

Imperial College
London



HUAWEI



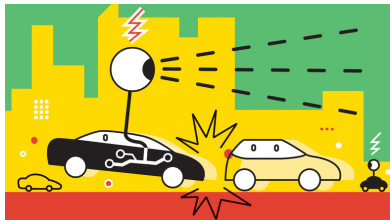
Settling the Variance of Multi-Agent Policy Gradients



Jakub Grudzien Kuba*, Muning Wen*, Yaodong Yang,
Linghui Meng, Shangding Gu, Haifeng Zhang,
David Henry Mguni, Jun Wang

Background: MARL

Multi-Agent Reinforcement Learning —to make everyone happy



(a) A self-driving with single-agent RL



(b) A system of MARL cars

Problem Formulation

At time step t , n agents are at state s_t



state s_t

Problem Formulation

The agents take actions $a_t^1 \sim \pi_\theta^1(\cdot^1|s_t), \dots, a_t^n \sim \pi_\theta^n(\cdot^n|s_t)$

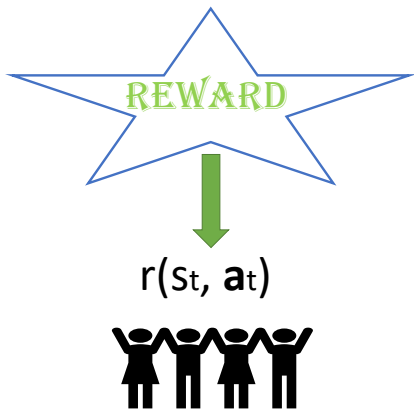


state s_t

Equivalently $a_t \sim \pi_\theta(\cdot|s_t)$

Problem Formulation

The environment emits the joint reward $r(s_t, a_t)$



Problem Formulation

The agents move to the next state $s_{t+1} \sim P(\cdot | s_t, a_t)$



state s_{t+1}

Problem Formulation

The agents want to maximise the joint return

$$\mathcal{J}(\theta) = \mathbb{E}_{s_0 \sim d^0, a_{0:\infty} \sim \pi_\theta, s_{1:\infty} \sim P} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

Multi-Agent Policy Gradient (MAPG) Estimators

Decentralised Training (DT)

$$g_D^i = \sum_{t=0}^{\infty} \gamma^t \hat{Q}^i(s_t, a_t^i) \nabla_{\theta^i} \log \pi_{\theta^i}^i(a_t^i | s_t)$$

Multi-Agent Policy Gradient (MAPG) Estimators

Decentralised Training (DT)

$$g_D^i = \sum_{t=0}^{\infty} \gamma^t \hat{Q}^i(s_t, a_t^i) \nabla_{\theta^i} \log \pi_{\theta^i}(a_t^i | s_t)$$

Centralised Training Decentralised Execution (CTDE)

$$g_C^i = \sum_{t=0}^{\infty} \gamma^t \hat{Q}(s_t, a_t^{-i}, a_t^i) \nabla_{\theta^i} \log \pi_{\theta^i}(a_t^i | s_t)$$

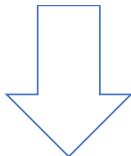
Multi-Agent Policy Gradient (MAPG) Estimators

Decentralised Training (DT)

$$g_D^i = \sum_{t=0}^{\infty} \gamma^t \hat{Q}^i(s_t, a_t^i) \nabla_{\theta^i} \log \pi_{\theta^i}^i(a_t^i | s_t)$$

Centralised Training Decentralised Execution (CTDE)

$$g_C^i = \sum_{t=0}^{\infty} \gamma^t \hat{Q}(s_t, a_t^{-i}, a_t^i) \nabla_{\theta^i} \log \pi_{\theta^i}^i(a_t^i | s_t)$$



same feedback

Credit Assignment Problem

$$\hat{Q}(s_t, a_t^{-i}, a_t^i)$$



Credit Assignment Problem

$$\hat{Q}(s_t, a_t^{-i}, a_t^i)$$



If the bias is not the problem, then what is it?

Credit Assignment Problem

$$\hat{Q}(s_t, a_t^{-i}, a_t^i)$$



If the bias is not the problem, then what is it?
Perhaps variance

Settling the Variance of Multi-Agent Policy Gradients

Settling the Variance of Multi-Agent Policy Gradients

begins with understanding how a subset of agents contributes to the return (the agents affect MAPG estimators through actions).

Settling the Variance of Multi-Agent Policy Gradients

begins with understanding how a subset of agents contributes to the return (the agents affect MAPG estimators through actions).

