

SUPPLEMENTARY MATERIALS: POLICY GRADIENT METHODS FOR THE NOISY LINEAR QUADRATIC REGULATOR OVER A FINITE HORIZON

BEN HAMBLY *, RENYUAN XU *, AND HUINING YANG *†

SM1. Market Simulator for Linear Price Dynamics. We estimate the parameters for the LQR model using NASDAQ ITCH data taken from Lobster¹.

Permanent Price Impact and Volatility. The model in (2.7) implies that prices changes are proportional to the market-order flow imbalances (MFI). We adopt the framework from [1], namely that the price change ΔS is given by

$$(SM1.1) \quad \Delta S = \gamma \text{MFI} + \sigma \epsilon,$$

with $\text{MFI} = M^b - M^s$ where M^s and M^b are the volumes of market sell orders and market buy orders respectively during a time interval $\Delta T = 5\text{mins}$ and $\epsilon \sim \mathcal{N}(0, 1)$. We then estimate γ and σ from the data.

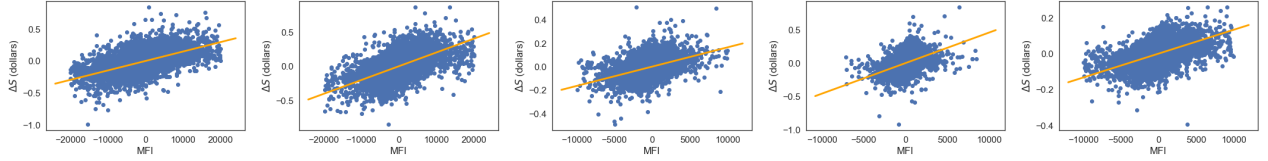


Fig. SM1: Relationship between MFI and ΔS . (Example (from left to right): AAP, FB, JPM, IBM and AAL, 10:00AM-11:00AM 01/01/2019-08/31/2019, $\Delta T = 1\text{min}$)

Temporary Price Impact. We assume the LOB has a flat shape with constant queue length l for the first few levels. Figure SM2 shows the average queue lengths for the first 5 levels so that our assumption is not too unreasonable. Therefore the following equation, on the amount received when we liquidate u shares with best bid price S , holds

$$u(S - \beta u) = \int_{S - \frac{u\Delta}{l}}^S l v dv.$$

Therefore we have $\beta = \frac{\Delta}{2l}$, where Δ is the tick size and l is the average queue length.

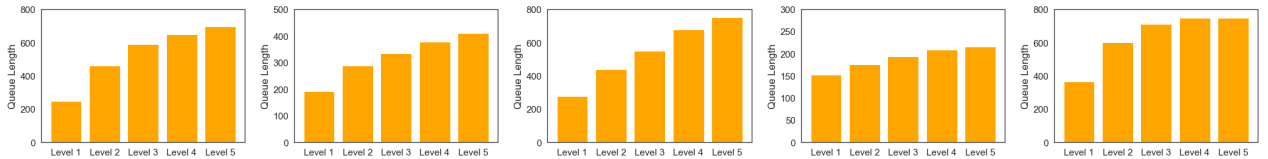


Fig. SM2: Average queue length (volume) of the first five levels on the limit buy side (Example (from left to right): AAP, FB, JPM, IBM and AAL, 10:00AM-11:00AM 01/01/2019-08/31/2019 with 5000 samples uniformly sampled with natural time clock in each trading day.)

Parameter Estimation. See the estimates for AAPL, FB, IBM, JPM, and AAL in Table SM1.

*Mathematical Institute, University of Oxford. **Email:** {hambly, xur, yang}@maths.ox.ac.uk

† Supported by the EPSRC Centre for Doctoral Training in Industrially Focused Mathematical Modelling (EP/L015803/1) in collaboration with BP plc.

¹<https://lobsterdata.com/>

Paramters/Stock	AAPL	FB	IBM	JPM	AAL
β	1.03×10^{-5}	1.30×10^{-5}	2.65×10^{-5}	9.28×10^{-6}	3.27×10^{-5}
γ	7.27×10^{-6}	1.40×10^{-5}	4.60×10^{-5}	1.65×10^{-5}	1.3310×10^{-5}
σ	0.107	0.115	0.082	0.059	0.042

Table SM1: Parameter estimation from NASDAQ ITCH Data (10:00AM-11:00PM 01/01/2019-08/31/2019).

SM2. Comparison between the Policy Gradient Method and Q-learning. The computational complexity of Q-learning is highly dependent on the size of the set of the (discrete) states and actions. Therefore Q-learning is typically less suited to problems with continuous and unbounded states and actions. In order to apply Q-learning for such problems, we need to discretize the continuous state and action space. Intuitively speaking, Q-learning suffers from low accuracy when the discretization scheme is less refined (see Figures SM3 and SM5). On the other hand, the computational complexity grows quadratically when increasing the level of granularity of discretization (see Figures SM4 and SM5).

To demonstrate this view point, we compare the performance of the Q-learning algorithm with the policy gradient method on a one-dimensional LQR problem with finite horizon as suggested by the reviewer. (We would expect the deep Q-learning algorithm and the deep policy gradient method to have similar comparison results.)

Q learning update. We initialize the Q table $\{q_t^{(0)}(x, u)\}_{x,u,t}$ with all zeros. In the i -th iteration, we update the Q table for $t = 0, 1, \dots, T - 1$,

$$(SM2.1) \quad q_t^{(i)}(x, u) = (1 - \tilde{\eta}) q_t^{(i-1)}(x, u) + \tilde{\eta} \left[c_t(x, u) + \min_{u'} q_{t+1}^{(i)}(x', u') \right],$$

with terminal condition $q_T^{(i)}(x, u) = x^2 Q_T$. Here $c_t(x, u) = x^2 Q_t + u^2 R_t$ is the instantaneous cost at time t ; x' is the next state simulated from the system when the agent takes an action u in state x at time t ; and $\tilde{\eta} \in (0, 1)$ is the learning rate.

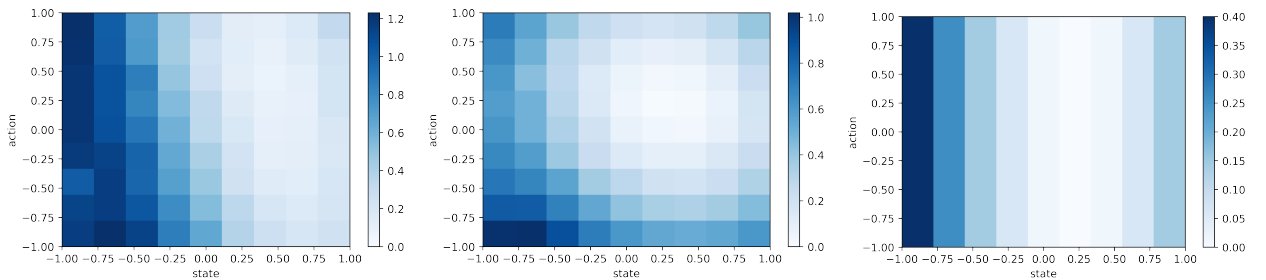
Model set-up. We set $d = k = 1$, $T = 5$, $A = 1.0$, $B = 0.2$, $Q_t = 0.2$ for $t = 0, 1, 2, 3, 4$, $Q_T = 0.4$, $R_t = 0.1(t + 1)$ for $t = 0, 1, \dots, 4$, $w_t \sim \mathcal{N}(0, 0.1)$, and $x_0 \sim \mathcal{N}(0, 0.1)$.

Parameter set-up. To perform Q-learning, we uniformly partition the states and actions in $[-1, 1]$. We set the learning rate for Q-learning as $\tilde{\eta} = 0.1$. For the policy gradient method, we set the learning rate $\eta = 0.2$ and the number of trajectories in the zero-th optimization as $m = 50$.

Conclusion. From Figure SM5, we observe that

- For LQR with finite horizon, the policy gradient method outperforms Q-learning algorithms (with the size of actions and states varying from 10 to 100) in terms of both sample efficiency and accuracy.
- When increasing the size of the states and actions from 10 to 100, the accuracy of the Q-learning algorithm improves, however, it requires many more samples to converge.

To conclude, Q-learning is less suited to handling decision-making problems with continuous and unbounded states and actions. More advanced approximation techniques may be needed in this case [2].

Fig. SM3: Q tables with 10 states and 10 actions: $q_0(s, a)$, $q_4(s, a)$ and $q_5(s, a)$ (from left to right).

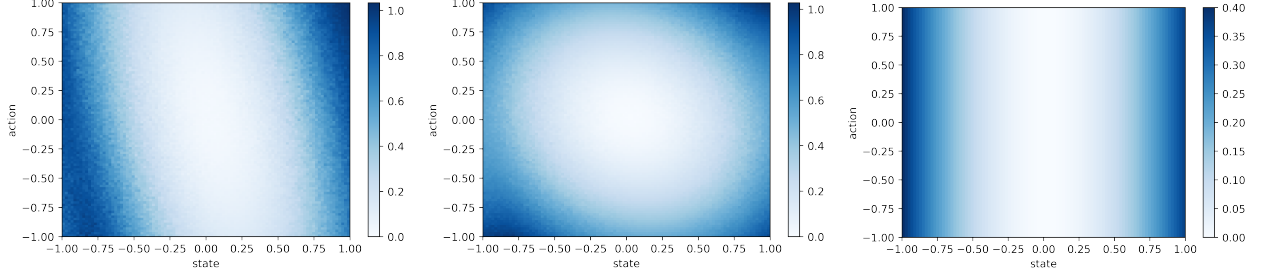


Fig. SM4: Q tables with 100 states and 100 actions: $q_0(s, a)$, $q_4(s, a)$ and $q_5(s, a)$ (from left to right).

