders, reconstructed by randomly scrambling the identity of traders, is identical to the impact of real metaorders. This is our major piece of evidence for a mechanical origin of the square-root law.

We also garnered further information shedding light and/or putting constraints on the interpretation of the square-root law:

• "Fast" traders, for which the holding period is less than a day, represent between 50% and 60% of the executed market orders. This includes market makers (HFT) and short term traders. Correspondingly, the fraction of market orders executed against "fast" traders represents nearly half of the exchanged volume, a number far too low to vindicate the standard interpretation of the square-root impact within the LLOB framework [76]. Since this empirical fact is inconsistent with the standard LLOB framework, it

Chapter 4.

motivates us to propose a new interpretation of the LLOB framework, as presented in this Chapter.

- One can define refill sequences for liquidity providers. One finds that the size of those sequences is also power-law distributed, as was found for metaorders in [16], following the suggestion of Lillo, Mike & Farmer [18].
- Once a buy (sell) market order is executed, liquidity providers tend to increase (decrease) their offered price. Such a price degradation however decreases as a power-law of the number of trades already executed, with a prefactor that separates aggressive and wary liquidity providers.

The outline of the Chapter is as follows: Section 4.2 describes the dataset and presents some general facts about execution. Section 4.3 investigates the cumulative price impact of child orders and proposes a non-linear propagator to rationalize the empirical results. In Section 4.4, we extend results from the previous Section to all market orders, and we introduce a method for generating synthetic metaorders that are found to follow exactly the same square root law as real metaorders. In Section 4.5, we scrutinize the opposite side of market orders by analyzing the behavior of liquidity providers, and we present our conclusions in Section 4.6.

4.2 Data description and preliminary observations

4.2.1 A unique dataset

Our study is based on a dataset from the Tokyo Stock Exchange (TSE), provided by the Japan Exchange Group (JPX) for academic purposes only and already used in [16]. The dataset contains all orders sent to the exchange, with a unique order ID, a virtual server ID, the price and type of the order, the volume and price of the best quotes, for all stocks available on the exchange from 2012 to 2018. Here the virtual server ID is the unit of trading accounts on the TSE. Technically it is not a membership ID (i.e., the corporate level ID) because any trader may have several virtual servers to avoid the submission-number limit during a fixed interval. However, one can reconstruct an effective trader ID, called the Trading Desk, by properly aggregating those virtual server IDs (see [16, 77] for the details). In this Chapter, the Trading Desks are referred to as trader IDs.

We focus on the top 100 liquid stocks of the exchange, including 10 ETFs. After anonymizing all assets names (for confidentiality reasons), we only kept orders submitted during continuous double auctions trading sessions: there are two sessions each day in the Japanese market, during 09:00 - 11:30 and 12:30 - 15:00. We discarded orders submitted during the 10 first and last minutes of each sessions,

as they might be affected by special conditions. In the following, we will refer to these two distinct periods as two separate days, with a slight abuse of language.

We define a *metaorder* as a sequence of consecutive market orders ("child orders") of same sign (buy or sell) submitted by the same trader during a given session. The dataset is unique for two reasons: (i) we have access to a colossal unbiased set of metaorders and (ii) we can analyze traders behaviors as we can assign each orders to a given participant.

Item (i) is particularly important regarding the claim that impact is an universal mechanism, independent of the type of traders and the information content of the trades [7, 36]. Indeed, access to metaorders data is rare (the Ancerno dataset being one exception [44, 78, 79]), and most of datasets used in the literature are proprietary and may be plagued by conditioning effects [7, 49]. ¹⁵

Item (ii) represents an interesting opportunity to categorize traders as market makers (MM), high frequency traders (HFT) or low frequency traders, see 4.2.3 and understand better their typical impact on the market. For example, the long term debate about the benefits of HFT for market stability has been dramatically improved with this kind of dataset [77, 81]. In addition, these identifiers enabled us to analyze the behavior of *liquidity providers*, which is only possible if we have access to the ownership of all limit or cancellation orders, see Section 4.5.

¹⁵Note however that CFM data [38] was acquired in such a way to minimize conditioning effects such as those discussed in [7], ch. 12.3, see also [80].

4.2.2 General stylized facts about metaorders execution

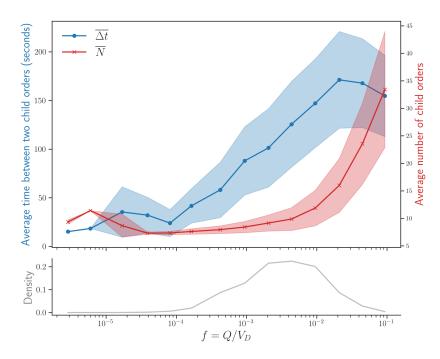


Figure 4.1: Top graph: The blue points represent the average time between two child orders of the same metaorder (in seconds), as a function of $f = Q/V_D$. The red points are the average number of child orders as a function of $f = Q/V_D$. The average red curve can be approximately fitted by a power law $f^{0.3}$, for $f \ge 10^{-3}$. The shaded areas represent the corresponding standard deviations. Bottom graph: Distribution of the size of the metaorders, showing a maximum between 0.1% and 1% of the daily volume. Results are averaged over the top 24 most liquid TSE stocks.

A metaorder ω is usually characterized by few metrics: $Q(\omega)$ is the total metaorder size in shares and $N(\omega)$ the total number of child orders, i.e. the number of consecutive market orders of the same sign from the same trader. $T(\omega)$ is the total duration of the execution, $q_i(\omega)$ is the size (in shares) of the *i*th child orders, and $p_i(\omega)$ is the log mid-price just before the time of execution $t_i(\omega)$.

Natural questions that arise (among many others) are:

- What is the typical volume Q of metaorders compared to the total daily volume V_D ?
- How does N and T depend on Q on average?
- What is the average execution schedule, i.e. how does the already executed

volume $\sum_{t_i < t_i} q_j$ depend on $t_i - t_0$?

We show in Fig. 4.1 (bottom graph) the distribution of executed volume fraction $f := Q/V_D$, which shows a broad maximum in the region $f \in [0.1\%, 1\%]$. Some metaorders correspond to 10% of the daily volume but they are relatively rare. Similarly, very small metaorders of size < 0.01% have a very small probability. The range $f \in [0.1\%, 10\%]$ for metaorders is typical of firms like, e.g., AQR or CFM [14, 80].

The average time between child orders Δt and the average number of child orders are plotted in Fig. 4.1 (top graph) as a function of $f=Q/V_D$, where V_D is the executed volume during the day in shares. We also show as a shaded region the standard deviation of these quantities. One sees that the average time between child orders mildly increases as a function of f, ranging from 25 secs. for f=0.01% to 150 secs. for f=1%, before saturating or even slightly decreasing for larger values of f. Hence the average execution time T increases slightly faster than Q itself, except perhaps for the largest metaorders.

The increase of Δt when Q increases is related to the fact that child orders are more and more aggressive in order to complete execution, so traders wisely wait longer before sending the next one, lest they are detected by market makers. This however becomes difficult for large fs, because traders are also attempting to execute their metaorders as quickly as possible. Note that translated into total execution time T, these results show that for f=0.01%, the typical value of T is ≈ 200 seconds, whereas for large metaorders with f=10%, $T\approx 90$ minutes. These numbers are however only indicative and the total duration of execution can rise to a full day.

The typical number N of child orders is around 10 for $f \lesssim 1\%$ before increasing steeply for larger f. This reflects the fact that the available volume at the best quotes is relatively small and if traders want to avoid "eating into the book", then the size q of child orders is also limited, which mechanically pushes the number of child orders up when Q increases. Note that typically the volume available at the best quote is around $10^{-4}V_D$ for the most liquid stocks.

Chapter 4.

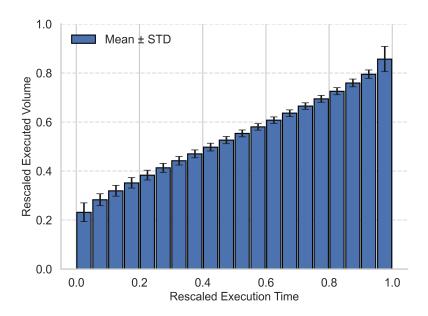


Figure 4.2: Mean fraction of executed volume of metaorders as a function of the rescaled execution time $(t_i - t_1)/T$. Slight shifts on the x-axis come from binning effects. Results from all metaorders in our dataset.

Let us finally turn to the average execution schedule. We plot in Fig. A.2 the average executed fraction $\sum_{t_j \leq t_i} q_j/Q$ as a function of the rescaled time since the start of execution $(t-t_1)/T$. One sees a nicely linear average execution profile, suggesting that metaorders are typically executed using a constant trading rate, except possibly at the beginning and at the end of the execution where trades are more aggressive – although this effect might be dominated by metaorders with a small number of child orders.

4.2.3 Time scales & Market ecology

One of the aspects that makes financial markets particularly complex is the wide range of time scales over which traders operate. Indeed, time horizons range from years or decades for institutional investors (pension funds, mutual funds etc) to sub-seconds for market makers.

These time scales are relevant to understand order flow and liquidity dynamics, and therefore price impact. Indeed, whereas market makers allow orderly trading by acting as intermediaries between final buyers and final sellers, their inventory constraint prevent them from offering "resistance" to large buy or sell metaorders.

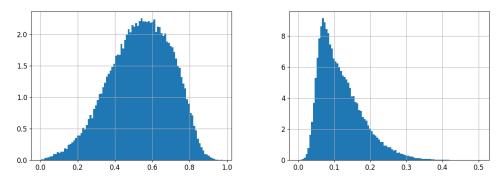
Only slow liquidity can counter-act persistent order imbalance [76].

In this Section, we leverage the identification of traders to classify them into four different classes. We denote by I_t the inventory of a given trader at time t, and by $\epsilon_t = \pm 1$ the sign (buy/sell) of the orders they submit at time t. We define the reversal time τ as the average time between two consecutive orders with different signs from the same trader:

- Long Term traders: These traders have long term trading horizon, and typically execute large metaorders. Orders submitted generally have the same sign as their inventories: $I_t \cdot \varepsilon_t > 0$ and the reversal time τ is typically longer than a trading session.
- Short Term traders: They can trade high frequency signals that flip sign during a trading session. To exploit their signals, they must however build significant positions I_t . Hence we expect them to trade in the same direction for a while, but with a reversal time shorter than a trading session, typically around 30 minutes. In this Chapter, a "fast" trader is broadly defined such that their τ is smaller than the session time as the broadest definition. Otherwise the trader is regarded as a "slow" trader.
- Market Markers: They provide liquidity to the order book, earning the spread but possibly suffering from price impact. In a way, these participants are the easiest to identify: they are responsible for a large part of trading activity while keeping their inventories I_t close to zero. Their reversal time is typically under the minute.
- Brokers: They are executing orders on behalf of their clients, so they are trading at high frequency and their inventories can vary greatly across sessions or stocks. It is difficult to assert with certainty that a trader is a broker, as they may resemble the other three categories in a given session.

Using the reversal time τ as a criterion to separate "fast" and "slow" traders, we compute the contribution of the two categories to the trading activity using either the fraction of the market order volume executed by fast traders, $V_{\rm fast}/V_D$ or the relative fraction of the number of fast traders, $N_{\rm fast}/N_D$. $N_{\rm fast}$ is the number of "fast" traders in a given session, whereas N_D is the total number of traders having traded at least once during the session. We show in Fig. 4.3 the histogram of these ratios, computed over all sessions and all stocks of our database. Whereas the most probable value of $N_{\rm fast}/N_D$ is around 8%, the most probable value of the ratio $V_{\rm fast}/V_D$ is between 50% and 60%. Therefore, the contribution of "slow" market orders to total volume is roughly one-half. This finding is in line with [82]. One can also determine the fraction of market orders executed against fast traders, which is found to be in the range 60% to 70%.

Chapter 4.



(a) $V_{\rm fast}/V_D$ per session for 100 stocks over 10 (b) $N_{\rm fast}/N_{\rm D}$ per session for 100 stocks over 10 years

Figure 4.3: Histogram of the participation of fast traders to the global trading activity. V_{fast} is the volume executed by fast trades during a session. V_D is the volume of market orders during a session. N_{fast} is the number of fast traders and N_D is the total number of traders participating to a given session. We used each session for each stock in our database (around 40,000 sessions).

The conclusion of this study is that while fast traders are dominant in terms of volume, the contribution of slow volumes to trading activity actually of the same order of magnitude. This finding is important since the standard interpretation of the square-root impact law in terms of latent liquidity [74, 76] assumes that slow volumes are a factor ~ 300 times smaller than V_D [79], which is certainly not the case here. In the next Section, we will revisit the empirical evidence for square-root impact and propose a new interpretation of the LLOB model.

4.3 Square-root impact: micro-scales & meso-scales

As recalled in the introduction, there is overwhelming empirical evidence for the square-root impact law for metaorders, which has again been confirmed in great detail in [72] for the TSE, using the very same dataset as here. More precisely, the square-root impact law states that

$$\mathcal{I}(Q) := \mathbb{E}[\Delta p \cdot \epsilon \mid Q] = Y \sigma_D \sqrt{\frac{Q}{V_D}}$$
(4.1)

where ϵ is the sign of the metaorder of total size Q, Δp is the log mid-price change between just before the first and just after last child order and σ_D and V_D are the contemporaneous daily volatility and exchanged volume. Note again that $\mathcal{I}(Q)$ is independent of the execution time T (see for example [70] and Fig. 4.6 below). In the rest of the Chapter, we will use $(p_{\text{high}} - p_{\text{low}})/p_{\text{open}}$ as a proxy for σ_D .

In this Section we attempt to dissect the square-root law into more microscopic components, which sheds further light into its origin and leads to a new interpretation of the Latent Liquidity model.

4.3.1 The "double" square-root impact of child orders

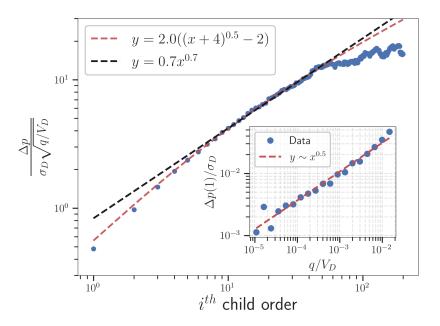


Figure 4.4: Average price profile during the execution of a metaorder. The vertical axis represents the cumulative impact of child orders, rescaled by the daily volatility and the square root of the relative volume of the child order q. These profiles are obtained averaging over the top 8 most liquid stocks of our dataset. We show with red dotted line the function $(\sqrt{i+i_0}-\sqrt{i_0})$ and in black a pure power-law fit $i^{0.7}$. Inset: Average impact of the first child order $\mathbb{E}[\Delta p(1) \cdot \epsilon]$ as a function of its size q rescaled by daily volume V_D , demonstrating the fact that the square-root law is in fact valid at the level of child orders.

Instead of measuring the impact \mathcal{I} of full metaorders, one can measure the partial impact $\mathcal{J}(q,i)$, i.e. the average price difference Δp_i between just before child order i+1 and just before child order 1, conditional on the size q of child orders and the rank i. We find that to a good approximation that impact is proportional to

Chapter 4.

both \sqrt{q} and \sqrt{i} :

$$\left| \mathcal{J}(q,i) := \mathbb{E}[\Delta p_i \cdot \epsilon | q, i] \propto \sigma_D \sqrt{\frac{q}{V_D}} \left(\sqrt{i + i_0} - \sqrt{i_0} \right) \right|$$
(4.2)

with $i_0 \approx 4$, see Fig. 4.4. We show in the inset that the \sqrt{q} dependence on the volume of child orders holds very well for i = 1, but we have checked such a dependence for other values of i as well. This demonstrates that the square-root law already holds at the level of child orders, provided one waits long enough after the execution, see below. Eq. (4.2) is the central result of this Chapter.

Several remarks are in order:

- The fit is very good up to i=50, beyond which another regime appears to set in, where impact saturates. This however only concerns a small fraction of large metaorders, for which conditioning effects should be taken into account (e.g. large prevailing liquidity at the opposite best, see [37]). An alternative interpretation is that such large metaorders are detected by the market, triggering the influx of opposing limit orders. Impact saturation for large Q has been reported elsewhere as well, see e.g. [78, 79].
- When fitting with a more general power-law $(i+i_0)^{1-\beta}-(i_0)^{1-\beta}$, the optimal value of β is found to be 0.48, i.e. very close to a square-root. Alternatively, imposing $i_0 = 0$ yields $1-\beta = 0.7$, but the fit is clearly worse for small times, see Fig. 4.4.
- When $i \gg i_0$, one finds that the temporal profile of the impacted price behaves as \sqrt{i} , a result already reported in [69, 78, 83].
- When i = N and using Q = qN one finds

$$\mathcal{I}(Q) \propto \sigma_D \sqrt{\frac{Q}{V_D}} \left(\sqrt{1 + i_0/N} - \sqrt{i_0/N} \right),$$
 (4.3)

allowing one to recover exactly the square-root impact law, Eq. (4.1), up to a weakly varying N-dependent factor that increases from 0.31 to 1 as N goes from 2 to ∞ (when $i_0 = 4$).

4.3.2 A non-linear propagator model

The above results can be summarized within the framework of a non-linear propagator model [7, 31], where the impact of child order j measured at time $t_i > t_j$

¹⁶The exponent β is defined as the decay exponent of the propagator, as $|t_i - t_j|^{-\beta}$, see [7, 31].

is proportional to $\sqrt{q_j/(t_i-t_j+s_0)}$, as predicted both by the LLOB model [74] and by the Bayesian theory of Ref. [30]. Indeed, from such an expression one gets:

$$\mathcal{J}(q,i) \propto \sqrt{q} \sum_{t_i < t_i} \frac{\sqrt{\Delta t}}{\sqrt{t_i - t_j + s_0}} \approx 2\sqrt{q} \left(\sqrt{i + i_0} - \sqrt{i_0} \right), \qquad s_0 \equiv i_0 \Delta t, \quad (4.4)$$

where we have assumed for simplicity that $q_j = q$, $\forall j$ and $t_j \approx j\Delta t + t_0$, with Δt the time between two consecutive child orders. We have furthermore approximated the discrete sum over j by an integral over t_j .

Such an interpretation however appears to violate the "diffusivity" condition derived in [31, 84], which relates the decay of the propagator with the decay of the autocorrelation of the sign of the trades. Superficially, a propagator decay in $(t_i-t_j)^{-1/2}$ should lead to strongly mean reverting prices, at odds with the diffusive nature of prices. A way out of this conundrum will be presented in a forthcoming paper [85]. Note that the above non-linear propagator model precisely saturates the no-arbitrage bound derived by Gatheral in [86].

The most striking result of the previous Section is that the square-root law appears to hold already at the level of child orders. This is not in line with the standard "mesoscopic" interpretation of the Latent Liquidity model [74], which, as we alluded to before, would require the fraction of slow volume to be much smaller than what we reported in Section 4.2.3.

We are thus led to the conclusion that the latent liquidity idea must in fact operate already at the micro level. Whereas the revealed order book contains primarily limit orders posted by market makers, the final sellers' or buyers' price would be distributed according to a locally linear profile, as predicted by the LLOB theory [74] – which, as a reminder, only relies on minimal assumptions, in particular on the diffusive nature of prices.

So the scenario would be as follows: once the incoming buy (sell) has been executed against a market-maker, a "hot potato" game starts between market-makers until the order is finally digested by a final seller (buyer). In order to find such a final seller (buyer) the price must on average move by an amount δp such that $\Gamma(\delta p)^2/2 = q$, where Γ is the slope of the latent liquidity. Once this is achieved, the price tends to revert back as $1/\sqrt{t-t_i}$, as predicted by the LLOB model, see [7, 74], but in disagreement with the predictions of the propagator model [31]. This discrepancy will be discussed further in [85].

If the above interpretation is correct, however, it should hold for arbitrary market orders. Since all orders are equivalent, they should lead to the same average impact on prices (as was indeed found in [36]). In other words, one should see a \sqrt{q} impact

Chapter 4.

for single market orders provided one waits long enough for the "hot potato" game to be completed. This is what we test in the next Section, which then opens up the question of reconstructing synthetic metaorders from a list of consecutive market orders that do not necessarily belong to the same ID. The anonymity of market orders suggests that, in certain conditions, the very same square-root impact law \mathcal{I} given by Eq. (4.1) should also hold for synthetic metaorders. This indeed turns out to be the case, as we discuss now.

4.4 From single market orders to synthetic metaorders

4.4.1 The impact of single public market orders

We want to understand the behavior of the price after a buy (sell) market order. Clearly, if the volume of the market order is less than the prevailing volume at the opposite best, the immediate impact is zero. However, as time goes by, one very quickly sees that impact grows and becomes approximately given by \sqrt{q} , whether or not immediate impact is zero (the full temporal aspects will be explored in more details in [85]).

More precisely, we show in Fig. 4.5 the impact of a single market order as a function of q for two typical stocks of the TSE, after waiting for a volume time equal to q itself – i.e. after the market has traded the same quantity as the initial market order. We see that independently of whether or not the initial market order has an immediate impact, the overall behaviour is compatible with an impact growing as \sqrt{q} . To remove intraday seasonality effects, we first determine the average intraday profile of volatility σ_b and executed volume V_b , using yearly data, computed on each 15 minutes bins. Then, we rescale impact and volume of market orders by, respectively, σ_b and V_b corresponding to the 15-minute bins of the day to which these orders belong.

For some stocks, a plateau regime however appears for very small q, perhaps related to tick size effects (see Fig. 4.5; right graph).

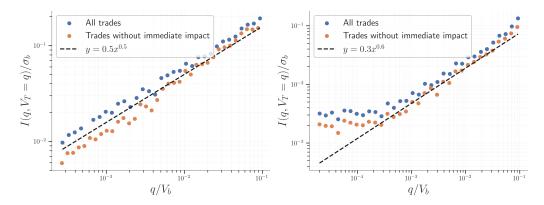


Figure 4.5: Impact of public market orders measured after a volume time equal to the size q of the market order itself. σ_b and V_b are the average volatility of volume of the 15 minute bin to which the order belongs. The blue points represents the average impact of all trades, while the orange points represents trades that have no immediate impact, i.e. such that q is smaller than the prevailing volume at the opposite best. In this way, we can see that the impact indeed builds up over time. Left and right graphs correspond to two typical stocks, one showing a nearly perfect \sqrt{q} behaviour for all q (see dashed black line). Note that the right graph exhibits a plateau for very small q's.

4.4.2 Synthetic metaorders

Since arbitrary market orders seem to all behave similarly, independently of the metaorder they belong to, we made the following numerical experiment that confirms the non-linear propagator interpretation of Eq. (4.2) for metaorders. We construct a new dataset of synthetic metaorders, by randomly shuffling traders ID and distributing them to market orders while we keep the historical market order flow, i.e. by randomly reordering real traders' IDs ¹⁷. This preserves the initial frequency distribution of traders: i.e. some of them appear many times whereas others are trading less frequently.

We then use the same method as in 4.2.1 to define metaorders as a consecutive sequence of trades of the same sign associated to the same new trader ID (that has been reshuffled). We obtain synthetic metaorders, that start and end at different times as the original (true) metaorders. Hence, any information associated with these metaorders is at least partially lost. Still, as shown in Fig. 4.6, we recover exactly the same square-root impact function as for the original metaorders! Note in passing that these graphs show once again that the square-root impact law

¹⁷The shuffling was based on the Fisher-Yates algorithm. In other words, we collect all the market orders within the same session for a specific stock and simply shuffle only trader IDs. Specifically, we apply the function numpy.random.shuffle in Python to the DataFrame column containing the trader IDs

Chapter 4.

only depends on the volume Q of the metaorder (either real or synthetic) but not of the execution time T, see [7, 45, 70]. We have tested different constructions of synthetic metaorders on different stocks, from the Paris and London Stock Exchange, with similar results – a more detailed discussion will be presented in [85]. Note that the preservation of the impact law under reshuffling provides further evidence that short-term impact should be decoupled from alpha (i.e., predictive signals). This is consistent with the fact that the typical time horizon over which traders seek alpha is significantly longer than the execution timescale of metaorders. This is particularly true in this study, as metaorders are defined within trading periods. Similar observations have been made, for instance with CFM's trades—where the square-root law remains valid—or in the ANcerno dataset, see [38, 40, 44]

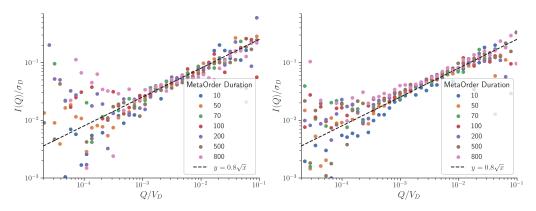


Figure 4.6: Left: Impact $\mathcal{I}(Q)$ of real metaorders as a function of their rescaled size Q/V_D for a typical Japanese Stock, using data from 2012 to 2018. The color of the dots correspond to different total execution time T, expressed in seconds. Right: Impact $\mathcal{I}(Q)$ of synthetic metaorders for the same stock, obtained via ID reshuffling. ID shuffling consists in a random permutation of historical trader IDs, preserving the frequency of apparition. Note that the vertical and horizontal scales are the same in the two plots: the square-root fit is exactly the same as for real metaorders. From the legend, one can also clearly see that $\mathcal{I}(Q)$ is independent of T [70]. However, it is worth noting that synthetic metaorders are generally smaller in size compared to real ones.

4.4.3 Discussion

The results of the previous two sub-sections strongly suggest a purely "mechanical" interpretation of the square-root impact law, based on a time decaying $\sqrt{q/t}$ impact of single market orders, independently of their association with a specific metaorder (since trader ID's can be scrambled without affecting the results).

These findings are difficult to reconcile with theories explaining the square-root law based on information, such as in [71], or on the detection by the market of the beginning of new metaorders, such as in [30, 43].

Those results are, on the other hand, perfectly in line with the fact that, due to anonymity, all market orders – even uninformed ones – should play an equivalent role and should on average impact prices similarly [87]. This was already noted in [36] by comparing the impact of CFM market orders with non-CFM market orders, and even more convincingly in an unpublished specifically designed 2010 experimental campaign with totally random market orders.

4.5 The other side of market orders: liquidity providers

4.5.1 Refill sequences

Whereas the long-term correlation of market orders is well documented [18, 31, 88] and mainly attributed to the order splitting of large metaorders [16, 18], our dataset also allows us to study the splitting strategy of liquidity providers.

To do so, we simplify the problem by restricting to the set of filled limit orders, i.e. limit orders that have been placed in the order book and subsequently executed by another participant. Thus, as for liquidity takers, one can aggregate those filled limit orders into "refill sequences", i.e. sequences of consecutive filled limit orders of same sign submitted by the same trader during a trading session. Given the splitting behavior of liquidity takers, market makers/liquidity providers face a sign-correlated succession of market orders. It is thus likely that the executed limit orders flow will be also be persistent, and lead to a power-law tail in the size distribution of the refill sequences, as we indeed confirm empirically, see Fig. 4.7. We show there a power-law fit of the distribution of the number n of child orders associated with refill sequences, as $\psi(n) \propto n^{-\mu_p}$ with μ_p in the range [1.4, 2.4] depending on the considered stock. This power-law decay echoes the Lillo-Mike-Farmer distribution of child market orders, although it is not expected to mirror it exactly since different traders provide liquidity to the same incoming metaorder – see [89] for a related discussion.

It should however be emphasized at this point that the boundary between liquidity provider and liquidity taker is somewhat blurred. Except for specific cases, the large majority of market participants use a mix of limit and market orders to acquire or sell shares. For example, it is quite common to see participants adding liquidity at the bid (ask) while sending market orders at the ask (bid). Large funds like AQR declare that most of their executed volume is though limit orders [14].

Chapter 4.

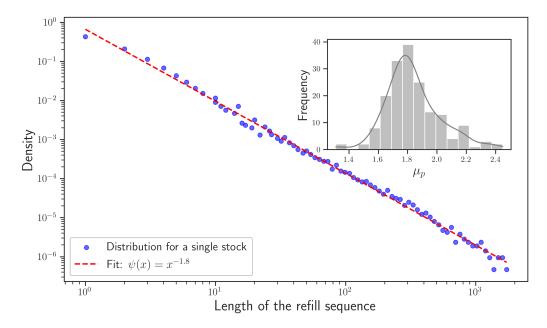


Figure 4.7: Distribution of the length n of refill sequences for a given stock of the TSE. A power-law distribution fits the data very well: $\psi(n) \sim n^{-\mu_p}$. **Inset:** Distribution of μ_p across the different stocks of our dataset, regressed independently.

4.5.2 Strategic behaviour of liquidity providers

As liquidity providers get executed on the ask (bid) side, they tend to increase (decrease) their next limit order such as to (i) control their inventory as the next trade will be biased towards the bid (ask), (ii) protect themselves against being picked up by making the next trade less favorable for the buyer (seller). This is often called "skewing" in the market making jargon, and/or (iii) ask for a better price in the case demand for liquidity is persistent.

Hence we expect the next limit order to be executed at a higher (lower) price, i.e.

$$\mathcal{K}(i) := \mathbb{E}\left[\epsilon \cdot \frac{p_{i+1} - p_i}{\sigma_D} \mid i\right] > 0,$$

where ϵ is the sign of the executed market order, and p_i is the log-price at which the *i*th child of a refill sequence is filled. We find that the "refill function" $\mathcal{K}(i)$ depends only weakly of the volume of the filled limit order, probably due to strategic liquidity provision.

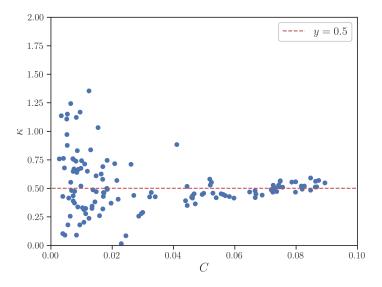


Figure 4.8: Coefficient of the refill function Eq. (4.5) when regressed separately. We used refill sequences from the top 4 most liquid stocks of our dataset, selecting the top 100 more active traders. Each dot is obtained by averaging two close-by data points, so that no individual data can be inferred from this graph.

We have found that the $\mathcal{K}(i)$ however depends on the liquidity provider ℓ , some being more aggressive than others. More precisely, we fitted $\mathcal{K}_{\ell}(i)$ as:

$$\mathcal{K}_{\ell}(i) = \frac{C_{\ell}}{i^{\kappa_{\ell}}},\tag{4.5}$$

where p is the label of the liquidity provider. The inverse dependence on i means that, as the number of previous executions increases, liquidity providers are more willing to post competitive quotes. One can observe two main types of traders, see Fig. 4.8:

- High $C_{\ell} \gtrsim 0.02$ traders are "wary" and place their next limit order quite a bit deeper in the book. The corresponding values of κ_{ℓ} cluster around 1/2.
- Low $C_{\ell} \lesssim 0.02$ traders, on the other hand, correspond to "aggressive" liquidity providers who compete for the spread. Corresponding values of κ_{ℓ} are also larger, meaning that even after being executed many times, they are still providing competitive quotes.

Low C_{ℓ} market makers are thus responsible for ensuring stable liquidity. Figure 4.9 confirms that low C_{ℓ} traders account for the largest fraction of consumed liquidity. Note that $C_{\ell} = 0.02$ corresponds to a price degradation of 2% of the

Chapter 4.

daily volatility σ_D after the first executed limit order, and of 0.2% after the 10th execution when $\kappa_{\ell} = 1$.

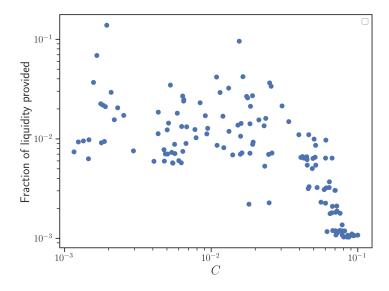


Figure 4.9: Plot of the fraction of liquidity provided by trader p as a function of parameter C_{ℓ} . As expected, traders responsible for most of liquidity have a low C_{ℓ} . We used refill sequences from the top 4 most liquid stocks of our dataset, selecting the top 100 more active liquidity providers. Each dot is obtained by averaging two close-by data points, so that no individual data can be inferred from this graph.

Interestingly, the fact that Eq. (4.5) decreases with i suggests that the available liquidity increases away from the best price, which is the fundamental ingredient leading to a concave impact function.

4.6 Conclusion

The JPX database provides a trove of interesting features, which have only started to be exploited by Sato and Kanazawa to understand the origin of the long memory of market order signs [16] and, more recently, to firmly establish the square-root law of market impact, Eq. (4.1) and rule out some of the proposed theories that predict a non-universal value for the concavity exponent δ [72].

Our aim in this study was to leverage the fact that *all* metaorders can be identified to shed light on the microscopic origin of the square-root law. Our central result, which we did not expect when starting this project, is that such a law has in fact microscopic roots, and applies already at the level of single child orders,

provided one waits long enough for the market to "digest" these orders. This is not consistent with the standard interpretation of the LLOB model [74], which assumed that the theory described the liquidity dynamics at a "mesoscopic" scale. The mesoscopic impact of metaorders rather arises from a "double" square-root effect at the level of child orders, see Eq. (4.2): square-root in volume of individual impact, followed by an inverse square root decay as a function of child order time, such that the cumulative impact of a metaorder yields back the partial (Eq. (4.4)) and total (Eq. (4.3)) square-root laws.

This finding however immediately suggests that since market orders are anonymous, the double square-root law Eq. (4.2) should apply to any market orders and the impact of synthetic metaorders, reconstructed by scrambling the identity of the issuers, should also be described by the square-root impact law, Eq. (4.1). We have provided empirical evidence that this both statements are indeed valid. In particular, synthetic metaorders behave exactly as real metaorders, see Fig. 4.6. We conclude that there is nothing special about child orders belonging to a given metaorder, at odds with theories that emphasize the information content of such trades to explain the square-root impact law, but in agreement with previous conjectures about the purely mechanical aspect of price impact [36, 87, 90].

Interestingly, our synthetic metaorder experiment suggests that it may be possible to reconstruct the impact of metaorders from the public tape only, without trader IDs. In the next chapter, we show that this is indeed the case, provided market orders are properly aggregated into synthetic metaorders, opening the path to a new wave of empirical studies, in particular concerning cross-impact [91, 92].

While our results show that the square-root impact law does not emerge at the meso-scale but is already present at the micro-scale, they trigger new unanswered questions. In particular, why is the impact of single market orders of volume q also a square-root? We have argued that this is because the latent order book is locally linear, such that after a "hot potato" game between liquidity providers, the final counterparty of the initial market order is on average at a distance \sqrt{q} from the initial mid-point. Although this scenario is intuitively plausible, we believe that a deeper dive into the JPX database (or a similar one) would allow one to (in-)validate such a picture. Another conundrum is the square-root time decay of individual market orders, which seems to violate the martingale constraint that relates the decay of the autocorrelation of the sign of market orders to impact decay [31, 84]. We will come back on this in Chapter 6.

Data Availability Statement

The data supporting our results were provided by Japan Exchange (JPX) Group, Inc. JPX Group is a third-party commercial company and provided their dataset through a non-disclosure agreement with Kyoto University, strictly for academic purposes. This non-disclosure agreement imposes legal restrictions on data availability, and therefore, we cannot make the data publicly accessible without approval from JPX Group.

Take Home Message

- The square-root law of market impact, long considered a mesoscopic phenomenon, already holds at the microscopic level of individual child orders.
- Impact follows a "double" square-root structure: a square-root dependence on volume per child order, combined with an inverse square-root decay over time.
- The order flow can be viewed as a succession of metaorders (as defined here), each of them following the SQL.
- The robustness of the square-root law under issuer reshuffling suggests that short-term impact is largely decoupled from alpha.
- Liquidity providers also seem to submit "provider" metaorders to the market, whose sizes are likewise power-law distributed.
- These results open the door to reconstructing metaorder impact from public data alone, with promising implications for cross-impact analysis.

Chapter 5

Metaorder Proxy: Examining the Puzzling Efficiency of Synthetic Metaorder Reconstruction

We have to remember that what we observe is not nature herself, but nature exposed to our method of questioning.

Werner Heiseinberg

This chapter builds directly on the previous one, where we demonstrated that reshuffling the identities of market order issuers preserves the Square Root Law. Motivated by this insight, we now tackle the problem of reconstructing realistic metaorders from public trade data and present our novel algorithm. This approach addresses the challenges arising from reliance on proprietary datasets in price impact research. We describe how the algorithm successfully recovers key stylized facts, including the Square Root Law, concave execution profiles, and post-execution decay. Finally, we discuss our findings, which suggest that average realized short-term—and even long-term—price impact is primarily mechanical rather than driven by information revelation, potentially explaining the universality of the SQL.

From: Generating Realistic metaorders from public data G. Maitrier, G. Loeper, JP. Bouchaud

85

Contents 5.1 Introduction 74 5.2**76** 5.3 Recovering metaorder stylized facts Peak impact: the Square Root Law 5.3.1 79 5.3.2 81 Concave profile during metaorder execution 82 5.3.3 5.3.4 83

5.1 Introduction

5.4

Although we have already introduced the Square-Root Law (SQL) several times, let us briefly recall its key features, mostly discovered using proprietary datasets.

- 1. the SQL is in a first approximation *independent* of the time T needed to execute the metaorder, and only depends on the volatility of the asset and the fraction of the total traded volume captured by Q [7, 49];
- 2. the SQL also holds "inside" each metaorder: the average price profile is itself a square-root as a function of the currently executed volume. This means that the last child orders impact less than the first ones [2, 39, 69];
- 3. the square-root impact decays post-execution over the time scale T of the metaorder itself [39, 43, 44, 83, 93], with a sharp decay at first and a very slow decay at long times with perhaps a small but non-zero permanent component [44, 53].

These empirical results are of primary importance for both academics and practitioners. The SQL indeed predicts that impact costs are extremely high even for small volumes Q, because of the infinite slope of the square-root function at the origin. Neglecting such costs can easily turn a profitable strategy on paper into a money losing machine once implemented, see e.g. [50]. From an academic point of view, the explanation of such a square-root dependence is far from trivial – is it due to information revelation, as many standard economic theories postulate [30, 43, 71, 75, 94, 95], or mostly "mechanical", as postulated by "latent liquidity" theories [38, 53, 62, 74, 90], see also [7].

Despite its significance, empirical research on market impact, specifically when it comes to "metaorders", often faces limitations due to data access constraints. Indeed, to track the impact of those metaorders, one should access proprietary datasets, typically held by private institutions, limiting the scope and reproducibility of academic research. Furthermore, such proprietary datasets are often not very large and possibly biased by the trading style of the managers: as emphasized in [7, 49] market impact and short term trading signals can be difficult to disentangle – in fact, mainstream economists would claim that "impact" is nothing but the correlation between the sign of informed trades and the subsequent price change [87].

Very recently, Sato and Kanazawa [42] have been able to access the records of all trades of the Tokyo Stock Exchange (TSE), with (anonymized) trader labels that allowed them to reconstruct all metaorders unambiguously. Their analysis allowed them to confirm once again the validity of the SQL with great precision, and to establish that such a law holds for all stocks individually, when some theories based on the volume distribution or on the autocorrelation of the sign of trades would have predicted systematic deviations from a square-root law [43, 71]. The same dataset has also been used to unveil further, more subtle properties of the SQL [2].

Such a unique dataset is however, quite unfortunately, inaccessible for open academic research. There have thus been many attempts to create proxies of metaorder impact using the public tape, i.e. the list of all buy and sell market orders executed on lit markets, but without any tags allowing one to track individual traders. To the best of our knowledge, these attempts have been unsuccessful. Identifying metaorder impact with the correlation between order imbalance in a time interval T is clearly completely wrong – impact is linear for small imbalance and saturates for large imbalance [37]. It is all but impossible to accurately identify real metaorders within the order flow [96]. It is also extremely difficult to generate data that recreate all the stylized facts mentioned above using VAR models – see Chapter 8 – or propagator models calibrated on real data, see e.g. [1, 21]. In recent years, machine learning has become a widely used tool for generating limit order book and understand impact, see [97, 98]. Nevertheless, all models still struggle to fully capture the characteristics of metaorder execution [99].

To address such challenges, we propose in this Chapter an algorithm that uses public trade data to generate synthetic, metaorders which lead to price impact indistinguishable from that observed using proprietary datasets. Our method not only circumvents the need for proprietary data but also facilitates the creation of larger and more robust datasets. By adequately aggregating public trade data, we demonstrate that the resulting synthetic metaorders preserve all major character-

Chapter 5.

istics found for real metaorders (SQL, concave execution profiles, post-execution decay), thus providing a valuable tool for both academics and practitioners.

This chapter is mostly algorithmic and empirical in nature, as we explain our procedure in a clear and reproducible way, and present a sample of the results we have obtained that fully validate our proposal. The theoretical justification for the success of our procedure is not yet completely understood and we will propose a modelling framework in the next Chapter. But we believe that the fundamental idea is the following: since the SQL is measured for all metaorders, independently of the trading firm, and since market orders executed in markets mostly originate from such metaorders and are anonymous, the emergence of the SQL cannot heavily rely on the precise matching between market orders and metaorders. This is indeed what was observed in our previous paper [2] using the detailed TSE data as a validation, and that we generalize in the present Chapter, which is structured as follows:

- Section 5.2 provides a detailed explanation of our synthetic metaorder generation algorithm, with an emphasis on the significance of what we call the "mapping" function.
- Section 5.3 presents several evidences of the method's effectiveness in replicating and validating the well-documented empirical facts about metaorders.

5.2 The algorithm

In this section, we present our algorithm designed to generate random metaorders from publicly available data, which lead to impact properties indistinguishable from actual proprietary data. We define a metaorder as a sequence of trades of the same sign initiated by a given trader within the same trading session. In order to generate *synthetic* metaorders, we propose the following algorithm, which requires only public trade data for any asset class (stock, futurs, options etc...). Although using aggregated order book data across multiple venues yields similar results, we recommend using data from a single exchange (ex: Euronext, Nasdaq, CME etc...).

Algorithm 1 Generating Synthetic Metaorders

Input: Trade data for a given stock and date

Output: Metaorder statistics

- 1. Load and clean trade data for given stock and date (e.g., remove opening and closing periods).
- 2. Compute daily traded volume V_D and intraday volatility σ_D , defined as:

$$V_D = \sum_{i} q_i, \quad \sigma_D = \frac{\max(p_t) - \min(p_t)}{p_0}$$
 (5.1)

- 3. Randomly assign trades to traders using a mapping function while preserving the chronological order of trades.
- 4. Sort trades by traders and timestamp.
- 5. Define a metaorder as a sequence of trades of the same sign from the same trader.
- 6. Compute metaorder features:
 - Log price at metaorder start and end.
 - Number of child orders in the metaorder.
 - Volume traded within the metaorder.
 - Any other relevant quantities.
- 7. Aggregate metaorder statistics and return only those with more than one child order.

The mapping function

The important feature of the algorithm is the mapping function, which assigns synthetic trader IDs to each market order executed on a particular day. This function is crucial: if one possesses a proprietary dataset [38] or an exhaustive one such as the TSE dataset [42], one knows this mapping exactly at least for a given set of market orders, and this allows one to measure the SQL in the usual manner, namely (see [7])

$$\frac{I(Q)}{\sigma_D} = Y \sqrt{\frac{Q}{V_D}}, \quad \text{with } \begin{cases} I(Q) = \mathbb{E}[\varepsilon \cdot (p_e - p_s)] \\ Y \in [0.5, 1], \end{cases}$$
 (5.2)

where p_s is the start mid price, just before the execution of the first child order and p_e is the end mid price, just before the execution of the market order immediately following the last child order. Now, as highlighted in [2], the introduction of random variations in the matching between real traders and orders, the square impact law is preserved, including its prefactor Y. Here, we show that even for a mapping function totally agnostic of the true mapping, we still recover the correct

Chapter 5.

impact.

We gave this mapping function only two degrees of freedom: the number of different traders in a given day and the distribution of their trading frequency, that is, the fraction of orders they participate to. We will show later that the impact law is only weakly dependent on those parameters, as expected. However, these parameters directly influence both the number and average length of the generated random metaorders. Therefore, considering the stock's liquidity (i.e. the number of trades per day and the average volume per trades), it may be necessary to set these parameters within an appropriate range to generate coherent metaorders.

Below is the pseudo-code for the mapping function. It is important to emphasize that this is merely a mapping function—quite simple in this case. While it performs well on the selected assets (and others tested), further fine-tuning may be necessary depending on the specific microstructural characteristics of the asset under study. A different, and potentially more natural mapping function that also yields satisfactory results will be presented in Chapter 7.

