

9. TRANSACTION COSTS

9.1. Limit Order Books. This is not a course on market microstructure, but we will need to understand the basics of trading in limit-order books to continue. For subsection 9.1, a good background reference is Gould et al. (2013).

The majority of the world’s transactions, by volume, are enacted in marketplaces which take the form of a continuous limit-order book (LOB). This means that almost anyone can participate in the market, by quoting prices at which they would buy or sell the asset. Quotes come flowing in, and an automated computer system listens to those quotes and uses them to build a representation of the supply and demand curves, more or less as you learned them in undergraduate micro-economics. The main difference is that the price levels are discrete, and the “curves” are continuously changing as new quotes come in and old ones are canceled.

Definition 9.1. A *limit order* or *quote* $x = (p_x, q_x, t_x)$ submitted at time t_x with price p_x and size $q_x > 0$ (respectively, $q_x < 0$) is a commitment to buy (respectively, sell) up to $|q_x|$ units of the traded asset at a price no greater than (respectively, no less than) p_x .

In this section, we will sometimes abbreviate “limit order” to simply “order” and assume all orders are limit orders.

Definition 9.2. The *lot size* ℓ of an LOB is the smallest amount of the asset that can be traded. All quotes must arrive with a size

$$q_x \in \{\pm k\ell \mid k = 1, 2, \dots\}$$

The *tick size* of an LOB is the smallest permissible price interval between different orders within it. All orders must arrive with a price that is specified to the accuracy of one tick.

The lot size and tick size are called *resolution parameters*.

When a buy (respectively, sell) order x is submitted, an LOB’s matching engine checks whether it is possible to match x to some other previously submitted sell (respectively, buy) order. If so, the matching occurs immediately. If not, x becomes *active*, and it remains active until either it becomes matched to another incoming sell (respectively, buy) order or it is canceled. Cancellation usually occurs because the owner of an order no longer wishes to offer a trade at the stated price, but rules governing a market can also lead to the cancellation of active orders.

Let $\mathcal{L}(t)$ denote the set of all active orders in a market at time t . The active orders in an LOB $\mathcal{L}(t)$ can be partitioned into the set of active buy orders $\mathcal{B}(t)$, for which $q_x > 0$, and the set of active sell orders $\mathcal{A}(t)$, for which $q_x < 0$. An

LOB can then be considered as a set of queues, each of which consists of active buy or sell orders at a specified price. Usually, but not always, in equity markets the queues are time-priority, or in other words they are FIFO queues. It follows that the first to join one of these queues (the trader who can react more quickly) has an advantage. Note that you cannot gain priority by posting a quote 10^{-9} cents above or below someone else's quote, due to the tick size.

Definition 9.3. The *bid price* at time t is the highest stated price among active buy orders at time t .

$$b(t) := \max_{x \in \mathcal{B}(t)} p_x.$$

The *ask price* at time t is the lowest stated price among active sell orders at time t .

$$a(t) := \min_{x \in \mathcal{A}(t)} p_x.$$

The bid-ask spread at time t is $s(t) := a(t) - b(t)$. The mid price at time t is $m(t) := [a(t) + b(t)]/2$.

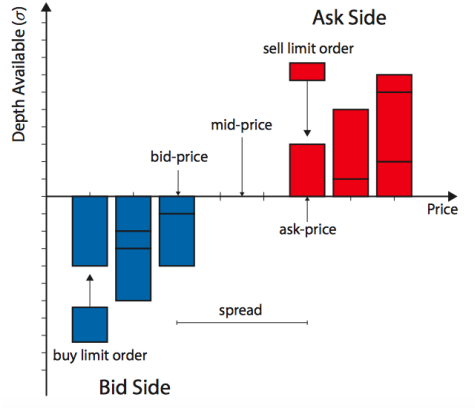


FIGURE 9.1. Schematic of an abstract limit-order book as depicted by Gould et al. (2013).

In the above schematic, note again the similarity to supply and demand curves from micro-economics.

Orders that result in immediate matching upon submission are known as *market orders*. This terminology is used only to emphasize whether an incoming order triggers an immediate matching or not. There is no fundamental difference between a limit order and a market order. Trading via limit orders which are not expected to immediately match is called *trading passively*. The opposite – submitting a sequence of market orders – is referred to as being *aggressive*.

In an LOB, traders are able to choose between submitting limit orders and submitting market orders. Limit orders stand a chance of matching at better prices than do market orders, but they also run the risk of never being matched. Conversely, market orders never match at prices better than $b(t)$ and $a(t)$, but they do not face the inherent uncertainty associated with limit orders. An LOB's bid-ask spread $s(t)$ is a measure of how highly the market values the immediacy and certainty associated with market orders versus the waiting and uncertainty associated with limit orders.

Foucault, Kadan, and Kandel (2005) argued that the popularity of LOBs is due in part to their ability to allow some traders to demand immediacy, while simultaneously allowing others to supply immediacy to those who later require it. Most traders use a combination of both limit orders and market orders; they select their actions for each situation based on their individual needs at that time.

9.2. Parent Orders and Slippage.

Definition 9.4. A *parent order* $\mathfrak{o} = (q, \tau, \tau')$ will be an instruction to buy or sell a fixed quantity q of a certain asset over the time window $[\tau, \tau']$. Per convention, an order to *sell* has $q < 0$ while an order to *buy* has $q > 0$. A parent order is also called a *metaorder*. A parent order may be split into *child orders*; the structure of a child order makes it mathematically equivalent to a parent order, but the language helps the user not to get confused when discussing order-splitting strategies. A *fill* is a statement that part of the buying or selling in a particular (parent or child) order has been completed at a certain time and a certain price, usually in the form of a single matching event generated by the exchange's matching engine.

This terminology is especially suited to cash equities, bonds, and equity options. For equities, the “quantity” is an *integer* in units of shares, while for options, the quantity is typically in units of contracts.

Thus the parent order \mathfrak{o} , if it were completed, would lead to a sequence of fills $\{f_i : i = 1 \dots n_f\}$ where each fill

$$f_i = (\mathfrak{o}, t_i, n_i, p_i)$$

is made up of the parent order, the time $t_i \in [\tau, \tau']$, the number n_i of shares filled, and the price p_i at which they were filled. Let $\text{pa}(f)$ denote the parent order that generated the fill.