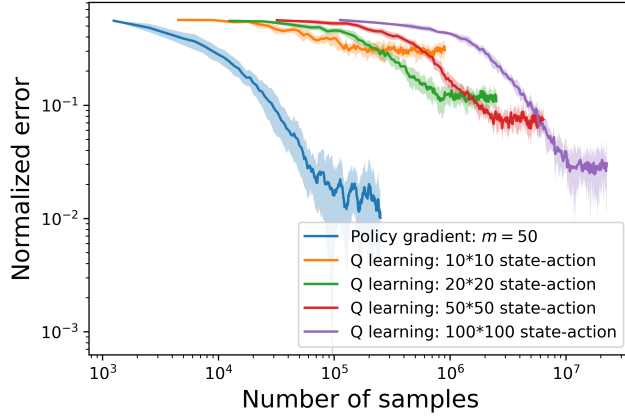


Fig. SM5: Comparison between Q-learning and the policy gradient method (log-log scale). (Average across 10 scenarios.)

SM3. Proofs of Technical Results. We now give the proofs that were omitted in the text.

SM3.1. Proofs in Section 3.1.

Proof of Lemma 3.2. Denote by $\{x_t\}_{t=0}^T$ the state trajectory induced by an arbitrary control \mathbf{K} . By Assumption 3.1 the matrix $\mathbb{E}[x_0 x_0^\top]$ is positive definite. For $t \geq 1$, we have

$$\mathbb{E}[x_t x_t^\top] = (A - BK_{t-1})\mathbb{E}[x_{t-1} x_{t-1}^\top](A - BK_{t-1})^\top + \mathbb{E}[w_{t-1} w_{t-1}^\top].$$

Now $(A - BK_{t-1})\mathbb{E}[x_{t-1} x_{t-1}^\top](A - BK_{t-1})^\top$ is positive semi-definite and $\mathbb{E}[w_{t-1} w_{t-1}^\top]$ is positive definite. Hence $\mathbb{E}[x_t x_t^\top]$ is positive definite and as a result $\underline{\sigma}_{\mathbf{X}} > 0$. In this case, we can simply take $\underline{\sigma}_{\mathbf{X}} = \min(\mathbb{E}[x_0 x_0^\top], \sigma_{\min}(W))$. \square

Proof of Proposition 3.4. This can be proved by backward induction. For $t = T$, $P_T^K = Q_T$ is positive definite since Q_T is positive definite. Assume P_{t+1}^K is positive definite for some $t + 1$, then take any $z \in \mathbb{R}^d$ such that $z \neq 0$,

$$z^\top P_t^K z = z^\top Q_t z + z^\top K_t^\top R_t K_t z + z^\top (A - BK_t)^\top P_{t+1}^K (A - BK_t) z > 0.$$

The last inequality holds since $z^\top Q_t z > 0$, $z^\top K_t^\top R_t K_t z \geq 0$ and $z^\top (A - BK_t)^\top P_{t+1}^K (A - BK_t) z \geq 0$. By backward induction, we have P_t^K positive definite, $\forall t = 0, 1, \dots, T$. \square

To prove Lemma 3.6, let us start with a useful result for the value function. Define the value function $V_{\mathbf{K}}(x, \tau)$ for $\tau = 0, 1, \dots, T - 1$, as

$$V_{\mathbf{K}}(x, \tau) = \mathbb{E}_{\mathbf{w}} \left[\sum_{t=\tau}^{T-1} (x_t^\top Q_t x_t + u_t^\top R_t u_t) + x_T^\top Q_T x_T \middle| x_\tau = x \right] = x^\top P_\tau x + L_\tau,$$

62 with terminal condition

$$63 \quad V_{\mathbf{K}}(x, T) = x^\top Q_T x,$$

64 where L_τ is defined in (3.10). We then define the Q function, $Q_{\mathbf{K}}(x, u, \tau)$ for $\tau = 0, 1, \dots, T-1$ as

$$65 \quad Q_{\mathbf{K}}(x, u, \tau) = x^\top Q_\tau x + u^\top R_\tau u + \mathbb{E}_{w_\tau} [V_{\mathbf{K}}(Ax + Bu + w_\tau, \tau + 1)],$$

66 and the advantage function

$$67 \quad A_{\mathbf{K}}(x, u, \tau) = Q_{\mathbf{K}}(x, u, \tau) - V_{\mathbf{K}}(x, \tau).$$

68 Note that $C(\mathbf{K}) = \mathbb{E}_{x_0 \sim \mathcal{D}}[V(x_0, 0)]$. Then we can write the difference of value functions between \mathbf{K} and \mathbf{K}' in
69 terms of advantage functions.

70 **LEMMA SM3.1.** Assume \mathbf{K} and \mathbf{K}' have finite costs. Denote $\{x'_t\}_{t=0}^T$ and $\{u'_t\}_{t=0}^{T-1}$ as the state and control
71 sequences of a single trajectory generated by \mathbf{K}' starting from $x'_0 = x_0 = x$, then

$$72 \quad (\text{SM3.1}) \quad V_{\mathbf{K}'}(x, 0) - V_{\mathbf{K}}(x, 0) = \mathbb{E}_{\mathbf{w}} \left[\sum_{t=0}^{T-1} A_{\mathbf{K}}(x'_t, u'_t, t) \right],$$

73 and $A_{\mathbf{K}}(x, -K'_\tau x, \tau) = 2x^\top (K'_\tau - K_\tau)^\top E_\tau x + x^\top (K'_\tau - K_\tau)^\top (R_\tau + B^\top P_{\tau+1} B)(K'_\tau - K_\tau)x$, where E_τ is defined
74 in (3.11).

75 *Proof.* Denote by $c'_t(x)$ the cost generated by \mathbf{K}' with a single trajectory starting from $x'_0 = x_0 = x$. That
76 is, $c'_t(x) = (x'_t)^\top Q_t x'_t + (u'_t)^\top R_t u'_t$, $t = 0, 1, \dots, T-1$, and $c'_T(x) = (x'_T)^\top Q_T x'_T$, with $u'_t = -K'_t x'_t$, $x'_{t+1} =$
77 $Ax'_t + Bu'_t + w_t$, $x'_0 = x$.

78 Therefore,

$$\begin{aligned} V_{\mathbf{K}'}(x, 0) - V_{\mathbf{K}}(x, 0) &= \mathbb{E}_{\mathbf{w}} \left[\sum_{t=0}^T c'_t(x) \right] - V_{\mathbf{K}}(x, 0) = \mathbb{E}_{\mathbf{w}} \left[\sum_{t=0}^T (c'_t(x) + V_{\mathbf{K}}(x'_t, t) - V_{\mathbf{K}}(x'_t, t)) \right] - V_{\mathbf{K}}(x, 0) \\ 79 \quad &= \mathbb{E}_{\mathbf{w}} \left[\sum_{t=0}^{T-1} (c'_t(x) + V_{\mathbf{K}}(x'_{t+1}, t+1) - V_{\mathbf{K}}(x'_t, t)) \right] \\ &= \mathbb{E}_{\mathbf{w}} \left[\sum_{t=0}^{T-1} (Q_{\mathbf{K}}(x'_t, u'_t, t) - V_{\mathbf{K}}(x'_t, t)) \middle| x_0 = x \right] = \mathbb{E}_{\mathbf{w}} \left[\sum_{t=0}^{T-1} A_{\mathbf{K}}(x'_t, u'_t, t) \middle| x_0 = x \right], \end{aligned}$$

80 where the third equality holds since $c'_T(x) = V_{\mathbf{K}}(x'_T, T)$ with the same single trajectory. For $u = -K'_\tau x$,

(SM3.2)

$$\begin{aligned} A_{\mathbf{K}}(x, -K'_\tau x, \tau) &= Q_{\mathbf{K}}(x, -K'_\tau x, \tau) - V_{\mathbf{K}}(x, \tau) \\ &= x^\top (Q_\tau + (K'_\tau)^\top R_\tau K'_\tau)x + \mathbb{E}_{w_\tau} [V_{\mathbf{K}}((A - BK'_\tau)x + w_\tau, \tau + 1)] - V_{\mathbf{K}}(x, \tau) \\ &= x^\top (Q_\tau + (K'_\tau)^\top R_\tau K'_\tau)x + (x^\top (A - BK'_\tau)^\top P_{\tau+1} (A - BK'_\tau)x + \text{Tr}(WP_{\tau+1}) + L_{\tau+1}) \\ &\quad - (x^\top P_\tau x + L_\tau) \\ 81 \quad &= x^\top (Q_\tau + (K'_\tau - K_\tau + K_\tau)^\top R_\tau (K'_\tau - K_\tau + K_\tau))x \\ &\quad + x^\top (A - BK_\tau - B(K'_\tau - K_\tau))^\top P_{\tau+1} (A - BK_\tau - B(K'_\tau - K_\tau))x \\ &\quad - x^\top (Q_\tau + K_\tau^\top R_\tau K_\tau + (A - BK_\tau)^\top P_{\tau+1} (A - BK_\tau))x \\ &= 2x^\top (K'_\tau - K_\tau)^\top ((R_\tau + B^\top P_{\tau+1} B)K_\tau - B^\top P_{\tau+1} A)x \\ &\quad + x^\top (K'_\tau - K_\tau)^\top (R_\tau + B^\top P_{\tau+1} B)(K'_\tau - K_\tau)x. \end{aligned}$$

□

82 **Proof of Lemma 3.6.** First for any K'_τ , from (SM3.2),

$$\begin{aligned}
 A_{\mathbf{K}}(x, -K'_\tau x, \tau) &= Q_{\mathbf{K}}(x, -K'_\tau x, \tau) - V_{\mathbf{K}}(x, \tau) \\
 &= 2 \operatorname{Tr}(xx^\top (K'_\tau - K_\tau)^\top E_\tau) + \operatorname{Tr}(xx^\top (K'_\tau - K_\tau)^\top (R_\tau + B^\top P_{\tau+1} B)(K'_\tau - K_\tau)) \\
 &= \operatorname{Tr}(xx^\top (K'_\tau - K_\tau + (R_\tau + B^\top P_{\tau+1} B)^{-1} E_\tau)^\top (R_\tau + B^\top P_{\tau+1} B) \\
 &\quad (K'_\tau - K_\tau + (R_\tau + B^\top P_{\tau+1} B)^{-1} E_\tau)) - \operatorname{Tr}(xx^\top E_\tau^\top (R_\tau + B^\top P_{\tau+1} B)^{-1} E_\tau) \\
 &\geq -\operatorname{Tr}(xx^\top E_\tau^\top (R_\tau + B^\top P_{\tau+1} B)^{-1} E_\tau),
 \end{aligned}
 \tag{SM3.3}$$

