# ON DETECTING SPOOFING STRATEGIES IN HIGH FREQUENCY TRADING

XUAN TAO, ANDREW DAY, LAN LING, AND SAMUEL DRAPEAU

ABSTRACT. Spoofing is an illegal act of artificially modifying the supply to drive temporarily prices in a given direction for profit. In practice, detection of such an act is challenging due to the complexity of modern electronic platforms and the high frequency at which orders are channeled. We present a micro-structural study of spoofing in a simple static setting. A multilevel imbalance which influences the resulting price movement is introduced upon which we describe the optimization strategy of a potential spoofer. We provide conditions under which a market is more likely to admit spoofing behavior as a function of the characteristics of the market. We describe the optimal spoofing strategy after optimization which allows us to quantify the resulting impact on the imbalance after spoofing. Based on these results we calibrate the model to real Level 2 datasets from TMX, and provide some monitoring procedures based on the Wasserstein distance to detect spoofing strategies in real time.

Keywords: Spoofing, High Frequency Trading, Imbalance, Limit Order Book.

## 1. INTRODUCTION

The act of spoofing is a specific trading activity that aims at artificially modifying the supply on the market, without intend to trade, to move it away from its equilibrium. One might profit from the resulting short term price movement by canceling the previous supply while the market comes back to its equilibrium. Such a strategy implies that the spoofer should be able to act anonymously, fast and in a market where all the other agents react to offer and demand. In this regard, with the recent rise of centrally cleared venue and high frequency algo-trading, the ground for the existence of spoofing schemes is rising, see Shorter and Miller [22].[1] In a competitive market where many potential spoofers are present, spoofing behavior might cancel out, but most regulations consider it as illegal. For instance, the 2010 Dodd-Frank Act prohibits spoofing – defined as activity of bidding or offering with the intent to cancel before execution – that can be prosecuted as "a felony punishable by up to $1 million in penalties and up to ten years in prison for each spoofing count".[2]

Yet, for several reasons, detecting and prosecuting spoofing behavior is a challenging problem. First is the sheer amount of data produced from high frequency trading across many financial products and venues. Second, it is usually impossible to trace in real time who

[1]In 2010, trader Navinder Singh Sarao was accused of exacerbating a flash crash by placing thousands of E-mini S&P 500 stock index futures contract orders in one day and changed or moved those orders more than 20 million times before they were cancelled.
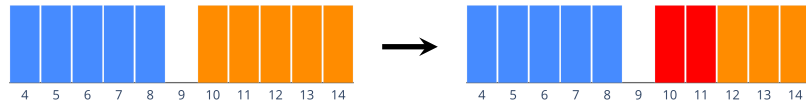
[2]In 2019, the high frequency company Tower Research Capital agreed to pay a fine of about $60 million over spoofing allegations. In 2020, JP Morgan settles spoofing lawsuit alleging fraud for about $920 million.

is behind every trade. From a CCP viewpoint, they mainly have access to the broker ID through which the trade has been channeled resulting only in aggregated informations. Furthermore, a potential spoofer might post those trades through different venues and brokers. Third, aside from a loose definition, it is unclear how a spoofing strategy differs quantitatively from other strategies and what is the resulting impact: Thus, the complexity of quantifying and discriminating spoofing strategies from legitimate ones. Finally, due to the lack of information from the few regulatory cases, the problem amount to an unsupervised classification problem. Based on the above points, it seems difficult to provide an efficient way to monitor the market for spoofing behavior.[3] However, a potential spoofer is also confronted to the constraints of modern electronic trading platforms. Indeed, from the basic spoofing description, the spoofer has to act rapidly in a complex and high frequency environment. Therefore, it must rely on fast – henceforth simple – algorithmic strategies which, due to the complexity of the dynamic structure of a limit order book, is based on aggregated signal.

Along these lines, and in view of this unsupervised classification task, we intentionally address a quantitative analysis of spoofing in a simple setup. As a basis for this study, let us consider a simple example. We suppose that in the next period the limit order book shifts up by one unit with a probability $\bar{p}$ and down by one unit otherwise. From the perspective of an agent whose objective is to purchase two shares, it faces the following three idealized situations.

1 - Immediately post a buy market order for a total cost of
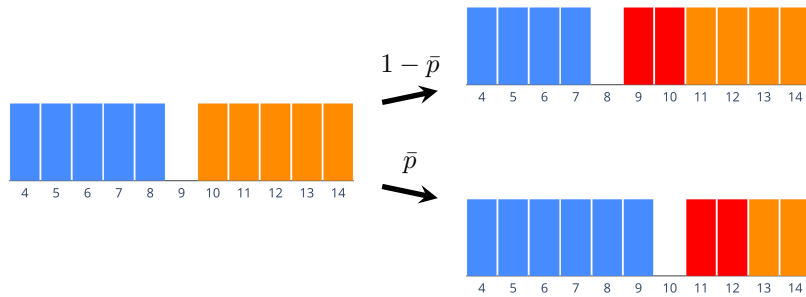
$$\hat{C} = 10 + 11 = 21$$



2 - Delay the buy market order for one period resulting in a total average cost of

$$\tilde{C} = (1 - \bar{p})(9 + 10) + \bar{p}(11 + 12) = (1 - \bar{p})19 + \bar{p}23$$

which is smaller than $\hat{C}$ if and only if $\bar{p} < 1/2$, hence a bearish market.



3 - Delay the buy market order and post a sell limit order of one share at a distance of one unit from the best ask price to artificially modify the offer and demand resulting in a temporary more bearish state $p < \bar{p}$.[4] Doing so, with a probability $q$, its sell limit order is executed through an incoming market order walking the limit order book beyond one unit. For this executed sell order, the agent receives an average price of $q((1-p)10 + p12)$

---

[3]Aside from obvious cases or exogenous approaches as for insider prosecution.

[4]We suppose that during this time period, if the market moves through limit order posting/canceling, the agent keeps its sell limit order one unit away from the best ask price by rapidly canceling and posting again.

while its inventory increases in average to $2 + q$. The cost of buying back this increased inventory minus the gain from selling its limit order results in an average net cost of

$$C = (1 - p)\,(9 + 10 + q11) + p\,(11 + 12 + q13) - q\,((1 - p)10 + p12)$$
$$= (1 - p)19 + p23 + q$$



From this simple situation, a profitable spoofing situation depends on the value of $\bar{p}$, how bearish $p < \bar{p}$ the market reacts to an increase of one share at a distance of one unit from the best ask price, as well as the probability $q$ of being adversarially executed through a market order during this time. The bottom line is this paper is to study the interplay between these different factors and the resulting impact in view of detection procedure, in particular the impact on the price movement as a function of the spoofing size as well as the depth in the limit order book. Based on this theoretical approach, we present some approaches to track spoofing behavior and calibrate those to real market data.

We model the impact of the offer and demand on the price through the volume imbalance often taken as the ratio of the volume of the best bid divided by the total volume on the best bid and ask. Since the spoofer never intends to have their orders executed, spoofing is more likely to happen beyond the top of the limit order book. Indeed, the probability of getting executed is too high, resulting in a negative payoff. To take this into account we weight the impact on the imbalance in terms of depth as follows

$$\bar{\imath} = \frac{\sum \bar{v}_k^- \, w_k}{\sum \left( \bar{v}_k^- + \bar{v}_k^+ \right) w_k}$$

where $\bar{v}_k^{\pm}$ represents the volume on the bid/ask $k$ units away from the best bid/ask and $w_k$ is the relative impact on the imbalance at level $k$. If the agent posts a sell limit order $v$ on the ask side at tick level $k_0$, the imbalance moves to

$$\imath(v) = \frac{\sum \bar{v}_k^- \, w_k}{\sum \left( \bar{v}_k^- + \bar{v}_k^+ \right) w_k + w_{k_0} v} \leq \bar{\imath}$$

If we denote by $dp_n$ the probability of a price deviation of $n$ units, the dependence on the imbalance $\imath$ is given as follows

$$dp_n = \imath dp_n^+ + (1 - \imath) dp_n^-, \quad n = \ldots, -2, -1, 0, 1, 2, \ldots$$

where $dp_n^{\pm}$ represents the price deviation when the imbalance is at its extreme. When the imbalance $\imath$ is low/high – the offer/demand dominates – the price distribution is biased downwards/upwards through $dp^{\pm}$. The agent can influence the price distribution in a non linear way through the volume it posts:

$$v \longmapsto dp(v) := \imath(v) dp^+ + (1 - \imath(v)) \, dp^-$$

Given the probability $dq$ of a sell limit order hitting the limit order book up to a given level the resulting net average cost of spoofing turns out to be

$$C(v) = pH + \underbrace{(1-Q)G(H)}_{\text{Cost for optimal situation}} + \underbrace{H\mu^+ (2\imath(v) - 1)}_{\text{Spoofing impact}} + \underbrace{QG(H+v) + v\nu}_{\text{Cost for being caught wrong way}}$$

where $H$ is the initial objective of shares to acquire, $G$ is the liquidity costs from walking through the limit order book, $\mu^+ > 0$ is the mean of $dp^+$, $Q$ (resp $\nu$) is the probability (resp mean beyond $k_0$) of being executed beyond $k_0$. From this expression, there is a competitive aspect between the risk of being caught on the wrong side and the fact of pushing $\imath(v)$ way beyond $1/2$ to get $H\mu^+(2\imath(v) - 1) \leq 0$.

In a theoretical part we first provide conditions for the limit order book to admit spoofing manipulation[5] as a function of the initial imbalance $\bar{\imath}$, local sensitivity of the imbalance $w$, overall price impact $\mu^+$, liquidity cost $G$, as well as the objective $H$. In short, this model confirms several intuitive facts when spoofing is more likely to occur

- if the probability $Q$ for the spoofing order to be executed is small;
- if the local sensitivity $w_{k_0}$ or the overall price impact $\mu^+$ is large;
- if the amount of share $H$ to buy is large with respect to the depth of the limit order book;
- if the initial market imbalance is close to $1/2$, that is, the market is equally balanced between offer and demand.
- away from the top of the limit order book;

As for the depth of the limit order book – how liquid the market is – the results are inconclusive. For this to be taken into account, one should model how the above mentioned parameters depend on the liquidity of the limit order book. The subsequent empirical study shows that it is the case, but we can not derive conclusions from this model as in Shorter and Miller [22] where illiquid markets seems more prone to spoofing. We then address the impact of spoofing on the resulting imbalance $\imath_{spoof} = \imath(v_{spoof})$ and discuss its dependence as a function of the aforementioned parameters. We characterize and discuss the deviation for the imbalance as a function of the different parameters. In particular as a function of the depth where the spoofing order is posted. We finally address the situation of a market maker using spoofing strategies for a positive round trip payoffs.

Based on this study, we can theoretically discriminate a spoofed imbalance $\imath_{spoof}$ from the legitimate one $\bar{\imath}$. Yet, from an outsider perspective, this is a hidden value since only the spoofer is aware of the actual imbalance. The main idea for the detection is to observe that a successful spoofing strategy requires the execution of a market order. We therefore compare the imbalance $\imath_-(t)$ before a market order at time $t$ with the imbalance $\imath_+(t)$ after this market order. If the market order is legitimate, the behavior of these two imbalances should follow some steady state distribution $(\bar{\imath}_-, \bar{\imath}_+)$. On the other hand, if the market order is the result of a spoofing strategy, the imbalance before the market order should be of the form

$$\imath_{spoof} \approx \frac{b}{b + a + wv_{spoof}} < \frac{b}{a+b} = \bar{\imath}_-$$

while returning to its steady state $\bar{\imath}_+$ as soon as the spoofed volumes are canceled. Hence, a quantification approach is to compare the distance from the instant imbalance $\imath_-(t)$ before a market order at time $t$ with, one the one hand, the long term legitimate one $\bar{\imath}_-$, and with, on the other hand, the theoretical spoofed one $\imath_{spoof}$. This measure is done according to the current market state situation $\imath_+(t)$. In other terms we measure and compare

$$\underbrace{d\left(\imath_-(t), \bar{\imath}_- | \imath_+(t)\right)}_{\substack{\text{distance of instant imbalance } \imath_-(t) \\ \text{before market order to legitimate imbalance } \bar{\imath}_- \\ \text{given current market conditions } \imath_+(t)}} \qquad \text{and} \qquad \underbrace{d\left(\imath_-(t), \imath_{spoof} | \imath_+(t)\right)}_{\substack{\text{distance of instant imbalance } \imath_-(t) \\ \text{before market order to spoofing imbalance } \imath_{spoof} \\ \text{given current market conditions } \imath_+(t)}}$$

---

[5] In other words, better than immediate or delayed market order.

For the distance, we adopt the non parametric Wasserstein distance. The technical design, in particular in terms of conditioning, the calibration with market data, implementation as well as the reason for such an approach are explained and illustrated for several stocks from TMX.

Before addressing the relevant literature, let us expose the shortcomings and modeling choices of this approach. The micro-structure dynamic of the limit order book at high frequency is complex. To excerpt some key impacts of spoofing behavior we deliberately focus on a static situation where the dynamic of the market is ignored.[6] For instance, we do not consider situations where compound spoofing behavior happens. We furthermore assume that there exists a single potential spoofer and that the market is infinitely reactive in the sense that it comes back to its steady state driven by the imbalance. There is no implicit game where the market acknowledges the existence of a potential spoofer, such as in Kyle [16] for instance. Finally, we assume that spoofing behavior is happening in a single market, while it has been documented and studied by van Kervel [24] that fast traders take advantage of multiple venues to post market orders in one while cancelling their limit orders in other venues. However, the present approach could also take into account an average imbalance over several venues. There is also no competitive game between two or more spoofers. Also, even though we shortly address the situation of a round trip for a market maker and the resulting optimal spoofing behavior, we take the viewpoint of a market taker willing to purchase/sell a given amount of shares. The overall goal being to understand the mechanism of spoofing in its most simple nature, quantify the resulting impact and derive potential detection procedures. Refinement of this approach, other take on, as well as more adequate quantification procedures are topics of further studies.