Algorithm 2 Mapping Function

Initialization: Let N be the number of traders, and F a probability law.

- 1. Generate $f_i \sim F$ for $i = 1, \ldots, N$.
- 2. Define $p_i = \frac{f_i}{\sum f_i}$ for each trader.
- 3. Compute cumulative probabilities:

$$c_i = \sum_{j=1}^{i} p_j, \quad c_0 = 0.$$

- 4. For each order in the market:
 - (a) Draw a random variable $U \sim U(0, 1)$.
 - (b) Find the trader i such that $c_{i-1} \leq U < c_i$.
 - (c) Assign the order to agent i.

We have evaluated the robustness of our algorithm using different values of N and two types of trader frequency f distributions: a power-law $P(f) \propto f^{-\alpha}$ and a homogeneous $f \equiv f_0$. In real markets, the distribution of trader participation is well approximated by a power-law distribution, where a small number of traders account for a significant portion of executed orders, while most traders participate less frequently, see [2].

Note that a key aspect of this mapping function is that it corresponds to *sampling* without replacement. We find this feature essential for recovering the SQL.

All source code used in this Chapter is publicly available at :https://github.com/glatouille/Metaorder_proxy

Now the procedure is established, we focus next on empirical results and demonstrate the robustness of our method with which we can generate an unlimited number of realistic metaorders.

5.3 Recovering metaorder stylized facts

In this section we report the results of our empirical investigations using synthetic metaorders. A series of sanity checks have been performed to rule out any trivial artifacts (see Appendix). For example we check that by randomly flipping the sign of market orders we measure zero impact, as it should be – impact is *not* merely related to volatility [70].

5.3.1 Peak impact: the Square Root Law

We first tested our algorithm on a very liquid asset: the EUROSTOXX futures contract from September 2016 to August 2018. We used public trade data, selecting only the Eurex exchange. We obtain a remarkably clean square root law over four decades, see Figure 5.1, with a noisy region for very small $Q/V_D \leq 5 \times 10^{-6}$, which might possibly be considered as a linear, as in [100]. Note that we do not only recover the square-root dependence on Q but also the correct a realistic prefactor in Eq. (5.2), with Y = 0.5, as found in [38, 42].

Chapter 5.

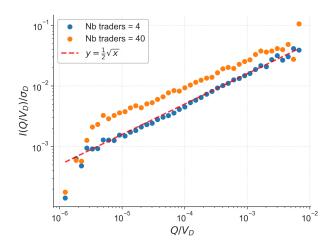


Figure 5.1: Retrieving the Square Root Law for futures on the EUROSTOXX, trade data from 2016 to 2018. We used a mapping function with 4 trades (resp. 40 trades) for the blue curve (resp. orange), and homogeneous distribution of their trading frequencies, resulting in approximately 3 million metaorders in both cases. We see that, by fine-tuning the number of traders, one may recover the correct prefactor $Y \approx 0.5$

We also tested these results on single stocks by selecting a basket of seven stocks traded on the Paris Stock Exchange, between January 2021 and December 2023. We used our algorithm to generate approximately 3 millions of random metaorders per stock, imposing 20 traders and a homogeneous trading frequency distribution. Once again, we obtained very precise results, with almost no variation in the prefactor $(Y \approx 0.5)$, see Figure 5.2. This could be explained by the fact that selected stocks are among the most liquid ones traded on the PSE, and thus may be quite similar. We also tested for a specific stock (BNP Paribas) the dependence of the SQL on the input parameters, i.e. the number of traders and the distribution of their trading frequencies. Again, we found no significant variations in the impact function. However, we do acknowledge that for some assets, particularly illiquid ones, one may have to fine-tune those parameters to recover the SQL. It is necessary to be statistically close enough to the real mapping function. The same remark also applies to the Y-ratio. While we mostly show impact functions with a realistic prefactor, i.e., close to 0.5, we have also encountered stocks where the prefactor was slightly higher of lower, but typically of order one, see Figure 5.1.

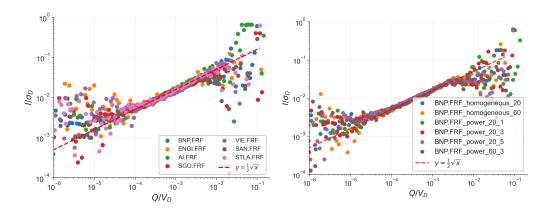


Figure 5.2: Left: Retrieving the Square Root Law for various stocks traded on the Paris Stock Exchange. Synthetic metaorders were generated using our algorithm, specifying 20 traders and a homogeneous distribution of their trading frequency. Stocks were traded on Euronext between 2021 and 2023. Right: Verifying the robustness of the algorithm in respect to variations in the mapping function parameters. Legend represents (Stock Name; Type of Distribution; Number of traders; Power law exponent). Data from Paris Stock Exchange, between 2021 and 2023.

5.3.2 Role of metaorder duration

One of the major enigmas of metaorder price impact is that, contrary to a priori expectations, the impact remains independent of the metaorder duration T. This phenomenon is a natural consequence of the SQL as written in Eq. (5.2): indeed as T is varied the volatility contribution scales as \sqrt{T} whereas the total traded volume scales as T, which means that the explicit T dependence cancels from I(Q). This property was extensively studied in [70] and confirmed in [2] both for real and synthetic metaorders based on the TSE dataset. We show in Fig. A.2 that this property also holds for our synthetic metaorders.

Chapter 5.

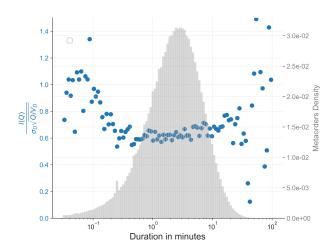


Figure 5.3: Approximate independence of metaorder impact I(Q) with respect to the metaorder duration T (expressed in minutes). The blue dots represent the impact divided by the volume contribution as a function of duration, which is approximately constant for $T \gtrsim 30$ seconds, and on average equal to Y = 0.6. The grey histogram shows the distribution of metaorders durations in minutes. Synthetic metaorders were generated with BNP Paribas share price between 2020 and 2023, for 20 homogeneous traders.

Figure 5.3 also reveals that, on average, synthetic metaorders constructed using our method are shorter than those typically found in proprietary datasets. For instance, it is not uncommon for firms like CFM to execute metaorders over an entire trading day. However, the durations of our synthetic metaorders are consistent with the average metaorder duration observed in the TSE dataset, which is typically also around 2-5 minutes, see [2]. In any case, one can adjust metaorder duration by tuning the mapping function's parameters. Furthermore, the other mapping function detailed in Chapter 7 will generate longer metaorders, more precisely featuring a power-law distribution of their durations.

5.3.3 Concave profile during metaorder execution

Beyond the peak impact I(Q), our method is also particularly effective in reproducing other stylized facts, such as the concave profile during metaorder execution, see [39, 69, 83, 93]. Indeed, as proposed in [7], the average price impact during the execution of the metaorder reads:

$$\mathcal{I}(\phi Q) = \sqrt{\phi}I(Q),\tag{5.3}$$

where $\phi \in [0,1]$ is the fraction of executed volume, with $\phi = 0$ at the start of the metaorder and $\phi = 1$ at the end.

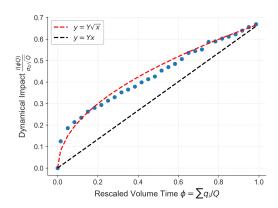


Figure 5.4: Concave profile during metaorder execution. Random metaorders constructed on BNP Paribas stock price, data from 2020 to 2023, with mapping function parameters of 20 traders and a power law distribution for their trading frequencies with an exponent of $\alpha = 2$. We selected only metaorders having more than 5 child orders.

Such a concave profile can be explained within the latent liquidity framework, which predicts that the latent limit order book is locally linear (LLOB) [7, 74]. The fact that this profile also holds for synthetic metaorders is another indication that, on average, the latent order book is indeed present and provides liquidity for all metaorders. This could be related to the action of market makers, as proposed in [2]. That said, the concavity is in fact crucial for market efficiency: it may be a key element to ensure the diffusivity of prices as argued recently in [101]. However, we believe that the framework developed in that paper is insufficient to accurately describe reality, as it lacks a crucial component: metaorder decay, to which we turn next. A unified theory of price impact that incorporate all known ingredients (autocorrelation of the sign of the trades, square-root impact, impact decay) is still under construction, a topic on which we hope to report soon [85].

5.3.4 Metaorder decay post execution

Impact decay has been subject to controversy, even when it is of crucial importance for optimal execution schedules. Indeed, assuming permanent impact or accounting for impact decay leads to radically different trading policies. What makes the empirical study of this problem particularly difficult is the fact that price variance increases linearly with the time elapsed since the end of execution, leading to large errors in the determination of impact decay. We know that metaorder impact during execution is generally small relative to the volatility (see Eq. (5.2) for small Q/V_D), this predicament is especially strong for metaorder decay: the signal-to-noise ratio significantly worsens when analyzing extended timescales after execution, highlighting the need for large metaorder datasets.

Chapter 5.

An initial line of research suggested that there is a permanent impact after execution, at approximately 2/3 of the peak impact, meaning that the price after execution is equal to the average paid during execution, see [43]. However, a later empirical study found that upon closer inspection the impact eventually decays to zero over a much longer timescale (several days). This can potentially be mistaken for a permanent impact of about 2/3 by the end of the trading day [44]. However, if the decay of the impact is evaluated over multiple days, a clear decay of impact is observed – although the long term fate is difficult to ascertain and it is plausible that a small permanent impact exists [52, 53, 76].

Interestingly, using our synthetic metaorders, we precisely replicate the fit observed in [44] with real data. Assuming a propagator decaying as $G(t) \approx t^{-\beta}$ with $\beta < 1$, the rescaled impact after execution for a metaorder of size Q and duration T can be expressed as:

$$\mathcal{I}(Q,z) = I(Q) \left(z^{1-\beta} - (z-1)^{1-\beta} \right), \tag{5.4}$$

with $z = t/T \ge 1$ and t = 0 corresponds to the start of the metaorder. Such a decay is fast at the beginning (with a sharp singularity $-(t-T)^{1-\beta}$ and a slow relaxation at long times (as $t^{-\beta}$).

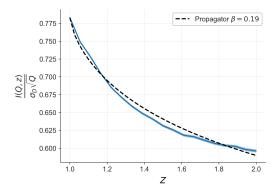


Figure 5.5: Price relaxation post metaorder execution, well fitted by the impact predicted by the propagator model, with the value $\beta \approx 0.2$. We used synthetic metaorders generated on BNP Paribas share price between 2021 and 2023, using a mapping function of parameters of 4 traders and power law distribution of exponent $\alpha = 2$. We kept only metaorders with more than 5 child orders.

Hence, we confirm using our synthetic metaorder database that impact not only decays after execution, but rather interestingly it decays exactly as found by Bucci et al. [44] using real metaorders, with a value of $\beta = 0.2$ very close to the one reported there ($\beta = 0.22$). This is also the value predicted by the propagator

model $\beta = \frac{1-\gamma}{2}$ [31], where γ describes the power-law decay of the autocorrelation of trades, which is typically found to be around 0.5 for most stocks, see [102].

It is indeed true that some may be interested in evaluating impact decay over much longer horizons—typically for asset managers holding positions for several months. This is a question we intend to explore in the near future.

5.4 Conclusion

In this Chapter, we introduced a straightforward yet surprisingly effective algorithm for generating realistic metaorders from public trade data. This approach offers a robust solution to a longstanding dataset challenge in price impact research, which has traditionally relied on proprietary data. Using this algorithm, we were able to recover all the salient stylized facts reported in the existing literature, specifically: the Square Root Law and its independence with respect to metaorder duration; the concave profile during metaorder execution; and the slow power-law decay after execution. We also confirm that this decay is effectively captured by the prediction of the propagator model, in line with previous studies. Of course, to generate even more realistic metaorders, some refinements could be made to this algorithm, particularly when it comes to the mapping function. However, our goal was to provide a highly reproducible yet effective algorithm in this area of research, which, by essence, is often not transparent when it comes to data. We therefore believe this could be a valuable tool for both practitioners and academics to enhance the quality of their empirical studies.

On the other hand, it also serves as further evidence that there is complete decorrelation between putative prediction signals (a.k.a. short term alpha) and the square-root law governing the average realized price impact (at least in the short term). By construction, a synthetic metaorder has no connection to a trader's intention to trade based on a predictive signal. Yet, its impact is indistinguishable from that of a real metaorder, which, by contrast, may be executed by a trader who follows such a signal. Hence, this result supports a mechanical origin of price impact, which may also explain its universal nature. This naturally leads to the next key question for future research: what are the theoretical foundations underlying our findings? Indeed, if price impact is purely mechanical, then one might expect a theory purely based on order flow dynamics, endogenously in a sense, could fully capture this phenomenon. Evidence seems to favour the theory of the latent limit order book [42], although some modifications may be required [2]. Finally, from an empirical perspective, a natural next step would be to extend this method to measure cross-impact, which is even more influenced by the scarcity of datasets than self-impact [50].

Take Home Message

- We introduce a simple and reproducible algorithm to reconstruct realistic metaorders from public trade data, bypassing the need for proprietary datasets.
- The algorithm successfully reproduces key stylized facts: the Square Root Law (independent of duration), concave execution profiles, and post-execution power-law decay.
- The synthetic metaorders have no alpha component, yet their impact matches that of real metaorders, supporting a mechanical —rather than informational —origin of price impact.
- Optimizing the mapping function and its parameters may also be necessary to generalize the procedure for any underlying asset —machine learning could be an ideal tool for this task.
- This mechanical nature may explain the universality of the Square Root Law and motivates the search for a fully endogenous, orderflow-based theoretical framework.
- Future research will focus on developing a theoretical framework to formalize this mechanism and expand the approach to estimate cross-impact, especially in contexts where public data is limited.
- Although this chapter reflects the chronological thought process inspired by Chapter 4, the mapping function and theoretical explanations in Chapter 7 effectively address the identified gaps. Therefore, revisiting this chapter with insights from Chapter 7 proves to be helpful.

Chapter 6

A Unified Framework for Market Microstructure: Reconciling the Square Root Law, Order Flow Dynamics & Price Dynamics

Extraordinary claims require extraordinary evidence.

Carl Sagan

After empirically reconstructing metaorders and examining the Square Root Law in the previous chapters, we now develop a theoretical framework aimed at reconciling several seemingly conflicting observations in market microstructure—namely, the Square Root Law for metaorders, the diffusive nature of prices, and the linear aggregated impact of order flow imbalance. Our model builds on a key insight established earlier: the order flow can be viewed as a superposition of metaorders, each following the Square Root Law (see Chapter 4).

We then derive theoretical predictions regarding the non-monotonic relationship between generalized volume imbalances and price changes, which we subsequently confirm through empirical analysis. Ultimately, we argue that these findings lend support to the "Order-Driven" theory of excess volatility, suggesting that price movements are primarily the result of mechanical trading impact rather than shifts in fundamental value.

Chapter 6.

From: The Subtle Interplay between Square-root Impact, Order Imbalance and Volatility: A Unifying Framework. G. Maitrier, JP. Bouchaud

Contents			
6.1	Introduction)
6.2	A continuous time description of order flow		L
	6.2.1	Model set-up	L
	6.2.2	Average number of metaorders	2
	6.2.3	Average trading activity and trading volume 92	2
6.3	\mathbf{Ord}	Order flow imbalance	
	6.3.1	Sign Imbalance	1
	6.3.2	Generalized Volume Imbalance 95	5
	6.3.3	The role of long-range correlations $between$ metaorders . 98	3
	6.3.4	Empirical observations)
6.4	The	Impact-Diffusivity puzzle and a generalized prop-	
	agat	$ ext{tor}$;
	6.4.1	Price diffusivity within the propagator model 106	3
	6.4.2	A generalized propagator model	3
	6.4.3	The role of metaorder autocorrelations	L
	6.4.4	The role of volume fluctuations	}
	6.4.5	The role of impact fluctuations	5
	6.4.6	Discussion	7
6.5	Cov	ariance between order flow imbalance and prices	
	char	$nges \dots \dots$	7
	6.5.1	Without volume fluctuations)
	6.5.2	With correlated metaorders)
	6.5.3	With correlated metaorders and volume fluctuations $$. $$ 119)
	6.5.4	With a random impact component	L
	6.5.5	With "informed" metaorders	L
	6.5.6	The correlation coefficient	2
	6.5.7	Empirical data	}
6.6	Con	clusion	L
On	Alpha	a Prediction and Permanent Impact 134	Į.

6.1 Introduction

Let us rephrase it here: *Price impact* refers to the fact that buyers push the price up and sellers push the price down [7, 87]. The traditional *Efficient Market* interpretation of this empirical fact is that buyers and sellers are on average informed, and lo and behold, the price moves according to their prediction [94, 103].

Another, very different interpretation of price impact is that it is a purely statistical reaction of the market to incoming order flow, where information plays little role. Prices move just because people trade, whatever the reason they are trading, and volatility is the result of people randomly buying and selling. This is the *Order-Driven* view of markets, explicitly spelled out in [7], chapter 20 – but see also [87, 104, 105] and in a different setting, [52, 53]. In this scenario, there is no "information revelation" but rather "self-fulfilling prophecies", as recently emphasized in [106].

Of course, reality should lie somewhere between these two extremes. There are surely some informed trades, and news do impact prices, but there is also overwhelming evidence for the presence of "noise traders", excess trading and excess volatility in financial markets, see e.g. [17, 57, 107–109]. Direct empirical estimates suggest that informed trades are a minority (see [7], chapter 16). While we are convinced that the Order-Driven view is a much closer approximation to the dynamics of markets, there are several empirical loose ends that need to be tied up. A major issue is how to reconcile two prominent stylized facts of market microstructure, namely, (i) the long-range memory in order signs (i.e. +1 for buy orders and -1 for sell orders) and (ii) the ubiquitous square-root law of market impact (that governs the average price move induced by the execution of a sequence of orders) with the random walk nature of prices, with a volatility σ that is directly proportional to the amplitude of the square-root law. More precisely, the square-root law states that the average price impact $\mathcal I$ of a metaorder of total volume Q is given by

$$\mathcal{I}(Q) = Y\sigma\sqrt{\frac{Q}{\phi}},\tag{6.1}$$

where Y is a O(1) numerical coefficient and ϕ is the average flow of orders executed in the market per unit time – see e.g. [3, 14, 38, 40, 42, 44, 47, 67–69, 93] and [7, 49] for reviews.

Eq. (6.1) may look familiar, and in fact trivial: superficially, it states that price changes (i.e. $\mathcal{I}(Q)$) grow as the square-root of execution time T, i.e. the law of random walks, *provided* one assumes that T and Q are proportional. But, as argued in [70], such an argument is completely misleading: not only $\mathcal{I}(Q)$ is an

Chapter 6.

average price change and not a standard deviation; but also $\mathcal{I}(Q)$ is found to depend *only* on the quantity executed Q and *not* on execution time T. Furthermore, it is known that the impact of metaorders *decays*, on average, after the end of the execution period (e.g. [44, 80, 83]). Note that all these features are at odds with the Kyle model, that predicts linear and permanent impact, resulting from information revelation [75].

We are thus confronted with three interrelated but separate problems:

- (a) What is the basic mechanism that explains the square-root law, Eq. (6.1), and its surprising universal character?
- (b) Can the volatility of prices σ be explained only in terms of the impact of intertwined, possibly uninformed metaorders? Or is it impact that is somehow slaved to some "Fundamental" volatility?
- (c) Can one reconcile the *square-root* law for metaorders with a *linear* relation between average price changes and order flow imbalance?

Although many theories have been proposed to explain the square-root law, there is no consensus on the issue. The Latent Liquidity Theory proposed in [74] (see also [7, 90]) seems to capture many features observed in data but seems inconsistent with others, in particular those reported in the Chapter 4 and here the recent [2]. We will not attempt to dwell further on this particular issue here but accept it as an incontrovertible empirical fact, still waiting for a fully convincing explanation. We will rather focus on points (b) and (c) above: knowing that the duration of metaorders is power-law distributed, can one reconcile the square-root impact law with the volatility of markets and with a locally linear aggregate impact law?

In order to answer these questions quantitatively, we introduce in section 6.2 a new theoretical framework to describe metaorders with different signs, sizes and durations, which impact prices as a square-root of volume but with a subsequent time decay. We show in section 6.4 that, as in the original propagator model, price diffusion is ensured by the long memory of cross-correlations between metaorders. In order to account for the effect of strongly fluctuating volumes q of individual trades, we need to further introduce two q-dependent exponents, which we justify empirically and allow us to account for the way the moments of generalized volume imbalance (section 6.3) and the correlation between price changes and generalized volume imbalance (section 6.5) scales with T. We predict in particular that the corresponding power-laws depend in a non-monotonic fashion on a parameter a that allows one to put the same weight on all child orders or overweight large orders, a behaviour clearly borne out by empirical data (section 6.5). We also predict that the correlation between price changes and volume imbalances should display a maximum as a function of a, which again matches observations (section

6.5). We conclude by arguing that our results support the "Order-Driven" theory of markets, and are at odds with the idea that a "Fundamental" component accounts for a large share of the volatility of financial markets.

6.2 A continuous time description of order flow

6.2.1 Model set-up

We posit that between t and t+dt and with probability νdt a new metaorder of random sign $\varepsilon(t)=\pm 1$ and duration s(t) is initiated. The volume of child orders is q (which might itself be random, see below), and during execution the probability that one of them gets executed is φdt , independently of the size q. We neglect throughout this Chapter activity fluctuations as well as intraday seasonalities, as these are not crucial for the effects we want to focus on. This means that μ, φ and the average duration \bar{s} are chosen to be time independent.

The total size of the metaorder is thus $Q = q\varphi s + O(\sqrt{s})$. The probability density of durations s is denoted $\Psi(s)$, which will typically has a power-law tail $\Psi(s) \propto s^{-1-\mu}$, such that the distribution of metaorder sizes Q inherits from this power-law, and decays as $Q^{-1-\mu}$ as suggested by empirical data [16, 43].

Such a power-law distribution of metaorder sizes is the basic mechanism proposed by Lillo, Mike and Farmer (LMF) [18] to explain the long memory of order signs, which is known to decay with lag τ as $\tau^{-\gamma}$ with $0 < \gamma < 1$ [88]. Within the LMF model, one has $\gamma = \mu - 1$, a result recently validated in great details by Sato and Kanazawa [16] using data from the Tokyo Stock Exchange. In a later stage, we will allow the exponent μ to depend on q, to account for the fact that large child orders tend to be less autocorrelated that small ones.

We will also allow the sign of different metaorders to be correlated, as indeed observed in data [69, 110]. More precisely, we will model the long-term decay of the autocorrelation of signs, $\mathbb{E}[\varepsilon(t)\varepsilon(t+\tau)]$ as a power-law $\tau^{-\gamma_{\times}}$, with an exponent γ_{\times} a priori such that $\gamma_{\times} \geq \gamma$ such not to contradict the LMF hypothesis.

We start warming up by computing two simple quantities, the total number of active metaorders and the average trading volume within a window of duration T. We will then turn to the distribution of volume imbalance in windows of different sizes.

6.2.2 Average number of metaorders

The total number of metaorders N_T that are active between t = 0 and t = T is given by

 $N_T = \int_{-\infty}^{T} dN_t \, \mathbb{I}(t + s(t) > 0), \tag{6.2}$

where $dN_t = 0$ if there is no new metaorder initiated between t and t + dt and $dN_t = 1$ otherwise. This equation means that to be active in [0, T], it must start before t = T and end at least after t = 0.

The average over the probability of initiating metaorders and over their duration gives

$$\overline{N}_T = \nu \int_{-\infty}^T dt \int_0^\infty ds \, \Psi(s) \, \mathbb{I}(t+s>0) = \nu(T+\bar{s}), \tag{6.3}$$

where we assume henceforth that the average size of metaorders is finite, i.e. $\bar{s} := \int_0^\infty \mathrm{d} s \, s \Psi(s) < +\infty$, which is tantamount to $\mu > 1$. Hence, for large $T \gg \bar{s}$, one finds $\overline{N}_T \approx \nu T$, as expected.

In the following, we will always make averages over the metaorder initiation process, and often replace dN_t by νdt whenever possible. When comparing with empirical data, we will work in trade time N_T but still call this quantity T. Translating our results is real time is, however, non-trivial because, as is well known, the activity rate ν shows strongly intermittent dynamics (often modeled using Hawkes processes, see e.g. [7], chapter 9) on top of a U-shaped intraday pattern.

6.2.3 Average trading activity and trading volume

Second warm-up question: what is the total activity A_T and total trading volume V_T executed between t=0 and t=T? In the following we assume that all child orders have the same size and denote $\kappa:=q\varphi\nu$, so that activity and trading volume are simply related by $\overline{V}_T=q\overline{A}_T$. More generally, the following results holds with $\kappa=\bar{q}\varphi\nu$.

There are two terms, corresponding to metaorders initiated within the period [0, T] or before t = 0 that are still active in [0, T]. We denote these two terms as V_T^1 and V_T^2 , with

$$V_T^1 = \int_0^T dN_t \left[\mathbb{I}(t+s(t) > T)q\varphi(T-t) + \mathbb{I}(t+s(t) < T)q\varphi s \right], \qquad (6.4)$$

which after averaging over dN_t gives

$$V_T^1 = \kappa \int_0^T dt \left[\mathbb{I}(t + s(t) > T)(T - t) + \mathbb{I}(t + s(t) < T)s \right]. \tag{6.5}$$

Similarly, for V_T^2 we get

$$V_T^2 = \kappa \int_{-\infty}^0 dt \left[\mathbb{I}(t + s(t) > T)T + \mathbb{I}(0 < t + s(t) < T)(s(t) + t) \right]$$
 (6.6)

Now let us compute the average over duration s, given by

$$\overline{V}_T^1 = \kappa \int_0^T dt \int_0^\infty ds \, \Psi(s) \left[\mathbb{I}(t+s > T)(T-t) + \mathbb{I}(t+s < T)s \right]$$
 (6.7)

and

$$\overline{V}_T^2 = \kappa \int_{-\infty}^0 dt \int_0^\infty ds \, \Psi(s) \left[\mathbb{I}(t+s > T)T + \mathbb{I}(0 < t+s < T)(s+t) \right]$$
 (6.8)

Carefully taking the derivative with respect to T one finds:

$$\frac{\partial \overline{V}_T^1}{\partial T} = \kappa \left[\int_0^T \mathrm{d}s \, s \Psi(s) + T \int_T^\infty \mathrm{d}s \, \Psi(s) \right]; \qquad \frac{\partial \overline{V}_T^2}{\partial T} = \kappa \int_T^\infty \mathrm{d}s \, (s - T) \Psi(s). \tag{6.9}$$

Hence, the total average executed volume $\overline{V}_T = \overline{V}_T^1 + \overline{V}_T^2$ is given, for large T, by

$$\overline{V}_T = \kappa \bar{s} T \approx q \varphi \bar{s} \, \overline{N}_T, \tag{6.10}$$

i.e. the average volume per metaorder $\overline{Q}=q\varphi\bar{s}$ times the average number of metaorders \overline{N}_T . So in this model the average volume flow per unit time is $\phi:=\nu\overline{Q}=\kappa\bar{s}$. Note that for large enough T it is dominated by V_T^1 .

The average activity (i.e. number of trades per unit time) is given by $\nu\varphi\bar{s}$. We will denote its inverse $\tau_0 := (\nu\varphi\bar{s})^{-1}$, which is the average time between two trades. Finally, note that the average number of child orders per metaorder is $\bar{n} := \varphi\bar{s}$.

6.3 Order flow imbalance

In this section we compute, within our model, the statistics of the flow imbalance during periods of duration T. In the case where all trades have the same size q, volume imbalance is trivially proportional to sign imbalance. It turns out that due to the heavy tail in metaorder durations, sign imbalance scales anomalously with T and has non-Gaussian fluctuations, even when the signs of metaorders are independent.

When single trade volumes are themselves strongly fluctuating, the statistics of volume imbalance can be very different from those of sign imbalance, a feature we actually observe on data and reproduce within our framework – see section 6.3.2 below.

6.3.1 Sign Imbalance

The sign imbalance I_T^0 in an interval of size T is given by a sum of two contributions, as for the traded volume above:

$$I_{T,1}^{0} = \varphi \int_{0}^{T} \varepsilon(t) dN_{t} \left[\mathbb{I}(t+s(t) > T)(T-t) + \mathbb{I}(t+s(t) < T)s \right], \tag{6.11}$$

and

$$I_{T,2}^{0} = \varphi \int_{-\infty}^{0} \varepsilon(t) dN_{t} \left[\mathbb{I}(t+s(t) > T)T + \mathbb{I}(0 < t+s(t) < T)(s(t)+t) \right]. \quad (6.12)$$

Because $\mathbb{E}[\varepsilon(t)] = 0$, these terms are of mean zero. In this subsection and the next, we assume metaorders to be independent, in particular one has $\mathbb{E}[\varepsilon(t)dN_t\varepsilon(t')dN_{t'}] = \delta(t-t')dN_t$.

The sign imbalance variance is then given by $\varphi^2 \nu$ times

$$\int_0^T dt \int_0^\infty ds \, \Psi(s) \left[\mathbb{I}(t+s>T)(T-t) + \mathbb{I}(t+s< T)s \right]^2$$

$$+ \int_{-\infty}^0 dt \int_0^\infty ds \, \Psi(s) \left[\mathbb{I}(t+s>T)T + \mathbb{I}(0 < t+s < T)(s+t) \right]^2$$
(6.13)

We note that because the indicator functions are non-overlapping, all cross-products are zero. Taking a derivative with respect to T of the previous expression and denoting the result D_2 we get

$$D_2 = \int_0^T ds \, s^2 \Psi(s) + 2T \int_T^\infty ds \, s \Psi(s) - T^2 \int_T^\infty ds \, \Psi(s)$$
 (6.14)

Suppose for definiteness that

$$\Psi(s) = \frac{\mu s_0^{\mu}}{s^{1+\mu}} \mathbb{I}(s > s_0), \qquad 1 < \mu < 2, \tag{6.15}$$

corresponding to a sign autocorrelation function decaying as $\tau^{-\gamma}$ with $\gamma = \mu - 1 < 1$, as found in the data [16, 18]. Then the previous expression becomes:

$$D_2 = \frac{2}{(2-\mu)(\mu-1)} s_0^{\mu} T^{2-\mu}$$
 (6.16)

Hence in this case the variance of the sign imbalance is given by

$$\Sigma^{2} = \frac{2\varphi^{2}\nu}{(3-\mu)(2-\mu)(\mu-1)} s_{0}^{\mu} T^{3-\mu}, \tag{6.17}$$

i.e. a growth faster than T but slower than T^2 . Note that when $\mu \nearrow 2$, one smoothly recovers the expected result for a short-range correlated order flow, namely $\Sigma^2 \propto T$.

One can also compute the fourth moment of the sign imbalance. Focusing on the $I_{T,1}^0$ contribution, one finds

$$\mathbb{E}[(I_{T,1}^{0})^{4}] = \varphi^{4} \nu \int_{0}^{T} dt \int_{0}^{\infty} ds \, \Psi(s) \left[\mathbb{I}(t+s > T)(T-t) + \mathbb{I}(t+s < T)s \right]^{4},$$
(6.18)

and taking the derivative with respect to T yields

$$D_4 = \int_0^T ds \, s^4 \Psi(s) + T^4 \int_T^\infty ds \, \Psi(s). \tag{6.19}$$

When $\mu < 4$, this behaves as $T^{4-\mu}$, so that the kurtosis of the sign imbalance distribution behaves as $T^{5-\mu}/(T^{3-\mu})^2 \sim T^{\mu-1}$ which grows with T! An important consequence is that within our model the sign imbalance does not become Gaussian for large T.

Generalizing to the 2n-th moment, one finds that it grows with T like $T^{2n+1-\mu}$ when $\mu < 2n$. This suggests that when $\mu < 2$, the sign imbalance converges at large T towards a truncated Lévy distribution of index μ for the rescaled variable $I^0/T^{1/\mu}$, where the truncation takes place for $|I^0|=T$ (see [111] for a very similar calculation, and the Appendix of [112] for a proof). Indeed, one can check that the moments of such a truncated Lévy distribution scale with T exactly as above. We will test this prediction in section 6.3.4.

6.3.2 Generalized Volume Imbalance

One can generalize the calculation to the volume imbalance, or in fact to any power a of the individual traded volume, $I^a(T)$, given again by the sum two contributions:

$$I_{T,1}^{a} = \int_{0}^{T} dN(t)\varepsilon(t)q^{a}(t) \left[\mathbb{I}(t+s(t) > T)(T-t) + \mathbb{I}(t+s(t) < T)s \right], \quad (6.20)$$

and

$$I_{T,2}^{a} = \int_{-\infty}^{0} dN(t)\varepsilon(t)q^{a}(t) \left[\mathbb{I}(t+s(t) > T)T + \mathbb{I}(0 < t+s(t) < T)(s(t)+t) \right].$$
(6.21)

Note that a = 0 corresponds to sign imbalances and a = 1 to volume imbalances. Obviously, if q(t) = q at all times, all these imbalances are equal, up to a trivial

factor, to the sign imbalance computed in the previous section. In the following, we will assume that metaorders differ not only by their duration s but also by the size of their child orders, with a joint distribution that we denote as $\Psi_q(s)\Xi(q)$. Inspired by empirical data (see below), we posit that metaorders that execute with larger child order sizes q still have a power-law distributed duration s, but with a tail exponent μ_q that increases with q – i.e. have a thinner tail. More precisely, the conditional distribution $\Psi_q(s)$ is of the form:

$$\Psi_q(s) = \frac{\mu_q s_0^{\mu_q}}{s^{1+\mu_q}}, \qquad \mu_q = \mu_1 + \lambda \log q, \tag{6.22}$$

where q = 1 is the lot size.

The generalized volume imbalance $I^a(T)$ still has mean zero and variance Σ_a^2 now given by $I_1^a + I_2^a$:

$$\nu \varphi^{2} \int_{0}^{T} dt \int_{0}^{\infty} dq \, q^{2a} \Xi(q) \int_{0}^{\infty} ds \, \Psi_{q}(s) \left[\mathbb{I}(t+s>T)(T-t) + \mathbb{I}(t+s
$$+ \nu \varphi^{2} \int_{-\infty}^{0} dt \int_{0}^{\infty} dq \, q^{2a} \Xi(q) \int_{0}^{\infty} ds \, \Psi_{q}(s) \left[\mathbb{I}(t+s>T)T + \mathbb{I}(0

$$(6.23)$$$$$$

Consider the I_1^a contribution (the I_2^a contribution does not change the conclusion below):

$$\Sigma_{a,1}^2 = \int_0^T du \int_0^\infty dq \, q^{2a} \Xi(q) \int_0^\infty ds \, \Psi_q(s) \left[\mathbb{I}(s > u) u^2 + \mathbb{I}(s < u) s^2 \right]$$
 (6.24)

The derivative of this quantity with respect to T gives

$$\partial_T \Sigma_{a,1}^2 = \int_0^\infty \mathrm{d}q \, q^{2a} \Xi(q) \left[T^2 \int_T^\infty \mathrm{d}s \, \Psi_q(s) + \int_0^T \mathrm{d}s \, \Psi_q(s) s^2 \right] \tag{6.25}$$

Now assume T is large and define q_2 such that $\mu_{q_2}=2$. Metaorders with smaller volumes $q< q_2$ thus have an infinite duration variance ($\mu_q\leq 2$), while larger volumes have a finite variance ($\mu_q>2$). The two contributions then read

$$\partial_T \Sigma_{a,1,1}^2 = \int_0^\infty dq \, q^{2a} \Xi(q) T^{2-\mu_q},$$
 (6.26)

and 18

$$\partial_T \Sigma_{a,1,2}^2 = \int_0^\infty \mathrm{d}q \, q^{2a} \Xi(q) \int_0^T \mathrm{d}s \, \Psi_q(s) s^2 = \int_0^{q_2} \mathrm{d}q \, q^{2a} \Xi(q) \frac{\mu_q T^{2-\mu_q}}{2-\mu_q} + \int_{q_2}^\infty \mathrm{d}q \, q^{2a} \Xi(q) \frac{\mu_q}{\mu_q - 2}.$$

$$(6.27)$$

A convenient mathematical description of the right tail of child order sizes is a log-normal:

$$\Xi(q) = \frac{1}{q\sqrt{2\pi\sigma_{\ell}^2}} e^{-\frac{(\ell-m)^2}{2\sigma_{\ell}^2}},$$
(6.28)

with $\ell := \log q$ and m is the most likely value of $\log q$. One then gets:

$$\partial_T \Sigma_{a,1,1}^2 \propto e^{2ma} T^{2-\mu_m} \int_0^\infty d\ell \ e^{(\ell-m)(2a-\lambda \log T) - (\ell-m)^2/2\sigma_\ell^2} \propto e^{2ma+2\sigma_\ell^2 a^2} T^{2-\widetilde{\mu}(a)}$$
(6.29)

with $\mu_m = \mu_1 + \lambda m$ and an effective exponent $\tilde{\mu}$ that reads

$$\widetilde{\mu}(a) = \mu_m + \lambda \sigma_\ell^2 (2a - \frac{1}{2}\lambda \log T)$$
(6.30)

Let us fix a range of T where the data is fitted, and assume for simplicity that we are in a case where $\frac{1}{2}\lambda \log T \ll 1$ while $\lambda \sigma_{\ell}^2 = O(1)$. Then the expression for the effective exponent $\widetilde{\mu}$ becomes very simple:

$$\widetilde{\mu}(a) = \mu_m + 2a\lambda\sigma_\ell^2 \tag{6.31}$$

So the effective exponent $\widetilde{\mu}$ increases with a, i.e. an exponent $2-\widetilde{\mu}$ that decreases with a. When a=0, the sign correlation is dominated by the most probable volumes $q=e^m$ and we recover the previous result with $\widetilde{\mu}=\mu_m$.

For the contribution $\partial_T \Sigma_{a,1,2}^2$, one has to separate the cases $\ell < \log q_2$ and $\ell > \log q_2$. Since the integral over ℓ is dominated by the region $\ell \approx \ell^* = m + 2a\sigma_\ell^2$, one can use the same expression as above for the first term in Eq. (6.27), when $\ell^* < \log q_2$. When $\ell^* > \log q_2$, $\partial_T \Sigma_{a,1,2}^2$ is dominated by the second term and becomes independent of T.

Putting everything together, the predictions of this simple model are thus that

$$\Sigma_a^2 \propto e^{2ma + 2\sigma_\ell^2 a^2} \times \begin{cases} T^{3 - \tilde{\mu}(a)}, & a < a_c(1) := \frac{2 - \mu_m}{2\lambda \sigma_\ell^2}; \\ T, & a \ge a_c(1). \end{cases}$$
(6.32)

$$\int_0^\infty dq \, q^{2a} \Xi(q) \int_0^T ds \, \Psi_q(s) s^2 = \int_0^\infty dq \, q^{2a} \Xi(q) \frac{\mu_q}{2 - \mu_q} \left(T^{2 - \mu_q} - s_0^{2 - \mu_q} \right),$$

which is finite for all μ_q .

 $^{^{18}}$ Note that the apparent divergence for $\mu_q=2$ is spurious. In fact, the correct expression should read:

In other words, one finds that the variance of the generalized volume imbalance scales anomalously with T when a is small enough (like the sign imbalance considered above), but becomes simply diffusive when a is large. Intuitively, it is because large child orders are much less auto-correlated than small child orders when $\lambda > 0.19$ We will compare these predictions with empirical data in the next section. Although the model is over-simplified, we will see that it captures the data semi-quantitatively. Typically, $a_c(1)$ is found to be around 2. With $\mu_m = 3/2$, we find $\lambda \sigma_\ell^2 \approx 1/8$, an estimate that will match other observations.

It is interesting to generalize these results to higher moments of the volume imbalance. Extending the calculation above, one finds

$$\Sigma_{a,1}^{(2n)} = \int_0^T du \int_0^\infty dq \, q^{2na} \Xi(q) \int_0^\infty ds \, \Psi_q(s) \left[\mathbb{I}(s > u) u^{2n} + \mathbb{I}(s < u) s^{2n} \right], \quad (6.33)$$

from which one derives the following result

$$\Sigma_{a,1}^{(2n)} \propto \begin{cases} T^{2n+1-\mu_m-2na\lambda\sigma_\ell^2}, & a < a_c(n); \\ T, & a \ge a_c(n), \end{cases}$$
 (6.34)

with $a_c(n) = (1 - \mu_m/2n)/\lambda \sigma_\ell^2$.

6.3.3 The role of long-range correlations between metaorders

It is known that the signs of metaorders initiated by different traders are also correlated, see [69, 110]. This may either be due to herding, or more plausibly to different traders following the same signal. As mentioned above, we assume that the sign cross-correlation $\mathbb{E}[\varepsilon(t)\varepsilon(t+\tau)]$ decays as $\Gamma(\tau_0/\tau)^{\gamma_{\times}}$, whereas the sizes q and q' are remain independent for simplicity.²⁰ When $\Gamma=0$, there is no correlation between successive metaorders.

$$\Psi_q(s'|s,\tau) = \frac{\tau^b}{1+\tau^b}\Psi_q(s') + \frac{1}{1+\tau^b}\frac{1}{s}F(s'/s),\tag{6.35}$$

where F(.) is a certain function and b a new exponent. The following results are unaffected provided b > 1.

¹⁹ Another mechanism that leads to dependence of the effective exponent of $\Sigma_a^2(T)$ is the presence of power-law tails in the distribution of q that can be a confounding factor. If the tail exponent is equal to ϖ (see section 6.3.4), one expects a crossover value $a_c(n)$ given by $\varpi/2n$, i.e. when $\mathbb{E}[q^{2na}]$ diverges. Although such a mechanism may certainly play a role, it comes in parallel with the dependence of μ_q on q for which there is direct evidence, see Fig. 6.3.

²⁰One can extend the following calculations to the case where conditional size distribution of a metaorder starting at $t + \tau$, knowing that one metaorder started at t is

When all order sizes q are equal, the variance of the sign imbalance again contains two terms, one of them reading

$$\Gamma \tau_0^{\gamma_{\times}} \iint_0^T \frac{\mathrm{d} u \, \mathrm{d} u'}{|u - u'|^{\gamma_{\times}}} \int_0^{\infty} \mathrm{d} s \, \Psi(s) \int_0^{\infty} \mathrm{d} s' \, \Psi(s') \left[\mathbb{I}(s > u)u + \mathbb{I}(s < u)s \right] \left[\mathbb{I}(s' > u')u' + \mathbb{I}(s' < u')s' \right]. \tag{6.36}$$

Taking the derivative with respect to T leads to

$$\Gamma \tau_0^{\gamma_{\times}} \int_0^T \frac{\mathrm{d}u'}{(T-u')^{\gamma_{\times}}} \int_0^{\infty} \mathrm{d}s \, \Psi(s) \int_0^{\infty} \mathrm{d}s' \, \Psi(s') \left[\mathbb{I}(s>T)T + \mathbb{I}(su')u' + \mathbb{I}(s'$$

The scaling of this expression with T is found to be $T^{1-\gamma_{\times}}$ provided $\mu > 1$, i.e. as soon as the mean size of metaorders is finite. Hence, we get an off-diagonal contribution to Σ^2 that scales as $T^{2-\gamma_{\times}}$, which must be compared to the "diagonal" contribution (i.e. for u = u' and s = s') given in Eq. (6.17), which scales as $T^{2-\gamma_{\times}}$.

In other words, the LMF model [18] that ascribes the main contribution to sign autocorrelation to long metaorders is only valid if such metaorders are not too strongly correlated between themselves, i.e. when

$$\gamma_{\times} \ge \gamma.$$
 (6.38)

In view of the empirical data supporting the LMF model, we stick to this assumption henceforth. In fact, one can measure γ_{\times} directly (G. Maitrier, unpublished, see also [69]) suggesting $\gamma_{\times} \approx \gamma$.

Let us now include volume fluctuations on top of long-range correlations between the sign of metaorders. Assuming that the sizes q, q' of the child orders of two different metaorders are independent, one finds that the off-diagonal (o.d.) contribution to Σ_a^2 reads:

$$(\Sigma_a^2)_{o.d.} \propto \Gamma e^{2ma + \sigma_\ell^2 a^2} T^{2 - \gamma_\chi}, \tag{6.39}$$

to be compared with Eq. (6.32).