84 with equality holds when $K'_\tau = K_\tau - (R_\tau + B^\top P_{\tau+1} B)^{-1} E_\tau$. Then,

$$\begin{aligned}
 C(\mathbf{K}) - C(\mathbf{K}^*) &= -\mathbb{E} \sum_{t=0}^{T-1} A_{\mathbf{K}}(x_t^*, u_t^*, t) \leq \mathbb{E} \sum_{t=0}^{T-1} \operatorname{Tr}(x_t^* (x_t^*)^\top E_t^\top (R_t + B^\top P_{t+1} B)^{-1} E_t) \\
 &\leq \|\Sigma_{\mathbf{K}^*}\| \sum_{t=0}^{T-1} \operatorname{Tr}(E_t^\top (R_t + B^\top P_{t+1} B)^{-1} E_t) \leq \frac{\|\Sigma_{\mathbf{K}^*}\|}{\underline{\sigma}_{\mathbf{R}}} \sum_{t=0}^{T-1} \operatorname{Tr}(E_t^\top E_t) \\
 &\leq \frac{\|\Sigma_{\mathbf{K}^*}\|}{4 \underline{\sigma}_{\mathbf{X}}^2 \underline{\sigma}_{\mathbf{R}}} \sum_{t=0}^{T-1} \operatorname{Tr}(\nabla_t C(\mathbf{K})^\top \nabla_t C(\mathbf{K})),
 \end{aligned}$$

88 where $\underline{\sigma}_{\mathbf{X}}$ is defined in (3.3) and $\underline{\sigma}_{\mathbf{R}}$ is defined in (3.4). For the lower bound, consider $K'_t = K_t - (R_t + B^\top P_{t+1} B)^{-1} E_t$ where the equality holds in (SM3.3). Using $C(\mathbf{K}^*) \leq C(\mathbf{K}')$

(SM3.4) □

$$\begin{aligned}
 C(\mathbf{K}) - C(\mathbf{K}^*) &\geq C(\mathbf{K}) - C(\mathbf{K}') = -\mathbb{E} \sum_{t=0}^{T-1} A_{\mathbf{K}}(x'_t, u'_t, t) = \mathbb{E} \sum_{t=0}^{T-1} \operatorname{Tr}(x'_t (x'_t)^\top E_t^\top (R_t + B^\top P_{t+1} B)^{-1} E_t) \\
 &\geq \underline{\sigma}_{\mathbf{X}} \sum_{t=0}^{T-1} \frac{1}{\|R_t + B^\top P_{t+1} B\|} \operatorname{Tr}(E_t^\top E_t)
 \end{aligned}$$

91 **Proof of Lemma 3.7.** By lemma SM3.1 we have

$$\begin{aligned}
 C(\mathbf{K}') - C(\mathbf{K}) &= \mathbb{E} \left[\sum_{t=0}^{T-1} A_{\mathbf{K}}(x'_t, -K'_t x'_t, t) \right] \\
 &= \sum_{t=0}^{T-1} (2 \operatorname{Tr}(\Sigma'_t (K'_t - K_t)^\top E_t) + \operatorname{Tr}(\Sigma'_t (K'_t - K_t)^\top (R_t + B^\top P_{t+1} B)(K'_t - K_t))).
 \end{aligned}$$

93 **Proof of Lemma 3.8.** For $t = 0, 1, \dots, T$,

$$C(\mathbf{K}) \geq \mathbb{E}[x_t^\top P_t x_t] \geq \|P_t\| \sigma_{\min}(\mathbb{E}[x_t x_t^\top]) \geq \underline{\sigma}_{\mathbf{X}} \|P_t\|,$$

$$C(\mathbf{K}) = \sum_{t=0}^{T-1} \operatorname{Tr}(\mathbb{E}[x_t x_t^\top](Q_t + K_t^\top R_t K_t)) + \operatorname{Tr}(\mathbb{E}[x_T x_T^\top] Q_T) \geq \underline{\sigma}_{\mathbf{Q}} \operatorname{Tr}(\Sigma_{\mathbf{K}}) \geq \underline{\sigma}_{\mathbf{Q}} \|\Sigma_{\mathbf{K}}\|.$$

97 Therefore the statement in Lemma 3.8 follows provided that $\underline{\sigma}_{\mathbf{X}} > 0$ and Assumption 2.1 holds. □

98 **Proof of Proposition 3.9.** Recall that $\Sigma_t = \mathbb{E}[x_t x_t^\top]$. Note that

$$\begin{aligned}
 \Sigma_1 &= \mathbb{E}[x_1 x_1^\top] = \mathbb{E} \left[((A - B K_0) x_0 + w_0) ((A - B K_0) x_0 + w_0)^\top \right] \\
 &= (A - B K_0) \Sigma_0 (A - B K_0)^\top + W = \mathcal{G}_0(\Sigma_0) + W.
 \end{aligned}$$

101 Now we first prove that

$$\Sigma_t = \mathcal{G}_{t-1}(\Sigma_0) + \sum_{s=1}^{t-1} D_{t-1,s} W D_{t-1,s}^\top + W, \quad \forall t = 2, 3, \dots, T.$$

102 (SM3.5)

103 When $t = 2$,

$$\begin{aligned}
 104 \quad \Sigma_2 &= \mathbb{E} [x_2 x_2^\top] = \mathbb{E} \left[((A - B K_1)x_1 + w_1)((A - B K_1)x_1 + w_1)^\top \right] \\
 105 \quad &= (A - B K_1)\Sigma_1(A - B K_1)^\top + W = \mathcal{G}_1(\Sigma_0) + (A - B K_1)W(A - B K_1)^\top + W,
 \end{aligned}$$

106 which satisfies (SM3.5). Assume (SM3.5) holds for $t \leq k$. Then for $t = k + 1$,

$$\begin{aligned}
 107 \quad \mathbb{E} [x_{t+1} x_{t+1}^\top] &= \mathbb{E} \left[((A - B K_t)x_t + w_t)((A - B K_t)x_t + w_t)^\top \right] \\
 108 \quad &= (A - B K_t)\Sigma_t(A - B K_t)^\top + W = \mathcal{G}_t(\Sigma_0) + \sum_{s=1}^t D_{t,s} W D_{t,s}^\top + W.
 \end{aligned}$$

109 Therefore (SM3.5) holds, $\forall t = 1, 2, \dots, T$. Finally,

$$110 \quad \Sigma_{\mathbf{K}} = \sum_{t=0}^T \Sigma_t = \Sigma_0 + \sum_{t=0}^{T-1} \mathcal{G}_t(\Sigma_0) + \sum_{t=1}^{T-1} \sum_{s=1}^t D_{t,s} W D_{t,s}^\top + TW = \mathcal{T}_{\mathbf{K}}(\Sigma_0) + \Delta(\mathbf{K}, W). \quad \square$$

111 SM3.2. Proofs in Section 3.2.

112 **Proof of Lemma 3.13.** By direct calculation,

$$113 \quad (\text{SM3.6}) \quad \|\mathcal{G}_t\| \leq \rho^{2(t+1)}, \quad \text{and} \quad \|\mathcal{G}'_t\| \leq \rho^{2(t+1)}.$$

114 Denote $\mathcal{F}_t = \mathcal{F}_{K_t}$ and $\mathcal{F}'_t = \mathcal{F}_{K'_t}$ to ease the exposition. Then for any symmetric matrix $\Sigma \in \mathbb{R}^{d \times d}$ and $t \geq 0$,

$$\begin{aligned}
 115 \quad \|(\mathcal{G}'_{t+1} - \mathcal{G}_{t+1})(\Sigma)\| &= \|\mathcal{F}'_{t+1} \circ \mathcal{G}'_t(\Sigma) - \mathcal{F}_{t+1} \circ \mathcal{G}_t(\Sigma)\| \\
 116 \quad &= \|\mathcal{F}'_{t+1} \circ \mathcal{G}'_t(\Sigma) - \mathcal{F}'_{t+1} \circ \mathcal{G}_t(\Sigma) + \mathcal{F}'_{t+1} \circ \mathcal{G}_t(\Sigma) - \mathcal{F}_{t+1} \circ \mathcal{G}_t(\Sigma)\| \\
 117 \quad &\leq \|\mathcal{F}'_{t+1} \circ \mathcal{G}'_t(\Sigma) - \mathcal{F}'_{t+1} \circ \mathcal{G}_t(\Sigma)\| + \|\mathcal{F}'_{t+1} \circ \mathcal{G}_t(\Sigma) - \mathcal{F}_{t+1} \circ \mathcal{G}_t(\Sigma)\| \\
 118 \quad &= \|\mathcal{F}'_{t+1} \circ (\mathcal{G}'_t - \mathcal{G}_t)(\Sigma)\| + \|(\mathcal{F}'_{t+1} - \mathcal{F}_{t+1}) \circ \mathcal{G}_t(\Sigma)\| \\
 119 \quad &\leq \|\mathcal{F}'_{t+1}\| \|(\mathcal{G}'_t - \mathcal{G}_t)(\Sigma)\| + \|\mathcal{G}_t\| \|\mathcal{F}'_{t+1} - \mathcal{F}_{t+1}\| \|\Sigma\| \\
 120 \quad &\leq \rho^2 \|(\mathcal{G}'_t - \mathcal{G}_t)(\Sigma)\| + \rho^{2(t+1)} \|\mathcal{F}'_{t+1} - \mathcal{F}_{t+1}\| \|\Sigma\|.
 \end{aligned}$$

121 Therefore,

$$122 \quad (\text{SM3.7}) \quad \|(\mathcal{G}'_{t+1} - \mathcal{G}_{t+1})(\Sigma)\| \leq \rho^2 \|(\mathcal{G}'_t - \mathcal{G}_t)(\Sigma)\| + \rho^{2(t+1)} \|\mathcal{F}'_{t+1} - \mathcal{F}_{t+1}\| \|\Sigma\|.$$