1.1. **Literature review.** There exists a solid stream of research showing that even rational speculative activities might destabilize prices, have an adverse effect on market efficiency or eventually lead to different forms of arbitrage; From market speculation based on various form of information asymmetries, for instance Hart and Kreps [14], Allen and Gale [3] or Jarrow [15], to price manipulation in limit order books using different market impact assumptions and trading strategies, as studied in Alfonsi and Schied [1], Alfonsi et al. [2], Gatheral [10], Gatheral and Schied [11]. The specific case of spoofing behavior has not yet been the subject of much theoretical study.

Although many high frequency trading strategies are legitimate, Shorter and Miller [22] point out that high frequency trading firms may engage in potentially manipulative strategies involving the usage of quote cancellations. Lee et al. [17] empirically study the change in spoofing behavior following a modification in volume disclosure rules on the Korean Exchange (KRX) at the start of 2002. Up to the end of 2001 the KRX disclosed the total volume of shares on both sides of the book and also the volumes at each tick up to 5 ticks from the best ask/bid. At the start of 2002 the KRX stopped disclosing the total volume on both sides in an effort to stop spoofing, but increased the disclosed volumes at each tick from the first 5 to the first 10 ticks from the best ask/bid. They show that spoofing is profitable and spoofers tend to prefer stocks with higher return volatility, lower market capitalization, lower price and lower managerial transparency. This study suggests the importance of the depth of book on spoofing strategies and potential price manipulation being carried out through a form of "volume imbalance". Wang [26] followed a similar methodology in empirically studying spoofing on Taiwan's index futures market. They found consistent results on the impact of spoofing on the market, but without the novel testing ground on changes in the disclosure of volumes deeper in the limit order book.

Some other studies to detect price manipulations are mainly based on learning algorithms. Among other, Cao et al. [5, 6] based on the definition of spoofing in [18] use K-nearest

---

[6]Since we consider the limit order book beyond its top, a dynamic version of the present approach would result into a fairly complex and high dimensional dynamic programing problem.

5

neighbour, one class support vector machine and adaptive hidden markov models to classify the data. Miranda et al. [20] characterize spoofing and pinging as full and partial observability of Markov decision processes. Under a reinforcement learning framework, they find that in order to maximise the investment growth, a trader will always employ spoofing or pinging orders except when market adds extra transaction costs or fines. In contrast to these empirical studies, our approach focuses on the micro economic features of spoofing behavior, in particular using our main stylized signal, the imbalance, which measures the difference between bid and ask side.

Concerning the impact of the imbalance on direction of the price movement: Lipton et al. [18] use the definition on the top of the book for the imbalance and study the impact on the trade arrival dynamic and resulting price movement. They fit a stochastic model for this behavior on real market data. Cartea et al. [7] employ volume imbalance as a signal to improve profits on the liquidation of a collection of shares in a dynamic high-frequency trading environment. Gould and Bonart [12] fit logistic regressions between the imbalance and the direction of the subsequent mid-price movement for each of 10 liquid stocks on Nasdaq, and illustrates the existence of a statistically significant relationship. Xu et al. [28] compute the imbalance at multi position in the limit order book and fit a linear relationship between this imbalance and the mid-price change. They find that the goodness-of-fit is considerably stronger for large-tick stocks than it is for small-tick stocks. The impact of order imbalance on prices has also been studied by Cont et al. [9] and Bechler and Ludkovski [4], for example. Bechler and Ludkovski also found that including characteristics of deeper parts of the book may be necessary for forecasting price impact. However, due to the nature of their dataset, they were only able to look at an aggregated form of the depth of book while we are able to use the exact volumes at all depths in the book. Sirignano [23] also used the book volumes beyond the touch to model price movements in a deep learning setting. Further suggesting the impact of depth of book on predicting future price movements.

Finally, the closest work to the present one in terms of quantitative analysis of spoofing behavior in relationship with imbalance is from Cartea et al. [8]. They adopt a dynamic approach where the trader influences the imbalance to derive the optimal strategy. They calibrate their model to market data and provide trading trajectories for the spoofer showing that spoofing considerably increases the revenues from liquidating a position. While being in a dynamic setting, in contrast to the present study, everything happens at the top of the limit order book for the imbalance to be manipulated. Furthermore, we do not focus here on the resulting gains from the spoofer, be rather on the impact on the imbalance from spoofer as for detection purposes from a regulatory viewpoint.

1.2. **Organization of the paper.** The first Section introduces the model, the imbalance and the dependence of the price movement on that imbalance. The second Section presents the spoofing strategies, addresses the theoretical conditions for spoofing behavior to happen and provide the resulting imbalance after spoofing together with numerical illustrations. The third Section is dedicated to the calibration procedure of the model on real Level 2 market data from TMX. The last Section discusses and introduces a quantitative approach to track spoofing behavior in real time illustrated on real datasets. Proofs, treatment of the round trip situation, calibration details and conditional distance specification are content of the Appendix.

## 2. Limit Order Book, Liquidity Costs and Imbalance

The ask price is denoted by $p$ and the limit order book on the ask side by $\bar{v} = (\bar{v}_0, \bar{v}_1, \ldots)$, that is, $\bar{v}_0$ is the volume posted at ask price $p$, $\bar{v}_1$ the volume posted at $p + \delta$, etc. where $\delta$ is the tick size. We denote by $\bar{v}^- = (\bar{v}_0^-, \bar{v}_1^-, \ldots)$ the limit order book on the bid side.[7]

---

[7]That is $\bar{v}_0^-$ is the volume posted at bid price $p^- < p$, $\bar{v}_1^-$ the volume posted at $p^- - \delta$, etc.

Given a limit order book inventory $\bar{v}$ on the ask side, we define for an amount $H \geq 0$ of shares the function

$$F(H) := \inf \left\{ x \in \mathbb{N}_0 \colon \sum_{k=0}^{x} \bar{v}_k \geq H \right\}$$

which represents how many positive price tick deviation an order of size $H$ generates. Given an amount of shares $H$ to buy, a bid price $p$ and an ask limit order $\bar{v}$, the resulting costs of the market order is

$$pH + \sum_{k=0}^{F(H)} k\delta\bar{v}_k - \delta F(H) \left( \sum_{k=0}^{F(H)} \bar{v}_k - H \right) = pH + \delta G(H)$$

The term $G$ on the right hand side represents the liquidity costs depending only on $\bar{v}$.

**Remark 2.1.** *Throughout the theoretical part of this work we assume that the limit order book is blocked shaped with an amount $a > 0$ of shares at each price level of the ask side. We then get the continuous approximation*

$$F(H) \approx \frac{H}{a} \quad and \quad G(H) \approx \frac{H^2}{2a}$$

As for the imbalance of the limit order book, measure of the difference between offer and demand, we proceed as follows. Let $w_0, w_1, \ldots$ with $\sum w_k = 1$ and $w_k \geq 0$, a weight for each tick level $k$, and a limit order book inventory $\bar{v}^-, \bar{v}$ on the bid and ask respectively, we denote by

$$\bar{B} = \sum w_k \bar{v}_k^- = \langle w, \bar{v}^- \rangle \qquad\qquad \bar{A} = \sum w_k \bar{v}_k = \langle w, \bar{v} \rangle$$

the weighted average bid and ask volumes. We define the imbalance as

$$\bar{\imath} := \frac{B}{B + A} \in (0, 1)$$

**Remark 2.2.** *In a blocked shaped setting, if $b$ denotes the amount of orders on every price level on the bid side, we get*

$$\bar{\imath} = \frac{\sum w_k b}{\sum w_k (a + b)} = \frac{b}{a + b}$$

*which yields $b = a\bar{\imath}/(1 - \bar{\imath})$.*

The price deviation in the next period can be triggered by two events. The posting and cancellation of incoming limit orders as well as the posting of market orders. We distinguish between both, since the former does not have an impact on the execution of existing limit orders while the latter has. We generically denote by

$$dp = \{\ldots, dp_{-1}, dp_0, dp_1, \ldots\} \qquad\qquad \mu = \sum x dp_x$$
$$dq = \{\ldots, dq_{-1}, dq_0, dq_1, \ldots\} \qquad\qquad \nu = \sum y dp_y$$

the distribution and mean, respectively, of the two possible price movement in the next period. For the sake of simplicity, we assume that the imbalance does not have an impact on incoming market orders and that the price deviation with respect to the market orders is neutral, that is $\nu = 0$. We furthermore assume that they are independent of each others.[8] To reflect the fact that the imbalance, as an indicator of the offer and demand on the market, has an impact on the market makers, we consider a parametrization $\imath \mapsto dp(\imath)$ of the price movement driven by limit orders as a function of the imbalance $\imath$. Since the imbalance

---

[8]The subsequent theoretical study adapt to eventual joint distribution of price movement due to limit and market orders also jointly dependent on the imbalance. The exposition of which is no longer explicit but can be solved numerically.

moves between 0 and 1, we assume that the distribution $dp$ moves as a convex combination of $\imath$ between the distribution $dp^-$ – distribution when the imbalance is close to 0, that is highly skewed to the left – and the distribution $dp^+$ – distribution when the imbalance is close to 1, that is highly skewed to the right. Mathematically:

$$dp(\imath) = \imath dp^+ + (1 - \imath)dp^-$$

From the skewness assumptions and symmetry of the imbalance indicator, we assume that

$$dp_x^+ \geq dp_x^- \quad \text{for every } x \geq 0 \quad \text{and} \quad dp_x^+ = dp_{-x}^- \quad \text{for every } x$$

which implies that

$$\mu(\imath) := \sum x dp_x(\imath) = (2\imath - 1) \sum x dp_x^+ = \mu^+(2\imath - 1)$$

By assumption, $\mu^+$ is positive, showing that $\mu(i)$ moves between $-\mu^+$ and $\mu^+$ and is equal to 0 for an imbalance of $1/2$ when offers equal demand.

## 3. Spoofing strategy

Suppose that at a given time we are given an ask price $p$ and a limit order book inventory $(\bar{v}^-, \bar{v})$. A trader willing to buy an amount $H$ of shares faces the following three options.

- **Immediate market order:** for a total costs of

$$pH + \delta G(H)$$

- **Delayed market order:** for an average total cost of

$$\sum_{x,y} [(p + \delta(x + y)) H + G(H)] dp_x(\bar{\imath}) dq_y = pH + \delta \left( G(H) + H\mu^+ (2\bar{\imath} - 1) \right)$$

Clearly if $\bar{\imath} < 1/2$, then this second option is better than a direct buy.
- **Spoofing and delayed market order:** Book first an ask limit order $v$ at a depth $k$ in $\{0, 1, \ldots, N\}$ on top of the ask limit order book $\bar{v}_k$ to increase the liquidity on the ask side and signal a surge in supply to the market. In the next period the price deviates from $p$ to $p + \delta(x + y)$ and two situations may happen:
  - $y \leq k$: no market order of sufficient magnitude hits the limit order book and therefore this limit order is not executed against an incoming market order. The previous limit order is canceled and the amount $H$ of shares is acquired for a cost of

$$(p + \delta(x + y)) H + \delta G(H)$$

  - $y > k$: the limit order is executed against an incoming market order at a price level $p + \delta(x + k)$. The new objective moves to $H + v$ resulting in a net cost of

$$(p + \delta(x + y)) (H + v) + \delta G (H + v) - (p + \delta(x + k)) v$$
$$= (p + x\delta) H + \delta G(H + v) + \delta(y - k)v$$

It follows that the spoofing net cost for a price deviation of $p + \delta(x + y)$ is given by

(3.1) $\qquad C_k(v, x, y) = (p + (x + y)\delta)H + \delta G \left( H + v1_{\{y>k\}} \right) + \delta(y - k)v1_{\{y>k\}}$

However, the posting of the selling limit order modifies the imbalance from $\bar{\imath}$ to

$$\imath_k(v) := \frac{B}{A + B + w_k v}$$

In other words, the imbalance will move downwards, shifting the distribution $dp$ to more favorable outcomes. Since we assume that $\nu = 0$, it follows that the average net costs are

given by

$$C_k(v) := \sum_{x,y} C_k(v,x,y) dp_x(\imath_k(v)) dq_y$$

$$= \sum_x (p + \delta x) H dp_x(\imath_k(v)) + \delta(1 - Q_k) G(H) + \delta Q_k G(H + v) + \delta v \sum_{y \geq k+1} (y - k) dq_y$$

$$= pH + \underbrace{\delta(1 - Q_k) G(H)}_{\text{Cost for optimal situation}} + \underbrace{\delta H \mu^+ (2\imath_k(v) - 1)}_{\text{Spoofing impact}} + \underbrace{\delta Q_k G(H + v) + \delta v \nu_k}_{\text{Cost for being caught wrong way}}$$

where

$$Q_k = \sum_{y \geq k+1} dq_y \quad \text{and} \quad \nu_k = \sum_{y \geq k+1} (y - k) dq_y$$

Note that this cost functional is convex in $v$ since $G$ and $\imath_k$ are convex functions. Note also that in order to take advantage of this spoofing impact, it is necessary to drive the imbalance $\imath_k(v)$ below $1/2$.

**Remark 3.1.** *Note that we implicitly assume that the spoofing only happens at a given depth $k$. It is possible to spoof simultaneously at different depths resulting in a slightly more complex cost function that can be solved numerically. The conclusions do not change qualitatively and we use the more general multi-depth spoofing for the analysis of data in the subsequent sections.*

3.1. **Existence of Spoofing Manipulation.** The question is whether it is possible to push the imbalance as much as possible to $0$ in order to offset the costs of posting selling orders, they being executed and paying the liquidity costs of buying them back.

