

## A Algorithms

The single-agent version of ICQ is shown in Algorithm 1. Its multi-agent version counterpart (ICQ-MA) is shown in Algorithm 2.

---

### Algorithm 1: Implicit Constraint Q-Learning in Single-Agent Tasks.

---

**Input:** Offline buffer  $\mathcal{B}$ , target network update rate  $d$ .

Initialize critic network  $Q^\pi(\cdot; \phi)$  and actor network  $\pi(\cdot; \theta)$  with random parameters.

Initialize target networks:  $\phi' = \phi, \theta' = \theta$ .

**for**  $t = 1$  **to**  $T$  **do**

    Sample trajectories from  $\mathcal{B}$ .

    Train policy according to  $\mathcal{J}_\pi(\theta) = \mathbb{E}_{\tau \sim \mathcal{B}} \left[ -\frac{1}{Z(\tau)} \log(\pi(a \mid \tau; \theta)) \exp \left( \frac{Q^\pi(\tau, a)}{\alpha} \right) \right]$ .

    Train critic according to

$$\mathcal{J}_Q(\phi) = \mathbb{E}_{\tau \sim \mathcal{B}} \left[ r + \gamma \frac{1}{Z(\tau')} \exp \left( \frac{Q(\tau', a'; \phi')}{\alpha} \right) Q(\tau', a'; \phi') - Q(\tau, a; \phi) \right]^2.$$

**if**  $t \bmod d = 0$  **then**

        Update target networks:  $\phi' = \phi, \theta' = \theta$ .

**end**

**end**

---



---

### Algorithm 2: Implicit Constraint Q-Learning in Multi-Agent Tasks.

---

**Input:** Offline buffer  $\mathcal{B}$ , target network update rate  $d$ .

Initialize critic networks  $Q^i(\cdot; \phi_i)$ , actor networks  $\pi^i(\cdot; \theta_i)$  and Mixer network  $M(\cdot; \psi)$  with random parameters.

Initialize target networks:  $\phi' = \phi, \theta' = \theta, \psi' = \psi$ .

**for**  $t = 1$  **to**  $T$  **do**

    Sample trajectories from  $\mathcal{B}$ .

    Train individual policy according to

$$\mathcal{J}_\pi(\theta) = \sum_i \mathbb{E}_{\tau^i, a^i \sim \mathcal{B}} \left[ -\frac{1}{Z^i(\tau^i)} \log(\pi^i(a^i \mid \tau^i; \theta_i)) \exp \left( \frac{w^i(\tau) Q^i(\tau^i, a^i)}{\alpha} \right) \right].$$

    Train critic according to  $\mathcal{J}_Q(\phi, \psi) =$

$$\mathbb{E}_{\mathcal{B}} \left[ \sum_{t \geq 0} (\gamma \lambda)^t \left[ r_t + \gamma \frac{\exp(\frac{1}{\alpha} Q(\tau_{t+1}, \mathbf{a}_{t+1}; \phi', \psi'))}{Z(\tau_{t+1}; \phi', \psi')} Q(\tau_{t+1}, \mathbf{a}_{t+1}; \phi', \psi') - Q(\tau_t, \mathbf{a}_t; \phi, \psi) \right] \right]^2.$$

**if**  $t \bmod d = 0$  **then**

        Update target networks:  $\phi' = \phi, \theta' = \theta, \psi' = \psi$ .

**end**

**end**

---

## B Detailed Proof

### B.1 Proof of Theorem 1

**Theorem 1.** *Given a deterministic MDP, the propagation of  $\epsilon_{\mathbf{b}}$  to  $\epsilon_{\mathbf{s}}$  is proportional to  $\|P_{\mathbf{s},\mathbf{u}}^\pi\|_\infty$ :*

$$\|\epsilon_{\mathbf{s}}\|_\infty \leq \frac{\gamma \|P_{\mathbf{s},\mathbf{u}}^\pi\|_\infty}{(1-\gamma) \left(1 - \gamma \|P_{\mathbf{s},\mathbf{s}}^\pi\|_\infty\right)} \|\epsilon_{\mathbf{b}}\|_\infty. \quad (21)$$