With $\mu_m = 3/2$ and $\gamma_{\times} = 1/2$, one therefore concludes that as soon as a > 0, the $T \to \infty$ behaviour of Σ_a^2 should, in principle, be dominated by the off-diagonal contribution. However, for a small the two exponents $3 - \widetilde{\mu}$ and $2 - \gamma_{\times}$ are indistinguishable, and the cross-over time T_{\times} beyond which $(\Sigma_a^2)_{o.d.}$ is dominant soon becomes unreachable when a grows. When $a > a_c(1)$ one finds, with $\Gamma = O(1)$

$$T_{\times} \approx e^{2\sigma_{\ell}^2 a^2}.\tag{6.40}$$

For a=2 and $\sigma_{\ell}^2=1$, this yields $T_{\times}\sim 10^4$ trades, beyond the range of times scales studied below.

6.3.4 Empirical observations

For this analysis (as well as the remainder of the Chapter), we have chosen four assets for which we have trade-by-trade prices and signed volumes. We have chosen two stocks from the LSE, one small tick stock (LLOY) with a small tick size, such that the average spread-to-tick ratio equal to ≈ 3 . The second is a medium tick stock (TSCO), with an average spread-to-tick ratio equal to ≈ 1.5 . We also selected two liquid futures contracts: the SPMINI, with a spread-to-tick ratio of approximately 1.1, and the EUROSTOXX, a large-tick asset with a ratio close to one (≈ 1.02). For equities, the dataset goes from 2012 to 2015, while for futures, it covers 2016–2018 for the EUROSTOXX and 2022 for the SPMINI. This selection allows us to cover two major asset classes actively traded in modern markets, a wide range of spread-to-tick ratios, and nearly a decade of market evolution.

Child volume distribution

As discussed in Section 6.3.2, we consider a log-normal distribution for the child order sizes, as defined by Eq. (6.28). It turns out to be a reasonable approximation of reality for large volumes, see Fig. 6.1, with values of σ_{ℓ} reported in the legend, around 1 for stocks and SPMINI and 1.2 for EUROSTOXX. A better fit of the tail of the distribution is, arguably, power-law $\propto q^{-1-\varpi}$, with ϖ found to be around 2.1-2.4. Such a choice makes the mathematical analysis more cumbersome, and we prefer the log-normal specification for our semi-quantitative discussion of the role of volume fluctuations. Nevertheless, a power-law tail can play a role similar to the coefficient λ in μ_q when it comes to the scaling of $\Sigma_a^{(2n)}$, crossing over to a linear behaviour when $2na \gtrsim \varpi$, see footnote 19.

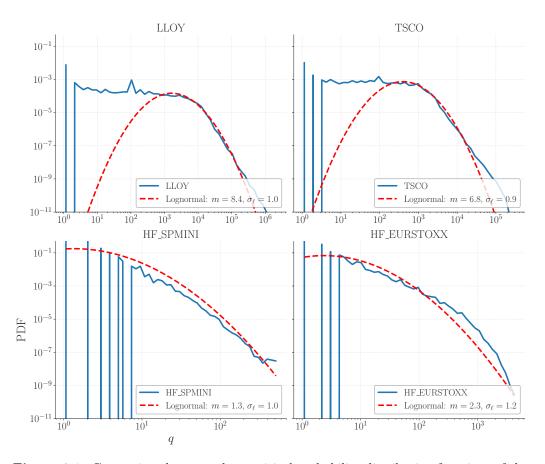


Figure 6.1: Comparison between the empirical probability distribution functions of the executed volume q for the four selected assets, along with a fitted lognormal distribution in the tail region.

Distribution of sign imbalances

We show in Fig. 6.2 the distribution of rescaled sign imbalances $I^0T^{-\chi}$, for different T in trade time and $\chi=0.72$. The theoretical analysis performed in the previous sections predicts that such distributions should collapse when χ is chosen to be $1/\mu$, where $\mu=1+\gamma$ is related to the autocorrelation of the sign of the trades, which gives $\mu\approx 1.4$. Although not perfect, the agreement is quite reasonable, in view of the fact that μ actually depends on the size of the child order q, see Fig. 6.3. The master curve is clearly non-Gaussian, with tails that become fatter as T increases, as expected from our theoretical prediction.

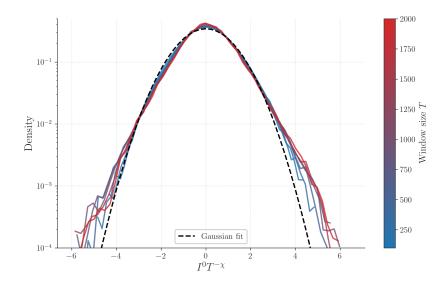


Figure 6.2: Distribution of sign imbalances I^0 as a function of the window size T (measured in number of trades) on EUROSTOXX, between 2016 and 2019. After a proper rescaling by $T^{-\chi}$, with $\chi=0.72$, distributions nicely collapse onto a single master curve, consistent with findings in [37]. The value of χ is not far from the theoretical prediction $\chi=1/\mu$ with $\mu=3/2$. The dotted line shows a Gaussian distribution with the same variance as red curve. As expected, the volume imbalance exhibits fat tails.

Sign autocorrelation function for different child volumes

Intuitively, the sign of large child orders should have shorter memory than small ones. Traders who have large quantities Q to execute is likely to trade small lots in order not to reveal information, whereas smaller Q might be possible to execute in a few shots.

In order to test this hypotheses, we define five logarithmic bins for the rescaled volume $\tilde{q} = q/\phi_D$ where ϕ_D is the daily traded volume. We then compute for each bin the autocorrelation function $C_{\mathcal{B}(\tilde{q})}(\tau) = \mathbb{E}[\varepsilon_{B(\tilde{q})}(t)\varepsilon_{B(\tilde{q})}(t+\tau)]$. We removed the largest bin, as it contains outliers, and present the four other autocorrelation functions in Fig. 6.3, in log-log, together with the unconditional autocorrelation function. As expected, the effective memory exponent γ increases with q, going from 0.6 to 1.3. We take this observation as a qualitative justification of the specification proposed above, i.e. that $\mu_q = 1 + \gamma_q = \mu_1 + \lambda \log q$, which as we now show, allows us to make quantitative predictions for the scaling of the generalized volume imbalances $I^a(T)$.

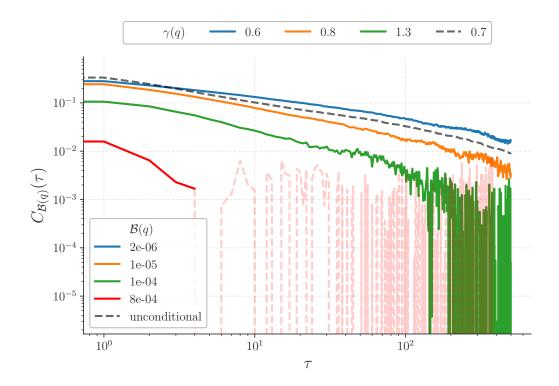


Figure 6.3: Evolution of the sign autocorrelation of market orders based on their corresponding volume bin $\mathcal{B}(q)$. We used data from the EUROSTOXX, between 2016 and 2019 rescaling market order volume q by the daily traded volume. Our findings indicate that larger market orders tend to be *less* correlated than smaller ones. Note that nearly 4% of market orders fall into the largest bin.

Scaling of the generalized volume imbalances

Our model predicts that the even moments $\Sigma_a^{(2n)}$ of the generalized volume imbalances I^a scale with T with exponents that depend on a, see Eq. (7.4). In order to test this prediction empirically, we first remark that trade-by-trade data typically exhibit numerous outliers (such as block trades, fat fingers etc...). These outliers can substantially influence the empirical estimation of the diffusion coefficient, particularly for large values of a. Consequently, trades quantities were clipped beyond 1% of the daily volume.

The moments $\Sigma_a^{(2n)}$ for n=1,2,3 are shown in Fig. 6.4 as a function of a. Remarkably, the theoretical predictions qualitatively reproduce the empirical data, in spite of the rather uncontrolled approximations made in the calculations. In particular, we do find that for large enough a, all these moments scale propor-

tionally to T, whereas super-linear behaviour in T is observed for small a, as a consequence of the long memory of order signs. Such behaviour is washed away when we look at large volumes only, i.e. when a is large enough.

Looking in particular at the curves for n=1, we see that $\widetilde{\mu}(a)$ decreases from $\widetilde{\mu}(a=0)\approx 3/2$ to $\widetilde{\mu}(a=a_c)\approx 2$ with $a_c\approx 1$ for large tick EUROSTOXX and $a_c\approx 2$ for smaller tick LLOY, TSCO and SPMINI. From Eq. (6.31), we deduce that $\lambda\sigma_\ell^2\approx 1/4$ for EUROSTOXX, and $\lambda\sigma_\ell^2\approx 1/8$ for LLOY, TSCO and SPMINI.

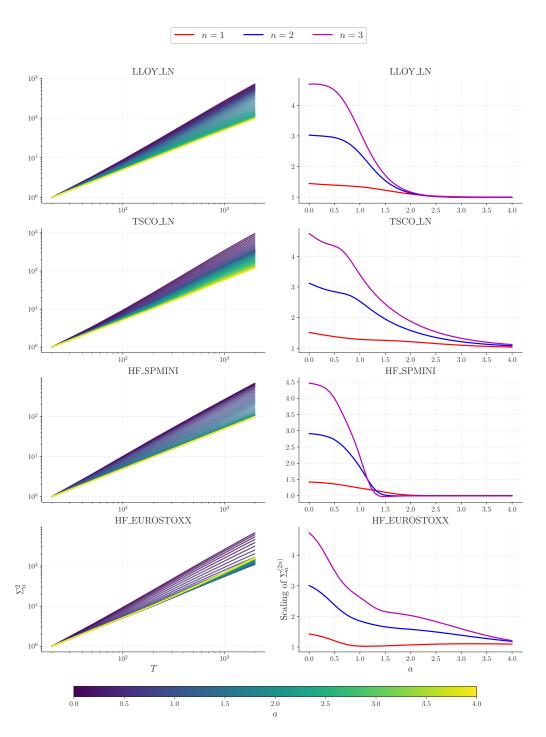


Figure 6.4: Data for LLOY, TSCO, SPMINI and EUROSTOXX. Right column: scaling of the different moments $\Sigma_a^{(2n)}$ as a function of trade time T, from which we extract the exponent from a regression of the data in log-log. Left column: scaling exponent as a function of a. As predicted by our model, increasing the value of a-1.05 giving more weight to orders with large volume – reduces the value of the exponent, which reaches unity for $a > a_c(n)$ with $a_c(1) \approx 2$ for LLOY, TSCO and SPMINI, and $a_c(1) \approx 1$ for EUROSTOXX.

6.4 The Impact-Diffusivity puzzle and a generalized propagator

We now discuss the problem of whether the volatility of price changes can be explained chiefly in terms of the impact of intertwined metaorders of different sizes and signs that get executed in the market. We will first recall how price diffusivity and long-memory of order flow are reconciled within the standard "propagator" framework, and then show how the issue becomes much more perplexing when the impact of metaorders obeys the square-root law. Three possible resolutions are proposed, together with their strengths and weaknesses.

6.4.1 Price diffusivity within the propagator model

We first assume the impact of single orders is given by a *deterministic* propagator model, i.e. the average price change due to an order of volume q executed at time t = 0 is $\theta(q)$ which then decays as [7, 31]

$$G(t) = \theta(q) \left(\frac{\tau_0}{t + \tau_0}\right)^{\beta}, \tag{6.41}$$

where τ_0 is the average time between two child orders. The average impact of a single metaorder of volume Q, duration $s \gg \tau_0$ and chosen participation rate $\widetilde{\varphi}$ (such that $Q = q\widetilde{\varphi}s$) is then given by²¹

$$\mathcal{I}(t \le s) = \theta(q) \int_0^t dt' \, \widetilde{\varphi} \left(\frac{\tau_0}{t - t'} \right)^{\beta} = \mathcal{I}_0(q, \widetilde{\varphi}) t^{1 - \beta}, \qquad \mathcal{I}_0(q, \widetilde{\varphi}) := \frac{\widetilde{\varphi} \theta(q) \tau_0^{\beta}}{1 - \beta}. \tag{6.42}$$

The peak impact $\mathcal{I}(s)$ reads

$$\mathcal{I}(s) = \frac{\widetilde{\varphi}\theta(q)s^{1-\beta}\tau_0^{\beta}}{1-\beta} = \frac{\theta(q)q^{\beta-1}(\widetilde{\varphi}\tau_0)^{\beta}}{1-\beta}Q^{1-\beta},\tag{6.43}$$

which reveals one problematic flaw of the propagator model: for $\beta=1/2$, peak impact does not only depend on Q, as empirically observed, but also on q and $\widetilde{\varphi}$. Although one can always choose $\theta(q) \propto q^{1-\beta}$ to get rid of the q dependence, one is still left with a square-root dependence on the participation rate $\widetilde{\varphi}$.

After the end of the metaorder, impact decays (see Fig. 6.5) and is given by

$$\mathcal{I}(t>s) = \theta(q) \int_0^s dt' \, \widetilde{\varphi} \left(\frac{\tau_0}{t-t'} \right)^{\beta} = \mathcal{I}_0(q, \widetilde{\varphi}) \left(t^{1-\beta} - (t-s)^{1-\beta} \right). \tag{6.44}$$

²¹Here we distinguish the participation rate of a specific metaorder, $\tilde{\varphi}$, from the average participation rate of the whole market, φ .

Note that this last expression behaves as $\mathcal{I}(s)(s/t)^{\beta}$ for $t \gg s$.

For simplicity, we assume for now that all market orders have the same volume q, and model the price variation Δ_T over time T as the superposition of the *average* impact of metaorders, neglecting fluctuations that will be considered in section 6.4.5 below.

Price variations are then given by the sum of two terms, $\Delta_{T,1}$ describing the impact of metaorders initiated within [0,T], and $\Delta_{T,2}$ the decaying impact of metaorders initiated before t=0:

$$\Delta_{T,1} = \mathcal{I}_0(q,\varphi) \int_0^T dN_t \,\varepsilon(t) \left[\mathbb{I}(t+s>T)(T-t)^{1-\beta} \right]$$
 (6.45)

$$+\mathbb{I}(t+s < T)\left((T-t)^{1-\beta} - (T-t-s)^{1-\beta}\right)$$
 (6.46)

and

$$\Delta_{T,2} = \mathcal{I}_{0}(q,\varphi) \int_{-\infty}^{0} dN_{t} \,\varepsilon(t) \left[\mathbb{I}(t+s>T) \left((T-t)^{1-\beta} - (-t)^{1-\beta} \right) + \mathbb{I}(t+s
(6.47)$$

The average of Δ_T over ε is of course nil, and its variance is given by two contributions:

$$\Sigma_{T,1}^{2} = \mathcal{I}_{0}^{2}(q,\varphi)\nu \int_{0}^{T} du \int_{0}^{\infty} ds \, \Psi(s) \left[\mathbb{I}(s>u)u^{2(1-\beta)} + \mathbb{I}(s< u) \left(u^{1-\beta} - (u-s)^{1-\beta} \right)^{2} \right]$$
(6.48)

and

$$\Sigma_{T,2}^{2} = \mathcal{I}_{0}^{2}(q,\varphi)\nu \int_{T}^{\infty} du \int_{0}^{\infty} ds \, \Psi(s) \left[\mathbb{I}(s>u) \left(u^{1-\beta} - (u-T)^{1-\beta} \right)^{2} + \mathbb{I}(s
(6.49)$$

All these contributions can be exactly computed for large T when $\Psi(s)$ decays as a power-law $s^{-(1+\mu)}$, but provided $\mu < 2$ the scaling can simply be obtained by the change of variables s = xT, u = yT, that yields

$$\Sigma_T^2 := \mathbb{E}[\Delta_T^2] \propto \mathcal{I}_0^2(q, \varphi) \, T^{3 - 2\beta - \mu}. \tag{6.50}$$

Hence, we see that metaorders contribute to volatility provided $3 - 2\beta - \mu = 1$. This equality coincides, as expected, with the critical condition $2\beta = 1 - \gamma$ derived within the propagator model (recall that $\gamma = \mu - 1$). When $\beta < (1 - \gamma)/2$, the price is super-diffusive (i.e. $\gg T$), whereas when $\beta > (1 - \gamma)/2$, the contribution of the average impact of metaorders to price variance is negligible, i.e. o(T).

In order to recover the square-root impact law, one should naively set $\beta=1/2$, such that $3-2\beta-\mu=1-\gamma$. But the contribution of metaorders to volatility Σ_T^2 would be then subdominant at long times, unless $\gamma\to 0$ (i.e. an hyper-slow decay of the sign autocorrelation function). Note that the choice $\beta=1/2^-, \gamma=0^+$ corresponds to the model advocated by Jusselin & Rosenbaum [51], but is difficult to reconcile with the empirically determined value $\gamma\approx 0.5$ [16]. This value of γ , in turn, means that a square-root impact appears to be unable to generate price diffusion, since in this case $3-2\beta-\mu\approx 0.5<1$.

One could then argue that volatility does not primarily come from the average impact of metaorders, but rather from its fluctuations, a possibility that we explore in section 6.4.4 below. But in any case, the propagator model with $\beta=1/2$ fails to account for two important stylized facts:

- Injecting $\beta = 1/2$ into Eq. (6.42), one finds, as already mentioned above, $\mathcal{I}(s) \propto \sqrt{\widetilde{\varphi}\tau_0}\sqrt{Q}$ when $\theta(q) = \sqrt{q}$ (as indeed suggested by the data of [2]). Hence, one recovers the square-root law $\mathcal{I}(Q) = Y\sqrt{Q}$ but with an extra square-root dependence of the prefactor Y on the participation rate $\widetilde{\varphi}$, when empirical data show that Y is all but independent of $\widetilde{\varphi}$.
- The decay of impact after the end of a metaorder, when fitted with Eq. (6.44), suggests a value $\beta \approx 0.2 < 1/2$ [2, 44, 93], i.e. a much faster short time decay and a much slower long time decay than predicted by $\beta = 1/2$.

We conclude that the propagator model, even with $\beta=1/2$ cannot fully explain the observed impact of metaorders, nor its post-execution decay. In the following sections, we explore different routes to reconcile metaorder impact with long-term volatility.

6.4.2 A generalized propagator model

As we just discussed, the square-root price profile *during* the execution of the metaorder and the subsequent impact decay cannot be captured within the standard propagator model. Here we propose a (somewhat ad-hoc) extension of this model that allows one to decouple these two profiles. We will not attempt to fully justify such a proposal from first principles, but use the resulting equations as a convenient way to capture the known phenomenology of metaorder impact.

Let us assume that once a metaorder has started, the market slowly adapts and, in the spirit of the LLOB model [7, 74], progressively provides more liquidity

to absorb the incoming flow. We represent this effect as a *two-time* propagator, describing the impact of a child order occurring at time t' after the start of the metaorder on the price at time t > t':

$$G(t' \to t) = \frac{\theta(q)}{(\widetilde{\varphi}t' + n_0)^{1/2 - \beta}} \left(\frac{\tau_0}{t - t' + \tau_0}\right)^{\beta}, \qquad (\beta < \frac{1}{2})$$
 (6.51)

where τ_0 is the average time between two trades and $\widetilde{\varphi} \times t'$ is *number* of child orders executed since the start of the metaorder, which have eaten into the LLOB and therefore reveal more hidden liquidity. n_0 is the number of trades after which the metaorder is statistically detected by liquidity providers. The immediate impact of a child order is thus

$$G(t' \to t') = \frac{\theta(q)}{(\widetilde{\varphi}t' + n_0)^{1/2 - \beta}}$$

$$(6.52)$$

which decreases with t', as liquidity adapts.

The impact of a metaorder of duration $s \gg \tau_0$ and $\widetilde{\varphi} s \gg n_0$ is now given by ²²

$$\mathcal{I}(t \le s) \approx \theta(q) \int_0^t dt' \frac{\widetilde{\varphi}^{1/(2+\beta)}}{t'^{1/(2-\beta)}} \left(\frac{\tau_0}{t-t'}\right)^{\beta} = \mathcal{I}_1(q, \widetilde{\varphi}) \sqrt{t}. \tag{6.53}$$

with

$$\mathcal{I}_1(q,\widetilde{\varphi}) := \mathcal{B}_\beta \sqrt{\widetilde{\varphi}} \theta(q) (\widetilde{\varphi} \tau_0)^\beta, \tag{6.54}$$

where $\mathcal{B}_{\beta} = 2\Gamma(1/2 + \beta)\Gamma(1 - \beta)/\sqrt{\pi}$.

After the end of the metaorder, impact now decays as

$$\mathcal{I}(t>s) = \mathcal{I}_1(q)\sqrt{s} \left[\left(\frac{t}{s}\right)^{1-\beta} - \left(\frac{t}{s} - 1\right)^{1-\beta} \right], \tag{6.55}$$

which reproduces the empirical decay of metaorders if one chooses $\beta \approx 0.2$. With $\theta(q) \propto \sqrt{q}$, the peak impact then reads

$$\mathcal{I}(Q) \propto (\widetilde{\varphi}\tau_0)^{\beta} \sqrt{Q},$$
 (6.56)

which still depends on the participation rate, but now with a smaller exponent $\beta = 0.2$, more difficult to exclude empirically.

²²Note that when $s \gg \tau_0$ but $\widetilde{\varphi}s \lesssim n_0$, impact behaves as in the standard propagator model as $t^{1-\beta}$. If $\beta \approx 0.2$, such a behaviour is much less concave than a square-root, in agreement with the results of [2, 79].

One could have hoped that the slower relaxation of impact in the post-execution regime would help recover a linear-in-T behaviour of $\Sigma_T^2 := \mathbb{E}[\Delta_T^2]$. Unfortunately, one finds that the contribution of metaorder impact to Σ_T^2 scales as

$$\Sigma_T^2 \propto_{T \to \infty} \begin{cases} T^{1-\gamma}, & \gamma < 2\beta; \\ T^{1-2\beta}, & \gamma > 2\beta, \end{cases}$$
 (6.57)

which is again sub-diffusive whenever $\beta > 0$ and $\gamma = \mu - 1 > 0$. In other words, price diffusion is only possible when impact is permanent, i.e. $\beta = 0$. This is in fact the assumption made by Sato & Kanasawa in their latest paper [101]. However, all empirical data known to us suggest that impact decays, at least over short to medium time scales [3, 44, 80, 83], which according to our calculation should lead to substantial price mean reversion on such time scales.

We now turn to three possible resolutions of the diffusion "paradox": one based on the autocorrelation of the sign of metaorders, a second one based on the permanent impact of large child orders, and the last one based on permanent impact fluctuations.

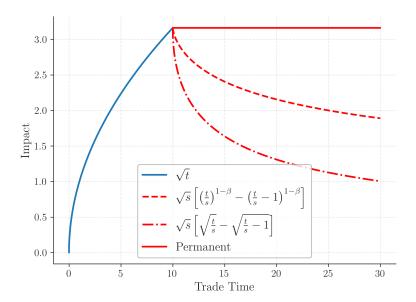


Figure 6.5: Comparison of the three main theories of impact decay for a metaorder of size s=10. During the execution phase (blue), impact follows a square-root impact growth, i.e. $\propto \sqrt{t}$, where t denotes the index of the child order [2, 83]. In the post-execution phase (red), impact may either remain permanent (solid line) or decay (dashed). Two distinct decay mechanisms are illustrated: one based on the generalized two-time propagator as $t^{-\beta}$ with $\beta=0.2$, and the other on the LLOB theory, corresponding to $\beta=\frac{1}{2}$.

6.4.3 The role of metaorder autocorrelations

What happens if we assume, as in section 6.3.3, that metaorders themselves are autocorrelated, with a new exponent γ_{\times} ? Extending the calculation of $\mathbb{E}[\Delta_T^2]$ the case $\gamma > \beta - 1/2$ (always satisfied when $\beta < \frac{1}{2}$), we find that these correlations contribute to volatility as

$$\mathbb{E}[\Delta_T^2]_{o.d.} \propto \Gamma T^{2-\gamma_{\times}-2\beta}. \tag{6.58}$$

This leads to diffusion provided $\gamma_{\times} = 1 - 2\beta$, which is, not surprisingly, the same condition as for the standard propagator model but with γ replaced by γ_{\times} . Indeed, after coarse graining metaorders into effective single orders, we are back to the usual propagator model.

For $\beta=0.2$, this gives $\gamma_{\times}=0.6$, which is close to γ itself and close to its empirical value [69] and G. Maitrier (unpublished). It is also compatible with the bound obtained in Eq. (6.38), which ensures that the autocorrelation of signs is dominated by the size distribution of metaorders, as postulated by [18] and firmly established in [16].

The resulting value of volatility σ^2 , defined as Σ_T^2/T , is given by

$$\sigma^2 = C(\beta, \gamma_{\times}) \Gamma \nu^2 \tau_0^{\gamma_{\times}} \mathcal{I}_1^2(q, \varphi) \left(\overline{s^{\frac{1}{2} + \beta}} \right)^2, \tag{6.59}$$

where $C(\beta, \gamma_{\times})$ is a numerical coefficient found to be ≈ 5.6 for $\gamma_{\times} = 0.6$ and $\beta = 0.2$.

With $\theta(q) = \theta_0 \sqrt{q}$ and $\phi = \nu q \varphi \bar{s}$, we finally find

$$\theta_0 = Y \frac{\sigma}{\sqrt{\phi}}, \qquad Y \propto \frac{\bar{n}^{\frac{1}{2} - \beta}}{\sqrt{C\Gamma}},$$
(6.60)

where we have assumed that $(s^{\frac{1}{2}+\beta})^2 \propto \bar{s}^{1+2\beta}$, which is justified in the present case since the $(1/2+\beta)$ -th moment of s converges (whenever $\mu=1+\gamma>\frac{1}{2}+\beta$). We recall that $\bar{n}=\varphi\bar{s}$ is the average number of child orders per metaorder.

This result is interesting since the peak impact is then precisely given by the standard square-root law, up to a weak participation rate dependence:

$$\mathcal{I}(Q) \propto (\widetilde{\varphi}\tau_0)^{\beta} \, \sigma \sqrt{\frac{Q}{\phi}}.$$
 (6.61)

Note that we find naturally that impact is proportional to volatility, simply because volatility is due to impact!

The above calculation can be generalized to higher moments of Δ_T . Assuming Wick-like factorisation of $\mathbb{E}[\varepsilon(t_1)\cdots\varepsilon(t_{2n})]$ as

$$\mathbb{E}[\varepsilon(t_1)\cdots\varepsilon(t_{2n})] \propto \underbrace{\prod_{i\neq j} |t_i - t_j|^{-\gamma_{\times}}}_{n \text{ pairs}},$$
(6.62)

it is easy to show that the 2n-th moment of Δ_T scales as:

$$\Sigma_T^{(2n)} \propto T^{2n(1-\beta)-n\gamma_{\times}} = T^n, \tag{6.63}$$

as indeed observed empirically (up to subleading corrections that can also be rationalized within our framework), see Fig. 6.6.

Note that we do not observe mutifractality (i.e. $\Sigma_T^{(2n)} \propto T^{\zeta_n}$ with $\zeta_n \neq n$) because we work in trade time and not in real time. As it is well known, multifractal effects come from intermittent fluctuations of the activity rate ν , see e.g. [51, 113]. An interesting extension of our model, which we leave for future work, would be to assume that ν itself has fractal properties.

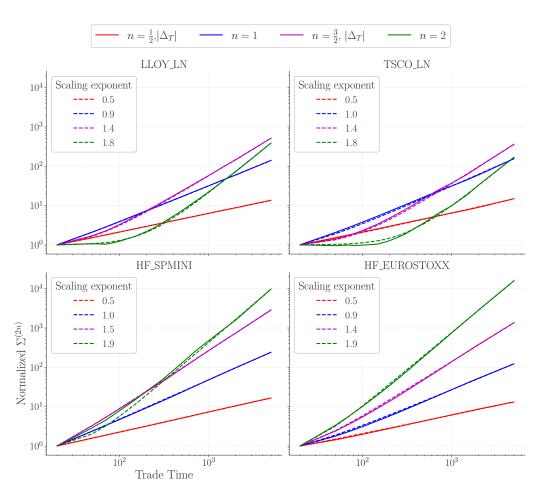


Figure 6.6: Scaling of the moments of price changes $|\Delta_T|^{2n}$ as a function of trade time T. We normalized the moment values such that all curves begin at 1 for T=1. As shown, a sublinear behavior is found at short times for n>1, likely attributable to price jumps that may dominate for short time scales when n increases. Therefore, we fitted the data as $\Sigma^{(2n)} = a_0 + a_1 T^{\zeta_n}$ and present the values of ζ_n in the legend. We find in all cases $\zeta_n \approx n$, up to subleading corrections that we attribute to volume fluctuations.

6.4.4 The role of volume fluctuations

Another possibility is to take seriously the fact that the size q of child orders is fluctuating and correlated with the duration s of metaorders. As we have shown in section 6.3.2, a dependence of the exponent μ_q on q explains how different moments of the volume imbalance depend on T, see Eqs. (6.32), (7.4). It is also in agreement with direct empirical results, see Fig. 6.3.

Now, similarly to μ_q , we might assume that the impact decay exponent β becomes q-dependent: $\beta_q = \beta_1 - \lambda' \log q$. One then observes that price impact becomes permanent for large child orders, with $q > q_0$ such that $\beta_q \leq 0$. When metaorder signs are independent, the non-vanishing contribution to volatility then reads:

$$\mathbb{E}[\Delta_T^2]_{q>q_0} = \nu \int_{q_0}^{\infty} \mathrm{d}q \,\Xi(q) \mathcal{I}_1^2(q,\varphi) \int_0^T \mathrm{d}u \, \int_0^u \mathrm{d}s \,\Psi_q(s) s \tag{6.64}$$

$$\approx_{T\to\infty} \nu T \int_{q_0}^{\infty} \mathrm{d}q \,\Xi(q) \mathcal{I}_1^2(q,\varphi) \bar{s}_q$$
 (6.65)

Taking $\theta(q) = \theta_0 \sqrt{q}$ one gets $\mathcal{I}_1(q > q_0, \varphi) = 2\theta_0 \sqrt{\varphi q}$, so the long-term price volatility is given by

$$\sigma^2 = \frac{\mathbb{E}[\Delta_T^2]_{q > q_0}}{T} \propto \theta_0^2 \phi_0, \tag{6.66}$$

where ϕ_0 is the average volume flow of the market, restricted to "large" child orders $q > q_0$. Interestingly, this relation can be read backwards as

$$\theta_0 \propto \frac{\sigma}{\sqrt{\phi_0}},$$
 (6.67)

allowing one to recover the full square-root impact law from the expression of \mathcal{I}_1 above:

$$\mathcal{I}(Q|q,\widetilde{\varphi}) = Y(q,\widetilde{\varphi})\sigma\sqrt{\frac{Q}{\phi}}, \qquad Y(q,\widetilde{\varphi}) \propto (\widetilde{\varphi}\tau_0)^{\beta_q}\sqrt{\frac{\phi}{\phi_0}}.$$
 (6.68)

We see that this result suggests a weak dependence of the prefactor of the square-root law in q and $\widetilde{\varphi}$, which disappears for large enough $q>q_0$, since $\beta_q\to 0$ in that case. However, in that case Y would be substantially larger than empirically observed, since from Fig. 6.3 we estimate $\phi_0\approx 0.1\phi$. Besides, since child orders of size $< q_0$ (which are the most numerous) only impact prices temporarily, this scenario would lead to strong visible mean-reversion in prices not observed in data – see e.g. Fig. 6.6 for n=1.

It is, however, important to discuss how volume fluctuations might affect the result Eqs. (6.57), (6.58) above, induced by the correlation between different metaorders. It is plain to extend the calculations of section 3.3 to get

$$\mathbb{E}[\Delta_T^2]_{o.d.} \propto \Gamma e^{m + \frac{\sigma_\ell^2}{4}} T^{2 - \gamma_{\times} - 2\beta_m + \lambda' \sigma_\ell^2}, \tag{6.69}$$

²³A sufficient condition for price diffusivity in the standard propagator model is $\beta_1 = 1 - \mu_1/2$ and $\lambda' = \lambda/2$, but we will not impose such constraints below and leave λ and λ' free.

whereas a similar calculation for the diagonal contribution gives

$$\mathbb{E}[\Delta_T^2]_{d.} \propto e^{m + \frac{\sigma_\ell^2}{2}} T^{1 - 2\beta_m + 2\lambda' \sigma_\ell^2}. \tag{6.70}$$

The ratio of the prefactors is now only $e^{\sigma_\ell^2/4}/\Gamma$ in favor of the diagonal term. This means that the off-diagonal contribution, with a larger power of T, becomes dominant beyond a reasonable small value of T when $\sigma_\ell^2 = 1$ and $\Gamma = O(1)$. We will therefore choose in the following $\gamma_{\times} = 0.6$ and $\widehat{\beta}(0) := \beta_m - \frac{1}{2}\lambda'\sigma_\ell^2 = 0.2$, such that the off-diagonal contribution is exactly diffusive.

6.4.5 The role of impact fluctuations

Up to now, we have assumed *deterministic* impact and neglected the role of price changes induced by "news" or other order book events that are not related to trades, with the ambition of recovering all the price volatility from the impact of metaorders. However, it is clear that:

- (i) news events do obviously exist (see e.g. [56] for a recent discussion) and should indeed contribute to volatility. In fact, Efficient Market Theory predicts that the *only* contribution to volatility comes from news!
- (ii) the impact of a given metaorder has no reason to be deterministic: it should depend on specific time-dependent market conditions and thus include a random component.

Such a random component was in fact indirectly observed by Bucci et al. [70], where it was found that the effect of a single metaorder on price changes reads

$$\Delta_T = p_T - p_0 \approx \varepsilon \mathcal{I}(Q) \left[1 + z \eta \right], \tag{6.71}$$

where η is a zero mean, unit variance, independent random variable, and z a coefficient measuring the relative fluctuations of impact, found to be around 3 in [70].²⁴

Let us postulate that, while the average impact decays to zero with time as per the propagator model, the random component $z\eta$ does not, or at least not completely. This assumption does not violate any known stylized facts about the average decay of impact. Following these ideas, we expect price changes Δ_T to include extra terms that read

$$\Delta_{T,1} = z_{\infty} \theta_0 \sqrt{q\varphi} \int_0^T \mathrm{d}N_t \, \varepsilon(t) \left[\mathbb{I}(t+s > T) \sqrt{T-t} + \mathbb{I}(t+s < T) \sqrt{s} \right] \eta_t + \sigma_F \xi \sqrt{T}, \tag{6.72}$$

 $^{^{24}}$ Note that there is an error in that paper, where z, called a there, was reported to be around 0.1.

where $z_{\infty} \leq z$ accounts for a possible time decay of the random component of impact, and the last contribution captures fundamental "news", with a volatility σ_F . (ξ is another zero mean, unit variance, independent random variable).

This gives rise the following contribution to volatility:²⁵

$$\Sigma_T^2 = \mathbb{E}[\Delta_T^2]_{\eta,\xi} = z_\infty^2 \theta_0^2 q \varphi \nu \int_0^T du \int_0^\infty ds \Psi(s) \left[\mathbb{I}(s > u) u + \mathbb{I}(s < u) s \right] + \sigma_F^2 T,$$
(6.73)

whence a long-term volatility given by

$$\sigma^2 = \sigma_F^2 + z_\infty^2 \theta_0^2 q \varphi \nu \bar{s} \equiv \sigma_F^2 + z_\infty^2 \theta_0^2 \phi, \tag{6.74}$$

where $\phi = q\varphi\nu\bar{s}$ is, again, the average volume flow in the market.

This expression is quite interesting: inserting $\theta_0 = Y\sigma/\sqrt{\phi}$, we get, provided $Yz_{\infty} < 1$:²⁶

$$\sigma^2 = \frac{\sigma_F^2}{1 - Y^2 z_\infty^2},\tag{6.75}$$

i.e. excess volatility induced by trading, independently of its information content. This is in line with the empirical results of [114] (section 4), [33] and [7] (Figs. 13.2 and 14.5), and is of course related to the well-known excess volatility/excess trading puzzle, see e.g. [57, 108, 109] and [7], chapters 2 & 20, see also [52, 87, 106] for related discussions.

Hence, even if the deterministic, decaying part of impact does not contribute to long-term volatility, its fluctuations might do the job. Of course, this somewhat contorted scenario relies on the assumption that the random component of impact has a permanent contribution to price changes, i.e. $z_{\infty} > 0$. Although this hypothesis is somewhat ad-hoc, the non-trivial result here is the relation between price impact and volatility given by

$$\mathcal{I}(Q) \propto \sqrt{\frac{\sigma^2 - \sigma_F^2}{z_\infty^2}} \times \sqrt{\frac{Q}{\phi}}.$$
 (6.76)

If, on the other hand, the permanent contribution z_{∞} vanishes, then trivially $\sigma^2 = \sigma_F^2$. This is the Efficient Market picture, where uninformed trades do not

²⁵Note that since $\mathbb{E}[\varepsilon\eta]$ is assumed to be zero, there is no particular role for metaorder correlations in this scenario. One could however wonder how the result given in Eq. (6.57) might change if we assume "informed metaorders", i.e. some correlations between the sign of the metaorder ε and the subsequent fondamental price change ξ , see section 6.5.5. The result is a contribution to Σ_T^2 proportional to $\rho T^{1-\beta}$, which is subdominant at large T.

 $[\]Sigma_T^2$ proportional to $\rho T^{1-\beta}$, which is subdominant at large T.

²⁶Note that the same expression would hold with $z_{\infty}=1$ and Y given by Eq. (6.60) if on top of the off-diagonal contribution to volatility one would add a fundamental contribution. The following discussion can thus be transposed to that case as well.

contribute to long-term volatility and our whole construction breaks down. However, this would leave θ_0 undetermined and would not allow rationalizing why θ_0 is proportional to $\sigma/\sqrt{\phi}$, as found empirically. Furthermore, the impact of uninformed metaorders (which is probably a large fraction of all metaorders) would generate large price reversion effects, which again are not observed.

6.4.6 Discussion

We thus have three possible scenarios for generating long term volatility from the impact of metaorders. The idea that metaorders are correlated between one another, with roughly the same long memory as within each one of them, seems plausible, is compatible with available data and provides a natural extension of the propagator scenario: diffusive prices emerge from the subtle interplay between decaying impact and autocorrelated order flow. Such a scenario predicts a weak dependence of the prefactor Y of the square-root law with the participation rate of the metaorder, as $\widetilde{\varphi}^{\beta}$ with $\beta \approx 0.2$.

The second scenario, which attributes long term to the non-decaying impact of metaorders executed with large child orders does not seem very credible to us, because it would lead to substantial mean-reversion effects due to the impact decay of small child orders, which is at odds with empirical decay.

Finally, long term volatility may result from the random component of impact, assumed to be permanent. This bypasses the paradox that average impact decays and, in the absence of correlations between metaorders, should lead to subdiffusion. In such a scenario, volatility is induced by trading activity alone, even if average impact was zero. Still, the fact that impact fluctuations are proportional to average impact, as postulated in Eq. (6.71), is important to recover the correct relation between peak impact $\mathcal{I}(Q)$ and $\sigma\sqrt{Q/\phi}$, Eq. (6.76), which is now independent of $\widetilde{\varphi}$.

We now turn to the study of the covariance of price changes and volume imbalances, to see whether we can constrain the theory further. We will indeed see that the data strongly favors an interpretation based on the average impact of correlated metaorders. The last scenario, where volatility arises from the fluctuating part of impact, does not pass the test.

6.5 Covariance between order flow imbalance and prices changes

Another quantity that can be computed within our model and easily measured empirically using the public tape of trades and prices is the "aggregated" impact

 $\mathbb{E}[\Delta|I^a]$, conditioned to a certain value of imbalance I^a , as studied in [37, 115]. We know that such a quantity behaves very differently from the square-root law, and has non-trivial scaling properties as a function of T, see [37] for a = 0 and a = 1. For a = 0, in particular, one finds that the initial slope of $\mathbb{E}[\Delta|I^0]$ as a function of I^0 scales like $T^{-\omega}$ with $\omega \approx 1/4$ [7, 37], a result we confirm in Fig. 6.7.

Note that if I^a and Δ were Gaussian variables one could use the following general relation to predict that slope:

$$\mathbb{E}[\Delta|I^a] = \frac{\mathbb{E}[\Delta \cdot I^a]}{\Sigma_a^2} I^a, \tag{6.77}$$

i.e. a linear aggregate impact for small imbalances, where Σ_a^2 was defined in section 6.3.2. However, in our model the Gaussian assumption does not hold since I^a is a truncated Lévy variable (see section 6.3.2). Hence the exact calculation of $\mathbb{E}[\Delta|I^a]$ is much more intricate, and we restrict the following analysis to the covariance $\mathbb{E}[\Delta \cdot I^a]$, which we compare to empirical data below. Still, naively applying Eq. (7.8) will predict a scaling in $T^{-\omega}$, albeit within an uncontrolled approximation.

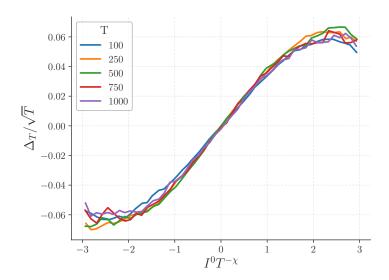


Figure 6.7: Aggregated impact versus sign imbalance for the EUROSTOXX. After appropriate rescaling, curves for different T nicely collapse onto a single master curve. The y-axis rescaling reflects the diffusive nature of the price, while the x-axis rescaling captures trade sign correlations, as detailed in section 6.3. In this case, we find $\omega = \chi - \frac{1}{2} \approx 0.22$.

6.5.1 Without volume fluctuations

Since impact fluctuations, considered in section 6.4.5, are independent of the order imbalance (i.e. $\mathbb{E}[\varepsilon\eta] = 0$), we can restrict the calculation to the deterministic part of impact, which is given by the generalized propagator model above. Neglecting any sign correlations between different metaorders, we thus get, when all child orders have the same size:

$$\mathbb{E}[\Delta_T \cdot I_T^a] = \nu \varphi q^a \mathcal{I}_1(q, \varphi) \int_0^T du \int_0^\infty ds \Psi(s) \left[\mathbb{I}(s > u) u \sqrt{u} \right]$$
 (6.78)

$$+\mathbb{I}(s < u)s\sqrt{s}\left(\left(\frac{u}{s}\right)^{1-\beta} - \left(\frac{u}{s} - 1\right)^{1-\beta}\right)\right].$$
 (6.79)

This quantity scales like $T^{1-\beta}$ whenever $\mu > 3/2 + \beta$ and like $T^{5/2-\mu}$ otherwise.²⁷ Note that if one disregards the non-Gaussian nature of I^0 and uses Eq. (7.8) to compute the initial slope of the aggregate impact, one finds (using Eq. (6.17)) $\omega = 1/2$ when $\mu = 3/2$ and $\beta > 0$, which is far from the empirical value $\omega \approx 1/4$. As we shall see in the following, empirical results suggest a strong dependence on a, meaning that volume fluctuations certainly cannot be neglected in this calculation.

6.5.2 With correlated metaorders

One can extend the result above to account for correlated metaorders. Assuming $\gamma > 0$ and $\beta < 1/2$, the "non-diagonal" contribution $t \neq t'$ gives a contribution that scales as

$$\mathbb{E}[\Delta_T \cdot I_T^a]_{o.d.} \propto \Gamma \, q^{a + \frac{1}{2}} \, T^{2 - \beta - \gamma_{\times}},\tag{6.80}$$

to be compared with the above results, i.e. $T^{1-\beta}$ or $T^{5/2-\mu}$. When $\gamma_{\times} < 1$, the non-diagonal contribution is always dominant at large T compared to $T^{1-\beta}$ and becomes subdominant compared to $T^{5/2-\mu} = T^{3/2-\gamma}$ when $\beta + \gamma_{\times} > 1/2 + \gamma$.

For $\gamma = \gamma_{\times}$ and $\beta < \frac{1}{2}$ we thus deduce that the "off-diagonal" contribution is always dominant for large enough times. The corresponding value of the slope exponent ω is now equal to $\beta + \gamma_{\times} - \gamma$, and is thus very close to the empirical value $\omega \approx 1/4$ when $\gamma_{\times} \approx \gamma$ and $\beta \approx 0.2$.

6.5.3 With correlated metaorders and volume fluctuations

Let us now add volume fluctuations, with a q-dependent value of μ_q as given by Eq. (6.22) and $\beta_q = \beta_1 - \lambda' \ell$, with $\ell = \log q$. Now, there exists a value $q = q_c'$ such that $5/2 - \mu_{q_c'} = 1 - \beta_{q_c'}$.

