

- Jazwinski, A. Stochastic process and filtering theory, academic press. *A subsidiary of Harcourt Brace Jovanovich Publishers*, 1970.
- Kamalapurkar, R., Klotz, J. R., Walters, P., and Dixon, W. E. Model-based reinforcement learning in differential graphical games. *IEEE Transactions on Control of Network Systems*, 5(1):423–433, 2016.
- Karatzas, I. and Shreve, S. Brownian motion and stochastic calculus (graduate texts in mathematics), 1991.
- Kushner, H. Numerical approximations for stochastic differential games. *SIAM J. Control Optim.*, 41:457–486, 2002.
- Kushner, H. and Chamberlain, S. On stochastic differential games: Sufficient conditions that a given strategy be a saddle point, and numerical procedures for the solution of the game. *Journal of Mathematical Analysis and Applications*, 26:560–575, 1969.
- Liu, I.-J., Yeh, R. A., and Schwing, A. G. Pic: permutation invariant critic for multi-agent deep reinforcement learning. In *Conference on Robot Learning*, pp. 590–602. PMLR, 2020.
- Lyle, C., van der Wilk, M., Kwiatkowska, M., Gal, Y., and Bloem-Reddy, B. On the benefits of invariance in neural networks. *arXiv preprint arXiv:2005.00178*, 2020.
- Ma, J., Zhang, J., et al. Representation theorems for backward stochastic differential equations. *Annals of Applied probability*, 12(4):1390–1418, 2002.
- Mahajan, A., Rashid, T., Samvelyan, M., and Whiteson, S. Maven: Multi-agent variational exploration. In *Advances in Neural Information Processing Systems*, pp. 7613–7624, 2019.
- Mataramvura, S. and Øksendal, B. Risk minimizing portfolios and hjbi equations for stochastic differential games. *Stochastics An International Journal of Probability and Stochastic Processes*, 80(4):317–337, 2008.
- Osborne, M. J. and Rubinstein, A. A course in game theory cambridge. MA: MIT Press [Google Scholar], 1994.
- Pereira, M. A., Wang, Z., Exarchos, I., and Theodorou, E. A. Learning deep stochastic optimal control policies using forward-backward sdes. In *Robotics: science and systems*, 2019.
- Prasad, A. and Sethi, S. P. Competitive advertising under uncertainty: A stochastic differential game approach. *Journal of Optimization Theory and Applications*, 123(1): 163–185, 2004.
- Raissi, M. Forward-backward stochastic neural networks: Deep learning of high-dimensional partial differential equations. *arXiv preprint arXiv:1804.07010*, 2018.
- Ramachandran, K. M. and Tsokos, C. P. *Stochastic differential games. Theory and applications*, volume 2. Springer Science & Business Media, 2012.
- Rashid, T., Samvelyan, M., De Witt, C. S., Farquhar, G., Foerster, J., and Whiteson, S. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. *arXiv preprint arXiv:1803.11485*, 2018.
- Rozen, O., Shwartz, V., Aharoni, R., and Dagan, I. Diversify your datasets: Analyzing generalization via controlled variance in adversarial datasets. *arXiv preprint arXiv:1910.09302*, 2019.
- Schmidt, L., Santurkar, S., Tsipras, D., Talwar, K., and Madry, A. Adversarially robust generalization requires more data. In *Advances in Neural Information Processing Systems*, pp. 5014–5026, 2018.
- Shi, G., Hönig, W., Shi, X., Yue, Y., and Chung, S.-J. Neural-swarm2: Planning and control of heterogeneous multi-rotor swarms using learned interactions. *arXiv preprint arXiv:2012.05457*, 2020.
- Simon, D. *Optimal state estimation: Kalman, H infinity, and nonlinear approaches*. John Wiley & Sons, 2006.
- Son, K., Kim, D., Kang, W. J., Hostallero, D. E., and Yi, Y. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:1905.05408*, 2019.
- Sunehag, P., Lever, G., Gruslys, A., Czarnecki, W. M., Zambaldi, V. F., Jaderberg, M., Lanctot, M., Sonnerat, N., Leibo, J. Z., Tuyls, K., et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *AAMAS*, pp. 2085–2087, 2018.
- Tan, M. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the tenth international conference on machine learning*, pp. 330–337, 1993.
- Theodorou, E., Tassa, Y., and Todorov, E. Stochastic differential dynamic programming. In *Proceedings of the 2010 American Control Conference*, pp. 1125–1132. IEEE, 2010.
- Wang, L., Yang, Z., and Wang, Z. Breaking the curse of many agents: Provable mean embedding q-iteration for mean-field reinforcement learning. *arXiv preprint arXiv:2006.11917*, 2020.

- Wang, Z., Lee, K., Pereira, M. A., Exarchos, I., and Theodorou, E. A. Deep forward-backward SDEs for min-max control. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 6807–6814. IEEE, 2019a.
- Wang, Z., Pereira, M. A., and Theodorou, E. A. Deep 2fbsdes for systems with control multiplicative noise. *arXiv preprint arXiv:1906.04762*, 2019b.
- Yang, Y., Tutunov, R., Sakulwongtana, P., Ammar, H. B., and Wang, J. α^α -rank: Scalable multi-agent evaluation through evolution. *arXiv preprint arXiv:1909.11628*, 2019.
- Yong, J. and Zhou, X. Y. *Stochastic controls: Hamiltonian systems and HJB equations*, volume 43. Springer Science & Business Media, 1999.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Póczos, B., Salakhutdinov, R. R., and Smola, A. J. Deep sets. In *Advances in neural information processing systems*, pp. 3391–3401, 2017.
- Zhou, M., Chen, Y., Wen, Y., Yang, Y., Su, Y., Zhang, W., Zhang, D., and Wang, J. Factorized q-learning for large-scale multi-agent systems. In *Proceedings of the First International Conference on Distributed Artificial Intelligence*, pp. 1–7, 2019.

Supplementary Material

A. Multi-agent HJB Derivation

Applying Bellman's principle to the value function (3) yields,

$$\begin{aligned}
 V^i(t, \mathbf{X}(t)) &= \inf_{U_i \in \mathcal{U}_i} \mathbb{E} \left[V^i(t+dt, \mathbf{X}(t+dt)) + \int_t^{t+dt} C^i d\tau \right] \\
 &= \inf_{U_i \in \mathcal{U}_i} \mathbb{E} \left[C^i dt + V^i(t, \mathbf{X}(t)) + V_t^i(t, \mathbf{X}(t))dt \right. \\
 &\quad \left. + V_x^{iT}(t, \mathbf{X}(t))d\mathbf{X} + \frac{1}{2}\text{tr}(V_{xx}^i(t, \mathbf{X}(t))\Sigma\Sigma^T)dt \right] \\
 &= \inf_{U_i \in \mathcal{U}_i} \mathbb{E} \left[C^i dt + V^i(t, \mathbf{X}(t)) + V_t^i(t, \mathbf{X}(t))dt \right. \\
 &\quad \left. + V_x^{iT}(t, \mathbf{X}(t))((f + G\mathbf{U})dt + \Sigma d\mathbf{W}) + \frac{1}{2}\text{tr}(V_{xx}^i(t, \mathbf{X}(t))\Sigma\Sigma^T)dt \right] \\
 &= \inf_{U_i \in \mathcal{U}_i} \left[C^i dt + V^i(t, \mathbf{X}(t)) + V_t^i(t, \mathbf{X}(t))dt \right. \\
 &\quad \left. + V_x^{iT}(t, \mathbf{X}(t))((f + G\mathbf{U})dt) + \frac{1}{2}\text{tr}(V_{xx}^i(t, \mathbf{X}(t))\Sigma\Sigma^T)dt \right] \\
 &\Rightarrow 0 = V_t^i(t, \mathbf{X}(t)) + \inf_{U_i \in \mathcal{U}_i} \left[C^i + V_x^{iT}(t, \mathbf{X}(t))(f + G\mathbf{U}) \right] + \frac{1}{2}\text{tr}(V_{xx}^i(t, \mathbf{X}(t))\Sigma\Sigma^T)
 \end{aligned} \tag{20}$$