In other words a fill $f_i = (\mathfrak{o}, t_i, n_i, p_i)$ is a statement that n_i shares have been exchanged for cash in the amount of $n_i p_i$ dollars (or other numeraire currency) at time t_i , as part of the parent order \mathfrak{o} . Also, we assume all of the fills associated to

a fixed parent order have $\text{sgn}(n_i) = \text{sgn}(q)$, as is logical. Thus necessarily

$$\sum_i |n_i| = |\sum_i n_i| \leq |q|.$$

The total amount filled is $\sum_{i=1}^{n_f} n_i$ which could be smaller in magnitude than $|q|$; in other words, not every parent order is completely filled. For example, a parent order could be an unrealistically large number of shares in a very illiquid stock. The execution algo could either decide not to fully fill based on certain agreed-upon limits, or it could simply fail to locate the shares. If the order involves taking a short position in a stock that is hard to borrow, failing to locate is a common occurrence.

Definition 9.5. An *execution algorithm* or, in common usage, an *algo*, is a means of creating a sequence of fills for any given parent order, ie. a mapping from $\mathfrak{o} \rightarrow \{f_i\}$. Equivalently it is a means of choosing n_i and t_i in the sequence $f_i = (\mathfrak{o}, t_i, n_i, p_i)$.

For any list of pairs $L = \{(n_i, p_i) : i = 1, \dots, N_L\}$ where $p_i > 0$ are prices and $n_i \in \mathbb{N}$, we define

$$\text{vwap}(L) = \frac{\sum_i n_i p_i}{\sum_i n_i}.$$

where vwap stands for “volume-weighted average price.” In one common example, L is the list of all trade prices, with the volume transacted at each price, for a particular (stock, day). For analysis of intraday patterns, one could take L to be a subset of this data over a finer-grained time period, such as a minute.

The vwap of an order is defined to be the vwap of the sequence of fills used to fill the order:

$$\text{vwap}(\mathfrak{o}) = \text{vwap}\{f : \text{pa}(f) = \mathfrak{o}\}$$

To the portfolio manager, the end result of executing the order is essentially the same as if the entire quantity q had been filled in one shot, at price $\text{vwap}(\mathfrak{o})$. The goal for a *buy* parent order is for $\text{vwap}(\mathfrak{o})$ to be as small as possible, and for a *sell* parent order, the goal is for $\text{vwap}(\mathfrak{o})$ to be as large as possible.

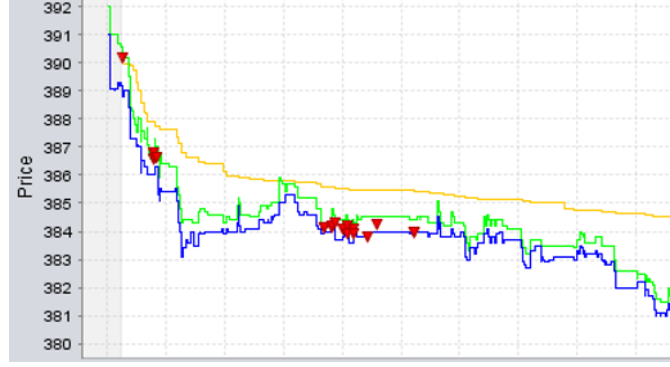


FIGURE 9.2. Schematic of a parent sell order being executed by an algo. Each red triangle denotes a child order being filled. The green (resp blue) line is the national best offer (resp bid). This order has high slippage to arrival mid, because the price moved so quickly.

Definition 9.6. A *benchmark pricing method* is a way of assigning a theoretical price $p_0(\mathfrak{o}) \in \mathbb{R}$ to any order, ie a mapping $p_0 : \mathfrak{D} \rightarrow \mathbb{R}$ where \mathfrak{D} denotes the space of all possible orders. This price need not represent the prices of any actual trades.

One popular benchmark is arrival (mid) price, or simply “arrival price.” This is the last midpoint price available before the order begins being executed. This seems to be the benchmark used in Almgren et al. (2005), for example.

We now come to the most important definition of this section:

Definition 9.7. Given a benchmark pricing method p_0 , the *slippage* of an order relative to this benchmark pricing method is defined as

$$\text{slip}(\mathfrak{o}) = \left(\sum_{i=1}^{n_f} n_i \right) [\text{vwap}(\mathfrak{o}) - p_0(\mathfrak{o})] \quad (9.1)$$

Note that $\text{slip}(\mathfrak{o})$ has units of whatever currency the prices are denominated in, and the sign is such that *positive* slippage denotes a *worse* result for the trader than transacting at the benchmark price. If \mathfrak{D} is an entire set of orders, then

$$\text{slip}(\mathfrak{D}) = \sum_{\mathfrak{o} \in \mathfrak{D}} \text{slip}(\mathfrak{o})$$

As stated, the benchmark pricing method is a mathematical construct and need not correspond to tradable prices. For example, if we take the zero mapping $p_0(\mathfrak{o}) = 0 \forall \mathfrak{o} \in \mathfrak{D}$ then the slippage is simply the total traded notional value. The arrival mid price is also clearly not achievable. Even the average of the midpoint price

over the lifetime of a parent order is not achievable without some special short-term alpha. How exactly are we supposed to consistently transact between the bid and the offer on a sequence of trades that are all in the same direction? On a sequence of trades which includes buys and sells in roughly equal amounts, the vwap of those trades might be closer to the midpoint.

Another interesting benchmark is “actual duration vwap.” This is computed after the algo has finished executing the parent order. Let $[t, t']$ be the interval over which, in the end, the order was executed. The actual-duration vwap is the vwap of all of the trades which occurred in the market over the same interval $[t, t']$. Low slippage to actual-duration VWAP is not necessarily a good thing! If you are causing huge market impact because you are a large percent of the volume, then the vwap of your order will be close to the actual-duration vwap.

For orders that are executed incrementally over the course of an entire day, the full-day vwap is a popular benchmark price. Like the other benchmark prices discussed above, this one is not exactly achievable for all orders. No algo can guarantee that the vwap of your order will equal the aggregate market vwap of the day.

A very simplistic implementation of an algo benchmarked to vwap might be as follows. If the length of the trading day is $5K$ minutes, divide the trading day up into K five-minute bins and let p_i be the fraction of the day’s volume that you predict will occur within each bin, so $\sum_{i=1}^K p_i = 1$ and all $p_i \geq 0$. For an order with total quantity $N = \sum_j n_j$, plan to execute $p_i N$ within the i -th bin. When there is no “interesting” activity going on in a given stock, then algos of this sort can achieve the vwap plus noise, where the noise mostly comes from the difference between the predicted intraday volume pattern and the realized one.

