

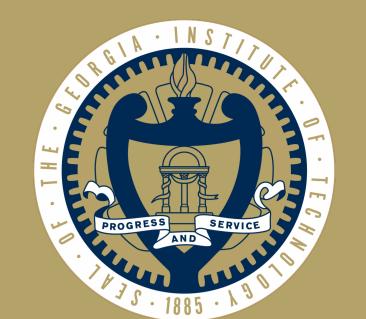
Large Scale Multi-Agent Deep FBSDEs

Tianrong Chen ¹ Ziyi Wang ² Ioannis Exarchos ³ Evangelos A. Theodorou ⁴

¹School of ECE of GaTech ²Center for Machine Learning of GaTech

²Center for Machine Learning of GaTech ³Department of Computer Science of Stanford University

⁴School of Aerospace Engineering of Gatech



Estimating Nash Equilibrium in Stochastic Differential Games

- HJB in Stochastic Differential Game (SDG) can be solve by a PDE solver but it does not scale beyond few dimensions.
- The deep learning model [2] suffers from limited exploration region since the dynamics is driven by Brownian motion.
- The learning algorithms might suffer from curse of many agents [3].

Scalable Fictitious Play Deep FBSDEs

The HJB-PDE of individual agent in SDG is,

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

Scalable Fictitious Play Deep FBSDEs (SFPD-FBSDE) is a deep learning model which can estimate Nash Equilibrium motivated by theoretic analysis of classical FBSDEs system [?],

$$d\mathbf{X}_{t} = (f + G\mathbf{U}_{0,*})dt + \Sigma d\mathbf{W}_{t}, \ \mathbf{X}_{t_{0}} = \mathbf{x}_{t_{0}} \quad (FSDE),$$

$$dY_{t}^{i} = -h_{t}^{i}dt + Z_{t}^{i}dW_{t}, \ Y_{T}^{i} = g(\mathbf{X}_{T}), \quad (BSDE),$$
(2)