Given the cost function assumption (the cost function is quadratic w.r.t control variable.), the infimum can be obtained explicitly using optimal control $U_i^* = -R^{-1}(G_i^T V_x^i + Q_i^T \mathbf{X})$. With that we can obtain the final form of the HJB PDE as

$$V_t^i + h + V_x^{iT}(f + G\mathbf{U}_{0,*}) + \frac{1}{2}\text{tr}(V_{xx}^i \Sigma \Sigma^T) = 0, \quad V^i(T, \mathbf{X}) = g(\mathbf{X}(T)). \tag{21}$$

B. Multi-agent FBSDE Derivation

Given the HJB PDE in equation 5, one can apply the non-linear Feynman-Kac lemma (Karatzas & Shreve, 1991) to obtain a set of FBSDE as

$$\begin{aligned}
 d\mathbf{X}(t) &= (f + G\mathbf{U}_{0,*})dt + \Sigma d\mathbf{W}, \quad \mathbf{X}(0) = \mathbf{x}_0 \quad (\text{FSDE}) \\
 dV^i &= -h dt + V_x^{iT} \Sigma d\mathbf{W}, \quad V(\mathbf{X}(T)) = g(\mathbf{X}(T)). \quad (\text{BSDE})
 \end{aligned} \tag{22}$$

The backward process is derived by applying Ito's lemma on V^i

$$\begin{aligned}
 dV^i &= V_t^i dt + V_x^{iT} d\mathbf{X} + \frac{1}{2}\text{tr}(V_{xx}^i \Sigma \Sigma^T)dt \\
 &\text{Plug in eq.21 into } V_t^i \text{ term, and eq.1 into } d\mathbf{X} \text{ term} \\
 &= (-h - V_x^{iT}(f + G\mathbf{U}_{0,*}) - \frac{1}{2}\text{tr}(V_{xx}^i \Sigma \Sigma^T))dt + V_x^{iT}((f + G\mathbf{U}_{0,*})dt + \Sigma d\mathbf{W}) + \frac{1}{2}\text{tr}(V_{xx}^i \Sigma \Sigma^T)dt \\
 &= -h dt + V_x^{iT} \Sigma d\mathbf{W}.
 \end{aligned}$$

C. Metrics

In the test set, we randomly select B initial states, and $B \times T$ noise \mathbf{W} , where B is Batch size and T is time horizon. We evaluate the performance of models based on three different metrics. All losses (for comparison) are computed by averaging the last 10 stages and over 3 random seeds for fair comparison.

C.1. Relative Square Error(RSE)

The RSE is a metric applied in test phase which is defined as following:

$$\mathcal{L}_{RSE} = \frac{\sum_{1 \leq j \leq B}^{i \in \mathbb{I}} (\hat{Y}^i(0, \mathbf{X}^j(0)) - Y^i(0, \mathbf{X}^j(0)))^2}{\sum_{1 \leq j \leq B}^{i \in \mathbb{I}} (\hat{Y}^i(0, \mathbf{X}^j(0)) - \bar{Y}^i(0, \mathbf{X}^j(0)))^2}, \quad (23)$$

Where Y^i is the analytical solution of value function for i th agents at initial state $\mathbf{X}^j(0)$. The initial state $\mathbf{X}^j(0)$ is new batch of data sampled from same distribution as $\mathbf{X}(0)$ in the training phase. The batch size B is 256 for all inter-bank simulations. \hat{Y}^i is the approximated value function for i th agent by FBSDE controller, and \bar{Y}^i is the average of analytical solution for i th agent over the entire batch.

C.2. Evaluation/training Loss

The evaluation loss is same as training loss which is defined as the mean square error between true terminal value evaluated on the terminal state and the value propagated by the BSDE,

$$\mathcal{L}(\hat{Y}_T^i, Y_T^i) = \frac{1}{B} \|\hat{Y}_T^i - Y_T^i\|_2^2 \quad (24)$$

C.3. Cumulative Loss

The cumulative loss is computed explicitly from the objective function of optimal control (2) in the test phase. Here we restate it for completeness.

$$\begin{aligned} \mathcal{L}_{cum} &:= J_t^i(\mathbf{X}, U_{i,m}; \mathbf{U}_{-i,m-1}) \\ &= \mathbb{E} \left[g(\mathbf{X}_T) + \int_0^T C^i(\mathbf{X}_\tau, U_i(\mathbf{X}_\tau); \mathbf{U}_{-i}) d\tau \right], \end{aligned} \quad (25)$$

D. Missing Derivation and Proof in Section 3

D.1. Assumption 2

Here we state the assumption 2 in detail.

Consider a general FBSDE system,

$$\begin{aligned} \mathbf{X}_T^{t,\mathbf{x}} &= \mathbf{x} + \int_t^T \mu_s ds + \int_t^T \Sigma_s d\mathbf{W}_s \quad (\text{FSDE}), \\ Y_t^{T,\mathbf{x}} &= g(\mathbf{X}_T^{t,\mathbf{x}}) - \int_t^T H_s ds + \int_t^T Z_s d\mathbf{W}_s \quad (\text{BSDE}), \end{aligned} \quad (26)$$

We use $|\cdot|$ and $\|\cdot\|_F$ to denote the L^2 norm and Frobenius norm respectively. For terminal loss $g(\cdot)$, drift function $\mu(\cdot, \cdot, \cdot)$, $H(\cdot, \cdot, \cdot)$, diffusion function $\Sigma(\cdot, \cdot)$, and control function $\mathbf{u}(\cdot, \cdot)$, they are Lipschitz continuous with respect to their arguments:

$$\begin{aligned} |g(t, \mathbf{x}_1) - g(t, \mathbf{x}_2)|^2 &\leq g_x |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ |H(t, \mathbf{x}_1, \mathbf{z}_1) - H(t, \mathbf{x}_2, \mathbf{z}_2)|^2 &\leq H_x |\mathbf{x}_1 - \mathbf{x}_2|^2 + H_z \|\mathbf{z}_1 - \mathbf{z}_2\|_F^2, \\ \|\Sigma(t, \mathbf{x}_1) - \Sigma(t, \mathbf{x}_2)\|_F^2 &\leq \Sigma_x |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ \|\mathbf{u}(t, \mathbf{x}_1) - \mathbf{u}(t, \mathbf{x}_2)\|^2 &\leq u_x |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ |\mu(t, \mathbf{x}_1, \mathbf{u}_1) - \mu(t, \mathbf{x}_2, \mathbf{u}_2)|^2 &\leq \mu_x |\mathbf{x}_1 - \mathbf{x}_2|^2 + \mu_u |\mathbf{u}_1 - \mathbf{u}_2|^2, \end{aligned} \quad (27)$$

Here $\mu_x, \mu_u, g_x, H_x, H_z, \Sigma_x, u_x$ are all positive constants.