123 Summing (SM3.7) up for $t \in \{1, 2, \dots, T-2\}$ with $\|\mathcal{G}'_0 - \mathcal{G}_0\| = \|\mathcal{F}'_0 - \mathcal{F}_0\|$, we have

$$124 \quad \sum_{t=0}^{T-1} \|(\mathcal{G}_t - \mathcal{G}'_t)(\Sigma)\| \leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\sum_{t=0}^{T-1} \|\mathcal{F}_t - \mathcal{F}'_t\| \right) \|\Sigma\|. \quad \square$$

125 SM3.3. Proofs in Section 3.3.

126 **Proof of Lemma 3.15.** Given (3.22) and condition (3.23), we have $\|K'_t - K_t\| = \eta \|\nabla_t C(\mathbf{K})\| \leq \frac{\sigma_{\mathbf{Q}} \sigma_{\mathbf{X}}}{2C(\mathbf{K})\|B\|}$.
 127 Therefore,

$$128 \quad (\text{SM3.8}) \quad \|B\| \|K'_t - K_t\| \leq \frac{\sigma_{\mathbf{Q}} \sigma_{\mathbf{X}}}{2C(\mathbf{K})} \leq \frac{1}{2}.$$

129 The last inequality holds since $\sigma_{\mathbf{X}} \leq \frac{C(\mathbf{K})}{\sigma_{\mathbf{Q}}}$ given by Lemma 3.8. Therefore, by Lemma 3.12,

$$130 \quad (\text{SM3.9}) \quad \sum_{t=0}^{T-1} \|\mathcal{F}_{K_t} - \mathcal{F}_{K'_t}\| \leq (2\rho + 1) \|B\| \left(\sum_{t=0}^{T-1} \|K_t - K'_t\| \right).$$

131 By Lemmas 3.5 and 3.7,

(SM3.10)

$$\begin{aligned}
 C(\mathbf{K}') - C(\mathbf{K}) &= \sum_{t=0}^{T-1} \left[2 \operatorname{Tr} \left(\Sigma'_t (K'_t - K_t)^\top E_t \right) + \operatorname{Tr} \left(\Sigma'_t (K'_t - K_t)^\top (R_t + B^\top P_{t+1} B) (K'_t - K_t) \right) \right] \\
 &= \sum_{t=0}^{T-1} \left[-4\eta \operatorname{Tr} \left(\Sigma'_t \Sigma_t E_t^\top E_t \right) + 4\eta^2 \operatorname{Tr} \left(\Sigma'_t \Sigma_t E_t^\top (R_t + B^\top P_{t+1} B) E_t \Sigma_t \right) \right] \\
 &= \sum_{t=0}^{T-1} \left[-4\eta \operatorname{Tr} \left((\Sigma'_t - \Sigma_t + \Sigma_t) \Sigma_t E_t^\top E_t \right) + 4\eta^2 \operatorname{Tr} \left(\Sigma'_t \Sigma_t E_t^\top (R_t + B^\top P_{t+1} B) E_t \Sigma_t \right) \right] \\
 132 \quad &\leq \sum_{t=0}^{T-1} \left[-4\eta \operatorname{Tr} \left(\Sigma_t E_t^\top E_t \Sigma_t \right) + 4\eta \operatorname{Tr} \left((\Sigma'_t - \Sigma_t) \Sigma_t E_t^\top E_t \Sigma_t \Sigma_t^{-1} \right) \right. \\
 &\quad \left. + 4\eta^2 \operatorname{Tr} \left(\Sigma'_t \Sigma_t E_t^\top (R_t + B^\top P_{t+1} B) E_t \Sigma_t \right) \right] \\
 &\leq \sum_{t=0}^{T-1} \left[-4\eta \operatorname{Tr} \left(\Sigma_t E_t^\top E_t \Sigma_t \right) + 4\eta \frac{\|\Sigma'_t - \Sigma_t\|}{\sigma_{\min}(\Sigma_t)} \operatorname{Tr} \left(\Sigma_t E_t^\top E_t \Sigma_t \right) \right. \\
 &\quad \left. + 4\eta^2 \|\Sigma'_t (R_t + B^\top P_{t+1} B)\| \operatorname{Tr} \left(\Sigma_t E_t^\top E_t \Sigma_t \right) \right] \\
 &\leq -\eta \left(1 - \frac{\sum_{t=0}^{T-1} \|\Sigma'_t - \Sigma_t\|}{\underline{\sigma}_{\mathbf{X}}} - \eta \|\Sigma_{\mathbf{K}'}\| \sum_{t=0}^{T-1} \|R_t + B^\top P_{t+1} B\| \right) \sum_{t=0}^{T-1} \left[\operatorname{Tr} (\nabla_t C(\mathbf{K})^\top \nabla_t C(\mathbf{K})) \right].
 \end{aligned}$$

133 By Lemma 3.6, we have

(SM3.11)

$$134 \quad C(\mathbf{K}') - C(\mathbf{K}) \leq -\eta \left(1 - \frac{\sum_{t=0}^{T-1} \|\Sigma'_t - \Sigma_t\|}{\underline{\sigma}_{\mathbf{X}}} - \eta \|\Sigma_{\mathbf{K}'}\| \sum_{t=0}^{T-1} \|R_t + B^\top P_{t+1} B\| \right) \left(\frac{4\underline{\sigma}_{\mathbf{X}}^2 \underline{\sigma}_{\mathbf{R}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}) - C(\mathbf{K}^*))$$

135 provided that

$$136 \quad (\text{SM3.12}) \quad 1 - \frac{\sum_{t=0}^{T-1} \|\Sigma'_t - \Sigma_t\|}{\underline{\sigma}_{\mathbf{X}}} - \eta \|\Sigma_{\mathbf{K}'}\| \sum_{t=0}^{T-1} \|R_t + B^\top P_{t+1} B\| > 0.$$

137 By (3.21), (3.22), and (SM3.8),

$$138 \quad \sum_{t=0}^{T-1} \|\Sigma'_t - \Sigma_t\| \leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}} + T\|W\| \right) \left(\eta(2\rho + 1)\|B\| \sum_{t=0}^{T-1} \|\nabla_t C(\mathbf{K})\| \right).$$

139 Given the step size condition in (3.23), we have

(SM3.13)

$$140 \quad \eta(2\rho + 1)\|B\| \sum_{t=0}^{T-1} \|\nabla_t C(\mathbf{K})\| \leq \eta(2\rho + 1)\|B\| \left(T \cdot \max_t \{\|\nabla_t C(\mathbf{K})\|\} \right) \leq \frac{(\rho^2 - 1) \underline{\sigma}_{\mathbf{Q}} \underline{\sigma}_{\mathbf{X}}}{2(\rho^{2T} - 1)(C(\mathbf{K}) + \underline{\sigma}_{\mathbf{Q}} T\|W\|)}.$$

141 Then, by Corollary 3.14 and (SM3.9),

$$\begin{aligned}
 142 \quad \frac{\|\Sigma_{\mathbf{K}'} - \Sigma_{\mathbf{K}}\|}{\underline{\sigma}_{\mathbf{X}}} &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\sum_{t=0}^{T-1} \|\mathcal{F}_{K_t} - \mathcal{F}_{K'_t}\| \right) \frac{\|\Sigma_0\| + T\|W\|}{\underline{\sigma}_{\mathbf{X}}} \\
 143 \quad &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1)\|B\| \left(\sum_{t=0}^{T-1} \eta \|\nabla_t C(\mathbf{K})\| \right) \frac{C(\mathbf{K}) + \underline{\sigma}_{\mathbf{Q}} T\|W\|}{\underline{\sigma}_{\mathbf{Q}} \underline{\sigma}_{\mathbf{X}}} \leq \frac{1}{2},
 \end{aligned}$$

144 where the last step holds by (SM3.13). Therefore, the bound of $\|\Sigma_{\mathbf{K}'}\|$ in (SM3.12) is given by

$$145 \quad (\text{SM3.14}) \quad \|\Sigma_{\mathbf{K}'}\| \leq \|\Sigma_{\mathbf{K}'} - \Sigma_{\mathbf{K}}\| + \|\Sigma_{\mathbf{K}}\| \leq \frac{1}{2} \underline{\sigma}_{\mathbf{X}} + \frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}} \leq \frac{1}{2} \|\Sigma_{\mathbf{K}'}\| + \frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}},$$

146 which indicates that $\|\Sigma_{\mathbf{K}'}\| \leq \frac{2C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}}$. Therefore, (SM3.12) gives

$$\begin{aligned}
 & 1 - \frac{\sum_{t=0}^{T-1} \|\Sigma'_t - \Sigma_t\|}{\underline{\sigma}_{\mathbf{X}}} - \eta \|\Sigma_{\mathbf{K}'}\| \sum_{t=0}^{T-1} \|R_t + B^\top P_{t+1} B\| \\
 147 & \geq 1 - \frac{(\rho^{2T} - 1)}{(\rho^2 - 1) \underline{\sigma}_{\mathbf{X}}} \left(\frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}} + T \|W\| \right) \left(\eta(2\rho + 1) \|B\| \sum_{t=0}^{T-1} \|\nabla_t C(\mathbf{K})\| \right) - \eta \frac{2C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}} \sum_{t=0}^{T-1} \|R_t + B^\top P_{t+1} B\| \\
 & = 1 - C_1 \eta,
 \end{aligned}$$

148 where C_1 is defined in (3.24). So if $\eta \leq \frac{1}{2C_1}$, then,

$$149 \quad 1 - \frac{\sum_{t=0}^{T-1} \|\Sigma'_t - \Sigma_t\|}{\underline{\sigma}_{\mathbf{X}}} - \eta \|\Sigma_{\mathbf{K}'}\| \sum_{t=0}^{T-1} \|R_t + B^\top P_{t+1} B\| \geq 1 - C_1 \eta \geq \frac{1}{2} > 0.$$

150 Hence, $C(\mathbf{K}') - C(\mathbf{K}) \leq -\frac{\eta}{2} \left(\frac{4 \underline{\sigma}_{\mathbf{X}}^2 \underline{\sigma}_{\mathbf{R}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}) - C(\mathbf{K}^*))$, and

$$151 \quad C(\mathbf{K}') - C(\mathbf{K}^*) = (C(\mathbf{K}') - C(\mathbf{K})) + (C(\mathbf{K}) - C(\mathbf{K}^*)) \leq \left(1 - 2\eta \frac{\underline{\sigma}_{\mathbf{X}}^2 \underline{\sigma}_{\mathbf{R}}}{\|\Sigma_{\mathbf{K}^*}\|} \right) (C(\mathbf{K}) - C(\mathbf{K}^*)).$$

□

152 **SM3.4. Proofs in Section 4.** Before proceeding to the proof of Theorem 4.5, we show two important
 153 Lemmas which provide the intermediate steps. We first show the optimality condition for the projection operator
 154 in Lemma SM3.2. We then show the one-step convergence result in Lemma SM3.3.