*Proof.* Based on the Remark 1 in BCQ [16], the exact form of  $\epsilon_{\text{MDP}}(\tau, a)$  is:

$$\begin{aligned} \epsilon_{\text{MDP}}(\tau, a) &= Q_M^\pi(\tau, a) - Q_{\mathcal{B}}^\pi(\tau, a) \\ &= \sum_{\tau'} (P_M(\tau' | \tau, a) - P_{\mathcal{B}}(\tau' | \tau, a)) \left( r(\tau, a, \tau') + \gamma \sum_{a'} \pi(a' | \tau') Q_{\mathcal{B}}^\pi(\tau', a') \right) \\ &\quad + P_M(\tau' | \tau, a) \gamma \sum_{a'} \pi(a' | \tau') \epsilon_{\text{MDP}}(\tau', a'), \end{aligned} \quad (22)$$

where  $P_{\mathcal{B}} = \frac{\mathcal{N}(\tau, a, \tau')}{\sum_{\tilde{\tau}} \mathcal{N}(\tau, a, \tilde{\tau})}$  and  $\mathcal{N}$  is the number of times the tuple  $(\tau, a, \tau')$  is observed in  $\mathcal{B}$ . If  $\sum_{\tilde{\tau}} \mathcal{N}(\tau, a, \tilde{\tau}) = 0$ , then  $P_{\mathcal{B}}(\tau_{\text{init}} | \tau, a) = 1$ . Since the considered MDP is deterministic, we have  $P_M(\tau' | \tau, a) - P_{\mathcal{B}}(\tau' | \tau, a) = 0$  for  $P_{\mathbf{s},\mathbf{s}}^\pi$  and  $P_{\mathbf{s},\mathbf{u}}^\pi$ . For notational simplicity, the error generated by  $P_M(\tau' | \tau, a) - P_{\mathcal{B}}(\tau' | \tau, a)$  in  $P_{\mathbf{u},\mathbf{s}}^\pi$  and  $P_{\mathbf{u},\mathbf{u}}^\pi$  is attributed to  $\epsilon_{\mathbf{b}}$  as they have the same dimension. Then, based on the extrapolation error decomposition assumption, we rewrite Equation 22 in the matrix form:

$$\begin{bmatrix} \epsilon_{\mathbf{s}} \\ \epsilon_{\mathbf{u}} \end{bmatrix} = \gamma \begin{bmatrix} P_{\mathbf{s},\mathbf{s}}^\pi & P_{\mathbf{s},\mathbf{u}}^\pi \\ P_{\mathbf{u},\mathbf{s}}^\pi & P_{\mathbf{u},\mathbf{u}}^\pi \end{bmatrix} \begin{bmatrix} \epsilon_{\mathbf{s}} \\ \epsilon_{\mathbf{u}} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \epsilon_{\mathbf{b}} \end{bmatrix}. \quad (23)$$

The result indicates that the error is the solution of a linear program with  $[0, \epsilon_{\mathbf{b}}]^T$  as the reward function. Thus, we solve this linear program and arrive at

$$\begin{bmatrix} \epsilon_{\mathbf{s}} \\ \epsilon_{\mathbf{u}} \end{bmatrix} = (I - \gamma P^\pi)^{-1} \begin{bmatrix} 0 \\ \epsilon_{\mathbf{b}} \end{bmatrix} = \begin{bmatrix} I - \gamma P_{\mathbf{s},\mathbf{s}}^\pi & -\gamma P_{\mathbf{s},\mathbf{u}}^\pi \\ -\gamma P_{\mathbf{u},\mathbf{s}}^\pi & I - \gamma P_{\mathbf{u},\mathbf{u}}^\pi \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ \epsilon_{\mathbf{b}} \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ \epsilon_{\mathbf{b}} \end{bmatrix}. \quad (24)$$

With the block matrix inverse formula, we have

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1} & -A^{-1}B(D - CA^{-1}B)^{-1} \\ -(D - CA^{-1}B)^{-1}CA^{-1} & (D - CA^{-1}B)^{-1} \end{bmatrix}. \quad (25)$$

Since  $(D - CA^{-1}B)^{-1}$  is just the lower right block of  $(I - \gamma P^\pi)^{-1}$ , we have

$$\|(D - CA^{-1}B)^{-1}\|_\infty \leq \|(I - \gamma P^\pi)^{-1}\|_\infty \leq \frac{1}{1-\gamma}. \quad (26)$$