D.2. Proof of Lemma 1

Proof. Denote $(\mathbf{X}_s^{t,\mathbf{x}}, Y_s^{t,\mathbf{x}}, Z_s^{t,\mathbf{x}})_{t \leq s \leq T}$ as the solution for the FBSDE system for the i th agent:

$$\begin{aligned} \mathbf{X}_T^{t,\mathbf{x}} &= \mathbf{x} + \int_t^T \mu_s ds + \int_t^T \Sigma_s dW_s & (\text{FSDE}), \\ Y_t^{T,\mathbf{x}} &= g(\mathbf{X}_T^{t,\mathbf{x}}) - \int_t^T H_s ds + \int_t^T Z_s dW_s & (\text{BSDE}), \end{aligned} \quad (28)$$

for any $(t_0, \mathbf{x}) \in [t_0, T] \times \mathcal{X}$. Then for any $t \in [t_0, T]$ and $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$, let $(\mathbf{X}_t^j, Y_t^j, Z_t^j)$ be the short notation of $(\mathbf{X}_t^{t_0, \mathbf{x}_j}, Y_t^{t_0, \mathbf{x}_j}, Z_t^{t_0, \mathbf{x}_j})$, where $j \in \{1, 2\}$, $t \in [0, T]$.

Here we define $\delta \mathbf{X}_t, \delta Y_t, \delta Z_t, \delta \Sigma_t, \delta H_t$ as:

$$\begin{aligned} \delta \mathbf{X}_t &= \mathbf{X}_t^1 - \mathbf{X}_t^2, \\ \delta Y_t &= Y_t^1 - Y_t^2, \\ \delta Z_t &= Z_t^1 - Z_t^2, \\ \delta \mu_t &= \mu_t(t, \mathbf{X}_t^1, \mathbf{U}_t^1) - \mu_t(t, \mathbf{X}_t^2, \mathbf{U}_t^2), \\ \delta H_t &= H(t, \mathbf{X}_t^1, Z_t^1) - H(t, \mathbf{X}_t^2, Z_t^2), \\ \delta \Sigma_t &= \Sigma(t, \mathbf{X}_t^1) - \Sigma(t, \mathbf{X}_t^2), \end{aligned} \quad (29)$$

Then we have:

$$\begin{aligned} d\delta \mathbf{X}_t &= d\mathbf{X}_t^1 - d\mathbf{X}_t^2 \\ &= (\mu_t(t, \mathbf{X}_t^1, \mathbf{U}_t^1)dt + \Sigma(t, \mathbf{X}_t^1)dW_t) - (\mu_t(t, \mathbf{X}_t^2, \mathbf{U}_t^2)dt + \Sigma(t, \mathbf{X}_t^2)dW_t) \\ &= (\delta \mu_t)dt + \delta \Sigma_t dW_t, \\ d\delta Y_t &= -\delta H_t dt + \delta Z_t^\top dW_t. \end{aligned} \quad (30)$$

By applying Itô lemma to $d\delta \mathbf{X}_t$ and $d\delta Y_t$,

$$\begin{aligned} d|\delta \mathbf{X}_t|^2 &= (2\delta \mu_t \delta \mathbf{X}_t + \frac{1}{2} \cdot 2\delta \|\Sigma_t\|_F^2)dt + 2(\delta \mathbf{X}_t)^\top \delta \Sigma_t dW_t \\ &= (2\delta \mu_t \delta \mathbf{X}_t + \delta \|\Sigma_t\|_F^2)dt + 2(\delta \mathbf{X}_t)^\top \delta \Sigma_t dW_t \\ d|\delta Y_t|^2 &= (-2\delta H_t \delta Y_t + \delta \|Z_t^\top\|_F^2)dt + 2(\delta Z_t \delta Y_t)^\top dW_t \end{aligned} \quad (31)$$

By taking the expectation on both sides of $d|\delta \mathbf{X}_t|^2$, it will yield:

$$\begin{aligned} \mathbb{E}[|\delta \mathbf{X}_t|^2] &= |\mathbf{x}_1 - \mathbf{x}_2|^2 + \int_{t_0}^t \mathbb{E} [2\delta \mu_s \delta \mathbf{X}_s + \delta \|\Sigma_s\|_F^2] ds \\ &\leq |\mathbf{x}_1 - \mathbf{x}_2|^2 + \int_{t_0}^t \mathbb{E} [(\mu_x + \mu_u u_x)^{-1} |\delta \mu_s|^2 + (\mu_x + \mu_u u_x) |\delta \mathbf{X}_s|^2 + \delta \|\Sigma_s\|_F^2] ds \\ &\leq |\mathbf{x}_1 - \mathbf{x}_2|^2 + \int_{t_0}^t \mathbb{E} [(\mu_x + \mu_u u_x)^{-1} (\mu_x |\delta \mathbf{X}_s|^2 + \mu_u |\delta u|^2) + (\mu_x + \mu_u u_x) |\delta \mathbf{X}_s|^2 + \delta \|\Sigma_s\|_F^2] ds \\ &\leq |\mathbf{x}_1 - \mathbf{x}_2|^2 + \int_{t_0}^t \mathbb{E} [(\mu_x + \mu_u u_x)^{-1} (\mu_x |\delta \mathbf{X}_s|^2 + \mu_u u_x |\delta \mathbf{X}_s|^2) + (\mu_x + \mu_u u_x) |\delta \mathbf{X}_s|^2 + \delta \|\Sigma_s\|_F^2] ds \\ &\leq |\mathbf{x}_1 - \mathbf{x}_2|^2 + \int_{t_0}^t \mathbb{E} [|\delta \mathbf{X}_s|^2 + (\mu_x + \mu_u u_x) |\delta \mathbf{X}_s|^2 + \Sigma_x |\delta \mathbf{X}_s|^2] ds \\ &= |\mathbf{x}_1 - \mathbf{x}_2|^2 + (I + \mu_x + \mu_u u_x + \Sigma_x) \int_{t_0}^t \mathbb{E} |\delta \mathbf{X}_s|^2 ds \end{aligned}$$

by Gronwall's inequality

$$\begin{aligned} &\leq e^{(I + \mu_x + \mu_u u_x + \Sigma_x)(t - t_0)} |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= e^{\xi(t - t_0)} |\mathbf{x}_1 - \mathbf{x}_2|^2 \end{aligned} \quad (32)$$

Where $\xi = I + \mu_x + \mu_u u_x + \Sigma_x$.

Similarly, we can have,

$$\begin{aligned}
 \mathbb{E}[|\delta \mathbf{Y}_t|^2] &= \mathbb{E}|\delta \mathbf{Y}_T|^2 + \int_t^T \mathbb{E} [2\delta H_s \delta Y_s - \delta \|\mathbf{Z}_s^T\|_F^2] ds \\
 &= \mathbb{E}|\mathbf{g}(T, \mathbf{x}_1) - \mathbf{g}(T, \mathbf{x}_2)|^2 + \int_t^T \mathbb{E} [2\delta H_s \delta Y_s] - \mathbb{E}\|\mathbf{Z}_s\|_F^2 ds \\
 &\leq g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T H_z \mathbb{E}|\delta Y_s|^2 + H_z^{-1} \mathbb{E}|\delta H_s|^2 - \mathbb{E}\|\mathbf{Z}_s\|_F^2 ds \\
 &\leq g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T H_z \mathbb{E}|\delta Y_s|^2 + H_z^{-1} \mathbb{E} [H_x |\delta \mathbf{X}_s|^2 + H_z |\delta \mathbf{Z}_s|^2] - \mathbb{E}\|\mathbf{Z}_s\|_F^2 ds \\
 &= g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T H_z \mathbb{E}|\delta Y_s|^2 + H_z^{-1} H_x \mathbb{E}|\delta \mathbf{X}_s|^2 ds \\
 &\leq \left[g_x e^{\xi(t-t_0)} + H_x \frac{e^{\xi(T-t)} - e^{\xi(t-t_0)}}{H_z \Sigma_x} \right] |\mathbf{x}_1 - \mathbf{x}_2|^2 + H_z \int_t^T \mathbb{E}|\delta Y_s|^2 ds
 \end{aligned} \tag{33}$$

by Gronwall's inequality

$$|\delta \mathbf{Y}_t|^2 \leq e^{H_z(T-t)} \left[g_x e^{\xi(T-t_0)} + H_x \frac{e^{\xi(T-t_0)} - e^{\xi(t-t_0)}}{H_z \Sigma_x} \right] |\mathbf{x}_1 - \mathbf{x}_2|^2$$

