

OTC market making

Zhou Fang¹

¹Department of Mathematics, The University of Texas at Austin

April 24, 2023

Abstract

The over-the-counter (OTC) market is known for its unique characteristic of allowing market makers to set different bid-ask spreads based on the size of the order. However, the market-making problems in the OTC market is a challenging high-dimensional stochastic control problems. This paper proposes a stochastic policy approach for setting bid-ask prices and using reinforcement learning to optimize the strategy. Under some stylized assumptions, we demonstrate that the optimal stochastic policy follows a normal distribution.

1 Model

an OTC market-maker will set quotes on different price ladders, and can also choose whether or not to hedge part of its inventory by externalization. Let the dynamics of the underlying asset be

$$\frac{dS_t}{S_t} = \sigma dW_t \quad (1)$$

Modeling the successful deals of size z_k as Poisson processes $N_t^+(k)$, and $N_t^-(k)$, and denote the intensity of those Poisson processes as $\lambda^+(k)$, and $\lambda^-(k)$. Define a function $h(q)$ as follows

$$h(q) = \begin{cases} 0 & q < Q \\ 1 & q \geq Q \end{cases}$$

Let $\epsilon_t = (\epsilon_t^\pm(k))_{k=1}^N$ be the bid-ask spreads posted by the market maker at time t , and let $\pi(\epsilon_t|t, S, q)$ be the probability density for posting spreads ϵ_t . If the market maker posts the bid-ask spreads at time t as ϵ_t , then the inventory has the following dynamics

$$dq_t = \sum_{k=1}^N z_k (dN_t^+(k) - dN_t^-(k)) \quad (2)$$

If the market maker chooses to externalize dq_t inventory, then she needs to pay an additional transaction fee δ . The wealth process is

$$dX_t = \sum_{k=1}^N z_k [\epsilon_t^b(k) dN_t^+(k) + \epsilon_t^a(k) dN_t^-(k)] + d(q_t S_t) - \delta h(q_t) dq_t \quad (3)$$

Value Function

Given a policy π . Let q_t^π be the inventory process under policy π , and the initial condition at time t be $S_t = S$, $q_t^\pi = q$. Then the value function under policy π is

$$\begin{aligned}
& V^\pi(t, S, q) \\
&= \mathbb{E} \left[\int_t^T \int_{\epsilon_u} \left[\sum_{k=1}^N z_k [\epsilon_u^b(k) dN_u^+(k) + \epsilon_u^a(k) dN_u^-(k)] + d(q_u S_u) - \delta h(q_u) dq_t \right] \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u \right. \\
&\quad \left. - \gamma \int_t^T \int_{\epsilon_u} \pi(\epsilon_u | u, S_u, q_u^\pi) \log \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u du \mid S_t = S, q_t^\pi = q \right] \\
&= \mathbb{E} \left[\int_t^T \int_{\epsilon_u} \sum_{k=1}^N \left[z_k (S_u + \epsilon_u^b(k) - \delta h(q_u)) dN_u^+(k) - z_k (S_u - \epsilon_u^a(k) - \delta h(q_u)) dN_u^-(k) \right] \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u \right. \\
&\quad \left. - \gamma \int_t^T \int_{\epsilon_u} \pi(\epsilon_u | u, S_u, q_u^\pi) \log \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u du \mid S_t = S, q_t^\pi = q \right] \tag{4}
\end{aligned}$$