155 **LEMMA SM3.2 (Optimality Condition).** Fix a policy matrix \mathbf{L}^1 and write $\mathbf{L}^* = \Pi_{\mathcal{S}}(\mathbf{L}^1)$. Then for any
 156 $\mathbf{L}^0 \in \mathcal{S}$, we have

$$157 \quad (\text{SM3.15}) \quad \sum_{t=0}^{T-1} \text{Tr}((L_t^0 - L_t^*)(L_t^* - L_t^1)^\top) \geq 0.$$

158 *Proof of Lemma SM3.2.* We show condition (SM3.15) by contradiction. Assume condition (SM3.15) does
 159 not hold, then there exist some $\mathbf{L}^3 \in \mathcal{S}$ and some constant $b > 0$ such that

$$160 \quad (\text{SM3.16}) \quad \sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^* - L_t^1)^\top) = -b < 0.$$

Let

$$M = 1 + \frac{\sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^3 - L_t^*)^\top)}{-2 \sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^* - L_t^1)^\top)} = 1 + \frac{\sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^3 - L_t^*)^\top)}{2b} > 1,$$

161 and take

$$162 \quad \bar{\mathbf{L}} = \frac{1}{M} \mathbf{L}^3 + \left(1 - \frac{1}{M} \right) \mathbf{L}^*.$$

163 By the convexity of \mathcal{S} , we have $\bar{\mathbf{L}} \in \mathcal{S}$. Hence from definition (4.3), we have

$$164 \quad (\text{SM3.17}) \quad \sum_{t=0}^{T-1} \text{Tr}((L_t^* - L_t^1)(L_t^* - L_t^1)^\top) \leq \sum_{t=0}^{T-1} \text{Tr}((\bar{L}_t - L_t^1)(\bar{L}_t - L_t^1)^\top).$$

On the other hand,

$$\begin{aligned}
 & \sum_{t=0}^{T-1} \text{Tr}((\bar{L}_t - L_t^1)(\bar{L}_t - L_t^1)^\top) \\
 &= \sum_{t=0}^{T-1} \text{Tr} \left(\left(\frac{1}{M}(L_t^3 - L_t^*) + (L_t^* - L_t^1) \right) \left(\frac{1}{M}(L_t^3 - L_t^*) + (L_t^* - L_t^1) \right)^\top \right) \\
 &= \sum_{t=0}^{T-1} \frac{1}{M^2} \text{Tr}((L_t^3 - L_t^*)(L_t^3 - L_t^*)^\top) + \sum_{t=0}^{T-1} \text{Tr}((L_t^* - L_t^1)(L_t^* - L_t^1)^\top) \\
 & \quad + \sum_{t=0}^{T-1} \frac{2}{M} \text{Tr}((L_t^3 - L_t^*)(L_t^* - L_t^1)^\top).
 \end{aligned}
 \tag{SM3.18}$$

By the definition of M we have,

$$\begin{aligned}
 & \sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^3 - L_t^*)^\top) + 2M \sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^* - L_t^1)^\top) \\
 &= \sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^3 - L_t^*)^\top) - 2Mb \\
 &= \sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^3 - L_t^*)^\top) - 2b - \sum_{t=0}^{T-1} \text{Tr}((L_t^3 - L_t^*)(L_t^3 - L_t^*)^\top) = -2b < 0.
 \end{aligned}$$

Thus substituting this in (SM3.18) contradicts (SM3.17) which completes the proof.

□

LEMMA SM3.3. Assume Assumption 2.1 holds, $\underline{\sigma}_{\mathbf{X}} > 0$, $\mathbf{K} \in \mathcal{S}$ and that

$$\text{(SM3.19)} \quad K'_t = K_t - \eta \nabla_t C(\mathbf{K}), \quad \text{where}$$

$$\text{(SM3.20)} \quad \eta \leq \min\{C_1, C_2\}, \quad \text{with}$$

$$\text{(SM3.21)} \quad C_1 = \frac{(\rho^2 - 1) \underline{\sigma}_{\mathbf{Q}} \underline{\sigma}_{\mathbf{X}}}{4dT^2 \sqrt{d+k}(\rho^{2T} - 1)(2\rho + 1)(C(\mathbf{K}) + \underline{\sigma}_{\mathbf{Q}} T \|W\|) \|B\| \max_t \{\|\nabla_t C(\mathbf{K})\|\}}$$

$$\text{(SM3.22)} \quad C_2 = \frac{\underline{\sigma}_{\mathbf{Q}}}{8C(\mathbf{K}) \sum_{t=0}^{T-1} \|R_t + B^\top P_{t+1} B\|}.$$

Take

$$\text{(SM3.23)} \quad \tilde{\mathbf{K}} = \Pi_{\mathcal{S}}(\mathbf{K}'),$$

with $\mathbf{K}' = (K'_0, \dots, K'_{T-1})$ and K'_t defined in (SM3.19) ($t = 0, 1, \dots, T-1$). Then we have

$$2\eta \sum_{t=0}^T \|G_t(\mathbf{K})\|_F^2 = 2\eta \sum_{t=0}^T \text{Tr}(G_t(\mathbf{K})^\top G_t(\mathbf{K})) \leq C(\mathbf{K}) - C(\tilde{\mathbf{K}}).$$

Proof. By definition of $\tilde{\mathbf{K}}$, and as $\mathbf{K} \in \mathcal{S}$, we have

$$\text{(SM3.24)} \quad \sum_{t=0}^{T-1} \text{Tr}((\tilde{K}_t - K'_t)(\tilde{K}_t - K'_t)^\top) \leq \sum_{t=0}^{T-1} \text{Tr}((K_t - K'_t)(K_t - K'_t)^\top).$$

Take $\mathbf{L}^1 = \mathbf{K}'$ and $\mathbf{L}^0 = \mathbf{K}$ in Lemma SM3.2, we have

$$(SM3.25) \quad \sum_{t=0}^{T-1} \text{Tr} \left((K_t - \tilde{K}_t)(\tilde{K}_t - K'_t)^\top \right) \geq 0.$$

Combining (SM3.24) and (SM3.25) leads to

$$(SM3.26) \quad \sum_{t=0}^{T-1} \text{Tr} \left((K_t - K'_t)(K_t - \tilde{K}_t)^\top \right) \geq \sum_{t=0}^{T-1} \text{Tr} \left((K_t - \tilde{K}_t)(K_t - \tilde{K}_t)^\top \right).$$

Given the definition (SM3.23), we have $G(\mathbf{K}) = \frac{\tilde{\mathbf{K}} - \mathbf{K}}{2\eta}$ and $G_t(\mathbf{K}) = \frac{\tilde{K}_t - K_t}{2\eta}$. By Lemmas 3.5 and 3.7,

$$(SM3.27) \quad \begin{aligned} C(\tilde{\mathbf{K}}) - C(\mathbf{K}) &= \sum_{t=0}^{T-1} \left[2 \text{Tr} \left(\tilde{\Sigma}_t (\tilde{K}_t - K_t)^\top E_t \right) + \text{Tr} \left(\tilde{\Sigma}_t (\tilde{K}_t - K_t)^\top (R_t + B^\top P_{t+1} B) (\tilde{K}_t - K_t) \right) \right] \\ &= \sum_{t=0}^{T-1} \left[4\eta \text{Tr} \left((\tilde{\Sigma}_t - \Sigma_t) (G_t(\mathbf{K}))^\top E_t \right) + 4\eta \text{Tr} \left(\Sigma_t (G_t(\mathbf{K}))^\top E_t \right) \right. \\ &\quad \left. + 4\eta^2 \text{Tr} \left(\tilde{\Sigma}_t (G_t(\mathbf{K}))^\top (R_t + B^\top P_{t+1} B) G_t(\mathbf{K}) \right) \right], \end{aligned}$$

with $\tilde{\Sigma}_t := \mathbb{E}[\tilde{x}_t \tilde{x}_t^\top]$ and $\{\tilde{x}_t\}_{t=0}^{T-1}$ is the trajectory under policy $\tilde{\mathbf{K}}$.

First, we have

$$(SM3.28) \quad \begin{aligned} \sum_{t=0}^{T-1} \text{Tr} \left(\Sigma_t (G_t(\mathbf{K}))^\top E_t \right) &= \sum_{t=0}^{T-1} \text{Tr} \left((G_t(\mathbf{K}))^\top E_t \Sigma_t \right) = \frac{1}{4\eta^2} \sum_{t=0}^{T-1} \text{Tr} \left((\tilde{K}_t - K_t)^\top (K_t - K'_t) \right) \\ &\leq -\frac{1}{4\eta^2} \sum_{t=0}^{T-1} \text{Tr} \left((\tilde{K}_t - K_t)(\tilde{K}_t - K_t)^\top \right) = -\text{Tr} \left((G_t(\mathbf{K}))^\top (G_t(\mathbf{K})) \right), \end{aligned}$$

in which the last inequality holds by (SM3.26).