When $t_0 = 0$, one can have:

$$\begin{aligned}
 |\delta Y_T|^2 &\leq g_x e^{\xi T} |\mathbf{x}_1 - \mathbf{x}_2|^2 \\
 &= L_1 |\mathbf{x}_1 - \mathbf{x}_2|^2 \\
 |\delta Y_0|^2 &\leq e^{H_z T} \left[g_x e^{\xi T} + H_x \frac{e^{\xi T} - 1}{H_z \Sigma_x} \right] |\mathbf{x}_1 - \mathbf{x}_2|^2 \\
 &= L_2 |\mathbf{x}_1 - \mathbf{x}_2|^2
 \end{aligned} \tag{34}$$

Where

$$\begin{aligned}
 L_1 &= g_x e^{\xi T} \\
 L_2 &= e^{H_z T} \left[g_x e^{\xi T} + H_x \frac{e^{\xi T} - 1}{H_z \Sigma_x} \right] \\
 \xi &= I + \mu_x + \mu_u u_x + \Sigma_x
 \end{aligned} \tag{35}$$

□

D.3. Lemma.2 with Proof

Lemma 2. Denote $(\mathbf{X}_s^{t,\mathbf{x}}, Y_s^{t,\mathbf{x}}, Z_s^{t,\mathbf{x}})_{t \leq s \leq T}$ as the solution for the FBSDE system with importance sampling (12, 13) satisfying assumptions 1 and 2. Denote the difference of Y component at two different states \mathbf{x}_1 and \mathbf{x}_2 as:

$$\delta \mathbf{X}_t = \mathbf{X}_t^{t_0, \mathbf{x}_1} - \mathbf{X}_t^{t_0, \mathbf{x}_2}, \delta Y_t = Y_t^{t_0, \mathbf{x}_1} - Y_t^{t_0, \mathbf{x}_2}. \tag{36}$$

Then we can have:

$$\begin{aligned}
 |\delta Y_T|^2 &\leq \tilde{L}_1 |\mathbf{x}_1 - \mathbf{x}_2|^2, \\
 |\delta Y_{t_0}|^2 &\leq \tilde{L}_2 |\mathbf{x}_1 - \mathbf{x}_2|^2,
 \end{aligned} \tag{37}$$

Where L_1 and L_2 are defined as:

$$\begin{aligned}
 \tilde{L}_1 &= g_x e^{\tilde{\xi}} \\
 \tilde{L}_2 &= e^{2(H_z + k_z)(T-t_0)} \left[g_x e^{\tilde{\xi}(T-t_0)} + H_x (H_z^{-1} + k_z^{-1}) \frac{e^{\tilde{\xi}(T-t_0)} - 1}{2\Sigma_x} \right], \\
 \tilde{\xi} &= I + \tilde{\mu}_x + \tilde{\mu}_u u_x + \Sigma_x,
 \end{aligned} \tag{38}$$

Where $\mu_x, \mu_u, \Sigma_x, H_x, H_z, g_x, u_x$ are Lipschitz constant defined in Assumption.2. The definition of Lipschitz constant $\tilde{\mu}_x, \tilde{\mu}_u, k_z$ and proof can be found in the following proof.

Proof. Similar to the proof of lemma.1, now we first analyze the forward process with IS. Inspired by the success of (Exarchos & Theodorou, 2018), we select the control computed from the last run as the importance sampling term. Then the IS term in FSDE is defined as $m_s := \Sigma K = GU_{*,0}$. New drift term is modified as $\tilde{\mu}_s = \mu_s + m_s$ with Lipschitz constant $\tilde{\mu}_x$ and $\tilde{\mu}_u$. In the BSDE, Then IS term is written as $k_s = Z_s K_s = Z_s \Gamma U_{*,0} = V_x G U_{*,0}$, and the modified $\tilde{H}_s = H_s + k_s$. Here we formally write the Lipschitz constant for IS terms in FSDE and BSDE are:

$$\begin{aligned} |m_s(t, \mathbf{x}_1, \mathbf{u}_1) - m_s(t, \mathbf{x}_2, \mathbf{u}_2)|^2 &\leq m_x |\mathbf{x}_1 - \mathbf{x}_2|^2 + m_u |\mathbf{u}_1 - \mathbf{u}_2|^2, \\ |k_s(t, \mathbf{x}_1, \mathbf{z}_1) - k_s(t, \mathbf{x}_2, \mathbf{z}_2)|^2 &\leq k_x |\mathbf{x}_1 - \mathbf{x}_2|^2 + k_z |\mathbf{z}_1 - \mathbf{z}_2|^2, \end{aligned} \quad (39)$$

Similar to the proof (D.2), we can have,

$$\begin{aligned} |\delta \mathbf{X}_t|^2 &\leq e^{(I + \tilde{\mu}_x + \tilde{\mu}_u u_x + \Sigma_x)(t-t_0)} |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= e^{\tilde{\xi}(t-t_0)} |\mathbf{x}_1 - \mathbf{x}_2|^2 \end{aligned} \quad (40)$$

Where $\tilde{\xi} = I + \tilde{\mu}_x + \tilde{\mu}_u u_x + \Sigma_x$.

And for Y term we will have,

$$\begin{aligned} \mathbb{E}[|\delta \mathbf{Y}_t|^2] &= \mathbb{E}|\delta \mathbf{Y}_T|^2 + \int_t^T \mathbb{E} [2(\delta H_s \delta Y_s + \delta k_s \delta Y_s) - \delta \|Z_s^T\|_F^2] \mathbf{d}s \\ &\leq g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T 2H_z \mathbb{E}|\delta Y_s|^2 + (2H_z)^{-1} \mathbb{E}|\delta H_s|^2 + 2k_z \mathbb{E}|\delta Y_s|^2 + (2k_z)^{-1} \mathbb{E}|\delta k_s|^2 - \mathbb{E}\|Z_s\|_F^2 \mathbf{d}s \\ &\leq g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T (2H_z + 2k_z) \mathbb{E}|\delta Y_s|^2 + (2H_z)^{-1} \mathbb{E} [H_x |\delta \mathbf{X}_s|^2 + H_z \|\delta Z_s\|_F^2] \\ &\quad + (2k_z)^{-1} \mathbb{E} [k_x |\delta \mathbf{X}_s|^2 + k_z \|\delta Z_s\|_F^2] - \mathbb{E}\|Z_s\|_F^2 \mathbf{d}s \end{aligned}$$