Then the value function under the optimal policy is

$$\begin{aligned}
& V(t, S, q) \\
&= \max_{\pi} \mathbb{E} \left[\int_t^T \int_{\epsilon_u} \sum_{k=1}^N \left[z_k (S_u + \epsilon_u^b(k) - \delta h(q_u)) dN_u^+(k) - z_k (S_u - \epsilon_u^a(k) - \delta h(q_u)) dN_u^-(k) \right] \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u \right. \\
&\quad \left. - \gamma \int_t^T \int_{\epsilon_u} \pi(\epsilon_u | u, S_u, q_u^\pi) \log \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u du \mid S_t = S, q_t^\pi = q \right] \\
&= \max_{\pi} \mathbb{E} \left[\int_t^{t+\Delta t} \int_{\epsilon_u} \sum_{k=1}^N \left[z_k (S_u + \epsilon_u^b(k) - \delta h(q_u)) dN_u^+(k) - z_k (S_u - \epsilon_u^a(k) - \delta h(q_u)) dN_u^-(k) \right] \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u \right. \\
&\quad \left. - \gamma \int_t^{t+\Delta t} \int_{\epsilon_u} \pi(\epsilon_u | u, S_u, q_u^\pi) \log \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u du + V(t + \Delta t, S_t + \Delta S_t, q_t + \Delta q_t^\pi) \mid S_t = S, q_t^\pi = q \right] \\
&= \max_{\pi} \left\{ \int_{\epsilon_t} \sum_{k=1}^N \left[z_k (S_t + \epsilon_t^b(k) - \delta h(q_t)) \lambda_t^+(k) - z_k (S_t - \epsilon_t^a(k) - \delta h(q_t)) dN_t^-(k) \right] \pi(\epsilon_t | t, S_t, q_t^\pi) d\epsilon_t \Delta t \right. \\
&\quad \left. - \gamma \int_{\epsilon_t} \pi(\epsilon_t | t, S_t, q_t^\pi) \log \pi(\epsilon_t | t, S_t, q_t^\pi) d\epsilon_t \Delta t + \mathbb{E} \left[V(t + \Delta t, S_t + \Delta S_t, q_t^\pi + \Delta q_t^\pi) \mid S_t = S, q_t^\pi = q \right] \right\} \tag{5}
\end{aligned}$$

Dynamic Programming

To make notation simpler, denote $\mathcal{L}V(t, S_t, q_t)$ as

$$\mathcal{L}V(t, S_t, q_t) = V(t, S_t, q_t) + (\partial_t V(t, S_t, q_t) + \frac{1}{2} \sigma^2 \partial_{SS} V(t, S_t, q_t) \Delta t + \sigma \partial_S V(t, S_t, q_t) dW_t) \tag{6}$$

Since $dS_t = \sigma S_t dW_t$, and $dq_t = \sum_k z_k (dN_t^+(k) - dN_t^-(k))$, by the Ito formula, we have the following,

$$\begin{aligned}
& V(t + \Delta t, S_t + \Delta S_t, q_t + \Delta q_t) \\
&= V(t + \Delta t, S_t + \Delta S_t, q_t) \prod_k (1 - dN_t^+(k)) (1 - dN_t^-(k)) \\
&\quad + \sum_k \left[V(t + \Delta t, S_t + \Delta S_t, q_t + z_k) dN_t^+(k) + V(t + \Delta t, S_t + \Delta S_t, q_t - z_k) dN_t^-(k) \right] \\
&= \mathcal{L}V(t, S_t, q_t) \prod_k (1 - dN_t^+(k)) (1 - dN_t^-(k)) + \sum_k \mathcal{L}V(t, S_t, q_t + z_k) dN_t^+(k) + \mathcal{L}V(t, S_t, q_t - z_k) dN_t^-(k) \tag{7}
\end{aligned}$$

Notice that the above Ito formula is based on the assumption that the inventory process is $dq_t = \sum_k z_k (dN_t^+(k) - dN_t^-(k))$. Since the intensities of Poisson processes are determined by the quoted bid-ask spreads. So, the inventory process in the above Ito formula assumes the bid-ask spreads are already determined. Thus, when computing conditional expectation, $\mathbb{E}[V(t + \Delta t, S_t + \Delta S_t, q_t^\pi + \Delta q_t^\pi) | S_t = S, q_t^\pi = q]$, one should average over all possibilities. Then the conditional expectation is