**Definition 3.2.** *We say that the limit order book $(\bar{v}^-, \bar{v})$ admits a (market taker) spoofing manipulation if there exists $v > 0$ and $k \in \{0, 1, \ldots, N\}$ such that*

$$(3.2) \qquad \begin{cases} C_k(v) < pH + \delta \left( G(H) + \mu^+ H (2\bar{\imath} - 1) \right) & \text{if } \bar{\imath} \leq 1/2 \\[2mm] C_k(v) < pH + \delta G(H) & \text{if } \bar{\imath} > 1/2 \end{cases}$$

According to the average costs of spoofing, these two inequalities turn into

$$(3.3) \qquad Q_k \left[ G(H + v) - G(H) \right] < 2\mu^+ H \left( \bar{\imath} \wedge \frac{1}{2} - \imath_k(v) \right) - v\nu_k$$

The following results concerns the existence of spoofing manipulation in a blocked shaped setting where the volume on the ask side of the limit order book is $a$ everywhere. For scaling reasons let $H = \rho a$, where $\rho$ represents the ratio of the shares to purchase to the depth of the limit order book. Our first result concerns the existence of spoofing manipulation.

**Proposition 3.3.** *In a block shaped setting, where the volume on the ask side of the limit order book is $a$ everywhere. The following assertions hold*

• *If $\bar{\imath} \leq 1/2$, the limit order book admits no spoofing manipulation if and only if (3.4) holds;*
• *If $\bar{\imath} > 1/2$, the limit order book admits no spoofing manipulation if (3.4) holds.*
*Where $H = \rho a$ and*

$$(3.4) \qquad 2\rho\mu^+ (1 - \bar{\imath}) \bar{\imath} w_k \leq Q_k \rho + \nu_k \quad \text{for all } k$$

For the proof, see Appendix A. From this proposition, we deduce that price manipulation is more likely to occur

 • if $Q_k$ is small – and as a byproduct $\nu_k$. If the probability to get a spoofing order executed is small, there is relatively no downsize at spoofing.

- if $\bar{\imath}$ is close to $1/2$. If the imbalance is close to $1/2$, then $\bar{\imath}(1-\bar{\imath})$ is maximum. The impact of moving the price in ones favor is maximal there.
- if $\mu^+$ is large: $\mu^+$ represents the mean deviation sensitivity as a function of the imbalance. The more sensitive the price movement is with respect to the imbalance, the more likely spoofing strategies may occur.
- If $w_k$ is large: $w_k$ represents the relative impact at tick level $k$ of a spoofing volume to the imbalance. If one of $w_k$ is large with respect to the corresponding $Q_k$, then spoofing is more likely to occur there.
- if $\rho$ is relatively large. If the amount of order to buy relative to the overall offer is very large, spoofing is more likely to happen.

Figure 1 represents the existence of spoofing condition 3.4 in terms of the initial imbalance $\bar{\imath}$ with varying market parameters.
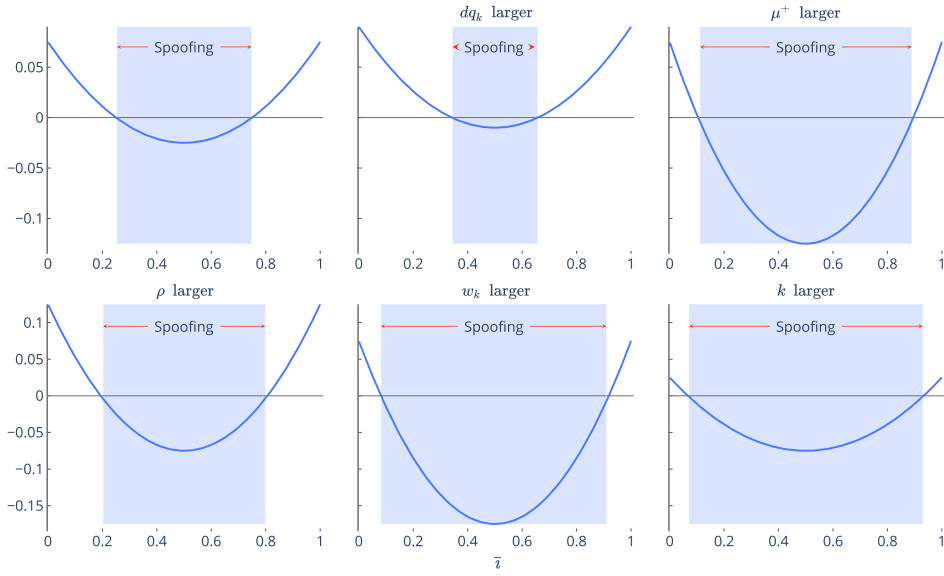


**Figure 1.** Spoofing condition (3.4) as a function of $\bar{\imath}$ and in $(a)$ $\mu^+ = 1, \rho = 1, k = 3, w_k = 0.2, dq_y = 0.025$ for all $y \geq k$. One parameter is increased each time with respect to $(a)$ where $(b)$: $dq_y = 0.03$ for all $y \geq k$; $(c)$: $\mu^+ = 2$; $(d)$: $\rho = 2$; $(e)$: $w_k = 0.5$; $(f)$: $k = 4$.

3.2. **Optimal Spoofing and Resulting Imbalance Impact.** Let us now address the problem of finding the optimal spoofing strategy. In particular as a function of the depth at which the spoofing order is placed.

**Proposition 3.4.** *The optimal spoofing volume $v_{spoof}$ at a given level $k$ and resulting imbalance $\imath_{spoof}$ – adopting the notations $w := w_k$, $Q := Q_k$ and $\nu := \nu_k$ – are given by:*

$$v_{spoof} = \frac{a}{Q}\left[2\rho w \mu^+ \frac{1-\bar{\imath}}{\bar{\imath}}\imath_{spoof}^2 - (Q\rho + \nu)\right]^+$$

*where $\imath_{spoof}$ is the unique cubic root solution in $(0, \bar{\imath}]$ of*

$$\frac{\bar{\imath}}{\imath} = 1 + \frac{(1-\bar{\imath})w}{Q}\left[2\rho w \mu^+ \frac{1-\bar{\imath}}{\bar{\imath}}\imath^2 - (Q\rho + \nu)\right]^+$$

For the proof, see Appendix A. Though the solution is implicit, we can inspect the spoofing behavior as a function of the distance to the top of the limit order book. Note first that $\imath_{spoof} = \bar{\imath}$ if and only if $2\rho w \mu^+(1-\bar{\imath})\bar{\imath} \leq Q\rho + \nu$ which results into $v_{spoof} = 0$. This

coincide with the no spoofing condition of the previous proposition. We are interested in the relative spoofing size as a function of these parameters. From the definition of the imbalance, $\imath(v)$ increases if and only if $v$ decreases, so we get more spoofing volume as $\imath_{spoof}$ gets smaller. Now from the implicit function it holds that

$$\frac{1}{\imath} = \frac{1}{\bar{\imath}} + \frac{w}{Q}\frac{1-\bar{\imath}}{\bar{\imath}}\left[2\rho w\mu^{+}\frac{1-\bar{\imath}}{\bar{\imath}}\imath^{2} - (Q\rho+\nu)\right]^{+} =: f\left(w,\mu^{+},\nu,Q,\bar{\imath},\rho,\imath\right)$$

where the function $f$ is an increasing function of $\imath$ greater than $1/\bar{\imath}$.

- Since $f$ is increasing as a function of $w$ and $\mu^{+}$, it follows that $\imath_{spoof}$ is decreasing as a function of $w$ and $\mu^{+}$. Hence, spoofing behavior increases as a function of the impact $w$ at level $k$ on the imbalance as well as a function of the overall sensitivity $\mu^{+}$ of the price movement with respect to the imbalance.
- Since $f$ is decreasing as a function of $Q$ and $\nu$, it follows that $\imath_{spoof}$ is increasing as a function of $Q$ and $\nu$. From an empirical viewpoint, $Q = Q_{k}$ as well as $\nu = \nu_{k}$ decreases as a function of the depth $k$. It follows that spoofing behavior is more likely to happen and increase deeper in the limit order book. However, this conclusion is short of the fact that local sensitivity of the imbalance $w = w_{k}$ also depends on the depth with an inverse impact. According to empirical analysis, it turns out that $w$ does not exert this decreasing behavior as a function of $k$ at least within a reasonable depth in the limit order book. It seems that spoofing behavior is more likely to happen at a reasonable distance from the top of the limit order book.

The behavior of the resulting imbalance $\imath_{spoof}$ as a function of the initial imbalance $\bar{\imath}$ is more difficult to stress out. We know that $\imath_{spoof} \leq \bar{\imath}$ and for the same reasons as before it is increasing as a function of $\bar{\imath}$. The same holds for the dependence on the relative number of shares to purchase $\rho$. Figure 2 represents the curves of the spoofed imbalance $\imath_{spoof}(\bar{\imath})$ as a function of the initial imbalance $\bar{\imath}$ for different depths with varying market parameters.
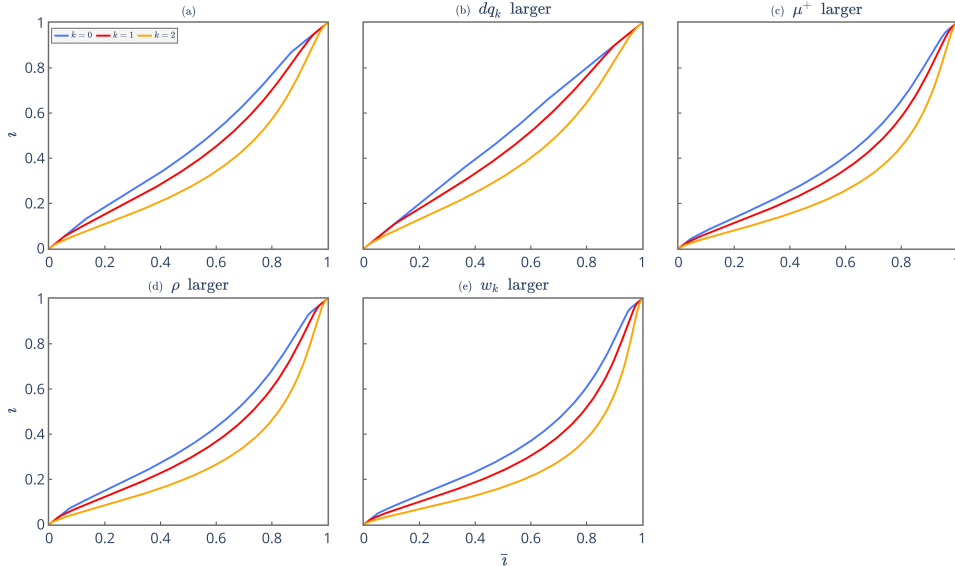


**Figure 2.** $\imath$ as a function of $\bar{\imath}$ and in $(a)$ $\mu^{+} = 1, \rho = 1, w_{k} = 0.2, dq_{y} = 0.003$ for all $y \geq k$. One parameter is increased each time with respect to $(a)$ where $(b)$: $dq_{y} = 0.006$ for all $y \geq k$; $(c)$: $\mu^{+} = 3$; $(d)$: $\rho = 3$; $(e)$: $w_{k} = 0.5$; Blue line :$k = 0$; red line: $k = 2$; orange line: $k = 4$.

**Remark 3.5.** *Throughout, we mainly focus on the spoofing behavior from a market taker's viewpoint. As for a market maker, spoofing behavior might be rewarding as well. It turns out that the rewards from spoofing are intertwined with the ones from pure market making. The resulting impact on the imbalance is however quite similar, up to the fact that the bid ask spread is an additional factor in the spoofing opportunity, since the market maker will have to cross the spread. In Appendix B, we derive the spoofing strategy in the round trip context, discuss the spoofing impact on the imbalance and numerical analysis in the same context as the present situation.*

## 4. CALIBRATION

According to this model, we calibrate the imbalance as well as $dp$ and $dq$ on real data provided by TMX. These datasets consists of level 2 data from June to September 2017. Among the 1500 available equities we selected 10, varying in company background, market capitalization as well as trading frequency. The level 2 datasets include time, order price, volume, type,[9] order ID and counterpart order ID in case of a trade, see Figure 3

| time | side | price | order_id | reason | change |
|------|------|-------|----------|--------|--------|
| 09:30:00.003886 | Sell | 62 | S20170417000000020 | Partial Fill | -100 |
| 09:30:00.042634 | Sell | 62.33 | S20170417000000031 | New Order | 200 |
| 09:30:00.081079 | Buy | 61.66 | B20170417000000017 | Cancelled | -100 |

**Figure 3.** Original Level 2 dataset of stock AEM provided by TMX.

Since these are provided in diff form, therefore a cumulative aggregation allows to construct the full limit order book at any time as in Figure 4

| time | side | price | change | order_id | reason | $p_b$ | $p_a$ | vb5 | vb4 | vb3 | vb2 | vb1 | va1 | va2 | va3 | va4 | va5 |
|------|------|-------|--------|----------|--------|-------|-------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 09:31:01.279206 | Sell | 62.1 | -100 | S20170417000000117 | Cancelled | 61.95 | 62.1 | 100 | 200 | 100 | 200 | 100 | 600 | 500 | 100 | 100 | 200 |
| 09:31:01.475385 | Sell | 62.08 | 100 | S20170417000000024 | New Order | 61.95 | 62.1 | 100 | 200 | 100 | 200 | 100 | 500 | 500 | 100 | 100 | 200 |
| 09:31:01.475612 | Sell | 62.08 | 100 | S20170417000000021 | New Order | 61.95 | 62.08 | 100 | 200 | 100 | 200 | 100 | 100 | 0 | 500 | 500 | 100 |

**Figure 4.** Generation of the full limit order book out of the Level 2 data.

This operation is computationally very intensive and therefore has been realized on a distributed data cluster of TMX with spark.
With the full limit order books at hand we divide the calibration into three steps:

- Find a normalized sample frequency and choose a maximal depth for the support of the distributions $dp^{\pm}$ and $dq$;
- Calibrate the imbalance generically, that is, as a function of the weights $w = (w_1, w_2, \ldots)$;
- Estimate $dq$, $dp^{\pm}$ and weights $w = (w_1, w_2, \ldots)$.

