

Reinforcement Learning for Optimal Execution with the Queue Reactive Model

Edoardo Vittori

Based on a joint work with Tomas Espana, Yadh Hafsi and Fabrizio Lillo

March 27, 2026

1. Introduction

2. Market simulation

- Stochastic models
- Agent based models
- Queue reactive model

3. RL for optimal execution

- Intro to RL
- Embedding optimal execution
- Experimental results

4. Conclusions

Problem statement

Assume the trader must buy (sell) V units of a security over $[0, T]$. The order is completed in N trades at times t_0, t_1, \dots, t_{N-1} , with $t_0 = 0$ and $t_{N-1} = T$. Let v_{t_n} denote the trade size at time t_n , then we have: $\sum_{n=0}^{N-1} v_{t_n} = V$. For a buy problem, $V > 0$, and for a sell problem, $V < 0$.

Execution cost

Assume P_0 is the arrival price, and \bar{P}_k is the execution price for trade v_k , then the execution cost is given by:

$$C(\mathbf{v}) = \sum_{k=0}^{N-1} v_k \bar{P}_k - VP_0 = \sum_{k=0}^{N-1} v_k (\bar{P}_k - P_0)$$

This expression is also referred to as implementation shortfall

Objective

$$\arg \min_{\mathbf{v}} \mathbb{E}[C(\mathbf{v})]$$

Price simulation:

$$\bar{P}_k = P_k + \rho v_k + \text{sign}(v_k) \frac{S}{2}$$

$$P_{k+1} = P_k + \theta v_k + \eta_k$$

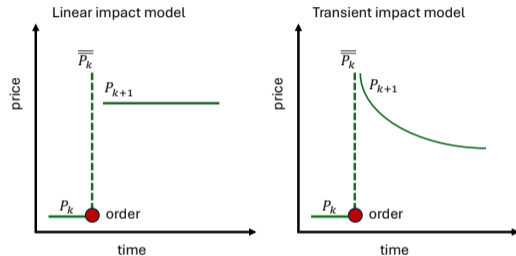
where:

- $\eta_k \sim \text{i.i.d. } \mathcal{N}(0, \sigma^2)$
- θ is the permanent impact coefficient
- ρ is the temporary impact coefficient
- S is the constant bid-ask spread

Considerations

- Further realism can be added by using a transient impact model like in [Obizhaeva and Wang, 2005]
- These models can be calibrated to real data
- **These models only simulate the price, not the limit order book**

Figure 1: Example of linear and transient impact models



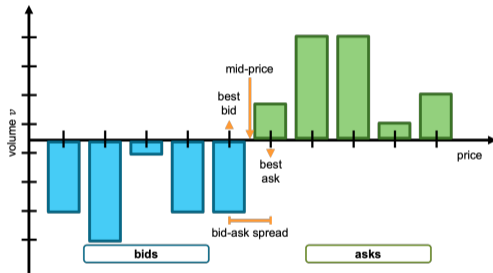
Order types

- **Market order** is an order to execute immediately at the best available price in the order book
- **Limit order** is an order that specifies both the price and volume of a trade
- A limit order sits in the order book until it is either **executed** against a matching market order or **canceled**

Examples of features of the LOB

- Volume imbalance $\frac{v_b - v_a}{v_b + v_a}$
- Volume at best bid and best ask

Figure 2: Illustration of Limit Order Book



ABIDES

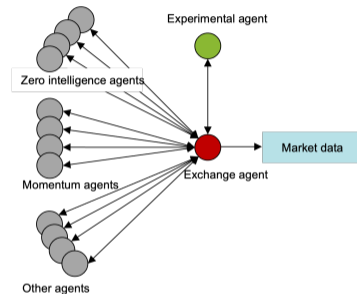
The Agent-Based Interactive Discrete Event Simulation (ABIDES) [Byrd et al., 2019] realistically replicates characteristics of electronic markets such as:

- Continuous double-auction trading
- Network latency and agent computation delays
- Communication solely by means of standardized message protocols