Multi-agent state-action value function

$$Q_{\theta}^{i: k}(s, \mathbf{a}^{i: k}) = \mathbb{E}_{\mathbf{a}^{-i: k} \sim \pi_{\theta}^{-i: k}} [Q_{\theta}(s, \mathbf{a}^{i: k}, \mathbf{a}^{-i: k})]$$

Settling the Variance of Multi-Agent Policy Gradients

begins with understanding how a subset of agents contributes to the return (the agents affect MAPG estimators through actions).

Multi-agent state-action value function

$$Q_{\theta}^{i:k}(s, \mathbf{a}^{i:k}) = \mathbb{E}_{\mathbf{a}^{-i:k} \sim \pi_{\theta}^{-i:k}} [Q_{\theta}(s, \mathbf{a}^{i:k}, \mathbf{a}^{-i:k})]$$

Multi-agent advantage function

$$A_{\theta}^{i:k}(s, \mathbf{a}^{1:m}, \mathbf{a}^{i:k}) = Q_{\theta}^{j_1:m, i:k}(s, \mathbf{a}^{j_1:m}, \mathbf{a}^{i:k}) - Q_{\theta}^{j_1:m}(s, \mathbf{a}^{j_1:m})$$

Multi-Agent Advantage Decomposition

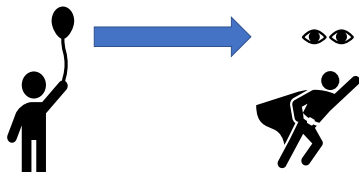
$$A_{\theta}^{i: m}(s, \mathbf{a}^{i: m}) = \sum_{j=1}^m A_{\theta}^{j_j}(s, \mathbf{a}^{1: j-1}, a^j)$$



$$A^{1,2}(s, \text{stick figure}, \text{superhero})$$

Multi-Agent Advantage Decomposition

$$A_{\theta}^{i:m}(s, \mathbf{a}^{i:m}) = \sum_{j=1}^m A_{\theta}^{ij}(s, \mathbf{a}^{1:j-1}, a^j)$$



$$A^{1,2}(s, \text{balloon}, \text{superhero}) \\ = A^1(s, \text{balloon}) + A^2(s, \text{balloon}, \text{superhero})$$

$$\text{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, a_{0:\infty} \sim \pi_{\theta}} [g_C^i] - \text{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, a_{0:\infty} \sim \pi_{\theta}} [g_D^i] = \mathcal{O}\left(\sum_{j \neq i} \epsilon_j^2\right)$$

$$\text{where } \epsilon_j = \max_{s, \mathbf{a}^{-j}, a^j} |A_{\theta}^j(s, \mathbf{a}^{-j}, a^j)|$$

$$\text{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, a_{0:\infty} \sim \pi_{\theta}} [g_C^i] - \text{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, a_{0:\infty} \sim \pi_{\theta}} [g_D^i] = \mathcal{O}\left(\sum_{j \neq i} \epsilon_j^2\right)$$

$$\text{where } \epsilon_j = \max_{s, \mathbf{a}^{-j}, a^j} |A_{\theta}^j(s, \mathbf{a}^{-j}, a^j)|$$

- The more agents, the worse.

MAPG Variance

$$\text{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, a_{0:\infty} \sim \pi_{\theta}} [g_C^i] - \text{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, a_{0:\infty} \sim \pi_{\theta}} [g_D^i] = \mathcal{O}\left(\sum_{j \neq i} \epsilon_j^2\right)$$

$$\text{where } \epsilon_j = \max_{s, \mathbf{a}^{-j}, a^j} |A_{\theta}^j(s, \mathbf{a}^{-j}, a^j)|$$

- The more agents, the worse.
- The more they explore, the worse.

