

Appendices

A Preliminary Remarks	15
B Proofs of the Theoretical Results	16
B.1 Proofs of Lemmas 1, 2, and 3	16
B.2 Proofs of Theorems 1 and 2	18
C Proofs of the Results about Optimal Baselines	21
C.1 Proof of Theorem 3	21
C.2 Remarks about the surrogate optimal baseline	22
C.3 Proof of Theorem 4	23
D Pytorch Implementations of the Optimal Baseline	25
E Computations for the Numerical Toy Example	26
F Detailed Hyper-parameter Settings for Experiments	29

A Preliminary Remarks

Remark 1. *The multi-agent state-action value function obeys the bounds*

$$\left| Q_{\theta}^{i_1, \dots, i_k} \left(s, \mathbf{a}^{(i_1, \dots, i_k)} \right) \right| \leq \frac{\beta}{1-\gamma}, \quad \text{for all } s \in \mathcal{S}, \mathbf{a}^{(i_1, \dots, i_k)} \in \mathcal{A}^{(i_1, \dots, i_k)}.$$

Proof. It suffices to prove that, for all t , the total reward satisfies $|R_t| \leq \frac{\beta}{1-\gamma}$, as the value functions are expectations of it. We have

$$|R_t| = \left| \sum_{k=0}^{\infty} \gamma^k r_{t+k} \right| \leq \sum_{k=0}^{\infty} |\gamma^k r_{t+k}| \leq \sum_{k=0}^{\infty} \gamma^k \beta = \frac{\beta}{1-\gamma}$$

□

Remark 2. *The multi-agent advantage function is bounded.*

Proof. We have

$$\begin{aligned} & \left| A_{\theta}^{i_1, \dots, i_k} \left(s, \mathbf{a}^{(j_1, \dots, j_m)}, \mathbf{a}^{(i_1, \dots, i_k)} \right) \right| \\ &= \left| Q_{\theta}^{j_1, \dots, j_m, i_1, \dots, i_k} \left(s, \mathbf{a}^{(j_1, \dots, j_m, i_1, \dots, i_k)} \right) - Q_{\theta}^{j_1, \dots, j_m} \left(s, \mathbf{a}^{(j_1, \dots, j_m)} \right) \right| \\ &\leq \left| Q_{\theta}^{j_1, \dots, j_m, i_1, \dots, i_k} \left(s, \mathbf{a}^{(j_1, \dots, j_m, i_1, \dots, i_k)} \right) \right| + \left| Q_{\theta}^{j_1, \dots, j_m} \left(s, \mathbf{a}^{(j_1, \dots, j_m)} \right) \right| \leq \frac{2\beta}{1-\gamma} \end{aligned}$$

□

Remark 3. *Baselines in MARL have the following property*

$$\mathbb{E}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} \left[b \left(s, \mathbf{a}^{-i} \right) \nabla_{\theta^i} \log \pi_{\theta}^i \left(\mathbf{a}^i | s \right) \right] = \mathbf{0}.$$

Proof. We have

$$\mathbb{E}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} \left[b \left(s, \mathbf{a}^{-i} \right) \nabla_{\theta^i} \log \pi_{\theta}^i \left(\mathbf{a}^i | s \right) \right] = \mathbb{E}_{s \sim d_{\theta}^t, \mathbf{a}^{-i} \sim \pi_{\theta}^{-i}} \left[b \left(s, \mathbf{a}^{-i} \right) \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\nabla_{\theta^i} \log \pi_{\theta}^i \left(\mathbf{a}^i | s \right) \right] \right],$$

which means that it suffices to prove that for any $s \in \mathcal{S}$

$$\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\nabla_{\theta^i} \log \pi_{\theta}^i \left(\mathbf{a}^i | s \right) \right] = \mathbf{0}.$$

We prove it for continuous \mathcal{A}^i . The discrete case is analogous.

$$\begin{aligned} \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\nabla_{\theta^i} \log \pi_{\theta}^i \left(\mathbf{a}^i | s \right) \right] &= \int_{\mathcal{A}^i} \pi_{\theta}^i \left(\mathbf{a}^i | s \right) \nabla_{\theta^i} \log \pi_{\theta}^i \left(\mathbf{a}^i | s \right) d\mathbf{a}^i \\ &= \int_{\mathcal{A}^i} \nabla_{\theta^i} \pi_{\theta}^i \left(\mathbf{a}^i | s \right) d\mathbf{a}^i = \nabla_{\theta^i} \int_{\mathcal{A}^i} \pi_{\theta}^i \left(\mathbf{a}^i | s \right) d\mathbf{a}^i = \nabla_{\theta^i} (1) = \mathbf{0} \end{aligned}$$

□

B Proofs of the Theoretical Results

B.1 Proofs of Lemmas 1, 2, and 3

In this subsection, we prove the lemmas stated in the paper. We realise that their application to other, very complex, proofs is not always immediately clear. To compensate for that, we provide the stronger versions of the lemmas; we give a detailed proof of the strong version of Lemma 1 which is supposed to demonstrate the equivalence of the normal and strong versions, and prove the normal versions of Lemmas 2 & 3, and state their stronger versions as remarks to the proofs.

Lemma 1 (Multi-agent advantage decomposition). *For any state $s \in \mathcal{S}$, the following equation holds for any subset of m agents and any permutation of their labels,*

$$A_{\theta}^{1,\dots,m} \left(s, \mathbf{a}^{(1,\dots,m)} \right) = \sum_{i=1}^m A_{\theta}^i \left(s, \mathbf{a}^{(1,\dots,i-1)}, a^i \right).$$

Proof. We prove a slightly **stronger**, but perhaps less telling, version of the lemma, which is

$$A_{\theta}^{k+1,\dots,m} \left(s, \mathbf{a}^{(1,\dots,k)}, \mathbf{a}^{(k+1,\dots,m)} \right) = \sum_{i=k+1}^m A_{\theta}^i \left(s, \mathbf{a}^{(1,\dots,i-1)}, a^i \right). \quad (13)$$

The original form of the lemma will follow from the above by taking $k = 0$.

By the definition of the multi-agent advantage, we have

$$\begin{aligned} & A_{\theta}^{k+1,\dots,m} \left(s, \mathbf{a}^{(1,\dots,k)}, \mathbf{a}^{(k+1,\dots,m)} \right) \\ &= Q_{\theta}^{1,\dots,k,k+1,\dots,m} \left(s, \mathbf{a}^{(1,\dots,k,k+1,\dots,m)} \right) - Q_{\theta}^{1,\dots,k} \left(s, \mathbf{a}^{(1,\dots,k)} \right) \end{aligned}$$

which can be written as a telescoping sum

$$\begin{aligned} & Q_{\theta}^{1,\dots,k,k+1,\dots,m} \left(s, \mathbf{a}^{(1,\dots,k,k+1,\dots,m)} \right) - Q_{\theta}^{1,\dots,k} \left(s, \mathbf{a}^{(1,\dots,k)} \right) \\ &= \sum_{i=k+1}^m \left[Q_{\theta}^{(1,\dots,i)} \left(s, \mathbf{a}^{(1,\dots,i)} \right) - Q_{\theta}^{(1,\dots,i-1)} \left(s, \mathbf{a}^{(1,\dots,i-1)} \right) \right] \\ &= \sum_{i=k+1}^m A_{\theta}^i \left(s, \mathbf{a}^{(1,\dots,i-1)}, a^i \right) \end{aligned}$$

□

Lemma 2. *For any state $s \in \mathcal{S}$, we have*

$$\mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [A_{\theta}(s, \mathbf{a})] = \sum_{i=1}^n \mathbb{E}_{a^1 \sim \pi_{\theta}^1, \dots, a^{i-1} \sim \pi_{\theta}^{i-1}} \left[\mathbf{Var}_{a^i \sim \pi_{\theta}^i} \left[A_{\theta}^i \left(s, \mathbf{a}^{(1,\dots,i-1)}, a^i \right) \right] \right].$$