$$\begin{aligned} & \mathbb{E} \left[V(t + \Delta t, S_t + \Delta S_t, q_t^\pi + \Delta q_t^\pi) \mid S_t = S, q_t^\pi = q \right] \\ &= V(t, S, q) + \int_{\epsilon_t} \pi(\epsilon_t | t, S, q) \left[- \sum_k (\lambda_t^+(k) + \lambda_t^-(k)) V(t, S, q) + \partial_t V(t, S, q) + \frac{1}{2} \sigma^2 \partial_{SS} V(t, S, q) \right. \\ & \quad \left. + \sum_k [\lambda_t^+(k) V(t, S, q + z_k) + \lambda_t^-(k) V(t, S, q - z_k)] \right] d\epsilon_t \Delta t \end{aligned} \quad (8)$$

Then one can get the HJB equation,

$$\begin{aligned} & \max_{\pi} \left\{ \int_{\epsilon_t} \sum_k [\lambda_t^+(k) V(t, S, q + z_k) + \lambda_t^-(k) V(t, S, q - z_k) - (\lambda_t^+(k) + \lambda_t^-(k)) V(t, S, q)] \pi(\epsilon_t | t, S, q) d\epsilon_t \right. \\ & \quad + \int_{\epsilon_t} \sum_{k=1}^N [z_k \lambda_t^+(k) (S + \epsilon_t^b(k) - \delta h(q)) - z_k \lambda_t^-(k) (S - \epsilon_t^a(k) - \delta h(q))] \pi(\epsilon_t | t, S, q) d\epsilon_t \\ & \quad \left. - \gamma \int_{\epsilon_t} \pi(\epsilon_t | t, S, q) \log \pi(\epsilon_t | t, S, q) d\epsilon_t \right\} + \partial_t V(t, S, q) + \frac{1}{2} \sigma^2 \partial_{SS} V(t, S, q) \\ &= 0 \end{aligned} \quad (9)$$

Optimal Stochastic Policy

In order to find the maximizer for the quantity inside the max bracket of the HJB equation, we apply calculus of variation

$$\begin{aligned} 0 &= \int_{\epsilon_t} \sum_k [\lambda_t^+(k) V(t, S, q + z_k) + \lambda_t^-(k) V(t, S, q - z_k) - (\lambda_t^+(k) + \lambda_t^-(k)) V(t, S, q)] \delta \pi d\epsilon_t \\ & \quad + \int_{\epsilon_t} \sum_{k=1}^N [z_k \lambda_t^+(k) (S + \epsilon_t^b(k) - \delta h(q)) - z_k \lambda_t^-(k) (S - \epsilon_t^a(k) - \delta h(q))] \delta \pi d\epsilon_t \\ & \quad - \gamma \int_{\epsilon_t} \pi \frac{\delta \pi}{\pi} d\epsilon_t - \gamma \int_{\epsilon_t} \delta \pi \log \pi d\epsilon_t \end{aligned} \quad (10)$$

Since π is probability density distribution, then

$$\int_{\epsilon_t} \delta \pi d\epsilon_t = 0 \quad (11)$$

Then equation (10) becomes

$$\begin{aligned} 0 &= \int_{\epsilon_t} \delta \pi \left(\sum_k [\lambda_t^+(k) V(t, S, q + z_k) + \lambda_t^-(k) V(t, S, q - z_k) - (\lambda_t^+(k) + \lambda_t^-(k)) V(t, S, q)] \right. \\ & \quad \left. + z_k \lambda_t^+(k) (S + \epsilon_t^b(k) - \delta h(q)) - z_k \lambda_t^-(k) (S - \epsilon_t^a(k) - \delta h(q)) \right] - \gamma (\delta \pi) \log \pi d\epsilon_t \end{aligned} \quad (12)$$