Thus, we obtain

$$\begin{aligned} \|-A^{-1}B(D - CA^{-1}B)^{-1}\|_\infty &\leq \|A^{-1}\|_\infty \|-B\|_\infty \|(D - CA^{-1}B)^{-1}\|_\infty \\ &\leq \frac{1}{1-\gamma} \|A^{-1}\|_\infty \|-B\|_\infty \\ &= \frac{1}{1-\gamma} \|(I - \gamma P_{\mathbf{s},\mathbf{s}}^\pi)^{-1}\|_\infty \|\gamma P_{\mathbf{s},\mathbf{u}}^\pi\|_\infty \\ &\leq \frac{\gamma \|P_{\mathbf{s},\mathbf{u}}^\pi\|_\infty}{(1-\gamma) \left(1 - \gamma \|P_{\mathbf{s},\mathbf{s}}^\pi\|_\infty\right)}. \end{aligned} \quad (27)$$

Plugging the result into Equation 25, we finish our proof at

$$\|\epsilon_{\mathbf{s}}\|_\infty \leq \|-A^{-1}B(D - CA^{-1}B)^{-1}\|_\infty \|\epsilon_{\mathbf{b}}\|_\infty \leq \frac{\gamma \|P_{\mathbf{s},\mathbf{u}}^\pi\|_\infty}{(1-\gamma) \left(1 - \gamma \|P_{\mathbf{s},\mathbf{s}}^\pi\|_\infty\right)} \|\epsilon_{\mathbf{b}}\|_\infty. \quad (28)$$

□

## B.2 Proof of Theorem 2

The proof of our Theorem 2 is based on the Theorem 3 in [46]. The main difference is that we consider a behavior policy to regularize the softmax operation. All the actions considered in the analysis are batch-constrained, thus  $\mu(a \mid \tau) > 0, \forall \tau, a$  in the proof.

**Lemma 1.** *By assuming  $f_\alpha^T(Q(\tau, \cdot))Q(\tau, \cdot)$  as target value of the Implicit Constraint Q-learning operator, we have  $\forall Q, \max_{a \sim \mathcal{B}} Q(\tau, a) - f_\alpha^T(Q(\tau, \cdot))Q(\tau, \cdot) \leq (|A_\tau| - 1) \max\{\frac{1}{(\frac{1}{\alpha} + 1)C + 1}, \frac{2Q_{\max}}{1 + C \exp(\frac{1}{\alpha})}\}$ , where  $Q_{\max} = \frac{R_{\max}}{1 - \gamma}$  represents the maximum Q-value in Q-iteration with  $\mathcal{T}_{\text{ICQ}}$ .*

*Proof.* The target value operation of Implicit Constraint Q-learning is defined as:

$$f_\alpha(\tau \mid \mu) = \frac{[\mu_1 \exp(\frac{1}{\alpha} \tau_1), \mu_2 \exp(\frac{1}{\alpha} \tau_2), \dots, \mu_{|A_\tau|} \exp(\frac{1}{\alpha} \tau_{|A_\tau|})]^T}{\sum_{i=1}^{|A_\tau|} \mu_i \exp(\frac{1}{\alpha} \tau_i)}, \quad (29)$$

We first sort the sequence  $\{Q(\tau, a_i)\}$  such that  $Q(\tau, a_{[1]}) \geq \dots \geq Q(\tau, a_{[|A_\tau|]})$ . Then,  $\forall Q$  and  $\forall \tau$ , we have that the distance between optimal Q-value and Implicit Constraint Q-value is:

$$\begin{aligned} & \max_{a \sim \mathcal{B}} Q(\tau, a) - f_\alpha^T(Q(\tau, \cdot) \mid \mu(\cdot \mid \tau))Q(\tau, \cdot) \\ &= Q(\tau, a_{[1]}) - \frac{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[\frac{1}{\alpha} Q(\tau, a_{[i]})] Q(\tau, a_{[i]})}{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[\frac{1}{\alpha} Q(\tau, a_{[i]})]} \\ &= \frac{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[\frac{1}{\alpha} Q(\tau, a_{[i]})] [Q(\tau, a_{[1]}) - Q(\tau, a_{[i]})]}{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[\frac{1}{\alpha} Q(\tau, a_{[i]})]}. \end{aligned} \quad (30)$$