Proof. The trick of this proof is to develop a relation on the variance of multi-agent advantage which is recursive over the number of agents. We have

$$\begin{aligned} \mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [A_{\theta}(s, \mathbf{a})] &= \mathbf{Var}_{a^1 \sim \pi_{\theta}^1, \dots, a^n \sim \pi_{\theta}^n} \left[A_{\theta}^{1,\dots,n} \left(s, \mathbf{a}^{(1,\dots,n)} \right) \right] \\ &= \mathbb{E}_{a^1 \sim \pi_{\theta}^1, \dots, a^n \sim \pi_{\theta}^n} \left[A_{\theta}^{1,\dots,n} \left(s, \mathbf{a}^{(1,\dots,n)} \right)^2 \right] \\ &= \mathbb{E}_{a^1 \sim \pi_{\theta}^1, \dots, a^{n-1} \sim \pi_{\theta}^{n-1}} \left[\mathbb{E}_{a^n \sim \pi_{\theta}^n} \left[A_{\theta}^{1,\dots,n} \left(s, \mathbf{a}^{(1,\dots,n)} \right)^2 \right] \right. \\ &\quad \left. - \mathbb{E}_{a^n \sim \pi_{\theta}^n} \left[A_{\theta}^{1,\dots,n} \left(s, \mathbf{a}^{(1,\dots,n)} \right) \right]^2 + \mathbb{E}_{a^n \sim \pi_{\theta}^n} \left[A_{\theta}^{1,\dots,n} \left(s, \mathbf{a}^{(1,\dots,n)} \right) \right]^2 \right] \\ &= \mathbb{E}_{a^1 \sim \pi_{\theta}^1, \dots, a^{n-1} \sim \pi_{\theta}^{n-1}} \left[\mathbf{Var}_{a^n \sim \pi_{\theta}^n} \left[A_{\theta}^{1,\dots,n} \left(s, \mathbf{a}^{(1,\dots,n)} \right) \right] \right] \\ &\quad + \mathbb{E}_{a^1 \sim \pi_{\theta}^1, \dots, a^{n-1} \sim \pi_{\theta}^{n-1}} \left[A_{\theta}^{1,\dots,n-1} \left(s, \mathbf{a}^{(1,\dots,n-1)} \right)^2 \right] \end{aligned}$$

which, by the stronger version of Lemma 1, given by Equation 13, applied to the first term, equals

$$\begin{aligned} & \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{n-1} \sim \pi_\theta^{n-1}} \left[\mathbf{Var}_{\mathbf{a}^n \sim \pi_\theta^n} \left[A_\theta^n \left(s, \mathbf{a}^{(1, \dots, n-1)}, \mathbf{a}^n \right) \right] \right] \\ & + \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{n-1} \sim \pi_\theta^{n-1}} \left[A_\theta^{1, \dots, n-1} \left(s, \mathbf{a}^{(1, \dots, n-1)} \right)^2 \right] \end{aligned}$$

Hence, we have a recursive relation

$$\begin{aligned} & \mathbf{Var}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^n \sim \pi_\theta^n} \left[A_\theta^{1, \dots, n} \left(s, \mathbf{a}^{(1, \dots, n)} \right) \right] \\ & = \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{n-1} \sim \pi_\theta^{n-1}} \left[\mathbf{Var}_{\mathbf{a}^n \sim \pi_\theta^n} \left[A_\theta^n \left(s, \mathbf{a}^{(1, \dots, n-1)}, \mathbf{a}^n \right) \right] \right] \\ & + \mathbf{Var}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{n-1} \sim \pi_\theta^{n-1}} \left[A_\theta^{1, \dots, n-1} \left(s, \mathbf{a}^{(1, \dots, n-1)} \right) \right] \end{aligned}$$

from which we can obtain

$$\begin{aligned} & \mathbf{Var}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^n \sim \pi_\theta^n} \left[A_\theta^{1, \dots, n} \left(s, \mathbf{a}^{(1, \dots, n)} \right) \right] \\ & = \sum_{i=1}^n \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}} \left[\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^i \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^i \right) \right] \right] \end{aligned}$$

□

Remark 4. Lemma 2 has a **stronger** version, coming as a corollary to the above proof; that is

$$\begin{aligned} & \mathbf{Var}_{\mathbf{a}^{k+1} \sim \pi_\theta^{k+1}, \dots, \mathbf{a}^n \sim \pi_\theta^n} \left[A_\theta^{k+1, \dots, n} \left(s, \mathbf{a}^{(1, \dots, k)}, \mathbf{a}^{(k+1, \dots, n)} \right) \right] \\ & = \sum_{i=k+1}^n \mathbb{E}_{\mathbf{a}^{k+1} \sim \pi_\theta^{k+1}, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}} \left[\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^i \left(s, \mathbf{a}^{(1, \dots, k)}, \mathbf{a}^{(k+1, \dots, i-1)}, \mathbf{a}^i \right) \right] \right]. \quad (14) \end{aligned}$$

We think of it as a corollary to the proof of the lemma, as the fixed joint action $\mathbf{a}^{1, \dots, k}$ has the same algebraic properties, throughout the proof, as state s .

Lemma 3. For any state $s \in \mathcal{S}$, we have

$$\mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [A_\theta(s, \mathbf{a})] \leq \sum_{i=1}^n \mathbf{Var}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}, \mathbf{a}^i \sim \pi_\theta^i} [A_\theta^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)].$$

Proof. By Lemma 2, we have

$$\mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [A_\theta(s, \mathbf{a})] = \sum_{i=1}^n \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}} \left[\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^i \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^i \right) \right] \right] \quad (15)$$

Take an arbitrary i . We have

$$\begin{aligned} & \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}} \left[\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^i \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^i \right) \right] \right] \\ & = \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}} \left[\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^i \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^i \right)^2 \right] \right] \\ & = \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}} \left[\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbb{E}_{\mathbf{a}^{i+1} \sim \pi_\theta^{i+1}, \dots, \mathbf{a}^n \sim \pi_\theta^n} \left[A_\theta^{i, \dots, n} \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^{(i, \dots, n)} \right) \right]^2 \right] \right] \\ & \leq \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}} \left[\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbb{E}_{\mathbf{a}^{i+1} \sim \pi_\theta^{i+1}, \dots, \mathbf{a}^n \sim \pi_\theta^n} \left[A_\theta^{i, \dots, n} \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^{(i, \dots, n)} \right) \right]^2 \right] \right] \\ & = \mathbb{E}_{\mathbf{a}^1 \sim \pi_\theta^1, \dots, \mathbf{a}^{i-1} \sim \pi_\theta^{i-1}, \mathbf{a}^{i+1} \sim \pi_\theta^{i+1}, \dots, \mathbf{a}^n \sim \pi_\theta^n} \left[\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^{i, \dots, n} \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^{(i, \dots, n)} \right) \right]^2 \right] \end{aligned}$$

The above can be equivalently, but more tellingly, rewritten after permuting (cyclic shift) the labels of agents, in the following way

$$\begin{aligned} & \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^{i+1, \dots, n, i} \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^{(i+1, \dots, n, i)} \right) \right]^2 \right] \\ & = \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} \left[A_\theta^{i+1, \dots, n, i} \left(s, \mathbf{a}^{(1, \dots, i-1)}, \mathbf{a}^{(i+1, \dots, n, i)} \right) \right] \right] \end{aligned}$$

which, by the strong version of Lemma 1, equals

$$\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_{\theta}^{-i}} \left[\mathbf{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[A_{\theta}^i (s, \mathbf{a}^{-i}, \mathbf{a}^i) \right] \right]$$

which can be further simplified by

$$\begin{aligned} \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_{\theta}^{-i}} \left[\mathbf{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[A_{\theta}^i (s, \mathbf{a}^{-i}, \mathbf{a}^i) \right] \right] &= \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_{\theta}^{-i}} \left[\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[A_{\theta}^i (s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right] \right] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi_{\theta}} \left[A_{\theta}^i (s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right] = \mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} \left[A_{\theta}^i (s, \mathbf{a}^{-i}, \mathbf{a}^i) \right] \end{aligned}$$

which, combined with Equation 15, finishes the proof. \square

Remark 5. Again, subsuming a joint action $\mathbf{a}^{(1, \dots, k)}$ into state in the above proof, we can have a **stronger** version of Lemma 3,

$$\begin{aligned} &\mathbf{Var}_{\mathbf{a}^{k+1} \sim \pi_{\theta}^{k+1}, \dots, \mathbf{a}^n \sim \pi_{\theta}^n} \left[A_{\theta}^{k+1, \dots, n} (s, \mathbf{a}^{(1, \dots, k)}, \mathbf{a}^{(k+1, \dots, n)}) \right] \\ &\leq \sum_{i=k+1}^n \mathbf{Var}_{\mathbf{a}^{k+1} \sim \pi_{\theta}^{k+1}, \dots, \mathbf{a}^n \sim \pi_{\theta}^n} \left[A_{\theta}^i (s, \mathbf{a}^{(k+1, \dots, i-1, i+1, \dots, n)}, \mathbf{a}^i) \right] \end{aligned} \quad (16)$$

B.2 Proofs of Theorems 1 and 2

Let us recall the two assumptions that we make in the paper.