Note that for the standard propagator model, one finds that the scaling is always $T^{3-\beta-\mu}$ when $\mu > 1$.

Then, using $\mathcal{I}_1(q,\varphi) \propto \sqrt{q}$, one gets, for the diagonal contribution

$$\mathbb{E}_{q}[\Delta_{T} \cdot I_{T}^{a}]_{d.} \propto \int_{0}^{q_{c}'} dq \,\Xi(q) \, q^{a+\frac{1}{2}} T^{5/2-\mu_{q}} + \int_{q_{c}'}^{\infty} dq \,\Xi(q) \, q^{a+\frac{1}{2}} T^{1-\beta_{q}}, \qquad (6.81)$$

and for the off-diagonal contribution

$$\mathbb{E}_{q}[\Delta_{T} \cdot I_{T}^{a}]_{o.d.} \propto \Gamma \,\mathbb{E}[q^{a}] \int_{0}^{\infty} \mathrm{d}q \,\Xi(q) \,q^{\frac{1}{2}} T^{2-\gamma_{\times}-\beta_{q}}, \tag{6.82}$$

Repeating the same calculations as in section 6.3.2, we now get:

$$\mathbb{E}_{q}[\Delta_{T} \cdot I_{T}^{a}]_{d.} \propto T^{5/2 - \mu_{m}} \int_{0}^{\ell_{c}'} d\ell \ e^{(\ell - m)(a + \frac{1}{2} - \lambda \log T) - (\ell - m)^{2}/2\sigma_{\ell}^{2}}
+ T^{1 - \beta_{m}} \int_{\ell_{c}'}^{\infty} d\ell \ e^{(\ell - m)(a + \frac{1}{2} + \lambda' \log T) - (\ell - m)^{2}/2\sigma_{\ell}^{2}}.$$
(6.83)

When $(\lambda, \lambda') \log T \ll 1$, the Gaussian integrals are dominated by the region around $\ell^* = m + \sigma_\ell^2 (a+1/2)$. So, schematically, when $\ell^* < \ell_c' := \log q_c'$ the first integral dominates, while for $\ell^* > \ell_c'$ the second integral dominates. Hence, the dominant term scales as:

$$\mathbb{E}_{q}[\Delta_{T} \cdot I_{T}^{a}]_{d.} \propto e^{m(a+\frac{1}{2})+\frac{1}{2}\sigma_{\ell}^{2}(a+\frac{1}{2})^{2}} \begin{cases} T^{5/2-\widehat{\mu}(a)}, & \widehat{\mu}(a) = \mu_{m} + (a+\frac{1}{2})\lambda\sigma_{\ell}^{2} & a < a'_{c}; \\ T^{1-\widehat{\beta}(a)}, & \widehat{\beta}(a) = \beta_{m} - \left(a+\frac{1}{2}\right)\lambda'\sigma_{\ell}^{2} & a > a'_{c}; \end{cases}$$

$$(6.84)$$

where a'_c is such that $\widehat{\mu}(a'_c) = \mu_{q'_c}$.

The power of T coming from the diagonal contribution is thus predicted to decrease with a when $a < a'_c$ and then to increase with a for larger a. Hence we expect an interesting non-monotonic behaviour of the effective exponent as a function of a, i.e. with the relative weight given to child orders with large volumes.

The off-diagonal contribution, on the other hand, gives:

$$\mathbb{E}_{q}[\Delta_{T} \cdot I_{T}^{a}]_{o.d.} \propto \Gamma e^{m(a+\frac{1}{2})+\frac{1}{2}\sigma_{\ell}^{2}(a^{2}+\frac{1}{4})} T^{2-\gamma_{\times}-\widehat{\beta}(0)}. \tag{6.85}$$

Note that the coefficient in front of the power-law is $\Gamma e^{-\sigma_\ell^2 a/2}$ smaller than the one corresponding to the diagonal contribution. Other numerical prefactors may however contribute as well, that were neglected in the rough estimate of the above integrals. The final theoretical prediction is that $\mathbb{E}_q[\Delta_T \cdot I_T^a]$ is the sum of three power-law contributions:

• $T^{5/2-\widehat{\mu}(a)}$, with an exponent equal to 1 for a=0 and decreasing as a increases,

- $T^{1-\widehat{\beta}(a)}$, with an exponent equal to ≈ 0.8 for a=0 and increasing as a increases,
- $T^{2-\gamma_{\times}-\widehat{\beta}(0)}$, with an exponent independent of a and equal to ≈ 1.2 for the default values of γ_{\times} , $\widehat{\beta}(0)$.

In section 6.5.7 below, we will show that empirical data can be fitted as a power-law of T, with an effective exponent that indeed behaves non-monotonically with a, which suggests $(\lambda, \lambda')\sigma_{\ell}^2$ in the range 0.1–0.2, also in line with the condition already obtained in section 6.3.2.

Finally, note that the slope exponent ω , naively predicted from Eq. (7.8) for a=0 is

$$\omega_{d.} = \frac{1}{2} + \widehat{\mu}(0) - \widetilde{\mu}(0) = \frac{1}{2}(1 + \lambda \sigma_{\ell}^2), \qquad \omega_{o.d.} = 1 - \widetilde{\mu}(0) + \gamma_{\times} + \widehat{\beta}(0), \quad (6.86)$$

depending on whether the diagonal or off-diagonal contribution dominates. Numerically, with $\lambda \sigma_{\ell}^2 = 1/8$ and $\gamma_{\times} + \widehat{\beta}(0) = 0.8$, one finds $\omega_{d.} \approx 0.56$ and $\omega_{n.d.} \approx 0.30$, close to the empirical value 1/4 in the second case.

6.5.4 With a random impact component

Quite interestingly, if the random component of impact (see Eq. (6.71)) is assumed to be such that $\mathbb{E}[\varepsilon\eta] = 0$, it will not contribute to the covariance between price changes and order imbalance. Indeed, by definition such a contribution to price changes does *not* contribute to the covariance between price changes and order imbalance. In such a scenario, only the fundamental component can contribute to the covariance, provided some metaorders are "informed", as we show next.

6.5.5 With "informed" metaorders

Up to now, we assumed that there are no correlations between the sign of metaorders ε and the "Fundamental" component of price changes on time scale T, $\sigma_F \xi \sqrt{T}$, see Eq. (6.71). If we rather assume that $\mathbb{E}[\varepsilon \xi] = \rho/\sqrt{\nu T}$, where ρ measures the average amount of information of individual metaorders, ²⁸ we find an extra contribution to $\mathbb{E}[\Delta_T \cdot I_T^a]$ that reads:

$$\mathbb{E}[\Delta_T \cdot I_T^a] = \rho \sqrt{\nu} \varphi q^a \sigma_F \int_0^T du \int_0^\infty ds \Psi(s) \left[\mathbb{I}(s > u) u + \mathbb{I}(s < u) s \right]. \tag{6.87}$$

²⁸For a detailed discussion of the scaling with T, see [87], section 16.1.3. The idea is that out of $N_T = \nu T$ metaorders of random signs, an excess fraction $\sim \sqrt{\nu T}$ is possibly informed. It is important to stress that, in the spirit of the Kyle model [75], the fundamental component $\sigma_F \xi \sqrt{T}$ is not a mechanical consequence of impact.

Extending the calculation above, we now find that this covariance scales as $T^{2-\mu}$ for $\mu < 1$ and as T for $\mu > 1$, which is the case we focus on here. Note that the naive prediction for the slope exponent ω (from Eq. (7.8)) is $\omega = 2 - \mu$.

Hence, we find that such a fundamental contribution predicts a linear scaling of $\mathbb{E}[\Delta_T \cdot I_T^a]$ as a function of T, independently of a.

6.5.6 The correlation coefficient

Finally, we turn our attention with a natural description of the interplay between (generalized) volume imbalance and price changes, namely the following correlation coefficient

$$R_a(T) := \frac{\mathbb{E}[\Delta_T \cdot I_T^a]}{\Sigma_T \Sigma_a}, \qquad \Sigma_T := \sqrt{\mathbb{E}[\Delta_T^2]}, \qquad \Sigma_a := \sqrt{\mathbb{E}[I_T^{a2}]}$$
 (6.88)

In order to simplify the discussion, we assume that the crossovers between the two regimes for Σ_a^2 (Eq. (6.32)) and for $\mathbb{E}[\Delta_T \cdot I_T^a]$ (Eq. (7.5)) occur for the same value of $a = a_c = a'_c$. Although this is not precisely true, the following conclusions will be qualitatively correct.

In the case where the sign of metaorders and the fundamental component of price changes are independent ($\rho = 0$) and the off-diagonal contribution can be neglected ($\Gamma = 0$), we find that for $T \gg 1$:

$$R_a(T)_{d.} \propto e^{\frac{\sigma_{\ell}^2}{2}a(1-a)} \times \begin{cases} T^{(1-\mu_m-\lambda\sigma_{\ell}^2)/2}, & a < a_c; \\ T^{-\widehat{\beta}(a)}, & a > a_c. \end{cases}$$
 (6.89)

A more refined analysis would be needed for small values of T, but from such an analysis we conclude that $R_a(T)$ should be decreasing with T for small values of a (since $\mu_m \approx 3/2$) and saturating for large values of a (since $\widehat{\beta}(a > a'_c) = 0$). Furthermore, note the non-monotonic behaviour (in a) of the prefactor in Eq. (6.89).

If we now consider the contribution coming from the correlation between metaorders, we get:

$$R_a(T)_{o.d.} \propto e^{-\frac{\sigma_\ell^2 a^2}{2}} \times T^{\mu_m/2 + a\lambda \sigma_\ell^2 - \gamma_\times - \widehat{\beta}(0)}.$$
 (6.90)

For $a \to 0$, the power of T is very close to zero for our default choice of parameters. For higher values of a, the exponent becomes positive and therefore this off-diagonal contribution adds an increasing function of T.²⁹

²⁹In the $T \to \infty$ limit, one should take into account the fact that Σ_a itself becomes dominated by the off-diagonal contribution, see the discussion around Eq. (6.39). Hence the correlation does saturate when $T \to \infty$, as it should be.

If we assume instead that $R_a(T)$ is dominated by informed metaorders, we obtain, in the regime where $\mu_q > 1 + \beta_q$, $\forall q$,

$$R_a(T)_F \propto \rho e^{-\frac{\sigma_\ell^2 a^2}{2}} \times \begin{cases} T^{\mu_m/2 + a\lambda \sigma_\ell^2 - 1}, & a < a_c; \\ T^0, & a > a_c. \end{cases}$$
 (6.91)

The scaling with T indicates that $R_a(T)$ should decrease with T for small values of a, and become independent of T for large values of a. Note that if volatility is mostly due to fundamentals and not due to impact, one finds that $R_a(T)$ is directly proportional to ρ .³⁰

Hence, we see that the three possible contributions to $R_a(T)$ have different monotonicity properties as functions of T, suggesting that different shapes of $R_a(T)$ might be observed empirically. In the following, we will confirm that this is indeed the case. The behaviour of $R_a(T)$ for a given T as a function of a is simpler to describe. One finds that for $a < a_c$

$$R_a(T) = e^{-\frac{\sigma_{\ell}^2 a^2}{2}} \left(A(T) e^{\frac{\sigma_{\ell}^2 a}{2}} + B(T) e^{\lambda \sigma_{\ell}^2 a \log T} \right), \tag{6.92}$$

where A, B are functions of T, the second contribution B(T) coming from the a dependence of the exponents. Hence we expect a behaviour of $R_a(T)$ that always increases for small a, independently of the dominant contribution (diagonal, off-diagonal or fundamental), but with a slope that increases with T in the last two cases. These predictions will be tested against empirical data in the next section.

6.5.7 Empirical data

The theoretical analysis laid out in the previous sections makes several non-trivial predictions:

- 1. Provided the main source of price moves is the average impact of random metaorders, the covariance $\mathbb{E}[\Delta_T \cdot I_T^a]$ behaves as a power-law of T with an effective exponent that is non-monotonic in a, reaching a minimum for some value of a, see Eq. (7.5). If volatility is dominated by the random component of impact (Eq. (6.71)) with a "Fundamental" component unrelated to trading, a linear behaviour in T independent of a should be observed;
- 2. The correlation coefficient $R_a(T)$ contains an off-diagonal contribution growing with T and two contributions (diagonal and fundamental) decaying with

 $^{^{30}}$ It may be realistic to assume that well informed metaorders are larger, and therefore that $\rho \propto q^{\psi}$ where $\psi \geq 0$. In this case, the scaling with a reads $R_a(T)_F \propto e^{\sigma_\ell^2 a(2\psi - a)/2}$, which reaches a maximum for $a^* = \psi$.

Chapter 6	Cha	pter	6
-----------	-----	------	---

T before saturating. Depending on the relative amplitude of these contribu-
tions, different shapes can be expected;

3. For a given T, the correlation coefficient $R_a(T)$ is predicted to be a humped shape function of a.

We will show below that, quite remarkably, all these predictions are in qualitative agreement with empirical data. This will enable us to estimate the parameters λ and λ' . We will also see a marked difference between large tick assets like EUROSTOXX and smaller tick assets (LLOY, TSCO and SPMINI).

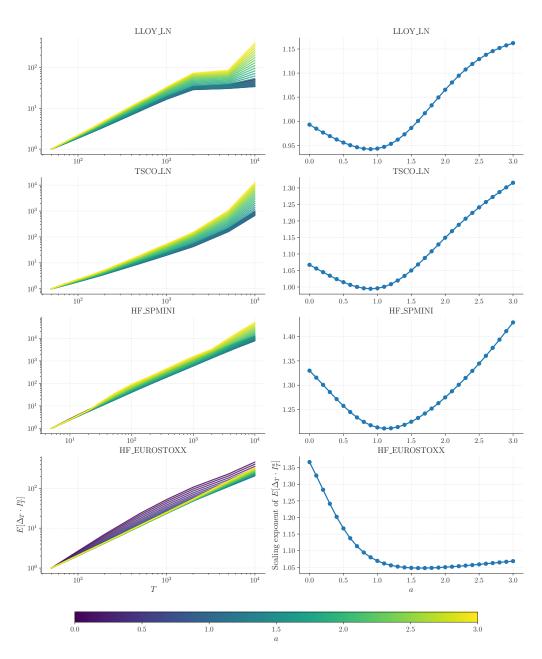


Figure 6.8: Covariance (Δ_T, I_T^a) as a function of (T, a) for the four considered assets. **Left:** Log-log plot of $\mathbb{E}[\Delta_T \cdot I_T^a]$ vs. T for different values of a. **Right:** Scaling exponents as a function of a, obtained by fitting the initial regime $(T < 10^3)$.

Power-law behaviour of the covariance

We first investigate point 1. above. As shown in Fig. 6.8 (left), a power-law behaviour of $\mathbb{E}[\Delta_T \cdot I_T^a]$ as a function of T is approximately verified for all a, albeit with some amount of convexity (for TSCO) or concavity (for LLOY or EUROSTOXX). When plotted as a function of a, the effective exponent of T displays the predicted non-monotonic behaviour, see Fig. 6.8 (right), reaching a minimum for $a \approx 1$ for LLOY and TSCO, $a \approx 1.1$ for SPMINI and $a \approx 1.5$ for EUROSTOXX. The left slope is predicted to be equal to $-\lambda \sigma_\ell^2$ and the right slope to $+\lambda'\sigma_\ell^2$. For LLOY, TSCO and SPMINI we thus find $\lambda\sigma_\ell^2 \approx 0.1$ (not far from the estimate derived from Fig. 6.4) and $\lambda'\sigma_\ell^2 \approx 0.15$ –0.2. For EUROSTOXX, we estimate $\lambda\sigma_\ell^2 \approx 0.5$, a factor 2 larger than from the behaviour of Σ_a (Fig. 6.4), but a very small (but positive) value for $\lambda'\sigma_\ell^2$.

The value of the effective exponent is in the range 0.95-1.35, as expected since the diagonal contributions give an exponent slightly below 1 (Eq. (7.5)) and the off-diagonal contribution yield an exponent ≈ 1.2 for the default values of γ_{\times} and $\widehat{\beta}(0)$ (see Eq. (6.85)).

Note that, as mentioned above, a volatility model based on fundamentals only, leads to a linear behaviour $\mathbb{E}[\Delta_T \cdot I_T^a] \propto T$, independently of the value a, clearly at odds with Fig. 6.8 (right). For the EUROSTOXX, such a linear regime can perhaps be observed for $a \gtrsim 1.5$, but also compatible with Eq. (7.5) if λ' is small.

Correlation vs. T and a

Turning to point 2., Fig. 6.9 shows the correlation coefficient $R_a(T)$ vs. T for different a in two different representations: standard plot and heatmap. A first immediate observation is that these correlations are O(1) for all T values, and peak around 0.45 for stocks and 0.7 for futures. This means that order flow and returns are indeed strongly correlated, as was emphasized many times (see e.g. [37, 88, 105, 116, 117]).

We observe that for LLOY and TSCO, $R_a(T)$ is a mildly decreasing function of T for small a, which becomes mildly increasing for larger a, as expected from the discussion in section 6.5.6, assuming that the diagonal contribution dominates for small a and the off-diagonal contribution kicks in for larger a, or larger values of T (as indeed suggested by the two upper plots in Fig. 6.9), where an upturn of $R_a(T)$ is observed at large T.

For EUROSTOXX and SPMINI, we observe a non-monotonic behaviour of $R_a(T)$ vs. T for small a, with a maximum reached for rather large values of T. This suggests that the non-diagonal contribution is dominant for small a, with an ex-

ponent that is already positive for a=0, i.e. values of μ_m and $\widehat{\beta}(0)$ larger than the default values quoted throughout the Chapter to be compatible with our admittedly rough theoretical analysis. The qualititative behaviour of both futures contracts thus appears to be quite different from that of stocks. Note in particular that the level of the correlation is markedly higher for EUROSTOXX, reaching a maximum value ≈ 0.75 , compared to ≈ 0.45 for stocks (see also Fig. 6.9) and ≈ 0.6 for SPMINI.

Finally, note that the saturation regime where $R_a(T)$ should become independent of T, as predicted by either the "diagonal" hypothesis ((6.89)) or the "Fundamental" hypothesis (Eq. (6.91)), is hardly observable in the data, at least up to 10^4 trades (roughly one trading day). This observation appears to confirm the predominant role of *impact* of mostly uninformed (but correlated) metaorders in the genesis of volatility [7].

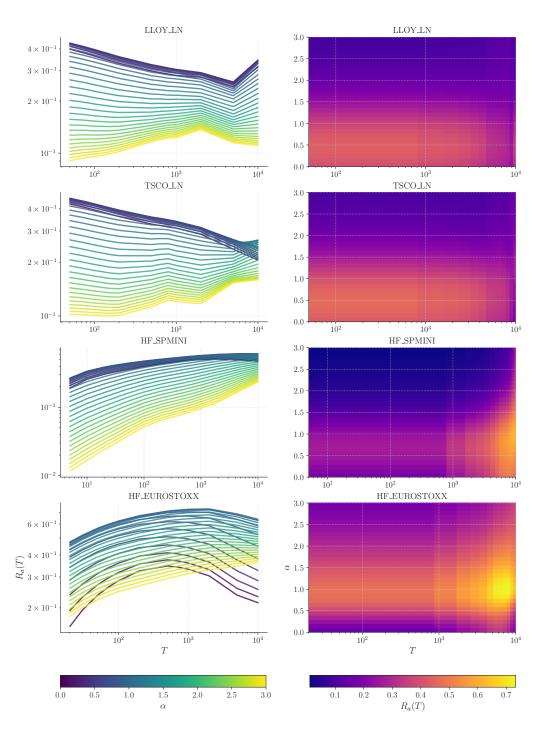


Figure 6.9: Analysis of the correlation function $R_a(T)$ for LLOY, TSCO, SPMINI and EUROSTOXX. Left column: Evolution of the correlation for different values of a, showing the non monotonic behavior **Right column:** Heatmap illustrating the distribution of correlation values within the (a,T) space, indicating that the correlation reaches its peaks for $a \approx 0.5 - 1$, regardless of the T values.

Non-monotonic behaviour of R_a vs. a

Finally, for point 3., Fig. 6.10 shows that the correlation between price returns and generalized order imbalance is non-monotonic as a function of a, reaching a maximum for $a^* \approx 0.5$ for LLOY, TSCO and SPMINI, and $a^* \approx 1$ for EU-ROSTOXX.

Such a non-monotonic behaviour is predicted by our theoretical analysis. Interestingly, when the diagonal contribution to R_a dominates, we expect that the maximum correlation is reached precisely for $a^* = 1/2$, with a peak amplitude that decreases with T, see Eq. (6.89). The data for the two stocks is therefore compatible with the fact that in the small a regime, $R_a(T)_d$ is a decreasing function of T. In this regime, one can also predict that

$$\frac{R_{a=\frac{1}{2}}(T)}{R_{a=0}(T)} = e^{\sigma_{\ell}^2/8},\tag{6.93}$$

to be compared with the data for which this ratio is ≈ 1.1 . The inferred value of σ_{ℓ}^2 is thus around 1, comparable to the direct estimate of σ_{ℓ}^2 from the variance of log-volumes, see section 6.3.4. The full fit of $R_a(T)$ vs. a neglecting that the B term in Eq. (6.89) is given in Fig. 6.10.

For the EUROSTOXX, on the other hand, the maximum is reached for larger values of a, and the ratio defined in (6.93) is much larger (3 – 5), suggesting that the B(T) term in Eq. (7.7) is now dominant, as demonstrated by a fit to the data, see Fig. 6.10. This is consistent with our remark above – that the off-diagonal term dominates $R_a(T)$ for small T. In this scenario, the initial positive slope of $R_a(T)$ vs. a should increase with T, in agreement with data. Neglecting A(T) in Eq. (6.89) and setting $a^* = 1$, we now obtain

$$\frac{R_{a=1}(T)}{R_{a=0}(T)} = e^{\sigma_{\ell}^2(\lambda \log T - \frac{1}{2})},\tag{6.94}$$

which can indeed become large: if we take $\sigma_\ell^2=2$ and $\lambda\sigma_\ell^2=1/2$, as suggested above, one finds that for T=100 the ratio above is ≈ 3.7 . For $\lambda\sigma_\ell^2=1/4$, that ratio is ≈ 1.1 .

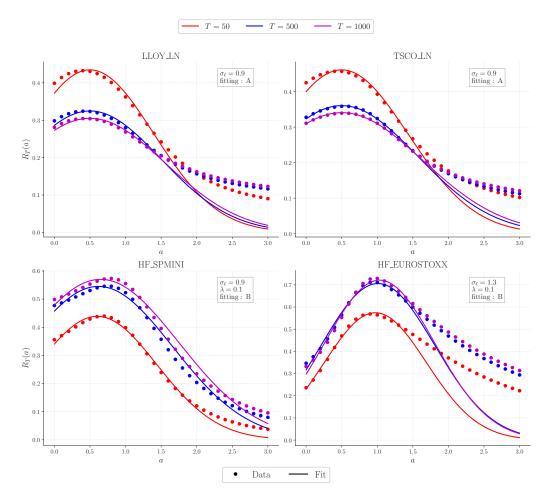


Figure 6.10: Fit of the correlation function $R_a(T)$, for several T. Since Eq. (7.7) holds only for $a < a_c$, we restrict the fit to a < 1.5. The empirical estimates of σ_{ℓ} obtained with this fit turn out to be surprisingly close to the ones obtained in Fig. 6.1.

Empirical covariance: conclusion

All empirical data appear to confirm the qualitative validity of our predictions, some of them being rather non-trivial. One of the main assumptions of our model is that the exponents describing the autocorrelation of child orders (μ_q) and the decay of their impact (β_q) depend on the size q of these child orders. This, in turn, leads to a non-monotonic dependence of the scaling of the covariance $\mathbb{E}[\Delta_T \cdot I_T^a]$ with a, a1 and of the correlation $R_a(T)$ with both T and a2. Although some of our

 $^{^{31}}$ Such a non-monotonic behaviour cannot be explained simply from the power-tail distribution of volumes q, see G. Maitrier et al., in preparation.

predictions fail to explain the data quantitatively, we are tempted to ascribe these discrepancies on the bluntness of our approximations, which, we argue, correctly capture the mechanisms at play.

Perhaps the two most important conclusions of this section, beyond the success of our model in capturing the main trends of the covariance data, are:

- Stocks and futures seem to differ quantitatively when it comes to the *correlation* between order imbalance and price changes. In particular, the correlation between the two is stronger for futures contract, and is reached for longer time intervals T and with more weight given on large child order volumes (i.e. $a^* = 1$ instead of $a^* = 1/2$). This points towards a stronger role of metaorder correlations for futures than for stocks.
- The hypothesis according to which most of the volatility comes from fundamental information is hard to reconcile with the data. First, as discussed in section 6.4.5, this assumption does not enable one to understand why the square-root law, which applies to all metaorders informed or not [42, 44, 118] is proportional to volatility. Second, it predicts that the covariance of price returns and volume imbalances is proportional to *T independently of a*, at odds with empirical data. As argued in [7, 33, 104], the most plausible hypothesis is that volatility stems from trading alone. In other words, the excess volatility puzzle has a microstructural origin.

6.6 Conclusion

The aim of this Chapter was to reconcile several apparently contradictory observations: is a square-root law of metaorder impact that decays with time compatible with the random-walk nature of prices and the linear impact of order imbalances? Can one entirely explain the volatility of prices as resulting from a "soup" of indistinguishable, randomly intertwined and uninformed metaorders?

In order to answer these questions, we have introduced a new theoretical framework to describe metaorders with different signs, sizes and durations, possibly correlated between themselves, which *all* impact prices as a square-root of volume (which we assume as an input) but with a subsequent time decay characterized by an exponent $\beta \neq \frac{1}{2}$, i.e. different from the one suggested by the classic propagator model [7, 31] or the LLOB model [74]. We proposed a generalized propagator model to account for such a feature.

We then established that the power-law tailed distribution of metaorder durations is not sufficient to counteract impact decay, leading to price sub-diffusion. Rather, as in the original propagator model, price diffusion is ensured by the long memory

Chapter 6.

of cross-correlations between metaorders, which is indeed present in data. In fact, we conjecture that the intra- and cross-correlations between child orders decay roughly in the same manner, a feature that may be crucial for explaining the success of the construction of synthetic metaorders from public data [3] - Chapter 5.

The existence of correlations between metaorders is therefore a crucial ingredient to recover price diffusion. The old debate between order splitting and "herding" that seemed to have been closed by several papers in favor of splitting [16, 119], is perhaps not so clear-cut. Such correlations could be due either to the fact that many participants use the same trading signals, or that copy-cat metaorders follow past order flow, or else that some traders successfully predict the future behaviour of other participants. Note that within our story any predictive alpha signal manifests itself through autocorrelated metaorders, since only trades can move the price (see also the discussion in [7], chapter 20).

In view of the strongly fluctuating volumes q of child orders, one quickly realizes that one needs to account for heterogeneity in the distribution of metaorder durations, and in the resulting decay of their impact, which we parametrized by two q-dependent exponents, μ_q and β_q , assumed to depend linearly on $\ell = \log q$. In a nutshell, this is needed to account for the fact that metaorders with large q are less autocorrelated than those with small q, as shown in Fig. 6.3. This feature allowed us to account semi-quantitatively for the way the moments of generalized volume imbalance scale with time T, and more importantly how the correlation between price changes and generalized volume imbalance scales with T. We predicted that the corresponding power-law should depend in a non-monotonic fashion on the parameter a that allows one to put the same weight on all child orders (a = 0)or overweight large orders (a large), a behaviour clearly borne out by empirical data, see Fig. 6.8. We also predicted that the correlation between price changes and volume imbalances should display a maximum as a function of a for fixed T, which again matches observations, see Fig. 6.10, with fitting parameters fixed with previously determined values. We found that stocks and futures appear to differ quite markedly in terms of these metrics, which could provide an interesting new way to characterize price formation mechanisms in different markets.

Such noteworthy agreement between theory and data suggests that our framework correctly captures the basic mechanism at the heart of price formation, namely the average impact of metaorders. We claim that our results strongly support the "Order-Driven" theory of markets, according to which it is the mechanical impact of trades, independently of any notion of "fundamental information", that generates volatility in financial markets, a picture advocated in [7, 53, 104, 105] and, in a different context, in [52]. In particular, the Efficient Market Hypothesis,

which posits that volatility is mostly due to the variation of "fundamental value", cannot easily explain our results on the correlations between price changes and order flow imbalance.

Take Home Messages

- We developed a generalized propagator model in which all metaorders follow the square-root law, with a decaying impact governed by an exponent $\beta = \frac{1-\gamma}{2}$.
- By introducing metaorder cross-correlations, we showed that such interactions are essential to recover price diffusion, challenging the prevailing focus on order splitting alone.
- Our framework predicts a non-monotonic relationship between pricevolume imbalance correlations and the weighting parameter a, confirmed by empirical results.
- We provided strong evidence supporting the "Order-Driven" theory of markets, where volatility originates from the mechanical impact of trades rather than fundamental information.

On alpha prediction and permanent impact

Let me pause here briefly and return to the questions raised in Section 2.5, now viewed in light of the results obtained so far. As the reader has likely understood, permanent impact is a crucial concept and yet it remains the subject of ongoing debate. This is partly because it is extremely difficult to estimate, but also because it touches on deeper matters of belief and perspective—across both academia and industry, as I've had the chance to witness firsthand.

While the framework developed in this thesis offers what I believe to be a new interpretation of these questions, I want to emphasize that what follows are more conjectures than definitive claims. These ideas still require careful empirical validation.

So far, we have seen that both the volatility and the trending behavior of prices can be largely explained by trading activity alone, consistent with the so-called excess volatility puzzle. In this view, price moves because trades are made. But under our framework, each trade belongs to a metaorder. Suppose now that one possesses a perfect signal—say, insider information—that a given asset will double in price tomorrow. For that to happen, it cannot be enough to *know* this: the price must be *pushed* in that direction by executed metaorders.

In this light, alpha is not fundamentally about discovering some deep economic truth. Rather, it becomes about being able to anticipate how others will trade. There may well be an economic reason that justifies the price doubling—but if no one else shares or acts upon that belief, the price will not move.

One might argue that other participants eventually acquire the same information and follow your lead. In that case, your alpha lies not just in having the information, but in acting on it *first*. Their trades—whether due to direct imitation, information diffusion—amplify the move you began.

Your profit, then, arises from their impact. They push the price further in your direction, even after you have finished executing your trade. And just as the accumulation of correlated market orders with decaying impact can produce a permanent price shift, the aggregation of correlated metaorders —even if each has no permanent impact individually —can likewise give rise to an effective permanent impact. This could align with the foundational work of Farmer [120].

In the current literature—see, for instance, the debate between Gabaix et al. [52] and Bouchaud [53]—permanent market impact is typically modeled as linear. This has long raised a conceptual issue: how can impact be dynamically concave yet result in a linear permanent effect?

Here, I would suggest that this is because the two phenomena are not governed by very distinct mechanisms. The concavity of dynamical impact, ie the SQL, is a consequence of microstructural factors—liquidity constraints, as captured for instance by the latent liquidity theory of Donier et al. [15], to which I strongly subscribe (though I believe it could still benefit from refinement). By contrast, permanent impact arises from the autocorrelation of investment decisions at the level of investors: the persistence of belief, of information, of allocation decisions, of trading behavior.

I fully acknowledge that this line of reasoning is still in its early stages and needs further time to mature. I hope to develop it more rigorously in future work.

Chapter 6.

Chapter 7

From Abstraction to Animation: An Artificial Market Simulator

Nothing is more practical than a good theory.

Ludwig Boltzmann

In the next chapter, we will confirm the ideas from the previous one using numerical simulations of the model, in order to (i) confirm the qualitative validity of our theoretical analysis, (ii) propose an efficient way to generate an exhaustive dataset, capturing the subtle interplay between order flows and prices, and (iii) investigate whether the metaorder proxy introduced in Chapter 5 can be used to construct synthetic metaorders directly from simulated data.

From:

The Subtle Interplay between Square-root Impact, Order Imbalance and Volatility II: An Artificial Market Simulator G. Maitrier, G Loeper, JP. Bouchaud

Contents					
7.1	Introduction				
7.2	A brief reminder of the generalized propagator model 140				
7.3	How to simulate our model?				
7.4	Empirical stylized facts vs. simulations 145				
	7.4.1 The q -dependence of the autocorrelation of trades \dots 145				
	7.4.2 The scaling of the order flow imbalance $\dots \dots 146$				
	7.4.3 Recovering a diffusive price $\dots \dots 149$				
	7.4.4 Aggregated impact and anomalous rescaling 151				
	7.4.5 The covariance coefficient $\dots \dots \dots$				
	7.4.6 The correlation coefficient				
7.5	The puzzling effectiveness of proxy metaorders 157				
	7.5.1 How the acceptance window drives mapping accuracy . 161				
7.6	Conclusion				

7.1 Introduction

Market microstructure —and more specifically, limit order books —constitutes the microscopic environment in which prices are formed. It can be viewed as a black box: orders, submitted by various market participants, enter as inputs, and the resulting output is the observed transaction price. We describe it as a black box not because it is fundamentally opaque or inaccessible, but because the interactions that occur within it are governed by a multitude of heterogeneous agents, operating at different timescales, with diverse objectives and information sets. These interactions generate a highly nonlinear and noisy environment, making it extremely challenging to disentangle cause and effect, or to isolate the fundamental mechanisms driving price formation. For these reasons, understanding the inner workings of this black box —i.e. constructing models that faithfully reproduce both order flow patterns and price behavior —remains one of the central challenges in market microstructure research. Indeed, several known stylized facts about price impact (the famous square-root law), order flow (with its long-memory properties) and volatility (i.e. that prices are diffusive) appear to be disconnected and, at least at first sight, hard to accomodate.

In our previous paper [4], we introduced a theoretical framework that allows one to reconcile the statistical properties of order flow and price dynamics. Our model makes detailed and somewhat non-trivial predictions about the cross-correlations between order flow and price variations that appear to all be borne out by empirical data on stocks and futures.

In order to derive such predictions, we made several assumptions and simplification that may appear somewhat strong and uncontrolled [4]. Whereas our theoretical model is challenging to solve in complete generality, it has the notable advantage of being straightforward to simulate numerically. The present follow up Chapter serves a dual purpose. First, it offers additional evidence for the robustness of our theoretical model by showing that the approximate analytical treatment proposed in [4] actually correctly describes the key empirical phenomena. Second, we introduce what we believe to be a versatile and realistic simulation tool that captures the intricate interplay between order flow and price formation.

This latter contribution could be of significant interest to the industry. Generating realistic market dynamics —encompassing both prices and order flow —remains a notoriously challenging task. It is fair to say that many existing approaches, including those based on neural networks, see [98, 121], often fall short of capturing the full complexity of market behavior. Yet, such generative models are essential for several practical applications: they enable robust strategy backtesting, and they provide enhanced fitting capabilities in situations where real financial data is limited or unavailable. Our framework is based on a direct, mechanistic description of order flow and price impact that abstracts away from the infinite complexities of the full order book dynamics, and surely suffers from some short-comings, but is transparent and computationally trivial. Hybridizing our model with higher frequency, data driven generative model would be very interesting.

This Chapter is divided in three parts, and contains:

- A detailed framework for generating synthetic data based on our assumptions. This synthetic dataset closely resembles the ideal one (similar to the TSE dataset, for instance) and includes all relevant information about order flow, metaorder ids, execution time, and impact.
- The reproduction of empirical graphes showed previously, but for simulated price, using parameters with fitted on real data. In this section, we aim to reproduce the behavior observed for TSCO by setting $\ell=6.5$ and $\sigma_\ell=1$ -empirical values obtained by fitting Ξ on TSCO trades quantities.
- A discussion of the puzzling possibility of reconstructing metaorder proxies from public data introduced in [3] and Chapter 5 that we confirm within our artificial market, thereby validating the procedure.

7.2 A brief reminder of the generalized propagator model

The present study is based on the unified framework proposed in [122] and in the previous Chapter. To succinctly recall the context, we summarize the model as follows:

- The order flow is composed of a succession of metaorders, initiated with rate ν per unit time. The size (i.e., the number of child orders per metaorder) is distributed according to a power law, $\Psi_q(s)$, with a q-dependent tail exponent μ_q , where q is the size of the child orders, assumed to be constant within each metaorder. Such child volumes are distributed according to a lognormal distribution with parameters (m, σ_ℓ) . To account for the empirical sign autocorrelations (see Fig. 3 of [122] and Fig. 7.1 below), we set $\mu_q = \mu_1 + \lambda \log(q)$.
- We also introduced the possibility of correlating the sign of different metaorders, starting respectively at time t and $t+\tau$. If ε_t is the sign of the t^{th} metaorder of the day, we assume that for $\tau \gg 1$

$$\mathbb{E}[\varepsilon_t \varepsilon_{t+\tau}] = \Gamma \tau^{-\gamma_{\times}}. \tag{7.1}$$

• Finally, to understand price formation from order flow, we introduced a generalized propagator model. This instrument is crafted to incorporate the three main stylized characteristics of the impact of metaorders (see [2, 7] and refs therein): (i) impact grows on average as the square-root of the number of child orders being executed, (ii) average peak impact at the end of the execution solely depends on the square root of the traded volume, and (iii) average impact subsequently decays as a slow power-law of time after the end of execution.

We posited that the impact of a child order of volume q, executed at time t' on the price at time t > t', knowing that the metaorder started at t = 0, is given by

$$G_q(t' \to t) = \frac{\theta \sqrt{q}}{(\varphi t' + n_0)^{1/2 - \beta_q}} \left(\frac{\tau_0}{t - t' + \tau_0}\right)^{\beta_q}, \qquad (\beta_q < \frac{1}{2})$$
 (7.2)

with θ , n_0 , τ_0 are constants—see section 7.3 for details—and φ the participation rate of the metaorder. Empirical observations led us to the following specification $\beta_q := \beta_1 - \lambda' \log(q)$, meaning that impact decay is slower for large child orders, as intuitively meaningful.

The entire framework is motivated and explained more thoroughly in [122], and leads to the following predictions:

1. The generalized order flow imbalance: We defined the weighted order flow imbalance, where ε_t is the sign of *child orders*:

$$I_T^a = \int_0^T \mathrm{d}t \,\varepsilon_t q_t^a,\tag{7.3}$$

and its moments $\Sigma_{I^a}^{(2n)}:=\mathbb{E}[(I_T^a)^{2n}],$ for which our theory predicts a non-trivial behavior:

$$\Sigma_{a,1}^{(2n)} \propto \begin{cases} T^{2n+1-\mu_m-2na\lambda\sigma_\ell^2}, & a < a_c(n); \\ T, & a \ge a_c(n), \end{cases}$$
 (7.4)

with $a_c(n) = (1 - \mu_m/2n)/\lambda \sigma_\ell^2$.

2. The time-dependent covariance function: Armed with the order flow description and the generalized propagator, we describe the interplay between price returns Δ_T and order flow by computing the covariance $\mathbb{E}[\Delta_T \cdot I_T^a]$. Our model tells us that such a quantity should behave as a power-law of T with an exponent that is a *non-monotonic* function of a:

$$\mathbb{E}_{q}[\Delta_{T} \cdot I_{T}^{a}] \propto \begin{cases} T^{5/2 - \widehat{\mu}(a)}, & \widehat{\mu}(a) = \mu_{m} + (a + \frac{1}{2})\lambda\sigma_{\ell}^{2} & a < a_{c}'; \\ T^{1 - \widehat{\beta}(a)}, & \widehat{\beta}(a) = \beta_{m} - \left(a + \frac{1}{2}\right)\lambda'\sigma_{\ell}^{2} & a > a_{c}'; \end{cases}$$
(7.5)

where $\mu_m = \mu_1 + \lambda m$, $\beta_m = \beta_1 - \lambda' m$ and a'_c such that $\widehat{\mu}(a'_c) = \mu_{q'_c}$, with $5/2 - \mu_{q'_c} = 1 - \beta_{q'_c}$.

3. The correlation coefficient: Finally, our model also allows one to predict the behavior of the following correlation coefficient:

$$R_a(T) := \frac{\mathbb{E}[\Delta_T \cdot I_T^a]}{\Sigma_T \Sigma_{I^a}}, \qquad \Sigma_T := \sqrt{\mathbb{E}[\Delta_T^2]}, \qquad \Sigma_{I^a} := \sqrt{\mathbb{E}[(I_T^a)^2]}$$
 (7.6)

The following prediction fits surprisingly well empirical data:

$$R_a(T) = e^{-\frac{\sigma_{\ell}^2 a^2}{2}} \left(A(T) e^{\frac{\sigma_{\ell}^2 a}{2}} + B(T) e^{\lambda \sigma_{\ell}^2 a \log T} \right), \tag{7.7}$$

for $a < a_c$, and A, B two functions of T. In particular, for a given T, $R_a(T)$ is non-monotonic in a and reaches a maximum for $a \approx 1/2$ for stocks and $a \approx 1$ for futures.

Although the model is based on only a few assumptions, the theoretical predictions above are not straightforward, and some uncontrolled approximations needed to be made. Still, the empirical data we analyzed in [122] agree surprisingly well

Chapter 7.

with our predictions. By simulating numerically the very same model, our goal is to replicate these stylized facts without any analytical approximations, and demonstrate that we have identified the correct mechanism. This will confirm that such good fits are not merely coincidental and that uncontrolled approximations are not, unwittingly, responsible for the success of our theory.

7.3 How to simulate our model?

Whereas generating order imbalances is relatively straightforward, simulating realistic price dynamics is more delicate. In our model, child orders from different metaorders can in principle be executed simultaneously, which complicates the price formation process. Furthermore, while the execution of a child order is clearly a discrete event, its impact decays continuously over time and should be taken into account at each timestep.

After testing several approaches, we found that using actual timestamps yields the most transparent and realistic simulations. The simulation procedure is thus divided into three main steps:

- Generating metaorders: We specify the average number of metaorders and the total trading period for the simulation (e.g., 10000 metaorders over an 8-hour trading day). This defines the rate ν at which new metaorders start. For each metaorder, we define the following parameters: a volume q, distributed as a log-normal truncated below q = 1, a size s (distributed according to Ψ_q), a sign (either randomly assigned or generated with crossmetaorder correlations), and a starting time, randomly chosen within the trading day with density ν . We ensure that starting times are unique, as they will later serve as metaorder identifiers. It is possible to control the trading rate and liquidity by modifying $\nu, \varphi, m, \sigma_{\ell}$, as described in section 7.2.
- Deriving the corresponding order flow: The order flow is then generated by iterating over all time-sorted metaorders. For each one, we store the execution time of child orders by generating time intervals δt thanks to a Poisson process: $\delta t \sim e^{-\varphi \delta t}$. For example, the second child order is executed at time $t = t_{\text{start}} + dt_1$. This approach allows us to sort all child orders by their execution time, thereby constructing an order flow that closely resembles real trade-by-trade data (or more precisely that from the TSE dataset). Each event (here execution) includes the timestamp, volume, sign, child order rank, and the time elapsed τ since the start of the corresponding metaorder. In addition, and specific to our model, we store the value of β_q associated with each metaorder.

timestamp	volume	sign	rank	$\mathbf{timestamp}_{\mathrm{start_meta}}$	β_q
10:32:01.35	10	+1	1	10:32:01.35	0.31
10:32:01.57	150	-1	7	09:15:03.86	0.20
10:32:02.15	80	+1	2	09:34:43.12	0.25
10:32:02.76	120	-1	1	10:32:02.76	0.22
10:32:02.78	90	-1	5	09:47:52.27	0.28

Table 7.1: Simulated order flow data with metaorder decomposition. Each row represents the execution of a child order, with 'rank' column indicating the position of the child order within its metaorder. The 'timestamp(start_meta)' column records the start time of the metaorder and also acts as an identifier, as it is uniquely assigned to each metaorder

• Reconstructing the mid-price: Armed with this simulated order flow and our generalized propagator, reconstructing the price dynamics becomes straightforward. We define the price p_t as the mid-price just before the execution occurring at time t. To compute this price, we aggregate the contributions from all child orders executed prior to t, ie $t_{\rm exec} < t$. We use the generalized propagator to compute for their respective impacts and sum them. Although not computationally optimized (with complexity $\sim \mathcal{O}(N^2)$), this algorithm appears to be the most rigorous. It also preserves a key property of price impact observed in real markets: most of real market orders have zero immediate impact (as their volume is smaller than the prevailing best), but their impact builds up over time (on this point, see e.g. [2, 7, 31]).