Let  $\delta_i(\tau) = Q(\tau, a_{[1]}) - Q(\tau, a_{[i]})$ . The distance in the Equation 30 can be rewritten as:

$$\begin{aligned} & \frac{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[\frac{1}{\alpha} Q(\tau, a_{[i]})] [Q(\tau, a_{[1]}) - Q(\tau, a_{[i]})]}{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[\frac{1}{\alpha} Q(\tau, a_{[i]})]} \\ &= \frac{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)] \delta_i(\tau)}{\sum_{i=1}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)]} \\ &= \frac{\sum_{i=2}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)] \delta_i(\tau)}{\mu(a_{[1]} \mid \tau) + \sum_{i=2}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)]} \end{aligned} \quad (31)$$

First note that for any two non-negative sequences  $\{x_i\}$  and  $\{y_i\}$ ,

$$\frac{\sum_i x_i}{1 + \sum_i y_i} \leq \sum_i \frac{x_i}{1 + y_i}. \quad (32)$$

We have the following conclusion by applying the Equation 32 to Equation 31:

$$\begin{aligned} & \frac{\sum_{i=2}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)] \delta_i(\tau)}{\mu(a_{[1]} \mid \tau) + \sum_{i=2}^{|A_\tau|} \mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)]} \leq \sum_{i=2}^{|A_\tau|} \frac{\mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)] \delta_i(\tau)}{\mu(a_{[1]} \mid \tau) + \mu(a_{[i]} \mid \tau) \exp[-\frac{1}{\alpha} \delta_i(\tau)]} \\ &= \sum_{i=2}^{|A_\tau|} \frac{\mu(a_{[i]} \mid \tau) \delta_i(\tau)}{\mu(a_{[i]} \mid \tau) + \mu(a_{[1]} \mid \tau) \exp[\frac{1}{\alpha} \delta_i(\tau)]} \\ &= \sum_{i=2}^{|A_\tau|} \frac{\delta_i(\tau)}{1 + \frac{\mu(a_{[1]} \mid \tau)}{\mu(a_{[i]} \mid \tau)} \exp[\frac{1}{\alpha} \delta_i(\tau)]} \\ &\leq \sum_{i=2}^{|A_\tau|} \frac{\delta_i(\tau)}{1 + C \exp[\frac{1}{\alpha} \delta_i(\tau)]}, \end{aligned} \quad (33)$$

where  $C = \inf_{\tau \in S} \inf_{2 \leq i \leq |A_\tau|} \frac{\mu(a_{[1]}|\tau)}{\mu(a_{[i]}|\tau)}$ .

If  $\delta_i(\tau) > 1$ , we have

$$\frac{\delta_i(\tau)}{1 + C \exp\left[\frac{1}{\alpha}\delta_i(\tau)\right]} \leq \frac{\delta_i(\tau)}{1 + C \exp\left(\frac{1}{\alpha}\right)} \leq \frac{2Q_{\max}}{1 + C \exp\left(\frac{1}{\alpha}\right)}. \quad (34)$$

else  $0 \leq \delta_i(\tau) \leq 1$ :

$$\frac{\delta_i(\tau)}{1 + C \exp\left[\frac{1}{\alpha}\delta_i(\tau)\right]} = \frac{1}{\frac{1+C}{\delta_i(\tau)} + \frac{1}{\alpha}C + 0.5\frac{1}{\alpha^2}\delta_i(\tau)C + \dots} \leq \frac{1}{\left(\frac{1}{\alpha} + 1\right)C + 1}. \quad (35)$$

By combining these two cases with Equation 33, we complete the proof.  $\square$

**Theorem 2.** Let  $\mathcal{T}_{\text{ICQ}}^k Q_0$  denote that the operator  $\mathcal{T}_{\text{ICQ}}$  are iteratively applied over an initial state-action value function  $Q_0$  for  $k$  times. Then, we have  $\forall(\tau, a)$ ,  $\limsup_{k \rightarrow \infty} \mathcal{T}_{\text{ICQ}}^k Q_0(\tau, a) \leq Q^*(\tau, a)$ ,

$$\liminf_{k \rightarrow \infty} \mathcal{T}_{\text{ICQ}}^k Q_0(\tau, a) \geq Q^*(\tau, a) - \frac{\gamma(|A| - 1)}{(1 - \gamma)} \max \left\{ \frac{1}{\left(\frac{1}{\alpha} + 1\right)C + 1}, \frac{2Q_{\max}}{1 + C \exp\left(\frac{1}{\alpha}\right)} \right\}, \quad (36)$$