MAPG Variance Reduction

Modify the estimator with the baseline trick

$$\begin{aligned} g_C^i(b) &= \sum_{t=0}^{\infty} \gamma^t [\hat{Q}(s_t, a_t^{-i}, a_t^i) - b(s_t, a_t^{-i})] \nabla_{\theta^i} \log \pi_{\theta^i}(a_t^i | s_t) \\ &= \sum_{t=0}^{\infty} \gamma^t g_{C,t}^i(b) \end{aligned}$$

MAPG Variance Reduction

Modify the estimator with the baseline trick

$$\begin{aligned}g_C^i(b) &= \sum_{t=0}^{\infty} \gamma^t [\hat{Q}(s_t, a_t^{-i}, a_t^i) - b(s_t, a_t^{-i})] \nabla_{\theta^i} \log \pi_{\theta^i}(a_t^i | s_t) \\ &= \sum_{t=0}^{\infty} \gamma^t g_{C,t}^i(b)\end{aligned}$$

which brings no bias

$$\mathbb{E}_{a_t^{-i} \sim \pi_{\theta^{-i}}, a_t^i \sim \pi_{\theta^i}} [g_{C,t}^i(b)] = \mathbb{E}_{a_t^{-i} \sim \pi_{\theta^{-i}}, a_t^i \sim \pi_{\theta^i}} [g_{C,t}^i(0)]$$

MAPG Variance Reduction

Modify the estimator with the baseline trick

$$\begin{aligned}g_C^i(b) &= \sum_{t=0}^{\infty} \gamma^t [\hat{Q}(s_t, a_t^{-i}, a_t^i) - b(s_t, a_t^{-i})] \nabla_{\theta^i} \log \pi_{\theta^i}(a_t^i | s_t) \\ &= \sum_{t=0}^{\infty} \gamma^t g_{C,t}^i(b)\end{aligned}$$

which brings no bias

$$\mathbb{E}_{a_t^{-i} \sim \pi_{\theta^{-i}}, a_t^i \sim \pi_{\theta^i}} [g_{C,t}^i(b)] = \mathbb{E}_{a_t^{-i} \sim \pi_{\theta^{-i}}, a_t^i \sim \pi_{\theta^i}} [g_{C,t}^i]$$

MAPG Variance Reduction: Decomposition

$$\begin{aligned} \text{Var}_{s_t \sim d_\theta^t, a_t \sim \pi_\theta} [g_{C,t}^i(b)] &= \underbrace{\text{Var}_{s_t \sim d_\theta^t} [\mathbb{E}_{a_t \sim \pi_\theta} [g_{C,t}^i(b)]]}_{\text{Variance from state}} \\ &+ \underbrace{\mathbb{E}_{s_t \sim d_\theta^t} [\text{Var}_{a_t^{-i} \sim \pi_\theta^{-i}} [\mathbb{E}_{a_t^i \sim \pi_\theta^i} [g_{C,t}^i(b)]]]}_{\text{Variance from other agents' actions}} \\ &+ \underbrace{\mathbb{E}_{s_t \sim d_\theta^t, a_t^{-i} \sim \pi_\theta^{-i}} [\text{Var}_{a_t^i \sim \pi_\theta^i} [g_{C,t}^i(b)]]}_{\text{Variance from agent } i\text{'s action}} \end{aligned}$$

MAPG Variance Reduction: Decomposition

$$\begin{aligned} \text{Var}_{s_t \sim d_{\theta}^t, a_t \sim \pi_{\theta}} [g_{C,t}^i(b)] &= \underbrace{\text{Var}_{s_t \sim d_{\theta}^t} [\mathbb{E}_{a_t \sim \pi_{\theta}} [g_{C,t}^i]]}_{\text{Variance from state}} \\ &+ \underbrace{\mathbb{E}_{s_t \sim d_{\theta}^t} [\text{Var}_{a_t^{-i} \sim \pi_{\theta}^{-i}} [\mathbb{E}_{a_t^i \sim \pi_{\theta}^i} [g_{C,t}^i]]]}_{\text{Variance from other agents' actions}} \\ &+ \underbrace{\mathbb{E}_{s_t \sim d_{\theta}^t, a_t^{-i} \sim \pi_{\theta}^{-i}} [\text{Var}_{a_t^i \sim \pi_{\theta}^i} [g_{C,t}^i(b)]]}_{\text{Variance from agent } i\text{'s action}} \end{aligned}$$

MAPG Variance Reduction: Objective

$$\min_{b(s, \mathbf{a}^{-i})} \text{Var}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, a^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s) \right]$$

MAPG Variance Reduction: Optimal Baseline

$$\min_{b(s, \mathbf{a}^{-i})} \text{Var}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, a^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s) \right]$$

The optimal baseline is given by

$$b^{\text{optimal}}(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, a^i) \|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}$$

MAPG Variance Reduction: Optimal Baseline

$$\min_{b(s, \mathbf{a}^{-i})} \text{Var}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, a^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s) \right]$$

The optimal baseline is given by

$$b^{\text{optimal}}(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, a^i) \|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}$$