To complement this description, we provide the following pseudo-code:

Algorithm 3 Simulation of Price Impact from Correlated Metaorders

Require: Number of metaorders N and base parameters γ_{\times} , μ_1 , β_1 , (m, σ_{ℓ}) , (λ, σ_{ℓ}) λ'

Ensure: Time series of executed orders with associated impact prices

- 1: Set ν , the Poisson rate for metaorders initiation and φ the participation rate within a metaorder.
- 2: Draw N metaorder start times $\{t_i^{\text{start}}\}_{i=1}^N$, with $t_{i+1}^{\text{start}} t_i^{\text{start}} \sim \text{Exp}(-\nu dt)$
- 3: **for** i = 1 to N **do**
- Sample metaorder volume $q_i \sim LN(m, \sigma_\ell)$ and size $\mu_i = \mu(q_i, \mu_1, \lambda)$
- Compute impact exponent $\beta_i = \beta(q_i, \beta_1, \lambda')$ 5:
- Sample metaorder sign ε_i autocorrelated sign time serie, if γ_{\times} 6:
- Sample number of child orders $s_i \sim \Psi_{q_i}$ 7:
- 8:
- Generate inter-arrival times $\{\delta t_k^{(i)}\}_{k=1}^{s_{i-1}} \sim \operatorname{Exp}(-\varphi \delta t)$ Compute execution times $t_k^{(i)} = t_i^{\operatorname{start}} + \sum_{j=1}^k \delta t_j^{(i)}$ Store each child order as $(t_k^{(i)}, \varepsilon_i, q_i, t_i^{\operatorname{start}}, \beta_i, \mu_i)$ 9:
- 10:
- 11: end for
- 12: Sort all child orders by execution time $\{t_k\}$
- 13: Initialize price impact array $p_k \leftarrow 0$
- 14: **for** each execution time t_k **do**
- 15: Identify past orders j < k
- Apply the generalized propagator: 16:

$$p_k = \sum_{j \le k} \varepsilon_j \cdot \sqrt{q_j} \left(\varphi(t_j - t_j^{\text{start}}) + n_0 \right)^{-\frac{1}{2} + \beta_j} \cdot \left(\frac{\tau_0}{t_k - t_j + \tau_0} \right)^{\beta_j}$$

- 17: end for
- 18: Convert timestamps to realistic time
- 19: **return** DataFrame of child orders with $\{t_k, p_k, q_k, t_k^{\text{start}}, \varepsilon_k, \beta_k\}$

This simple model relies on only a few parameters that require fine-tuning. To stay as close as possible to [122], we set $m \in \{3,6\}$, $\sigma_{\ell} = 1$, $\lambda \sigma_{\ell}^2 = \frac{1}{8}$, and $\lambda' = 2\lambda$. We also set $\mu_m = 1.5$, and $\beta_m = 0.25$. For consistency, we ensure that $0 < \beta_q < 1$.

Finally, we fix $n_0 = 3$, based on empirical observations in [2], after verifying that this parameter has only a mild influence on the rest of the system. The average time between two child orders, denoted τ_0 , is theoretically given by $\tau_0 := (\nu \varphi \bar{s})^{-1}$ [122]. For simplicity, we assume a uniform participation rate φ across metaorders. By adjusting ν and φ , one can control the average number of concurrently active metaorders. In the rest of the paper, we will typically impose $\nu = 1.5 \cdot 10^{-3}$ and

$$\varphi = 2 \cdot 10^{-3}.$$

7.4 Empirical stylized facts vs. simulations

We simulated the system under five different scenarios in order explore the relative importance of metaorder correlation, child volume fluctuations and the q-dependence of exponents β and μ . We summarize the different names of these specifications in Table 7.2.

Name	Metaorder Correlation	q-Dependence	q-Fluctuations
NC-NVD-NVF	$\Gamma = 0$	$\lambda, \lambda' = 0$	$q \equiv 1$
NC-NVD-VF	$\Gamma = 0$	$\lambda, \lambda' = 0$	$LN(m,\sigma_\ell)$
NC-VD-VF	$\Gamma = 0$	$\lambda, \lambda' \neq 0$	$LN(m,\sigma_\ell)$
C-NVD-VF	$\Gamma > 0$	$\lambda, \lambda' = 0$	$LN(m,\sigma_\ell)$
C-VD-VF	$\Gamma > 0$	$\lambda, \lambda' \neq 0$	$LN(m,\sigma_\ell)$

Table 7.2: Summary of the five simulated configurations. Each model is named using a triplet notation, with C = correlation (described by parameter Γ), VD = volume dependence of μ_q , β_q , VF = volume fluctuations. Here, "N" indicates negation, such as ND = no metaorder correlation (Γ = 0, see Eq. (7.1)), NVD = no volume dependence ($\lambda, \lambda' = 0$), NVF = no volume fluctuations ($\sigma_{\ell} = 0$). The fully realistic case corresponds to the last line C-VD-VF.

7.4.1 The q-dependence of the autocorrelation of trades

We begin by examining the relationship between child order volume and their autocorrelation in the C-VD-VF scenario, which captures all effects we purport are important. To this end, we partition the simulated rescaled volume $\tilde{q} = q/\phi_D$, where ϕ_D denotes the daily traded volume, into four logarithmic bins \mathcal{B} . For each bin, we compute the sign autocorrelation function defined as

$$C_{\mathcal{B}(\tilde{q})}(\tau) = \mathbb{E}[\varepsilon_{\mathcal{B}(\tilde{q})}(t)\varepsilon_{\mathcal{B}(\tilde{q})}(t+\tau)] \propto \tau^{-\gamma(q)}$$

The autocorrelation functions are displayed in Fig. 7.1, in log-log scale, along with the unconditional autocorrelation function (dotted line). As observed in the data [122], the effective memory exponent $\gamma(q)$ systematically increases with volume, ranging from 0.4 (long memory) to 1.3 (short memory). This graph is strikingly similar to the one obtained for the EUROSTOXX, see Fig 3. in [122].

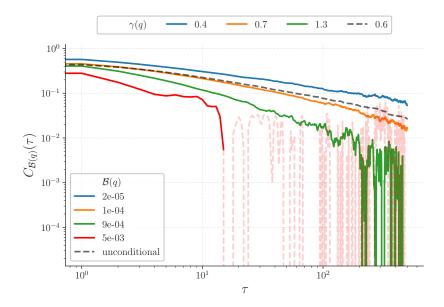


Figure 7.1: Evolution of the sign autocorrelation of market orders based on their corresponding volume bin $\mathcal{B}(q)$. Simulation were done in the C-VD-VF case, with $m=3, \sigma_\ell=1$ and $\lambda\sigma_\ell^2=1/8$. The dotted line corresponds to the *unconditional* autocorrelation function. Compare with Fig 3. in [122].

It is straightforward to verify numerically that the q-dependence of μ_q is indeed responsible for this phenomenon. If the order flow is simulated without incorporating this dependence, the stylized fact completely disappears, with $\gamma(q) \approx 0.5$ independently of q (data not shown).

7.4.2 The scaling of the order flow imbalance

We now turn to the scaling behavior of the moments of the generalized order imbalance $\Sigma_{I^a}^{(2n)}$, which is one of the main successes of the theoretical framework introduced in [122]. When q is constant, the dependence on a disappears trivially, and the imbalance was shown in [122] to follow a truncated Lévy distribution, entirely driven by the long memory of trade signs, thereby justifying the scaling $\Sigma_{I^a}^2 \sim T^{3-\mu}$ with $\mu=1.5$ for the NC-NVD-NVF simulation. However, by introducing a q-dependent μ_q (i.e. when $\lambda>0$), we retrieve scalings that resemble very closely empirical ones, see Fig. 7.2, both with (C) and without (NC) metaorder correlations, as expected.

Note that volume fluctuations alone can induce a spurious dependence of the scaling exponent on a (see NC-NVD-VF in Fig. 7.2) which is due to finite size effects, for which extreme events are artificially amplified as a increases, with a

mechanism similar to the Random Energy Model (REM) in spin glass theory [123, 124]. Indeed, we only simulated 100 trading days with approximately 50000 trades each day such that these finite-size effects are noticeable.

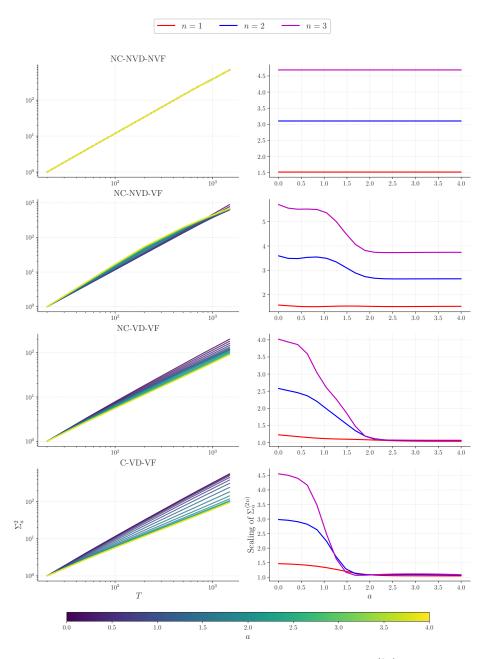


Figure 7.2: Right column: Scaling behavior of the moments $\Sigma_a^{(2n)}$ as a function of trade time T, from which the scaling exponent is extracted via a log-log regression. Left column: Scaling exponent plotted as a function of a. As predicted by our model, increasing a—which gives greater weight to large-volume orders—reduces the scaling exponent. We set $m=6,\sigma_\ell=1$ and $\lambda\sigma_\ell^2=1/8$.

7.4.3 Recovering a diffusive price

A well known puzzle in the literature is the compatibility of decaying square-root impact, long-memory of trade signs and the diffusivity of prices —see [7, 31, 33, 101, 125]. Several solutions to this conundrum were proposed in Section 4 of Ref. [122]. In particular (i) the sign of metaorders themselves should be long-range correlated, as in Eq. (7.1) and (ii) large child orders tend to have a permanent impact, i.e. beyond some value called q_0 in [122], the decay exponent β_q becomes zero.

These two scenarii are both confirmed by numerical simulations: we indeed find that the generalized propagator model leads to a sub-diffusive price in the absence of metaorder correlations ($\Gamma=0$) and without volume effects. Introducing metaorder autocorrelations with the correct exponent γ_{\times} or incorporating a volume dependence β_q restores price diffusivity at long times. By correctly tuning Γ and λ' , one can control the full signature plot and not only the long time diffusive behaviour, and reproduce empirical results that show a variety of possibly short time behaviour, from locally trending to locally mean-reverting —although tick size effects, not modeled here, are expected to play an important role at short times.

Chapter 7.

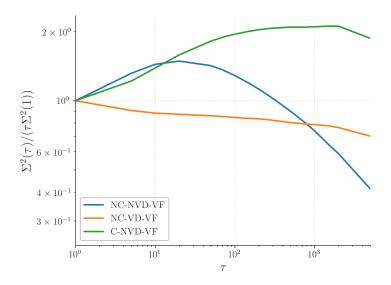


Figure 7.3: Signature plot Σ^2/τ of the simulated price as a function of the trade lag τ . Diffusion corresponds to a flat, horizontal signature plot. The generalized propagator model NC-NVD-VF (blue curve) results in sub diffusive behavior, as expected, while the two other impact models exhibit diffusive behavior after an initial trending phase (C-NVD-VF, green line) or mean-reverting phase (NC-VD-VF, orange line). Simulations were conducted for $\Gamma=0.1$ in the C-NVD-VF case, and $\lambda=\lambda'=1/6$ for the NC-VD-VF case. In both cases, we used $\varphi=2\cdot 10^{-3}, \mu_m=1.5, m=3$ and $\sigma_\ell^2=1$

As in [122], we can also investigate 2n-moments of price changes, and check that all moments scale asymptotically as T^n , as for empirical data, see Fig. 7.4. We insist again that we work here in trade time, so that multifractal effects due to intermittent activity bursts, are not present.

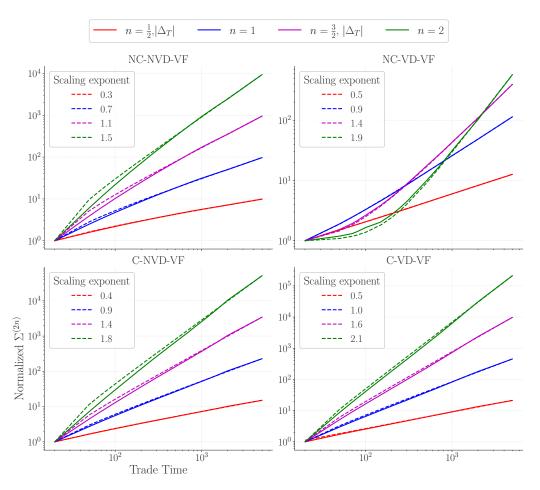


Figure 7.4: Scaling of the moments of price changes $|\Delta_T|^{2n}$ as a function of trade time T. We normalized the moment values such that all curves begin at 1 for T=1. To account for short term, we fitted the data as $\Sigma^{(2n)}=a_0+a_1T^{\zeta_n}$ and present the values of ζ_n in the legend.

7.4.4 Aggregated impact and anomalous rescaling

Aggregated impact is a very natural observable to investigate, but it also turns out to be highly non trivial. It is defined as the conditional expectation $\mathbb{E}[\Delta|I^a]$ of price change Δ given an imbalance I^a , is a natural and empirically accessible observable [37, 115]. However, it exhibits non-trivial behavior that departs from the standard square-root law, with scaling properties that vary significantly with the time horizon T.

In particular, for a=0, the initial slope of $\mathbb{E}[\Delta|I^0]$ scales as $T^{-\omega}$ with $\omega\approx 1/4$, a

result documented in [7, 37]. While a Gaussian assumption would suggest a linear relation

 $\mathbb{E}[\Delta|I^a] = \frac{\mathbb{E}[\Delta \cdot I^a]}{\Sigma_{Ia}^2} I^a, \tag{7.8}$

such an approximation has a priori no reason to hold within in our setting, where I^a is a truncated Lévy variable. Despite this, Eq. (7.8) still captures the correct T-scaling.

We now revisit this observable using simulations based on our model and confirm that the anomalous rescaling $\sim T^{-\omega}$ is precisely recovered, validating the theoretical prediction. However, Fig. 7.5 shows that the concavity seen in empirical curves for large imbalances is absent in our simulations. As demonstrated in Ref. [37], such a concavity is due to a selection bias, not described in our model: large orders tend to be executed when large limit orders are available on the other side, limiting the impact of these market orders.

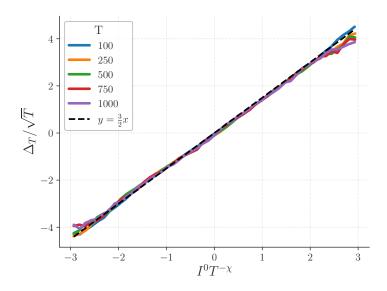


Figure 7.5: Aggregated impact as a function of sign imbalance for C-ND-V simulations. As in real market data, curves corresponding to different values of T collapse onto a single master curve after appropriate rescaling. We find a scaling exponent $\chi=0.75$, in close agreement with the theoretical prediction $1/\mu$, as we simulated with $\mu=1.5$. The slope exponent $\omega=0.25$ is also consistent with empirical observations.

7.4.5 The covariance coefficient

We now focus on the covariance coefficient. Our theoretical predictions suggest that the non-monotonic shape as a function of a originates from volume fluctua-

tions —particularly in the upward branch, which depends on the parameter λ' in the relation $\beta(q) = \beta_1 - \lambda' q$ (see Eq. (7.5)). We clearly confirm this phenomenon in Fig. 7.7, case C-VD-VF. Some aspects still require further investigation, in particular why the NC-VD-VF configuration exhibits a monotonically increasing pattern, when Eq. (7.5) predicts no dependence on metaorder correlations. Nevertheless, we believe that the difference between C-NVD-VF (or NC-NVD-VF) and C-VD-VF supports and reinforces our claim that volume fluctuations coupled to volume dependence of the impact decay is key to account for such a non monotonic behaviour.

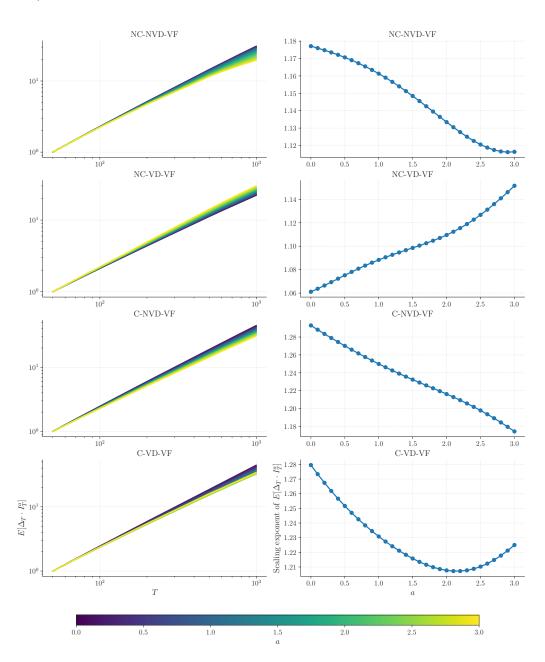


Figure 7.6: Covariance (Δ_T, I_T^a) as a function of (T,a) for simulated markets, with $m=3, \sigma_\ell=1, \lambda\sigma_\ell^2=1/8$ and $\lambda'=\lambda$, for the four configurations considered here. From top to bottom NC-NVD-VF, NC-VD-VF, C-NVD-VF and C-VD-VF. Left: Log-log plot of $\mathbb{E}[\Delta_T \cdot I_T^a]$ vs. T for different values of a. Right: Scaling exponents as a function of a, obtained by fitting the initial regime $(T<10^3)$.

7.4.6 The correlation coefficient

Finally, an important quantity is the correlation coefficient $R_a(T)$, for which our theoretical model also predicts a non-trivial behavior for fixed T as a function of a. Once again, the resulting curves show quite a remarkable agreement with empirical data, as illustrated in Fig. 7.7. Moreover, by fitting Eq. (7.7) to the simulated data, we can extract the values of σ_{ℓ} and λ , which are very close to the parameters originally used in the simulation, see Fig. 7.8.

By fitting $R_a(T)$ as a function of a for specific values of T, one can assess which term -A or B —is dominant, see Eq. (7.7). This is done by successively fitting only one term at a time, i.e., either setting B=0 and fitting A, or setting A=0 and fitting B. Our theoretical framework also predicts which term should dominate depending on whether $\lambda \neq 0$.

Not only does the model show good qualitative agreement with the data, but the fits presented in Fig. 7.8 are also remarkably convincing from a quantitative point of view. In particular, we observe a clear match between:

- the NVD case and fits using only the A-term (with negligible B),
- the VD case and fits where B dominates (with A negligible).

Moreover, the fits yield realistic estimates for both the input values of σ_{ℓ} and λ , further validating the consistency of the model.

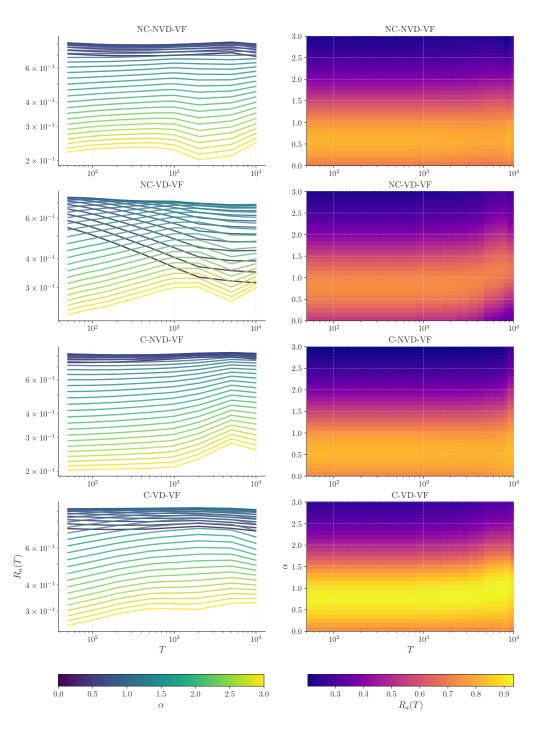


Figure 7.7: Simulations were done for $m=3, \sigma_\ell=1$ and $\lambda=1/(8\sigma_\ell)$. Left column: Evolution of the correlation for different values of a, showing the non monotonic behavior. Right column: Heatmap illustrating the distribution of correlation values within the (a,T) space, indicating that the correlation reaches its peaks for $a\approx 0.5-1$, regardless of the T values.

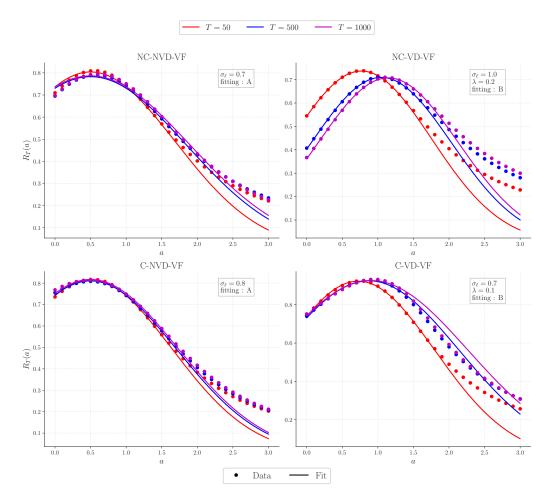


Figure 7.8: Fit of the correlation function $R_a(T)$ for several values of T. Since Eq. (7.7) is valid only for $a < a_c$, the fit is restricted to a < 1.5. The letters (A, B) indicate which term of Eq. (7.7) is being fitted. The empirical estimates of σ_{ℓ} and λ obtained through this procedure are remarkably close to the values used as input in the simulations: $\sigma_{\ell}^2 = 1$, $\lambda = 1/8$.

7.5 The puzzling effectiveness of proxy metaorders

In this final section, we address a central puzzle in the study of price impact: the surprising effectiveness of "proxy metaorders" introduced in [3]. Our algorithm constructs synthetic metaorders while preserving the exact trade history and sampling trades without replacement, two conditions that turn out to be essential. This algorithm originates from a study of metaorder impact using the TSE dataset, which includes real trading identifiers. A striking initial finding was

Chapter 7.

that randomly shuffling trading IDs and reconstructing synthetic metaorders still preserves the square-root impact law (SQL): see [2] section 4.2 for details.

However, one might argue that obtaining such a result relies on the prior knowledge of the original trading IDs. The shuffling process may preserve hidden information—such as the distribution of trading frequencies—which could, in turn, explain the impact function observed for the synthetic metaorders. Although appealing and somehow intuitive, this hypothesis was refuted in [3] through the construction of synthetic metaorders using public data. Yet the justification of the success of that method in reproducing the SQL remained somewhat mysterious.

The framework we introduce here allows one to justify further our proposal using purely simulated data. Although we have not yet been able to compute exactly the impact of proxy metaorders within our model, we believe that our numerical results are convincing enough to believe that the procedure proposed in Ref. [3] is warranted.

In Section 7.3, we introduced a detailed procedure for generating a dataset that closely approximates the ideal case (such as the TSE dataset), which provides trade-by-trade data along with metaorder identifiers across the entire market. Building on this, we conduct a numerical experiment where we pretend we do not know the mapping between trades and metaorders, and construct a proxy in the spirit of [3]. For the purpose of such an experiment, we assume no volume dependence, i.e, $\lambda = \lambda' = 0^{32}$. Each metaorder can thus be characterized by only three parameters: its size s drawn from a distribution $\Psi(s) \sim s^{-1-\mu}$, its execution rate which we choose to be the same for all metaorders $\tilde{\varphi} = \varphi$ and its average child order volume q, with $q \sim LN(m, \sigma_{\ell}^2)$.

The core challenge in designing a reliable proxy for metaorders lies in aggregating market orders in a way that statistically approximates the true (yet usually unobservable) matching between traders and trades. A natural method to reconstruct realistic metaorders from the observed order flow is to first separate buy and sell orders and then, for each list, iterate through the orders while performing the following: if an order is already part of an existing metaorder, we skip it; otherwise, we draw a size $s \sim \psi(s)$ and group the next s orders that occur within intervals of duration φ into a new metaorder. Since splitting and grouping orders can bring together orders with the same sign that were actually executed far apart in time, we introduce an inter-time threshold between two child orders. If the inter-time is above the threshold, we consider that the two child orders belong to different metaorders. This inter-time constraint proves essential for reproducing the SQL,

The proposed study and code can be readily extended to scenarios where $(\lambda, \lambda') \neq (0, 0)$ and $q \sim LN(m, \sigma_q)$ by separating buy and sell orders, binning the volume q, and applying the algorithm using the corresponding value of μ_q .

and corresponds to usual execution schemes where child orders tend to be relatively close to one another. A long pause in the execution schedule is tantamount to starting a new metaorder.

This procedure is summarized in Algorithm 4, where C is a constant, which we arbitrarily set to 4φ , as it provides satisfactory results, see Fig. 7.9, where we compare the numerical evaluation of the square-root law using the known exact matching between child orders and metaorders generated by our simulation (blue line) and the impact law estimated using proxy (or synthetic) metaorders (orange line). One sees that the agreement is almost perfect when Q/V_D is not too small, whereas the effective behaviour of the reconstructed impact becomes more linear. This is expected, since the start of short proxy metaorders have a higher probability to miss the start of "real" metaorders, for which the impact is most concave. We also confirm that the derivation of the prefactor Y of the SQL in [122] is correct, namely $I(Q) = Y \sigma \sqrt{Q/V_D}$. We believe that this additional quantitative validation is important, since the prefactor is usually less studied in the literature, although it remains of significant interest for the estimate of actual impact costs.

These simulation results therefore bolster the claim made in [3] that a realistic estimate of the impact of metaorders can be obtained using anonymous trade by trade data, provided the mapping function that generates proxy metaorders is chosen adequately. In fact as shown in [3] (Appendix), this mapping function, based on the theoretical framework developed here, also performs well on real data.

Algorithm 4 Generate Metaorder Identifiers with Time Threshold

```
1: function GENERATEMETAIDS(t\_execs, \varphi, \mu, s_{\text{max}})
         n \leftarrow \text{len}(t \ execs)
 2:
 3:
         ids \leftarrow zeros(n, dtype=int)
        id\_meta \leftarrow 1
 4:
        i \leftarrow 0
 5:
         while i < n do
 6:
             size \leftarrow \Psi(\mu, s_{\text{max}})
 7:
             count \leftarrow 0
 8:
             current\_time \leftarrow t\_execs[i]
 9:
             while count < size and i < n do
10:
11:
                 ids[i] \leftarrow id\_meta
                 count \leftarrow count + 1
12:
                 next\_time \leftarrow current\_time + \text{Exp}(\varphi)
13:
                 i\_next \leftarrow searchsorted(t\_execs, next\_time, left)
14:
                 if i\_next \ge n or (t\_execs[i\_next] - next\_time) > C/\varphi then
15:
                      break
16:
                 end if
17:
18:
                 i \leftarrow i\_next
19:
                 current\_time \leftarrow t\_execs[i]
             end while
20:
             id\_meta \leftarrow id\_meta + 1
21:
22:
             while i < n and ids[i] \neq 0 do
23:
                 i \leftarrow i + 1
24:
             end while
         end while
25:
        {f return}\ ids
26:
27: end function
```

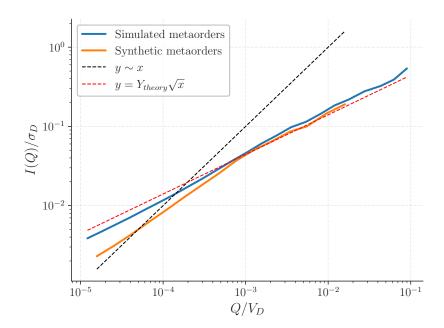


Figure 7.9: Comparison of the impact of simulated metaorders in the C-NVD-VF setup and synthetic metaorders generated using the metaorder proxy and constructed from simulated prices. For small Q/V_D , synthetic metaorders tend to have less concave, but after a crossover value around 10^{-3} , it nicely converges to the expected SQL, which is an input of our simulation. Note that we also recover the exact theoretical prefactor Y_{theory} computed in [122]. Both the simulation algorithm and the mapping function use $\varphi = 2.10^{-3}$, $\mu = 1.5$, m = 3 and $\sigma_{\ell}^2 = 1$. A total of 1,000,000 simulated metaorders were generated. We obtain similar results in others simulations cases

A crucial aspect of price impact lies in the decay post execution, which remains the subject of active debate within the community. We verify that the post-execution dynamics of synthetic metaorders closely match those of simulated metaorders. Since analyzing post-execution impact is subtle [44], we check here that the method produces realistic outputs. A more refined investigation, both theoretical and empirical, is left for future work. Nevertheless, we confirm that the decay is consistent with propagator theory, i.e. a power-law decay $t^{-\beta}$ with $\beta \in [0.1, 0.5]$

7.5.1 How the acceptance window drives mapping accuracy

Let's try to explain more theoretically why this mapping function works that well, and what is the role of the cutoff. In simulations, each metaorder i is generated

Chapter 7.

as a Poisson process with intensity φ :

$$X_k^{(i)} \sim \text{Exp}(\varphi), \qquad t_j^{(i)} = t_1^{(i)} + \sum_{k=1}^{j-1} X_k^{(i)}.$$

The propagator kernel depends on offsets

$$\Delta t_i = t_i - t_1$$

and the SQL holds then by construction.

In the reconstruction setting, we observe only the split order flow—i.e., individual buy or sell orders—which results from the superposition of multiple metaorders executed simultaneously. After observing a child at time t, the next candidate can come from:

- the same metaorder, after a gap $\sim \text{Exp}(\varphi)$,
- from the rest of the split flow

We denote by m the number of active metaorders of same sign. We readily obtain $m = \frac{\nu \mathbb{E}[s]}{2\varphi}$ - up to a correction factor when metaorder signs are autocorrelated. With it comes the rate of execution of the split order flow (say buy orders here): $\lambda^+ = \frac{\nu \mathbb{E}[s]}{2}$

Let us assume that we accept the next order if it occurs within a window w. Conditioning on at least one event in [0, w], the probability that this event is from the same metaorder is:

$$P_{\text{true}}(w) = \frac{\varphi(1 - e^{-\varphi w})}{\varphi(1 - e^{-\varphi w}) + \lambda^{+}(1 - e^{-\lambda^{+}w})}$$
(7.9)

As $w \to \infty$, $P_{\text{true}}(w) \to \frac{\varphi}{\varphi + \lambda^+} = \frac{1}{1+m}$. Thus, in the regime of low density $m \ll 1$ (ie a small number of metaorder is active at the same time), this mapping function allows one to reconstruct the real simulated metaorder.

Unfortunately, in real financial markets as in our simulation, m is typically larger than 1, as different investors execute their metaorder simultaneously. Thus, it is highly probable that synthetic metaorders integrate also child orders from other metaorders. Then, the generalized propagator for synthetic metaorders reads:

$$I(\hat{Q}) = q \sum_{j=1}^{\hat{s}} \left(\frac{t_j - t_1(t_j)}{\varphi} + n_0\right)^{-(1/2 - \beta)} \left(\frac{\tau_0}{t_{\hat{Q}} - t_j + \tau_0}\right)^{-\beta}$$
(7.10)

where $t_1(t_j)$ is the starting time of the simulated metaorder to which the child order executed at t_j belongs. As expected, the decay part is not impacted by the reconstruction principle, but a distortion of the SQL could arise from the mismatching $t_1(t_j)! = t_1$, as we aggregated child orders that may not come from metaorders initiated at the same time. Luckily, fine-tuning the threshold C has two main effects:

First, in (7.9), it increase the probability of selecting the correct next child order, and more importantly, it prevent $\Delta_j = t_j - t_1(t_j)$ to be too large. Indeed, even when linkage is correct, large gaps inside a metaorder distort impact. For s children drawn i.i.d. from $\text{Exp}(\varphi)$:

$$\mathbb{E}[\max_{1 \le i \le s-1} X_i] \sim \frac{\log s}{\varphi},$$

which grows unboundedly with s. These large offsets inflate the kernel argument

$$\left(\frac{\Delta t_j}{\varphi} + n_0\right)^{-\left(\frac{1}{2} - \beta\right)},$$

causing the last trades to contribute almost zero impact, breaking SQL at the metaorder level. Imposing a maximum gap $X_i < C/\varphi$ ensures that Δt_j remain within realistic values. Even if some reconstructed children are slightly misaligned with respect to their true metaorder start, compensation effects occur when averaging over a large number of metaorders. This self-averaging restores then the SQL observed in Fig. 7.9.

To fully comprehend this phenomenon, a more extended mathematical derivation is necessary.

7.6 Conclusion

This work extends and complements our previous theoretical paper on the subtle interplay between impact, order flow and volatility [122]. In that work, most of our predictions turned out to be in rather remarkable agreement with empirical observations, despite the simplifying mathematical approximations that we had to make. In the present paper, we show using numerical simulations that these approximations are actually quantitatively justified, which provides further support for the validity of our theoretical framework, and bolsters our conclusion that price volatility can be fully explained by the superposition of correlated metaorders that all impact prices, on average, as a square-root of executed volume. One of the most striking predictions of our model is the structure of the correlation between generalized order flow and returns, which is observed empirically and reproduced using our synthetic market generator.

Chapter 7.

Furthermore, we were able to construct proxy metaorders from simulated order flow that reproduce the square-root law of market impact —a law that has long been, and in some circles still is, attributed to information revelation; see e.g. [30, 71, 94]. Our model, on the other hand, makes the assumption that impact is purely mechanical and a result of the random dynamics of latent liquidity that creates a buffer for price moves, see [7, 74, 90]. The possibility of measuring the impact of metaorders from tape data (i.e. anonymized trades) was long thought to be impossible. However, Ref. [3] showed that a suitable mapping between market orders and proxy metaorders allows one to reconstruct many statistical features of real metaorders. We confirm that this is indeed the case within our purely synthetic market as well, lending further credence to our proposal [3].

Take Home Message

- We present an algorithm designed to create artificial markets in line with the theoretical framework established in the previous chapter.
- Our results successfully replicate the complex stylized facts emphasized by the theory, demonstrating a strong alignment between data, theory, and simulation and thus reinforcing the robustness of this framework.
- Building on our framework, we propose a more intuitive mapping function compared to the one introduced in Chapter 5. This advancement enables us to generate realistic metaorders on a *simulated price*.

Part III Market Stability

Chapter 8

Microstructure modes in the Limit Order Book: Signature of Marginal Instability

Tout doit tendre au bon sens: mais, pour y parvenir, Le chemin est glissant et pénible à tenir; Pour peu qu'on s'en écarte, aussitôt l'on se noie. La raison pour marcher n'a souvent qu'une voie.

Nicolas Boileau, L'Art poétique, Chant I (1674)

While prices are ultimately the outcome of supply and demand interactions, the high-frequency nature of order book activity makes it difficult to disentangle meaningful patterns from noise. In this chapter, we investigate the joint dynamics of order flow and price movements using an order-by-order dataset from the $HF_EUROSTOXX$. We introduce a coarse-graining method that reveals statistically meaningful patterns at the minute scale. Using Principal Component Analysis, we extract dominant "microstructure modes" that decompose market activity into symmetric and anti-symmetric components. We then calibrate a Vector Auto-Regressive (VAR) model to describe the dynamics of these modes, finding remarkably stable parameters and strong predictive power for symmetric liquidity patterns. As the number of lags increases, the model approaches marginal instability, echoing the persistent memory of order flow and hinting at the endogenous nature of liquidity crises.

From:
"Microstructure Modes":
Disentangling the Joint Dynamics of Prices & Order Flow
S. Elomari-Kessab, G. Maitrier, J.Bonart, JP. Bouchaud

Contents 8.1 Introduction 8.2 Variables of interest and intraday profile 171 8.2.2 8.2.3 8.3PCA Analysis I: Raw data PCA Analysis II: Binned data 177 A VAR model for flow dynamics 178 8.4.2 8.5 An attempt to model price impact 184 Conclusion and further discussions 188

8.1 Introduction

The micro-dynamics of asset prices is intricate, resulting from a subtle interplay between market orders, limit orders and cancellations happening at an amazingly fast pace in modern electronic order books. The mathematical description of the succession of these different events, the volume in the order book, and the occasional price changes when the queue at the best bid or ask empties out, is extremely difficult. This is due both to the high dimensionality of the problem, and to the presence of long-range correlations in the sign of the market/limit orders, which makes it necessary to have strong enough feedback loops. For instance, Zero Intelligence models [126, 127], where agents make decisions without any strategic reasoning or foresight, fail for exactly this reason at creating coherent sequences in time.

One possible avenue, which has led to interesting results recently, is to train generative neural networks on large datasets [97, 98], interpreting each event as a word and trying to guess the series of event following a given word history. Learning the underlying statistical structure of the order book dynamics would allow one to generate realistic synthetic limit order books. This would in turn offer valuable

opportunities to enhance market making strategies, or its dual problem: optimal execution. It would also allow one to simulate the counter-factual impact of additional orders, that are not present in the public tape, by understanding how the market digests such orders [88]. Indeed, inferring the impact of – say – market orders based only on the public tape is marred with conditioning problems.

Although some success of using the analogue of Large Language Models was reported [97, 128, 129], the prediction horizon for the order book dynamics appears to be limited to a few tens of events. But since such events happen at extremely high frequencies, the time horizon of these predictions is shorter that one second for electronic liquid markets, during which the price itself seldom changes. Although possibly useful for High Frequency Trading [97, 130], one would like to develop tools that account for the joint dynamics of prices and order flows on somewhat longer time scales, say minutes.

One of the main problems faced by "complete" models where all events are taken into account is that the high frequency dynamics of order books contains what one would like to call "jitter", i.e. orders that are placed and immediately cancelled, providing little information on the longer term fate of the order book. Another source of "jitter" are market orders that empty a queue at the best only to be immediately refilled by limit orders, creating high-frequency mid-point bounces.

Our main idea in this Chapter is to coarse-grain and simplify the dynamics in such a way that only "significant" price changes (more precisely defined below) are retained. The flow of market orders, limit orders and cancellations, both at the bid and at the ask, are aggregated between two price changes and used as the relevant dynamical variables we want to focus on and predict, together with the time elapsed between two price changes and the corresponding return itself. These variables define an 8-dimensional space on which we project, in a sense, the full joint dynamics of prices and order flow.

We then perform a Principal Component Analysis of the fluctuations, which defines "liquidity modes" that turn out to be stable in time and have a clean interpretation of market dynamics. This allows us to define a VAR model for predicting such modes one lag ahead, with a very significant \mathbb{R}^2 score.

One of our key findings is that one should actually distinguish between two natural coarse-graining procedures. The first one is to exclude price changes that are immediately reverted, and define other price changes as significant. However, we still see very strong mean-reversion (or "bounce") effects for the resulting price changes, that we call "raw" henceforth. We therefore define a second coarse-graining scale by aggregating N successive raw price changes, constructing what we will call "binned" returns, choosing N in such a way that the autocorrelation

of successive binned returns is below 0.01. On longer time scales, the series of price returns is thus closer to white noise, such that mechanical microstructure effects are smoothed out. For such binned data, our VAR model predicts flows with a substantial R^2 score ($\sim 25\%$) whereas, not unexpectedly, the prediction for returns is smaller but still significant, both in-sample and out-of-sample.

Both the "raw" scale and the "binned" scale are important for applications, but for different end users. The raw scale is presumably most useful for market makers and HFT, whereas the binned time scale is relevant for optimal execution and even, possibly, fast alpha signals. Our reduced model allows us to generate realistic time series of price changes and order flow. It also allows us to detect regime changes, when residuals with respect to the VAR model become anomalously high.

Interestingly, when our VAR model is extended to multi-lags, we detect clear signs related to known long memory effects, i.e. several activity directions correspond to eigenvalues tending to one and become marginally stable under the dynamics. A similar effect is known to occur when one fits linear Hawkes processes to financial data [59]: the only way to capture long memory is to bring the model close to instability [60, 131]. If taken at face value, the marginally stable eigenvectors of our VAR model would suggest incipient liquidity crises, a scenario advocated in various contexts, see e.g. [55, 64, 132–134]. An alternative interpretation of such marginal stability is the effect of changing activity levels across different periods, which is in a sense another manifestation of the long range correlations of the flows.

Finally, we can also use our model to simulate the impact of additional flows and see how far we can recover the various stylized facts reported in the literature, i.e. impact concavity and relaxation of impact after the trade is completed.

The outline of the Chapter is as follows. Section 8.2 introduces the variable of interest in our modeling. It describes the dataset, and the chosen pre-processing of the data. Section 8.3 suggests an analysis of microstructure modes based on Principal Component Analysis (PCA) of our data. In Section 8.4, we present the VAR model applied to our data, with an analysis of the stability of the resulting dynamics. Measures of the price impact under our model can be found in Section 8.5, and we conclude in Section 8.6.

8.2 Data presentation

The dataset used in this study consists of 545 days of the futures contract on EU-ROSTOXX from September 2016 to August 2019. The original data was obtained at the tick level, capturing detailed information about each price change. During

the analysis period, just under 4 million price changes were observed, with on average 7264 price changes per day.

It is noteworthy that EUROSTOXX is a liquid large-tick asset, with a spread almost always equal to one tick. This feature puts strong constraints on possible price changes: very often the mid-point mechanically bounces back because one of the best queues is immediately refilled after its depletion. Considering that these price changes represent microstructure noise, we filtered them out of the data and retained only "significant" price changes, defined as follows:

A significant price change corresponds to cases when the new bid corresponds to the old ask, or when the new ask corresponds to the old bid.

In other words, most spread-opening events correspond to a mid-point change of half a tick. If the following spread-closing event corresponds to another half-tick move in the *same* direction, we consider the price change to be significant. This definition allows us to remove some of the "jitter" that we deem insignificant in the dynamics that we want to capture, and to reduce the dimensionality of the problem by focusing on the total flux of orders between two successive price changes.

8.2.1 Variables of interest and intraday profile

For the n^{th} significant price change of the day, occurring at time t_n , we define the following variables:

 $\Delta t_n = t_n - t_{n-1}$: Time duration between the $(n-1)^{th}$ and n^{th} price changes $V_n^{\text{ex, a}}, V_n^{\text{ex, b}}$: Volume executed at the ask and bid, respectively, between t_{n-1} and t_n $V_n^{\text{lo, a}}, V_n^{\text{lo, b}}$: Volume posted to the first levels of the LOB at the bid and the ask, respectively $V_n^{\text{c, a}}, V_n^{\text{c, b}}$: Volume cancelled at the first levels of the LOB at the bid and the ask, respectively t_n : The return generated by the price change

All variables except returns are, by definition, positive. Returns can take positive and negative values, and are equal to ± 1 tick in most cases. For later use, we stack these 8 variables into the following 8-dimensional dynamical vector

$$\mathbf{X}_{n} = \left(\Delta t_{n}, V_{n}^{\text{lo, b}}, V_{n}^{\text{lo, a}}, V_{n}^{\text{c, b}}, V_{n}^{\text{c, a}}, V_{n}^{\text{ex, b}}, V_{n}^{\text{ex, a}}, r_{n}\right)$$
(8.1)

A consequence of focusing on "significant" price changes is that when the prices move up (twice, to be deemed significant), the first inserted volume at the new, higher bid is counted in $V_n^{\text{lo, b}}$ whereas the pre-existing volume at the new ask is not - and equivalently, when the prices moves down. In other words: queues that

move from a second-best to a best position are not considered as new placement flows.

Fig. 8.1 depicts the normalized average shape of the bid volume variables $V_n^{\star, b}$ throughout the day, binned in 1-minute intervals. The observed peak around 15:00 coincides with the opening of the US market. Since our modelling approach does not incorporate intraday volume patterns, all volumes are scaled by a smoothed average profile, fitted using two distinct exponential decay functions $A \exp(-t/\tau) + B$ with 3 parameters each: amplitude (A), on the decay time constant (τ) , and baseline (B), see table 8.1.

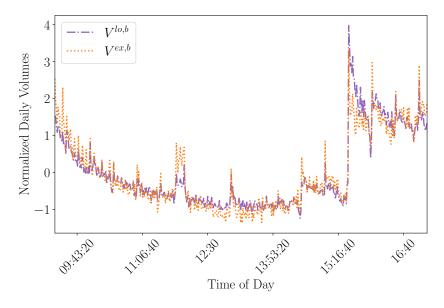


Figure 8.1: Normalized intraday profile of LOB placement and trade flows from the futures on EUROSTOXX data. The cancellation flows have similar profiles as the placement flows and are not presented in the figure for clarity. The activity level is high at the beginning of the day and decreases until a surge of activity at the open of the U.S. market. The intraday profile is the same for all activity flows and we characterize it by a unique set of parameters given in table 8.1.

