
Blockhouse Work Trial Test

Amos Anderson

28 July, 2025

1 PROBLEM 1

We want to model the temporary impact function $g_t(X)$ which measures the slippage per share for a buy market order of X shares at time t , based on the limit order book data in the provided files for the 3 stocks namely **SOUN**, **FROG** and **CRWW**.

1.1 SET-UP

The key variables for this problem are defined as follows:

- $S \in R$ is the total number of shares to buy by the end of each trading day, specific to each stock.
- N is the total trading periods within the trading day. We are given that $N = 390$
- $t = 1, \dots, N$ is the time index for each trading period.
- $x_t = [x_1, x_2, \dots, x_N]$ is the vector of allocated shares of a stock to buy at time t using market order.
- $i = 00, 01, \dots, 09$ represents the 10 levels of orders in the limit order book. In simpler notation, we shall sometimes use $i = 0, 1, \dots, 9$.

From the ask-side of the limit order book of each stock, we define

$$\begin{aligned} \text{Prices: } p_{t,i} &= \text{ask_px_i}_t, \quad \text{for } i = 0, 1, \dots, 9 \quad \text{such that } p_{t,1} \leq p_{t,2} \leq \dots \leq p_{t,9} \\ \text{Quantities: } q_{t,i} &= \text{ask_sz_i}_t \\ \text{Mid prices: } m_t &= \frac{\text{ask_px_00}_t + \text{bid_px_00}_t}{2} \end{aligned}$$

1.2 DERIVING A DEFINITION FOR $g_t(X)$

Let $(p_{t,0}, q_{t,0})$ be the best ask price and quantity pair. Then for a buy market order of X shares,

- If $X \leq q_{t,0}$, then all X shares are bought at $p_{t,0}$ for a cost of $X \cdot p_{t,0}$.
- If $X > q_{t,0}$, then we deliver all $q_{t,0}$ shares at $p_{t,0}$, and then move to level 1 and take up to $q_{t,1}$ at $p_{t,1}$ and if there is still a remaining unsettled stocks, we move up to level 2 and so on.

In general, suppose k is the smallest index such that the cumulative quantity available up to level k is at least X . That is, k is such that

$$\sum_{i=0}^k q_{t,i} \geq X$$

Then for any level k , we take all shares available at all levels before k , that is, $i = 0, \dots, k-1$ which will cost $q_{t,i} \cdot p_{t,i}$ at each level i . At level k , we take only the remaining shares needed which is given by $X - \sum_{i=0}^{k-1} q_{t,i}$ at the price at level k , $p_{t,k}$. The total cost of this trade is:

$$C_t(X) = \sum_{i=0}^{k-1} q_{t,i} \cdot p_{t,i} + \left(X - \sum_{i=0}^{k-1} q_{t,i} \right) \cdot p_{t,k}$$

The average cost is

$$AC_t(X) = \frac{1}{X} C_t(X)$$

The slippage, or temporary impact, is the difference between the average cost and the mid-price and is given by

$$g_t(X) = AC_t(X) - m_t$$

1.3 MODELING THE TEMPORARY IMPACT $g_t(X)$

To model the temporary impact function $g_t(X)$ using the 1-minute snapshots of the limit order book data for the three stocks (**SOUN**, **FROG**, and **CRWW**), we compute $g_t(X)$ across a grid of orders sizes (from $X = 10$ to $X = 600$). We then fit three candidate models (linear, power law and quadratic) to the empirical $g_t(X)$ curves and evaluate the results based on fit statistics, including Mean Squared Error (MSE), R-squared, Adjusted R-squared, and Akaike Information Criterion (AIC). The analysis justifies our selection of the quadratic model, showing that it is superior to the other

Table 1.1: Fit Statistics for Linear, Power-Law, and Quadratic Models

Model	Metric	SOUN	FROG	CRWW
Linear	MSE	0.00001	0.00142	0.00958
	R-squared	-100.11927	-7.81274	0.14158
	Adjusted R-squared	-100.11927	-7.81274	0.14158
	AIC	-1403.58801	-771.72045	-546.45600
Power Law	MSE	0.00000	0.00016	0.01050
	R-squared	0.79879	0.01550	0.05957
	Adjusted R-squared	0.79705	0.00702	0.05146
	AIC	-2135.51196	-1028.35549	-533.68907
Quadratic	MSE	0.00000	0.00001	0.00035
	R-squared	0.99949	0.94564	0.96869
	Adjusted R-squared	0.99948	0.94469	0.96815
	AIC	-2838.53357	-1368.13391	-933.17931

candidate models in capturing the piecewise convex behavior of $g_t(X)$.