Noticing that the drift term in BSDE w/o IS is $H_s = C^{i*} + V_x G U_{*,0}$ which is a Lipschitz continuous function, while IS term is $k_s = V_x G U_{*,0}$. then we can have $k_x \leq H_x$. By replacing k_x by H_x , it yields:

$$\begin{aligned} &\leq g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T (2H_z + 2k_z) \mathbb{E}|\delta Y_s|^2 + (2H_z)^{-1} \mathbb{E} [H_x |\delta \mathbf{X}_s|^2 + H_z \|\delta Z_s\|_F^2] \\ &\quad + (2k_z)^{-1} \mathbb{E} [H_x |\delta \mathbf{X}_s|^2 + k_z \|\delta Z_s\|_F^2] - \mathbb{E}\|Z_s\|_F^2 \mathbf{d}s \\ &\leq g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T (2H_z + 2k_z) \mathbb{E}|\delta Y_s|^2 + \frac{H_x}{2} (H_z^{-1} + k_z^{-1}) \mathbb{E}|\delta \mathbf{X}_s|^2 + \mathbb{E}\|\delta Z_s\|_F^2 - \mathbb{E}\|Z_s\|_F^2 \mathbf{d}s \\ &= g_x \mathbb{E}|\delta \mathbf{X}_T|^2 + \int_t^T (2H_z + 2k_z) \mathbb{E}|\delta Y_s|^2 + \frac{H_x}{2} (H_z^{-1} + k_z^{-1}) \mathbb{E}|\delta \mathbf{X}_s|^2 \mathbf{d}s \end{aligned}$$

by Gronwall's inequality

$$|\delta \mathbf{Y}_t|^2 \leq e^{2(H_z + k_z)(T-t)} \left[g_x e^{\tilde{\xi}(T-t_0)} + H_x (H_z^{-1} + k_z^{-1}) \frac{e^{\tilde{\xi}(T-t_0)} - e^{\tilde{\xi}(t-t_0)}}{2\Sigma_x} \right] |\mathbf{x}_1 - \mathbf{x}_2|^2 \quad (41)$$

When $t_0 = 0$, we have,

$$\begin{aligned} |\delta Y_T|^2 &\leq g_x e^{\tilde{\xi}T} |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= \tilde{L}_1 |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ |\delta Y_0|^2 &\leq e^{2(H_z + k_z)T} \left[g_x e^{\tilde{\xi}T} + H_x (H_z^{-1} + k_z^{-1}) \frac{e^{\tilde{\xi}T} - 1}{2\Sigma_x} \right] |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= \tilde{L}_2 |\mathbf{x}_1 - \mathbf{x}_2|^2 \end{aligned} \quad (42)$$

Where $\tilde{\xi} = I + \tilde{\mu}_x + \tilde{\mu}_u u_x + \Sigma_x$. Following arguments in (Ma et al., 2002), one further has,

$$\|Z_t\|_S^2 \leq \|\Sigma\|_S^2 \|\nabla_x Y_t\|_S^2 \leq M_\Sigma \tilde{L}_2 \quad (43)$$

□

D.4. Proof of Theorem 2

According to the result in Lemma.1 with the assumption that the initial state dataset \mathcal{D} are identical for FBSDE w/ and w/o importance sampling.

$$\begin{aligned} |\delta Y_T|^2 &\leq g_x e^{\xi T} |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= L_1 |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ |\delta Y_0|^2 &\leq e^{H_z T} \left[g_x e^{\xi T} + H_x \frac{e^{\xi T} - 1}{H_z \Sigma_x} \right] |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= L_2 |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ \xi &= I + \mu_x + \mu_u u_x + \Sigma_x \end{aligned} \quad (44)$$

Similarly, According to Lemma.2, one have,

$$\begin{aligned} |\delta Y_T|^2 &\leq g_x e^{\tilde{\xi} T} |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= \tilde{L}_1 |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ |\delta Y_0|^2 &\leq e^{2(H_z + k_z)T} \left[g_x e^{\tilde{\xi} T} + H_x (H_z^{-1} + k_z^{-1}) \frac{e^{\tilde{\xi} T} - 1}{2\Sigma_x} \right] |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ &= \tilde{L}_2 |\mathbf{x}_1 - \mathbf{x}_2|^2 \\ \tilde{\xi} &= I + \tilde{\mu}_x + \tilde{\mu}_u u_x + \Sigma_x \end{aligned} \quad (45)$$

We have $\tilde{\mu}_x = \mu_x + m_x \geq \mu_x$ and $\tilde{\mu}_u = \mu_u + m_u \geq \mu_u$, where m_x and m_u are the Lipschitz constants for m_s w.r.t. \mathbf{x} and \mathbf{u} defined in equation.39. Then we have $\tilde{\xi} \geq \xi$ which leads to $\tilde{L}_1 \geq L_1$. Noticing that the drift term in BSDE w/o IS is $H_s = C^{i*} + V_x G U_{*,0}$, while IS term is $k_s = V_x G U_{*,0}$. then we can have $k_x \leq H_x$, and $k_z \leq H_z$ which leads to

$$\frac{1}{2} \left(\frac{1}{H_z} + \frac{1}{k_z} \right) > \frac{1}{H_z} \quad (46)$$

We have known that $\tilde{\mu}_x \geq \mu_x$ and $\frac{1}{2} \left(\frac{1}{H_z} + \frac{1}{k_z} \right) > \frac{1}{H_z}$. Then we can have $\tilde{L}_1 \geq L_1$ and $\tilde{L}_2 \geq L_2$ strictly.

E. Invariant Layer Introductions and Implementation Techniques

E.1. Invariant Mapping

A function f maps its domain from \mathcal{X} to \mathcal{Y} . Domain \mathcal{X} is a vector space \mathbb{R}^d and \mathcal{Y} is a continuous space \mathbb{R} . Assume the function takes a set as input: $\mathbb{X} = \{x_1 \dots x_N\}$, then the function f is invariant if it satisfies property E.1.

Property 1. A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ defined on sets is permutation invariant to the order of objects in the set. i.e. For any permutation function π : $f(\{x_1 \dots x_N\}) = f(\{x_{\pi(1)} \dots x_{\pi(N)}\})$

In this paper, we discuss the case when f is a neural network only.

Theorem 3. (Zaheer et al., 2017) \mathbf{X} has elements from countable universe. A function $f(\mathbf{X})$ is a valid permutation invariant function, i.e invariant to the permutation of \mathbf{X} , iff it can be decomposed in the form $\rho(\sum_{x \in \mathbf{X}} \phi(x))$, for appropriate functions ρ and ϕ .

In the symmetric multi-agent system, each agent is not distinguishable. This property gives some hints about how to extract the features of the $-i$ th agents using a neural network. The states of the $-i$ th agents can be represented as a set: $\mathbf{X} = \{X_1, X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_N\}$. We want to design a neural network f which has the property of permutation invariance. Specifically, ϕ is represented as a one layer neural network and ρ is a common nonlinear activation function, and the invariant layer module is shown in Fig.11.

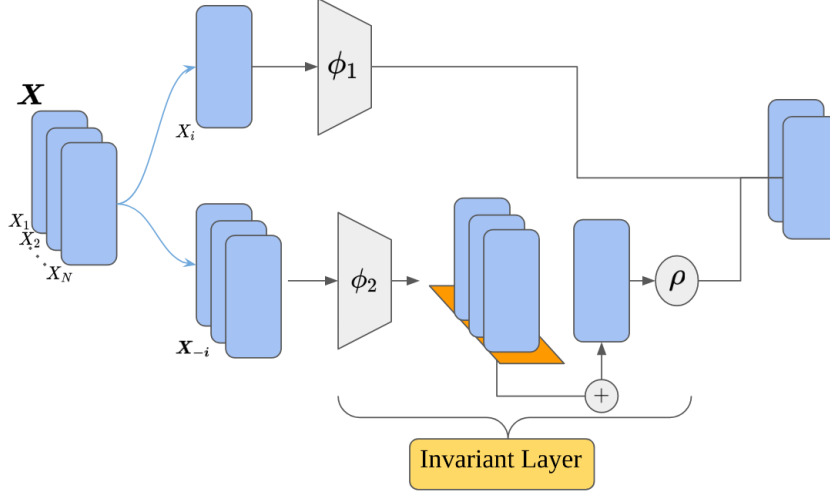


Figure 11. Feature extractor architecture. ϕ represents fully connected neural network. ρ is the ReLU activation function.