The sample frequency should be large enough such that there exists enough variance in price change, see Figure 5

---
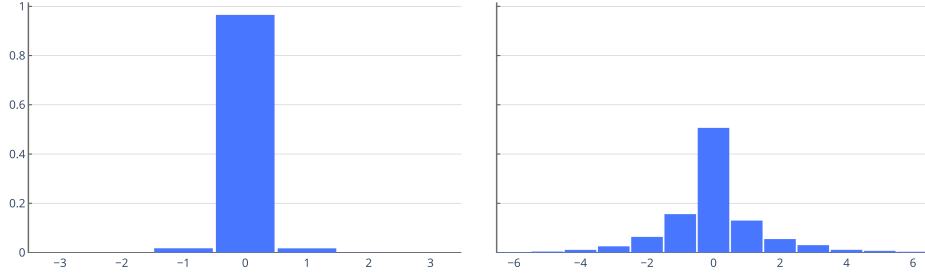
[9]Buy/sell; booked/cancelled/traded

**Figure 5.** Left panel: Histogram of original AEM price change. Right panel: Histogram of AEM price change after sampling.

To compare across markets with different trading activity—and eventually time during the day—we fix a target variance of $\sigma^2$ for the price movement and select the optimal frequency $f$ for each stock as to minimize the square distance between $\sigma_f$ and $\sigma$. For a target variance $\sigma^2 = 2$, Table 1 is the sample frequency for different stocks with a summary statistics of the average volume and arrival rate for Market/Limit Orders on the bid and ask side. As for the maximal depth for the support of the distribution, we take the 99% quantile of empirical sampled price change distribution. The depths are around 4 since we use the same $\sigma^2 = 2$ to determine the sampling frequency. The sampling frequency is related to how fast limit orders arrive, not market orders. For most stocks, $f$ is small when the arrival rate of limit orders is high.

| Stock | $f$ | Depth | Market Orders | | | | Limit Orders | | | |
| | | | Buy | | Sell | | Buy | | Sell | |
| | | | Vol | Rate | Vol | Rate | Vol | Rate | Vol | Rate |
|-------|-----|-------|------|-------|------|-------|------|-------|------|-------|
| AEM | 4 | 4 | 146 | 0.072 | 144 | 0.065 | 144 | 5.145 | 142 | 5.232 |
| BB | 38 | 4 | 567 | 0.107 | 613 | 0.085 | 3257 | 3.034 | 3245 | 3.008 |
| BMO | 11 | 4 | 171 | 0.107 | 168 | 0.134 | 149 | 2.482 | 151 | 2.553 |
| CNR | 6 | 4 | 141 | 0.104 | 134 | 0.1 | 138 | 2.397 | 141 | 2.399 |
| CPG | 53 | 5 | 356 | 0.13 | 349 | 0.115 | 866 | 1.762 | 879 | 1.741 |
| FNV | 3 | 5 | 123 | 0.059 | 121 | 0.056 | 126 | 2.881 | 140 | 2.794 |
| FR | 60 | 3 | 243 | 0.072 | 272 | 0.063 | 795 | 2.115 | 842 | 2.049 |
| PPL | 26 | 4 | 142 | 0.101 | 152 | 0.11 | 165 | 2.31 | 183 | 2.469 |
| TD | 20 | 4 | 223 | 0.205 | 213 | 0.218 | 277 | 3.87 | 278 | 3.866 |
| VET | 6 | 5 | 127 | 0.07 | 152 | 0.069 | 130 | 2.474 | 137 | 2.439 |

**Table 1.** Stock data from June 5, 2017 to June 9, 2017. The Vol columns corresponds to the average volume of a single order during inspected time interval and the Rate columns corresponds to the number of orders per second.

With the sampling frequency $f$ and depth $N$, we define the average imbalance at time $t$ as

$$\hat{i}(w, t) = \frac{\sum\limits_{k \leq N} \sum\limits_{t-f \leq s < t} \bar{v}_k^-(s) w_k}{\sum\limits_{k \leq N} \sum\limits_{t-f \leq s < t} \left(\bar{v}_k^-(s) + \bar{v}_k(s)\right) w_k}$$

which sums up order book volumes within a certain time interval weighted by time difference $\Delta s_i = s_{i+1} - s_i$. The weighting parameter $w$ impacts the average imbalance distribution, see Figure 6.
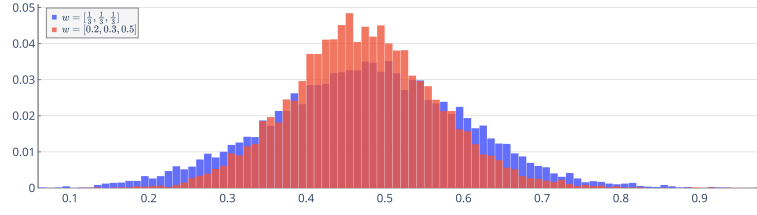
13

**Figure 6.** Distribution of the average imbalance for different weights for the stock BMO.

Nevertheless for each weight vector, the resulting distribution is close to a skewed normal distribution. For a given weight $w$, using maximum likelihood, we fit the empirical distribution to the corresponding skew normal distribution $\mathcal{SN}(\alpha(w), \xi(w), \omega(w))$, the fit of which is particularly good, see figure 7 for an example.
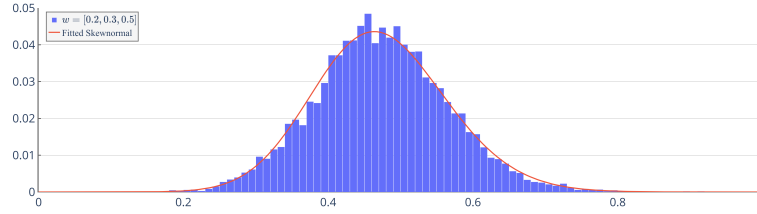


**Figure 7.** Histogram of BMO average imbalance and fitted skewnormal distribution: $w = [0.2, 0.2, 0.2, 0.2, 0.2]$.

The third step is to determine $dq$, $dp^{\pm}$ with the optimal weights $w$. As for $dq$, it is the probability that the price moves by $k$ ticks triggered by market orders. Thus for each market order, compute

$$F(H) = \inf\{x \in \mathbb{N}_0 : \sum_{k=0}^{x} v_k \geq H\}$$

where $H$ is the volume of the market order. This represents exactly how many positive tick price deviation an order of size $H$ will produce. We derive $dq$ from the empirical distribution.

As for $dp^{\pm}$ and $w$, a maximum likelihood estimation is implemented to solve

$$(4.1) \qquad dp^*, w^* = \underset{dp^+, w}{\arg\min} \left[ -\frac{1}{M} \sum_{m=1}^{M} \log p\left(x_m, \hat{\imath}_m\right) \right]$$

where $x_m$ is the empirical price change, $\hat{\imath}_m$ the average imbalance for a given weight $w$,

$$p\left(x_m, \hat{\imath}_m\right) = dp_{x_m}\left(\hat{\imath}_m\right) p(\hat{\imath}_m)$$

where $dp_{x_m}\left(\hat{\imath}_m\right) = \hat{\imath}_m dp_{x_m}^+ + (1 - \hat{\imath}_m) dp_{x_m}^-$ represents the conditional probability of price change equal to $x_m$ given $\hat{\imath}_m$, and $p(\hat{\imath}_m)$ is the density of the fitted skewnormal distribution evaluated at $\hat{\imath}_m$.

Figure 8, illustrating the value of the optimal weights $w$ for selected stocks, shows different patterns. Overall, it turns out that the relative impact of the imbalance to the price distribution is more important away from the top of the limit order book.

We also performed this calibration procedure on stock BMO weekly from June 5th to June 30th, as well as for the first hour of trading monthly from June to September. Figure 9 provide the optimal weights in each case for BMO.
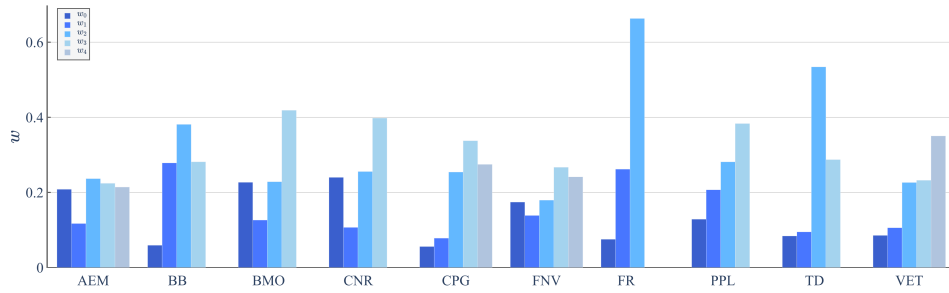
**Figure 8.** $w$ for stock AEM, BB, BMO, CNR, CPG, FNV, FR, PPL, TD, VET from June 5th, 2017 to June 9th, 2017.
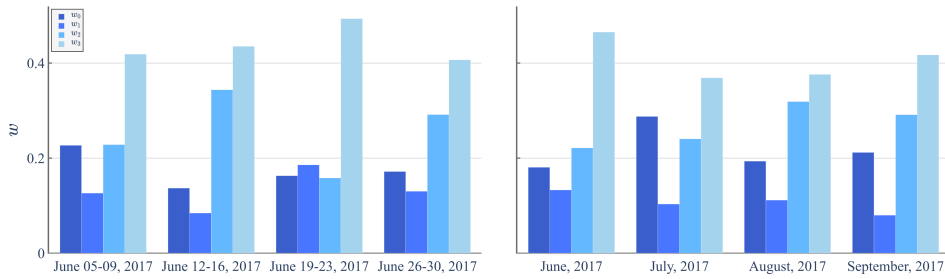


**Figure 9.** Left panel: $w$ for stock BMO each week in June 2017. Right panel: $w$ for stock BMO each month in June 2017, only using first hour trading data.

Notice that for the first hour of trading the optimal weights are more consistent across time, but all show that the weight impact on the price movement happens deeper in the limit order book.

As for the corresponding $dp^+$ and $dq$, they are represented in Figure 10 for stock BMO from June 5th, 2017 to June 9th, 2017. As expected, $dp^+$, representing the price movement as the imbalance is large, is skewed to the right.



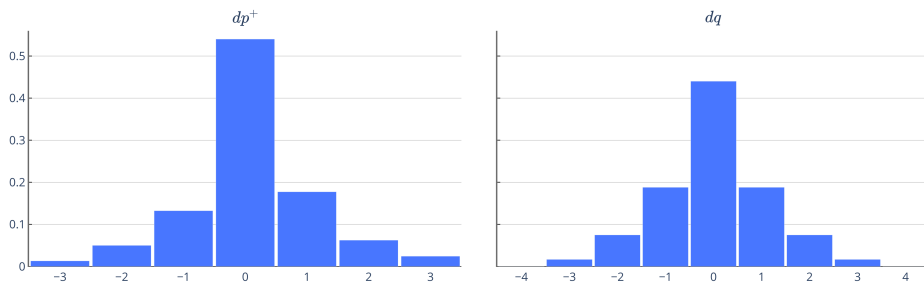**Figure 10.** $dp^+$ and $dq$ for BMO from June 5th to June 9th.

Table 2 provides the moments of $dp^+$ for each stock – in tick values. Skewness shows how much $dp^+$ is skewed to right. When it is large, $\mu^+$ is also large and spoofing has a larger impact according to theoretical part. Except for CPG which is relatively small, all of the other stocks under study excerpt this pattern of right-skewness.

15

| Stock | $\mu^+$ | Variance | Skewness | Kurtosis |
|-------|---------|----------|----------|----------|
| AEM | 0.411 | 0.881 | 1.175 | 4.46 |
| BB | 0.467 | 0.766 | 0.957 | 3.786 |
| BMO | 0.103 | 1.087 | 0.095 | 4.263 |
| CNR | 0.398 | 0.636 | 1.447 | 4.928 |
| CPG | 0.100 | 1.171 | -0.068 | 5.275 |
| FNV | 0.404 | 0.850 | 2.038 | 7.169 |
| FR | 0.209 | 1.033 | 0.532 | 2.203 |
| PPL | 0.076 | 1.177 | 0.645 | 3.83 |
| TD | 0.118 | 1.389 | 0.633 | 4.471 |
| VET | 0.315 | 1.253 | 1.899 | 7.393 |

**Table 2.** Moments of $dp^+$.

## 5. APPROACHES TO SPOOFING DETECTION

For reasons mentioned in the introduction, it is difficult from a regulatory viewpoint to figure out whether or not spoofing happened a-posteriori. According to the theoretical part, the act of spoofing will influence the resulting imbalance. However, to monitor the imbalance is akin to contemplate pure noise as shown in Figure 11.
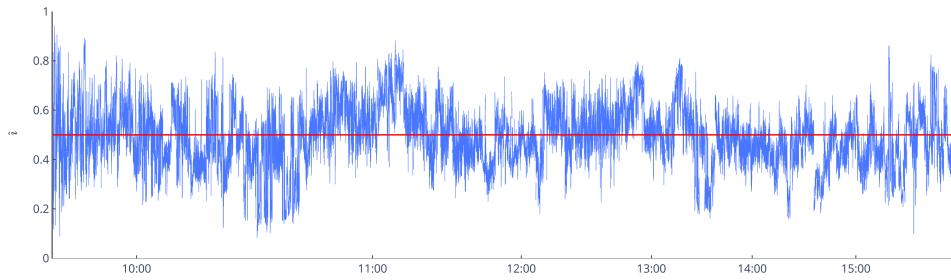


**Figure 11.** Imbalance of stock BMO from 09:30 to 16:00 on June, 7, 2017.

In the following, we propose some possible ways to perform such a monitoring based on the theoretical results. The strategy comes from the following observation: For a spoofing strategy to be successfully fulfilled, a market order has to be executed.[10] Hence when observing an executed market order two situations may happen:

*1-* The market order is a legitimate one. In that case, the imbalance before this market order $\imath_-$ and after $\imath_+$ should follow statistically the classical long run behavior. In other words, in a legitimate situation, we should observe statistically the pair

$$(\imath_-, \imath_+)$$

for each market order.

*2-* The market order is the result of a spoofing behavior. The implicit equilibrium without spoofing would be the pair $(\imath_-, \imath_+)$. After the market order is executed, the market imbalance should be back to its equilibrium $\imath_+$. However before the market order, the spoofer observes the implicit imbalance $\imath_-$ and decides to spoof according to this information, sending to the market $\imath_{spoof}(\imath_-)$ instead of $\imath_-$. The resulting observation for those spoofed market orders is therefore the pair

$$(\imath_{spoof}(\imath_-), \imath_+)$$

---

[10]In this paper, we do not consider spoofing strategies involving only limit orders.