Assumption 1. The state space \mathcal{S} , and every agent i 's action space \mathcal{A}^i is either discrete and finite, or continuous and compact.

Assumption 2. For all $i \in \mathcal{N}$, $s \in \mathcal{S}$, $\mathbf{a}^i \in \mathcal{A}^i$, the map $\theta^i \mapsto \pi_{\theta}^i(\mathbf{a}^i | s)$ is continuously differentiable.

These assumptions assure that the supremum $\sup_{s, \mathbf{a}^i} \|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|$ exists for every agent i . We notice that the supremum $\sup_{s, \mathbf{a}^{-i}, \mathbf{a}^i} |A_{\theta}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)|$ exists regardless of assumptions, as by Remark 2, the multi-agent advantage is bounded from both sides.

Theorem 1. The CTDE and DT estimators of MAPG satisfy

$$\mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{C}}^i] - \mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{D}}^i] \leq \frac{B_i^2}{1 - \gamma^2} \sum_{j \neq i} \epsilon_j^2 \leq (n-1) \frac{(\epsilon B_i)^2}{1 - \gamma^2}$$

where $B_i = \sup_{s, \mathbf{a}} \|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|$, $\epsilon_i = \sup_{s, \mathbf{a}^{-i}, \mathbf{a}^i} |A_{\theta}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)|$, and $\epsilon = \max_i \epsilon_i$.

Proof. It suffices to prove the first inequality, as the second one is a trivial upper bound. Let's consider an arbitrary time step $t \geq 0$. Let

$$\begin{aligned} \mathbf{g}_{\mathcal{C}, t}^i &= \hat{Q}(s_t, \mathbf{a}_t) \nabla_{\theta^i} \log \pi_{\theta}^i(s_t, \mathbf{a}_t^i) \\ \mathbf{g}_{\mathcal{D}, t}^i &= \hat{Q}^i(s_t, \mathbf{a}_t^i) \nabla_{\theta^i} \log \pi_{\theta}^i(s_t, \mathbf{a}_t^i) \end{aligned}$$

be the contributions to the centralised and decentralised gradient estimators coming from sampling $s_t \sim d_{\theta}^t$, $\mathbf{a}_t \sim \pi_{\theta}$. Note that

$$\mathbf{g}_{\mathcal{C}}^i = \sum_{t=0}^{\infty} \gamma^t \mathbf{g}_{\mathcal{C}, t}^i \quad \text{and} \quad \mathbf{g}_{\mathcal{D}}^i = \sum_{t=0}^{\infty} \gamma^t \mathbf{g}_{\mathcal{D}, t}^i$$

Moreover, let $\mathbf{g}_{\mathcal{C}, t, j}^i$ and $\mathbf{g}_{\mathcal{D}, t, j}^i$ be the j^{th} components of $\mathbf{g}_{\mathcal{C}, t}^i$ and $\mathbf{g}_{\mathcal{D}, t}^i$, respectively. Using the law of total variance, we have

$$\begin{aligned} &\mathbf{Var}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{C}, t, j}^i] - \mathbf{Var}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{D}, t, j}^i] \\ &= \left(\mathbf{Var}_{s \sim d_{\theta}^t} [\mathbb{E}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{C}, t, j}^i]] + \mathbb{E}_{s \sim d_{\theta}^t} [\mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{C}, t, j}^i]] \right) \\ &\quad - \left(\mathbf{Var}_{s \sim d_{\theta}^t} [\mathbb{E}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{D}, t, j}^i]] + \mathbb{E}_{s \sim d_{\theta}^t} [\mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{\mathcal{D}, t, j}^i]] \right) \end{aligned} \quad (17)$$

Noting that \mathbf{g}_C^i and \mathbf{g}_D^i have the same expectation over $\mathbf{a} \sim \pi_\theta$, the above simplifies to

$$\begin{aligned} & \mathbb{E}_{s \sim d_\theta^i} [\mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{C,t,j}^i]] - \mathbb{E}_{s \sim d_\theta^i} [\mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{D,t,j}^i]] \\ &= \mathbb{E}_{s \sim d_\theta^i} [\mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{C,t,j}^i] - \mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{D,t,j}^i]] \end{aligned} \quad (18)$$

Let's fix a state s . Using (again) the fact that the expectations of the two gradients are the same, we have

$$\begin{aligned} & \mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{C,t,j}^i] - \mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{D,t,j}^i] \\ &= \left(\mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\mathbf{g}_{C,t,j}^i \right)^2 \right] - \mathbb{E}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{C,t,j}^i]^2 \right) - \left(\mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\mathbf{g}_{D,t,j}^i \right)^2 \right] - \mathbb{E}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{D,t,j}^i]^2 \right) \\ &= \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\mathbf{g}_{C,t,j}^i \right)^2 \right] - \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\mathbf{g}_{D,t,j}^i \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\mathbf{g}_{C,t,j}^i \right)^2 - \left(\mathbf{g}_{D,t,j}^i \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\frac{\partial \log \pi_\theta^i(\mathbf{a}^i | s)}{\partial \theta^i} \hat{Q}(s, \mathbf{a}) \right)^2 - \left(\frac{\partial \log \pi_\theta^i(\mathbf{a}^i | s)}{\partial \theta^i} \hat{Q}^i(s, \mathbf{a}^i) \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\left(\frac{\partial \log \pi_\theta^i(\mathbf{a}^i | s)}{\partial \theta^i} \right)^2 \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\hat{Q}(s, \mathbf{a}^i, \mathbf{a}^{-i})^2 - \hat{Q}^i(s, \mathbf{a}^i)^2 \right] \right]. \end{aligned}$$

The inner expectation is the variance of $\hat{Q}(s, \mathbf{a}^i, \mathbf{a}^{-i})$, given \mathbf{a}^i . We rewrite it as

$$\begin{aligned} &= \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\left(\frac{\partial \log \pi_\theta^i(\mathbf{a}^i | s)}{\partial \theta^i} \right)^2 \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\left(\hat{Q}(s, \mathbf{a}^i, \mathbf{a}^{-i}) - \hat{Q}^i(s, \mathbf{a}^i) \right)^2 \right] \right] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\frac{\partial \log \pi_\theta^i(\mathbf{a}^i | s)}{\partial \theta^i} \right)^2 \left(\hat{Q}(s, \mathbf{a}) - \hat{Q}^i(s, \mathbf{a}^i) \right)^2 \right]. \end{aligned}$$

Now, recalling that the variance of the total gradient is the sum of variances of the gradient components, we have

$$\begin{aligned} & \mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{C,t}^i] - \mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{D,t}^i] = \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left\| \nabla_{\theta^i} \log \pi_\theta^i(\mathbf{a}^i | s) \right\|^2 \left(\hat{Q}(s, \mathbf{a}) - \hat{Q}^i(s, \mathbf{a}^i) \right)^2 \right] \\ & \leq B_i^2 \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\left(\hat{Q}(s, \mathbf{a}) - \hat{Q}^i(s, \mathbf{a}^i) \right)^2 \right] = B_i^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\left(\hat{Q}(s, \mathbf{a}) - \hat{Q}^i(s, \mathbf{a}^i) \right)^2 \right] \right] \\ & = B_i^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\left(\hat{Q}(s, \mathbf{a}^i, \mathbf{a}^{-i}) - \hat{Q}^i(s, \mathbf{a}^i) \right)^2 \right] \right] \\ & = B_i^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\hat{A}^{-i}(s, \mathbf{a}^i, \mathbf{a}^{-i})^2 \right] \right] = B_i^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbf{Var}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\hat{A}^{-i}(s, \mathbf{a}^i, \mathbf{a}^{-i}) \right] \right] \end{aligned}$$

which by the strong version of Lemma 3, given in Equation 16, can be upper-bounded by

$$B_i^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\sum_{j \neq i} \mathbf{Var}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\hat{A}^j(s, \mathbf{a}^{-j}, \mathbf{a}^j) \right] \right] = B_i^2 \sum_{j \neq i} \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbf{Var}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\hat{A}^j(s, \mathbf{a}^{-j}, \mathbf{a}^j) \right] \right]$$

Notice that, for any $j \neq i$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbf{Var}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\hat{A}^j(s, \mathbf{a}^{-j}, \mathbf{a}^j) \right] \right] \\ &= \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_\theta^{-i}} \left[\hat{A}^j(s, \mathbf{a}^{-j}, \mathbf{a}^j)^2 \right] \right] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi_\theta} \left[\hat{A}^j(s, \mathbf{a}^{-j}, \mathbf{a}^j)^2 \right] \leq \epsilon_j^2 \end{aligned}$$