E.2. Feature Extractor with Invariant Layer Architecture

The architecture of feature extractor with invariant layer is described in Fig. 11.

E.3. Invariant Layer Techniques

Noticing that all the agents has the access to the global states, we define the state input features of invariant layer for the i th agent as:

$$\mathbf{X}_t = \{X_i, X_1, X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_N\}, \quad (47)$$

with shape of $[B, N, N_x]$, where B is the batch size, N_x is the dimension of the observed states. In the other word, **we always put own feature in the first position**. For each agent i , there will exist such a feature tensor, then for the feature extractor, the shape of input is $[BS, N, N, N_x]$. Therefore, the shape of input tensor will become $[BS, N, N - 1, N_x]$ for invariant layer. where N is the number of agents. First, we could use neural network to map the observed states to the feature space with dimension N_f . Then the shape of the tensor will become $[BS, N, N - 1, N_f]$. After summing up the features of all the element in the set, the dimension of the tensor would reduce to $[BS, N, 1, N_f]$, and we denote this feature tensor as F . However, the memory complexity is $\mathcal{O}(N^2 \times N_f)$ which is not tolerable when the number of agent N increases. Alternatively, we can simply map the tensor \mathbf{X}_t whose dimension is $[BS, N, N_f]$ into the desired feature dimension N_f , then the shape of the tensor would become $[BS, N, N_f]$, and we denote this tensor as F' . Now we create another tensor which is the average of features of element in set with size $[BS, 1, N_f]$ and we denote it to be \bar{F}' . Then we compute $\hat{F} = (\bar{F}' \times N - F') / (N - 1)$ which has size of $[BS, N, N_f]$. We can find that $\hat{F} = F$, and the memory complexity of computing \hat{F} is just $\mathcal{O}(N \times N_f)$. The derivation only holds when the system is symmetric and the agents are not distinguishable. The technique can be extended to higher state dimension for individual agent.

F. Experiment Configurations

This section elaborates the experiment configurations for §4. For all the simulation, the number of SGD iteration is fixed as $N_{SGD} = 100$. We are using Adam as optimizer with 1E-3 learning rate for all simulations.

F.1. Inter-bank Experiments

In section §4.1, For the prediction of initial value function V_0^i , all frameworks are using 2 layers feed forward network with 128 hidden dimension. For the baseline framework, we followed the suggested configuration motioned in (Han et al., 2018). At each time steps, V_x^i is approximated by three layers of feed forward network with 64 hidden dimensions. We add batch norm (Ioffe & Szegedy, 2015) after each affine transformation and before each nonlinear activation function. For Deep

FBSDE with LSTM backbone, we are using two layer LSTM parametrized by 128 hidden state. If the framework includes the invariant layer, the number of mapping features is chosen to be 256. The hyperparameters of the dynamics is listed as following:

$$a = 0.1, q = 0.1, c = 0.5, \epsilon = 0.5, \rho = 0.2, \sigma = 1, T = 1. \quad (48)$$

In the simulation, the time horizon is separated into 40 time-steps over 1 second by Euler method. Learning rate is chosen to be 1E-3 which is the default learning rate for Adam optimizer. The initial state for each agents are sampled from the uniform distribution $[-\delta_0, \delta_0]$. Where δ_0 is the constant standard deviation of state $\mathbf{X}(t)$ during the process as described in (Han & Hu, 2019). In the evaluation, we are using 256 new sampled trajectory which are different from training trajectory to evaluate the performance. The number of stage is set to be 100 which is enough for all framework to converge.

F.2. Belief Space Car Racing

In §4.2, the hyperparameter is listed as following:

$$c_{drag} = 0.01, L = 0.1, c = 0.5, T = 10.0 \quad (49)$$

The observation noise is sampled from Gaussian noise $m \sim \mathcal{N}(0, 0.1\mathbf{I})$. The time horizon is enrolled into 100 time-steps by Euler method. In this experiments, the initial value V_i is approximated a single trainable scale and $V_{x,i}(t)$ is approximated by two layers of LSTM parametrized with 32 hidden dimensions. The number of stage is set to be 10.

G. FBSDEs and Analytical Solution for Inter-Bank Borrowing/Lending Problem

G.1. FBSDEs for Inter-Bank Borrowing/Lending Problem

By plugging the running cost (17) to the HJB (5) function, one can have,

$$V_t^i + \inf_{U_i \in \mathcal{U}_i} \left[\sum_{j=1}^N [a(\bar{X} - X_j) + U_j^2] V_{x_j} + \frac{1}{2} U_i^2 - q U_i (\bar{X} - X_i) + \frac{\epsilon}{2} (\bar{X} - X_i)^2 \right] + \frac{1}{2} \text{tr}(V_{xx}^i \Sigma \Sigma^T) = 0. \quad (50)$$

By computing the infimum explicitly, the optimal control of player i is: $U_i^*(\mathbf{X}, t) = q(\bar{X} - X_i) - V_x^i(\mathbf{X}, t)$. The final form of HJB can be obtained as

$$V_t^i + \frac{1}{2} \text{tr}(V_{xx}^i \Sigma \Sigma^T) + a(\bar{X} - X_i) V_x^i + \sum_{j \neq i} [a(\bar{X} - X_j) + U_j] V_x^j + \frac{\epsilon}{2} (\bar{X} - X_i)^2 - \frac{1}{2} (q(\bar{X} - X_i) - V_x^i)^2 = 0 \quad (51)$$

Applying Feynman-Kac lemma to equation 51, the corresponding FBSDE system is

$$\begin{aligned} d\mathbf{X}(t) &= (f(\mathbf{X}(t), t) + G(\mathbf{X}(t), t)\mathbf{u}(t))dt + \Sigma(t, \mathbf{X}(t))d\mathbf{W}_t, \quad \mathbf{X}(0) = \mathbf{x}_0 \\ dV^i &= -\left[\frac{\epsilon}{2} (\bar{X} - X_i)^2 - \frac{1}{2} (q(\bar{X} - X_i) - V_x^i)^2 + U_i \right] dt + V_x^{iT} \Sigma dW, \quad V(T) = g(\mathbf{X}(T)). \end{aligned} \quad (52)$$

G.2. Analytical solutions for Inter-Bank Borrowing/Lending Problem

The analytical solution for linear inter-bank problem was derived in (Carmona et al., 2013). We provide them here for completeness. Assume the ansatz for HJB function is described as:

$$V_i(t, \mathbf{X}) = \frac{\eta(t)}{2} (\bar{X} - X_i)^2 = \mu(t) \quad i \in \mathbb{I} \quad (53)$$

Where $\eta(t), \mu(t)$ are two scalar functions. The optimal control under this ansatz is:

$$U_i^*(t, \mathbf{X}) = \left[q + \eta(t) \left(1 - \frac{1}{N} \right) \right] (\bar{X} - X_i) \quad (54)$$

By plugging the ansatz into HJB function derived in equation (51), one can have,

$$\begin{aligned}\dot{\eta}(t) &= 2(a+q)\eta(t) + \left(1 - \frac{1}{N^2}\right)\eta^2(t) - (\epsilon - q^2), \quad \eta(T) = c, \\ \dot{\mu}(t) &= -\frac{1}{2}\sigma^2(1 - \rho^2)\left(1 - \frac{1}{N}\right)\eta(t), \quad \mu(T) = 0.\end{aligned}\tag{55}$$

There exists the analytical solution for the Riccati equation described above as,

$$\eta(t) = \frac{-(\epsilon - q^2)(e^{(\delta^+ - \delta^-)(T-t)} - 1) - c(\delta^+ e^{(\delta^+ - \delta^-)(T-t)} - \delta^-)}{(\delta^- e^{(\delta^+ - \delta^-)(T-t)} - \delta^+) - c(1 - 1/N^2)(e^{(\delta^+ - \delta^-)(T-t)} - 1)}.\tag{56}$$

Where $\delta^\pm = -(a+q) \pm \sqrt{R}$ and $R = (a+q)^2 + (1 - 1/N^2)(\epsilon - q^2)$

H. Additional Tables and Figures

H.1. Evaluation Loss with different number of agents

Fig.12 shows the comparison of SDFP-FBSDE and baseline by evaluation loss.