Furthermore, spoofing strategies are supposed to happen sporadically but intensively within a short time horizon. Before presenting some strategies, let us fix some notations:

- $\Pi = \{t_1 < t_2 < \ldots < t_M\}$ represents the time stamps of each (buy) market orders in a long sample (several weeks).
- $\hat{\imath}_-(t)$ and $\hat{\imath}_+(t)$ represents the imbalance before and after the market order happening at time $t$ in $\Pi$.
- $(\imath_-, \imath_+)$ represents the overall joint distribution of the imbalance right before and after each market orders fitted to the overall data. We assume that these represents the stable behavior of the market without spoofing, and therefore representative of legitimate market orders.
- $(\hat{\imath}_-^N(t), \hat{\imath}_+^N(t))$ represents the (short span) empirical distribution at time $t$ in $\Pi$ generated by the last $N$ market orders observed imbalances $(\hat{\imath}_-(s), \hat{\imath}_+(s))$, where $N \ll M$ is a short horizon sample size (in our case about 100).
- The previous theoretical part, even if not explicit in terms of solution allows us to compute numerically $\imath_{spoof}(\imath_-)$ for a given implicit imbalance $\imath_-$.

5.1. **Monitoring $\hat{\imath}_-^N$.** A first idea is to monitor the behavior of the short term imbalance $\hat{\imath}_-^N(t)$ as times passes to test whether it is statistically different from the equilibrium $\imath_-$. This is however not adequate for the following reasons. First, this is not related to spoofing behavior and might reflects some other market patterns. Second, and more importantly, the sequence of $\hat{\imath}_-(t)$ for each market order is highly dependent. Indeed, there might exists market conditions – bullish/bearish, etc – such that a short horizon sample $\hat{\imath}_-^N$ differs strongly from the long term behavior. Figure 12 provides empirical evidence about the sequential dependence of the imbalance $\hat{\imath}_-$ as well as $\hat{\imath}_+$ over time.
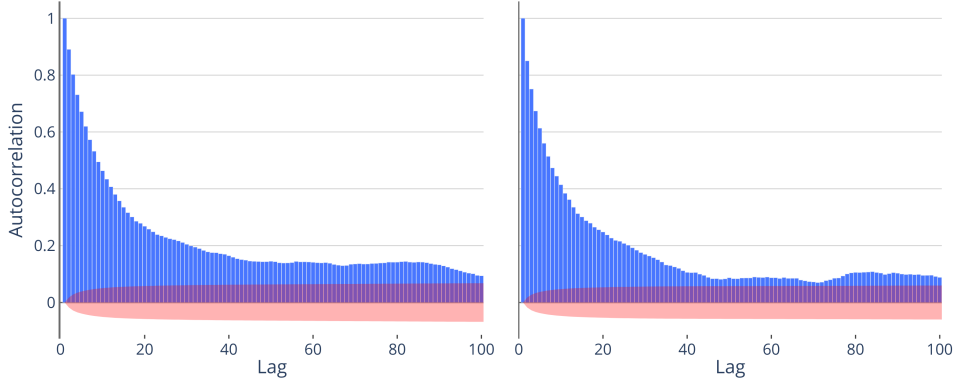


**Figure 12.** Left panel: $\hat{\imath}_-$ autocorrelation of stock BMO on June 7, 2017. Right panel: $\hat{\imath}_+$ autocorrelation of stock BMO on June 7, 2017. Red area is the 95% confidence interval of the autocorrelation.

5.2. **Monitoring $(\hat{\imath}_-^N, \hat{\imath}_+^N)$.** The statistical link towards discrimination of $\imath_{spoof}(\imath_-)$ from $\imath_-$ is the additional observation of the imbalance after the spoofing happen. This provides statistical a-posteriori information about the implicit market equilibrium before spoofing which in case of spoofing can not be directly observed. Figure 13 shows on the left panel the joint distribution $(\imath_-, \imath_+)$ while the right panel represents, based on the model of the theoretical part and calibration, the joint distribution $(\imath_{spoof}(\imath_-), \imath_+)$ in the case of spoofing. The spoofed joint distribution is skewed to the left in comparison to the non-spoofed one, in accordance to the theoretical analysis that spoofing decreases the imbalance – in the buy order case – before a market order.
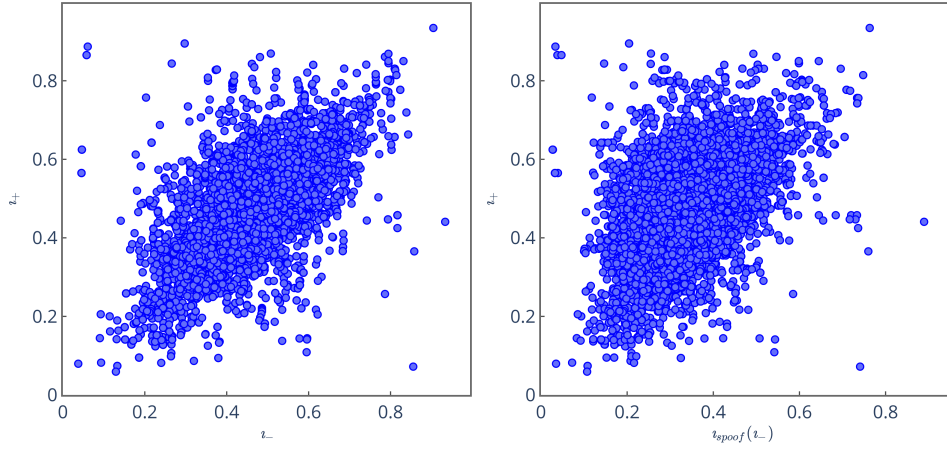
**Figure 13.** Left panel: Empirical joint distribution of $(\imath_-, \imath_+)$. Right panel: Joint distribution of $(\imath_{spoof}(\imath_-), \imath_+)$

A possible way to detect spoofing is therefore to compare the long run distribution $(\imath_-, \imath_+)$ with the short term empirical distribution $(\hat{\imath}_-^N, \hat{\imath}_+^N)$. These two distributions encode the possibility to disentangle legitimate market behavior from spoofed ones. However, as in the previous approach, the sequence of joint observation is once again not iid. For short time horizon, the market may be legitimate, though far away from the long run distribution.

5.3. **Monitoring $\hat{\imath}_-^N$ conditioned on $\hat{\imath}_+^N$.** To overcome the previous shortcomings, the next approach is to monitor $\hat{\imath}_-^N(t)$ conditioned on the current market state $\hat{\imath}_+^N(t)$. From our hypothesis, $\imath_+$ represents the steady state of the market at equilibrium after a market order. It turns out that conditioned on $\hat{\imath}_+(t)$ the sequence of $\hat{\imath}_-(t)$ is closer to iid.
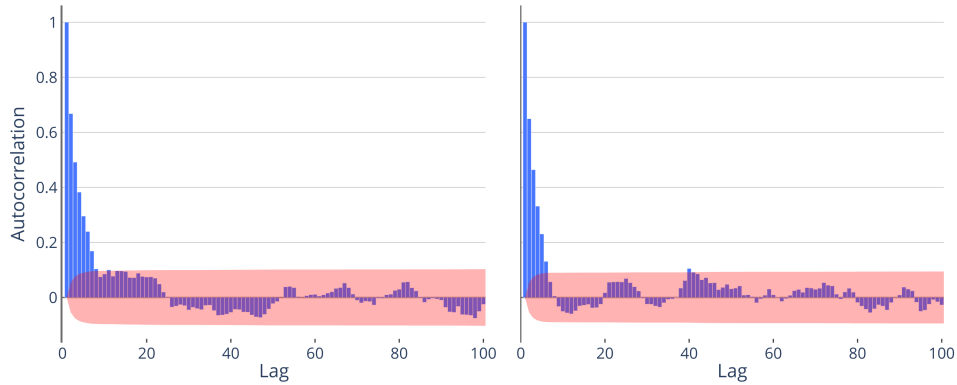


**Figure 14.** Left panel: Autocorrelation for $0.329 \leq \hat{\imath}_-^N | \hat{\imath}_+^N \leq 0.378$ of stock BMO on June 7, 2017. Right panel: Autocorrelation for $0.561 \leq \hat{\imath}_-^N | \hat{\imath}_+^N \leq 0.61$ of stock BMO on June 7, 2017. Red area is the 95% confidence interval of the autocorrelation.

In order to detect spoofing behavior, instead of adopting a statistical test for which some parametric assumptions on the distribution has to be made, we measure the distance between $\imath_-$ and $\hat{\imath}_-^N$ conditioned on the current observed imbalance $\hat{\imath}_+^N$ using a non-parametric distance, the Wasserstein distance, see Appendix for precise definition of this distance.

On the one hand, we know the conditional distribution $\imath_- | \imath_+$ as well as $\imath_{spoof}(\imath_-)$. Hence, we can deduce the conditional distribution of $\imath_{spoof} | \imath_+$. We monitor the following two quantities

$$\underbrace{t \longmapsto d\left(\hat{\imath}_-^N(t), \imath_- | \imath_+^N(t)\right)}_{\substack{\text{Distance from the short term imbalance } \hat{\imath}_-^N(t) \\ \text{to the \textbf{legitimate} imbalance } \imath_- \\ \text{given that } \imath_+ \sim \imath_+^N(t)}} \quad \text{and} \quad \underbrace{t \longmapsto d\left(\hat{\imath}_-^N(t), \imath_{spoof} | \imath_+^N(t)\right)}_{\substack{\text{Distance from the short term imbalance } \hat{\imath}_-^N(t) \\ \text{to the \textbf{spoofed} imbalance } \imath_{spoof} \\ \text{given that } \imath_+ \sim \imath_+^N(t)}}$$

In Figure 15 and Figure 16, are the plots thereof for the selected stocks. We mark in red the area where the distance of the short term imbalance to the spoofed one is smaller than the distance to the legitimate one. Since the spoofed imbalance is computed regardless whether spoofing is rewarding or not, we additionally mark in blue the area where spoofing is not worth according to Proposition 3.3, that is when

$$(5.1) \qquad \sup_k \left\{ 2\bar{\rho}\mu^+(1-\bar{\imath})\bar{\imath}w_k - Q_k\bar{\rho} - \nu_k \right\} \le 0$$

where $\bar{\rho}$ represent the mean of the total volume of market orders purchased divided by the total volume available within $N$ ticks on the limit order book, $\bar{\imath}$ is the mean of $\imath_+^N(t)$.
According to Table 1, different stocks have different frequencies, depth of order book as well as volume of incoming market orders. Our approach does take these different factors into account and are reflected into the different plots. According to the previous results, spoofing is more likely to happen if either $\mu^+$ – the overall price impact – or if $\rho$ – the relative size of the market orders with respect to the liquidity present in the limit order book – is large. We recap for the stocks under study some of their key aspects as well as $\bar{\rho}$, the total amount of market orders within the observed time window relative to the average liquidity available within the observed depth.

| Stock | $f$ | Depth | $\bar{\rho}$ | $\mu^+$ |
|-------|-----|-------|--------------|---------|
| AEM | 4 | 4 | 54.1% | 0.411 |
| BB | 38 | 4 | 17.8% | 0.467 |
| BMO | 11 | 4 | 62.9% | 0.103 |
| CNR | 6 | 4 | 89.5% | 0.398 |
| CPG | 53 | 5 | 26.4% | 0.100 |
| FNV | 3 | 5 | 81.8% | 0.404 |
| FR | 60 | 3 | 21.4% | 0.209 |
| PPL | 26 | 4 | 54.1% | 0.076 |
| TD | 20 | 4 | 48.2% | 0.118 |
| VET | 6 | 5 | 75.4% | 0.315 |

**Table 3.** Frequency, depth, market order volume relative to liquidity available $\bar{\rho}$, price impact $\mu^+$ for the studied Stocks.