This gives

$$\mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{C,t}^i] - \mathbf{Var}_{\mathbf{a} \sim \pi_\theta} [\mathbf{g}_{D,t}^i] \leq B_i^2 \sum_{j \neq i} \epsilon_j^2$$

and combining it with Equations 17 and 18 for entire gradient vectors, we get

$$\mathbf{Var}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{C,t}^i] - \mathbf{Var}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{D,t}^i] \leq B_i^2 \sum_{j \neq i} \epsilon_j^2 \quad (19)$$

Noting that

$$\begin{aligned} \mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}^i] &= \mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{g}_{t}^i \right] \\ &= \sum_{t=0}^{\infty} \mathbf{Var}_{s_t \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\gamma^t \mathbf{g}_{t}^i] = \sum_{t=0}^{\infty} \gamma^{2t} \mathbf{Var}_{s_t \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{t}^i] \end{aligned}$$

Combining this series expansion with the estimate from Equation 19, we finally obtain

$$\begin{aligned} \mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_C^i] - \mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_D^i] \\ \leq \sum_{t=0}^{\infty} \gamma^{2t} \left(B_i^2 \sum_{j \neq i} \epsilon_j^2 \right) \leq \frac{B_i^2}{1 - \gamma^2} \sum_{j \neq i} \epsilon_j^2 \end{aligned}$$

□

Theorem 2. *The COMA and DT estimators of MAPG satisfy*

$$\mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_{COMA}^i] - \mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_D^i] \leq \frac{(\epsilon_i B_i)^2}{1 - \gamma^2}$$

Proof. Just like in the proof of Theorem 1, we start with the difference

$$\mathbf{Var}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{COMA,t,j}^i] - \mathbf{Var}_{s \sim d_{\theta}^t, \mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{D,t,j}^i]$$

which we transform to an analogue of Equation 18:

$$\mathbb{E}_{s \sim d_{\theta}^t} [\mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{COMA,t,j}^i] - \mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{D,t,j}^i]]$$

which is trivially upper-bounded by

$$\mathbb{E}_{s \sim d_{\theta}^t} [\mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{COMA,t,j}^i]]$$

Now, let us fix a state s . We have

$$\begin{aligned} \mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{COMA,t,j}^i] &= \mathbf{Var}_{\mathbf{a}^{-i} \sim \pi_{\theta}^{-i}, \mathbf{a}^i \sim \pi_{\theta}^i} \left[\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} A^i(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right] \\ &\leq \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_{\theta}^{-i}, \mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 A^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right] \\ &\leq \epsilon_i^2 \mathbb{E}_{\mathbf{a}^{-i} \sim \pi_{\theta}^{-i}, \mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] \end{aligned} \quad (20)$$

which summing over all components of θ^i gives

$$\mathbf{Var}_{\mathbf{a} \sim \pi_{\theta}} [\mathbf{g}_{COMA,t}^i] \leq (\epsilon_i B_i)^2$$

Now, applying the reasoning from Equation 19 until the end of the proof of Theorem 1, we arrive at the result

$$\mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_{COMA}^i] - \mathbf{Var}_{s_{0:\infty} \sim d_{\theta}^{0:\infty}, \mathbf{a}_{0:\infty} \sim \pi_{\theta}} [\mathbf{g}_D^i] \leq \frac{(\epsilon_i B_i)^2}{1 - \gamma^2}$$

□

C Proofs of the Results about Optimal Baselines

In this section of the Appendix we prove the results about optimal baselines, which are those that minimise the CTDE MAPG estimator's variance. We rely on the following variance decomposition

$$\begin{aligned}
\text{Var}_{s_t \sim d_{\theta}^t, \mathbf{a}_t \sim \pi_{\theta}} [\mathbf{g}_{C,t}^i(b)] &= \text{Var}_{s_t \sim d_{\theta}^t} [\mathbb{E}_{\mathbf{a}_t \sim \pi_{\theta}} [\mathbf{g}_{C,t}^i(b)]] + \mathbb{E}_{s_t \sim d_{\theta}^t} [\text{Var}_{\mathbf{a}_t \sim \pi_{\theta}} [\mathbf{g}_{C,t}^i(b)]] \\
&= \text{Var}_{s_t \sim d_{\theta}^t} [\mathbb{E}_{\mathbf{a}_t \sim \pi_{\theta}} [\mathbf{g}_{C,t}^i(b)]] + \mathbb{E}_{s_t \sim d_{\theta}^t} [\text{Var}_{\mathbf{a}_t^{-i} \sim \pi_{\theta}^{-i}} [\mathbb{E}_{\mathbf{a}_t^i \sim \pi_{\theta}^i} [\mathbf{g}_{C,t}^i(b)]] + \mathbb{E}_{\mathbf{a}_t^{-i} \sim \pi_{\theta}^{-i}} [\text{Var}_{\mathbf{a}_t^i \sim \pi_{\theta}^i} [\mathbf{g}_{C,t}^i(b)]]] \\
&= \underbrace{\text{Var}_{s_t \sim d_{\theta}^t} [\mathbb{E}_{\mathbf{a}_t \sim \pi_{\theta}} [\mathbf{g}_{C,t}^i(b)]]}_{\text{Variance from state}} + \underbrace{\mathbb{E}_{s_t \sim d_{\theta}^t} [\text{Var}_{\mathbf{a}_t^{-i} \sim \pi_{\theta}^{-i}} [\mathbb{E}_{\mathbf{a}_t^i \sim \pi_{\theta}^i} [\mathbf{g}_{C,t}^i(b)]]]}_{\text{Variance from other agents' actions}} + \underbrace{\mathbb{E}_{s_t \sim d_{\theta}^t, \mathbf{a}_t^{-i} \sim \pi_{\theta}^{-i}} [\text{Var}_{\mathbf{a}_t^i \sim \pi_{\theta}^i} [\mathbf{g}_{C,t}^i(b)]]}_{\text{Variance from agent } i\text{'s action}}.
\end{aligned} \tag{21}$$

This decomposition reveals that baselines impact the variance via the local variance $\text{Var}_{\mathbf{a}_t^i \sim \pi_{\theta}^i} [\mathbf{g}_{C,t}^i(b)]$. We rely on this fact in the proofs below.

C.1 Proof of Theorem 3

Theorem 3 (Optimal baseline for MAPG). *The optimal baseline (OB) for the MAPG estimator is*

$$b^{\text{optimal}}(s, \mathbf{a}^{-i}) = \frac{\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|^2]}{\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|^2]} \tag{7}$$

Proof. From the decomposition of the estimator's variance, we know that minimisation of the variance is equivalent to minimisation of the local variance

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

For a baseline b , we have

$$\begin{aligned}
&\text{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b \right) \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right) \right] \\
&= \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b \right)^2 \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] \\
&\quad - \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b \right) \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right) \right]^2 \\
&= \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b \right)^2 \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] \\
&\quad - \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right) \right]^2
\end{aligned} \tag{22}$$

as b is a baseline. So in order to minimise variance, we shall minimise the term 22.

$$\begin{aligned}
& \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b \right)^2 \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] \\
&= \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(b^2 - 2b \hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) + \hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right) \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] \\
&= b^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] - 2b \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] \\
&\quad + \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right]
\end{aligned}$$

which is a quadratic in b . The last term of the quadratic does not depend on b , and so it can be treated as a constant. Recalling that the variance of the whole gradient vector $\mathbf{g}^i(b)$ is the sum of variances of its components $g_j^i(b)$, we obtain it by summing over j

$$\begin{aligned}
& \mathbf{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - b \right) \nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s) \right] \\
&= \sum_j \left(b^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] \right. \\
&\quad \left. - 2b \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \left(\frac{\partial \log \pi_{\theta}^i(\mathbf{a}^i | s)}{\partial \theta_j^i} \right)^2 \right] + const \right) \\
&= b^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|^2 \right] \\
&\quad - 2b \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|^2 \right] + const \tag{23}
\end{aligned}$$

As the leading coefficient is positive, the quadratic achieves the minimum at

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

□

C.2 Remarks about the surrogate optimal baseline

In the paper, we discussed the impracticality of the above baseline. To handle this, we noticed that the policy $\pi_{\theta}^i(\mathbf{a}^i | s)$, at state s , is determined by the output layer, $\psi_{\theta}^i(s)$, of an actor neural network. With this representation, in order to handle the impracticality of the above optimal baseline, we considered a minimisation objective, the *surrogate local variance*, given by

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

As a corollary to the proof, the surrogate version of the optimal baseline (which we refer to as OB) was proposed, and it is given by

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

Remark 6. The $x_{\psi_{\theta}^i}^i$ measure, for which $b^*(s, \mathbf{a}^{-i}) = \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)]$, is generally given by

$$x_{\psi_{\theta}^i}^i(\mathbf{a}^i | s) = \frac{\pi_{\theta}^i(\mathbf{a}^i | s) \|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|^2}{\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\|\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i | s)\|^2]} \quad (24)$$

Let us introduce the definition of the softmax function, which is the subject of the next definition. For a vector $\mathbf{z} \in \mathbb{R}^d$, we have $\text{softmax}(\mathbf{z}) = \left(\frac{e^{z_1}}{\eta}, \dots, \frac{e^{z_d}}{\eta}\right)$, where $\eta = \sum_{j=1}^d e^{z_j}$. We write $\text{softmax}(\psi_{\theta}^i(s))(\mathbf{a}^i) = \frac{\exp(\psi_{\theta}^i(s)(\mathbf{a}^i))}{\sum_{\tilde{\mathbf{a}}^i} \exp(\psi_{\theta}^i(s)(\tilde{\mathbf{a}}^i))}$.