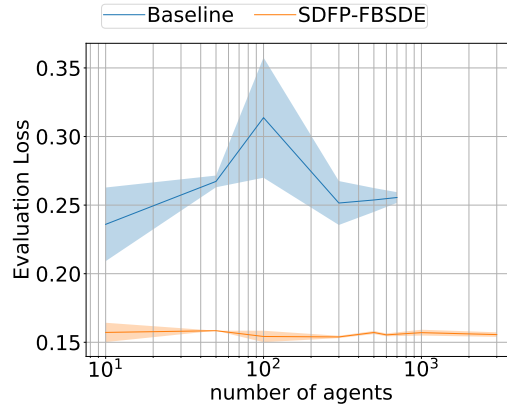


Figure 12. Comparison of SDFP-FBSDE and Baseline for inter-bank problem with different number of agents evaluated on evaluation loss(24).

H.2. Superlinear Inter-Bank Plots

Fig.13 demonstrates the performance difference between Baseline and our algorithm. One can find that our algorithm convergence faster and better than baseline. Since in the superlinear case, the influence of control term in the forward dynamics is mitigated, then the final performances are similar.

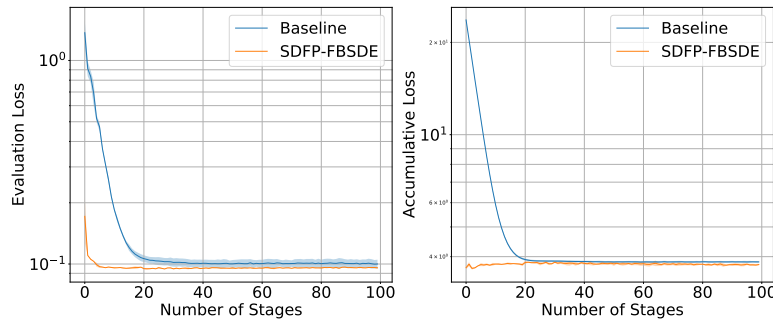


Figure 13. Comparison of SDFP-FBSDE and Baseline for inter-bank problem in superlinear case with evaluation loss (24) and cumulative loss(25)

H.3. Posterior Plot of Car Racing

Fig.14 illustrates the trajectory of single game with posterior estimated by each car. One can find that the variance does not blow up, and both of two cars are staying in the track.

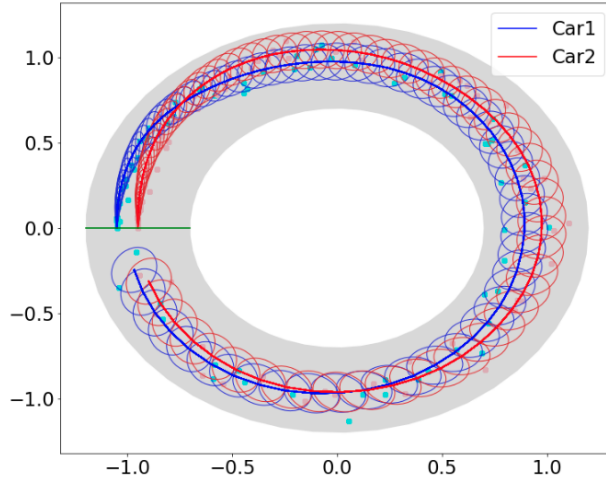


Figure 14. Car racing plot with posterior trained with SDFP-FBSDE. The competition loss is turned off

H.4. DFP-FBSDE Framework

In this subsection, we demonstrate the framework of Deep Fictitious Play FBSDE (DFP-FBSDE) in Fig.15. Each NN (blue box) represents for the FBSDE module shown in Fig.1.

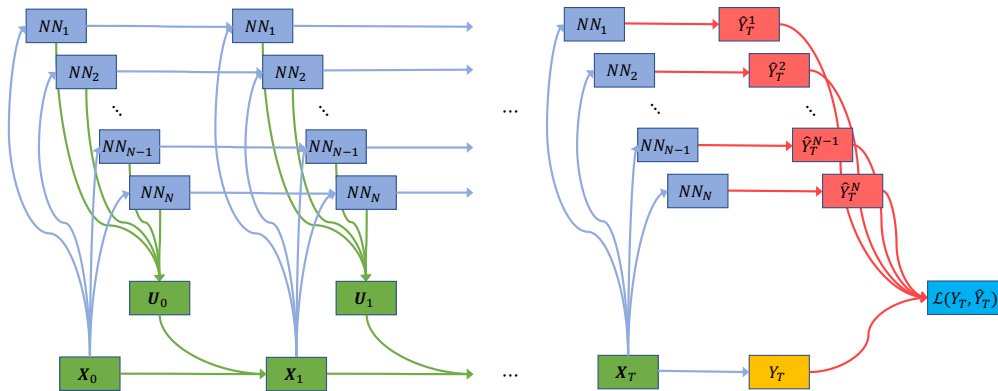


Figure 15. SDFP framework for N Players. Each NN block has the architecture in Fig. 1

I. Belief Car Racing

The framework for the racing problem is trained with batch size of 64, and 100 time steps over a time horizon of 10 seconds.

I.1. Continuous Time Extended Kalman Filter

The Partial Observable Markov Decision Process is generally difficult to solve within infinite dimensional space belief. Commonly, the Value function does not have explicit parameterized form. Kalman filter overcome this challenge by

presuming the noise distribution is Gaussian distribution. In order to deploy proposed Forward Backward Stochastic Differential Equation (FBSDE) model in the Belief space, we need to utilize extended Kalman filter in continuous time (Jazwinski, 1970) correspondingly. Given the partial observable stochastic system:

$$\frac{dx}{dt} = f(x, u, w, t), \quad \text{and} \quad z = h(x, v, t) \quad (57)$$

Where f is the stochastic state process featured by a Gaussian noise $w \sim \mathcal{N}(0, Q)$, h is the observation function while $v \sim \mathcal{N}(0, R)$ is the observation noise. Next, we consider the linearization of the stochastic dynamics in eq.(47) represented as follows:

$$A = \left. \frac{\partial f}{\partial x} \right|_{\hat{x}}, L = \left. \frac{\partial f}{\partial w} \right|_{\hat{x}}, C = \left. \frac{\partial h}{\partial x} \right|_{\hat{x}}, M = \left. \frac{\partial h}{\partial v} \right|_{\hat{x}}, \tilde{Q} = LQL^T, \tilde{R} = MRM^T \quad (58)$$

one can write the posterior mean state \hat{x} and prior covariance matrix P^- estimation update rule by (Simon, 2006):

$$\begin{aligned} \hat{x}(0)\mathbb{E}[x(0)], \quad P^-(0) &= \mathbb{E}[(x(0) - \hat{x})(x(0) - \hat{x})^T] \\ K &= PC^T \tilde{R}^{-1} \\ \dot{\hat{x}} &= b(\hat{x}, u, w, t) = f(\hat{x}, u, w_0, t) + K[z - h(\hat{x}, v_0, t)] \\ \dot{P}^- &= AP^- + P^-A^T + \tilde{Q} - P^-C^T \tilde{R}^{-1}CP^- \end{aligned} \quad (59)$$