The price process can be described by a **fundamental value** or by using **historical data**. It is possible to create a multi-agent composition using pre-defined agents such as:

- **exchange** agent
- **value** agents
- **momentum** agents
- **noise** agents
- **market maker** agents
- **custom made** agents

Figure 3: Illustration of agent based models



Considerations

- It is not possible to calibrate this simulator to real data
- It is not possible to calibrate a realistic impact model

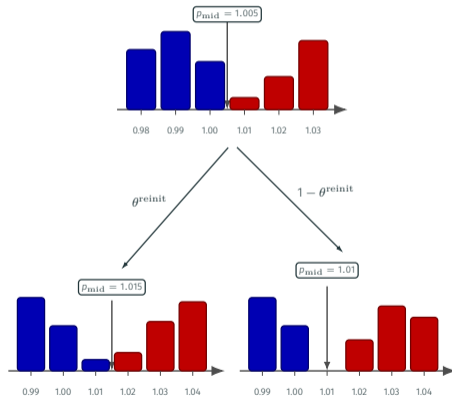
Core Components:

- LOB simulation model for **large tick assets**
- **Queue dynamics** (at fixed price) modeled as a continuous-time Markov process
- State: queue sizes at bid/ask levels
- For queue i :
 - Insertions (limit orders) with intensity $f_i(q)$
 - Removals (due to cancellations or market orders) with intensity $g_i(q)$
 - Queue sizes change by ± 1 at each event
- f_i, g_i calibrated on LOB data

Price Dynamics:

- Mid price updates occur when the best bid or ask queue is depleted
- With the θ and θ^{reinit} parameter you can control the market impact behavior
- θ^{reinit} causes sampling from empirical distributions

Figure 4: QRM response to consuming the best ask q_1 at time t



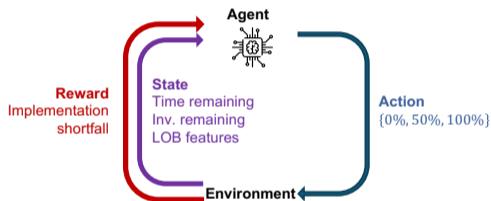
Reinforcement learning basics

- **MDP:** the Markov decision process describes the interaction between agent and environment
- **Objective:** find the policy π which maximizes the discounted sum of the rewards
- $J_\pi = \mathbb{E}_\pi[\sum_t \gamma^t r_t]$ with the reward at time t as r_t

Optimal Execution MDP

- **State:** time remaining, inventory remaining, LOB features
- **Action:** do nothing, market order for 50% or 100% of volume present in first level of LOB:
- **Reward:** $r_t = v_t(P_0 - \bar{P}_t) - \alpha \mathbf{1}_{\{t_k=T\}} V_T$

Figure 5: Illustration of MDP flow



Q-learning

- Q-function

$$Q_{\pi} = \mathbb{E}_{\pi} \left[\sum \gamma^t R_t \mid s_0, a_0 \right]$$

- Bellman Equation

$$Q_{\pi} = r(s, a) + \gamma \mathbb{E}_{s', a'} [Q_{\pi}(s', a')]$$

- Q-learning algorithm

$$Q_t(s, a) = r(s, a) + \gamma \max_{a'} Q_t(s', a')$$

- Q-learning is a tabular algorithm which can be generalized using function approximators

Algorithm examples

- DDQN [Van Hasselt et al., 2016]
- FQI [Ernst et al., 2005]

MDP

- **State:** time remaining, inventory remaining, volume at best ask, volume at best bid and best ask price
- **Action:** do nothing, market order for 50% or 100 % of the volume present in first level of lob
- **Reward:** $r_t = v_t(P_0 - \bar{P}_t) - \alpha \mathbf{1}_{\{t_k=T\}} V_T$

Execution setup

- **Market simulator:** QRM
- **Target:** Buy 25 shares
- **Horizon:** 600 seconds
- **Timestep:** 25-second intervals
- **RL algorithm:** DDQN