Remark 7. When the action space is discrete, and the actor's policy is $\pi_{\theta}^i(\mathbf{a}^i | s) = \text{softmax}(\psi_{\theta}^i(s))(\mathbf{a}^i)$, then the $x_{\psi_{\theta}^i}^i$ measure is given by

$$x_{\psi_{\theta}^i}^i(\mathbf{a}^i | s) = \frac{\pi_{\theta}^i(\mathbf{a}^i | s) \left(1 + \|\pi_{\theta}^i(s)\|^2 - 2\pi_{\theta}^i(\mathbf{a}^i | s)\right)}{1 - \|\pi_{\theta}^i(s)\|^2}$$

Proof. As we do not vary states s and parameters θ in this proof, let us drop them from the notation for π_{θ}^i , and $\psi_{\theta}^i(s)$, hence writing $\pi^i(\mathbf{a}^i) = \text{softmax}(\psi^i)(\mathbf{a}^i)$. Let us compute the partial derivatives:

$$\begin{aligned} \frac{\partial \log \pi^i(\mathbf{a}^i)}{\partial \psi^i(\tilde{\mathbf{a}}^i)} &= \frac{\partial \log \text{softmax}(\psi^i)(\mathbf{a}^i)}{\partial \psi^i(\tilde{\mathbf{a}}^i)} = \frac{\partial}{\partial \psi^i(\tilde{\mathbf{a}}^i)} \left[\log \frac{\exp(\psi^i(\mathbf{a}^i))}{\sum_{\hat{\mathbf{a}}^i} \exp(\psi^i(\hat{\mathbf{a}}^i))} \right] \\ &= \frac{\partial}{\partial \psi^i(\tilde{\mathbf{a}}^i)} \left[\psi^i(\mathbf{a}^i) - \log \sum_{\hat{\mathbf{a}}^i} \exp(\psi^i(\hat{\mathbf{a}}^i)) \right] \\ &= \mathbf{I}(\mathbf{a}^i = \tilde{\mathbf{a}}^i) - \frac{\exp(\psi^i(\tilde{\mathbf{a}}^i))}{\sum_{\hat{\mathbf{a}}^i} \exp(\psi^i(\hat{\mathbf{a}}^i))} = \mathbf{I}(\mathbf{a}^i = \tilde{\mathbf{a}}^i) - \pi^i(\tilde{\mathbf{a}}^i) \end{aligned}$$

where \mathbf{I} is the indicator function, taking value 1 if the statement input to it is true, and 0 otherwise. Taking \mathbf{e}_k to be the standard normal vector with 1 in k^{th} entry, we have the gradient

$$\nabla_{\psi^i} \log \pi^i(\mathbf{a}^i) = \mathbf{e}_{\mathbf{a}^i} - \pi^i \quad (25)$$

which has the squared norm

$$\begin{aligned} \|\nabla_{\psi^i} \log \pi^i(\mathbf{a}^i)\|^2 &= \|\mathbf{e}_{\mathbf{a}^i} - \pi^i\|^2 = (1 - \pi^i(\mathbf{a}^i))^2 + \sum_{\tilde{\mathbf{a}}^i \neq \mathbf{a}^i} (-\pi^i(\tilde{\mathbf{a}}^i))^2 \\ &= 1 + \sum_{\tilde{\mathbf{a}}^i} (-\pi^i(\tilde{\mathbf{a}}^i))^2 - 2\pi^i(\mathbf{a}^i) = 1 + \|\pi^i\|^2 - 2\pi^i(\mathbf{a}^i). \end{aligned}$$

The expected value of this norm is

$$\begin{aligned} \mathbb{E}_{\mathbf{a}^i \sim \pi^i} \left[1 + \|\pi^i\|^2 - 2\pi^i(\mathbf{a}^i) \right] &= 1 + \|\pi^i\|^2 - \mathbb{E}_{\mathbf{a}^i \sim \pi^i} [2\pi^i(\mathbf{a}^i)] \\ &= 1 + \|\pi^i\|^2 - 2 \sum_{\tilde{\mathbf{a}}^i} (\pi^i(\tilde{\mathbf{a}}^i))^2 = 1 - \|\pi^i\|^2 \end{aligned}$$

which combined with Equation 24 finishes the proof. \square

C.3 Proof of Theorem 4

Theorem 4. The excess surrogate local variance for baseline b satisfies

$$\Delta \text{Var}(b) = (b - b^*(s, \mathbf{a}^{-i}))^2 \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\left\| \nabla_{\psi_{\theta}^i} \log \pi^i(\mathbf{a}^i | \psi_{\theta}^i(s)) \right\|^2 \right]$$

In particular, the excess variance of the vanilla MAPG and COMA estimators satisfy

$$\Delta \text{Var}_{\text{MAPG}} \leq D_i^2 \left(\text{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} [A_{\theta}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)] + Q_{\theta}^{-i}(s, \mathbf{a}^{-i})^2 \right) \leq D_i^2 \left(\epsilon_i^2 + \left[\frac{\beta}{1-\gamma} \right]^2 \right)$$

$$\Delta \text{Var}_{\text{COMA}} \leq D_i^2 \text{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} [A_{\theta}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)] \leq (\epsilon_i D_i)^2$$

where $D_i = \sup_{\mathbf{a}^i} \left\| \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i | \psi_{\theta}^i(s)) \right\|$, and $\epsilon_i = \sup_{s, \mathbf{a}^{-i}, \mathbf{a}^i} |A_{\theta}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)|$.

Proof. The first part of the theorem (the formula for excess variance) follows from Equation 23. For the rest of the statements, it suffices to show the first of each inequalities, as the later ones follow directly from the fact that $|Q_\theta(s, \mathbf{a})| \leq \frac{\beta}{1-\gamma}$, $\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} [A_\theta^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)] = \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} [A_\theta^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2]$, and the definition of ϵ_i . Let us first derive the bounds for $\Delta \mathbf{Var}_{\text{MAPG}}$. Let us, for short-hand, define

$$c_\theta^i := \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\left\| \nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i) \right\|^2 \right]$$

We have

$$\begin{aligned} \Delta \mathbf{Var}_{\text{MAPG}} &= \Delta \mathbf{Var}(0) = c_\theta^i b^*(s, \mathbf{a}^{-i})^2 = c_\theta^i \mathbb{E}_{\mathbf{a}^i \sim x_{\psi_\theta^i}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right]^2 \\ &\leq c_\theta^i \mathbb{E}_{\mathbf{a}^i \sim x_{\psi_\theta^i}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right] = c_\theta^i \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\frac{\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \left\| \nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i) \right\|^2}{c_\theta^i} \right] \\ &= \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \left\| \nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i(s)) \right\|^2 \right] \\ &\leq \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 D_i^2 \right] \\ &= D_i^2 \left(\mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right] - \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right]^2 + \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right]^2 \right) \\ &= D_i^2 \left(\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right] + \hat{Q}^{-i}(s, \mathbf{a}^{-i})^2 \right) \\ &= D_i^2 \left(\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{A}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right] + \hat{Q}^{-i}(s, \mathbf{a}^{-i})^2 \right) \end{aligned}$$