Second given (SM3.19) and condition (SM3.21), we have

$$(SM3.29) \quad \|K'_t - K_t\| = \eta \|\nabla_t C(\mathbf{K})\| \leq \frac{\underline{\sigma}_{\mathbf{Q}} \underline{\sigma}_{\mathbf{X}}}{2T\sqrt{d+k}C(\mathbf{K})\|B\|}.$$

Therefore,

$$(SM3.30) \quad \begin{aligned} \sum_{t=0}^{T-1} \|B\| \|\tilde{K}_t - K_t\| &\leq \sum_{t=0}^{T-1} \|B\| \|\tilde{K}_t - K_t\|_F \leq \sum_{t=0}^{T-1} \|B\| \|K'_t - K_t\|_F \\ &\leq \sqrt{d+k} \|B\| \sum_{t=0}^{T-1} \|K'_t - K_t\| \leq \frac{\underline{\sigma}_{\mathbf{Q}} \underline{\sigma}_{\mathbf{X}}}{2C(\mathbf{K})} \leq \frac{1}{2}. \end{aligned}$$

The second inequality holds by (SM3.26) and the last inequality holds since $\underline{\sigma}_{\mathbf{X}} \leq \frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}}$ given by Lemma 3.8.

By (3.21),

$$(SM3.28) \quad \begin{aligned} \sum_{t=0}^{T-1} \|\tilde{\Sigma}_t - \Sigma_t\| &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}} + T\|W\| \right) \left(2\rho \|B\| \|\mathbf{K} - \tilde{\mathbf{K}}\| + \|B\|^2 \|\mathbf{K} - \tilde{\mathbf{K}}\|^2 \right) \\ &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left(\frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}} + T\|W\| \right) \left(2(2\rho + 1) \|B\| \sum_{t=0}^{T-1} \eta \|G_t(\mathbf{K})\| \right) \\ &\leq \frac{\underline{\sigma}_{\mathbf{X}}}{2dT^2\sqrt{d+k} \cdot \max_t \|\nabla_t C(\mathbf{K})\|} \sum_{t=0}^{T-1} \|G_t(\mathbf{K})\|, \end{aligned}$$

where the last inequality holds by step size condition $\eta \leq C_1$. Hence

$$\begin{aligned}
 (\text{SM3.27}) &\leq \sum_{t=0}^{T-1} \left[2\eta \frac{d\|\tilde{\Sigma}_t - \Sigma_t\|}{\sigma_{\min}(\Sigma_t)} \|G_t(\mathbf{K})\| \|\nabla_t C(\mathbf{K})\| - 4\eta \text{Tr}((G_t(\mathbf{K}))^\top (G_t(\mathbf{K}))) \right. \\
 &\quad \left. + 4\eta^2 \|\Sigma_{\tilde{\mathbf{K}}}\| \|R_t + B^\top P_{t+1} B\| \text{Tr}((G_t(\mathbf{K}))^\top G_t(\mathbf{K})) \right] \\
 &\leq 2\eta \frac{d}{\underline{\sigma}_{\mathbf{X}}} \left(\sum_{t=0}^{T-1} \|\tilde{\Sigma}_t - \Sigma_t\| \right) \left(\sum_{t=0}^{T-1} \|G_t(\mathbf{K})\| \right) \left(\sum_{t=0}^{T-1} \|\nabla_t C(\mathbf{K})\| \right) - 4\eta \sum_{t=0}^{T-1} \text{Tr}((G_t(\mathbf{K}))^\top (G_t(\mathbf{K}))) \\
 &\quad + 4\eta^2 \sum_{t=0}^{T-1} \|\Sigma_{\tilde{\mathbf{K}}}\| \|R_t + B^\top P_{t+1} B\| \text{Tr}((G_t(\mathbf{K}))^\top G_t(\mathbf{K})) \\
 &\leq \frac{\eta}{T} \left(\sum_{t=0}^{T-1} \|G_t(\mathbf{K})\| \right)^2 - 4\eta \sum_{t=0}^{T-1} \text{Tr}((G_t(\mathbf{K}))^\top (G_t(\mathbf{K}))) \\
 &\quad + 4\eta^2 \sum_{t=0}^{T-1} \|\Sigma_{\tilde{\mathbf{K}}}\| \|R_t + B^\top P_{t+1} B\| \text{Tr}((G_t(\mathbf{K}))^\top G_t(\mathbf{K})) \\
 &\leq \sum_{t=0}^{T-1} [\eta \|G_t(\mathbf{K})\|^2 - 4\eta \text{Tr}((G_t(\mathbf{K}))^\top (G_t(\mathbf{K}))) \\
 &\quad + 4\eta^2 \|\Sigma_{\tilde{\mathbf{K}}}\| \|R_t + B^\top P_{t+1} B\| \text{Tr}((G_t(\mathbf{K}))^\top G_t(\mathbf{K}))] \\
 &\leq \sum_{t=0}^{T-1} \eta [-3 + 4\eta \|\Sigma_{\tilde{\mathbf{K}}}\| \|R_t + B^\top P_{t+1} B\|] \text{Tr}((G_t(\mathbf{K}))^\top G_t(\mathbf{K})),
 \end{aligned}$$

where the third inequality holds by (SM3.28) and the fourth inequality holds by Cauchy-Schwarz inequality. By (SM3.26) we have $\sqrt{d+k} \sum_{t=0}^{T-1} \|\nabla_t C(\mathbf{K})\| \geq \sum_{t=0}^{T-1} \|G_t(\mathbf{K})\|$ and thus (SM3.28) $\leq \frac{\underline{\sigma}_{\mathbf{X}}}{2}$ and

$$\|\Sigma_{\tilde{\mathbf{K}}}\| \leq \|\Sigma_{\tilde{\mathbf{K}}} - \Sigma_{\mathbf{K}} + \Sigma_{\mathbf{K}}\| \leq \frac{\underline{\sigma}_{\mathbf{X}}}{2} + \|\Sigma_{\mathbf{K}}\| \leq \frac{\|\Sigma_{\tilde{\mathbf{K}}}\|}{2} + \frac{C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}}.$$

Thus $\|\Sigma_{\tilde{\mathbf{K}}}\| \leq \frac{2C(\mathbf{K})}{\underline{\sigma}_{\mathbf{Q}}}$. Therefore when $\eta \leq C_2$, we have

$$(\text{SM3.27}) \leq -2\eta \sum_{t=0}^T \text{Tr}((G_t(\mathbf{K}))^\top (G_t(\mathbf{K}))).$$

□

Proof of Theorem 4.5. The key step in this proof is Lemma SM3.3 and it suffices to show that the projected policy gradient method enjoys sublinear convergence rate in the setting of known parameters. This is because moving from the analysis for the case of known parameters to that for the case of unknown parameters follows the same procedure of policy gradient descent (without projection). In particular, the zeroth order estimation of the gradient term $\eta \nabla_t C(\mathbf{K})$ in (SM3.19) is the same for policy gradient method and projected policy gradient method.

We now show that the projected policy gradient method with known parameters enjoys a sublinear convergence rate. Since the step size conditions (SM3.20)-(SM3.22) are independent of the term $G_t(\mathbf{K})$, the existence of η follows the analysis in Theorem 3.3. Hence when $\eta \in \mathcal{H}(\frac{1}{C(\mathbf{K}^0)+1})$ is an appropriate polynomial in $\frac{1}{C(\mathbf{K}^0)+1}$ and model parameters, by Lemma SM3.3, we have for any $N \in \mathbb{N}^+$,

$$(\text{SM3.29}) \quad \sum_{n=1}^N \left(\sum_{t=0}^{T-1} \text{Tr}((G_t^{\mathbf{K}^n})^\top (G_t^{\mathbf{K}^n})) \right) \leq \frac{\sum_{n=1}^N C(\mathbf{K}^{n-1}) - C(\mathbf{K}^N)}{2\eta} \leq \frac{C(\mathbf{K}^0) - C(\mathbf{K}^*)}{2\eta},$$

Therefore $\left\{ \frac{1}{N} \sum_{n=0}^{N-1} \left(\sum_{t=0}^{T-1} \|G_t(\mathbf{K}^n)\|_F^2 \right) \right\}_{N \geq 1}$ converges at rate $\mathcal{O}(\frac{1}{N})$, which thus completes the proof for the case of known parameters. □

239 **Proof of Lemma 4.7.** Under Assumption 4.3, we have $\mathbb{E}[x_0 x_0^\top] = \widetilde{W}_0 \mathbb{E}[z_0 z_0^\top] \widetilde{W}_0^\top$, and $\|\mathbb{E}[x_0 x_0^\top]\| \leq$
 240 $\sigma_0^2 \|\widetilde{W}_0\|^2$. With the sub-Gaussian distributed noise, $W = \mathbb{E}[w_t w_t^\top] = \widetilde{W} \mathbb{E}[v_t v_t^\top] \widetilde{W}^\top$, then we have $\|W\| \leq$
 241 $\sigma_w^2 \|\widetilde{W}^2\|$.

242 Denote $S_t = Q_t + K_t^\top R_t K_t$, $\forall t = 1, \dots, T-1$. Thus, for $t = 0, 1, \dots, T-2$,

$$\begin{aligned} \mathbb{E}[x_{t+1}^\top Q_{t+1} x_{t+1} + u_{t+1}^\top R_{t+1} u_{t+1}] &= \mathbb{E}[x_{t+1}^\top S_{t+1} x_{t+1}] = \text{Tr}(\mathbb{E}[x_{t+1}^\top S_{t+1} x_{t+1}]) = \text{Tr}(\mathbb{E}[x_{t+1} x_{t+1}^\top] S_{t+1}) \\ &= \text{Tr} \left(\mathcal{G}_t(\Sigma_0) S_{t+1} + \sum_{s=1}^t D_{t,s} W D_{t,s}^\top S_{t+1} + W S_{t+1} \right). \end{aligned}$$

244 The last equality holds by (SM3.5). Therefore,

$$\begin{aligned} C(\mathbf{K}') - C(\mathbf{K}) &= \underbrace{\mathbb{E}[x_0^\top (K'_0)^\top R_0 K'_0 x_0 - x_0^\top K_0^\top R_0 K_0 x_0]}_{(I)} + \underbrace{\sum_{t=0}^{T-2} \text{Tr} \left(\mathcal{G}'_t(\Sigma_0) S'_{t+1} - \mathcal{G}_t(\Sigma_0) S_{t+1} \right)}_{(II)} \\ &\quad + \underbrace{\sum_{t=0}^{T-2} \text{Tr} \left(\sum_{s=1}^t (D'_{t,s} W (D'_{t,s})^\top S'_{t+1} - D_{t,s} W D_{t,s}^\top S_{t+1}) + W (S'_{t+1} - S_{t+1}) \right)}_{(III)} \\ &\quad + \underbrace{\text{Tr} \left(\mathcal{G}_{T-1}(\Sigma_0) Q_T - \mathcal{G}'_{T-1}(\Sigma_0) Q_T + \sum_{s=1}^{T-1} (D'_{T-1,s} W (D'_{T-1,s})^\top Q_T - D_{T-1,s} W D_{T-1,s}^\top Q_T) \right)}_{(IV)}. \end{aligned}$$