Overall, we do not observe many crossings, if ever. If a significant crossing happens, it is isolated showing some abnormal behavior. This is particularly obvious for CPG and PPL. An outlier in this series of observations is the stock TD where many crossings happens. TD is a stock which is particularly active with a high rate of market orders – the largest in our study group. Furthermore the price impact is particularly low $0.118$. However, from these market specificities, according to the equation (5.1), most of the time it is not worth spoofing. Hence, there remain only one significant crossing out of the area where spoofing would be worthwhile.
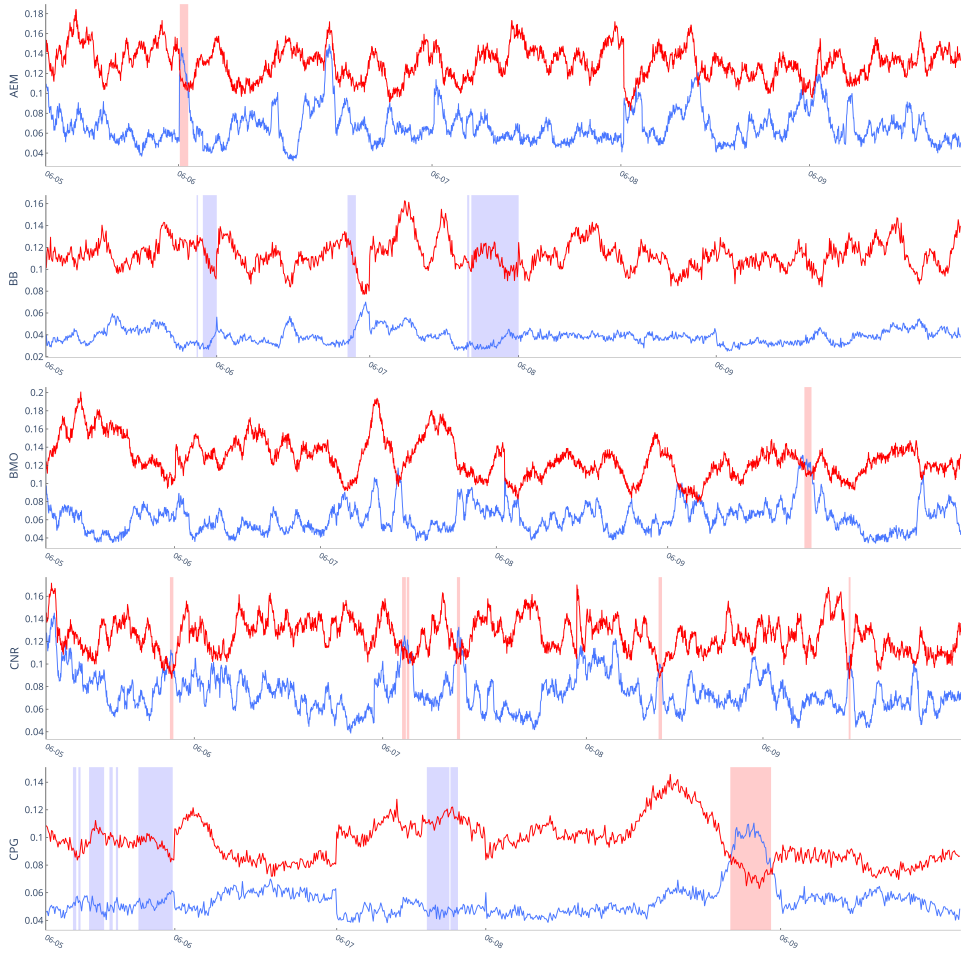
**Figure 15.** Time series of $t \mapsto d\left(\imath_{-}, \hat{\imath}_{-}^{N}(t)|\hat{\imath}_{+}^{N}(t)\right)$ (blue line) and $t_k \mapsto d\left(\imath_{spoof}, \hat{\imath}_{-}^{N}(t)|\hat{\imath}_{+}^{N}(t)\right)$ (red line) from June 5, 2017 to June 9, 2017. Red area is where $d\left(\imath_{spoof}, \hat{\imath}_{-}^{N}|\hat{\imath}_{+}^{N}\right) \leq d\left(\imath_{-}, \hat{\imath}_{-}^{N}|\hat{\imath}_{+}^{N}\right)$ over more than 10 consecutive times. Blue area is where $\sup_k\left\{2\bar{\rho}\mu^{+}(1-\bar{\imath})\bar{\imath}w_k - Q_k\bar{\rho} - \nu_k\right\} \leq 0$ and $\bar{\rho}, \bar{\imath}$ are the mean of $\rho, \imath_N^{+}$ in each window respectively. From top to bottom are stock AEM, BB, BMO, CNR, and CPG respectively.
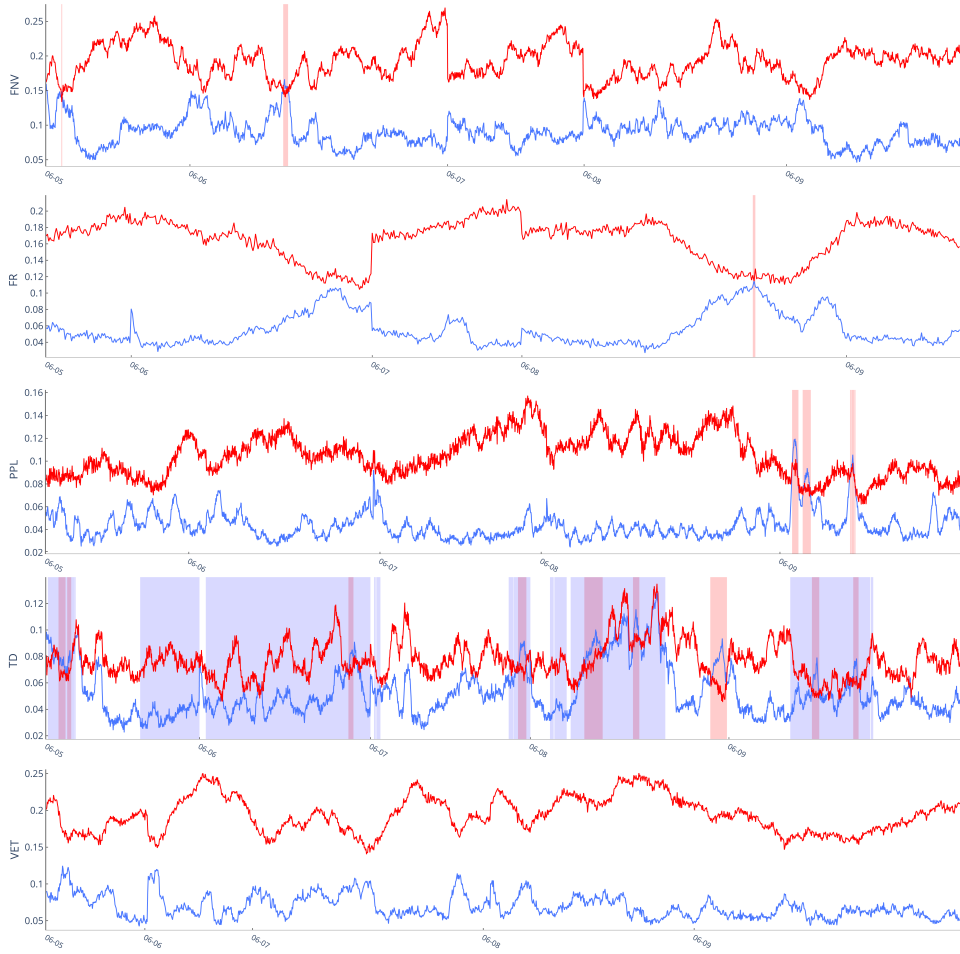
**Figure 16.** Time series of $t_k \mapsto d\left(\imath_-, \hat{\imath}_-^N(t) | \hat{\imath}_+^N(t)\right)$ (blue line) and $t_k \mapsto d\left(\imath_{spoof}, \hat{\imath}_-^N(t) | \hat{\imath}_+^N(t)\right)$ (red line) from June 5, 2017 to June 9, 2017. Red area is where $d\left(\imath_{spoof}, \hat{\imath}_-^N | \hat{\imath}_+^N\right) \leq d\left(\imath_-, \hat{\imath}_-^N | \hat{\imath}_+^N\right)$ over more than 10 consecutive times. Blue area is where $\sup_k \left\{2\bar{\rho}\mu^+(1-\bar{\imath})\bar{\imath}w_k - Q_k\bar{\rho} - \nu_k\right\} \leq 0$ and $\bar{\rho}, \bar{\imath}$ are the mean of $\rho, \imath_N^+$ in each window respectively. From top to bottom are stock FNV, FR, PPL, TD and VET respectively.

## 6. Conclusion

In this paper we address the question of assessing quantitatively eventual spoofing behavior in high frequency trading. In a stylised setting we present how a spoofing strategy from a market taker or maker is designed by manipulating the imbalance at different depth level to impact the subsequent price movement. We provide and discuss the conditions for the market to allow for spoofing manipulations. We subsequently solve the optimization problem from a spoofer perspective and derive/discuss the resulting imbalance after spoofing as a function of the market parameters. We calibrate the weighted imbalance and price movement impact to Level 2 data provided by TMX. Using these results we propose a quantification instrument to monitor in real time eventual spoofing behavior on the market using a conditional Wasserstein distance. We illustrate these results on the data provided by TMX.

This approach is by no means a definitive answer to spoofing detection but rather a first take on. The dynamic structure of the limit order book and strategy, the memory dependence

of the parameters over time, as well as the specificities of one market with respect to another one are left to further study. Also left to further studies is the consideration of multiple venue that could also be integrated into this framework. Furthermore, there might be alternative approaches subject to new research directions – monitoring arrival rates of orders, frequency of book/cancelling, etc. – that could complement such a monitoring approach.

## 7. CREDITS

## APPENDIX A. PROOFS

*Proof of Proposition 3.3.* Let $H = \rho a$, and the imbalance $\imath_k(v) = b/(b + a + w_k v)$ with $b = a\bar{\imath}/(1 - \bar{\imath})$. It follows that the gradient of $\imath_k(v)$ is given by

$$\nabla \imath_k(v) = -\frac{b}{(a + b + w_k v)^2} w_k = -\frac{(1 - \bar{\imath})}{a\bar{\imath}} \imath_k^2(v) w_k$$

• If $\bar{\imath} \leq 1/2$, from the previous equations, since $g(x) = x^2/(2a)$, it follows that there is no spoofing manipulation if and only if

$$f(v) := Q_k \frac{v^2}{2a} + Q_k \rho v + v \nu_k - 2\rho a \mu^+ (\bar{\imath} - \imath_k(v)) \geq 0$$

for any $v \geq 0$. Taking the gradient for this function yields

$$\nabla f(v) = Q_k \frac{v}{a} + Q_k \rho + \nu_k - 2\rho \mu^+ \frac{1 - \bar{\imath}}{\bar{\imath}} \imath_k^2(v) w_k$$

which is a monotone functional in $v$. Since $f(0) = 0$, it follows that $f(v) \geq 0$ for any $v$ if and only if $\nabla f(0) \geq 0$ which is equivalent to

$$Q_k \rho + \nu_k \geq 2\rho \mu^+ (1 - \bar{\imath}) \bar{\imath} w_k$$

• If $\bar{\imath} > 1/2$, there is no spoofing manipulation if and only if

$$f(v) := q_k \frac{v^2}{2a} + \rho q_k v - q_k k v - 2\rho a \mu^+ (1/2 - \imath_k(v)) \geq 0$$

for any $v \geq 0$ the gradient of which is given by

$$\nabla f(v) = q_k \frac{v}{a} + \rho q_k - q_k k v - 2\rho \mu^+ \frac{1 - \bar{\imath}}{\bar{\imath}} \imath_k^2(v) w_k$$

Since $f(0) > 0$, as previously argued, it follows that $f(v) \geq 0$ as soon as $\nabla f(0) \geqslant 0$, which yields the same conditions.

$\square$

*Proof of Proposition 3.4.* Adopting the notations $Q := Q_k$, $\nu = \nu_k$, $w = w_k$, $H = \rho a$, the goal is to optimize over $v \geq 0$ the objective function

$$f(v) = (1 - Q) \frac{(\rho a)^2}{2a} + Q \frac{(\rho a + v)^2}{2a} + \rho a \mu^+ (2\imath(v) - 1) + v \nu$$

$$= \frac{(\rho a)^2}{2a} + Q \frac{v^2}{2a} + (Q\rho + \nu) v + \rho a \mu^+ (2\imath(v) - 1)$$

First order condition with Lagrangian $\lambda$ yields

$$Q\frac{v}{a} + (Q\rho + \nu) - 2\rho w\mu^+ \frac{1 - \bar{\imath}}{\bar{\imath}}\imath^2 = \lambda$$

where $\imath := \imath(v)$. Solving as a function of $\imath$ in $(0, 1)$, we get

$$\lambda(\imath) = \left[(Q\rho + \nu) - 2\rho w\mu^+ \frac{1 - \bar{\imath}}{\bar{\imath}}\imath^2\right]^+$$

$$v(\imath) = \frac{a}{Q}\left[2\rho w\mu^+ \frac{1 - \bar{\imath}}{\bar{\imath}}\imath^2 - (Q\rho + \nu)\right]^+$$

Given now the optimal $v(\imath)$ as a function of $\imath$, we solve for $\imath$ such that

$$\frac{1}{\imath} = \frac{a + b + wv(\imath)}{b} = \frac{1}{\bar{\imath}} + w\frac{1 - \bar{\imath}}{a\bar{\imath}}v(\imath) = \frac{1}{\bar{\imath}} + \frac{w}{Q}\frac{1 - \bar{\imath}}{\bar{\imath}}\left[2\rho w\mu^+ \frac{1 - \bar{\imath}}{\bar{\imath}}\imath^2 - (Q\rho + \nu)\right]^+$$

Since the left hand side in strictly decreasing on from $\infty$ to $\frac{1}{\bar{\imath}}$ on $(0, \bar{\imath}]$ and the right hand side is increasing from $\frac{1}{\bar{\imath}}$, on $(0, \bar{\imath}]$, there exists a unique solution which is a cubic root. $\quad\square$

## APPENDIX B. ROUND TRIP SITUATION

In this paper we mainly focus on the spoofing behavior from a market taker's viewpoint. As for a market maker, spoofing behavior might be rewarding as well. However, as seen in the following subsection, the rewards from spoofing are intertwined with the ones from pure market making.

We present a simple situation together with the numerical analysis in a blocked shape setting with the same model assumptions as before. We assume that the potential market maker spoofer acts as follows: At the first stage it decides to spoof with a volume $v$ at depth $k$ on the ask side to drive the price down and acquire an amount $H$ of shares after this price movement. When the market comes back to its steady state, it liquidates $H$ and eventually buys back $v$ if it has been executed. We assume that $v$ and $H$ are decided at the very beginning.[11]

After spoofing a volume $v$ at level $k$, as soon as the price moves the spoofer executes its market order $H$ for a revenue of

$$-\left(H - v1_{\{y>k\}}\right)(p + \delta\Delta + \delta(x + y)) - \delta G_a(H) - \delta(y - k)v1_{\{y>k\}}$$

where $G_a(H) = H^2/(2a)$, $p = (p^+ + p^-)/2$, and $\Delta = (p^+ - p^-)/(2\delta)$ is the effective spread in ticks. The spoofer then waits for the market to return to its steady state and liquidate the resulting inventory with market orders. For ease of computation, we assume that it executes two market orders: One for $H$ and one for $v$ if it has been executed[12] for a revenue of

$$Hp - \delta\Delta H - G_b(H) - 1_{\{y>k\}}\left(v(p + \delta\Delta) + \delta G_a(v)\right)$$

Adding both and integrating yields an average net revenue of

$$\frac{R(H, v)}{\delta} = -H\left(2\Delta + \mu^+\left(2\imath_k(v) - 1\right)\right) - G_a(H) - G_b(H)$$
$$+ Q_k\left[\left(k + \mu^+\left(2\imath_k(v) - 1\right)\right)v - G_a(v)\right]$$
$$= -H\left(2\Delta + \mu^+\left(2\imath_k(v) - 1\right)\right) - \frac{1}{a\bar{\imath}}\frac{H^2}{2} + Q_kv\left[k + \mu^+\left(2\imath_k(v) - 1\right)\right] - \frac{Q_k}{a}\frac{v^2}{2}$$

---

[11] This stylised situation makes strong assumptions and simplifications. First $H$ is decided at time 0 even if it is executed after the price movement. This is to prevent conditional optimization. Second, the liquidation of the inventory $H$ and $v$ occurs separately. Once again, to provide simplified optimization problem, while we could numerically consider a liquidation of the net inventory $H - v$. Finally, a second spoofing could happen at the second stage as in the previous section to liquidate the inventory.