MAPG Variance Reduction: Optimal Baseline

$$\min_{b(s, \mathbf{a}^{-i})} \text{Var}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, a^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s) \right]$$

The optimal baseline is given by

$$b^{\text{optimal}}(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, a^i) \|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}$$

- No closed-form formula

MAPG Variance Reduction: Optimal Baseline

$$\min_{b(s, \mathbf{a}^{-i})} \text{Var}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, a^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s) \right]$$

The optimal baseline is given by

$$b^{\text{optimal}}(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, a^i) \|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}$$

- No closed-form formula
- The large dimension introduces extra variance in estimation

MAPG Variance Reduction: Optimal Baseline

$$\min_{b(s, \mathbf{a}^{-i})} \text{Var}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, a^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s) \right]$$

The optimal baseline is given by

$$b^{\text{optimal}}(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, a^i) \|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}{\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\|\nabla_{\theta^i} \log \pi_{\theta}^i(a^i | s)\|^2 \right]}$$

- No closed-form formula
- The large dimension introduces extra variance in estimation
- Multiple backpropagations are computationally expensive

Neural Network Policies

$$\pi_{\theta}^i(a^i|s) = \pi^i(a^i|\psi_{\theta}^i(s))$$

Neural Network Policies

Decomposition with the last layer

$$\pi_{\theta}^i(a^i|s) = \pi^i(a^i|\psi_{\theta}^i(s))$$

For example, ψ_{θ}^i can be the layer of logits

$$\pi_{\theta}^i(a^i|s) = \text{softmax}(\psi_{\theta}^i(s)) [a^i]$$

Neural Network Policies

Decomposition with the last layer

$$\pi_{\theta}^i(a^i|s) = \pi^i(a^i|\psi_{\theta}^i(s))$$

Gradient decomposition with chain rule

$$\nabla_{\theta^i} \log \pi_{\theta}^i(a^i|s) = \nabla_{\theta^i} \psi_{\theta}^i(a^i|s) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i|\psi_{\theta}^i(s))$$

Neural Network Policies

Decomposition with the last layer

$$\pi_{\theta}^i(a^i|s) = \pi^i(a^i|\psi_{\theta}^i(s))$$

Gradient decomposition with chain rule

$$\nabla_{\theta^i} \log \pi_{\theta}^i(a^i|s) = \nabla_{\theta^i} \psi_{\theta}^i(a^i|s) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i|\psi_{\theta}^i(s))$$

Neural Network Policies

Decomposition with the last layer

$$\pi_{\theta}^i(a^i|s) = \pi^i(a^i|\psi_{\theta}^i(s))$$

Gradient decomposition with chain rule

$$\nabla_{\theta^i} \log \pi_{\theta}^i(a^i|s) = \nabla_{\theta^i} \psi_{\theta}^i(a^i|s) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i|\psi_{\theta}^i(s))$$

MAPG Variance Reduction: Surrogate Objective

$$\begin{aligned} & \min_{b(s, \mathbf{a}^{-i})} \text{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s) \right] \\ \propto & \min_{b(s, \mathbf{a}^{-i})} \text{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b(s, \mathbf{a}^{-i}) \right) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i | \psi_{\theta}^i(s)) \right] \end{aligned}$$

MAPG Variance Reduction: Optimal Baseline (OB)

The optimal baseline (OB) minimising the surrogate variance is

$$b^*(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} [\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \|\nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i(s))\|^2]}{\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} [\|\nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i(s))\|^2]}$$

MAPG Variance Reduction: Optimal Baseline (OB)

The optimal baseline (OB) minimising the surrogate variance is

$$b^*(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{a^i \sim \pi_\theta^i} [\hat{Q}(s, \mathbf{a}^{-i}, a^i) \|\nabla_{\psi_\theta^i} \log \pi_\theta^i(a^i | \psi_\theta^i(s))\|^2]}{\mathbb{E}_{a^i \sim \pi_\theta^i} [\|\nabla_{\psi_\theta^i} \log \pi_\theta^i(a^i | \psi_\theta^i(s))\|^2]}$$

- it has an explicit formula in the discrete-action case

MAPG Variance Reduction: Optimal Baseline (OB)

The optimal baseline (OB) minimising the surrogate variance is

$$b^*(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} [\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \|\nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i(s))\|^2]}{\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} [\|\nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i(s))\|^2]}$$

- it has an explicit formula in the discrete-action case
- it is easy to approximate in the continuous-action case

MAPG Variance Reduction: Optimal Baseline (OB)