Then the quantity inside the bracket above is a constant

$$C = \sum_k \left[\lambda_t^+(k) V(t, S, q + z_k) + \lambda_t^-(k) V(t, S, q - z_k) - (\lambda_t^+(k) + \lambda_t^-(k)) V(t, S, q) \right. \\ \left. + z_k \lambda_t^+(k) (S + \epsilon_t^b(k) - \delta h(q)) - z_k \lambda_t^-(k) (S - \epsilon_t^a(k) - \delta h(q)) \right] - \gamma \log \pi \quad (13)$$

We assume the relation between the intensity and spreads is

$$\lambda_t^\pm(k) = A_k - B_k \epsilon_t^{a,b}(k) \quad (14)$$

To simplify the notations, let

$$\mathcal{H}_k^+(t, S, q, \pi) = V^\pi(t, S, q + z_k) - V^\pi(t, S, q) + z_k(S + \delta h(q)) \quad (15)$$

$$\mathcal{H}_k^-(t, S, q, \pi) = V^\pi(t, S, q - z_k) - V^\pi(t, S, q) - z_k(S - \delta h(q)) \quad (16)$$

So, under optimal policy π^* , we have

$$\mathcal{H}_k^+(t, S, q) = V(t, S, q + z_k) - V(t, S, q) + z_k(S + \delta h(q)) \quad (17)$$

$$\mathcal{H}_k^-(t, S, q) = V(t, S, q - z_k) - V(t, S, q) - z_k(S - \delta h(q)) \quad (18)$$

Then the optimal stochastic policy is

$$\begin{aligned} \pi^*(\epsilon_t | t, S, q) &\propto \exp \left\{ \frac{1}{\gamma} \sum_k (A_k - B_k \epsilon_t^{a,b}(k)) (z_k \epsilon_t^{a,b}(k) + \mathcal{H}_k^\pm(t, S, q)) \right\} \\ &\propto \prod_k \exp \left\{ - \frac{z_k B_k}{\gamma} \left[\epsilon_t^{a,b}(k) - \frac{A_k}{2B_k} + \frac{\mathcal{H}_k^\pm(t, S, q)}{2z_k} \right]^2 \right\} \\ &\propto \prod_k \mathcal{N} \left(\epsilon_t^{a,b} \mid \frac{A_k}{2B_k} - \frac{\mathcal{H}_k^\pm(t, S, q)}{2z_k}, \frac{\gamma}{2z_k B_k} \right) \end{aligned} \quad (19)$$

Therefore, we know that the optimal policy will be a multi-dimensional Gaussian distribution. In order to simplify the notation, let

$$\begin{aligned} \mu(t, S, q, \pi) &= \left(\frac{A_1}{2B_1} - \frac{\mathcal{H}_1^\pm(t, S, q, \pi)}{2z_1}, \dots, \frac{A_N}{2B_N} - \frac{\mathcal{H}_N^\pm(t, S, q, \pi)}{2z_N} \right) \\ \Sigma &= \begin{bmatrix} \frac{\gamma}{2z_1 B_1} & & & & \\ & \frac{\gamma}{2z_1 B_1} & & & \\ & & \ddots & & \\ & & & \frac{\gamma}{2z_k B_k} & \\ & & & & \frac{\gamma}{2z_k B_k} \end{bmatrix} \end{aligned}$$

Under the above notations, the optimal policy is

$$\pi^* \sim \mathcal{N}(\cdot \mid \mu(t, S, q, \pi^*), \Sigma) \quad (20)$$

Policy Improvement Theorem

Theorem 1.1 (policy improvement theorem). *Given any π , let the new policy π_{new} to be*

$$\pi_{new} \sim \mathcal{N}(\cdot \mid \mu(t, S, q, \pi), \Sigma) \quad (21)$$

then following inequality holds

$$V^\pi(t, S, q) \leq V^{\pi_{new}}(t, S, q) \quad (22)$$

Proof. Let $q_t^{\pi_{new}}$ be the inventory process under policy π_{new} . Let the initial condition at time t be $q_t^{\pi_{new}} = q$, and $S_t = S$. Then by the Ito formula, and averaging over all possibilities, we have the following