Suppose bad news comes out at close minus 5 minutes. This causes the volume to rise continuously for the next 5 minutes as the price is falling. The VWAP algo did not foresee this, so already executed a fraction of $1 - p_K$ of the order, where p_K is the typical/predicted fraction of volume in the last 5 minutes. Due to the news, the fraction of volume in the last 5 minutes today is much larger than p_K . It follows that, by the time the dust settles at the end of the day, the vwap algo will buy at a higher price than the full-day vwap (so it will underperform the full-day vwap benchmark).

This effect is symmetrical – if the news was good, your vwap algo would have bought at a lower (hence better) price than the full-day vwap. Over time, the aggregate slippage to full-day vwap is then related to how often the direction of your trades were “on the right side” of the news announcements which cause the most volume. This is probably also related to your P/L over the period!

9.3. Slippage as a Cost. Let's now consider a concrete example in which we buy a stock that is going up, and then later sell it, and our benchmark is arrival mid price. Assume that immediately before the buy order begins execution, the bid and the ask are $b_0 < a_0$ and immediately after the subsequent sell order begins execution, the bid and the ask are $b_1 < a_1$. We also assume that both transactions are for 100 shares, and also that $a_0 < b_1$ so the transaction will be profitable.

Hence the benchmark prices are the mids

$$m_0 = (a_0 + b_0)/2 \text{ and } m_1 = (a_1 + b_1)/2.$$

Assume that we are trading aggressively, hence we buy at the ask and sell at the bid, so

$$\text{vwap}(\mathfrak{o}_0) = a_0, \quad \text{vwap}(\mathfrak{o}_1) = b_1.$$

Let π denote our P&L, then

$$\pi = 100(b_1 - a_0) = 100(m_1 - m_0) - \text{slip}(\mathfrak{o}_0) - \text{slip}(\mathfrak{o}_1) \quad (9.2)$$

where $\text{slip}(\mathfrak{o}_0) = 100(a_0 - m_0)$ and $\text{slip}(\mathfrak{o}_1) = 100(m_1 - b_1)$. Note that as per our convention $\text{slip}(\mathfrak{o}_0)$ and $\text{slip}(\mathfrak{o}_1)$ are both positive. Eq (9.2) is easily verified by simple arithmetic.

Furthermore,

$$\pi = 100(b_1 - a_0) = hR - \text{slip}(\mathfrak{o}_0) - \text{slip}(\mathfrak{o}_1), \quad (9.3)$$

$$R := \frac{m_1 - m_0}{m_0}, \quad h := 100m_0.$$

We can interpret (9.3) as stating that if we price our intended holding of 100 shares at the arrival mid, so that our intended holding is worth $h = 100m_0$ dollars, then the P&L π can be represented as the holding value times the return R (which must be price-return using the benchmark price!) minus the total slippage from both orders.

We have proven the following.

Lemma 9.1. Over a sequence of trades which begin and end with zero holdings, the P/L can be represented as the holding value times the return R (defined as price-return using the benchmark price) minus the total slippage from all orders.

Eq.(9.3) generalizes to portfolios in the following way. Suppose that we hold portfolio $h_0 \in \mathbb{R}^n$ now and intend to trade into portfolio $h \in \mathbb{R}^n$. Suppose there are two times t_0 and t_1 , and at t_0 we will begin trading from h_0 to h at t_0 , we will reach h before t_1 , and then at t_1 we will begin liquidating h . Let \mathfrak{O}_i denote all orders started at t_i . Then the P&L can be written

$$\pi = h'R - \text{slip}(\mathfrak{O}_1) - \text{slip}(\mathfrak{O}_2) \quad (9.4)$$

where $R \in \mathbb{R}^n$ is the vector of returns over the interval $[t_0, t_1]$ computed with respect to benchmark price.

We could equivalently write

$$\pi = h'R - \text{slip}(h_0, h) - \text{slip}(h, 0) \quad (9.5)$$

where $\text{slip}(x, y)$ denotes the slippage of the orders needed to trade from portfolio x into portfolio y . The second term $\text{slip}(h, 0)$ is the slippage incurred from liquidating the final portfolio h . As in the simple example above, it is necessary to liquidate the final portfolio to actually realize all profits in dollars; otherwise some portion of the profits will be left as “unrealized” and any unrealized profits will be subject to slippage before they are “realized” or translated to dollars.

Definition 9.8. *Liquidation slippage* of a portfolio h is defined as $\text{slip}(h, 0)$, i.e. the slippage incurred on the full set of orders necessary to convert the holdings entirely to cash. The liquidation slippage of h will be denoted by

$$\text{liqslip}(h) := \text{slip}(h, 0).$$

For a perspective on optimal execution algos with fairly similar notation to ours, see Almgren and Chriss (1999) and Almgren and Chriss (2001).

Note that $\text{slip}(\mathbf{o})$ is not knowable at the order creation time τ (as it involves future prices). For the same reason (it involves future prices), it is hard to predict with high R^2 .

Definition 9.9. A *predictive slippage model* is a model for the conditional density $p(\text{slip}(\mathbf{o}) \mid I_\tau)$ where I_τ denotes the information set available at time τ . Many researchers simply model $\mathbb{E}[\text{slip}(\mathbf{o}) \mid I_\tau]$ directly without modeling the full distribution.

A number of prominent academics have studied the problem of predicting $\text{slip}(\mathbf{o})$ as function of the order quantity q and attributes of the asset being traded. Some attributes that have been found to be predictive include that asset’s volatility, volume, and the window $T = [\tau, \tau']$ over which the orders are filled. One of the most oft-cited such studies is Almgren et al. (2005).

Let’s now suppose that our prior holding in some asset is h_0 dollars and we are considering a trade of

$$\delta := h - h_0$$

so that our new holding will be h . Suppose we translate δ into a quantity of shares q using the arrival price, so that up to roundoff errors, $q = \delta/p_0$. If we assume that the order *will be fully executed* then we can algebraically manipulate the definition

(9.1) to express it in terms of price return and order value:

$$\begin{aligned} \text{slip}(\mathfrak{o}) &= \delta \cdot R_S(\mathfrak{o}) \quad \text{where} \\ R_S(\mathfrak{o}) &:= \frac{\text{vwap}(\mathfrak{o}) - p_0}{p_0} \end{aligned} \quad (9.6)$$

The quantity $R_S(\mathfrak{o})$ will be referred to by me as *slippage price return* – I do not know if this is standard terminology. The quantity $R_S(\mathfrak{o})$ is defined in such a way that the dollar slippage equals return on (signed) dollars traded, using this number as the return.