Benchmark execution algorithms

- **TWAPd:** Time-weighted average price — execute 1 share at each timestep
- **POPV1, POPV2:** execute 50% of available volume every 1, 2 steps respectively
- **POPV3, POPV4:** execute 100% of available volume every 3, 4 steps respectively

Figure 6: Learning curve

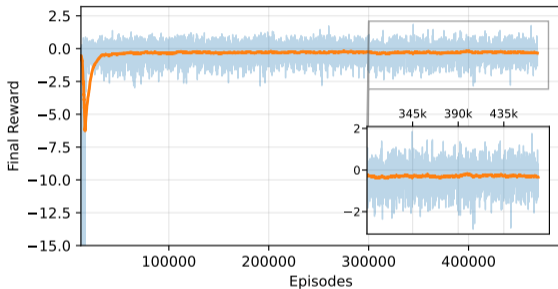


Figure 7: Execution trajectory

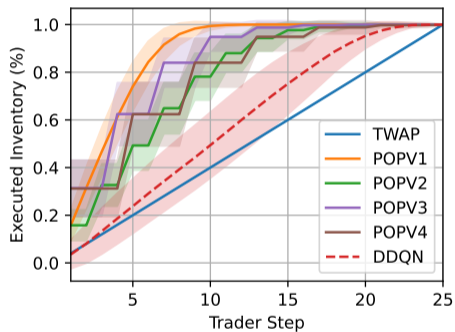


Figure 8: Distribution of the implementation shortfall

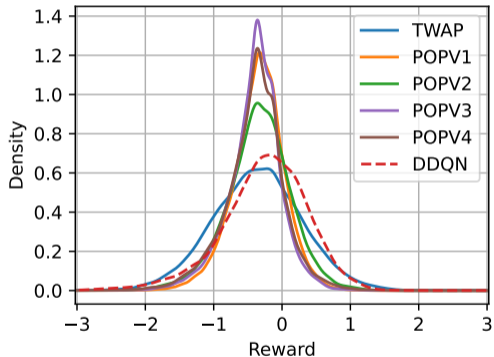
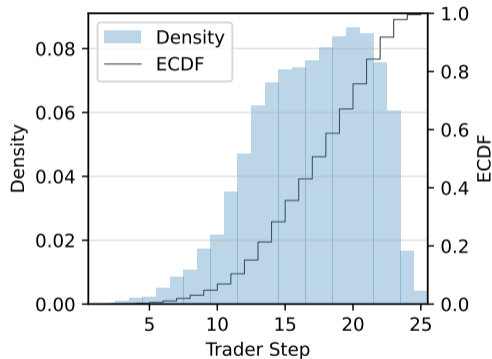


Figure 9: Distribution of the number of execution steps



Takeaways from the work:

- In-depth analysis of the QRM model, verifying that market impact is simulated realistically
- Learned an optimal execution strategy using RL
- Obtained execution strategies with a superior performance with respect to the benchmarks

Contact: edoardo.vittori@intesasanpaolo.com

Figure 10: Download the paper!



Figure 11: Open source code



- [Byrd et al., 2019] Byrd, D., Hybinette, M., and Balch, T. H. (2019).
Abides: Towards high-fidelity market simulation for ai research.
arXiv preprint.
- [Ernst et al., 2005] Ernst, D., Geurts, P., and Wehenkel, L. (2005).
Tree-based batch mode reinforcement learning.
JMLR, 6(Apr):503–556.
- [Huang et al., 2015] Huang, W., Lehalle, C.-A., and Rosenbaum, M. (2015).
Simulating and analyzing order book data: The queue-reactive model.
Journal of the American Statistical Association, 110(509):107–122.
- [Obizhaeva and Wang, 2005] Obizhaeva, A. A. and Wang, J. (2005).
Optimal trading strategy and supply/demand dynamics.
Journal of Financial markets, 16(1):1–32.
- [Van Hasselt et al., 2016] Van Hasselt, H., Guez, A., and Silver, D. (2016).
Deep reinforcement learning with double q-learning.
In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.