where  $|A|$  is the action space,  $|A_\tau|$  is the action space for state  $\tau$ ,  $C \triangleq \inf_{\tau \in S} \inf_{2 \leq i \leq |A_\tau|} \frac{\mu(a_{[1]}|\tau)}{\mu(a_{[i]}|\tau)}$  and  $\mu(a_{[1]}|\tau)$  denotes the probability of choosing the expert action according to behavioral policy  $\mu$ . Moreover, the upper bound of  $\mathcal{T}_{\text{BCQ}}^k Q_0 - \mathcal{T}_{\text{ICQ}}^k Q_0$  decays exponentially fast in terms of  $\alpha$ .

*Proof.* We first conjecture that

$$\mathcal{T}_{\text{BCQ}}^k Q_0(\tau, a) - \mathcal{T}_{\text{ICQ}}^k Q_0(\tau, a) \leq \sum_{j=1}^k \gamma^j \zeta, \quad (37)$$

where  $\zeta = \sup_Q \max_\tau [\max_{a \sim \mathcal{B}} Q(\tau, a) - f_\alpha^T(Q(\tau, \cdot)) Q(\tau, \cdot)]$  denotes the supremum of the difference between the BCQ and ICQ operators, over all  $Q$ -functions that occur during  $Q$ -iteration, and state  $\tau$ . Equation 37 is proven using induction as follows:

- When  $i = 1$ , we start from the definitions for  $\mathcal{T}_{\text{BCQ}}$  and  $\mathcal{T}_{\text{ICQ}}$ , and proceed as

$$\begin{aligned} \mathcal{T}_{\text{BCQ}} Q_0(\tau, a) - \mathcal{T}_{\text{ICQ}} Q_0(\tau, a) &= \gamma \sum_{\tau'} P(\tau' | \tau, a) \left[ \max_{a' \sim \mathcal{B}} Q_0(\tau', a') - f_\alpha^T(Q_0(\tau', \cdot)) Q_0(\tau', \cdot) \right] \\ &\leq \gamma \sum_{\tau'} P(\tau' | \tau, a) \zeta = \gamma \zeta. \end{aligned} \quad (38)$$

- Suppose the conjecture holds when  $i = l$ , i.e.,  $\mathcal{T}_{\text{BCQ}}^l Q_0(\tau, a) - \mathcal{T}_{\text{ICQ}}^l Q_0(\tau, a) \leq \sum_{j=1}^l \gamma^j \zeta$ , then

$$\begin{aligned} \mathcal{T}_{\text{BCQ}}^{l+1} Q_0(\tau, a) - \mathcal{T}_{\text{ICQ}}^{l+1} Q_0(\tau, a) &= \mathcal{T}_{\text{BCQ}} \mathcal{T}_{\text{BCQ}}^l Q_0(\tau, a) - \mathcal{T}_{\text{ICQ}}^{l+1} Q_0(\tau, a) \\ &\leq \mathcal{T}_{\text{BCQ}} \left[ \mathcal{T}_{\text{ICQ}}^l Q_0(\tau, a) + \sum_{j=1}^l \gamma^j \zeta \right] - \mathcal{T}_{\text{ICQ}}^{l+1} Q_0(\tau, a) \\ &= \sum_{j=1}^l \gamma^{j+1} \zeta + (\mathcal{T}_{\text{BCQ}} - \mathcal{T}_{\text{ICQ}}) \mathcal{T}_{\text{ICQ}}^l Q_0(\tau, a) \\ &\leq \sum_{j=1}^l \gamma^{j+1} \zeta + \gamma \zeta = \sum_{j=1}^{l+1} \gamma^j \zeta. \end{aligned} \quad (39)$$