When a parent order (or metaorder) is finished executing, if you’ve had significant impact on the price, then typically that impact will revert somewhat once the price pressure that you were creating has been removed.

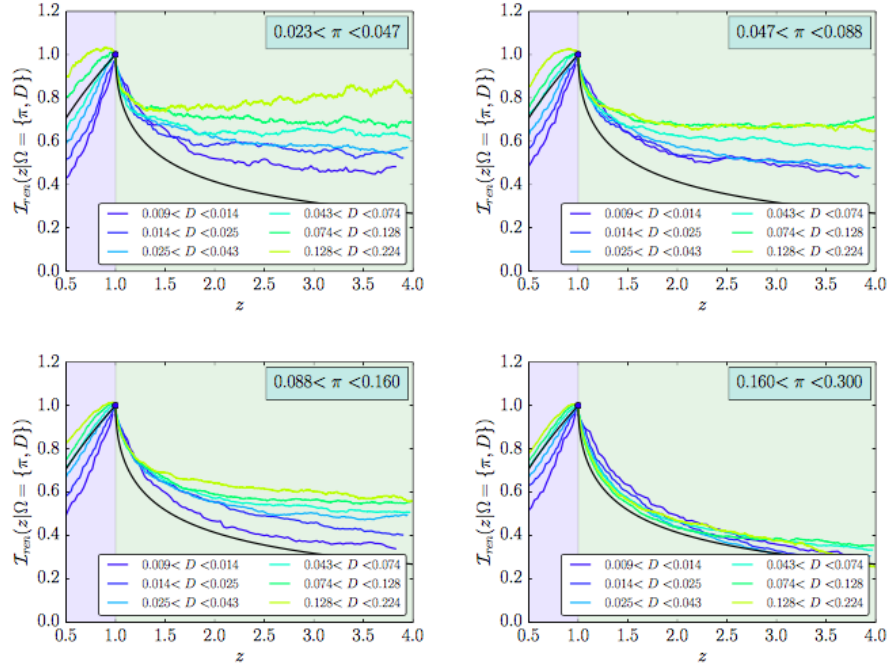


FIGURE 9.3. From Zarinelli et al. (2015). Decay of temporary market impact after the execution of a metaorder.

Within each panel the solid lines correspond to the average market impact trajectory for metaorders with different durations D ; the four panels correspond to different participation rates. The black line corresponds to the prediction of the transient impact model with $\delta = 0.5$.

9.4. Optimal Execution in the Almgren–Chriss Model. Suppose we hold a block of X units of a security that we want to completely liquidate before time

T . We divide T into N intervals of length $\tau = T/N$, and define the discrete times $t_k = k\tau$, for $k = 0, \dots, N$. We define a trading trajectory to be a list x_0, \dots, x_N , where x_k is the number of units that we plan to hold at time t_k . The boundary conditions are:

$$x_0 = X, \text{ and } x_N = 0.$$

There are certain simplifications associated to assuming all of the trades are in the same direction, which is always the case when considering executing a fixed order and without short-term alphas. We treat the case of selling shares; however, the order could just as well be a buy – essentially the same reasoning applies with only minor modifications. Let

$$n_k = x_{k-1} - x_k \geq 0$$

be the number of units that we will *sell* between times t_{k-1} and t_k . Clearly, x_k and n_k are related by

$$x_k = X - \sum_{j=1}^k n_j = \sum_{j=k+1}^N n_j$$

Almgren and Chriss (2001) define a *trading strategy* to be a rule for determining n_k in terms of information available at time t_{k-1} . Broadly speaking we distinguish two types of trading strategies: dynamic and static. Static strategies are determined in advance of trading, that is the rule for determining each n_k depends only on information available at time t_0 . Dynamic strategies, conversely, depend on all information up to and including time t_{k-1} .

We distinguish two kinds of market impact. Temporary impact refers to temporary imbalances in supply in demand caused by our trading leading to temporary price movements away from equilibrium. Permanent impact means changes in the equilibrium price due to our trading, which remain at least for the life of our liquidation. In other words, the definitions are such that temporary impact is assumed to revert instantaneously – it is completely undetectable even one period later. In practice, one of the largest determinants of temporary impact is simply spread pay.

The (midpoint) price dynamics are taken to be an arithmetic random walk:

$$S_k = S_{k-1} + \sigma\tau^{1/2}\xi_k - \tau g(n_k/\tau), \quad (9.7)$$

Here σ represents the volatility of the asset, $\xi_k \sim N(0, 1)$ are i.i.d. normal, and the permanent impact is a function of the average rate of trading $v_k = n_k/\tau$ during the interval. Note that the innovation in (9.7) is an i.i.d. random term and a deterministic function of the execution, so we can telescope the sum all the way back:

$$S_k = S_0 + \sum_{k=1}^N [\sigma\tau^{1/2}\xi_k - \tau g(n_k/\tau)]. \quad (9.8)$$

In this model, we do not explicitly model the bid and the ask as separate processes, but the model does allow for costs associated to the spread as we shall see.

Temporary impact is modeled by assuming the actual price received on the k -th transaction is

$$\tilde{S}_k = S_{k-1} - h(v_k), \quad v_k = n_k/\tau$$

but the effect of $h(v)$ does not appear in S_k . Thus temporary impact literally just means we get a worse price (than the midpoint price) on each of our child fills.

The full trading revenue upon completion of all trades is:

$$\sum_{k=1}^N n_k \tilde{S}_k = \sum_{k=1}^N n_k [S_{k-1} - h(v_k)] = \sum_{k=1}^N n_k S_{k-1} - \sum_{k=1}^N n_k h(v_k)$$

The $n_k S_{k-1}$ term needs further simplification. Note that

$$\begin{aligned} \sum_{k=1}^N S_{k-1} n_k &= \sum_{k=1}^N S_{k-1} (x_{k-1} - x_k) = \sum_{k=0}^{N-1} S_k x_k - \sum_{k=1}^N S_{k-1} x_k \\ &= S_0 X + \sum_{k=1}^N (S_k - S_{k-1}) x_k \end{aligned}$$

where in the last line we used the boundary conditions $x_0 = X, x_N = 0$.