Table 8.1: Exponential decay fit parameters for the average of the 6 normalized intraday flow profiles.

Parameter	Before 15:30	After 15:30
Amplitude (A)	2.20 ± 0.14	1.95 ± 0.42
Decay Time (τ)	50.0 ± 5.7 minutes	6.85 ± 2.4 minutes
Baseline (B)	1.79 ± 0.04	3.95 ± 0.07

8.2.2 A second coarse-graining

Even after filtering out price changes deemed not significant, the returns r_n still show very strong anti-correlations. Fig. 8.2 shows that the empirical correlation function $C_r(\ell) := \langle r_n r_{n+\ell} \rangle$ can be approximated as

$$C_r(\ell) \approx (-\gamma)^{\ell}; \qquad (\gamma \approx 0.8 < 1).$$
 (8.2)

These correlations only become small (< 0.01) beyond lag $\ell = 20$. To wit, strong microstructure effects still affect our "significant" price changes: the mid-price only becomes approximately diffusive for lags ≥ 20 .

In view of these persistent anti-correlations, we have introduced a second coarse-graining scale by further binning consecutive significant price changes into groups of 20. Throughout this Chapter, we will refer to our initial definition of significant price changes as "raw" and the aggregated (in batches of 20) price changes as "binned". Due to the very short time scales of the Raw Price Change data, one observes very many null flows $V_n^{\star, a}$ or $V_n^{\star, b}$ between t_n and t_{n+1} , an effect that completely disappears when the data is binned.

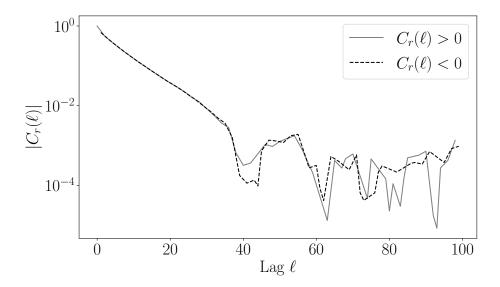


Figure 8.2: Absolute value of autocorrelation of returns. The alternation of positive and negative values in the autocorrelation indicates bounces in the return, which essentially disappear after binning 20 successive price changes.

8.2.3 Box-Cox transformation

Because of the strongly non-Gaussian nature of the variables V_n and Δt_n , even after binning, we start by applying a "Box-Cox" transformation $f(x; \lambda)$ to the binned variables, with

$$f(x;\lambda) := \begin{cases} \frac{(ax)^{\lambda} - 1}{\lambda}, & \text{if } \lambda \neq 0\\ \log(ax), & \text{if } \lambda = 0 \end{cases}.$$
 (8.3)

and a parameter λ possibly different for the volume variables λ_v (for chosen to be the same for all such variables) and λ_t for the time variable. These parameters are chosen to maximise the likelihood of the Gaussian distribution of the transformed variables, which yields $\lambda_v = 0.20$ and $\lambda_t = 0.14$. The scale parameter a can be set to unity without loss of generality.

We will henceforth work with a series of 8-dimensional vectors \mathbf{T}_n defined as:

$$\mathbf{T}_{n} = \left(f\left(\Delta t_{n}; \lambda_{\Delta t}\right), f\left(V_{n}^{\text{lo, b}}; \lambda_{f}\right), f\left(V_{n}^{\text{lo, a}}; \lambda_{f}\right), f\left(V_{n}^{\text{c, b}}; \lambda_{f}\right), f\left(V_{n}^{\text{c, b}}; \lambda_{f}\right), f\left(V_{n}^{\text{c, a}}; \lambda_{f}\right), f\left(V_{n}^{\text{ex, b}}; \lambda_{f}\right), f\left(V_{n}^{\text{ex, a}}; \lambda_{f}\right), r_{n} \right).$$

$$(8.4)$$

We further normalize the Binned Price Change data using a moving window spanning the days preceding the day of interest. Let w be the width of the time window

used for the computation of the means and the scales of the variables, with w = 20. Let d be a day in the data set, and N_k the number of observed price changes in day k. We write \mathbf{T}_n^d to indicate that the vector is observed at day d and define a causal local mean and the scale as follows:

$$\mu_j^d = \frac{1}{w} \sum_{k=d-w}^{d-1} \frac{1}{N_k} \sum_{n=1}^{N_k} (\mathbf{T}_n^k)_j, \quad j = 1, 2, \dots, 8,$$
(8.5)

$$(\sigma_j^d)^2 = \frac{1}{w} \sum_{k=d-w}^{d-1} \frac{1}{N_k} \sum_{n=1}^{N_k} \left((\mathbf{T}_n^k)_j - \mu_j^k \right)^2, \quad j = 1, 2, \dots, 8,$$
 (8.6)

with which we normalize each component of the T_n vectors as:

$$T_n^{\prime d} = \frac{T_n^d - \mu^d}{\sigma^d}. (8.7)$$

8.3 Microstructure modes

As expected intuitively, the volume V_n and time Δt_n variables are strongly correlated. For example, a large flux of market orders might trigger more limit orders and vice-versa. It is thus natural to use a Principal Component Analysis (PCA) to understand the structure of these (same bin) correlations, and define a set of uncorrelated principal components. These vectors, ordered by their associated eigenvalues, represent the dominant microstructure modes of the market. It turns out that all these modes exhibit near-perfect bid-ask symmetry (or anti-symmetry), especially when computed using a large number of days. Since there is no reason for this symmetry to be broken at high frequencies, we manually removed all remaining spurious bid-ask asymmetry in the results presented below. Note that the PCA analysis is always performed on the Box-Cox transformed variables T'_n , with the averaging window w chosen to be 20 days.

8.3.1 PCA Analysis I: Raw data

The eigenvectors decomposition of the raw data is given in Fig. 8.3, the corresponding eigenvalues λ_{α} ranging from $\lambda_1 = 4.07$ to $\lambda_8 = 0.02$, with $\sum_{\alpha=1}^{8} \lambda_{\alpha} = 8$ from the normalisation of the covariance.

Each eigenmode U_{α} has a rather intuitive and transparent interpretation, on which we comment below. Three of them are bid/ask symmetric, four are bid/ask anti-symmetric and the last one only contains duration, which appears to be independent variable at such high frequencies. Note that the sign of these eigenvectors is arbitrary; each direction is equally explored by the dynamics, with an intensity given by the square root of the corresponding eigenvalue.

- Mode 1 only contains volumes, with all coefficients positive. This represents an increase (or decrease) of general activity in the order book, with more (or less) market orders, limit orders and cancellations. It represents 51 % of the total variance.
- Mode 2 mixes market order imbalance with the contemporaneous return.
 As expected, more executions at the ask lead to a positive return and vice versa.
- Mode 3 is a pure duration mode.
- Mode 4 is bid/ask symmetric and describes situations where the aggressive flow becomes more active, whereas the passive flows (limit orders, cancellations) slows down or vice-versa.
- Mode 5 is anti-symmetric: more market orders at the ask than at the bid (and slightly less limit orders at the ask than at the bid), but resulting to a negative return, opposite to Mode 2. This counter-intuitive result is in this case due to the initial imbalance in the size of the queues. With the sign convention here, the bid side is less populated than the ask side, indicating net sell pressure overall. Still, higher liquidity at the ask attracts more buy market orders, explaining the excess of market orders at the ask.
- Mode 6, 7 and 8 are liquidity modes, since market order activity is absent from these directions. These modes represent 6.25 % of the total variance. Mode 6 and 7 and bid/ask anti-symmetric, and Mode 8 is symmetric. Mode 7 corresponds to a growing imbalance of the available liquidity at the bid and at the ask, since we see more limit orders and less cancellations at the ask and less limit orders and more cancellations at the bid (or vice-versa). Mode 8 has a very small intensity, and corresponds to a simultaneous loss (or increase) of liquidity on both sides of the book.

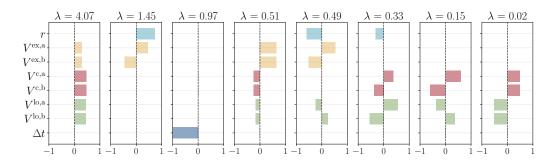


Figure 8.3: Normalized eigenvectors \mathbf{U}_{α} of the PCA decomposition (Raw data). For clarity purposes, the amplitudes lower than 0.15 (corresponding to weights less than $0.15^2 \approx 2\%$) have been set to zero. The directions should be interpreted as Box-Cox transforms of the original directions (except return r).

8.3.2 PCA Analysis II: Binned data

We now conduct exactly the same PCA analysis but now for binned data, aggregating volume flows and returns across 20 successive price changes. The emerging eigenmodes have very much the same structure as for the raw data: the PCA yield two categories of modes, one capturing symmetric activity between the bid and ask, and the other anti-symmetric activity and non-zero price changes.

Mode 1 again correspond to a global rise (or decline) of activity and mode 2 to a market order imbalance leading to a price change in the same direction as the imbalance. The total weight of these two modes $\lambda_1 + \lambda_2$ now reaches ≈ 6.80 , i.e. 85 % of the total variance, compared to 69 % for the raw data. Modes 3 and 4 are essentially the same as for raw data, apart from a permutation of their rank. The exact same thing happens for modes 5 and 6, and again for modes 7 and 8. Mode 4 (ex mode 3 for raw data) now associates shorter time duration with more volume added and cancelled in the order book. Interestingly, bid-ask symmetric fluctuations capture 82 % of the total variance, leaving only 18 % of the variance to asymmetric, price changing fluctuations.

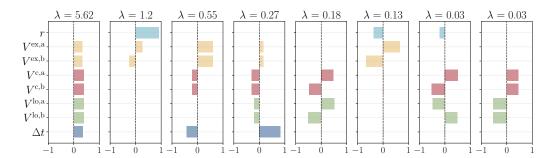


Figure 8.4: Normalized eigenvectors of the PCA decomposition (Binned data). For clarity purposes, the amplitudes lower than 0.15 (corresponding to weights less than $0.15^2 \approx 2\%$) have been set to zero. We observe 4 symmetric modes (1, 3, 4 and 8) and 4 anti-symmetric modes (2, 5, 6 and 7). The directions should be interpreted as Box-Cox transforms of the original directions (except return r).

8.4 A VAR model for flow dynamics

In this section, we present the mathematical framework underlying our modeling approach. We adopt a Vector Autoregression (VAR) model to capture the dynamic relationships among the variables associated with each price change. Regressing on the projection of the data onto eigenvectors rather than directly on the data itself helps handling collinearity issues and eases the interpretability of the results.

We will be interested in understanding the evolution of \mathbf{X}_n defined in Eq. (8.1), for the binned data. (For the raw data, the large fraction of zero entries would require a specific treatment, following for example [135–137]. We leave this for later investigations). In order to do so, we transform the data using Box-Cox and set up an Auto-Regressive Vector Model in the space of the 8 principal components (or eigenmodes) described in the previous section. For every n, the Box-Coxed vector \mathbf{X}_n is projected onto the j^{th} eigenmode \mathbf{U}_{α} , and the resulting projection is further demeaned and normalized to have unit variance, finally defining an 8-vector in the eigenmode space \mathbf{Y}_n .

The p-lag VAR model is then specified by the following evolution equation

$$\mathbf{Y}_n = \mathbf{\Phi}_1 \mathbf{Y}_{n-1} + \mathbf{\Phi}_2 \mathbf{Y}_{n-2} + \ldots + \mathbf{\Phi}_p \mathbf{Y}_{n-p} + \boldsymbol{\epsilon}_n, \tag{8.8}$$

where ϵ_n represents a vector of white noise innovations and Φ_k are 8 × 8 transition matrices capturing the inter-dependencies and temporal dynamics in the eigenmode space. The VAR model is calibrated using standard regression methods, except that we add by hand an additional constraint that the model has

to respect the bid-ask symmetry. This means that all coefficients $(\Phi_k)_{\alpha\beta}$ relating bid-ask symmetric modes $(\alpha = 1, 3, 4, 8)$ to bid-ask anti-symmetric modes $(\beta = 2, 5, 6, 7)$ must be zero. Without this constraint, all symmetry-breaking coefficients are found to be very small anyway.

8.4.1 1-lag VAR model

We first focus on the p = 1 lag VAR model:

$$\mathbf{Y}_n = \mathbf{\Phi}_1 \mathbf{Y}_{n-1} + \boldsymbol{\epsilon}_n. \tag{8.9}$$

The transition matrix Φ_1 is presented in the table 8.2. The most significant elements, i.e., such that $|(\Phi_1)_{\alpha\beta}| > 0.1$, are highlighted in bold and correspond mostly to diagonal elements (except 22 and 66). However, a better description of the transition matrix is in terms of its eigenvalues and eigenvectors. 6 eigenvectors correspond to real eigenvalues, 5 positive and one negative, and 2 eigenvectors correspond to a pair of complex conjugate eigenvalues, with a very small modulus. The five eigenvectors with largest norm are shown in Fig. 8.5. The fact that all eigenvalues within the unit circle means that the lag-1 VAR model is stable, with fluctuations dampening instead of getting amplified. Notice that the top eigenvalue is equal to 0.68 and corresponds to a symmetric cancellation mode, mostly reflecting the activity of market makers.

The second mode, with eigenvalue 0.56, is also symmetric and corresponds to more limit orders, less market orders and less inter price change time, or vice-versa. The largest anti-symmetric mode has eigenvalue $\lambda_5 = -0.23$ and is the imbalance level for all flows, which is seen to be mean-reverting (since $\lambda_5 < 0$).

Mode	1S	2A	3S	4S	5A	6A	7A	8S
1S	0.56	0.00	-0.02	-0.03	0.00	0.00	0.00	0.09
2A	0.00	-0.06	0.00	0.00	-0.00	0.04	-0.09	0.00
3S	0.05	0.00	0.53	0.00	0.00	0.00	0.00	-0.09
4S	-0.06	0.00	-0.02	0.59	0.00	0.00	0.00	-0.03
5A	0.00	0.01	0.00	0.00	0.15	-0.02	-0.04	0.00
6A	0.00	0.08	0.00	0.00	-0.09	-0.05	0.09	0.00
7A	0.00	-0.08	0.00	0.00	-0.08	0.09	-0.11	0.00
8S	0.12	0.00	-0.05	-0.04	0.00	0.00	0.00	0.52

Table 8.2: The transition matrix for microstructure modes, where values exceeding a significance threshold of 0.05 in the corresponding p-value have been set to zero. Columns correspond to input modes from time n-1, rows to predicted modes at time n. Symbol S (or A) refers to the bid-ask symmetry (anti-symmetry) of the modes.

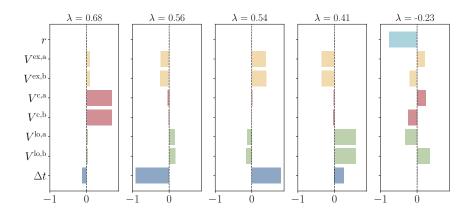


Figure 8.5: 5 eigenvectors with the largest eigenvalue norm from the decomposition of the transition matrix. The 4 first eigenvectors have positive eigenvalues and describe symmetric scenarios in the bid and the ask, the fifth one is anti-symmetric and mean-reverting.

The success of the lag-1 VAR model can be quantified in terms of the predictive R^2 scores, presented in table 8.3, both across modes of the transformed variables \mathbf{Y}_n and for the original variables \mathbf{X}_n . Note that, as expected, R^2 scores are much higher ($\sim 0.28-0.32$) for symmetric modes, which carry no information on returns, than for anti-symmetric modes ($\sim 0.01-0.03$). However, the in-sample R^2 scores

and out-of-sample scores are close, highlighting the fact that the predictive value of the VAR model is statistically significant. We have used the first 465 days of the data for model calibration and computing the in-sample scores. The remaining 80 days are allocated for computing the out-of-sample scores.

Mode	1S	2A	3S	4S	5A	6A	7A	8S
In Sample (%)	32.4	1.2	29.1	35.1	2.43	2.39	3.07	28.8
Out Of Sample (%)	28.0	1.11	21.8	36.1	4.33	1.68	2.24	32.5
Variable	Δt	V ^{lo, b}	V ^{lo, a}	<i>V</i> ^{c, b}	V ^{c, a}	V ^{ex, b}	V ^{ex, a}	r
Variable In Sample (%)	Δt 21.3	V ^{lo, b} 29.8	V ^{lo, a} 29.8	V ^{c, b} 36.4	V ^{c, a} 36.0	V ^{ex, b} 25.3	V ^{ex, a} 24.8	r 1.60

Table 8.3: R^2 scores in % both in mode space (top) and in the original space (bottom). Symbol S (or A) refer to the bid-ask symmetry (anti-symmetry) of the modes.

The R^2 score of 1.6% for return r is of particular interest. It is in particular significantly higher than the score of 0.49% obtained when predicting returns using the past return as the only feature. This shows that flow variables add useful predictive power to the return variable.

8.4.2 Multi-lag VAR model

In this subsection, we extend our modeling approach to the multi-lag Vector Autoregression VAR(p) model, specified by Eq. (8.8) with p>1. Interestingly, adding more lags reduces auto- correlation of residuals and increases the out of sample R^2 score of all the modes, by $\sim 25\%$ both for the symmetric and antisymmetric ones when p increases from 1 to 10 – see tables 8.4 and 8.5.

Lags	1	2	3	4	5	6	7	8	9	10
In Sample S (%)	31.4	35.8	37.3	38.1	38.5	38.8	39.0	39.2	39.3	39.4
Out Of Sample S (%)	29.6	34.3	36.2	37.0	37.3	37.6	37.9	38.1	38.3	38.4
In Sample A (%)	2.29	2.56	2.69	2.77	2.86	2.94	3.00	3.05	3.10	3.15
Out Of Sample A (%)	2.35	2.56	2.65	2.73	2.79	2.85	2.92	2.99	3.03	3.04

Table 8.4: Average R^2 scores for symmetric (S) and asymmetric (A) modes in-sample and out-of-sample for different number of lags p.

Mode	1S	2A	3S	4 S	5A	6A	7A	8S
In Sample (%)	36.0	1.46	35.5	43.5	3.47	2.29	4.00	36.3
Out Of Sample (%)	36.7	1.30	27.5	45.4	5.53	2.22	2.89	42.6
Variable	Δt	V ^{lo, b}	V ^{lo, a}	<i>V</i> ^{c, b}	V ^{c, a}	V ^{ex, b}	V ^{ex, a}	r
Variable In Sample (%)	Δt 30.4	V ^{lo, b} 37.6	V ^{lo, a} 37.9	V ^{c, b} 44.0	V ^{c, a} 43.8	V ^{ex, b} 32.1	V ^{ex, a} 31.5	1.87

Table 8.5: R^2 scores using the VAR(8) in % both in mode space (top) and in the original space (bottom). Note that the out-of-sample R^2 score of the returns increases from 1.46 for p=1 to 1.7 for p=8.

Another interesting question is whether adding memory to the system makes it less stable. In order to discuss this point, let us look for a vector \mathbf{Z} such that at long times the p-VAR model in the absence of innovations would yield

$$\mathbf{Y}_n \approx_{n\gg 1} \gamma^n \mathbf{Z}$$
.

Injecting in eq. (8.8) and dividing by γ^n , we find the following condition:

$$\mathbf{Z} = \mathbb{M}_p \mathbf{Z}, \qquad \mathbb{M}_p(\gamma) := \left[\sum_{k=1}^p \gamma^{-k} \mathbf{\Phi}_k \right].$$
 (8.10)

In other words, one should look for a value of γ such that the matrix $\mathbb{M}_p(\gamma)$ has one eigenvalue exactly equal to unity, the corresponding eigenvector defining **Z**. The least stable direction of the p-VAR model is associated with the largest possible value of $|\gamma|$ (with γ possibly complex).

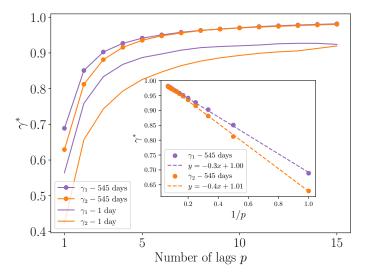


Figure 8.6: The two largest values of $\gamma_{1,2}(p)$, such as the matrix $\mathbb{M}_p(\gamma)$ has an eigenvalue equal to unity, as a function of lag p. The plot compares methods for normalizing the data: one using all 545 available days and the other where each day is normalized independently. Inset: same results, plotted as a function of 1/p showing a near perfect linear behaviour extrapolating to unity when $p \to \infty$.

Quite interestingly, we observe in Fig. 8.6 that both $\gamma_1(p)$ and $\gamma_2(p)$ can be fitted as $1 - C_{1,2}/p$ and therefore appear converge to unity as the number of lags increases. This means that the dominant eigenvectors, shown in Fig. 8.7, become more and more persistent as we increase the number of lags p. This suggests that the flow dynamics is in fact marginally stable, which is in line with the well-known stylized fact that order flow has power law, long memory correlations [88], corresponding to a unit root within a VAR description, or to marginal stability within a Hawkes process description [131]. Marginal stability could however result from the inadequacy of the VAR model to represent the data, since the only way to represent long memory correlations within a VAR framework is to have unit roots.

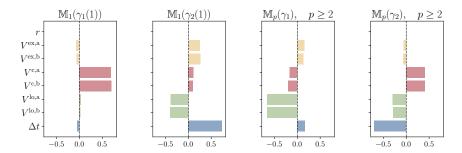


Figure 8.7: Dominant eigenvectors of $\mathbb{M}_p(\gamma_{1,2})$. For p=1 we recover eigenvectors from Fig. 8.5. When $p \geq 2$, all dominant eigenvectors are essentially independent of p and are associated with important liquidity fluctuations. For example γ_2 describes a persistent mode with less order placements and more cancellations, which can lead to liquidity crises.

Dominant eigenvectors are identical for all $p \geq 2$ and describe liquidity fluctuations. The mode associated with $\gamma_1(p)$ predicts less (or more) placements than usual. The second one, with rate $\gamma_2(p)$, describes a persistent mode with less order placements and more cancellations, which can lead to liquidity crises, as argued in [64]. Even if $\gamma_2(p)$ is below unity, the system appears to be very close to this stability boundary, and therefore be prone to endogenous liquidity crisis. In this context, recall that we chose a particularly stable, large tick contract (the EUROSTOXX); it would be interesting to perform the same analysis with small tick single name stocks.

8.5 An attempt to model price impact

In this section, we are interested in understanding the impact of the trading of one agent on the market and its future states within the VAR framework established above. For the rest of the section, we focus on the impact of perturbations of the flows at the ask without any loss of generality since the regression matrix is symmetric between the bid and the ask.

A phenomenon commonly studied in the literature is price impact [49, 127, 138], or by how much a trader modifies the price of an asset by buying or selling it. This metric is crucial for practitioners, but also from an academic point of view. Price impact exhibits interesting theoretical properties, such as the so-called square root law (for a review, see [49, 127]).

In principle, the mechanical impact of market orders (i.e. the part that is inde-

pendent of any information motivating the trades) is defined as [49, 127]

$$\mathcal{I}(\ell) := \mathbb{E}[m_{t+\ell} - m_t | \operatorname{exec}_t] - \mathbb{E}[m_{t+\ell} - m_t | \operatorname{no-exec}_t], \tag{8.11}$$

where m(t) is the mid-price at time t, when a market order is executed. In other words, one should compare the price change between time t and $t+\ell$ with and without order execution. Of course, such a measurement is impossible, since these two states of the market are mutually exclusive. Therefore, in practice one assumes that for short enough time scales, market orders issued by slow traders have little short term predictability such that the second term in Eq. (8.11) is negligible. Hence the observable impact is defined as

$$\mathcal{I}^{\text{obs}}(\ell|exec_t) := \mathbb{E}[m_{t+\ell} - m_t]. \tag{8.12}$$

The whole idea of constructing a faithful generating model for prices and order flow is to be able to perform numerically the "do-operation" [49] described in Eq. (8.11).

We have performed such a numerical experiment using the VAR model calibrated above on binned data – which, we recall, aggregates together 20 successive significant price changes. The procedure is as follows: we add to the observed flow of market orders at the ask a specific quantity corresponding to our extra buyer, between coarse-grained time n and time n + k. At each time step, the instantaneous impact is calculated using the average impact curves obtained in [37], that are reproduced in the B.

However, there is a subtlety related to the execution flows predicted by the VAR model. Rotating the matrix Φ_1 into real flows' space, we obtain the matrix shown in table 8.6.

Variables	Δt	$V^{ m lo,\ b}$	V ^{lo, a}	$V^{ m c,\ b}$	V ^{c, a}	$V^{ m ex,b}$	V ^{ex, a}	r
Δt	0.56	0.05	-0.05	-0.00	0.00	-0.02	-0.03	-0.0
$V^{ m lo,\;b}$	-0.02	0.21	0.22	0.06	-0.05	0.02	-0.01	0.06
$V^{\mathrm{lo, a}}$	-0.02	0.22	0.21	-0.05	0.06	-0.01	0.02	-0.05
$V^{\mathrm{c, b}}$	-0.00	0.08	-0.06	0.39	0.28	-0.00	0.02	-0.04
$V^{\mathrm{c, a}}$	-0.00	-0.06	-0.08	0.28	0.39	0.06	-0.00	0.04
$V^{\mathrm{ex, b}}$	0.02	0.02	0.04	-0.02	0.04	0.26	0.28	-0.05
$V^{ m ex,\ a}$	0.02	0.04	0.02	0.04	-0.02	0.28	0.26	0.05
r	-0.00	0.06	-0.05	-0.06	0.04	-0.02	0.02	-0.14

Table 8.6: An approximation of the transition matrix in the real flows' space obtained as a rotation of Φ_1 back to the real variables space. Columns correspond to input variables at time step n-1, rows correspond to an estimation of the variables at time step n. Note: This is not strictly speaking a transition matrix because of the non-linear Box-Cox operation.

This matrix reveals that an increase in the market order flow at the ask is most likely followed at the next time step by an increase in market order flows in both the ask and the bid, with slightly higher values observed at the opposite side. The succession of market orders at the same side is a manifestation of the well-known long range correlation of the flows in the market [88], which is primarily due to metaorder splitting, with very little contribution from herding [119].

Our model has been trained on real-world price and flow data, whose causal structure includes but cannot be reduced to a perturbation-response mechanism. Single market participants do not act in isolation and they may, through complex trading strategies, influence the market dynamics on long time scales, and even cross-sectionally. To the extent that the exogenous perturbation, whose impact we wish to simulate, is not representative of the average market participants' trading schedule, the model cannot fully distinguish whether correlations are due to market response or individual complex trading strategies. In order to model consistently the impact of a specific exogenous metaorder, we must avoid double counting such contributions. Thus, as an approximation, within our simulation framework we disregard subsequent execution orders predicted by the model on the same side and only take into account induced effects. The perturbed flows and returns are then propagated forward in time using the VAR model. The total price impact is then obtained by subtracting the unperturbed observed price trajectory and averaging over time. However, it is important to note that this approach is, even on a conceptual level, an approximation whose accuracy is difficult



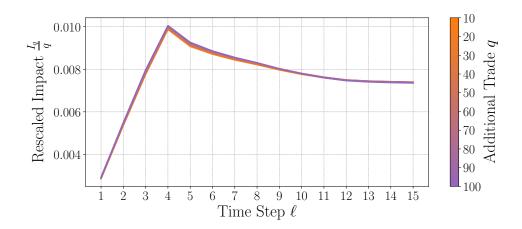


Figure 8.8: Simulations of impact of metaorders of length k=4. For the first 4 time steps, a market order of size q is added to the observed execution flow. Starting from the 5^{th} time step the market is no longer perturbed. We rescale the measured impact by the size of the added trade flow.

Empirically, as mentioned above, impact is strongly concave, and shows a square-root dependence both in time (within a metaorder) and in total size (at peak impact), see [127], chapter 12. Furthermore, such an impact strongly mean-reverts at the end of a metaorder. Our methodology, however, yields impact behaviour that differs in notable ways. What we observe in Fig. 8.8 is that the generated impact is only slightly concave within the metaorder, and then decays back down once the metaorder is completed. Despite this, the peak impact is linear in the size of the metaorder, contrarily to the concave behavior in observed market data. This linear shape is in fact expected within our perturbation approach, where the added trade flow is small enough to be absorbed by the market, leading to a linear behavior of the non-linear Box-Cox transformation. Furthermore, the level of reversion of the price between the moment we stop perturbing the market and when the price stabilizes is around 75% of the peak impact, whereas impact decay is much steeper in real data, with a significantly lower plateau value [44].

Note that Fig.8.8, as well as the overall approach, was independently reproduced and validated in [21]. This reinforces the conclusion that capturing metaorder impact using linear models is nearly impossible. This very limitation motivates the generalized propagator framework and underpins the work developed throughout PartII.

The conclusion of this section is that although our VAR framework offers a good benchmark for modelling the impact of metaorders, a crucial element appears to be missing since the strongly concave, mean-reverting nature impact is missed. We conjecture that such a missing element is an explicit reference to recent price changes, in a way to incorporate the idea of asymmetric latent liquidity, as argued in [15, 127]. Additionally, due to the use of the impact curves from [37] for computing the instantaneous return, the return produced by our model is diluted in scale. To address these limitations, future work could benefit from exploring models with enhanced non-linearity, such as neural networks.

8.6 Conclusion and further discussions

The frantic and noisy order book dynamics at the highest frequency hamper modelling attempts based on order by order activity. In this work, we have devised a specific coarse-graining procedure to extract meaningful information from such erratic flow data. First, in order to remove "flickering" bid-ask bounce noise, we have proposed a definition of significant price changes, and defined the flow variables of interest as aggregates of market orders, limit orders and cancellations between two such significant price changes.

However, we have found it necessary to introduce a second coarse-graining time scale in order to (i) smooth out strong price mean-reversion that survives until ~ 20 significant price changes and (ii) eliminate the large quantity of zeros in the flow variables that make linear analysis difficult to interpret.

One of our most interesting novel result is the appearance of what we called "microstructure modes", i.e. principal components of the joint, coarse-grained dynamics of price and order flow. These modes are extremely stable over time and all have an intuitive interpretation. They fall into two categories: bid-ask symmetric and bid-ask anti-symmetric. The first category describes, for example, an increase/decrease of cancellations and a decrease/increase of limit orders on both sides of the book simultaneously, associated to the dynamics of liquidity. The second category describes, for example, an increase of market orders at the ask and a decrease of market orders at the bid, associated to a positive price return.

Using these microstructure modes as inputs, we built and calibrated a multi-lag VAR model that captures their dynamics. The model is stable in time and leads to high R^2 scores $\sim 30-40\%$ for symmetric modes and, as expected, lower but significant R^2 scores $\sim 2-3\%$ for anti-symmetric (directional) modes. Non-linear, neural network models that take our microstructure modes as features should improve further the quality of the prediction.

We have found that the VAR model becomes marginally stable as the number of lags increases. This reflects the well-known long memory nature of the order flow in financial markets. The analysis of the flow directions that become unstable gives further credence to the "endogenous liquidity crisis" scenario suggested in [55, 64, 132–134].

Finally, we have used our VAR formalism to measure the impact of metaorders on the price. Although we observe some price mean-reversion at the end of the metaorder, similar to real data, we failed to reproduce the concave square-root dependence of impact on time and volume. We conjectured that an explicit conditioning of the VAR transition matrix on the recent returns is needed to capture "latent liquidity" effects that are thought to be at the origin of impact concavity [15, 127].

When working on the "raw", unbinned data, we were confronted with the fact that at short time scales, most of the observed flow volumes are null, making our linear VAR model unsuitable. One could address this problem using recent statistical techniques [135–137], or using more complex neural network architectures combining recurrent neural networks and attention techniques. It would also be interesting to revisit the price impact problem within this framework.

Take Home Message

- We introduced a coarse-graining procedure to extract meaningful signals from noisy, high-frequency order book data.
- Our analysis revealed robust and interpretable "microstructure modes" —principal components that capture the joint dynamics of price and order flow, split into bid-ask symmetric (liquidity) and antisymmetric (directional) patterns.
- A VAR model built on these modes successfully captures their temporal structure, with high predictive power for symmetric modes.
- The marginal instability of the VAR model as lags increase reflects the long-memory nature of order flow and supports the endogenous liquidity crisis scenario.
- This framework fails to capture the square-root impact law, suggesting the need for nonlinear dynamics or return-conditioned transitions.

Chapter 9

A generalized Santa Fe-like model to study liquidity crisis in the Limit Order Book

Done properly, computer simulation represents a kind of "telescope for the mind", multiplying human powers of analysis and insight just as a telescope does our powers of vision.

Mark Buchanan

In the previous chapter, we showed that precursor signs of endogenous liquidity crises can already be observed at the microstructure level, even at high frequency. Motivated by this further evidence challenging market stability, we now turn to a Santa-Fe inspired model —an extension of the original zero-intelligence agent-based framework, augmented with several forms of feedback. With these feedback mechanisms in place, we find that the order book dynamics undergo second-order phase transitions. Crucially, the system can shift from a stable to an unstable regime purely as a result of internal dynamics —depending on the degree of endogeneity and the memory of agents —without any external shocks. This provides yet another rejection of the Efficient Market Hypothesis.

From: *In preparation* **G. Maitrier**, G. Loeper, M. Benzaquen, JP. Bouchaud

Contents	
9.1	Introduction
9.2	Motivation: The initial Santa Fe model with feeedback 194
9.3	A 4 degree of freedom Santa Fe model 195
9.4	Investigating phase transitions in the 4DF Sante Fe model
	9.4.1 Stability maps
	9.4.2 Finite-size scaling
	9.4.3 Overview of exponents values 199
	9.4.4 Study of spread explosion
9.5	Analytical Determination of the Stability Boundary $$. 203
9.6	Conclusion

9.1 Introduction

Among its key predictions, the EMH asserts that markets operate in equilibrium, with prices fully reflecting all available information. Price movements are assumed to result from external news, and any deviations from equilibrium are expected to be promptly corrected. Within this framework, crises—such as crashes or liquidity dry-ups—are interpreted as rare, exogenously triggered anomalies.

However, as discussed in Chapter 3, even a brief look at financial market data challenges this view. Markets display recurring episodes of excessive volatility and heavy-tailed return distributions that are hard to reconcile with the EMH. At the macroscopic level, this manifests as flash crashes; at the microstructural level, it appears as abrupt price jumps occurring in the absence of any identifiable news.

These phenomena point to two key conclusions. First, the EMH likely incorrect—as we' ve been trying to convince the reader throughout this thesis, with increasing insistence. Second, market instabilities need not be triggered by external shocks, they can emerge spontaneously from internal feedback mechanisms inherent to market dynamic. Another striking particularity that caught the eye when examining financial prices is their similarity to turbulent processes well-known in statistical physics. As Mandelbrot [24] first pointed out, prices exhibit intermittency, scale invariance, and closely resemble multifractal processes. In physics, fractal properties often serve as indicators of phase transitions—critical points

where a system undergoes a fundamental change. These transitions are typically characterized, at criticality, by scale invariance, self-similarity, and power-law distributions, all hallmarks of fractal structures.

The occurrence of large price jumps and the striking statistical similarities between financial time series and multifractal processes motivate a deeper investigation of the limit order book as a complex, interacting system—where orders behave like particles. While a full theoretical description remains elusive, agent-based models (ABMs) offer a powerful framework to explore the emergence of such phenomena from simple agent behavior. Among these, the zero-intelligence model [126] has been a foundational starting point, successfully reproducing several stylized facts. However, it falls short when it comes to explaining extreme events like liquidity crises or abrupt price dislocations.

A significant step forward was made in [64], where the authors introduced a feed-back mechanism that makes otherwise random agents react to past price trends. This single addition was sufficient to induce a second-order phase transition in the order book, leading to endogenous liquidity crises —without any need for external shocks. As in physical systems, this transition is governed by a small number of parameters: here, the degree of endogeneity and the memory span over which agents integrate past price movements. Near criticality, key observables such as the spread and liquidity display power-law behavior, with associated critical exponents extracted using finite-size scaling techniques.

Building on this result, the present chapter aims to extend the model by incorporating a broader class of feedback mechanisms and testing the robustness of the transition. We also introduce a new method to extract critical exponents more reliably and propose a framework for analytically deriving the stability frontier —a step toward bridging the gap between microscopic agent behavior and macroscopic market instability.

The outline of the Chapter is as follow:

- We begin by providing a more in-depth presentation of the Santa Fe model introduced in [64]
- We then propose several extensions to make the model more realistic.
- Next, we perform a numerical analysis to investigate the presence of a secondorder phase transition for each type of feedback.
- Finally, we introduce a simple framework to recover the stability frontier, solved under the mean-field approximation.

9.2 Motivation: The initial Santa Fe model with feeed-back

The efficient operation of financial markets relies heavily on the order book, which plays a crucial role in facilitating interactions between buyers and sellers. A widely used framework for simulating market dynamics is the Santa Fe model (see [139], [126]). This agent-based model is built on minimal assumptions: orders of unit size are submitted according to Poisson processes, and the tick size, representing the smallest possible price increment, is set to one.

To simulate such a market, one has to specify a few parameters:

- \bullet N is the size of the limit orderbook
- T is the number of time steps for the simulation
- λ Poisson rate for limit order deposition, per time per tick.
- 2μ the Poisson rate for market order deposition per unit time, additive, falling with a probability 1/2 at the best bid or the best ask quote.
- ν_0 Poisson rate for constant cancellation per time per order.
- dP_t a modification of the mid price : $dP_t = \pm \psi dN_t$ with ψ the tick-size and N_t a point process

While this simple model has not been able to recover stylized facts such as diffusive prices or spread-volatility relationship, it can nevertheless account for the basic mechanisms of trading. Furthermore, the Santa Fe model is interesting from a physical point of view, as analogy can easily be made with statistical physics and complex systems. One can interpret orders as particles falling on a 1D grid. The are two kinds of particles, corresponding to buy and sell orders. If two particles of the same nature fall on the same site, they pile up, but if they are of a different nature, they annihilate each other. This classic framework, with fixed Poisson rates is known to be at equilibrium, as long as $\lambda > \nu$, see [7]. Several extensions of the original Santa Fe model have been proposed to better capture stylized facts. For example, [140] introduces market order aggressiveness to generate diffusive prices, while [63] explores a path-dependent refill probability, to recover price diffusivity or the famous Square Root Law impact of metaorders.

To study endogenous liquidity crises, a simple feedback on the cancellation rate

was introduced in [64]. This feedback reads:

$$\nu_t^Z = \nu_0 + \alpha \underbrace{\left(\int_0^t \sqrt{2\beta} e^{-\beta(t-s)} dP_s\right)^2}_{\text{Zumbach Kernel}}$$
(9.1)

where α is the endogeneity ratio and β the memory of the system. This feedback represents the impact of past price trends on future activity. Indeed, for market makers, prolonged trends correspond to significant losses, as they trade against the market. Consequently, the longer the trend persists, the larger their absolute share inventory grows, but the lower its value. This kernel is also natural in the sense that it has been calibrated on empirical data, see [64, 131].

Incorporating this simple feedback across all levels of the order book yields intriguing results, as the system appears to follow a second-order phase transition \cdot

- For α superior to a critical value α^* , an infinitely large orderbook $(N \to \infty)$ will empty in finite time
- Close to the critical point, several functions such as susceptibility or spread distribution exhibit a power-law scaling behavior.
- The critical exponents from the finite-size scaling procedure turns out to be robust regarding fixed values of the system, such as initial conditions.

This result is crucial for market stability, as it demonstrates that a simple agent-based model can account for the anomalously high frequency of non news-related price jumps observed in real financial markets. However, we aim to extend this model to incorporate more realistic feedback mechanisms, to assess whether the phase transition persists when the simulated order book more closely resembles real markets.

9.3 A 4 degree of freedom Santa Fe model

To be closer to reality, we gave this model four different degrees of freedom:

• Kernel: We introduce the Hawkes kernel:

$$\nu_t^H = \nu_0 + \alpha \underbrace{\int_0^t \beta e^{-\beta(t-s)} dP_s^2}_{\text{Hawkes Kernel}}$$
(9.2)

The Hawkes kernel can be interpreted as the effect of volatility on future activity: a period of high volatility represents a significant risk for traders, who are therefore inclined to cancel their orders.

Chapter 9.

- Spatial range: In real markets, orders placed close to the mid-price are more likely to be canceled, as they are often submitted by high-frequency traders. These HFTs are also more responsive to endogenous signals, as they rely on them to earn money. Thus, they are the most affected by these feedback effects. We call local feedback a feedback that occurs only at the best quote, while the cancellation rate remains constant for all other quotes, and global feedback the one that occurs on the whole orderbook.
- In-spread rules The aggressiveness of order placement may also be a key factor in the existence of a phase transition. In the original Santa Fe specification, limit orders could only be placed within one tick of the best quote, referred to as mild deposition. Alternatively, orders could be allowed to fall at any price between the best quote and the mid price. We will call it wild deposition. This could significantly influence the phase transition, as during a liquidity crisis, the spread may grow to infinity, which in turn increases the probability of an order falling within the spread. This compensation phenomenon could potentially mitigate the crisis, thus potentially destroying the phase transition property.
- Triggering Events: Initially, all feedback mechanisms were based on midprice variations, dP_t . However, allowing orders to fall within the spread can lead to artificially large dP_t values, driven by random, insignificant events where a single order falling in the middle of the spread. To resolve this issue, we introduced an additional triggering event, assigning a value of one whenever the spread opens. Then, the dP_t becomes $\mathbb{I}(dS_t > 0)$ with S_t being the size of the spread:

$$\nu_t^S = \nu_0 + \int_0^t \beta e^{-\beta t} \mathbb{I}(dS_t > 0)$$
 (9.3)

We can then simulate this model with 16 different configurations:

Spatial Range	In-spread rules	Kernel	Triggering Event
Global	Mild	Zumbach	Mid Price Change
Local	Wild	Hawkes	Spread opening

Table 9.1: Summary of all possible feedback configurations

From now on, we will refer to feedbacks by their acronyms. To simplify the notation, in the rest of the Chapter we will refer to each feedback mechanism by its acronym. For instance, the Global Mild Zumbach feedback is denoted GMZ. By default, the

triggering event is a change in the mid-price. If we want to refer to feedback triggered by a spread opening, we will denote this by adding an 'S' at the end of the name.

9.4 Investigating phase transitions in the 4DF Sante Fe model

9.4.1 Stability maps

A first step is to verify whether instabilities exist for those configurations, by computing the probability that the time of crisis τ is inferior to the simulation duration T, see Figure 9.1.

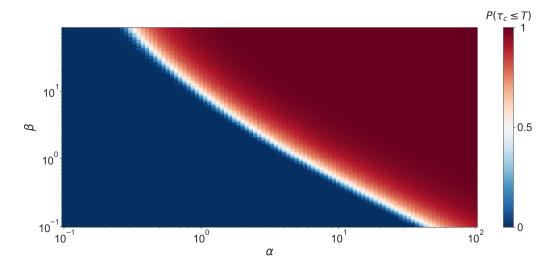


Figure 9.1: Evolution of $\mathbb{P}_N[\tau_c \leq T | \alpha, \beta]$ for the Hawkes global mild retroaction, clearly dividing the (α, β) space in two parts: a blue stable region and a red unstable one. This map is obtained for: $\lambda = 10, \mu = 10, \nu_0 = 1, T = 100, N = 240$. See Appendix C for other stability maps.

These diagrams clearly distinguish two phases in the parameter space (α, β) , separated by the frontier $f(\alpha, \beta)$, which will be the focus of Section 9.5. However, since (N, T) are fixed in these simulations, these stability maps are not conclusive enough to claim that the system exhibits a phase transition. As pointed out in [64], to answer this questions one should study the behavior of $P_N[\tau_c \leq T]$ in the double limit $(N, T) \to \infty$. Taking $T \to \infty$ before N does not provide information on the system, as the probability if crisis will be always one for a given N, if one wait long enough. Thus, mathematically speaking, one should find an $\alpha^*(\beta)$ such

Chapter 9.

that:

$$\lim_{T \to \infty} \lim_{N \to \infty} \mathbb{P}_N[\tau_c \le T, \alpha] = \begin{cases} 1, & \text{when } \alpha \ge \alpha^* \\ 0, & \text{when } \alpha \le \alpha^* \end{cases}$$
(9.4)

This equation is solvable as, for $\alpha > \alpha^*$, $\lim_{N\to\infty} \mathbb{P}_N[\tau_c \leq T, \alpha] > 0$. Thus, the stability of the system lies in a fragile competition between the effect of those two limits. Before going into more mathematical considerations, we can perform a finite size and time scaling method to empirically question the existence of a phase transition.