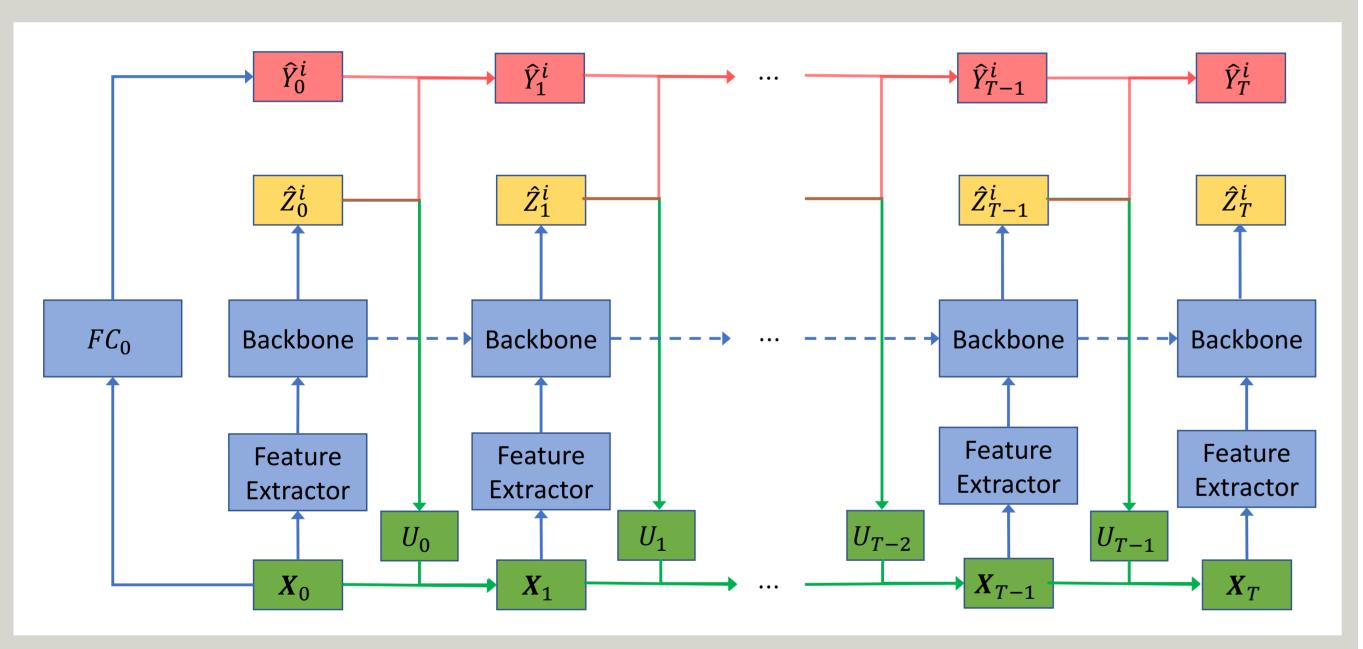


Figure 1. FBSDE Network for a single agent. The dashed arrow indicates hidden states propagation if LSTM is chosen as backbone. The dash arrow would disappear when FC is chosen. Red arrow represents for the BSDE propagation, Green arrow represents for the FSDE propagation

The solution of the BSDE is the solution of HJB-PDE almost surely and can thus be used to estimate the HJB PDE. Besides, in this paper we tackle the aforementioned difficulties:

- Efficient Exploration By incorporating importance Sampling technique, SFPD-FBSDE is guaranteed to have more efficient exploration compared to baselines from theoretical analysis. (modify the green arrow and the red arrow)
- Mitigate Curse of Many Agents SFPD-FBSDE integrates invariant layer [4] to mitigate the curse of many agents. (modify the feature extractor blue boxes)

Importance Sampling

The Importrance Sampling of FBSDEs is known as,

$$d\tilde{\boldsymbol{X}}_{s} = [\mu_{s} + \Sigma K_{s}]ds + \Sigma d\boldsymbol{W}_{s}, \ \tilde{\boldsymbol{X}}_{t} = \boldsymbol{x}_{t},$$

$$d\tilde{Y}_{s}^{i} = [-h_{s}^{i} + \tilde{Z}_{s}K_{s}]ds + \tilde{Z}_{s}^{i}dW_{s}, \ Y_{T}^{i} = g(\boldsymbol{X}_{T}).$$
(3)

The solution of modified BSDE is still aligned with the original HJB-PDE [1].

Exploration and bound for FBSDE system with drift term

Assumption 1 (informal): For a general FBSDE system,

$$\begin{split} \boldsymbol{X}_{T}^{t,\boldsymbol{x}} &= \boldsymbol{x} + \int_{t}^{T} \mu_{s} \mathrm{d}s + \int_{t}^{T} \Sigma_{s} \mathrm{d}\boldsymbol{W}_{s} & \text{(FSDE)}, \\ Y_{t}^{T,\boldsymbol{x}} &= g(\boldsymbol{X}_{T}^{t,\boldsymbol{x}}) - \int_{t}^{T} H_{s} \mathrm{d}s + \int_{t}^{T} Z_{s} \mathrm{d}\boldsymbol{W}_{s} & \text{(BSDE)}, \end{split}$$

The Lipschitz function are bounded.

Lemma 1 (informal): FBSDE system satisfying assumption 1, then we have

$$|\delta Y_T|^2 \le L_1 |\boldsymbol{x}_1 - \boldsymbol{x}_2|^2, |\delta Y_{t_0}|^2 \le L_2 |\boldsymbol{x}_1 - \boldsymbol{x}_2|^2,$$

$$||Z_t||_S^2 \le ||\Sigma||_S^2 ||\nabla_x Y_t||_S^2 \le M_\Sigma L_2$$
(5)

Where L_1 and L_2 are the function of Lipschitz constants.

Theorem 2 (informal) Denote \tilde{Y}_T is the solution of FBSDE with importance sampling, one has,

$$\max |\delta Y_T|^2 \le \max |\delta \tilde{Y}_T|^2,$$

$$\max |\delta Y_0|^2 \le \max |\delta \tilde{Y}_0|^2$$
(6)

Following table shows the similarity of training SFPD-FBSDE and traditional supervised learning task.