Hence the full trading revenue is

$$\begin{aligned} \sum_{k=1}^N n_k \tilde{S}_k &= S_0 X + \sum_{k=1}^N (S_k - S_{k-1}) x_k - \sum_{k=1}^N n_k h(v_k) \\ &= S_0 X + \sum_{k=1}^N [\sigma \tau^{1/2} \xi_k - \tau g(n_k/\tau)] x_k - \sum_{k=1}^N n_k h(v_k) \end{aligned}$$

where we have used eq. (9.7) to obtain an expression for the increment $S_k - S_{k-1}$.

The total cost of trading is the difference $X S_0 - \sum_k n_k \tilde{S}_k$ between the initial book value and the revenue. This is the standard ex-post measure of transaction costs used in performance evaluations, and is also called *implementation shortfall* or *slippage*.

Prior to trading, implementation shortfall is a random variable. We write $E(x)$ for the expected shortfall and $V(x)$ for the variance of the shortfall. These are calculated as

$$E(x) = \sum_{k=1}^N \tau g(n_k/\tau) x_k + \sum_{k=1}^N n_k h(v_k) \quad (9.9)$$

$$V(x) = \sigma^2 \sum_{k=1}^N \tau x_k^2 \quad (9.10)$$

For each value of $\kappa > 0$ there corresponds a unique trading trajectory x such that

$$E(x) + \frac{\kappa}{2} V(x)$$

is minimal. As we know, this trajectory is optimal from the point of view of an investor with Arrow (1971)–Pratt (1964) constant absolute risk aversion parameter $\kappa > 0$. Note that the sign is flipped; usually we would maximize $E[\pi] - (\kappa/2)V[\pi]$ where π is profit. In this case slippage (or shortfall) is like the negative of profit (subject to the boundary conditions).

Computing optimal trajectories is significantly easier if we take the permanent and temporary impact functions to be linear in the rate of trading. For linear permanent impact,

$$g(v) = \gamma v$$

Eq. (9.8) then yields

$$S_k = S_0 + \sigma \sum_{j=1}^k \tau^{1/2} \xi_j - \gamma(X - x_k)$$

Lemma 9.2. With $g(v) = \gamma v$, the permanent impact term from (9.9) equals

$$\sum_{k=1}^N \tau g(n_k/\tau) x_k = \frac{1}{2} \gamma X^2 - \frac{1}{2} \gamma \sum_k n_k^2$$

Proof. The left side is

$$\gamma \sum_{k=1}^N n_k x_k = \gamma \sum_{k=1}^N n_k \left(\sum_{j=k+1}^N n_j \right) = \gamma \sum_{j>k} n_k n_j$$

Noting that $X = \sum_k n_k$, the right side is

$$\frac{1}{2} \gamma \left[\left(\sum_k n_k \right)^2 - \sum_k n_k^2 \right]$$

The desired result now follows easily. \square

Similarly, for the temporary impact we take

$$h(v) = \epsilon \operatorname{sgn}(v) + \eta v.$$

The units of ϵ are \$/share, and those of η are (\$/share)/(share/time). A reasonable estimate for ϵ is the fixed costs of selling, such as half the bid-ask spread plus fees. It is more difficult to estimate η since it depends on internal and transient aspects of the market microstructure. The linear model above is often called a *quadratic* cost model because the total cost incurred by buying or selling n units in a single unit of time is

$$nh(n/\tau) = \epsilon |n| + (\eta/\tau) n^2$$

With both linear cost models,

$$E(x) = \frac{1}{2}\gamma X^2 + \epsilon \sum_{k=1}^N |n_k| + \frac{\tilde{\eta}}{\tau} \sum_{k=1}^N n_k^2, \quad \tilde{\eta} = \eta - \frac{1}{2}\gamma\tau$$

If all $n_k \geq 0$ as we have assumed, then $\sum_k |n_k| = X$ which means the L^1 term is irrelevant for optimization purposes.

In reality, it's naive to assume the bid-ask spread will be constant over the entire execution path (and one must be wary of simply using formulas from papers without questioning all of the various assumptions that are being made). Sophisticated practitioners would surely use a version of this in which ϵ depends on the time of day.

We need to minimize $U(x) = E(x) + (\kappa/2)V(x)$, which we do by enforcing the first-order condition:

$$\frac{\partial}{\partial x_j}(E + \frac{\kappa}{2}V) = 0 \quad (9.11)$$

Keeping only the relevant terms from E , and using $n_k = x_{k-1} - x_k$ we have

$$\begin{aligned} \frac{\partial}{\partial x_j}(E + \frac{\kappa}{2}V) &= \frac{\partial}{\partial x_j} \left[\frac{\tilde{\eta}}{\tau} \sum_k (x_{k-1} - x_k)^2 + \frac{1}{2} \tau \kappa \sigma^2 \sum_k x_k^2 \right] \\ &= \frac{\tilde{\eta}}{\tau} \frac{\partial}{\partial x_j} [(x_{j-1} - x_j)^2 + (x_j - x_{j+1})^2] + \tau \kappa \sigma^2 x_j \end{aligned}$$

Setting this equal to zero (and dividing by -2) leads to

$$\frac{\tilde{\eta}}{\tau} [(x_{j-1} - x_j) - (x_j - x_{j+1})] = \frac{1}{2} \tau \kappa \sigma^2 x_j$$

Thus we are led to a linear difference equation

$$\frac{1}{\tau^2} (x_{j-1} - 2x_j + x_{j+1}) = \tilde{\psi}^2 x_j, \quad \tilde{\psi}^2 := \frac{\kappa \sigma^2}{2\tilde{\eta}} \quad (9.12)$$

Solutions of such equations may be written as a combination of exponentials. There exists a constant ψ defined as the solution to

$$\frac{2}{\tau^2} (\cosh(\psi\tau) - 1) = \tilde{\psi}^2.$$

In terms of this, the precise solution to (9.12) which respects the boundary conditions is

$$x_j = \frac{\sinh(\psi(T - t_j))}{\sinh(\psi T)} X \quad \text{where} \quad \frac{2}{\tau^2} (\cosh(\psi\tau) - 1) = \tilde{\psi}^2. \quad (9.13)$$

Note that as $\tau \rightarrow 0$ we have $\tilde{\eta} \rightarrow \eta$ and $\tilde{\psi} \rightarrow \psi$.