$$\begin{aligned}
& V^\pi(t, S, q) \\
&= \mathbb{E} \left[V^\pi(s, S_s, q_s^{\pi_{new}}) + \int_t^s \int_{\epsilon_u} \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) V^\pi(u, S_u, q_u^{\pi_{new}}) \sum_k [\lambda_u^+(k) + \lambda_u^-(k)] d\epsilon_u du \right. \\
&\quad - \int_t^s \int_{\epsilon_u} \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) \sum_k \left[V^\pi(u, S_u, q_u^{\pi_{new}} + z_k) \lambda_u^+(k) + V^\pi(u, S_u, q_u^{\pi_{new}} - z_k) \lambda_u^-(k) \right] d\epsilon_u du \\
&\quad \left. - \int_s^t \left(\partial_t V^\pi(u, S_u, q_u^{\pi_{new}}) + \frac{1}{2} \sigma^2 \partial_{SS} V^\pi(u, S_u, q_u^{\pi_{new}}) \right) du \mid S_t = S, q_t^{\pi_{new}} = q \right] \tag{23}
\end{aligned}$$

Since at time t , under policy π , the following equality holds,

$$\begin{aligned}
& \int_{\epsilon_t} \sum_k \left[\lambda_t^+(k) V^\pi(t, S, q + z_k) + \lambda_t^-(k) V^\pi(t, S, q - z_k) - (\lambda_t^+(k) + \lambda_t^-(k)) V^\pi(t, S, q) \right] \pi(\epsilon_t | t, S, q) d\epsilon_t \\
&+ \int_{\epsilon_t} \sum_{k=1}^N \left[z_k \lambda_t^+(k) (S + \epsilon_t^b(k) - \delta h(q)) - z_k \lambda_t^-(k) (S - \epsilon_t^a(k) - \delta h(q)) \right] \pi(\epsilon_t | t, S, q) d\epsilon_t \\
&- \gamma \int_{\epsilon_t} \pi(\epsilon_t | t, S, q) \log \pi(\epsilon_t | t, S, q) d\epsilon_t + \partial_t V^\pi(t, S, q) + \frac{1}{2} \sigma^2 \partial_{SS} V^\pi(t, S, q) \\
&= 0 \tag{24}
\end{aligned}$$

For π_{new} , based on its construction, and by the same calculus of variation arguments as in equations (10) – (13), π_{new} is the maximizer for the following quantity,

$$\begin{aligned}
& \max_{\pi} \left\{ \int_{\epsilon_t} \sum_{k=1}^N \left[z_k \lambda_t^+(k) (S + \epsilon_t^b(k) - \delta h(q)) - z_k \lambda_t^-(k) (S - \epsilon_t^a(k) - \delta h(q)) \right] \tilde{\pi}(\epsilon_t | t, S, q) d\epsilon_t \right. \\
&+ \int_{\epsilon_t} \sum_k \left[\lambda_t^+(k) V^\pi(t, S, q + z_k) + \lambda_t^-(k) V^\pi(t, S, q - z_k) - (\lambda_t^+(k) + \lambda_t^-(k)) V^\pi(t, S, q) \right] \tilde{\pi}(\epsilon_t | t, S, q) d\epsilon_t \\
&\left. - \gamma \int_{\epsilon_t} \tilde{\pi}(\epsilon_t | t, S, q) \log \tilde{\pi}(\epsilon_t | t, S, q) d\epsilon_t \right\} \tag{25}
\end{aligned}$$