The fit statistics for the three models are presented in Table 1.1, derived from the observed $g_t(X)$ values across the specified order sizes and timestamps.

From Table 1.1, we observe the following:

- The linear model fit given by $g_t(X)_{\text{linear}} = \beta_t x$ is not a good fit for modeling the temporary impact function across SOUN, FROG, and CRWW due to its consistently poor performance across all metrics, as evidenced by negative R-squared values (-100.11927 for SOUN, -7.81274 for FROG) and a marginal 0.14158 for CRWW as well as recording the highest AIC values among all three candidate models. These poor fit statistics indicate that the linear model performs worse than a simple mean predictor and fails to capture the nonlinear relationship between order size X and slippage $g_t(X)$.
- The power law model, given by $g_t(X)_{\text{power}} = \gamma_t x^{\delta_t}$ for $\delta_t > 0$, offers some improvement with positive R-squared values (like 0.79879 for SOUN), but still underperforms with low values (0.01550 for FROG, 0.05957 for CRWW) and gives a relatively higher MSE values (0.01050 for CRWW for instance), suggesting it struggles to adapt to the varying curvature across stocks and timestamps.
- The quadratic model, given by $g_t(X)_{\text{quadratic}} = ax^2 + bx + c$, performs robustly across all fit metrics for all three stocks and outperforms the linear and power law models in all metrics. The quadratic model yields Adjusted R-squared values higher than 0.94 and close to 1 across all three stocks. This shows that the model captures the non-linear dynamics between the order size and the temporary impact. In fact, the structure of the quadratic model makes it easier to capture the convex nature of temporary impact. The ax^2 term acts like a gentle slope that steepens over time, matching the natural ups and downs seen in trading data over different time points. Moreover, the quadratic fit also gives lower prediction errors since its MSE values are closer to zero. In particular, the highest MSE value here which is 0.00035 for CRWW is significantly lower than the linear (0.00958) and power-law (0.01050) models. Moreover, the quadratic model gives the best optimal complexity tradeoff. We see that through its AIC values (being the least among all candidate models for all stocks), outperforming the simpler linear model and the moderately complex power-law. The balance between high Adjusted R-squared and low AIC helps mitigate overfitting risks which partly addresses any overfitting or underfitting issues that may arise.

Based on the above analysis, we recommend the quadratic model as the preferred choice for modeling the temporary market impact function for all three stocks under consideration.

Code Availability

Please use the link below to access all code and notebook file that contains all the code, data processing, and modeling steps applied in this project:

https://github.com/Amos-Anderson/Modeling-Temporary-Impact-Function/blob/main/main_work.ipynb

2 PROBLEM 2

Given our choice model:

$$g_t(x_i) = ax_i^2 + bx_i + c$$

and that a, b, c are known coefficients from fitting the quadratic model to the temporary impacts (in Problem 1), we want to determine the optimal order sizes $\{x_i\}_{i=1}^N$ to execute a total of S shares over N discrete time points $\{t_1, t_2, \dots, t_N\}$ that minimizes the total cost :

$$\sum_{i=1}^N x_i g_t(x_i)$$

such that $\sum_{i=1}^N x_i = S$.

The trading cost becomes

$$\begin{aligned} \text{Total Cost} &= \sum_{i=1}^N x_i (ax_i^2 + bx_i + c) \\ &= \sum_{i=1}^N (ax_i^3 + bx_i^2 + cx_i) \end{aligned}$$

2.1 OPTIMIZATION PROBLEM

We want to solve the following constrained optimization problem:

$$\begin{aligned} \min_{x_1, \dots, x_N} \quad & \sum_{i=1}^N (ax_i^3 + bx_i^2 + cx_i) \\ \text{s.t} \quad & \sum_{i=1}^N x_i = S \end{aligned}$$

We apply constrained convex optimization, specifically the Lagrangian method, to derive the optimality conditions. Since the objective is strictly convex, we know any local minimum is global. Additionally, due to the piecewise constant structure of the solution, we reduce the problem to a discrete optimization over a single parameter n_1 , which we solve using grid search.