Some comments and discussion are in order. First note we can re-write (9.12) as

$$x_{j+1} = (2 + \tau^2 \tilde{\psi}^2) x_j - x_{j-1},$$

thus representing it in recursive form. You could then, for instance, calculate x_2 from x_1 and x_0 , x_3 from x_2 and x_1 , etc, but to get the iteration started you need to set the values of the “free parameters” x_0 and x_1 using the boundary conditions, which are $x_0 = X$ and $x_N = 0$. Since the initial boundary condition determines x_0 , we are left solving for x_1 such that $x_N = 0$. This is certainly, but a little annoying and potentially prone to floating-point errors. For this reason, (9.13) (and more generally, continuous-time solutions) are to be preferred when available.

Equation (9.12) resembles the simplest non-trivial second-order differential equation,

$$x''(t) = c x(t)$$

and it is natural to wonder whether one might not have gotten to this equation directly, without having to discretize time. This is, indeed, possible and leads to an illuminating parallel with the Lagrangian formulation of classical mechanics. One ends up representing the trading path as a twice-differentiable function $x(t)$ from the outset, and applying a method for minimizing “functionals” $\mathcal{F}[x(t)]$ which are functions of the entire path. This technique is known as the “calculus of variations” and is treated in my paper with Jerome: <https://ssrn.com/abstract=3057570>.

9.5. Single-period optimization with costs. Let $h_0 \in \mathbb{R}^n$ denote our current portfolio holdings in dollars and let $h \in \mathbb{R}^n$ denote a hypothetical set of portfolio holdings that we could trade into. Let $\pi(h)$ denote the random variable which represents our P/L if we do these trades, and then liquidate the portfolio h so that we hold all cash. As discussed previously (9.5),

$$\pi(h) = h'R - [\text{slip}(h_0, h) + \text{liqslip}(h)]$$

We want to optimize $\mathbb{E}[u(w_T)]$ where T is the final time (eg. after the liquidation has been completed). We assume asset returns follow a multivariate elliptical distribution, and hence there exists some $\kappa > 0$ such that we can equivalently maximize

$$\mathbb{E}[w_T] - \frac{\kappa}{2} \mathbb{V}[w_T]$$

The final wealth is $w_T = w_0 + \pi(h)$ and at this point, dependence on the initial wealth w_0 drops out of the problem (it adds a constant term).

We are thus left with the problem

$$\max_h \left\{ \mathbb{E}[\pi(h)] - \frac{\kappa}{2} \mathbb{V}[\pi(h)] \right\}$$

Suppose we assume that $\mathbb{V}[\pi(h)]$ is well approximated by the variance of the term $h'R$ in $\pi(h)$, in other words

$$\mathbb{V}[\pi(h)] \approx h' \Sigma h \text{ where } \Sigma := \text{cov}(R).$$

If we aren't planning to liquidate the portfolio h after the next period, and we are happy to identify unrealized (ie. paper, mark-to-market) P/L with other P/L realized as cash, then we can ignore the liquidation slippage term.

Then the mean-variance problem becomes

$$\max_h \left\{ h' \mathbb{E}[R] - \frac{\kappa}{2} h' \Sigma h - \mathbb{E}[\text{slip}(h_0, h)] \right\} \quad (9.14)$$

Note that (9.14) can equivalently be considered as an optimization over the *trade list*

$$\delta := h - h_0 \in \mathbb{R}^n$$

since h_0 is fixed, and is not a parameter in the optimization. The trade list δ is the more natural variable for the largest sources of cost; slippage $\text{slip}(h_0, h)$ is a function of δ , as are commissions.

Financing costs, on the other hand, are functions of h rather than δ . You pay to finance the portfolio you end up with. On the short side, financing costs are usually called *borrow costs*, and can be quite high for stocks that are hard to borrow. On the long side, you will pay to finance a portfolio that is larger (in notional terms) than your capital, usually at the rate of 25bps per year times financed notional value.

Note that there are (at least) two functional forms for the latter two terms in (9.14) which allow for easy solution of the mean-variance maximization problem: (1) purely quadratic, and (2) quadratic plus absolute-value type penalty terms. In the first case, the entire problem remains quadratic, while in the second case, the problem becomes mathematically equivalent to a Lasso regression. The Almgren et al. (2005) form does not lead to such a well-known procedure as Lasso, but the associated problem is convex and differentiable, hence standard optimization routines can be expected to perform well.

9.6. General Multiperiod Problems. Gârleanu and Pedersen (2013) studied the multiperiod quantitative-trading problem under the somewhat restrictive assumptions that the alpha models follow mean-reverting dynamics and that the only source of trading frictions are purely linear market impacts (leading to purely quadratic impact-related trading costs). We are going to do something similar, but not so restrictive and general enough to apply to real trading scenarios.

We now place ourselves into the position of a rational agent planning a sequence of trades beginning presently and extending into the future. Specifically, a *trading plan* for the agent is modeled as a specific portfolio sequence

$$\mathbf{h} = (h_1, h_2, \dots, h_T),$$

where h_t is the portfolio the agent plans to hold at time t in the future. If r_{t+1} is the vector of asset returns over $[t, t+1]$, then the trading profit (ie. difference between initial and final wealth) associated to the trading plan \mathbf{h} is given by

$$\pi(\mathbf{h}) = \sum_t [h_t \cdot r_{t+1} - c_t(h_{t-1}, h_t)] \quad (9.15)$$

where $c_t(h_{t-1}, h_t)$ is the total cost (including but not limited to market impact, spread pay, borrow costs, ticket charges, financing, etc.) associated with holding portfolio h_{t-1} at time $t-1$ and ending up with h_t at time t .

Trading profit $\pi(\mathbf{h})$ is a random variable, since many of its components are future quantities unknowable at time $t=0$. Thus the problem we treat initially is that of maximizing $u(\mathbf{h})$, where

$$u(\mathbf{h}) := \mathbb{E}[\pi(\mathbf{h})] - (\kappa/2)\mathbb{V}[\pi(\mathbf{h})] \quad (9.16)$$

Our task is to find the maximum-utility path $\mathbf{h}^* = \operatorname{argmax}_{\mathbf{h}} u(\mathbf{h})$.