[12] Combining both in terms of $H - v1_{\{y>k\}}$ is cost effective but complicates the exposition of the result.

From this equation, we can derive the following remarks concerning the decision of the spoofer:

- If $v = 0$: This corresponds to the classical situation where a market maker takes advantage of the temporary market movement to execute a market order and cash out at a later time when the market comes back to its steady state. Clearly, it gets a positive gain if and only if

$$\bar{\imath} \leq \frac{1}{2} - \frac{\Delta}{\mu^+}$$

In particular, if the effective spread $\Delta$ is large, or if $\mu^+$ is small, then it is impossible or the initial imbalance should be very small. In the case where this happens, then $\bar{H}^*$ is given by

$$\bar{H}^* = a\bar{\imath} \left( 2\Delta + \mu^+ \left( 2\bar{\imath} - 1 \right) \right)^-$$

with corresponding revenue of

$$\bar{R}^* = \frac{1}{2} \left[ a\bar{\imath} \left( 2\Delta + \mu^+ \left( 2\bar{\imath} - 1 \right) \right)^- \right]^2$$

- If $H = 0$: This corresponds to the classical situation where a market maker posts limit orders at a given depth to gain from possible fluctuations. This results in corresponding average revenue given $v$ of

$$\hat{R}(v) = Q_k v \left( k + \mu^+ \left( 2\imath_k(v) - 1 \right) \right) - \frac{Q_k}{a} \frac{v^2}{2}$$

From this equation, even if the spoofer gets a positive gain of $k$ ticks buy executing its order, it will drive the imbalance $\imath_k(v)$ below $1/2$ and face adverse price movement that will offset its gains. The optimal $\hat{v}^* = \hat{v}^*(\bar{\imath})$ in that situation is not explicit, but can be easily numerically implemented and corresponds to an optimal revenue of

$$\hat{R}^* = Q_k \hat{v}^* \left( k + \mu^+ \left( 2\imath_k(\hat{v}^*) - 1 \right) \right) - \frac{Q_k}{a} \frac{\hat{v}^2}{2}$$

In general, solving for the optimal $H$ is straightforward with

$$H^* = a\bar{\imath} \left( 2\Delta + \mu^+ \left( 2\imath_k(v) - 1 \right) \right)^-$$

and corresponding average revenue:

$$\frac{R(v)}{\delta} = \frac{1}{2} \left[ a\bar{\imath} \left( 2\Delta + \mu^+ \left( 2\imath_k(v) - 1 \right) \right)^- \right]^2 + Q_k v \left( k + \mu^+ \left( 2\imath_k(v) - 1 \right) \right) - \frac{Q_k}{a} \frac{v^2}{2}$$

These two effects are difficult to disentangle from a truly spoofing gain when $H$ as well as $v$ are strictly positive. However, this can be done numerically and the results are presented in Figure 17, where the spoofing region – $H > 0$ as well as $v > 0$ – is indicated.

We can however draw some stylised facts about the spoofing behavior from this market maker viewpoint. The impact of the different parameters – initial imbalance $\bar{\imath}$, probability of getting executed $Q_k$, local sensitivity of imbalance on the price impact $w = w_k$ as well as overall price deviation $\mu^+$ are similar to the previous case. However, in addition to the previous part, the effective spread $\Delta$ acts negatively on the spoofing opportunity in that context. Indeed, a positive market order $H$ is only triggered if $\mu^+ (2\imath^* - 1) \leq -2\Delta$, which requires a spoofed imbalance satisfying

$$\imath_{spoof} \leq -\frac{\Delta}{\mu^+} + \frac{1}{2}$$

If $\Delta$ is too large or $\mu^+$ too low, a spoofing strategy is no longer rewarding.
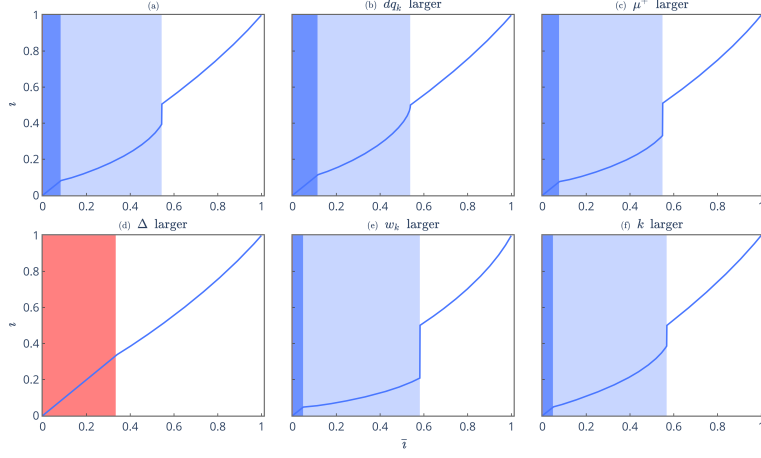
**Figure 17.** $\imath_{spoof}$ as a function of $\bar{\imath}$ and in $(a)$ $\mu^+ = 3, \Delta = 0, k = 1, w_k = 0.1, dq_y = 0.001$ for all $y \geq k$. One parameter is increased each time with respect to $(a)$ where $(b)$: $dq_y = 0.002$ for all $y \geq k$; $(c)$: $\mu^+ = 4$; $(d)$: $\Delta = 2$; $(e)$: $w_k = 0.4$; $(f)$: $k = 4$. Red area: $v^* = 0, H^* = 0$; Dark blue area: $v^* = 0, H^* > 0$; Light blue area: $v^* > 0, H^* > 0$; White area: $v^* > 0, H^* = 0$.

## APPENDIX C. GOODNESS OF FIT

The model for the price movement is conditioned on the imbalance. So, as for the goodness of fit, we compare the empirical price change distribution with the fitted one conditioned on different level of imbalance.
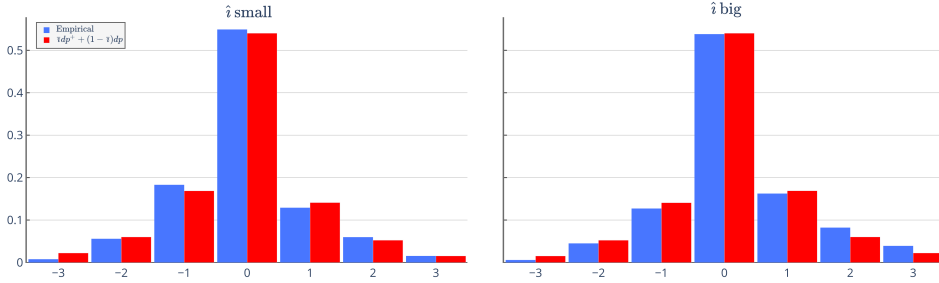


**Figure 18.** Empirical price change conditioned on the imbalance (blue) and $\bar{\imath} dp^+ + (1 - \bar{\imath}) dp^-$ (red). Left panel, $\hat{\imath} < 0.298$ and right panel, $\hat{\imath} \geq 0.648$, which correspond to the low and top 5% quantiles.

Since we fit a conditional distribution, we bucket the observed imbalance into 20 equidistant quantile intervals and perform a $\mathcal{X}^2$ test of goodness of fit between the empirical distribution conditioned on the imbalance within this bucket and $\bar{\imath} dp^+ + (1 - \bar{\imath}) dp^-$ where $\bar{\imath}$ is taken as the mid value of the imbalance for the corresponding bucket. For each bucket, we have about 500 data points which we test against the null hypothesis that the empirical price change distribution is consistent with $\bar{\imath} dp^+ + (1 - \bar{\imath}) dp^-$. Table 4 shows the $\mathcal{X}^2$-statistics and $p$-values in each bucket.[13] 3 buckets exceed the bound (marked in black) where the null hypothesis is rejected while for the 17 others, the null hypothesis sustains.

---

[13]Significance level 0.05 corresponds to a $\mathcal{X}_6^2(0.95) = 12.592$.

| $l$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $\bar{\imath}_l$ | 0.245 | 0.321 | 0.356 | 0.381 | 0.4 | 0.416 | 0.43 | 0.443 | 0.456 | 0.468 |
| $\mathcal{X}^2$-statistic | 7.112 | 7.635 | 12.476 | 6.253 | 6.658 | 8.041 | 4.089 | 2.814 | 12.56 | 4.177 |
| $p$-value | 0.311 | 0.266 | 0.052 | 0.395 | 0.354 | 0.235 | 0.665 | 0.832 | 0.051 | 0.653 |

| $l$ | 11 | 12 | 13 | 14 | 15 | **16** | 17 | **18** | 19 | **20** |
|---|---|---|---|---|---|---|---|---|---|---|
| $\bar{\imath}_l$ | 0.48 | 0.493 | 0.506 | 0.52 | 0.534 | **0.549** | 0.568 | **0.591** | 0.624 | **0.702** |
| $\mathcal{X}^2$-statistic | 9.372 | 8.775 | 5.219 | 3.014 | 7.634 | **16.046** | 2.334 | **17.835** | 6.949 | **17.953** |
| $p$-value | 0.154 | 0.187 | 0.516 | 0.807 | 0.266 | **0.014** | 0.887 | **0.007** | 0.326 | **0.006** |

**Table 4.** Chisquare test on stock BMO from June 5, 2017 to June 9, 2017.

## APPENDIX D. COMPUTATION OF $\imath_{spoof}$

For $t$ in $\Pi$, we compute $\hat{\imath}_-(t)$ and $\hat{\imath}_+(t)$, the imbalance before and after a market order, as follows:

$$\hat{\imath}_-(t) = \frac{\sum\limits_{k \leq N} \sum\limits_{t-f \leq s < t} w_k \bar{v}_k^-(s)}{\sum\limits_{k \leq N} \sum\limits_{t-f \leq s < t} w_k \left(\bar{v}_k^-(s) + \bar{v}_k(s)\right)}$$

$$\hat{\imath}_+(t) = \frac{\sum\limits_{k \leq N} \sum\limits_{t \leq s < t+1} w_k \bar{v}_k^-(s)}{\sum\limits_{k \leq N} \sum\limits_{t \leq s < t+1} w_k \left(\bar{v}_k^-(s) + \bar{v}_k(s)\right)}$$

For each $\hat{\imath}_-(t)$, the optimal spoofing strategy $v_{spoof}$ can be solved explicitly from

$$v_{spoof,k} = 1 + \frac{(1 - \hat{\imath}_-(t))\, w_k}{Q} \left[2\rho_t w_k \mu^+ \frac{1 - \hat{\imath}_-(t)}{\hat{\imath}_-(t)} \imath^2 - (Q_k \rho_t + \nu)\right]^+$$

where $\imath = \frac{b_t}{b_t + a_t + w_k v_{spoof,k}}$, $a_t$ is the average size of the limit order book $f$ seconds before a market order

$$a_t = \frac{\sum\limits_{k=1}^{N} \sum\limits_{t-f \leq s < t} \bar{v}_k(s) \Delta s}{Nf}$$

$$\rho_t = \frac{\sum\limits_{k=1}^{N} \sum\limits_{t-f \leq s < t} H_s}{a_t}$$

where $H_s$ is the market order volume at time $s$. In the same way, we can define $b_t$ and

$$\imath_{spoof}(\hat{\imath}_-(t)) = \frac{b_t}{b_t + a_t + \sum_k w_k v_{spoof,k}}$$

## APPENDIX E. WASSERSTEIN DISTANCE, KERNEL APPROXIMATION AND CONDITIONAL ESTIMATION

For two distributions $\mu$ and $\nu$ the 2-Wasserstein distance is defined as

$$W_2(\mu, \nu) = \left(\inf\left\{\int (x-y)^2 \pi(dx, dy) : \pi_1 \sim \mu, \pi_2 \sim \nu\right\}\right)^{1/2}$$

$$= \left(\int_0^1 (q_\mu(\alpha) - q_\nu(\alpha))^2 \, d\alpha\right)^{1/2}$$

From a generic perspective, the conditional distance we consider is as follows: If we assume that

$$(\imath_-, \imath_+) \sim K(y, dx) \otimes \mu(dy)$$

where $\mu \sim \imath_+$ and $K(y, \cdot) \sim \imath_- | \imath_+ = y$, it follows that

$$(\imath_{spoof}(\imath_-), \imath_+) \sim K_{spoof}(y, dx) \otimes \mu(dy) \quad \text{where} \quad K_{spoof}(y, \cdot) = K(y, \cdot) \circ \imath_{spoof}^{-1}$$

Hence given $\imath_+ = y$, we have

$$W_2(K(y, \cdot), K_{spoof}(y, \cdot)) = \left( \int_0^1 \left( q_{K(y,\cdot)}(\alpha) - \imath_{spoof}\left( q_{K(y,\cdot)}(\alpha) \right) \right)^2 d\alpha \right)^{1/2}$$

Heuristically we wish to monitor the following two quantities

$$\underbrace{W_2\left(\imath_-, \hat{\imath}_-^N\right) | \imath_+^N}_{\substack{\text{Distance from the short term imbalance } \hat{\imath}_-^N \\ \text{to the \textbf{equilibrium} imbalance } \imath_- \\ \text{given that } \imath_+ \sim \imath_+^N}} \qquad \text{and} \qquad \underbrace{W_2\left(\imath_{spoof}(\imath_-), \hat{\imath}_-^N\right) | \imath_+^N}_{\substack{\text{Distance from the short term imbalance } \hat{\imath}_-^N \\ \text{to the \textbf{spoofed} imbalance } \imath_{spoof} \\ \text{given that } \imath_+ \sim \imath_+^N}}$$

From the data, we can calibrate the joint distribution $(\imath_-, \imath_+)$ as well as $(\imath_{spoof}(\imath_-), \imath_+)$.[14] Hence, we have a parametrization of $K(y, dx)$ and $K_{spoof}(y, dx)$ for every $y$. However for each value $\hat{\imath}_+(l)$ from the discrete distribution $\hat{\imath}_+^N$ we only have a single sample point $\hat{\imath}_-(l)$ at hand. In order to overcome this problem we bucket the values of $\hat{\imath}_+(l)$ in the sample of $\hat{\imath}_+^N$ into several equidistant quantile intervals to get a Kernel approximation of $\hat{\imath}_-^N$.