246 For the first term, $(I) \leq \text{Tr}(\mathbb{E}[x_0 x_0^\top]) \|(K'_0)^\top R_0 K'_0 - K_0^\top R_0 K_0\|$. For the second term (II) , since

$$247 \sum_{t=0}^{T-2} (\text{Tr}(\mathcal{G}_t(\Sigma_0) S_{t+1})) = \mathbb{E} \left[\sum_{t=0}^{T-2} (\text{Tr}(\Pi_{i=0}^t (A - BK_i) x_0 x_0^\top \Pi_{i=0}^t (A - BK_{t-i})^\top S_{t+1})) \right] \leq \text{Tr}(\mathbb{E}[x_0 x_0^\top]) \left\| \sum_{t=0}^{T-2} \mathcal{G}_t(S_{t+1}) \right\|,$$

248 we have, $(II) \leq \text{Tr}(\mathbb{E}[x_0 x_0^\top]) \left\| \sum_{t=0}^{T-2} (\mathcal{G}'_t(S'_{t+1}) - \mathcal{G}_t(S_{t+1})) \right\|$.

249 We denote $\mathcal{G}_d := \sum_{t=0}^{T-2} (\mathcal{G}'_t(S'_{t+1}) - \mathcal{G}_t(S_{t+1}))$, then

$$\begin{aligned} \|\mathcal{G}_d\| &\leq \sum_{t=0}^{T-2} \left\| \mathcal{G}'_t(Q_{t+1} + (K'_{t+1})^\top R_{t+1} K'_{t+1}) - \mathcal{G}_t(Q_{t+1} + (K'_{t+1})^\top R_{t+1} K'_{t+1}) - \right. \\ &\quad \left. \mathcal{G}_t \circ (K_{t+1}^\top R_{t+1} K_{t+1} - (K'_{t+1})^\top R_{t+1} K'_{t+1}) \right\| \\ &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left((2\rho + 1) \|B\| \sum_{t=0}^{T-2} \|K_t - K'_t\| \right) \left(\sum_{t=1}^{T-1} \|Q_t + (K'_t)^\top R_t K'_t\| \right) \\ &\quad + \sum_{t=0}^{T-2} \|\mathcal{G}_t\| \|(K'_{t+1})^\top R_{t+1} K'_{t+1} - K_{t+1}^\top R_{t+1} K_{t+1}\| \\ 250 \quad (\text{SM3.30}) \quad &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} \left((2\rho + 1) \|B\| \sum_{t=0}^{T-2} \|K_t - K'_t\| \right) \left(\sum_{t=1}^{T-1} \|Q_t + (K'_t)^\top R_t K'_t - K_t^\top R_t K_t + K_t^\top R_t K_t\| \right) \\ &\quad + \frac{\rho^2(\rho^{2(T-1)} - 1)}{\rho^2 - 1} \sum_{t=1}^{T-1} \|(K'_t)^\top R_t K'_t - K_t^\top R_t K_t\| \\ &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|\mathbf{K}' - \mathbf{K}\| \left(\|\mathbf{Q}\| + \|\mathbf{K}\|^2 \|\mathbf{R}\| \right) \\ &\quad + \left(\frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|\mathbf{K}' - \mathbf{K}\| + \frac{\rho^2(\rho^{2(T-1)} - 1)}{\rho^2 - 1} \right) \sum_{t=1}^{T-1} \|(K'_t)^\top R_t K'_t - K_t^\top R_t K_t\|. \end{aligned}$$

where the second inequality holds by Lemma 3.13 and (SM3.9), and the third inequality holds by (SM3.6). For the first term in (III), we have

$$\begin{aligned}
 & \sum_{t=0}^{T-2} \text{Tr} \left(\sum_{s=1}^t D'_{t,s} W(D'_{t,s})^\top S'_{t+1} - D_{t,s} W D_{t,s}^\top S_{t+1} \right) \\
 &= \sum_{t=0}^{T-2} \text{Tr} \left(\sum_{s=1}^t D'_{t,s} W(D'_{t,s})^\top (S'_{t+1} - S_{t+1}) + (D'_{t,s} W(D'_{t,s})^\top - D_{t,s} W D_{t,s}^\top) S_{t+1} \right) \\
 &\leq \left(\sum_{t=0}^{T-2} \sum_{s=1}^t \text{Tr}(W) \|D'_{t,s}\|^2 \right) \left\| \sum_{t=1}^{T-1} (K'_t)^\top R_t K'_t - K_t^\top R_t K_t \right\| \\
 &\quad + \sum_{t=0}^{T-2} \left\| \sum_{s=1}^t D'_{t,s} W(D'_{t,s})^\top - D_{t,s} W D_{t,s}^\top \right\| \left(\sum_{t=1}^{T-1} \text{Tr}(Q_t) + \|K_t\|^2 \text{Tr}(R_t) \right) \\
 &\leq \text{Tr}(W) \frac{(T-1)(\rho^{2(T-1)} - 1)}{\rho^2 - 1} \left\| \sum_{t=1}^{T-1} (K'_t)^\top R_t K'_t - K_t^\top R_t K_t \right\| \\
 &\quad + T \frac{(\rho^{2T} - 1)}{\rho^2 - 1} (2\rho + 1) \|B\| \|W\| \|\mathbf{K}' - \mathbf{K}\| \left(\text{Tr} \left(\sum_{t=1}^{T-1} Q_t \right) + \|\mathbf{K}\|^2 \text{Tr} \left(\sum_{t=1}^{T-1} R_t \right) \right),
 \end{aligned}$$

where the last step holds by (3.20). The second term in (III) is bounded by

$$\sum_{t=0}^{T-2} \text{Tr} \left(W(S'_{t+1} - S_{t+1}) \right) \leq \text{Tr}(W) \sum_{t=1}^{T-1} \left\| (K'_t)^\top R_t K'_t - K_t^\top R_t K_t \right\|.$$

Similarly, by (3.20) and (SM3.9), (IV) is bounded by

$$\begin{aligned}
 (IV) &\leq \text{Tr}(\mathbb{E}[x_0 x_0^\top]) \sum_{t=0}^{T-1} \left\| (\mathcal{G}'_t - \mathcal{G}_t)(Q_T) \right\| + \text{Tr} \left(\sum_{s=1}^{T-1} D'_{T-1,s} W(D'_{T-1,s})^\top Q_T - D_{T-1,s} W D_{T-1,s}^\top Q_T \right) \\
 &\leq \text{Tr}(\mathbb{E}[x_0 x_0^\top]) \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|Q_T\| \|\mathbf{K}' - \mathbf{K}\| + \text{Tr}(Q_T) \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|W\| \|\mathbf{K}' - \mathbf{K}\|.
 \end{aligned}$$

Now we bound the term $\sum_{t=1}^{T-1} \left\| (K'_t)^\top R_t K'_t - K_t^\top R_t K_t \right\|$, which appears several times in previous inequalities:

$$\begin{aligned}
 \sum_{t=1}^{T-1} \left\| (K'_t)^\top R_t K'_t - K_t^\top R_t K_t \right\| &= \sum_{t=1}^{T-1} \left\| (K'_t - K_t + K_t)^\top R_t (K'_t - K_t + K_t) - K_t^\top R_t K_t \right\| \\
 &\leq \sum_{t=1}^{T-1} \|K'_t - K_t\|^2 \|R_t\| + 2 \|K_t\| \|R_t\| \|K'_t - K_t\| \leq 3 \|\mathbf{K}\| \|\mathbf{R}\| \|\mathbf{K}' - \mathbf{K}\|.
 \end{aligned}$$

The last step holds since $\|K'_t - K_t\| \leq \|K_t\|$ by assumption.

268 Therefore,

$$\begin{aligned}
|C(\mathbf{K}') - C(\mathbf{K})| &\leq \text{Tr}(\mathbb{E}[x_0 x_0^\top]) \left\{ 3 \|\mathbf{K}\| \|R_0\| \|\mathbf{K}' - \mathbf{K}\| + \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|Q_T\| \|\mathbf{K}' - \mathbf{K}\| \right. \\
&\quad + \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|\mathbf{K}' - \mathbf{K}\| \left(\|\mathbf{Q}\| + \|\mathbf{K}\|^2 \|\mathbf{R}\| \right) \\
&\quad + \left(\frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|\mathbf{K}' - \mathbf{K}\| + \frac{\rho^2(1 - \rho^{2(T-1)})}{\rho^2 - 1} \right) 3 \|\mathbf{K}\| \|\mathbf{R}\| \|\mathbf{K}' - \mathbf{K}\| \Big\} \\
&\quad + 3 \text{Tr}(W) \left(\frac{(T-1)(\rho^{2(T-1)} - 1)}{\rho^2 - 1} + 1 \right) \|\mathbf{K}\| \|\mathbf{R}\| \|\mathbf{K}' - \mathbf{K}\| \\
&\quad + \left(T \frac{(\rho^{2T} - 1)}{\rho^2 - 1} (2\rho + 1) \|B\| \|W\| \|\mathbf{K}' - \mathbf{K}\| \right) \left(\text{Tr} \left(\sum_{t=1}^{T-1} Q_t \right) + \|\mathbf{K}\|^2 \text{Tr} \left(\sum_{t=1}^{T-1} R_t \right) \right) \\
&\quad + \text{Tr}(Q_T) \frac{\rho^{2T} - 1}{\rho^2 - 1} (2\rho + 1) \|B\| \|W\| \|\mathbf{K}' - \mathbf{K}\|.
\end{aligned}$$