Which results in the following inequality,

$$\begin{aligned}
& \int_{\epsilon_t} \sum_k \left[\lambda_t^+(k) V^\pi(t, S, q + z_k) + \lambda_t^-(k) V^\pi(t, S, q - z_k) - (\lambda_t^+(k) + \lambda_t^-(k)) V^\pi(t, S, q) \right] \pi_{new}(\epsilon_t | t, S, q) d\epsilon_t \\
&+ \int_{\epsilon_t} \sum_{k=1}^N \left[z_k \lambda_t^+(k) (S + \epsilon_t^b(k) - \delta h(q)) - z_k \lambda_t^-(k) (S - \epsilon_t^a(k) - \delta h(q)) \right] \pi_{new}(\epsilon_t | t, S, q) d\epsilon_t \\
&- \gamma \int_{\epsilon_t} \pi_{new}(\epsilon_t | t, S, q) \log \pi_{new}(\epsilon_t | t, S, q) d\epsilon_t + \partial_t V^\pi(t, S, q) + \frac{1}{2} \sigma^2 \partial_{SS} V^\pi(t, S, q) \\
&\geq 0 \tag{26}
\end{aligned}$$

Then equation (23) yields

$$\begin{aligned}
& V^\pi(t, S, q) \\
& \leq \mathbb{E} \left[\int_t^s \int_{\epsilon_t} \sum_{k=1}^N \left[z_k \lambda_u^+(k) (S_u + \epsilon_u^b(k) - \delta h(q_u^{\pi_{new}})) - z_k \lambda_u^-(k) (S_u - \epsilon_u^a(k) - \delta h(q_u^{\pi_{new}})) \right] \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) d\epsilon_u du \right. \\
& \quad \left. - \gamma \int_t^s \int_{\epsilon_u} \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) \log \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) d\epsilon_u du + V^\pi(s, S_s, q_s^{\pi_{new}}) \mid S_t = S, q_t^{\pi_{new}} = q \right] \quad (27)
\end{aligned}$$

Set $s = T$, we have $V^\pi(T, S_T, q_T^{\pi_{new}}) = V^{\pi_{new}}(T, S_T, q_T^{\pi_{new}})$ then we have

$$\begin{aligned}
& V^\pi(t, S, q) \\
& \leq \mathbb{E} \left[\int_t^T \int_{\epsilon_t} \sum_{k=1}^N \left[z_k \lambda_u^+(k) (S_u + \epsilon_u^b(k) - \delta h(q_u^{\pi_{new}})) - z_k \lambda_u^-(k) (S_u - \epsilon_u^a(k) - \delta h(q_u^{\pi_{new}})) \right] \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) d\epsilon_u du \right. \\
& \quad \left. - \gamma \int_t^T \int_{\epsilon_u} \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) \log \pi_{new}(\epsilon_u | u, S_u, q_u^{\pi_{new}}) d\epsilon_u du + V^\pi(T, S_T, q_T^{\pi_{new}}) \mid S_t = S, q_t^{\pi_{new}} = q \right] \\
& \leq V^{\pi_{new}}(t, S, q) \quad (28)
\end{aligned}$$

□

Martingale Loss

Given a policy π , and q_t^π is inventory process under policy π . Let the initial condition at time t to be $S_t = S$, $q_t^\pi = q$, the value function under policy π is

$$\begin{aligned}
& V^\pi(t, S, q) \\
& = \mathbb{E} \left[\int_t^s \int_{\epsilon_u} \sum_{k=1}^N \left[z_k (S_u + \epsilon_u^b(k) - \delta h(q_u^\pi)) dN_u^+(k) - z_k (S_u - \epsilon_u^a(k) - \delta h(q_u^\pi)) dN_u^-(k) \right] \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u \right. \\
& \quad \left. - \gamma \int_t^s \int_{\epsilon_u} \pi(\epsilon_u | u, S_u, q_u^\pi) \log \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u du + V(s, S_s, q_s^\pi) \mid S_t = S, q_t^\pi = q \right] \quad (29)
\end{aligned}$$