Since they will come up over and over again, let us define the shorthand notations

$$\alpha_t := \mathbb{E}[r_{t+1}] \quad \text{and} \quad \Sigma_t := \mathbf{V}[r_{t+1}], \quad y_t := (\kappa \Sigma_t)^{-1} \alpha_t. \quad (9.17)$$

Then combining (9.15) with (9.16), one has

$$u(\mathbf{h}) = \sum_t \left[h'_t \alpha_t - \frac{\kappa}{2} h'_t \Sigma_t h_t - c_t(h_{t-1}, h_t) \right] \quad (9.18)$$

Note that any symmetric, positive-definite matrix Q defines a bilinear form

$$b_Q(x, y) = N_Q(x - y).$$

Lemma 9.3. One has

$$-b_{\kappa \Sigma_t}(h_t, y_t) = O(y^2) + h'_t \alpha_t - \frac{1}{2} h'_t (\kappa \Sigma_t) h_t$$

where $O(y^2)$ denotes a collection of terms which is a quadratic function of y_t and which doesn't contain h_t .

The proof is left as an exercise.

Therefore the first two terms in the utility calculation (9.18) (ie. all the terms not dealing with costs) are given by

$$-b_{\kappa \Sigma_t}(h_t, y_t)$$

Then up to \mathbf{h} -independent terms,

$$u(\mathbf{h}) = - \sum_t [b_{\kappa \Sigma_t}(h_t, y_t) + c_t(h_{t-1}, h_t)] \quad (9.19)$$

Note that maximizing (9.19) is naturally a tracking problem. We are finding the sequence h_t that minimizes tracking error $b_{\kappa\Sigma_t}(h_t, y_t)$ but also minimizes cost $c_t(h_{t-1}, h_t)$.

If you had written a computer program that maximizes (9.19), you could also apply it to other cases where y_t was something other than (9.17). For example, applying it to the sequence $y_t = 0$, and adding the appropriate constraint, we recover the Almgren-Chriss problem. For hedging exposure to derivatives, y_t should be our expectation of the offsetting replicating portfolio at all future times until expiration.

Tracking the portfolios of Black and Litterman (1992) is also a special case of our framework in which y_t is the solution to a mean-variance problem with a Bayesian posterior distribution for the expected returns. Since the posterior is Gaussian in the original Black-Litterman model, the two-moment approximation to utility is exact, and one simply replaces α_t and Σ_t with the appropriate quantities.

We are now starting to think that a computer program that maximizes (9.19) would be pretty useful. Next we describe how to write such a program, in practical terms. It's surprisingly easy, and you'll do it on your homework.

9.7. Non-differentiable Optimization. Given a convex, differentiable map $f : \mathbb{R}^n \rightarrow \mathbb{R}$, if we are at a point x such that $f(x)$ is minimized along each coordinate axis, have we found a global minimizer? In other words, does

$$f(x + d \cdot e_i) \geq f(x) \quad \text{for all } d, i$$

imply that $f(x) = \min_z f(z)$? Here $e_i = (0, \dots, 1, \dots, 0) \in \mathbb{R}^n$, the i -th standard basis vector.

The answer is: Yes!

Now consider the same question, but without the differentiability assumption.

The answer changes to no, and Fig. 9.4 below gives a counterexample.

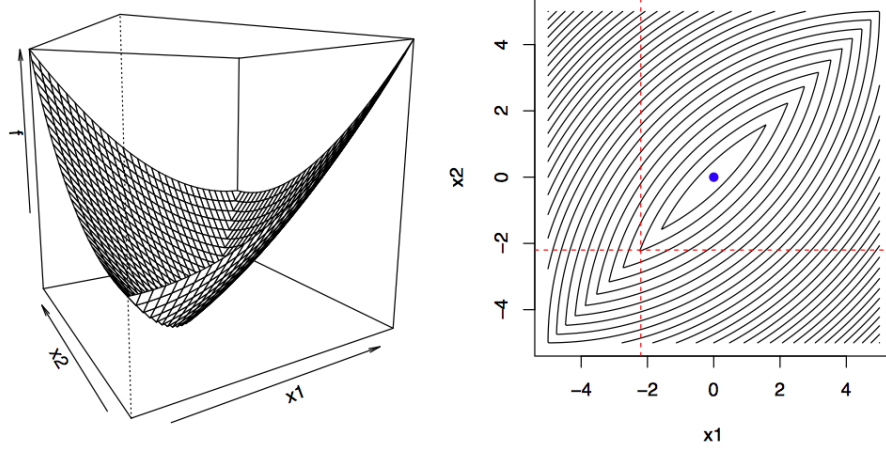


FIGURE 9.4. A convex function for which coordinate descent will get “stuck” before finding the global minimum.

Consider the same question again: “if we are at a point x such that $f(x)$ is minimized along each coordinate axis, have we found a global minimizer?” only now

$$f(x) = g(x) + \sum_{i=1}^n h_i(x_i)$$

with g convex, differentiable and each h_i convex? (In this case, we say the non-differentiable part is *separable*.)

If the non-differentiable term is separable, the answer is yes once again. This is a special case of a deep general result proved by Tseng (2001), which we will call “Tseng’s theorem.” The main take-away is: we can easily optimize

$$f(x) = g(x) + \sum_{i=1}^n h_i(x_i)$$

with g convex, differentiable and each h_i convex, by coordinate-wise optimization.

Tseng’s results also suggest an algorithm, called *blockwise coordinate descent* (BCD).

Algorithm 9.1. Chose an initial guess for x . Repeatedly iterate cyclically through $i = 1, \dots, N$, and perform the following optimization and update:

$$x_i = \underset{\omega}{\operatorname{argmin}} f(x_1, \dots, x_{i-1}, \omega, x_{i+1}, \dots, x_N)$$

Tseng (2001) shows that for functions of the form above, any limit point of the BCD iteration is a minimizer of f . The order of cycling through coordinates is arbitrary; and we can use any scheme that visits each of $\{1, 2, \dots, n\}$ every M

steps for fixed constant M . We can also everywhere replace individual coordinates with blocks of coordinates.

9.8. Single-asset trading paths. Now let us consider the multiperiod problem for a single asset, in which case the ideal sequence $\mathbf{y} = (y_t)$ and the holdings (or equivalently, hidden states) $\mathbf{h} = (h_t)$ are both univariate time series. Since the multiperiod many-asset problem can be reduced to iteratively solving a sequence of single-asset problems, the methods we develop in this section are important even if our main interest is in multi-asset portfolios.