270 By (3.27), Lemma 3.8, and Lemma 3.16, ρ is bounded above by polynomials in $\|A\|$, $\|B\|$, $\|\mathbf{R}\|$, $\frac{1}{\sigma_{\mathbf{X}}}$, $\frac{1}{\sigma_{\mathbf{R}}}$ and
 271 $C(\mathbf{K})$, or a constant $1 + \xi$. Therefore, we rewrite the above inequality by

$$272 \quad (\text{SM3.31}) \quad |C(\mathbf{K}') - C(\mathbf{K})| \leq h_{CK} \|\mathbf{K}' - \mathbf{K}\| + h'_{CK} \|\mathbf{K}' - \mathbf{K}\|^2,$$

273 where $h_{CK} \in \mathcal{H}(C(\mathbf{K}))$ and $h'_{CK} \in \mathcal{H}(C(\mathbf{K}))$ are polynomials in $C(\mathbf{K})$ and model parameters. Given assumption
 274 (4.7), we have $\|\mathbf{K}' - \mathbf{K}\| \leq 1$ and hence

$$275 \quad \|\mathbf{K}' - \mathbf{K}\| \geq \|\mathbf{K}' - \mathbf{K}\|^2.$$

276 Define $h_{cost} = h_{CK} + h'_{CK}$, then (SM3.31) gives

$$277 \quad |C(\mathbf{K}') - C(\mathbf{K})| \leq h_{cost} \|\mathbf{K}' - \mathbf{K}\|,$$

278 with $h_{cost} \in \mathcal{H}(C(\mathbf{K}))$. □

279 **Proof of Lemma 4.8.** Recall $\nabla_t C(\mathbf{K}) = 2E_t \Sigma_t$ and $W = \mathbb{E}[w_t w_t^\top] = \widetilde{W} \mathbb{E}[v_t v_t^\top] \widetilde{W}^\top$. We have,

$$280 \quad (\text{SM3.32}) \quad \|\nabla_t C(\mathbf{K}') - \nabla_t C(\mathbf{K})\| = \|2E'_t \Sigma'_t - 2E_t \Sigma_t\| \leq 2\|E'_t - E_t\| \|\Sigma'_t\| + 2\|E_t\| \|\Sigma'_t - \Sigma_t\|,$$

281 For the second term, by Lemma 3.6 and Cauchy-Schwarz inequality,

$$282 \quad (\text{SM3.33}) \quad \|E_t\| \leq \sum_{t=0}^{T-1} \|E_t\| \leq \sum_{t=0}^{T-1} \sqrt{\text{Tr}(E_t^\top E_t)} \leq \sqrt{T \cdot \frac{\max_t \|R_t + B^\top P_{t+1} B\|}{\sigma_{\mathbf{X}}}} (C(\mathbf{K}) - C(\mathbf{K}^*)).$$

283 By (SM3.7) and direct calculation, we have

$$284 \quad \|(\mathcal{G}'_{t+1} - \mathcal{G}_{t+1})(\Sigma_0)\| \leq \rho^{2(t+1)} \left(\sum_{s=0}^{t+1} \|\mathcal{F}_{K'_s} - \mathcal{F}_{K_s}\| \|\Sigma_0\| \right).$$

285 By (SM3.9) and (3.20), for $t = 1, 2, \dots, T-1$,

$$\begin{aligned}
286 \quad (\text{SM3.34}) \quad \|\Sigma'_t - \Sigma_t\| &\leq \|(\mathcal{G}'_t - \mathcal{G}_t)(\Sigma_0)\| + \left\| \sum_{s=0}^{t-1} D_{t-1,s} W D_{t-1,s}^\top - D'_{t-1,s} W (D'_{t-1,s})^\top \right\| \\
&\leq \rho^{2t} (2\rho + 1) \|B\| \|\Sigma_0\| \|\mathbf{K}' - \mathbf{K}\| + \frac{(\rho^{2T} - 1)}{\rho^2 - 1} (2\rho + 1) \|B\| \|W\| \|\mathbf{K}' - \mathbf{K}\|.
\end{aligned}$$

287 Therefore the second term in (SM3.32) is bounded by the product of (SM3.33) and (SM3.34).
 288

Next we bound the first term in (SM3.32). Similar to (SM3.14), $\|\Sigma'_t\| \leq \|\sum_{t=0}^T \Sigma'_t\| = \|\Sigma_{\mathbf{K}'}\| \leq \|\Sigma'_{\mathbf{K}} - \Sigma_{\mathbf{K}}\| + \|\Sigma_{\mathbf{K}}\| \leq \frac{C(\mathbf{K})}{\sigma_Q} + \|\Sigma_{\mathbf{K}}\|$. For $\|E'_t - E_t\|$, we first need a bound on $\|P'_t - P_t\|$. Since $P_0 = S_0 + \sum_{t=0}^{T-2} \mathcal{G}_t(S_{t+1}) + \mathcal{G}_{T-1}(Q_T)$, by (SM3.30), we have

(SM3.35)

$$\begin{aligned} \|P'_t - P_t\| &\leq \|P'_0 - P_0\| \leq 3\|K_0\|\|R_0\|\|K'_0 - K_0\| + \|\mathcal{G}_d\| + \frac{\rho^{2T} - 1}{\rho^2 - 1}(2\rho + 1)\|B\|\|Q_T\| \left(\sum_{t=0}^{T-1} \|K_t - K'_t\| \right) \\ &\leq \frac{\rho^{2T} - 1}{\rho^2 - 1}(2\rho + 1)\|B\|\|\mathbf{K}' - \mathbf{K}\| \left(\|\mathbf{Q}\| + \|\mathbf{K}\|^2\|R\| \right) \\ &\quad + 3 \left(1 + \frac{\rho^{2T} - 1}{\rho^2 - 1}(2\rho + 1)\|B\|\|\mathbf{K}' - \mathbf{K}\| + \frac{\rho^2(1 - \rho^{2(T-1)})}{\rho^2 - 1} \right) \cdot \|\mathbf{K}\|\|R\|\|\mathbf{K}' - \mathbf{K}\| \\ &\quad + \frac{\rho^{2T} - 1}{\rho^2 - 1}(2\rho + 1)\|B\|\|Q_T\|\|\mathbf{K}' - \mathbf{K}\|. \end{aligned}$$

Thus,

$$\begin{aligned} \|E'_t - E_t\| &= \|R_t(K'_t - K_t) - B^\top(P'_{t+1} - P_{t+1})A + B^\top(P'_{t+1} - P_{t+1})BK'_t + B^\top P_{t+1}B(K'_t - K_t)\| \\ &\leq (\|R_t\| + \|B\|^2\|P_0\|)\|\mathbf{K}' - \mathbf{K}\| + \|B\|\|P'_0 - P_0\|\|A\| + 2\|B\|^2\|P'_0 - P_0\|\|\mathbf{K}\|. \end{aligned}$$

Given the bound on $\|\mathbf{K}\| = \sum_{t=0}^{T-1} \|K_t\|$ in Lemma 3.16 and the bound on $\|P_t\|$ in Lemma 3.8, all the terms in (SM3.32) can be bounded by polynomials of related parameters multiplied by $\|\mathbf{K}' - \mathbf{K}\|$ and $\|\mathbf{K}' - \mathbf{K}\|^2$. Similarly to the proof of Lemma 4.7, we have $\|\mathbf{K}' - \mathbf{K}\| \leq 1$ and

$$\|\nabla_t C(\mathbf{K}') - \nabla_t C(\mathbf{K})\| \leq h_{grad}\|\mathbf{K}' - \mathbf{K}\|, \quad \square$$

for some polynomial $h_{grad} \in \mathcal{H}(C(\mathbf{K}))$.

SM3.5. Proofs in Section 5.

Proof of Proposition 5.2. Denote $H_t := \begin{pmatrix} 1 + \gamma k_t^1 & \gamma k_t^2 \\ k_t^1 & 1 + k_t^2 \end{pmatrix}$. Since H_t has two eigenvalues 1 and $\gamma k_t^1 + k_t^2 + 1$, H_t is positive definite when $\gamma k_t^1 + k_t^2 > -1$ ($0 \leq t \leq T-1$).

Then let us show the first claim by induction. Assume $\mathbb{E}[x_s x_s^\top]$ is positive definite for all $s \leq t$, then

$$\begin{aligned} \mathbb{E}[x_{t+1} x_{t+1}^\top] &= \mathbb{E}[(A - BK_t)x_t + w_t][(A - BK_t)x_t + w_t]^\top = \mathbb{E}[(H_t x_t + w_t)(H_t x_t + w_t)^\top] \\ &= \mathbb{E}[H_t x_t x_t^\top H_t^\top + w_t w_t^\top + w_t w_t^\top + 2H_t x_t w_t^\top] = H_t \mathbb{E}[x_t x_t^\top] H_t^\top + \begin{pmatrix} \sigma & 0 \\ 0 & 0 \end{pmatrix}. \end{aligned}$$

Hence $\mathbb{E}[x_{t+1} x_{t+1}^\top]$ is positive definite since $\mathbb{E}[x_t x_t^\top]$ is positive definite and H_t is positive definite. Therefore $\underline{\sigma}_{\mathbf{X}} > 0$.

The second claim can be proved by backward induction. For $t = T$, $P_T^K = Q_T$ is positive definite since Q_T is positive definite. Assume P_{t+1}^K is positive definite for some $t+1$, then take any $z \in \mathbb{R}^d$ such that $z \neq 0$,

$$z^\top P_t^K z = z^\top Q_t z + z^\top K_t^\top R_t K_t z + z^\top H_t^\top P_{t+1}^K H_t z > 0.$$

Note that H_t is positive definite when $\gamma k_t^1 + k_t^2 > -1$ and $1 + \gamma k_t^1 > 0$. The last inequality holds since Q_t and $H_t^\top P_{t+1}^K H_t$ are positive definite, and $K_t^\top R_t K_t$ is positive semi-definite. Hence we have P_t^K positive definite for all $t = 0, 1, 2, \dots, T$. \square

References.

- [1] R. CONT, A. KUKANOV, AND S. STOIKOV, *The price impact of order book events*, Journal of Financial Econometrics, 12 (2014), pp. 47–88.
- [2] S. TU AND B. RECHT, *Least-squares temporal difference learning for the linear quadratic regulator*, in International Conference on Machine Learning, 2018, pp. 5005–5014.