The optimal baseline (OB) minimising the surrogate variance is

$$b^*(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{a^i \sim \pi_\theta^i} [\hat{Q}(s, \mathbf{a}^{-i}, a^i) \|\nabla_{\psi_\theta^i} \log \pi_\theta^i(a^i | \psi_\theta^i(s))\|^2]}{\mathbb{E}_{a^i \sim \pi_\theta^i} [\|\nabla_{\psi_\theta^i} \log \pi_\theta^i(a^i | \psi_\theta^i(s))\|^2]}$$

- it has an explicit formula in the discrete-action case
- it is easy to approximate in the continuous-action case
- in each case, it is cheap to compute

MAPG Variance Reduction: Optimal Baseline (OB)

The optimal baseline (OB) minimising the surrogate variance is

$$b^*(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \|\nabla_{\psi^i} \log \pi_{\theta}^i(\mathbf{a}^i | \psi_{\theta}^i(s))\|^2]}{\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\|\nabla_{\psi^i} \log \pi_{\theta}^i(\mathbf{a}^i | \psi_{\theta}^i(s))\|^2]}$$

- it has explicit formula in the discrete-action case
- it is easy to approximate in the continuous-action case
- in each case, it is cheap to compute

To use it in an MAPG estimator

$$\hat{X}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i) \triangleq \hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b^*(s, \mathbf{a}^{-i})$$

Empirical Results

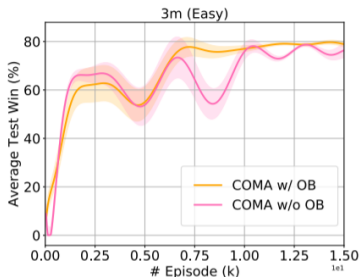
OB significantly reduces the variance of MAPG estimators

a^i	$\psi_{\theta}^i(a^i)$	$\pi_{\theta}^i(a^i)$	$x_{\psi_{\theta}^i}^i(a^i)$	$\hat{Q}(a^{-i}, a^i)$	$\hat{A}^i(a^{-i}, a^i)$	$\hat{X}^i(a^{-i}, a^i)$	Method	Variance
1	$\log 8$	0.8	0.14	2	-9.7	-41.71	MAPG	1321
2	0	0.1	0.43	1	-10.7	-42.71	COMA	1015
3	0	0.1	0.43	100	88.3	56.29	OB	673

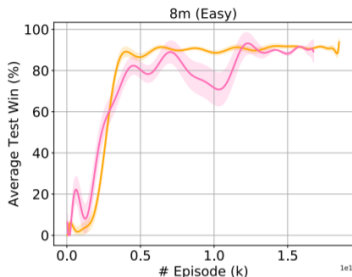
Figure 2: Toy Example

Empirical Results

Lower variance leads to more stable training (example of COMA).



(a) 3 marines

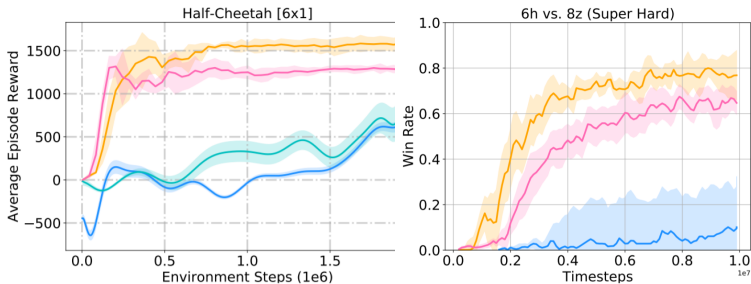


(b) 8 marines

Figure 3: StarCraftII: COMA with OB vs COMA

Empirical Results

Lower variance leads to better performance (example of MAPPO).



(a) MAMuJoCo: MAPPO with OB
vs MAPPO, COMIX, MADDPG

(b) StarCraftII: MAPPO with OB
vs MAPPO, QMIX

Thank you for your attention!

- Jakub Grudzien Kuba (*Imperial College London, Huawei R&D UK*)
- Muning Wen (*Shanghai Jiao Tong University*)
- Yaodong Yang (*King's College London*)
- Linghui Meng (*Institute of Automation, Chinese Academy of Science*)
- Shangding Gu (*Institute of Automation, Chinese Academy of Science*)
- Haifeng Zhang (*Institute of Automation, Chinese Academy of Science*)
- David Henry Mguni (*Huawei R&D UK*)
- Jun Wang (*University College London*)