The monitoring strategy at a given time $t_k$ in $\Pi$ is given as follows

1 - Consider the discrete short term joint distribution $(\hat{\imath}_-^N(t_k), \hat{\imath}_+^N(t_k))$ given by the sample

$$(\imath_-(s), \imath_+(s)), \quad s = t_k, \ldots, t_{k-N+1}$$

of the last $N$ pairs of imbalances before time $t$.

2 - We define the following $L$ buckets of equal cardinality $N/L$

$$J_l = \left\{ s: s = t_k, \ldots, t_{k-N+1}, q_{\hat{\imath}_+^N}\left( \frac{l-1}{L} \right) \leq \imath_+(s) < q_{\hat{\imath}_+^N}\left( \frac{l}{L} \right) \right\}, \quad l = 1 \ldots, L$$

as well as the mid point of each

$$\imath_l = \frac{L}{N} \sum_{s \in J_l} \imath_+(s)$$

3 - For each $l$, we generate a random sample $\imath_-^{N,L}$ and $\imath_{spoof}^{N,L}$ of $N/L$ points each drawn from $K(\imath_l, \cdot)$ and $K_{spoof}(\imath_l, \cdot)$, respectively.

4 - For each $l$ we compute the Wasserstein distances

$$W_2\left( \imath_-^{N,l}, \hat{\imath}_-^{N,L} \right) \quad \text{and} \quad W_2\left( \imath_{spoof}^{N,l}, \hat{\imath}_-^{N,L} \right)$$

where $\hat{\imath}_-^{N,L}$ is the discrete distribution out of the sample $\hat{\imath}_-(s)$ for $s$ in $J_l$. This is an approximation for the Wasserstein distance

$$W_2\left( \imath_-, \hat{\imath}_n^N \right) | \imath_+^N \approx \imath_l \quad \text{and} \quad W_2\left( \imath_{spoof}, \hat{\imath}_n^N \right) | \imath_+^N \approx \imath_l$$

---

[14]The former fits well with a joint normal distribution, while the second one with a skewed normal distribution, see Figure 19. Other parametrization could eventually be used too.

5 - we aggregate all together and define the indicators

$$d\left(\imath_-, \hat{\imath}_-^N | \hat{\imath}_+^N\right) := \frac{1}{L}\sum_{l=1}^{L} W_2\left(\imath_-^{N,l}, \hat{\imath}_-^{N,L}\right)$$

$$d\left(\imath_{spoof}, \hat{\imath}_-^N | \hat{\imath}_+^N\right) := \frac{1}{L}\sum_{l=1}^{L} W_2\left(\imath_{spoof}^{N,l}, \hat{\imath}_-^{N,L}\right)$$

**Remark E.1.** *To enhance the accuracy of this indicator, we run step 3 to 5 a couple of times with different samples and average again.*

As for the Kernel approximation, we fit $(\imath_-, \imath_+)$ to a bivariate normal distribution

$$(\imath_-, \imath_+) \sim \mathcal{N}\left(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho\right)$$

thus the conditional distribution is also a normal distribution and $K(y, \cdot)$ is its density function

$$\imath_- | \imath_+ = y \sim \mathcal{N}\left(\mu_1 + \frac{\sigma_1}{\sigma_2}\rho(y - \mu_2), (1 - \rho^2)\sigma_1^2\right)$$

Similarly, we fit $(\imath_{spoof}(\imath_-), \imath_+)$ to a bivariate skewnormal distribution

$$(\imath_{spoof}(\imath_-), \imath_+) \sim \mathcal{SN}\left(\alpha, \xi, \Omega\right)$$

where $\alpha = [\alpha_1, \alpha_2]^\top$, $\xi = [\xi_1, \xi_2]^\top$, $\Omega = \begin{bmatrix} w_1 & w \\ w & w_2 \end{bmatrix}$. It can be derived that

$$K_{spoof}(y, \cdot) = \phi\left(\frac{\cdot - \xi_1^c}{\sqrt{w_{11.2}}}\right)\frac{\Phi\left(\alpha_1\sqrt{\omega_1}(\cdot - \xi_1^c) + x_0'\right)}{\Phi(x_0)}$$

where $\phi$, $\Phi$ are the density function and cumulative distribution function of a standard normal distribution, and

$$\xi_1^c = \xi_1 + \frac{\omega}{\omega_2}(y - \xi_2), \quad w_{11.2} = w_1 - \frac{w^2}{w_2}$$

$$\bar{\alpha}_2 = \frac{\alpha_2 + \sqrt{\frac{w^2}{w_1 w_2}}\alpha_1}{\sqrt{1 + \frac{w_{11.2}}{w_1}\alpha_1^2}}, \quad x_0 = \frac{\bar{\alpha}_2}{\sqrt{w_2}}(y - \xi_2)$$

$$x_0' = \sqrt{1 + \frac{w_{11.2}}{w_1}\alpha_1^2}\, x_0$$

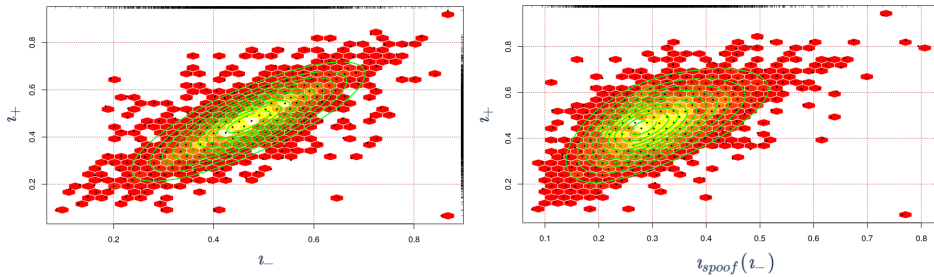Figure 19 shows that $K$ and $K_{spoof}$ fit well with real data in Figure 13.



**Figure 19.** Left panel: Empirical joint distribution of $(\imath_-, \imath_+)$ and $K$ (green contours) . Right panel: Joint distribution of $(\imath_{spoof}(\imath_-), \imath_+)$ and $K_{spoof}$ (green contours).

REFERENCES

[1] A. Alfonsi and A. Schied. Optimal trade execution and absence of price manipulations in limit order book models. *SSRN Electronic Journal*, 02 2010.

[2] A. Alfonsi, A. Fruth, and A. Schied. Optimal execution strategies in limit order books with general shape functions. *Quantitative Finance*, 10, 08 2007.

[3] F. Allen and D. Gale. Stock-price manipulation. *Review of Financial Studies*, 5: 503–29, 02 1992.

[4] K. Bechler and M. Ludkovski. Order flows and limit order book resiliency on the meso-scale. *Market Microstructure and Liquidity*, 3(03n04):1850006, 2017.

[5] Y. Cao, Y. Li, S. Coleman, A. Belatreche, and T. M. McGinnity. Detecting price manipulation in the financial market. In *2014 IEEE Conference on Computational Intelligence for Financial Engineering Economics (CIFEr)*, pages 77–84, 2014.

[6] Y. Cao, Y. Li, S. Coleman, A. Belatreche, and T. M. McGinnity. Adaptive hidden markov model with anomaly states for price manipulation detection. *IEEE Transactions on Neural Networks and Learning Systems*, 26(2):318–330, 2015.

[7] Á. Cartea, R. Donnelly, and S. Jaimungal. Enhancing trading strategies with order book signals. *Applied Mathematical Finance*, 25(1):1–35, 2018.

[8] A. Cartea, S. Jaimungal, and Y. Wang. Spoofing and price manipulation in order-driven markets. *Applied Mathematical Finance*, pages 1–32, 02 2020. doi: 10.1080/1350486X.2020.1726783.

[9] R. Cont, A. Kukanov, and S. Stoikov. The price impact of order book events. *Journal of Financial Econometrics*, 12(1):47–88, 2014.

[10] J. Gatheral. No-dynamic-arbitrage and market impact. *Quantitative Finance*, 10(7): 749–759, 2010.

[11] J. Gatheral and A. Schied. Dynamical models of market impact and algorithms for order execution. *SSRN Electronic Journal*, 01 2013.

[12] M. Gould and J. Bonart. Queue imbalance as a one-tick-ahead price predictor in a limit order book. *Market Microstructure and Liquidity*, 12 2015.

[13] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. ISSN 1476-4687.

[14] O. Hart and D. Kreps. Price destabilizing speculation. *Journal of Political Economy*, 94(5):927–52, 1986.

[15] R. Jarrow. Market manipulation, bubbles, corners, and short squeezes. *Journal of Financial and Quantitative Analysis*, 27:311–336, 09 1992.

[16] A. Kyle. Continuous auctions and insider trading. *Econometrica*, 53:1315–1336, 1985.

[17] E. J. Lee, K. S. Eom, and K. S. Park. Microstructure-based manipulation: Strategic behavior and performance of spoofing traders. *Journal of Financial Markets*, 16(2): 227–252, 2013.

[18] A. Lipton, U. Pesavento, and M. Sotiropoulos. Trade arrival dynamics and quote imbalance in a limit order book. *Preprint*, 12 2013.

[19] M. McKerns, L. Strand, T. Sullivan, A. Fang, and M. Aivazis. Building a framework for predictive science. *Proceedings of the 10th Python in Science Conference, 2011*, 2011.

[20] E. M. Miranda, P. McBurney, and M. J. Howard. Learning unfair trading: A market manipulation analysis from the reinforcement learning perspective. *2016 IEEE Conference on Evolving and Adaptive Intelligent Systems (EAIS)*, pages 103–109, 2016.

[21] J. Reback, W. McKinney, jbrockmendel, J. V. den Bossche, T. Augspurger, P. Cloud, gfyoung, Sinhrks, A. Klein, M. Roeschke, S. Hawkins, J. Tratner, C. She, W. Ayd, T. Petersen, M. Garcia, J. Schendel, A. Hayden, MomIsBestFriend, V. J. L. Rechenzentrum, P. Battiston, S. Seabold, chris b1, h vetinari, S. Hoyer, W. Overmeire, alimcmaster1, K. Dong, C. Whelan, and M. Mehyar. pandas-dev/pandas: Pandas. Software, Feb. 2020.

[22] G. Shorter and R. Miller. High-frequency trading: Background, concerns, and regulatory developments. Technical report, Congressional Research Service, 04 2015.

[23] J. A. Sirignano. Deep learning for limit order books. *Quantitative Finance*, 19(4): 549–570, 2019.

[24] V. van Kervel. Competition for order flow with fast and slow traders. *The Review of Financial Studies*, 28(7):2094–2127, 2015.

[25] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, İ. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, A. Vijaykumar, A. P. Bardelli, A. Rothberg, A. Hilboll, A. Kloeckner, A. Scopatz, A. Lee, A. Rokem, C. N. Woods, C. Fulton, C. Masson, C. Häggström, C. Fitzgerald, D. A. Nicholson, D. R. Hagen, D. V. Pasechnik, E. Olivetti, E. Martin, E. Wieser, F. Silva, F. Lenders, F. Wilhelm, G. Young, G. A. Price, G.-L. Ingold, G. E. Allen, G. R. Lee, H. Audren, I. Probst, J. P. Dietrich, J. Silterra, J. T. Webber, J. Slavič, J. Nothman, J. Buchner, J. Kulick, J. L. Schönberger, J. V. de Miranda Cardoso, J. Reimer, J. Harrington, J. L. C. Rodríguez, J. Nunez-Iglesias, J. Kuczynski, K. Tritz, M. Thoma, M. Newville, M. Kümmerer, M. Bolingbroke, M. Tartre, M. Pak, N. J. Smith, N. Nowaczyk, N. Shebanov, O. Pavlyk, P. A. Brodtkorb, P. Lee, R. T. McGibbon, R. Feldbauer, S. Lewis, S. Tygier, S. Sievert, S. Vigna, S. Peterson, S. More, T. Pudlik, T. Oshima, T. J. Pingel, T. P. Robitaille, T. Spura, T. R. Jones, T. Cera, T. Leslie, T. Zito, T. Krauss, U. Upadhyay, Y. O. Halchenko, Y. Vázquez-Baeza, and S. . Contributors. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods*, 17(3):261–272, March 2020. ISSN 1548-7105.

[26] Y. Wang. Strategic spoofing order trading by different types of investors in the futures markets. *Wall Street Journal*, 2015.

[27] Wes McKinney. Data Structures for Statistical Computing in Python. In Stéfan van der Walt and Jarrod Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56 – 61, 2010. doi: 10.25080/Majora-92bf1922-00a.

[28] K. Xu, M. Gould, and S. Howison. Multi-level order-flow imbalance in a limit order book. *Preprint*, 07 2019.

SHANGHAI JIAO TONG UNIVERSITY, SHANGHAI, CHINA
*Email address*: taoxuan@sjtu.edu.cn

WESTERN UNIVERSITY, CANADA
*Email address*: aday46@uwo.ca

SHANGHAI JIAO TONG UNIVERSITY (NOW AT PING AN TECHNOLOGY), SHANGHAI, CHINA
*Email address*: lan.ling@sjtu.edu.cn
*Email address*: linglan987@pingan.com.cn

SHANGHAI JIAO TONG UNIVERSITY, SHANGHAI, CHINA
*Email address*: sdrapeau@saif.sjtu.edu.cn
*URL*: http://www.samuel-drapeau.info