Items	Definitions
Training Data	random generated $oldsymbol{x}_0$
Label	terminal cost $Y^* = \boldsymbol{g}(\boldsymbol{x}_T)$
Training Loss Function	$ \hat{oldsymbol{Y}}^T - oldsymbol{Y}^* _2^2$
Trainable Parameters	parameters of backbones and FCs

Thus, Theorem 2 indicates that the framework with importance sampling can provide richer training data and better exploration efficiency which is critical for training a deep learning model.

Invariant Layers

Because of the symmetric problem setup, we incorporate the Invariant Layers (IL) [4] to extract permutation invariant features to improve the performance SFPD-FBSDE. It helps to mitigate the curse of many agents when the number of agent is increasing.

Experiments

We first conduct the ablation experiments on the Inter-Bank Game lending/borrowing problem to verify the effectiveness of IS and IL in fig.2

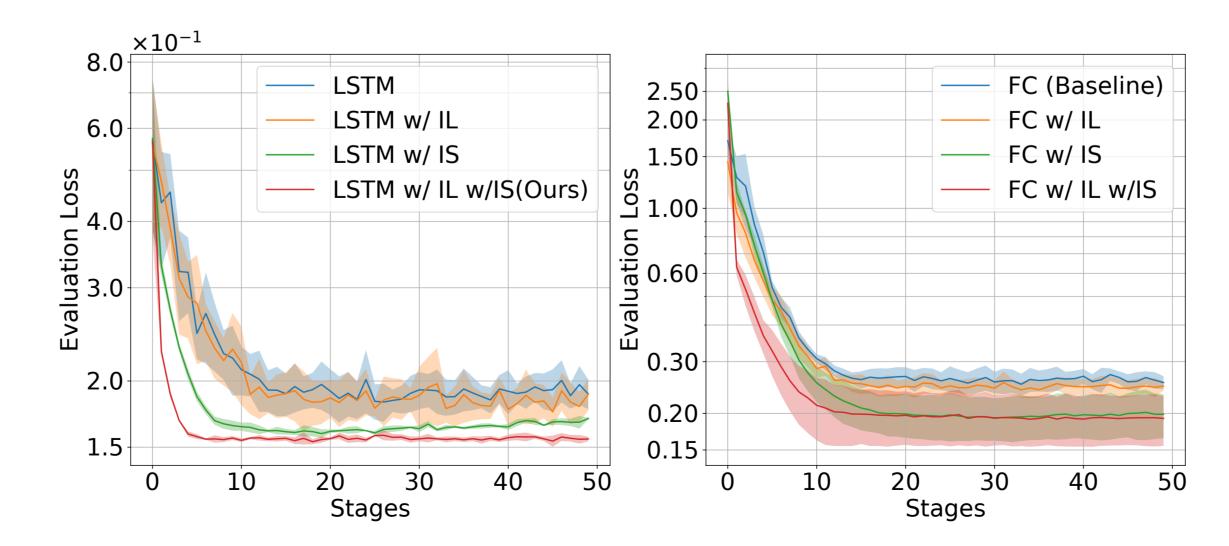


Figure 2. Ablation experiments on LSTM and FC backbone.