9.4.2 Finite-size scaling

Following the approach outlined in [64], we employ the finite size scaling method to determine the critical exponents. Finite size scaling, introduced by [141], is a useful technique for extrapolating finite systems to the thermodynamic limit. Phase transition theory assumes an infinite system without boundaries; however, our simulations are limited to finite values of (N,T). In this context, finite size scaling allows us to examine how system responses are affected by this constraint. Near the critical point, the fluctuations due to finite size effects should be negligible in comparison to those arising from the phase transition itself. In our analysis, we utilize the same scaling functions as those described by Fosset et al. [64]:

$$\mathbb{P}_{N}[\tau_{c} \leq T, \alpha_{K}] = F\left(T(\alpha_{K} - \alpha_{m}(N, T))^{\zeta}\right)$$
(9.5)

Here, F(u) is a monotonic regular function that approaches 0 as $u \to -\infty$ and approaches 1 as $u \to +\infty$. The term $\alpha_m(N,T)$ is defined as:

$$\alpha_m(N,T) = \alpha^* - \frac{1}{T^{1/\zeta}} g\left(\frac{N^\eta}{T}\right) \tag{9.6}$$

where g(v) is another function that converges to a constant g_{∞} as $v \to \infty$ and tends to $+\infty$ as $v \to 0$.

We then introduce the function $\chi(\alpha_K, T, N)$ as:

$$\chi(\alpha_K, T, N) = T^{\gamma} G\left(T(\alpha_K - \alpha_m(T, N))^{\zeta}\right) = T^{\gamma} \mathcal{G}\left(NT^{-1/\eta}, T^{1/\zeta}(\alpha_K - \alpha^*)\right)$$
(9.7)

Where χ is the susceptibility and reads :

$$\chi(\alpha_K, T, N) = Var(\min(\tau_c, T)) \tag{9.8}$$

This scaling form has the following interpretation:

- When $1 \ll T \ll N^{\eta}$, $\alpha_m \approx \alpha^*$. As α_K increases, $P_N[\tau_c \leq T, \alpha_K]$ evolves from 0 (no crises) to 1 (crises) in a region of width $T^{-1/\zeta}$ around α^* .
- When $T \gg N^{\eta}$, α_m becomes negative, meaning that $P_N[\tau_c \leq T, \alpha_K]$ is close to 1 for any α_K if one waits long enough.

By applying the finite size scaling procedure detailed in [64], we can extract γ and ζ for each configuration, as shown in 9.2. However, for determining η , the exponent relating N and T, we propose a new procedure, as the previously suggested one focused on the evolution of the maximum of the spread, which may introduce some bias.

9.4.3 Overview of exponents values

The first parameter of the phase transition is γ , which represents the scaling parameter of the crisis time variance. For most feedback configurations, $\gamma \approx 2$, as expected, since the variance naturally scales as T^2 .

Chapter 9.

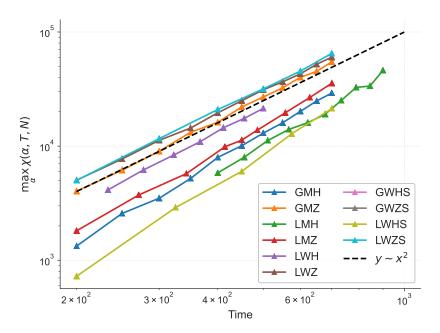


Figure 9.2: Derivation of γ governing the scaling of the maximum of the susceptibility : $\lim_{T\to\infty} \max \chi(\alpha, T, N) \sim T^{\gamma}$. Simulation were done for $\lambda = 10$, $\nu_0 = 1$, N = 140 and $\beta = 1$.

For the second exponent ζ we obtain :

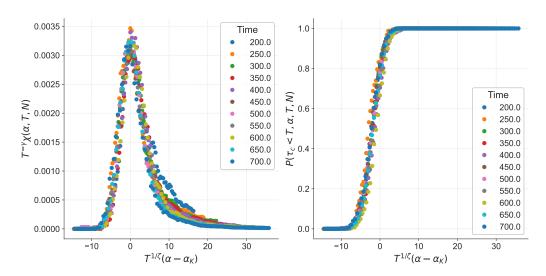


Figure 9.3: Derivation of ζ for the global Hawkes mild retroaction, for $\lambda = 10$, $\nu_0 = 1$, N = 140 and $\beta = 1$. Left: Scaling of the suceptibility, for different T. The x axis is rescaled by the exponent ζ where the y axis is the rescaled by γ . Right: Rescaling of the transition probability function, with the x-axis rescaled by ζ . Additional finite-size scaling plots are provided in Appendix C.

We performed this finite size scaling for most of retroaction:

	GMH	GMZ	LMH	LMZ	GWH	LWH	LWHS	LWZ	LWZS
γ	2	2	2	2	2	2	3	2	2
ζ	2	1	14	99	7	-3	1	10	10

Table 9.2: Summary of the values of γ and ξ for different feedbacks. For reference, GMH stands for "Global Mild Hawkes" feedback. Parameters values used in simulations are given in the appendix, along with the corresponding graphs.

As can be readily observed, some of the exponents obtained from the finite-size scaling appear unrealistic (e.g., -3). We suspect this behavior is related to spread variations: allowing orders to fall within the spread can trigger large mid-price movements, which in turn lead to significant kernel fluctuations. This issue motivated the introduction of new types of triggering events, such as spread openings (see Section 9.3). To reduce the resulting noise, we focus exclusively on spread-opening events and count only one event per spread variation. As for the extreme value 10, 14, 99, it likely corresponds to the limit $\zeta \to \infty$, where the rescaled quantity $T^{1/\zeta}(\alpha - \alpha_K)$ becomes independent of T.

Chapter 9.

These findings are still under investigation, as additional simulations are being conducted to better understand the underlying retroactions. The ultimate—and admittedly optimistic—goal is to classify these retroactive effects into distinct universality classes.

9.4.4 Study of spread explosion

One of the noteworthy facts to emerge from the simulations is the behaviour of the expected value of the spread : $\mathbb{E}[S_t|\alpha,\beta]$, that scales as a power law of time at criticality :

- $\alpha(\beta) < \alpha^*(\beta) \Rightarrow \mathbb{E}[S_t] \sim \bar{S}(\alpha, \beta)$
- $\alpha(\beta) = \alpha^*(\beta) \Rightarrow \mathbb{E}[S_t] \sim t^{1/\eta}, \quad \eta \in \mathbb{N}$
- $\alpha(\beta) > \alpha^*(\beta) \Rightarrow \mathbb{E}[S_t] \sim e^t$

Therefore, near criticality, we propose the following finite-time scaling functions to rescale the spread dynamics for a given β :

$$\mathbb{E}[S(t, \alpha < \alpha^*)] \sim t^{1/\eta} (|\alpha - \alpha^*|)^{\zeta}$$
(9.9)

We believe this method is more robust than the one proposed in [64], as focusing on $\max_{[0,t]} S_t$ may introduce bias.

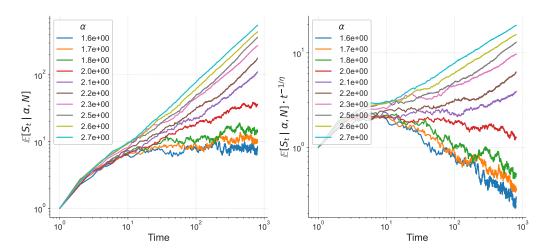


Figure 9.4: Explosion of the spread close at the critical α^* . Simulation was done for the LMH retroaction, with $\lambda=10$, $\nu_0=1$, $\beta=1$ and N=1000. Once rescaled by $t^{-1/\eta}$ with $\eta=2$, it clearly reveals the value of α at which the transition occurs. Further spread scaling plots can be found in Appendix C.

To obtain more accurate results, here is a convenient method to determine the parameters (α^*, η) . First one can find the value of η by plotting $\frac{\mathbb{E}[S(t, \alpha > \alpha)]}{t^{\xi}}$ for $\xi \in]0, 1]$ and select $\xi = 1/\eta$ such as one of the curve is a horizontal straight line. Then, one should fix $x = t^{1/\eta}$ and fit $\frac{\mathbb{E}[S(t, \alpha > \alpha)]}{x} = (a_0 + (a_1(\alpha)x)^q)^{1/q}$ with same optimal a_0, q for all curves. Then, by taking optimal $\alpha_1(\alpha)$ for each α , one could regress $\alpha_1(\alpha) = (\alpha - \alpha^*)^{2\zeta}$

	GMH	GMZ	LMH	LMZ
η	2	3	2	3

Table 9.3: Values of the scaling exponent η estimated for selected feedback types.

Thus, interestingly, we can see that going from a local retroaction to a global one doesn't change the value of η , but this value seems to be determined by the kind of retroaction.

9.5 Analytical Determination of the Stability Boundary

In this section, we want to analytically derive the equation for the frontier $\alpha(\beta)$, that separates the stable and unstable phases. The overall approach can be summarized as follows: there is a competition between queue depletion at the best price, driven by an increasing cancellation rate, and the placement of orders within the spread. From a mean-field perspective, the problem reduces to the following question: if $\nu^{Z,H,S}(t)$ eventually reaches a level at which the mean time of queue depletion equals the mean time for order placement within the spread, what conditions on (α, β) ensure this balance is maintained indefinitely? For simplicity, we focus on the ask side, as the system is symmetric. Furthermore, we provide here the method and derived it for two of the feedback, even though it could be done for all of them. This section is currently under development.

Dynamic of the cancellation:

Let's begin by explicitly describe the dynamic of $\nu^{Z,H}(t)$. We denote by t_i , the time just after the i-th price jump. For the Hawkes kernel we get, for $0 < t < t_{i+1} - t_i$:

$$\nu^{H}(t_{i}+t) = \nu_{0} + \alpha \int_{0}^{t_{i}+t} \beta e^{-\beta(t_{i}+t-s)} dP_{s}^{2}$$
(9.10)

$$= (\nu^{H}(t_i) - \nu_0)e^{-\beta t} + \nu_0 \tag{9.11}$$

Chapter 9.

and, at t_{i+1} :

$$\nu^{H}(t_{i+1}) = (\nu^{H}(t_i) - \nu_0)e^{-\beta(t_{i+1} - t_i)} + \nu_0 + \alpha\beta dP_{t_{i+1}}^2$$
(9.12)

For the Zumbach kernel, cancellation then reads:

$$\nu^{Z}(t_{i}+t) = \nu_{0} + \alpha \left(\int_{0}^{t_{i}+t} \sqrt{2\beta} e^{-\beta(t_{i}+t-s)} dP_{s} \right)^{2}$$

$$= \nu_{0} + \alpha \left(e^{-\beta t} \int_{0}^{t_{i}} \sqrt{2\beta} e^{-\beta(t_{i}-s)} dP_{s} + \int_{t_{i}}^{t_{i}+t} \sqrt{2\beta} e^{-\beta(t_{i}+t-s)} dP_{s} \right)^{2}$$

$$= (\nu^{Z}(t_{i}) - \nu_{0})e^{-2\beta t} + \nu_{0}$$

$$(9.13)$$

Then, at t_{i+1} :

$$\nu^{Z}(t_{i+1}) = (\nu^{Z}(t_i) - \nu_0)e^{-2\beta(t_{i+1} - t_i)} + \nu_0 + 2\alpha\beta dP_{t_{i+1}}^2$$
(9.16)

+
$$\sqrt{2\beta}e^{-\beta(t_{i+1}-t_i)}\sqrt{(\nu_i-\nu_0)}\text{sign}(dP_{t_{i+1}}\int_0^{t_i}\sqrt{2\beta}e^{-\beta(t_i-s)}dP_s$$
 (9.17)

Dynamics of queue under Poissonian order flow

The second ingredient important to understand spread dynamic is the behavior of the best quote queue size V_t^b , given by :

$$dV_t^b = -\nu_t V_t^b dN_t^\nu + dN_t^\lambda - dN_t^\mu$$
(9.18)

$$V_t^b = V_0^b e^{-\int_0^t dN_u^{\nu}} + \int_0^t e^{\int_s^t dN_u^{\nu}} (dN_s^{\lambda} - dN_s^{\mu})$$
 (9.19)

$$\mathbb{E}[V_t^b] = V_0^b e^{-\int_0^t \nu^{Z,H}(s)ds} + \int_0^t e^{\int_s^t \nu^{Z,H}(u)du} (\lambda - \mu)ds$$
 (9.20)

(9.21)

If $\nu^{Z,H}(t) = \nu_0$, we recover a classic queue dynamic, in line with [7]:

$$\mathbb{E}[V_t^b] = V_0^b e^{-\nu_0 t} + \frac{\lambda - \mu}{\nu_0} (1 - e^{-\nu_0 t})$$
(9.22)

Under the assumption $\lambda = \mu$, we obtain the two following equations describing the queue behavior at the best quote V^b and in the rest of the orderbook, V:

$$\mathbb{E}[V_t^b] = V_0 e^{-\int_0^t \nu^{Z,H}(s)ds}$$
(9.23)

$$\mathbb{E}[V_t] = V_0 e^{-\int_0^t \nu^{Z,H}(s)ds} + \lambda \int_0^t e^{\int_s^t \nu^{Z,H}(u)du} ds$$
 (9.24)

Then, by taking as a reference time t=0 the time at which a queue became the best quote, the average depletion time is then given by:

$$t^* = \min_{t} (\mathbb{E}[V_t^b] < 1) \iff \int_0^{t^*} \nu^{H,Z,S}(s) ds = \log(V_0^b)$$
 (9.25)

In the mean field approximation, it can be easily shown that the spread will widen if $t^* < 1/\lambda$. If $t^* > 1/\lambda$, a deposition occurs within the spread, and the system remains stable. Even if this approximation seems very gross, we will show that it's enough to recover the analytical frontier in certain cases. We believe that by fine tuning this approach, it may also be possible to derive analytical frontier for each feedback.

Solving for local Zumbach The LMZ stability frontier can be recovered in the limit of small β . With retroaction, an additional term emerges that captures the price trend. This term amplifies jump sizes when the *i*-th price change has the same sign as the memory kernel. However, under the small β approximation—corresponding to a long-memory regime—this kernel is significantly weakened. Furthermore, due to bid-ask symmetry, the price tends to mean-revert. Thus, in the mean-field approximation, this term can be neglected. We then obtain:

$$\nu_{i+1} = (\nu_i - \nu_0)e^{-2\beta/\lambda} + \nu_0 + 2\alpha\beta \tag{9.26}$$

$$\nu_i = \lambda \log(V_0) \tag{9.27}$$

Leading to a stability frontier given by:

$$\alpha = \frac{\lambda \log(V_0)(1 - e^{-2\beta/\lambda}) - \nu_0(1 - e^{-\beta/\lambda})}{2\beta}$$
 (9.28)

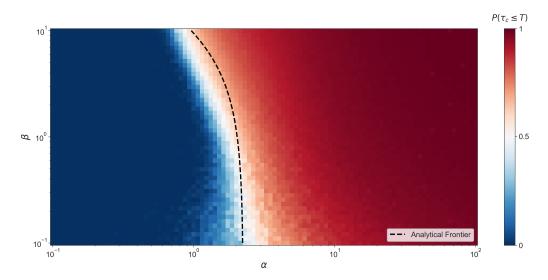


Figure 9.5: Derivation and Mean Field Validation of the LMZ Feedback Frontier: Simulations conducted with parameters $\lambda = 10$, $\mu = 20$, $\nu_0 = 1$, and N = 140.

Solving for the local Hawkes For the LMH feedback model, things are more messy due to the absence of a closed-form solution. We are still interested in identifying the conditions under which a queue is depleted before a new order is placed inside the spread.

Asymptotic regime $\beta \to 0$: In the limit of small β , the decay of the kernel is very slow. The queue will be consumed before any new deposition in the spread occurs if the following condition holds:

$$\nu_0/\lambda + \alpha(1 - e^{-\beta/\lambda}) > \log(V_0) \tag{9.29}$$

This can be rewritten as a constraint on α :

$$\alpha < \frac{\log(V_0) - \nu_0/\lambda}{1 - e^{-\beta/\lambda}} \tag{9.30}$$

Asymptotic regime $\beta \to \infty$: In the opposite limit, where β tends to infinity, the feedback behaves more like a jump process: it decays rapidly between events but exhibits sharp increases at each price change. In this regime, the condition for queue depletion becomes, after the *ith* price change:

$$\nu_0/\lambda + \frac{\nu_0 + i\alpha\beta}{\beta} (1 - e^{-\beta/\lambda}) > \log(V_0)$$
(9.31)

which can be rewritten as:

$$i\alpha\beta > \frac{\beta(\log(V_0) - \nu_0/\lambda)}{1 - e^{-\beta/\lambda}} - \nu_0 \tag{9.32}$$

These two expressions allow us to characterize the asymptotic behavior of the system under fast and slow feedback regimes, and to identify the parameter ranges where queue depletion dominates over order placement. However, Figure 9.6 shows that these analytical formulas are accurate only in the small β regime.

This discrepancy can be attributed to the simulation's implementation. For computational efficiency, the evolution of $\nu(t)$ is modeled as a step function, updated only at discrete event times. Specifically, we approximate the continuous-time integral as:

$$\int_0^t e^{-\beta(t-s)} dP_s \approx \sum_{T_n \le t} e^{-\beta(t-T_n)} \Delta P_{T_n}$$
(9.33)

Thus, with this discretization, equation (30) becomes:

$$\begin{cases} \nu_i e_0 + \sum_{i>1} (\nu_i - \nu_0) e^{-\beta e_i} (e_i - e_{i-1}) \\ e_i - e_{i-1} = \frac{1}{\lambda + \mu + \nu(e_i) V_{e_i}} \end{cases}$$
(9.34)

where $\{e_0, e_1, \dots\}$ denote the times at which events occur in the queue. The gap between the theoretical and discretized curves can then be explained as follows: for large β , the function $\nu(t)$ increases sharply by an amount proportional to $\alpha\beta$ following a price change, but also decays rapidly due to the strong exponential term $e^{-\beta t}$. As a result, the behavior of $\nu(t)$ approaches that of a Dirac delta function. However, once discretized, $\nu(t)$ retains its post-jump value over a longer interval—until the next event in the order book—thereby distorting the intended continuous decay. Then green dashed line represents the analytical frontier based on eq. (9.34), setting i=10, as it provides the best fit.

However, to get rid of this approximation, the next step is to calculate this frontier in the thermodynamic limit.

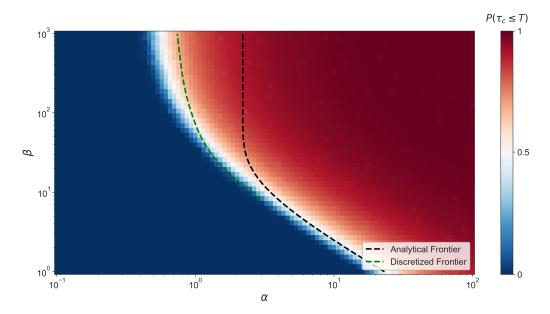


Figure 9.6: Analytical derivation of the frontier under LMH feedback. Simulations were performed with $\lambda=10,\,\mu=20,\,\nu_0=1,$ and N=140. The black dashed curve represents the analytical solution, which is asymptotically accurate in the small β regime but fails for large β . The green dashed curves show the large- β asymptotes, adjusted to account for discretization effects.

9.6 Conclusion

In this chapter, we explored the emergence of endogenous liquidity crises in a Santa Fe-inspired agent-based model of financial markets. Extending the zero-intelligence framework, we introduced various feedback mechanisms through which agents react to past market activity. These mechanisms included volatility-sensitive Hawkes-type kernels, price-trend-sensitive Zumbach-type kernels, spatially localized responses, in-spread order placement rules, and alternative triggering events such as spread openings. For the feedback previously investigated in [64], we obtained comparable results, thereby affirming the validity of their analysis.

We demonstrated that these feedback loops can lead to second-order phase transitions, where the order book dynamics shift from a stable to an unstable regime purely as a result of internal interactions—without any exogenous shocks. This finding provides further evidence against the Efficient Market Hypothesis and supports the view of financial markets as complex, self-organizing systems.

Using a combination of empirical simulations and finite-size scaling techniques, we

extracted critical exponents associated with the transitions and uncovered universal behaviors across feedback types. We also proposed a novel method based on the power-law explosion of the spread to estimate the critical exponent η more robustly. We must further interpret these exponents and investigate their possible connections to universality classes, should such relationships exist. Furthermore, we derived analytical approximations of the stability frontier $\alpha^*(\beta)$ for certain feedback configurations, providing insight into the interplay between memory effects, endogeneity, and market fragility.

Chapter 9.

Part IV Conclusion and Perspectives

Conclusion

The stock market is filled with individuals who know the price of everything,
but the value of nothing

Oscar Wild

Overview of the results

Let us briefly summarize the main results of this thesis. Our aim was to better understand the mechanisms behind price formation and to explore whether we could support an *Order-Driven view* of financial markets—one in which prices are primarily shaped by the mechanical impact of trades rather than by the flow of information.

Chapter 4: Exploring the Square Root Law Using the TSE Dataset

In Chapter 1, we conducted an extensive investigation of the so-called Square Root Law of metaorder impact. Thanks to a unique dataset, we were able to study this phenomenon at a very granular level and without the biases typically associated with proprietary data. We disentangled the impact due to the volume of child orders from the effect of their execution rank. After categorizing each participant, we also used the dataset to challenge the predictions of several theoretical models.

We further used this dataset to study the liquidity provider side, which is usually less in the spotlight. We showed that liquidity providers also tend to execute metaorders—in the sense that they submit successively multiple limit orders of the same sign. The lengths of these metaorders follow a power-law distribution as well, albeit with a larger exponent, indicating that providers' metaorders are typically shorter than those of liquidity takers.

Conclusion

Finally, and somewhat surprisingly, we found that metaorders could be reconstructed randomly—by breaking the link between traders and their executed trades—without significantly altering their empirical properties. This undermines what was initially considered as the main added value of the dataset.

From this investigation in the Japanese market, we drew two key observations, which we explore further in Chapters 2 and 3:

- (O1) The order flow can be viewed as a superposition of metaorders, each obeying a square-root impact law.
- (O2) Knowing the specific information behind trades—or at least identifying the traders who executed them—is not essential for reconstructing realistic metaorders.

Chapter 5: The metaorder proxy

The next natural step was to validate and extend this synthetic metaorder reconstruction method to a broader set of assets traded on various exchanges, using only publicly available trade data. To do so, we introduced a simple, easy-to-use algorithm—which could, and perhaps should, be further refined—to generate realistic synthetic metaorders. We call them realistic because they reproduce three well-established stylized facts from the literature:

- The peak impact scales with the square root of the volume.
- The impact is a concave function of the execution time—more precisely, it also follows a square-root law.
- After execution, the impact decays over a time scale longer than the execution itself.

Regarding the last point, although the nature of impact decay remains a topic of debate, several empirical studies showed that it should decay slowly, following a propagator kernel with a power-law form: $t^{-\beta}$, where $\beta \approx 0.25$. Once again, this observation aligns well with what we observe in our synthetic dataset.

The surprisingly effective performance of this synthetic metaorder proxy has two main implications:

• It provides a powerful tool for both practitioners and academics. Practitioners can enrich their datasets, while researchers can now study market impact without relying on proprietary data, which is often inaccessible or

biased by execution logic. Thanks to publicly available data—especially in crypto markets like Binance or Coinbase—realistic metaorders can now be reconstructed by anyone.

- This approach may also facilitate a more thorough investigation of phenomena that have been difficult to capture due to data scarcity, such as *cross impact*, see [50].
- It fully decouples information content from metaorders. Since synthetic metaorders are generated independently of any alpha or economic predictive signal, they are free from informational bias. One might argue they still reflect average information patterns, but given the consistency of impact across synthetic metaorders, I find this unlikely.

Chapter 6 & 7: A unified framework for market microstructure

Armed with the two empirical observations from Chapters 3 and 4, we found that these were sufficient to construct a unified framework for market microstructure. Using (O1) and the LMF hypothesis, we were able to model a realistic order flow. Then, drawing on (O2), we could understand how each metaorder impacts prices.

First, we had to slightly generalize the LMF hypothesis to reflect the empirical fact that large orders tend to be less correlated than small ones. Then, we introduce a parameter a, which allows one to give more weight to either small or large volumes. With this, we were able to predict and validate non-trivial scaling laws related to order flow imbalances. These findings were further confirmed with simulations in Chapter 4 7, reinforcing the validity of our theoretical results.

Turning to prices, the situation is more subtle. As repeatedly emphasized in this thesis, the black box that transforms order flow into prices is complex. Our goal was to reconcile two aspects that most (all) existing theories - unfortunately -don't manage to conciliate: (1) the scaling of price moments, i.e., the diffusive behavior that is crucial for financial markets, and (2) the square-root impact law for metaorders, which encompasses three key empirical features.

To achieve this, we proposed a generalized propagator —a mechanism which, when inserted into the black box, transforms a realistic order flow into a realistic price process. Once again, this led to non-trivial predictions, particularly for the correlation and covariance structure of the price and order flow. These predictions matched empirical data surprisingly well and were also reproduced in our simulations.

Conclusion

Finally, we demonstrated that our framework naturally explains the coexistence of a linear aggregated impact (modulated by an anomalous prefactor) and the square-root law for individual metaorders. We were even able to predict—and verify both empirically and in simulation—the exact (or very close) value of the exponent in the prefactor. Moreover, applying the synthetic metaorder proxy to our simulated price recovered the correct square-root law, although some fine-tuning was required.

In summary, this framework delivers two key messages:

- Prices can be fully explained by the order flow. We have shown that volatility
 arises mechanically from trading activity, rather than from news or informational shocks.
- We provide a simple yet powerful algorithm to simulate realistic market dynamics—both order flow and prices. This could be highly valuable to researchers and practitioners in market microstructure. It also supports the order-driven view of markets, demonstrating that information is not necessary to replicate realistic behaviors.

Chapter 8: The VAR model and microstructure modes

After focusing exclusively on trades and mid-price dynamics, we broadened our analysis of market microstructure to include a more comprehensive set of variables. Once the appropriate time granularity was chosen—specifically, aggregating every 20 price changes—we defined a state vector composed of eight variables: the full order flow (including limit orders, cancellations, and executions) at the best bid and ask, along with the volatility (measured as the time spent in each state) and the return.

Following a series of technical analyses, this framework enabled us to identify a few distinct microstructural modes, which proved to be highly significant from a phenomenological perspective. Not only did this approach allow us to predict the system's evolution—with greater accuracy than most traditional models, including AI-driven ones—but it also enabled us to test the system's stability. In particular, as we increased the memory of the system (i.e., the lag in the autoregressive model), we observed that the dominant mode of evolution tended toward an opening of the spread, driven by increased cancellations on both sides of the book. In the limit, this mode could even become critical, with its associated eigenvalue approaching one.

Two key findings emerge from this study:

- Markets appear to transition between a few distinct modes, and modeling them as such provides accurate predictions for both order flow and returns.
- The concept of endogenous liquidity crises is inherently embedded in market microstructure: the limit order book itself appears to operate near marginal instability.

Chapter 9: A generalized Agent Based Model model to study endogenous liquidity crisis

Seeing that liquidity crises are both predicted by market microstructure models and frequently make headlines of macro-finance journals provided additional motivation to revisit and extend the so-called Santa Fe model. Indeed, if these anomalies are both frequent and notoriously difficult to understand theoretically, agent-based models appear to be an efficient tool to gain deeper insight into them.

To this end, we developed a model in which traders interact randomly with the market, except for one crucial feature: they are subject to both leverage and the Zumbach effect—that is, feedback from past volatility and past price trends—two well-known mechanisms in equity markets. We proposed new realistic configurations in which this feedback occurs only at the best quotes, or where agents are allowed to place orders within the spread, to better match real-world behavior.

We showed that, in most cases, the finite-size scaling procedure is sufficient to confirm the existence of a second-order phase transition and to extract the critical exponents of this transition. Finally, we proposed a simple method to derive the analytical form of the stability frontier under a mean-field approximation.

The take home message is the following:

Assuming that agents trade randomly but react to past information is enough
to trigger endogenous liquidity crises—provided that their memory and the
level of endogeneity in the system exceed certain thresholds.

Future works and closing remarks

Many shadow areas still need to be understood, and I fully agree that some parts of this thesis raised more questions than they answered. I will therefore outline here few main research directions that, to my view, may deserve further clarifications:

- Empirical extension of the metaorder proxy: This may be the most puzzling aspect to understand —why does the proxy work so well for some assets, and how can it be generalized to a broader set of assets? So far, the proxy has performed remarkably well on diverse instruments, including a highly liquid future and various stocks on the Paris Stock Exchange, with little fine-tuning. However, preliminary investigations suggest that for some assets, more careful calibration is needed. This makes sense: the mapping function should resemble the true one as closely as possible, and that may depend on the specific microstructural features of each asset. A natural next step is to generalize the proxy to a wider range of assets and possibly use machine learning techniques to automatically identify optimal parameters. To the best of my knowledge, I have not encountered a single asset for which, after some manual adjustment, the square-root law could not be recovered —though this sometimes required significant effort. We argue that the mapping function is not unique: obtaining realistic metaorders does not necessarily require the exact underlying trader-trade correspondence.
- Extending the unified framework: Another promising direction would be to extend the framework and to explore whether modifying the rate at which new metaorders are initiated enables one to reproduce other well-known stylized facts of financial time series, such as volatility clustering, fat tails, the leverage effect and others. Furthermore, our framework also makes predictions regarding the Y-ratio, i.e., the prefactor of the square-root law. This prefactor is of particular interest to practitioners, so it would be valuable to investigate whether our prediction can be empirically validated. Another possible extension is to incorporate spread dynamics into the model, which we have completely neglected so far.
- The theoretical foundation of the metaorder proxy: A more theoretical task would be to use the unified framework to rigorously demonstrate why and under which conditions the metaorder proxy works.
- Beyond the propagator: While the (generalized) propagator model offers a convenient mathematical framework for describing market impact, it lacks a phenomenological foundation—that is, a deeper theory grounded in the behavior of market participants. The LLOB model comes close by introducing random agents who revise their "fundamental" price stochastically, but

- as we have shown, it does not fully align with empirical data. In my view, a promising direction would be to model liquidity providers with memory, building on the concept of provider metaorders introduced in Chapter 4.
- Self-organized criticality in market microstructure: Despite being a very intriguing and active area of research, the hypothesis that financial markets are self-organized critical systems remains unproven (to the best of my knowledge). A breakthrough would be to show that, under a realistic agent-based model such as the Santa Fe model with reaction terms and PnL optimization, markets *should* naturally evolve toward a state close to the unstable frontier to remain efficient.

More broadly, even if this work—hopefully—convinces the reader that the *Order-Driven view* of markets is the right one, I hope that the available codes (both the metaorder proxy and the market simulator) and the presented theories will inspire further research and lead to deeper insights into this fascinating field of market microstructure.

Conclusion

Bibliography

- Elomari-Kessab, S., Maitrier, G., Bonart, J. & Bouchaud, J. "Microstructure Modes" –Disentangling the Joint Dynamics of Prices & Order Flow. Wilmott 2024. ISSN: 1541-8286. http://dx.doi.org/10.54946/wilm.12074 (2024).
- 2. Maitrier, G., Loeper, G., Kanazawa, K. & Bouchaud, J.-P. The" double" square-root law: Evidence for the mechanical origin of market impact using Tokyo Stock Exchange data. arXiv preprint arXiv:2502.16246 (2025).
- 3. Maitrier, G., Loeper, G. & Bouchaud, J.-P. Generating realistic metaorders from public data 2025. arXiv: 2503.18199 [q-fin.TR]. https://arxiv.org/abs/2503.18199.
- 4. Maitrier, G. & Bouchaud, J.-P. The Subtle Interplay between Square-root Impact, Order Imbalance & Volatility: A Unifying Framework 2025. arXiv: 2506.07711 [q-fin.TR]. https://arxiv.org/abs/2506.07711.
- 5. Maitrier, G., Loeper, G. & Bouchaud, J.-P. The Subtle Interplay between Square-root Impact, Order Imbalance & Volatility: A Numerical Simulation In preparation. 2025.
- 6. Maitrier, G., Loeper, G., Benzaquen, M. & Bouchaud, J.-P. A Generalized Santa-Fe-like Model to Understand Liquidity Crises In preparation. 2025.
- 7. Bouchaud, J.-P., Bonart, J., Donier, J. & Gould, M. *Trades, quotes and prices: financial markets under the microscope* (Cambridge University Press, 2018).
- 8. Odlyzko, A. Newton's Financial Misadventures in the South Sea Bubble. Notes and Records: The Royal Society Journal of the History of Science 73, 29-59. http://doi.org/10.1098/rsnr.2018.0018 (2019).
- 9. Salek, M., Challet, D. & Toke, I. M. Price impact in equity auctions: zero, then linear 2023. arXiv: 2301.05677 [q-fin.ST]. https://arxiv.org/abs/2301.05677.

- Bonart, J. & and, M. D. G. Latency and liquidity provision in a limit order book. Quantitative Finance 17, 1601–1616. eprint: https://doi.org/10. 1080/14697688.2017.1296177. https://doi.org/10.1080/14697688. 2017.1296177 (2017).
- 11. Fabre, T. & Toke, I. M. High-Frequency Market Manipulation Detection with a Markov-modulated Hawkes process. arXiv preprint arXiv:2502.04027 (2025).
- 12. Laruelle, S., Rosenbaum, M. & Savku, E. Assessing MiFID 2 Regulation on Tick Sizes: A Transaction Costs Analysis Viewpoint tech. rep. Available at SSRN: https://ssrn.com/abstract=3256453 or http://dx.doi.org/10.2139/ssrn.3256453 (Sept. 2018). https://ssrn.com/abstract=3256453.
- 13. Menkveld, A. J. & Zoican, M. Need for Speed? Exchange Latency and Liquidity. Review of Financial Studies 30. Available at SSRN: https://ssrn.com/abstract=2442690, 1188-1228. http://dx.doi.org/10.2139/ssrn.2442690 (2017).
- 14. Frazzini, A., Israel, R. & Moskowitz, T. J. Trading costs. *Available at SSRN 3229719* (2018).
- 15. Donier, J., Bonart, J., Mastromatteo, I. & Bouchaud, J.-P. A fully consistent, minimal model for non-linear market impact. *Quantitative finance* **15**, 1109–1121 (2015).
- Sato, Y. & Kanazawa, K. Inferring microscopic financial information from the long memory in market-order flow: A quantitative test of the Lillo-Mike-Farmer model. *Physical Review Letters* 131, 197401 (2023).
- 17. Kurth, J. G., Majewski, A. A. & Bouchaud, J.-P. Revisiting the Excess Volatility Puzzle Through the Lens of the Chiarella Model. arXiv preprint arXiv:2505.07820 (2025).
- 18. Lillo, F., Mike, S. & Farmer, J. D. Theory for long memory in supply and demand. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics* **71**, 066122 (2005).
- Gabaix, X. Power Laws in Economics and Finance. Annual Review of Economics 1, 255-293. https://doi.org/10.1146/annurev.economics.050708.142940 (2009).
- 20. Tsaknaki, I.-Y., Lillo, F. & Mazzarisi, P. Online Learning of Order Flow and Market Impact with Bayesian Change-Point Detection Methods 2024. arXiv: 2307.02375 [q-fin.TR]. https://arxiv.org/abs/2307.02375.

- 21. Naviglio, M., Bormetti, G., Campigli, F., Rodikov, G. & Lillo, F. Why is the estimation of metaorder impact with public market data so challenging? 2025. arXiv: 2501.17096 [q-fin.TR]. https://arxiv.org/abs/2501.17096.
- 22. Bouchaud, J.-P. Statistical Mechanics of Financial Markets https://www.college-de-france.fr/chaire/jean-philippe-bouchaud-mecanique-statistique-des-marches-financiers. Lecture series at the Collège de France. 2023.
- 23. Bachelier, L. Théorie de la spéculation in Annales scientifiques de l'École normale supérieure 17 (1900), 21–86.
- 24. Mandelbrot, B. The variation of some other speculative prices. *The Journal of Business* **40**, 393–413 (1967).
- 25. Bacry, E., Delour, J. & Muzy, J.-F. Multifractal random walk. *Physical review E* **64**, 026103 (2001).
- 26. Zarhali, O., Aubrun, C., Bacry, E., Bouchaud, J.-P. & Muzy, J.-F. Why is the volatility of single stocks so much rougher than that of the S&P500? arXiv preprint arXiv:2505.02678 (2025).
- 27. Gatheral, J., Jaisson, T. & Rosenbaum, M. Volatility is rough. *Quantitative Finance* **18**, 933–949. eprint: https://doi.org/10.1080/14697688.2017.1393551. https://doi.org/10.1080/14697688.2017.1393551 (2018).
- 28. Blanc, P., Donier, J. & Bouchaud, J.-P. Quadratic Hawkes processes for financial prices 2015. arXiv: 1509.07710 [q-fin.TR]. https://arxiv.org/abs/1509.07710.
- Fosset, A., Bouchaud, J.-P. & and, M. B. Non-parametric estimation of quadratic Hawkes processes for order book events. The European Journal of Finance 28, 663–678. eprint: https://doi.org/10.1080/1351847X. 2021.1917441. https://doi.org/10.1080/1351847X.2021.1917441 (2022).
- 30. Saddier, L. & Marsili, M. A Bayesian theory of market impact. *Journal of Statistical Mechanics: Theory and Experiment* **2024**, 083404 (2024).
- 31. Bouchaud, J.-P., Gefen, Y., Potters, M. & Wyart, M. Fluctuations and response in financial markets: the subtle nature of random' price changes. *Quantitative finance* 4, 176 (2003).
- 32. Gerig, A. A Theory for Market Impact: How Order Flow Affects Stock Price 2008. arXiv: 0804.3818 [q-fin.ST]. https://arxiv.org/abs/0804.3818.

- 33. Taranto, D. E., Bormetti, G., Bouchaud, J.-P., Lillo, F. & Tóth, B. Linear models for the impact of order flow on prices. I. History dependent impact models. *Quantitative Finance* 18, 903–915 (2018).
- 34. Taranto, D. E., Bormetti, G., Bouchaud, J.-P., Lillo, F. & Toth, B. Linear models for the impact of order flow on prices II. The Mixture Transition Distribution model 2016. arXiv: 1604.07556 [q-fin.TR]. https://arxiv.org/abs/1604.07556.
- 35. Patzelt, F. & Bouchaud, J.-P. Nonlinear price impact from linear models. Journal of Statistical Mechanics: Theory and Experiment 2017, 123404. ISSN: 1742-5468. http://dx.doi.org/10.1088/1742-5468/aa9335 (Dec. 2017).
- 36. Toth, B., Eisler, Z. & Bouchaud, J.-P. The short-term price impact of trades is universal. *Market Microstructure and Liquidity* **3**, 1850002 (2017).
- 37. Patzelt, F. & Bouchaud, J.-P. Universal scaling and nonlinearity of aggregate price impact in financial markets. *Physical Review E* **97**, 012304 (2018).
- 38. Tóth, B. *et al.* Anomalous price impact and the critical nature of liquidity in financial markets. *Physical Review X* 1, 021006 (2011).
- Said, E., Ayed, A. B. H., Husson, A. & Abergel, F. Market impact: A systematic study of limit orders. Market microstructure and liquidity 3, 1850008 (2017).
- 40. Tóth, B., Eisler, Z. & Bouchaud, J.-P. The Square-Root Impace Law Also Holds for Option Markets. *Wilmott* **2016**, 70–73 (2016).
- 41. Eisler, Z. & Bouchaud, J.-P. Price Impact Without Order Book: A Study of the OTC Credit Index Market. *SSRN Electronic Journal*. Available at SSRN: https://ssrn.com/abstract=2840166 or http://dx.doi.org/10.2139/ssrn.2840166 (Sept. 2016).
- 42. Sato, Y. & Kanazawa, K. Does the square-root price impact law belong to the strict universal scalings?: quantitative support by a complete survey of the Tokyo stock exchange market. arXiv preprint arXiv:2411.13965 (2024).
- 43. Farmer, J. D., Gerig, A., Lillo, F. & Waelbroeck, H. How efficiency shapes market impact. *Quantitative Finance* **13**, 1743–1758 (2013).
- 44. Bucci, F., Benzaquen, M., Lillo, F. & Bouchaud, J.-P. Slow decay of impact in equity markets: insights from the ANcerno database. *Market Microstructure and Liquidity* 4, 1950006 (2018).
- 45. Kyle, A. S. & Obizhaeva, A. A. The Market Impact Puzzle. Anna A., The Market Impact Puzzle (February 4, 2018) (2018).

- 46. Torre, N. & Ferrari, M. Market Impact Model Handbook Available at http://www.barra.com/newsletter/nl166/miminl166.asp. Berkeley, 1997.
- 47. Grinold, R. C. & Kahn, R. N. Active portfolio management (2000).
- 48. Benzaquen, M. & Bouchaud, J.-P. Market impact with multi-timescale liquidity. *Quantitative Finance* **18**, 1781–1790 (2018).
- 49. Webster, K. T. Handbook of price impact modeling (Chapman and Hall/CRC, 2023).
- 50. Hey, N., Bouchaud, J.-P., Mastromatteo, I., Muhle-Karbe, J. & Webster, K. The Cost of Misspecifying Price Impact 2023. arXiv: 2306.00599 [q-fin.TR]. https://arxiv.org/abs/2306.00599.
- 51. Jusselin, P. & Rosenbaum, M. No-arbitrage implies power-law market impact and rough volatility 2018. arXiv: 1805.07134 [q-fin.ST]. https://arxiv.org/abs/1805.07134.
- 52. Gabaix, X. & Koijen, R. S. In search of the origins of financial fluctuations: The inelastic markets hypothesis tech. rep. (National Bureau of Economic Research, 2021).
- 53. Bouchaud, J.-P. The inelastic market hypothesis: a microstructural interpretation. *Quantitative Finance* **22**, 1785–1795 (2022).
- 54. U.S. Securities and Exchange Commission & U.S. Commodity Futures Trading Commission. Findings Regarding the Market Events of May 6, 2010 Staff Report SEC-2010-81. Report of the staffs of the CFTC and SEC to the Joint Advisory Committee on Emerging Regulatory Issues (U.S. Securities, Exchange Commission, and U.S. Commodity Futures Trading Commission, Washington, D.C., Sept. 2010). https://www.sec.gov/files/marketevents-report.pdf.
- 55. Marcaccioli, R., Bouchaud, J.-P. & Benzaquen, M. Exogenous and endogenous price jumps belong to different dynamical classes. *Journal of Statistical Mechanics: Theory and Experiment* **2022**, 023403 (2022).
- 56. Aubrun, C., Morel, R., Benzaquen, M. & Bouchaud, J.-P. Identifying new classes of financial price jumps with wavelets. *Proceedings of the National Academy of Sciences* **122**, e2409156121 (2025).
- 57. Shiller, R. Do Stock Prices Move Too Much to Be Justified by Subsequent Changes in Dividends? *American Economic Review* **71**, 421–36 (Jan. 1981).
- 58. Guyon, J. & Sopze, J. L. Volatility Is (Mostly) Path-Dependent Available at SSRN: https://ssrn.com/abstract=4174589. SSRN working paper. July 2022. http://dx.doi.org/10.2139/ssrn.4174589.