which finishes the proof for MAPG. For COMA, we have

$$\begin{aligned} \Delta \mathbf{Var}_{\text{COMA}} &= \Delta \mathbf{Var} \left(\hat{Q}^{-i}(s, \mathbf{a}^{-i}) \right) = c_\theta^i \left(b^*(s, \mathbf{a}^{-i}) - \hat{Q}^{-i}(s, \mathbf{a}^{-i}) \right)^2 \\ &= c_\theta^i \left(\mathbb{E}_{\mathbf{a}^i \sim x_{\psi_\theta^i}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right] - \hat{Q}^{-i}(s, \mathbf{a}^{-i}) \right)^2 \\ &= c_\theta^i \mathbb{E}_{\mathbf{a}^i \sim x_{\psi_\theta^i}^i} \left[\hat{Q}(s, \mathbf{a}^{-i}, \mathbf{a}^i) - \hat{Q}^{-i}(s, \mathbf{a}^{-i}) \right]^2 \\ &= c_\theta^i \mathbb{E}_{\mathbf{a}^i \sim x_{\psi_\theta^i}^i} \left[\hat{A}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i) \right]^2 \\ &\leq c_\theta^i \mathbb{E}_{\mathbf{a}^i \sim x_{\psi_\theta^i}^i} \left[\hat{A}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right] \\ &= c_\theta^i \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{A}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \frac{\left\| \nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i(s)) \right\|^2}{c_\theta^i} \right] \\ &= \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{A}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \left\| \nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i | \psi_\theta^i(s)) \right\|^2 \right] \\ &\leq \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[D_i^2 \hat{A}^i(s, \mathbf{a}^{-i}, \mathbf{a}^i)^2 \right] \leq (\epsilon_i D_i)^2 \end{aligned}$$

which finishes the proof. \square

D Pytorch Implementations of the Optimal Baseline

First, we import necessary packages, which are **PyTorch** [20], and its **nn.functional** sub-package. These are standard Deep Learning packages used in RL [1].

```
1 import torch, torch.nn.functional as F
```

We then implement a simple method that normalises a row vector, so that its (non-negative) entries sum up to 1, making the vector a probability distribution.

```
1 # x: batch of row vectors to normalise to probability mass
2 normalize = lambda x: F.normalize(x, p=1, dim=-1)
```

The **discrete** OB is an exact dot product between the measure $x_{\psi_{\theta}^i}^i$, and available values of \hat{Q} .

```
1 # q: Q values of actions of agent i
2 # pi: policy of agent i
3 def optimal_baseline(q, pi):
4     M = torch.norm(pi, dim=-1, keepdim = True) ** 2 + 1
5     xweight = normalize( (M - 2 * pi) * pi )
6     return (xweight * q).sum(-1)
```

In the **continuous** case, the measure $x_{\psi_{\theta}^i}^i$ and Q -values can only be sampled at finitely many points.

```
1 # a: sampled actions of agent i
2 # q: Q values of the sampled actions
3 # mu, std: parameters of the Gaussian policy of agent i
4 def optimal_baseline(a, q, mu, std):
5     mu_term = torch.norm((a - mu)/std**2, dim=-1)
6     std_term = torch.norm(((a - mu)**2 - std**2)/std**3, dim=-1)
7     xweight = normalize(mu_term**2 + std_term**2)
8     return (xweight * q).sum(-1)
```

We can incorporate it into our MAPG algorithm by simply replacing the values of advantage with the values of X , in the buffer. Below, we present a discrete example

```
1 # compute the policy and sample an action from it
2 a, pi = actor(obs)
3 q = critic(obs)
4
5 #compute OB
6 ob = optimal_baseline(q, pi)
7
8 # use OB to construct the loss
9 q = q.gather(-1, a)
10 pi = pi.gather(-1, a)
11 X = q - ob
12 loss = -(X * torch.log(pi)).mean()
```

and a continuous one

```
1 # normal sampling step, where log_pi is the log probability of a
2 a, log_pi = actor(obs, deterministic=False)
3 q = critic(obs, a)
4
5 # resample m (e.g., m=1000) actions for the observation
6 obs_m = obs.unsqueeze(0).repeat(m, 1)
7 a_m, mu_m, std_m = actor(obs, deterministic=False)
8
9 # approximate OB
10 q_m = critic(obs, a_m)
11 ob = optimal_baseline(a_m, q_m, mu_m, std_m)
12
13 # use OB to construct the loss
14 X = q - ob
15 loss = -(X * log_pi).mean()
```

E Computations for the Numerical Toy Example

Here we prove that the quantities in table are filled properly.

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

Proof. In this proof, for convenience, the multiplication and exponentiation of vectors is element-wise. Firstly, we trivially obtain the column $\pi_{\theta}^i(a^i)$, by taking softmax over $\psi_{\theta}^i(a^i)$. This allows us to compute the counterfactual baseline of COMA, which is

$$\begin{aligned}\hat{Q}^{-i}(a^{-i}) &= \mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\hat{Q}(a^{-i}, a^i) \right] = \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{Q}(a^{-i}, a^i) \\ &= 0.8 \times 2 + 0.1 \times 1 + 0.1 \times 100 = 1.6 + 0.1 + 10 = 11.7\end{aligned}$$

By subtracting this value from the column $\hat{Q}(a^{-i}, a^i)$, we obtain the column $\hat{A}^i(a^{-i}, a^i)$.

Let us now compute the column of $x_{\psi_{\theta}^i}^i$. For this, we use Remark 7. We have

$$\|\pi_{\theta}^i\|^2 = 0.8^2 + 0.1^2 + 0.1^2 = 0.66$$

and $1 + \|\pi_{\theta}^i\|^2 - 2\pi_{\theta}^i(a^i) = 1.66 - 2\pi_{\theta}^i(a^i)$, which is 0.06 for $a^i = 1$, and 1.46 when $a^i = 2, 3$. For $a^i = 1$, we have that

$$\pi_{\theta}^i(a^i) \left(1 + \|\pi_{\theta}^i\|^2 - 2\pi_{\theta}^i(a^i) \right) = 0.8 \times 0.06 = 0.048$$

and for $a^i = 2, 3$, we have

$$\pi_{\theta}^i(a^i) \left(1 + \|\pi_{\theta}^i\|^2 - 2\pi_{\theta}^i(a^i) \right) = 0.1 \times 1.46 = 0.146$$

We obtain the column $x_{\psi_{\theta}^i}^i(a^i)$ by normalising the vector (0.048, 0.146, 0.146). Now, we can compute OB, which is the dot product of the columns $x_{\psi_{\theta}^i}^i(a^i)$ and $\hat{Q}(a^{-i}, a^i)$

$$b^*(a^{-i}) = 0.14 \times 2 + 0.43 \times 1 + 0.43 \times 100 = 0.28 + 0.43 + 43 = 43.71$$

We obtain the column $\hat{X}^i(a^{-i}, a^i)$ after subtracting $b^*(a^{-i})$ from the column $\hat{Q}(a^{-i}, a^i)$.

Now, we can compute and compare the variances of vanilla MAPG, COMA, and OB. The surrogate local variance of an MAPG estimator $\mathbf{g}^i(b)$ is

$$\begin{aligned}\mathbf{Var}_{a^i \sim \pi_{\theta}^i} \left[\mathbf{g}^i(b) \right] &= \mathbf{Var}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}^i(a^{-i}, a^i) - b(a^{-i}) \right) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right] \\ &= \text{sum} \left(\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\left(\left[\hat{Q}^i(a^{-i}, a^i) - b(a^{-i}) \right] \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right)^2 \right] - \mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\left(\hat{Q}^i(a^{-i}, a^i) - b(a^{-i}) \right) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right]^2 \right) \\ &= \text{sum} \left(\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\left(\left[\hat{Q}^i(a^{-i}, a^i) - b(a^{-i}) \right] \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right)^2 \right] \right) - \text{sum} \left(\mathbb{E}_{a^i \sim \pi_{\theta}^i} \left[\hat{Q}^i(a^{-i}, a^i) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right]^2 \right)\end{aligned}$$

where ‘‘sum’’ is taken element-wise. The last equality follows by linearity of element-wise summing, and the fact that b is a baseline. We compute the variance of vanilla MAPG ($\mathbf{g}_{\text{MAPG}}^i$), COMA ($\mathbf{g}_{\text{COMA}}^i$),

and OB (\mathbf{g}_X^i). Let us derive the first moment, which is the same for all methods

$$\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(\mathbf{a}^{-i}, \mathbf{a}^i) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right] = \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{Q}(\mathbf{a}^{-i}, a^i) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i)$$

recalling Equation 25

$$\begin{aligned} &= \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{Q}(\mathbf{a}^{-i}, a^i) (\mathbf{e}_{a^i} - \pi_{\theta}^i) \\ &= 0.8 \times 2 \times \begin{bmatrix} 0.2 \\ -0.1 \\ -0.1 \end{bmatrix} + 0.1 \times 1 \times \begin{bmatrix} -0.8 \\ 0.9 \\ -0.1 \end{bmatrix} + 0.1 \times 100 \times \begin{bmatrix} -0.8 \\ -0.1 \\ 0.9 \end{bmatrix} \\ &= \begin{bmatrix} 0.32 \\ -0.16 \\ -0.16 \end{bmatrix} + \begin{bmatrix} -0.08 \\ 0.09 \\ -0.01 \end{bmatrix} + \begin{bmatrix} -8 \\ -1 \\ 9 \end{bmatrix} = \begin{bmatrix} -7.76 \\ -1.07 \\ 8.83 \end{bmatrix} \end{aligned}$$