Then we have

$$\begin{aligned}
& \mathbb{E} \left[\frac{1}{s-t} \int_t^s \int_{\epsilon_u} \sum_{k=1}^N \left[z_k (S_u + \epsilon_u^b(k) - \delta h(q_u^\pi)) dN_u^+(k) - z_k (S_u - \epsilon_u^a(k) - \delta h(q_u^\pi)) dN_u^-(k) \right] \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u \right. \\
& \quad \left. - \frac{\gamma}{s-t} \int_t^s \int_{\epsilon_u} \pi(\epsilon_u | u, S_u, q_u^\pi) \log \pi(\epsilon_u | u, S_u, q_u^\pi) d\epsilon_u du + \frac{V^\pi(s, S_s, q_s^\pi) - V^\pi(t, S_t, q_t^\pi)}{s-t} \mid S_t = S, q_t^\pi = q \right] = 0 \quad (30)
\end{aligned}$$

Let $s \rightarrow t$, and parametrize the value function V_θ^π , define the temporal difference error as,

$$\begin{aligned}
\delta_t^\theta & = \lim_{s \rightarrow t} \mathbb{E} \left[\frac{V_\theta^\pi(s, S_s, q_s^\pi) - V_\theta^\pi(t, S_t, q_t^\pi)}{s-t} \mid S_t = S, q_t^\pi = q \right] - \gamma \int_{\epsilon_t} \pi(\epsilon_t | t, S, q) \log \pi(\epsilon_t | t, S, q) d\epsilon_t \\
& \quad + \int_{\epsilon_t} \sum_{k=1}^N \left[z_k (S + \epsilon_t^b(k) - \delta h(q)) dN_t^+(k) - z_k (S - \epsilon_t^a(k) - \delta h(q)) dN_t^-(k) \right] \pi(\epsilon_t | t, S, q) d\epsilon_t \quad (31)
\end{aligned}$$

Using the Monte Carlo method to generate a set of sample paths $\mathcal{D} = \{(t_i, S_i^d, q_i^d)_{i=1}^T\}_{d=1}^D$, then define the martingale loss as

$$\begin{aligned} \mathbf{ML}(\theta) = & \frac{1}{2} \sum_{\mathcal{D}} \sum_i \left(\frac{V_{\theta}^{\pi}(t_{i+1}, S_{t_{i+1}}^d, q_{t_{i+1}}^d) - V(t_i, S_{t_i}^d, q_{t_i}^d)}{\Delta t} - \gamma \int_{\epsilon_{t_i}} \pi(\epsilon_{t_i} | t_i, S_{t_i}^d, q_{t_i}^d) \log \pi(\epsilon_{t_i} | t_i, S_{t_i}^d, q_{t_i}^d) d\epsilon_{t_i} \right. \\ & \left. \int_{\epsilon_{t_i}} \sum_{k=1}^N \left[z_k(S_{t_i}^d + \epsilon_{t_i}^b(k) - \delta h(q_{t_i}^d)) dN_{t_i}^+(k) - z_k(S_{t_i}^d - \epsilon_{t_i}^a(k) - \delta h(q_{t_i}^d)) dN_{t_i}^-(k) \right] \pi(\epsilon_{t_i} | t_i, S_{t_i}^d, q_{t_i}^d) d\epsilon_{t_i} \right)^2 \Delta t \end{aligned} \quad (32)$$

Algorithm 1 EMM: Exploratory Market Making

Require: Initialize hyperparameters

for $l = 1$ to L **do**

for $m = 1$ to M **do**

 Generate one sample path $\mathcal{D} = \{(t_i, S_{t_i}, q_{t_i})_{i=1}^T\}$ under policy π^{ϕ}

 Compute $\mathbf{ML}(\theta)$

 Updates $\theta \leftarrow \theta - \alpha \nabla_{\theta} \mathbf{ML}(\theta)$

end for

 Update $\pi^{\phi} \leftarrow \mathcal{N}\left(\epsilon \mid \left(\frac{A_k}{2B_k} - \frac{\mathcal{H}_{\theta}^{\pm}(t, S, q, \pi^{\phi})}{2z_k}\right), \Sigma\right)$

end for