- 59. Bacry, E., Mastromatteo, I. & Muzy, J.-F. Hawkes processes in finance. Market Microstructure and Liquidity 1, 1550005 (2015).
- 60. Hardiman, S. J., Bercot, N. & Bouchaud, J.-P. Critical reflexivity in financial markets: a Hawkes process analysis. *The European Physical Journal B* **86**, 1–9 (2013).
- 61. Knicker, M. S., Naumann-Woleske, K., Bouchaud, J.-P., et al. Post-COVID inflation and the monetary policy dilemma: an agent-based scenario analysis. Journal of Economic Interaction and Coordination 20, 141–195. https://doi.org/10.1007/s11403-024-00413-3 (2025).
- 62. Mastromatteo, I., Tóth, B. & Bouchaud, J.-P. Anomalous Impact in Reaction-Diffusion Financial Models. *Phys. Rev. Lett.* **113**, 268701. https://link.aps.org/doi/10.1103/PhysRevLett.113.268701 (26 Dec. 2014).
- 63. Ravagnani, A. & Lillo, F. Modeling metaorder impact with a Non-Markovian Zero Intelligence model 2025. arXiv: 2503.05254 [q-fin.TR]. https://arxiv.org/abs/2503.05254.
- 64. Fosset, A., Bouchaud, J.-P. & Benzaquen, M. Endogenous liquidity crises. Journal of Statistical Mechanics: Theory and Experiment 2020, 063401 (2020).
- 65. Taleb, N. N. Statistical Consequences of Fat Tails: Real World Preasymptotics, Epistemology, and Applications 2022. arXiv: 2001.10488 [stat.OT]. https://arxiv.org/abs/2001.10488.
- 66. Bochud, T. & Challet, D. Optimal approximations of power-laws with exponentials 2006. arXiv: physics/0605149 [physics.data-an]. https://arxiv.org/abs/physics/0605149.
- 67. Loeb, T. F. Trading cost: the critical link between investment information and results. *Financial Analysts Journal* **39**, 39–44 (1983).
- 68. Almgren, R., Thum, C., Hauptmann, E. & Li, H. Direct estimation of equity market impact. *Risk* 18, 58–62 (2005).
- 69. Donier, J. & Bonart, J. A million metaorder analysis of market impact on the Bitcoin. *Market Microstructure and Liquidity* 1, 1550008 (2015).
- 70. Bucci, F., Mastromatteo, I., Benzaquen, M. & Bouchaud, J.-P. Impact is not just volatility. *Quantitative Finance* **19**, 1763-1766. https://hal.science/hal-02323182 (July 2019).

- 71. Gabaix, X., Gopikrishnan, P., Plerou, V. & Stanley, H. E. Institutional Investors and Stock Market Volatility*. *The Quarterly Journal of Economics* 121, 461-504. ISSN: 0033-5533. eprint: https://academic.oup.com/qje/article-pdf/121/2/461/5324363/121-2-461.pdf. https://doi.org/10.1162/qjec.2006.121.2.461 (May 2006).
- 72. Sato, Y. & Kanazawa, K. Does the square-root price impact law belong to the strict universal scalings?: quantitative support by a complete survey of the Tokyo stock exchange market 2024. arXiv: 2411.13965 [q-fin.TR]. https://arxiv.org/abs/2411.13965.
- Glosten, L. R. & Milgrom, P. R. Bid, ask and transaction prices in a specialist market with heterogeneously informed traders. *Journal of financial economics* 14, 71–100 (1985).
- 74. Donier, J., Bonart, J., Mastromatteo, I. & Bouchaud, J.-P. A fully consistent, minimal model for non-linear market impact. *Quantitative finance* **15**, 1109–1121 (2015).
- 75. Kyle, A. S. Continuous auctions and insider trading. *Econometrica: Journal of the Econometric Society*, 1315–1335 (1985).
- 76. Benzaquen, M. & Bouchaud, J.-P. Market impact with multi-timescale liquidity. *Quantitative Finance* **18**, 1781–1790 (2018).
- 77. Goshima, K., Tobe, R. & Uno, J. Trader classification by cluster analysis: Interaction between HFTs and other traders. Waseda University Institute for Business and Finance, Working Paper Series (2019).
- 78. Zarinelli, E., Treccani, M., Farmer, J. D. & Lillo, F. Beyond the square root: Evidence for logarithmic dependence of market impact on size and participation rate. *Market Microstructure and Liquidity* 1, 1550004 (2015).
- 79. Bucci, F., Benzaquen, M., Lillo, F. & Bouchaud, J.-P. Crossover from linear to square-root market impact. *Physical review letters* **122**, 108302 (2019).
- 80. Brokmann, X., Serie, E., Kockelkoren, J. & Bouchaud, J.-P. Slow decay of impact in equity markets. *Market Microstructure and Liquidity* 1, 1550007 (2015).
- 81. Hosaka, G. Analysis of high-frequency trading at tokyo stock exchange. Japan Exchange Group, JPX Working Papers, Tokyo, Japan, May (2014).
- 82. Oyama, A. & Tsuda, H. Characterizing High Frequency Trading (HFT) Through Algorithmic Criteria. *JAFEE Journal* **20**, 55–69 (2022).
- 83. Moro, E. et al. Market impact and trading profile of hidden orders in stock markets. Physical Review E—Statistical, Nonlinear, and Soft Matter Physics 80, 066102 (2009).

Bibliography

- 84. Jusselin, P. & Rosenbaum, M. No-arbitrage implies power-law market impact and rough volatility. *Mathematical Finance* **30**, 1309–1336 (2020).
- 85. Maitrier, G., Loeper, G. & Bouchaud, J.-P. in preparation 2025.
- 86. Gatheral, J. No-dynamic-arbitrage and market impact. *Quantitative finance* **10**, 749–759 (2010).
- 87. Bouchaud, J.-P. Price Impact. Encyclopedia of Quantitative Finance (2010).
- 88. Bouchaud, J.-P., Farmer, J. D. & Lillo, F. in *Handbook of financial markets:* dynamics and evolution 57–160 (Elsevier, 2009).
- 89. Eisler, Z., Bouchaud, J.-P. & Kockelkoren, J. The price impact of order book events: market orders, limit orders and cancellations. *Quantitative Finance* 12, 1395–1419 (2012).
- 90. Donier, J. & Bouchaud, J.-P. From Walras' auctioneer to continuous time double auctions: A general dynamic theory of supply and demand. *Journal of Statistical Mechanics: Theory and Experiment* **2016**, 123406 (2016).
- 91. Benzaquen, M., Mastromatteo, I., Eisler, Z. & Bouchaud, J.-P. Dissecting cross-impact on stock markets: An empirical analysis. *Journal of Statistical Mechanics: Theory and Experiment* **2017**, 023406 (2017).
- Hey, N., Mastromatteo, I., Muhle-Karbe, J. & Webster, K. Trading with Concave Price Impact and Impact Decay-Theory and Evidence. *Available* at SSRN (2023).
- 93. Bacry, E., Iuga, A., Lasnier, M. & Lehalle, C.-A. Market impacts and the life cycle of investors orders. *Market Microstructure and Liquidity* 1, 1550009 (2015).
- 94. Hasbrouck, J. Empirical market microstructure: The institutions, economics, and econometrics of securities trading (Oxford University Press, 2007).
- 95. Durin, B., Rosenbaum, M. & Szymanski, G. The two square root laws of market impact and the role of sophisticated market participants 2023. arXiv: 2311.18283 [q-fin.MF]. https://arxiv.org/abs/2311.18283.
- 96. Tsaknaki, I.-Y., Lillo, F. & Mazzarisi, P. Online learning of order flow and market impact with Bayesian change-point detection methods. *Quantitative Finance*, 1–16 (2024).
- 97. Nagy, P. et al. Generative ai for end-to-end limit order book modelling: A token-level autoregressive generative model of message flow using a deep state space network in Proceedings of the Fourth ACM International Conference on AI in Finance (2023), 91–99.

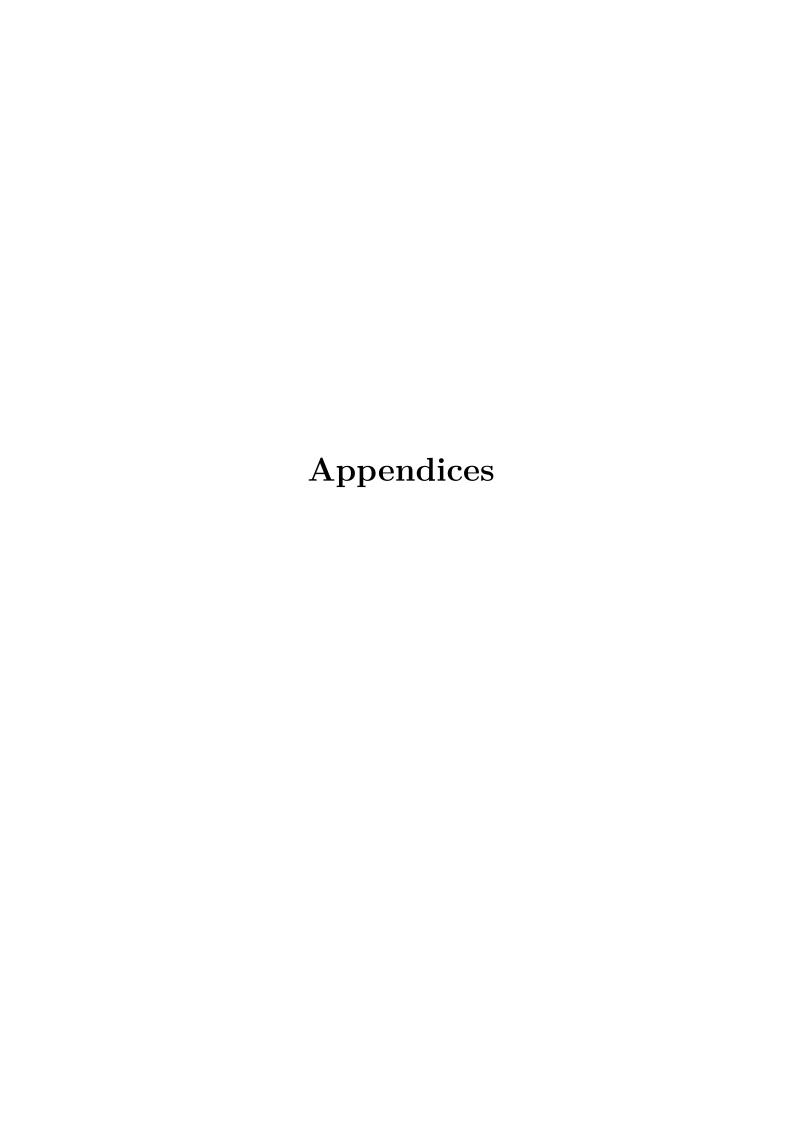
- 98. Coletta, A., Moulin, A., Vyetrenko, S. & Balch, T. Learning to simulate realistic limit order book markets from data as a world agent in Proceedings of the third acm international conference on ai in finance (2022), 428–436.
- 99. Cont, R., Cucuringu, M., Kochems, J. & Prenzel, F. Limit order book simulation with generative adversarial networks. *Available at SSRN 4512356* (2023).
- 100. Bucci, F., Benzaquen, M., Lillo, F. & Bouchaud, J.-P. Crossover from Linear to Square-Root Market Impact. *Physical Review Letters* 122. ISSN: 1079-7114. http://dx.doi.org/10.1103/PhysRevLett.122.108302 (Mar. 2019).
- 101. Sato, Y. & Kanazawa, K. Exactly solvable model of the square-root price impact dynamics under the long-range market-order correlation 2025. arXiv: 2502.17906 [q-fin.TR]. https://arxiv.org/abs/2502.17906.
- 102. Sato, Y. & Kanazawa, K. Quantitative statistical analysis of order-splitting behavior of individual trading accounts in the Japanese stock market over nine years. *Physical Review Research* 5. ISSN: 2643-1564. http://dx.doi.org/10.1103/PhysRevResearch.5.043131 (Nov. 2023).
- 103. Hasbrouck, J. Measuring the information content of stock trades. *The Journal of Finance* **46**, 179–207 (1991).
- 104. Bouchaud, J.-P., Kockelkoren, J. & Potters, M. Random walks, liquidity molasses and critical response in financial markets. *Quantitative finance* 6, 115–123 (2006).
- 105. Hopman, C. Do supply and demand drive stock prices? Quantitative Finance 7, 37–53 (2007).
- 106. Van der Beck, P., Bouchaud, J.-P. & Villamaina, D. Ponzi funds. arXiv preprint arXiv:2405.12768 (2024).
- 107. Black, F. Noise. The journal of finance 41, 528–543 (1986).
- 108. Odean, T. Do investors trade too much? American economic review 89, 1279–1298 (1999).
- 109. LeRoy, S. F. Excess volatility. The New Palgrave Dictionary of Economics, 2nd Edition. Palgrave Macmillan 13 (2006).
- 110. Bucci, F. et al. Co-impact: Crowding effects in institutional trading activity. Quantitative Finance 20, 193–205 (2020).
- 111. Wyart, M. & Bouchaud, J.-P. Statistical models for company growth. *Physica A: Statistical Mechanics and its Applications* **326**, 241–255. ISSN: 0378-4371. http://dx.doi.org/10.1016/S0378-4371(03)00267-X (Aug. 2003).

- 112. Moran, J., Secchi, A. & Bouchaud, J.-P. Revisiting Granular Models of Firm Growth 2024. arXiv: 2404.15226 [econ.GN]. https://arxiv.org/abs/ 2404.15226.
- 113. Muzy, J.-F., Delour, J. & Bacry, E. Modelling fluctuations of financial time series: from cascade process to stochastic volatility model. *The European Physical Journal B-Condensed Matter and Complex Systems* 17, 537–548 (2000).
- 114. Wyart, M., Bouchaud, J.-P., Kockelkoren, J., Potters, M. & Vettorazzo, M. Relation between bid-ask spread, impact and volatility in order-driven markets. *Quantitative finance* 8, 41–57 (2008).
- 115. Plerou, V., Gopikrishnan, P., Gabaix, X. & Stanley, H. E. Quantifying stock-price response to demand fluctuations. *Physical review E* **66**, 027104 (2002).
- 116. Evans, M. D. & Lyons, R. K. Order flow and exchange rate dynamics. *Journal of political economy* **110**, 170–180 (2002).
- 117. Chordia, T. & Subrahmanyam, A. Order imbalance and individual stock returns: Theory and evidence. *Journal of Financial Economics* **72**, 485–518 (2004).
- 118. Gomes, C. & Waelbroeck, H. Is market impact a measure of the information value of trades? Market response to liquidity vs. informed metaorders. *Quantitative Finance* **15**, 773–793 (2015).
- 119. Toth, B., Palit, I., Lillo, F. & Farmer, J. D. Why is equity order flow so persistent? *Journal of Economic Dynamics and Control* **51**, 218–239 (2015).
- 120. Farmer, J. D. Market force, ecology and evolution 2000. arXiv: adap-org/9812005 [adap-org]. https://arxiv.org/abs/adap-org/9812005.
- 121. Elomari, S. Modelling of the Limit Order Book: From Statistical Methods to Machine Learning Techniques Theses (Institut Polytechnique de Paris, Dec. 2024). https://theses.hal.science/tel-04966271.
- 122. Maitrier, G. & Bouchaud, J.-P. The Subtle Interplay between Square-root Impact, Order Imbalance\& Volatility: A Unifying Framework. arXiv preprint arXiv:2506.07711 (2025).
- 123. Derrida, B. Random-energy model: An exactly solvable model of disordered systems. *Physical Review B* **24**, 2613 (1981).
- 124. Bouchaud, J.-P. & Mézard, M. Universality classes for extreme-value statistics. *Journal of Physics A: Mathematical and General* **30**, 7997 (1997).
- 125. Lillo, F. & Farmer, J. D. The long memory of the efficient market. Studies in nonlinear dynamics & econometrics 8, 20123001 (2004).

- 126. Farmer, J. D., Patelli, P. & Zovko, I. I. The predictive power of zero intelligence in financial markets. *Proceedings of the National Academy of Sciences* **102**, 2254–2259 (2005).
- 127. Bouchaud, J.-P., Bonart, J., Donier, J. & Gould, M. in *Trades, Quotes and Prices: Financial Markets Under the Microscope* 208–228 (Cambridge University Press, 2018).
- 128. Hultin, H., Hult, H., Proutiere, A., Samama, S. & Tarighati, A. A generative model of a limit order book using recurrent neural networks. *Quantitative Finance* 23, 931–958 (2023).
- 129. Coletta, A. et al. Towards realistic market simulations: a generative adversarial networks approach in Proceedings of the Second ACM International Conference on AI in Finance (2021), 1–9.
- 130. Coletta, A., Jerome, J., Savani, R. & Vyetrenko, S. Conditional generators for limit order book environments: Explainability, challenges, and robustness in Proceedings of the Fourth ACM International Conference on AI in Finance (2023), 27–35.
- 131. Hardiman, S. J. & Bouchaud, J.-P. Branching-ratio approximation for the self-exciting Hawkes process. *Physical Review E* **90**, 062807 (2014).
- 132. Joulin, A., Lefevre, A., Grunberg, D. & Bouchaud, J.-P. Stock price jumps: news and volume play a minor role. *Wilmott Magazine* **46** (2008).
- 133. Bouchaud, J.-P. The endogenous dynamics of markets: Price impact, feed-back loops and instabilities. *Lessons from the credit crisis*, 345–74 (2011).
- 134. Aubrun, C., Morel, R., Benzaquen, M. & Bouchaud, J.-P. Riding Wavelets: A Method to Discover New Classes of Price Jumps. arXiv preprint arXiv:2404.16467 (2024).
- 135. Boulton, A. J. & Williford, A. Analyzing Skewed Continuous Outcomes With Many Zeros: A Tutorial for Social Work and Youth Prevention Science Researchers. *Journal of the Society for Social Work and Research* 9, 721–740. eprint: https://doi.org/10.1086/701235. https://doi.org/10.1086/701235 (2018).
- 136. Neelon, B., O'Malley, A. J. & Smith, V. A. Modeling zero-modified count and semicontinuous data in health services research Part 1: background and overview. Statistics in Medicine 35, 5070-5093. eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/sim.7050. https://onlinelibrary.wiley.com/doi/abs/10.1002/sim.7050 (2016).
- 137. Liu, L. *et al.* Statistical analysis of zero-inflated nonnegative continuous data. *Statistical Science* **34**, 253–279 (2019).

Bibliography

- 138. Bouchaud, J.-P. *Price impact. Encyclopedia of Quantitative Finance. 2010* (Wiley. AQ4).
- 139. Smith, E., Farmer, J. D., Gillemot, L. & Krishnamurthy, S. Statistical theory of the continuous double auction. *Quantitative finance* **3**, 481 (2003).
- 140. Mastromatteo, I., Tóth, B. & Bouchaud, J.-P. Agent-based models for latent liquidity and concave price impact. *Phys. Rev. E* **89**, 042805. https://link.aps.org/doi/10.1103/PhysRevE.89.042805 (4 Apr. 2014).
- 141. Fisher, M. Critical Phenomena Proc. 51st Enrico Fermi Summer School, ed. MS Green 1971.



Microstructure et lois du marché

Les marchés financiers jouent un rôle central dans les économies modernes. Ils assurent la rencontre entre offreurs et demandeurs de capitaux, attribuent une valeur aux actifs et permettent plus généralement le financement de l'activité économique. La théorie économique dominante les conçoit comme des marchés efficients, où les prix reflètent instantanément toute l'information disponible. Pourtant, l'expérience historique et l'observation empirique contredisent cette vision idéalisée. Les crises financières, bulles spéculatives et l'excès de volatilité sont autant de phénomènes qui surviennent régulièrement sans qu'un choc d'information exogène ne puisse les expliquer. Ces anomalies révèlent le rôle déterminant des mécanismes endogènes dans la formation des prix. C'est précisément le domaine de la microstructure de marché, qui étudie les règles et interactions élémentaires à l'échelle du carnet d'ordres électronique. Elle s'intéresse directement au flux d'ordres d'achat et de vente, et au rôle de la liquidité dans les dynamiques de prix. Cette approche s'inscrit dans une perspective d'*Econophysique*, qui mobilise les outils de la physique statistique pour rechercher des lois universelles. Et ici, bien que les marchés financiers soient souvent considérés comme chaotiques ou imprévisibles, ils semblent néanmoins obéir à certaines régularités, notamment la loi empirique dite "en racine carrée", que nous détaillerons par la suite. Cette loi montre que l'évolution des prix présente des comportements quantitatifs analogues à ceux des systèmes physiques.

Métaordres et impact mécanique

Le coeur de cette thèse est consacré à l'étude du phénomène d'impact de marché, c'est-à-dire la réponse des prix lorsque des transactions sont effectuées. Grâce à une base de données unique provenant de la Bourse de Tokyo, incluant les identifiants des traders, il a été possible d'analyser de façon fine le rôle des métaordres,

ces séquences d'ordres élémentaires émises par un même agent dans la même direction. Les résultats montrent que la loi en racine carrée reste valide même lorsque l'on zoom à l'intérieur d'un métaordre, ou que l'on détruit le lien entre identité des traders et trade effectué, et que l'on reconstruit des métaordres de manière aléatoire. L'impact apparaît ainsi comme un phénomène purement mécanique, indépendant de l'information, ce qui contredit fortement la vision traditionnelle. A partir de cette intuition, cette thèse propose une méthode de reconstruction de métaordres à partir de données publiques, permettant de reproduire de manière réaliste les faits stylisés documentés dans la littérature. Ces "métaordres synthétiques" ouvrent la voie à une recherche plus ouverte et reproductible, affranchie des contraintes liées aux données propriétaires. Sur cette base, un cadre théorique unifié a été développé combinant les caractéristiques statistiques des flux d'ordres, des prix et la loi en racine carrée et introduisant par là un propagateur généralisé qui pourrait remplacer la "boîte noire" de complexité qu'était auparavant la microstructure de marché. Les simulations confirment la validité de ce cadre, qui démontre que la volatilité émerge de manière endogène, comme conséquence mécanique des interactions d'ordres.

Dynamique critique et instabilité

Au-delà de l'impact des transactions, la thèse explore la question de la stabilité des marchés financiers. Un premier axe consiste à élargir la focale, considérer tous les types d'ordres et utiliser un modèle vectoriel autorégressif appliqué à plusieurs variables du carnet d'ordres. L'analyse révèle l'existence de modes de microstructure distincts, certains symétriques associés à la liquidité, d'autres antisymétriques associés aux déséquilibres directionnels. Ces modes permettent d'anticiper l'évolution du marché et mettent en évidence son caractère marginalement instable : le système opère souvent à une frontière où de petites perturbations peuvent entraîner un assèchement de la liquidité et un élargissement du spread. Dans un second temps, un modèle multi-agents inspiré du modèle de Santa Fe est développé et enrichi de rétroactions réalistes, telles que l'effet de levier et la dépendance à la volatilité passée. Ce modèle montre l'émergence de transitions de phase endogènes dans le carnet d'ordres. L'analyse numérique et analytique de ces transitions permet de caractériser une véritable frontière de stabilité du marché et de déterminer les exposants critiques liés à ces transitions de phase du second ordre. Ainsi, la thèse établit un pont entre la dynamique microstructurelle, les propriétés statistiques de l'impact et la stabilité globale du système financier. Elle suggère que les marchés, loin d'être des systèmes parfaitement efficients, sont des systèmes complexes qui fonctionnent parfois proche d'un régime critique. La stabilité des marchés n'est donc jamais acquise mais sans cesse menacée par leurs

propres dynamiques internes.

Appendix A

Chapter 4: the Metaorder Proxy

Sanity Checks and Validation of the Algorithm

We provide here a non-exhaustive list of sanity checks one may perform easily to validate the coherence of generated metaorders.

Validating the Consistency of Metaorder Size and Duration

We know that the Square Root Law is independent of the specific microscopic parameters of a market, such as the distributions of metaorder size and duration, as discussed in [42]. Their empirical analysis is particularly important to support the universalism of such theory. However, even though it may not modify the impact function, it may matter to verify that our generated metaorders are statistically coherent, see Figure A.1.

Appendix A. Chapter 4: the Metaorder Proxy

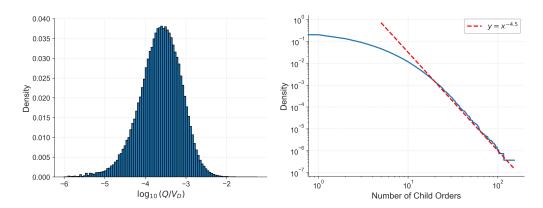


Figure A.1: Left: Distribution of the size of synthetic metaorders, for BNP Paribas, data from 2021 to 2023 generated with a mapping function of parameters 10 traders and homogeneous distribution. **Right:** On the same dataset, distribution of the number of child orders per metaorders, fitted by a power law of exponent $1 + \mu = 4.5$.

With this algorithm and those parameters, we obtained metaorders distributed around $10^{-3}V_D$, which is typically the order of magnitude one may expect for a metaorder. We also retrieve a power law distribution of the length of metaorders, in line with the Lillo-Mike-Farmer theory [18] and empirical studies, see [16]. However, one may note that the decay is much faster than expected, with $\mu = 3.5$ whereas theory suggests that $1 \le \mu_t \le 2$.

Modifying public trade data

To ensure that our algorithm is not merely generating random metaorders with the correct impact due to an unknown construction bias, we conducted several tests that removed information from the public data. One of the simplest consists in randomly shuffling trade signs, referring to buy or sell, in the public data, and then perform again the aggregation method.

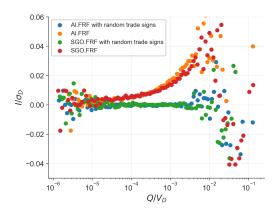


Figure A.2: Comparison of the impact law for random metaorders created using public data versus trade data where the signs of market orders have been randomly shuffled. As expected, the average impact is null when the real signs time series is modified. Data from 2020 to 2023, generated using a mapping function with parameters (20 traders, homogeneous distribution).

We also verified that randomly modifying the chronology of public trades execution, keeping real permanent impact for each trade affects the impact law. More generally, to obtain the correct impact, it is important to keep the original trade data as it is, indicating that some information is indeed present. Our argument is that this information is not associated with economic considerations but rather with flow reactions and liquidity responses.

Synthetic metaorders constructed on synthetic prices

A natural question arising from this empirical study is whether public trade data are actually needed to construct synthetic metaorders that verify the SQL. In other words, is there a specific element or piece of information in real trade data that distinguishes actual prices from random prices?

A possible experiment is to generate a basic synthetic price. This can be achieved by first creating a sequence of random volumes and random signs. Each order is then assigned a signed instantaneous impact, which can be a linear or concave function of the volume. Then, returns read:

$$r_t = \varepsilon_t q^{\chi}, \quad \chi \in \{0, \frac{1}{2}, 1\}, \quad q \sim \mathcal{U}(q_{min}, q_{max})$$
 (A.1)

In the following simulation we set $q_{min} = 1$ and $q_{max} = 100$, but the same results hold for $q = q_{min} = q_{max} = 1$, just generating synthetic metaorders with a lower average size and variance. Based on [2], we fixed $\chi = 0.5$.

Appendix A. Chapter 4: the Metaorder Proxy

This process produces synthetic trade data that can serve as a foundation for constructing random metaorders. However, metaorders constructed with this algorithm won't exhibit SQL, but a linear impact function, see Figure A.3.

Since trade signs are known to be autocorrelated, it is tempting to leverage this well-established stylized fact to recover the square root law. We generate a series of trade signs with an autocorrelation $C(l) = \langle \varepsilon_t \varepsilon_{t+l} \rangle = l^{-\gamma}$ and we naturally impose a propagator-type impact $G(t) \approx t^{-\beta}$, $\beta = \frac{1-\gamma}{2}$, see [7], to recover price diffusivity. But again, by doing so, one will still obtain a linear impact for synthetic metaorders, see Figure A.3.

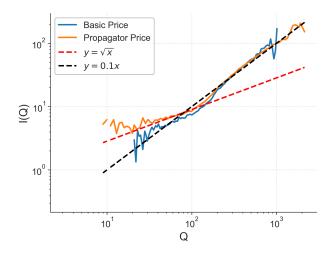


Figure A.3: Retrieving a linear impact when constructing synthetic metaorders on a synthetic price. The basic price is built using uniformly distributed volumes $(q_{min} = 1 \text{ and } q_{max} = 100 \text{ and a series of uncorrelated random signs.}$ The instantaneous impact reads $r_t = \varepsilon_t q^{\chi}$ with $\chi = 0.5$. Synthetic metaorders were constructed with 4 traders and a power law distribution of exponent $\chi = 2$. The propagator price refers to a synthetic price constructed with autocorrelated signs, $C(l) = \langle \varepsilon_t \varepsilon_{t+l} \rangle = l^{-\gamma}$ and a transient impact that decays according to the propagator model: $G(l) = t^{-\beta}$

Appendix B

Chapter 8: the VAR model

Aggregated Impact

Inspired by the work of F. Patzelt and one of us (JPB) [37], we quantify the relationship between the aggregated execution imbalance and its impact on the price for a bin size N. In a similar way, let us define the aggregate-imbalance impact for N consecutive observed price changes:

$$\mathcal{R}_N(\mathcal{I}_N) = \left\langle m_{t+i} - m_t | \mathcal{I}_N = \sum_{i=0}^{N-1} V_i^{ex,a} - V_i^{ex,b} \right\rangle.$$
 (B.1)

As in [37], we write:

$$\mathcal{R}_N(\mathcal{I}) = g(N)\mathcal{F}_{\alpha,\beta}\left(\frac{\mathcal{I}}{h(N)}\right),$$
 (B.2)

where g(N) and h(N) are the appropriate scaling of the return and the imbalance for a bin size N, and $\mathcal{F}_{\alpha,\beta}$ is a sigmoidal parametric function

$$\mathcal{F}_{\alpha,\beta}(x) = \frac{x}{(1+|x|^{\alpha})^{\frac{\alpha}{\beta}}}.$$
 (B.3)

Appendix B. Chapter 8: the VAR model

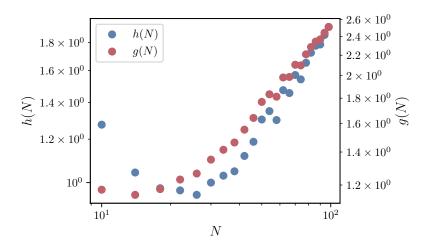


Figure B.1: Evolution of the scaling of the impact and of the return. Starting from N=20, the evolution of the scaling is stable.

After calibrating of the parameters of (B.2), the rescaled aggregated impact is the same for all the bin sizes N, as one can see in Figs. B.1 and B.2.

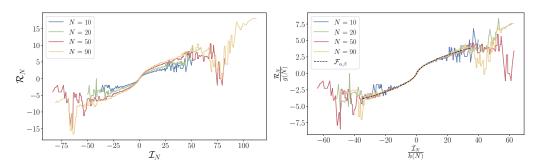


Figure B.2: Left: Aggregated imbalance impact on raw data, for different bin sizes N before the rescaling. Right: Aggregated imbalance impact on raw data, for different bin sizes N after the rescaling.

It is interesting to note that the universality of the aggregate impact holds even for price change-by-price change data, although the scaling of the returns and the impact no longer follows a pure power law.

Appendix C

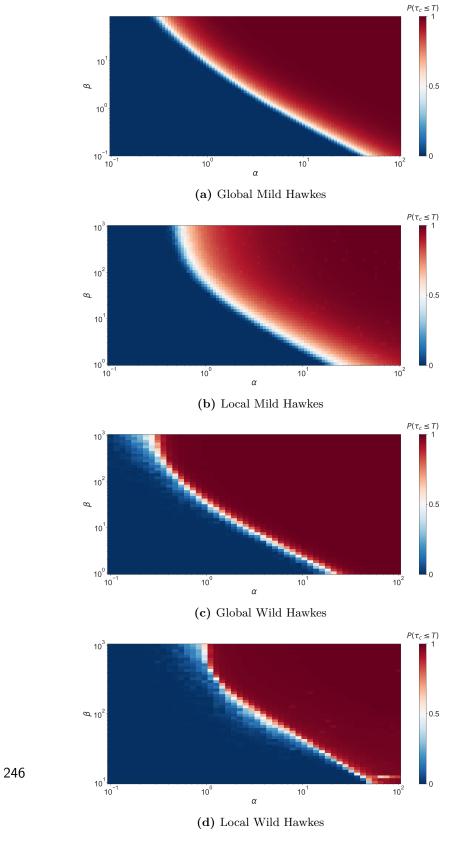
Chapter 9: the Santa-Fe like model

Stability and Scaling in Feedback Models

Overview. This appendix provides additional numerical results supporting the analysis of various feedback mechanisms. We focus on the stability landscape, finite-size susceptibility scaling, and spread behavior for models with Hawkes-type and Zumbach-type feedback kernels. All results are computed for varying model parameters and interaction strengths. These plots are part of ongoing work, and we plan to revisit and refine them shortly.

Stability Maps

We display the stability maps for four variants of the generalized Hawkes and Zumbach feedbacks. Simulations are performed with $N=140, T=500, \lambda=10, \mu=20, \nu_0=1$.



 ${\bf Figure~C.1:~Stability~maps~for~the~four~Hawkes~kernel~feedback~models.}$

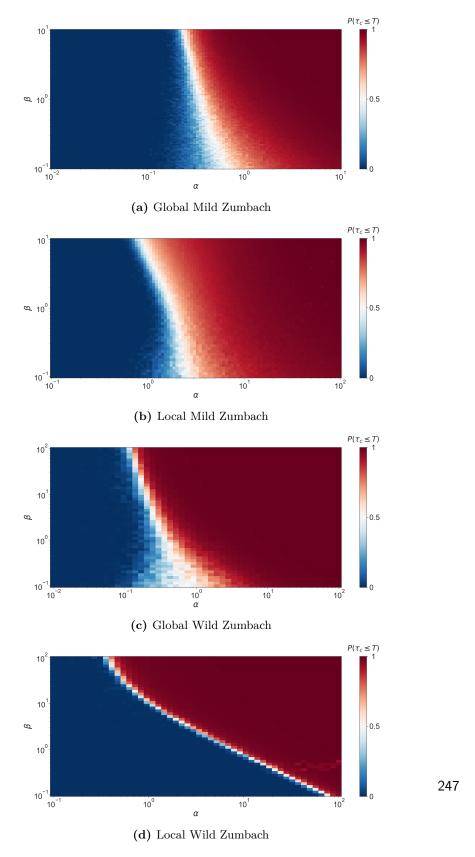


Figure C.2: Stability maps for the four Zumbach kernel feedback models.

Finite-Size Scaling of the Susceptibility

Hawkes Kernel

These plots show the scaling of the system susceptibility with system size for the Hawkes-based feedback models. Again, simulations are done with $\lambda=10,\,\mu=20,$ and $\nu_0=1.$

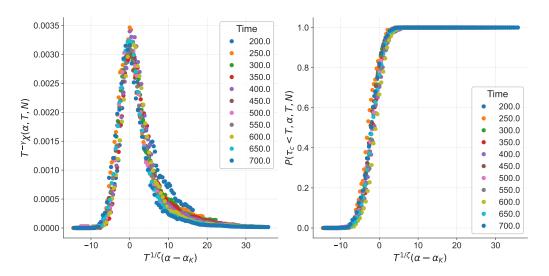


Figure C.3: Susceptibility scaling for GMH.

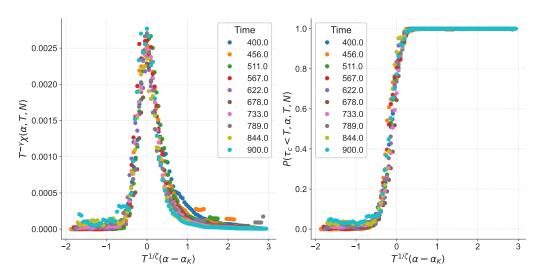


Figure C.4: Susceptibility scaling for LMH.

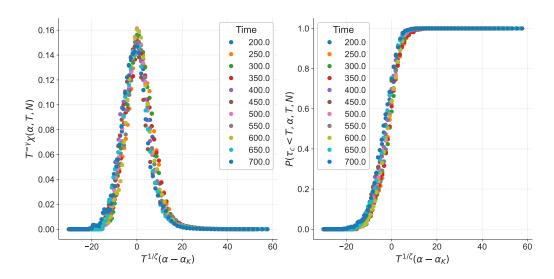


Figure C.5: Susceptibility scaling for GWH.

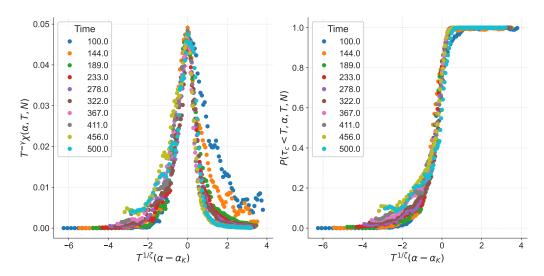


Figure C.6: Susceptibility scaling for LWH.

Appendix C. Chapter 9: the Santa-Fe like model

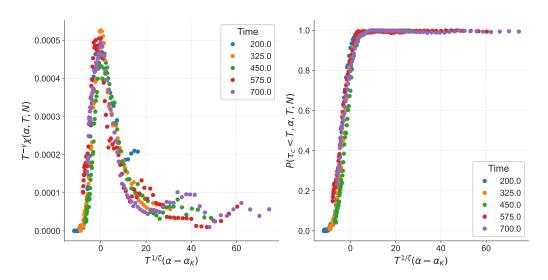


Figure C.7: Susceptibility scaling for LWS.

Zumbach Kernel

Same analysis for Zumbach-based models. Simulations are done with $\lambda=10,$ $\mu=20,$ and $\nu_0=1.$

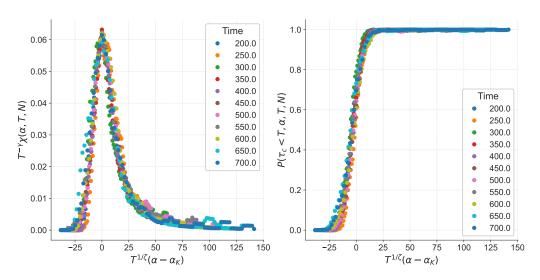


Figure C.8: Susceptibility scaling for GMZ.

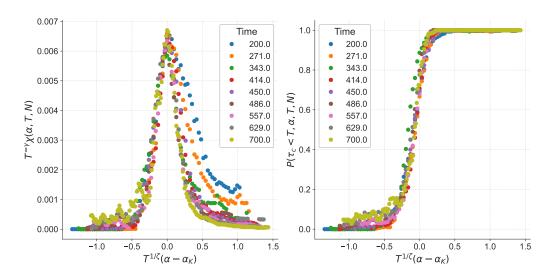


Figure C.9: Susceptibility scaling for LMZ.

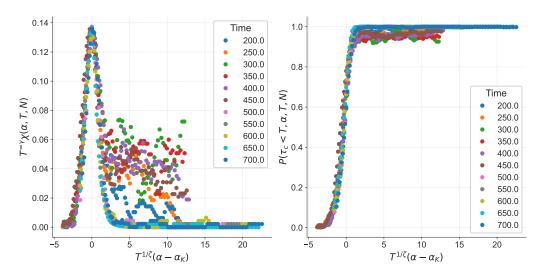
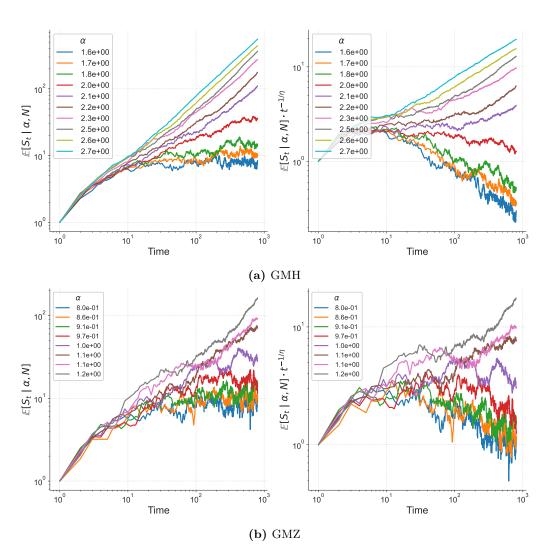


Figure C.10: Susceptibility scaling for LWZ.

Spread Scaling

Finally, we analyze how the average spread varies with system parameters and model type close to criticality. Each figure shows the scaling for a specific model. Simulations are performed for $\lambda = 10, \nu_0 = 1, \beta = 1$ and N = 1000

Appendix C. Chapter 9: the Santa-Fe like model





Titre : La Complexité Cachée de la Formation des Prix : Exploration des Mécanismes Microstructuraux des Marchés Financiers

Mots clés: Physique Statistique, Microstructure de Marché, Impact de Marché, Transition de Phase

Résumé: Les marchés financiers constituent une pièce maîtresse des économies modernes. Ils facilitent les échanges entre agents, le financement des institutions, et attribuent un prix, une valeur aux titres échangés. Si la théorie économique traditionnelle postule que ces marchés sont efficients – c'est-à-dire qu'ils intègrent instantanément et parfaitement toute l'information disponible – la réalité semble tout autre, en témoignent les nombreuses crises financières, bulles spéculatives, sauts brutaux de prix ou encore l'excès de volatilité observés fréquemment mais toujours incompris.

L'objectif de cette thèse est de mieux comprendre le processus de formation des prix, en adoptant une perspective microstructurelle, soit en étudiant les mécanismes les plus élémentaires à l'œuvre dans les carnets d'ordres électroniques. Nous commençons par analyser une base de données unique, incluant les identifiants des traders, issue de la Bourse de Tokyo. Cette étude met en lumière de nouveaux faits stylisés concernant l'impact: l'effet mécanique des ordres d'achat ou de vente sur les prix.

À la suite de cette analyse empirique, nous proposons un cadre théorique unifié, reposant sur des hypothèses minimales. Ce modèle parvient à concilier les propriétés statistiques du flux d'ordres, celles de la dynamique des prix, ainsi que la relation non triviale entre ces deux processus. Il permet également de formuler des prédictions testables

que nous validons sur données réelles. Par ailleurs, nous développons plusieurs algorithmes permettant de générer des marchés simulés réalistes, statistiquement indiscernables des marchés réels, et ce sans avoir à ajouter une quelconque information économique exogène dans le système. Ces outils non seulement appuient nos théories par leur efficacité surprenante, mais aussi ouvrent de nouvelles perspectives pour la recherche académique et industrielle, souvent freinée par l'inaccessibilité aux données, propriété privée des acteurs de marchés.

Nous introduisons ensuite un modèle de liquidité markovienne, révélant l'existence de "modes de microstructure". Ces modes permettent d'étudier à la fois la stabilité du marché, mais aussi de prédire les futurs flux d'ordres et variations de prix. Ils mettent en évidence le caractère marginalement stable des marchés et montrent que crises et sauts de prix anormaux sont déjà prédits à l'échelle microscopique, de manière purement endogène. Pour mieux comprendre ces instabilités, nous développons dans une dernière partie un modèle d'agents, inspiré du modèle de Santa Fe quadratique et enrichi par des rétroactions plus réalistes. Nous montrons numériquement l'émergence de transitions de phase dans le carnet d'ordres en déterminant numériquement leurs exposants critiques par une analyse de taille finie, et analytiquement leur frontière de stabilité.

Title: The Hidden Complexity of Price Formation: Exploring Microstructural Mechanisms in Financial Markets

Keywords: Statistical Physics, Market Microstructure, Price Impact, Phase Transition

Abstract : Financial markets are a cornerstone of modern economies. They facilitate exchanges between agents, enable institutional financing, and assign a price—a value—to traded securities. While traditional economic theory posits that these markets are efficient—meaning they instantaneously and perfectly incorporate all available information—reality often tells a different story, as evidenced by the many financial crises, speculative bubbles, abrupt price jumps, and excessive volatility that frequently occur yet remain poorly understood.

The objective of this thesis is to improve our understanding of the price formation process by adopting a microstructural perspective—that is, by studying the most elementary mechanisms at work within electronic order books. We begin by analyzing a unique dataset from the Tokyo Stock Exchange, which includes trader identifiers. This empirical study reveals new stylized facts about market impact—that is, the mechanical effect of buy and sell orders on prices.

Following this empirical analysis, we propose a unified theoretical framework based on minimal assumptions. This model successfully reconciles the statistical properties of order flow, those of price dynamics, and the nontrivial relationship between the two. It also yields testable predictions that we

validate using real market data. Moreover, we develop several algorithms capable of generating realistic simulated markets that are statistically indistinguishable from actual ones, without requiring any exogenous economic information. These tools not only support our theoretical claims through their surprising effectiveness, but also open new avenues for academic and industrial research, which is often hindered by the inaccessibility of privately owned market data.

We then introduce a Markovian liquidity model that reveals the existence of "microstructural modes." These modes enable the study of market stability as well as the prediction of future order flows and price movements. They highlight the marginal stability of markets and show that crises and abnormal price jumps are already foreshadowed at the microscopic scale, in a purely endogenous manner. To better understand these instabilities, we conclude by developing an agent-based model inspired by the quadratic Santa Fe model, enriched with more realistic feedback mechanisms. We numerically demonstrate the emergence of phase transitions within the order book, estimating their critical exponents through finite-size scaling analysis and deriving their stability boundaries analytically.