Now, let's compute the second moment for each of the methods, starting from vanilla MAPG

$$\begin{aligned} &\mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right)^2 \right] \\ &= \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{Q}(\mathbf{a}^{-i}, a^i)^2 \left(\nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right)^2 \\ &= \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{Q}(\mathbf{a}^{-i}, a^i)^2 (\mathbf{e}_{a^i} - \pi_{\theta}^i)^2 \\ &= 0.8 \times 2^2 \times \begin{bmatrix} 0.2 \\ -0.1 \\ -0.1 \end{bmatrix}^2 + 0.1 \times 1^2 \times \begin{bmatrix} -0.8 \\ 0.9 \\ -0.1 \end{bmatrix}^2 + 0.1 \times 100^2 \times \begin{bmatrix} -0.8 \\ -0.1 \\ 0.9 \end{bmatrix}^2 \\ &= 0.8 \times 4 \times \begin{bmatrix} 0.04 \\ 0.01 \\ 0.01 \end{bmatrix} + 0.1 \times \begin{bmatrix} 0.64 \\ 0.81 \\ 0.01 \end{bmatrix} + 0.1 \times 10000 \times \begin{bmatrix} 0.64 \\ 0.01 \\ 0.81 \end{bmatrix} \\ &= \begin{bmatrix} 0.128 \\ 0.032 \\ 0.032 \end{bmatrix} + \begin{bmatrix} 0.064 \\ 0.081 \\ 0.001 \end{bmatrix} + \begin{bmatrix} 640 \\ 10 \\ 810 \end{bmatrix} = \begin{bmatrix} 640.192 \\ 10.113 \\ 810.033 \end{bmatrix} \end{aligned}$$

We have

$$\begin{aligned} &\mathbf{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\mathbf{g}_{\text{MAPG}}^i] \\ &= \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right)^2 \right] \\ &\quad - \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{Q}(\mathbf{a}^{-i}, \mathbf{a}^i) \nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right]^2 \\ &= \begin{bmatrix} 640.192 \\ 10.113 \\ 810.033 \end{bmatrix} - \begin{bmatrix} -7.76 \\ -1.07 \\ 8.83 \end{bmatrix}^2 = \begin{bmatrix} 640.192 \\ 10.113 \\ 810.033 \end{bmatrix} - \begin{bmatrix} 60.2176 \\ 1.1449 \\ 77.9689 \end{bmatrix} = \begin{bmatrix} 579.9744 \\ 8.968 \\ 732.064 \end{bmatrix} \end{aligned}$$

So the variance of vanilla MAPG in this case is

$$\mathbf{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\mathbf{g}_{\text{MAPG}}^i] = 1321.007$$

Let's now deal with COMA

$$\begin{aligned}
& \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{A}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right)^2 \right] \\
&= \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{A}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right)^2 \\
&= \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{A}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 (e_{a^i} - \pi_{\theta}^i)^2 \\
&= 0.8 \times (-9.7)^2 \times \begin{bmatrix} 0.2 \\ -0.1 \\ -0.1 \end{bmatrix}^2 + 0.1 \times (-10.7)^2 \times \begin{bmatrix} -0.8 \\ 0.9 \\ -0.1 \end{bmatrix}^2 + 0.1 \times 88.3^2 \times \begin{bmatrix} -0.8 \\ -0.1 \\ 0.9 \end{bmatrix}^2 \\
&= 0.8 \times 94.09 \times \begin{bmatrix} 0.04 \\ 0.01 \\ 0.01 \end{bmatrix} + 0.1 \times 114.49 \times \begin{bmatrix} 0.64 \\ 0.81 \\ 0.01 \end{bmatrix} + 0.1 \times 7796.89 \times \begin{bmatrix} 0.64 \\ 0.01 \\ 0.81 \end{bmatrix} \\
&= \begin{bmatrix} 3.011 \\ 0.753 \\ 0.753 \end{bmatrix} + \begin{bmatrix} 2.327 \\ 9.274 \\ 0.114 \end{bmatrix} + \begin{bmatrix} 499.001 \\ 7.797 \\ 631.548 \end{bmatrix} = \begin{bmatrix} 504.339 \\ 17.824 \\ 632.415 \end{bmatrix}
\end{aligned}$$

We have

$$\begin{aligned}
& \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\mathbf{g}_{\text{COMA}}^i] \\
&= \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{A}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right)^2 \right] - \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{A}^i(\mathbf{a}^{-i}, \mathbf{a}^i) \nabla_{\theta^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right]^2 \\
&= \begin{bmatrix} 504.339 \\ 17.824 \\ 632.415 \end{bmatrix} - \begin{bmatrix} -7.76 \\ -1.07 \\ 8.83 \end{bmatrix}^2 = \begin{bmatrix} 504.339 \\ 17.824 \\ 632.415 \end{bmatrix} - \begin{bmatrix} 60.2176 \\ 1.1449 \\ 77.9689 \end{bmatrix} = \begin{bmatrix} 444.1214 \\ 16.6791 \\ 554.4461 \end{bmatrix}
\end{aligned}$$

and we have

$$\mathbf{Var}_{\mathbf{a}^i \sim \pi_{\theta}^i} [\mathbf{g}_{\text{COMA}}^i] = 1015.2466$$

Lastly, we figure out OB

$$\begin{aligned}
& \mathbb{E}_{\mathbf{a}^i \sim \pi_{\theta}^i} \left[\hat{X}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(\mathbf{a}^i) \right)^2 \right] \\
&= \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{X}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\psi_{\theta}^i} \log \pi_{\theta}^i(a^i) \right)^2 \\
&= \sum_{a^i=1}^3 \pi_{\theta}^i(a^i) \hat{X}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 (e_{a^i} - \pi_{\theta}^i)^2 \\
&= 0.8 \times (-41.71)^2 \times \begin{bmatrix} 0.2 \\ -0.1 \\ -0.1 \end{bmatrix}^2 + 0.1 \times (-42.71)^2 \times \begin{bmatrix} -0.8 \\ 0.9 \\ -0.1 \end{bmatrix}^2 + 0.1 \times 56.29^2 \times \begin{bmatrix} -0.8 \\ -0.1 \\ 0.9 \end{bmatrix}^2 \\
&= 0.8 \times 1739.724 \times \begin{bmatrix} 0.04 \\ 0.01 \\ 0.01 \end{bmatrix} + 0.1 \times 1824.144 \times \begin{bmatrix} 0.64 \\ 0.81 \\ 0.01 \end{bmatrix} + 0.1 \times 3168.564 \times \begin{bmatrix} 0.64 \\ 0.01 \\ 0.81 \end{bmatrix} \\
&= \begin{bmatrix} 55.6712 \\ 13.92 \\ 13.92 \end{bmatrix} + \begin{bmatrix} 116.7452 \\ 147.756 \\ 1.824 \end{bmatrix} + \begin{bmatrix} 202.788 \\ 3.169 \\ 256.654 \end{bmatrix} = \begin{bmatrix} 375.2044 \\ 164.845 \\ 272.398 \end{bmatrix}
\end{aligned}$$

We have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} [\mathbf{g}_X^i] \\
 &= \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{X}^i(\mathbf{a}^{-i}, \mathbf{a}^i)^2 \left(\nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i) \right)^2 \right] - \mathbb{E}_{\mathbf{a}^i \sim \pi_\theta^i} \left[\hat{X}^i(\mathbf{a}^{-i}, \mathbf{a}^i) \nabla_{\psi_\theta^i} \log \pi_\theta^i(\mathbf{a}^i) \right]^2 \\
 &= \begin{bmatrix} 375.2044 \\ 164.845 \\ 272.398 \end{bmatrix} - \begin{bmatrix} -7.76 \\ -1.07 \\ 8.83 \end{bmatrix}^2 = \begin{bmatrix} 375.2044 \\ 164.845 \\ 272.398 \end{bmatrix} - \begin{bmatrix} 60.2176 \\ 1.1449 \\ 77.9689 \end{bmatrix} = \begin{bmatrix} 314.987 \\ 163.7 \\ 194.429 \end{bmatrix}
 \end{aligned}$$

and we have

$$\mathbf{Var}_{\mathbf{a}^i \sim \pi_\theta^i} [\mathbf{g}_X^i] = 673.116$$

□

F Detailed Hyper-parameter Settings for Experiments

In this section, we include the details of our experiments. Their implementations can be found in the following codebase:

<https://github.com/morning9393/Optimal-Baseline-for-Multi-agent-Policy-Gradients>.