We follow the notation in (Simon, 2006), where x is the real state, \hat{x} is the mean of state estimated by Kalman filter based on the noisy sensor observation, P^- represents for the covariance matrix of the estimated state, nominal noise values are given as $w_0 = 0$ and $v_0 = 0$, where superscript + is the posterior estimation and - is the prior estimation. Then we can define a Gaussian belief dynamics as $\mathbf{b}(\hat{x}_k, P_k^-)$ by the mean state \hat{x} and variance P^- of normal distribution $\mathcal{N}(\hat{x}_k, P_k^-)$

The belief dynamics results in a decoupled FBSDE system as follows:

$$\begin{aligned} d\mathbf{b}_k &= g(\mathbf{b}_k, \mathbf{U}_k, 0)dt + \Sigma(\mathbf{b}_k, \mathbf{U}_k, 0)dW, \quad dW \sim \mathcal{N}(0, I) \\ dV &= -C^{i*} dt + V_x^T \Sigma dW \end{aligned} \quad (60)$$

where:

$$\begin{aligned} g(\mathbf{b}_k, \mathbf{U}_k) &= \begin{bmatrix} b(t, \mathbf{X}(t), U_{i,m}(t); \mathbf{U}_{-i,m}) \\ \text{vec}(A_k P_k^- + P_k^- A_k^T + \tilde{Q}_k - P_k^- C_k^T \tilde{R}_k^{-1} C_k P_k^-) \end{bmatrix} \\ \Sigma(\mathbf{b}_k, \mathbf{U}_k) &= \begin{bmatrix} \sqrt{K_k C_k P_k^- dt} \\ \mathbf{0} \end{bmatrix} \\ V(T) &= g(\mathbf{X}(T)) \\ \hat{\mathbf{X}}(0) &= \mathbb{E}[\mathbf{X}(0)] \\ P^-(0) &= \mathbb{E}[(\mathbf{X}(0) - \hat{\mathbf{X}})(\mathbf{X}(0) - \hat{\mathbf{X}})^T] \end{aligned} \quad (61)$$

In the car racing case, the dynamic function $f(\cdot, \cdot)$ in eq.57 is described as,

$$\begin{aligned} d\mathbf{X} &= (f(\mathbf{X}) + G(\mathbf{X})\mathbf{U})dt + \Sigma(\mathbf{X})dW, \quad \mathbf{z} = h(\mathbf{X}) + m \\ f(\mathbf{X}) &= \begin{bmatrix} v \cos \theta \\ v \sin \theta \\ -c_{\text{drag}}v \\ 0 \end{bmatrix}, \quad G(\mathbf{X}) = \Sigma(\mathbf{X}) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & v/L \end{bmatrix}, \quad h(\mathbf{X}) = \mathbf{x} \end{aligned} \quad (62)$$

Where dW is standard Brownian motion.

I.2. Cost Functions

We consider the problem of two cars racing on a circular track. The cost function of each car is designed as

$$J_t = \underbrace{\exp\left(\left|\frac{x^2}{a^2} + \frac{y^2}{b^2} - 1\right|\right)}_{\text{track cost}} + \underbrace{\text{ReLU}(-v)}_{\text{velocity cost}} + \underbrace{\exp(-d)}_{\text{collision cost}}$$

Where d is Euclidean distance between two cars. We use continuous time extended Kalman Filter to propagate belief space dynamics described in equation 61.

We introduce the concept competitive game by using an additional competition cost:

$$J_{competition} = \exp\left(-\begin{bmatrix} \cos(\theta) \\ \sin(\theta) \end{bmatrix}^T \begin{bmatrix} x_1 - x_2 \\ y_1 - y_2 \end{bmatrix}\right)$$

Where x_i, y_i is the x, y position of the i th car. When the i th car is leading, the competition loss will be minor, and it will increase exponentially when the car is trailing.

Thanks to the decoupled BSDE structure, each car can measure this competition loss separately and optimize the value function individually.

J. Hardware

All simulations are run on

1. Nvidia RTX TITAN
2. Nvidia GTX TITAN BLACK

K. Algorithm

Algorithm 1 Scalable Deep Fictitious Play FBSDE

- 1: **Hyper-parameters:** N : Number of players; T : Number of timesteps; M : Number of stages in fictitious play; N_{gd} : Number of gradient descent steps per stage; \mathbf{U}_0 : the initial strategies for players in set \mathbb{I} ; B : Batch size; Δt : time discretization (Total time/Number of timesteps); π : Permutation function (E.1).
 - 2: **Parameters:** ϕ : Network weights for Initial Value (IV) prediction $f_{IV}(\cdot)$; θ : Weights and bias of Backbone and Feature extractor (BF) $f_{BF}(\cdot)$.
 - 3: Initialize trainable parameters: θ^0, ϕ^0
 - 4: **for** $m \leftarrow 1$ to M **do**
 - 5: Generate B sample \mathbf{x}_0 and $B \times T$ Noise $\Delta \mathbf{w} \sim \mathcal{N}(0, \mathbf{I}\Delta t)$.
 - 6: **for** $l \leftarrow 0$ to $N_{gd} - 1$ **do**
 - 7: **for** $t \leftarrow 0$ to $T - 1$ **do**
 - 8: **if** $t=0$ **then**
 - 9: Predict value function for i th player: $\hat{y}_0^i = f_{IV}(\mathbf{x}_0; \phi^{m \times N_{gd} + l})$
 - 10: **else**
 - 11: Compute network prediction \hat{z}_i for i th player: $\hat{z}_i = \sum_i^T f_{BF}(\mathbf{x}_t; \theta^{m \times N_{gd} + l})$
 - 12: **end if**
 - 13: Compute i th optimal control: $u_i^* = -R^{-1}(\Gamma_i^T z_i + Q_i^T \mathbf{x}_t)$
 - 14: Infer $-i$ th players' network prediction and stop the gradient for them: $\hat{z}_{-i} = \sum_{-i}^T f_{BF}(\pi(\mathbf{x}_t); \theta^{m \times N_{gd}})$
 - 15: Compute $-i$ th optimal Control and stop the gradient for them: $\mathbf{u}_{-i}^* = -R_{-i}^{-1}(\Gamma_{-i}^T \hat{z}_{-i} + Q_{-i}^T \mathbf{x}_t)$
 - 16: Propagate FSDE: $\mathbf{x}_{t+1} = f_{FSDE}(\mathbf{x}_t, u_i^*, \mathbf{u}_{-i}^*, \Delta \mathbf{w}_t, t)$ (12)
 - 17: Propagate BSDE: $\hat{y}_{t+1}^i = f_{BSDE}(\hat{y}_t^i, \mathbf{x}_t, u_i^*, \mathbf{u}_{-i}^*, \hat{z}_i, \Delta \mathbf{w}_t, t)$ (13)
 - 18: **end for**
 - 19: **end for**
 - 20: Compute True terminal value $y_T^i = g^i(\mathbf{x}_T)$
 - 21: Compute loss: $\mathcal{L}(\hat{y}_T^i, y_T^i) = \frac{1}{B} \|\hat{y}_T^i - y_T^i\|_2^2$ (24)
 - 22: Gradient Update: θ^l, ϕ^l
 - 23: **end for**
-