2.2 SOLVING USING LAGRANGIAN METHOD

We define the Lagrangian function as follows

$$\mathcal{L}(x_1, \dots, x_N, \lambda) = \sum_{i=1}^N (ax_i^3 + bx_i^2 + cx_i) - \lambda \left(\sum_{i=1}^N x_i - S \right)$$

Taking partial derivatives with respect to x_i and setting to zero, we have

$$\frac{\partial \mathcal{L}}{\partial x_i} = 3ax_i^2 + 2bx_i + c - \lambda = 0 \implies 3ax_i^2 + 2bx_i + c = \lambda$$

We have used the fact that

$$\frac{\partial}{\partial x_k} \sum_{i=1}^T f(x_i) = \frac{\partial}{\partial x_k} f(x_k) = f'(x_k)$$

since all terms in the sum, $f(x_i), i \neq k$ are treated as constants and their derivative is zero.

So we have established that each x_i must satisfy the quadratic equation

$$3ax_i^2 + 2bx_i + c - \lambda = 0 \quad \text{for } i = 1, \dots, N$$

So the optimal x_i must be one of the roots of this quadratic equation.

We let $f(x) = 3ax^2 + 2bx + c - \lambda = 0$ and assume that the optimal Lagrange multiplier λ^* lies within a known bounded interval, and perform a grid search over that interval. Then for each λ in the predetermined interval, the solution to the quadratic equation are the roots of $f(x)$ given by

$$x = \frac{-2b \pm \sqrt{4b^2 - 12a(c - \lambda)}}{6a}$$

This yields two roots: $x^{(1)}$ and $x^{(2)}$. So the optimal x_i must be equal to one of the roots of that quadratic equation and the full vector $\{x_i\}$ must contain a combination of these roots.

Suppose $x^{(1)} \neq x^{(2)}$, then since each x_i must be either $x^{(1)}$ or $x^{(2)}$, we denote:

- n_1 as the number of x_i 's equal to $x^{(1)}$.
- $n_2 = N - n_1$ as the number of x_i 's equal to $x^{(2)}$.

Under this discrete allocation structure, the constraint becomes

$$n_1 x^{(1)} + (N - n_1) x^{(2)} = S$$

and the objective total cost becomes:

$$\sum_{i=1}^N h(x_i) = n_1 h(x^{(1)}) + (N - n_1) h(x^{(2)}) \quad \text{where} \quad h(x_i) = ax_i^3 + bx_i^2 + cx_i$$

Since n_1 is determined by the constraint, we solve for n_1 from the constraint to get

$$n_1 = \frac{S - Nx^{(2)}}{x^{(1)} - x^{(2)}}$$

This yields a closed-form expression for how many of the optimal x_i 's should be set to $x^{(1)}$ and the rest to $x^{(2)}$.

Next, for each λ , after getting n_1 and n_2 , we substitute into the total cost function to find its cost. The best λ with its corresponding pair (n_1, n_2) is the one for which the total cost value is the lowest. This gives the optimal allocation for x_i .

If $x^{(1)} = x^{(2)}$, then the constraint simplifies to :

$$Nx^{(1)} = S$$

and the optimal allocation becomes

$$x^{(1)} = \frac{S}{N}$$

Here is a summary of our proposed algorithm:

2.3 ALGORITHM FOR OPTIMAL TRADE ALLOCATION

1. **Input:** Parameters a, b, c , total shares S , number of time steps N , and a grid of candidate λ values.

2. **For each candidate λ :**

(a) Solve the quadratic equation:

$$3ax^2 + 2bx + c = \lambda$$

to obtain the two real roots $x^{(1)}$ and $x^{(2)}$. If the discriminant is negative, skip this λ .

(b) If $x^{(1)} = x^{(2)}$, set $x_i = x^{(1)} = \frac{S}{N}$ for all i and compute total cost:

$$\text{Cost} = N \cdot g(x^{(1)})$$

(c) Else, compute:

$$n_1 = \frac{S - Nx^{(2)}}{x^{(1)} - x^{(2)}}, \quad n_2 = N - n_1$$

Check if n_1 is in $[0, N]$ and is an integer. Then compute total cost:

$$\text{Cost} = n_1 \cdot g(x^{(1)}) + (N - n_1) \cdot g(x^{(2)})$$

(d) Store the configuration $(\lambda, x^{(1)}, x^{(2)}, n_1, n_2)$ and its corresponding cost.

3. **Select the configuration with the lowest total cost.**

4. **Output:** Optimal values $\{x_i\}$ given by $x^{(1)}$ repeated n_1 times and $x^{(2)}$ repeated n_2 times.