A very important class of examples arises when there are no constraints, but the cost function is a convex and non-differentiable function of the difference

$$\delta_t := h_t - h_{t-1}.$$

This allows for non-quadratic terms as in Almgren et al. (2005) and non-differentiable terms such as linear proportional costs.

In this case, we can use Tseng's theorem, applied to *trades* rather than *positions*. Writing

$$h_t = h_0 + \sum_{s=1}^t \delta_s,$$

the objective function becomes

$$u(\mathbf{h}) = - \sum_t \left[b\left(h_0 + \sum_{s=1}^t \delta_s, y_t\right) + c_t(\delta_t) \right] \quad (9.20)$$

Eq. (9.20) satisfies the convergence criteria of Tseng (2001) that the non-differentiable term is separable *across time*, while the non-separable term is differentiable. One then performs coordinate descent over the trades $\delta_1, \delta_2, \dots, \delta_T$. Almost any reasonable starting point will do to initialize the iteration, but if a warm start from a previous optimization is available, that may speed things along. This approach was introduced in Kolm and Ritter (2015).

Problem 9.1. Prove Lemma 9.3

Problem 9.2. Consider optimally trading a single stock over $T = 30$ days. Each period is one day, and you can trade once per day. The stock's daily return volatility is σ . Suppose your forecast is 50 basis points for the first period, and decays exponentially with half-life 5 days. This means that

$$\alpha_t := \mathbb{E}[r_{t,t+1}] = 50 \times 10^{-4} \times 2^{-t/5}. \quad (9.21)$$

Let $c(\delta)$ be the cost, in dollars, of trading δ dollars of this stock. For selling, $\delta < 0$. Following Almgren, we assume that

$$c(\delta) = PX \left(\frac{\gamma\sigma}{2} \frac{X}{V} \left(\frac{\Theta}{V} \right)^{1/4} + \text{sign}(X) \eta\sigma \left| \frac{X}{V} \right|^\beta \right), \quad X = \delta/P$$

where P is the current price in dollars, X is the signed trade size in shares, V is the daily volume in shares, Θ is the total number of shares outstanding, and finally $\gamma = 0.314$ and $\eta = 0.142$ and $\beta = 0.6$ are constants fit to market data. For concreteness, suppose the asset we are trading has

$$P = \$40, \quad V = 2 \times 10^6, \quad \Theta = 2 \times 10^8, \quad \sigma = 0.02.$$

For a trading path $\mathbf{x} = (x_0, x_1, \dots, x_T)$ where x_t denotes dollar holdings of the stock at time t , define the profit (also in dollars) as

$$\pi(\mathbf{x}) = \sum_{t=1}^t [x_t r_{t,t+1} - c(x_t - x_{t-1})]$$

This is a random variable due to the presence of $r_{t,t+1}$ which you can assume is Gaussian with mean $\alpha_t := \mathbb{E}[r_{t,t+1}]$ and variance $\mathbb{V}[r_{t,t+1}] = \sigma^2$. In this problem, always assume $x_0 = 0$ is fixed.

- (a) Find the sequence of positions x_1, x_2, \dots, x_T that maximizes

$$u(x_1, \dots, x_T) = \sum_{t=1}^T \left[x_t \alpha_t - \frac{\kappa}{2} \sigma^2 x_t^2 - c(x_t - x_{t-1}) \right] \quad (9.22)$$

with risk-aversion $\kappa = 10^{-7}$. Set tolerance so that your algorithm does not terminate unless each $x_t \in \mathbb{R}$ is within a distance of one dollar to the true optimal path. Plot the optimal path $\mathbf{x}^* := (x_0 = 0, x_1^*, \dots, x_T^*)$ and also report its values in a table. Also report the computation time.

Submit your code and a clear explanation of the algorithm you used, why you chose it over other possible algorithms, and how you know that it converges. For example, if you used a method that requires convexity, explain why the function you are optimizing is convex.

- (b) Use the program you wrote in part (a) to plot expected profit of the optimal path, $\mathbb{E}[\pi(\mathbf{x}^*)]$ and *ex ante* Sharpe ratio of the optimal path, defined as

$$\text{Sharpe}(\mathbf{x}^*) = \sqrt{252} \frac{\mathbb{E}[\pi(\mathbf{x}^*)]}{\sqrt{\mathbb{V}[\pi(\mathbf{x}^*)]}}$$

as a function of κ , as a function of the half-life (which was taken to be 5 in equation (9.21) above), as a function of the initial strength (taken to be 50 in equation (9.21)), and as a function of σ . So you need to do eight plots in all: profit and Sharpe ratio, each as a function of one of

four parameters (holding the others fixed). Choose appropriate intervals around the parameter values in part (a). Note that κ cannot be negative in reasonable models.

REFERENCES

- Almgren, Robert and Neil Chriss (1999). “Value under liquidation”. In: *Risk* 12.12, pp. 61–63.
- (2001). “Optimal execution of portfolio transactions”. In: *Journal of Risk* 3, pp. 5–40.
- Almgren, Robert et al. (2005). “Direct estimation of equity market impact”. In: *Risk* 18.7, pp. 58–62.
- Arrow, Kenneth J (1971). “Essays in the theory of risk-bearing”. In:
- Black, Fischer and Robert Litterman (1992). “Global portfolio optimization”. In: *Financial Analysts Journal*, pp. 28–43.
- Foucault, Thierry, Ohad Kadan, and Eugene Kandel (2005). “Limit order book as a market for liquidity”. In: *The review of financial studies* 18.4, pp. 1171–1217.
- Gârleanu, Nicolae and Lasse Heje Pedersen (2013). “Dynamic trading with predictable returns and transaction costs”. In: *The Journal of Finance* 68.6, pp. 2309–2340.
- Gould, Martin D et al. (2013). “Limit order books”. In: *Quantitative Finance* 13.11, pp. 1709–1742.
- Kolm, Petter N and Gordon Ritter (2015). “Multiperiod Portfolio Selection and Bayesian Dynamic Models”. In: *Risk* March 2015.
- Pratt, John W (1964). “Risk aversion in the small and in the large”. In: *Econometrica: Journal of the Econometric Society*, pp. 122–136.
- Tseng, Paul (2001). “Convergence of a block coordinate descent method for non-differentiable minimization”. In: *Journal of optimization theory and applications* 109.3, pp. 475–494.
- Zarinelli, Elia et al. (2015). “Beyond the square root: Evidence for logarithmic dependence of market impact on size and participation rate”. In: *Market Microstructure and Liquidity* 1.02, p. 1550004.