In COMA experiments, we use the official implementation in their codebase [7]. The only difference between COMA with and without OB is the baseline introduced, that is, the OB or the counterfactual baseline of COMA [7].

Hyper-parameters used for COMA in the SMAC domain.

Hyper-parameters	3m	8m	2s3z
actor lr	5e-3	1e-2	1e-2
critic lr	5e-4	5e-4	5e-4
gamma	0.99	0.99	0.99
epsilon start	0.5	0.5	0.5
epsilon finish	0.01	0.01	0.01
epsilon anneal time	50000	50000	50000
batch size	8	8	8
buffer size	8	8	8
target update interval	200	200	200
optimizer	RMSProp	RMSProp	RMSProp
optim alpha	0.99	0.99	0.99
optim eps	1e-5	1e-5	1e-5
grad norm clip	10	10	10
actor network	rnn	rnn	rnn
rnn hidden dim	64	64	64
critic hidden layer	1	1	1
critic hidden dim	128	128	128
activation	ReLU	ReLU	ReLU
eval episodes	32	32	32

As for Multi-agent PPO, based on the official implementation [45], the original V-based critic is replaced by Q-based critic for OB calculation. Simultaneously, we have not used V-based tricks like the GAE estimator, when either using OB or state value as baselines, for fair comparisons.

We provide the pseudocode of our implementation of Multi-agent PPO with OB. We highlight the novel components of it (those unrepresent, for example, in [45]) in colour.

Algorithm 1 Multi-agent PPO with Q-critic and OB

```

1: Initialize  $\theta$  and  $\phi$ , the parameters for actor  $\pi$  and critic  $Q$ 
2:  $episode_{max} \leftarrow step_{max}/batch\_size$ 
3: while  $episode \leq episode_{max}$  do
4:   Set data buffer  $D = \{\}$ 
5:   Get initial states  $s_0$  and observations  $o_0$ 
6:   for  $t = 0$  to  $batch\_size$  do
7:     for all agents  $i$  do
8:       if discrete action space then
9:          $a_t^i, p_{\pi,t}^i \leftarrow \pi(o_t^i; \theta)$  // where  $p_{\pi,t}^i$  is the probability distribution of available actions
10:      else if continuous action space then
11:         $a_t^i, p_{a,t}^i \leftarrow \pi(o_t^i; \theta)$  // where  $p_{a,t}^i$  is the probability density of action  $a_t^i$ 
12:      end if
13:       $q_t^i \leftarrow Q(s_t, i, a_t^i; \phi)$ 
14:    end for
15:     $s_{t+1}, o_{t+1}^n, r_t \leftarrow \text{execute} \{a_t^1 \dots a_t^n\}$ 
16:    if discrete action space then
17:      Append  $[s_t, o_t, a_t, r_t, s_{t+1}, o_{t+1}, q_t, p_{\pi,t}]$  to  $D$ 
18:    else if continuous action space then
19:      Append  $[s_t, o_t, a_t, r_t, s_{t+1}, o_{t+1}, q_t, p_{a,t}]$  to  $D$ 
20:    end if
21:  end for
  // from now all agents are processed in parallel in  $D$ 
22:  if discrete action space then
23:     $ob \leftarrow \text{optimal\_baseline}(q, p_{\pi})$  // use data from  $D$ 
24:  else if continuous action space then
25:    Resample  $a_{t,1\dots m}, q_{t,1\dots m} \sim \mu_t, \sigma_t$  for each  $s_t, o_t$ 
26:     $ob \leftarrow \text{optimal\_baseline}(a, q, \mu, \sigma)$  // use resampled data
27:  end if
28:   $X \leftarrow q - ob$ 
29:   $\text{Loss}(\theta) \leftarrow -\text{mean}(X \cdot \log p_a)$ 
30:  Update  $\theta$  with Adam/RMSProp to minimise  $\text{Loss}(\theta)$ 
31: end while

```

The critic parameter ϕ is trained with TD-learning [34].

Hyper-parameters used for Multi-agent PPO in the SMAC domain.

Hyper-parameters	3s vs 5z / 5m vs 6m / 6h vs 8z / 27m vs 30m
actor lr	1e-3
critic lr	5e-4
gamma	0.99
batch size	3200
num mini batch	1
ppo epoch	10
ppo clip param	0.2
entropy coef	0.01
optimizer	Adam
opti eps	1e-5
max grad norm	10
actor network	mlp
hidden layper	1
hidden layer dim	64
activation	ReLU
gain	0.01
eval episodes	32
use huber loss	True
rollout threads	32
episode length	100

Hyper-parameters used for Multi-agent PPO in the Multi-Agent MuJoCo domain.

Hyper-parameters	Hopper(3x1)	Swimmer(2x1)	HalfCheetah(6x1)	Walker(2x3)
actor lr	5e-6	5e-5	5e-6	1e-5
critic lr	5e-3	5e-3	5e-3	5e-3
lr decay	1	1	0.99	1
episode limit	1000	1000	1000	1000
std x coef	1	10	5	5
std y coef	0.5	0.45	0.5	0.5
ob n actions	1000	1000	1000	1000
gamma	0.99	0.99	0.99	0.99
batch size	4000	4000	4000	4000
num mini batch	40	40	40	40
ppo epoch	5	5	5	5
ppo clip param	0.2	0.2	0.2	0.2
entropy coef	0.001	0.001	0.001	0.001
optimizer	RMSProp	RMSProp	RMSProp	RMSProp
momentum	0.9	0.9	0.9	0.9
opti eps	1e-5	1e-5	1e-5	1e-5
max grad norm	0.5	0.5	0.5	0.5
actor network	mlp	mlp	mlp	mlp
hidden layper	2	2	2	2
hidden layer dim	32	32	32	32
activation	ReLU	ReLU	ReLU	ReLU
gain	0.01	0.01	0.01	0.01
eval episodes	10	10	10	10
use huber loss	True	True	True	True
rollout threads	4	4	4	4
episode length	1000	1000	1000	1000

For QMIX and COMIX baseline algorithms, we use implementation from their official codebases and keep the performance consistent with the results reported in their original papers [21, 24]. MADDPG is provided along with COMIX, which is derived from its official implementation as well [15].

Hyper-parameters used for QMIX baseline in the SMAC domain.

Hyper-parameters	3s vs 5z / 5m vs 6m / 6h vs 8z / 27m vs 30m
critic lr	5e-4
gamma	0.99
epsilon start	1
epsilon finish	0.05
epsilon anneal time	50000
batch size	32
buffer size	5000
target update interval	200
double q	True
optimizer	RMSProp
optim alpha	0.99
optim eps	1e-5
grad norm clip	10
mixing embed dim	32
hypernet layers	2
hypernet embed	64
critic hidden layer	1
critic hiddem dim	128
activation	ReLU
eval episodes	32

Hyper-parameters used for COMIX baseline in the Multi-Agent MuJoCo domain.

Hyper-parameters	Hopper(3x1) / Swimmer(2x1) / HalfCheetah(6x1) / Walker(2x3)
critic lr	0.001
gamma	0.99
episode limit	1000
exploration mode	Gaussian
start steps	10000
act noise	0.1
batch size	100
buffer size	1e6
soft target update	True
target update tau	0.001
optimizer	Adam
optim eps	0.01
grad norm clip	0.5
mixing embed dim	64
hypernet layers	2
hypernet embed	64
critic hidden layer	2
critic hiddem dim	[400, 300]
activation	ReLU
eval episodes	10

Hyper-parameters used for MADDPG baseline in the Multi-Agent MuJoCo domain.

Hyper-parameters	Hopper(3x1) / Swimmer(2x1) / HalfCheetah(6x1) / Walker(2x3)
actor lr	0.001
critic lr	0.001
gamma	0.99
episode limit	1000
exploration mode	Gaussian
start steps	10000
act noise	0.1
batch size	100
buffer size	1e6
soft target update	True
target update tau	0.001
optimizer	Adam
optim eps	0.01
grad norm clip	0.5
mixing embed dim	64
hypernet layers	2
hypernet embed	64
actor network	mlp
hidden layer	2
hidдем dim	[400, 300]
activation	ReLU
eval episodes	10