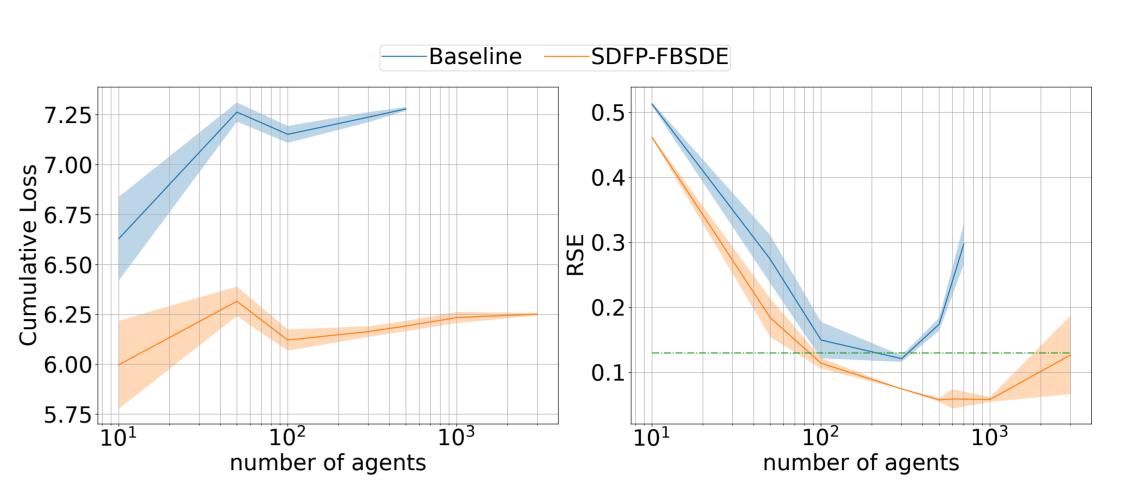


Figure 3. Comparison of SDFP-FBSDE and Baseline for inter-bank problem with different number of agents evaluated on cumulative loss and RSE.

Fig.3 shows the suprior performance of our model when the number of agents increases.

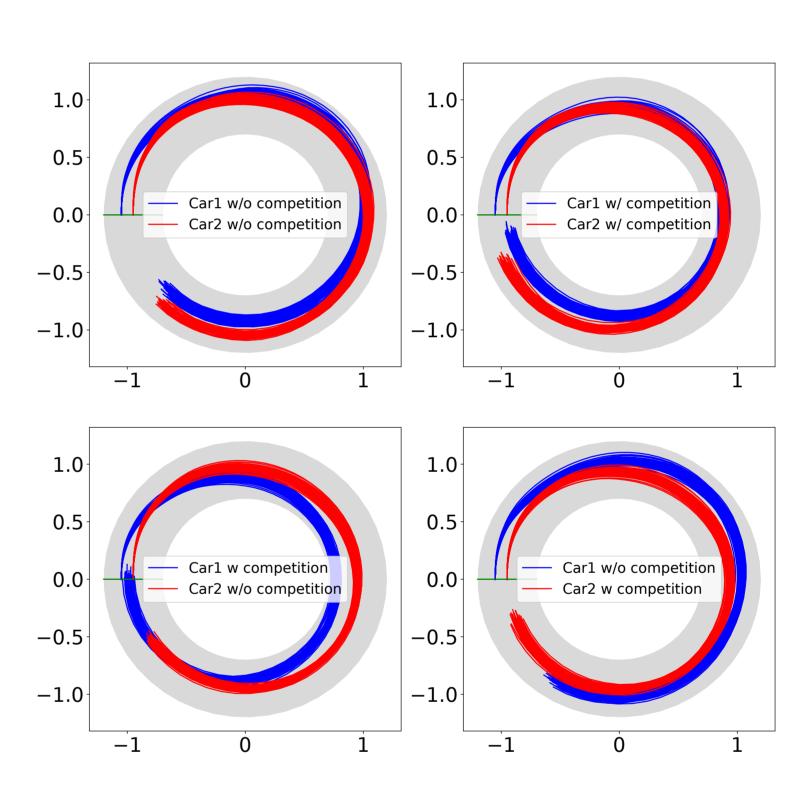


Figure 4. The plots contains 64 trials of racing. The performance varies with respect to the competition loss.

The algorithm is also executed on the partial observed autonomous racing cases in fig.4.

References

- [1] Christian Bender and Thilo Moseler. Importance sampling for backward sdes. Stochastic Analysis and Applications, 28(2):226–253, 2010.
- [2] Jiequn Han, Arnulf Jentzen, and E Weinan. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [3] Lingxiao Wang, Zhuoran Yang, and Zhaoran Wang. Breaking the curse of many agents: Provable mean embedding q-iteration for mean-field reinforcement learning. In *International Conference on Machine Learning*, pages 10092–10103. PMLR, 2020.
- [4] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov, and Alexander Smola. Deep sets. arXiv preprint arXiv:1703.06114, 2017.