By using the fact that  $\lim_{k \rightarrow \infty} \mathcal{T}_{\text{BCQ}}^k Q_0(\tau, a)$  and applying Lemma 1 to bound  $\zeta$ , we have  $\forall(\tau, a)$ ,  $\limsup_{k \rightarrow \infty} \mathcal{T}_{\text{ICQ}}^k Q_0(\tau, a) \leq Q^*(\tau, a)$  and  $\liminf_{k \rightarrow \infty} \mathcal{T}_{\text{ICQ}}^k Q_0(\tau, a) \geq Q^*(\tau, a) - \frac{\gamma(|A|-1)}{(1-\gamma)} \max\{\frac{1}{(\frac{1}{\alpha}+1)C+1}, \frac{2Q_{\max}}{1+C \exp(\frac{1}{\alpha})}\}$ . Based on the Equation 33, we can bound Equation 37 as:

$$\mathcal{T}_{\text{BCQ}}^k Q_0(\tau, a) - \mathcal{T}_{\text{ICQ}}^k Q_0(\tau, a) \leq \frac{\gamma(1-\gamma^k)}{1-\gamma} \sum_{i=2}^{|A|} \frac{\delta_i(\tau)}{1+C \exp[\frac{1}{\alpha}\delta_i(\tau)]}. \quad (40)$$

From the definition of  $\delta_i(\tau)$ , we have  $\delta_{|A_\tau|}(\tau) \geq \delta_{|A_\tau|-1}(\tau) \geq \dots \geq \delta_2(\tau) \geq 0$ . Furthermore, there must exist an index  $i^* \leq |A_\tau|$  such that  $\delta_i > 0, \forall i^* \leq i \leq |A_\tau|$ . Therefore, we can proceed from Equation 40 as

$$\begin{aligned} \frac{\gamma(1-\gamma^k)}{1-\gamma} \sum_{i=2}^{|A|} \frac{\delta_i(\tau)}{1+C \exp[\frac{1}{\alpha}\delta_i(\tau)]} &= \frac{\gamma(1-\gamma^k)}{1-\gamma} \sum_{i=i^*}^{|A|} \frac{\delta_i(\tau)}{1+C \exp[\frac{1}{\alpha}\delta_i(\tau)]} \\ &\leq \frac{\gamma(1-\gamma^k)}{1-\gamma} \sum_{i=i^*}^{|A|} \frac{\delta_i(\tau)}{C \exp[\frac{1}{\alpha}\delta_i(\tau)]} \leq \frac{\gamma(1-\gamma^k)}{1-\gamma} \sum_{i=i^*}^{|A|} \frac{\delta_i(\tau)}{C \exp[\frac{1}{\alpha}\delta_{i^*}(\tau)]} \\ &= \frac{\gamma(1-\gamma^k)}{1-\gamma} \exp\left[-\frac{1}{\alpha}\delta_{i^*}(\tau)\right] \sum_{i=i^*}^{|A|} \frac{\delta_i(\tau)}{C}, \end{aligned} \quad (41)$$

which implies an exponential convergence rate in terms of  $\alpha$ .  $\square$

### B.3 Proof of Remark 3.2

We analyze the MMDP experimental result in Section 3.2 from the perspective of the concentrability coefficient  $C(\Pi)$ , which illustrates the degree to which states and actions are out of distribution. In the MMDP case, we theoretically prove  $C(\Pi^i)$  satisfies:  $C(\Pi^1) < C(\Pi^2) < \dots < C(\Pi^n)$ , where  $\Pi^i$  denotes the set of joint policies including  $i$  agents. As illustrated in the above conclusion, the increase in the number of agents makes the distribution shift issue more severe in the MMDP case.

**Remark 1.** Let  $\varrho(s)$  denote the marginal distribution over  $S$ ,  $\rho_0$  indicate the initial state distribution, and  $\Pi^i$  represent the set of joint policies including  $i$  agents. Assume there exist coefficients  $c(k)$  satisfying  $\rho_0 P^{\pi_1} P^{\pi_2} \dots P^{\pi_k}(s) \leq c(k)\varrho(s)$ . We define the concentrability coefficient  $C(\Pi) \triangleq (1-\gamma)^2 \sum_{k=1}^{\infty} k\gamma^{k-1}c(k)$ , which illustrates the degree to which states and actions are out of distribution. Due to the limited datasets, the number of seen state-action pairs  $m$  is fixed. Then,  $C(\Pi^i)$  is monotonically increasing with the number of agents

$$C(\Pi^1) < C(\Pi^2) < \dots < C(\Pi^n) \quad (42)$$

*Proof.* We first note that  $c(k) \geq \frac{\rho_0 P^{\pi_1} P^{\pi_2} \dots P^{\pi_k}(s)}{\varrho(s)}$  and  $c(k)$  determines the value of  $C(\Pi^i)$ . To compare  $C(\Pi^i)$ , we just need to compare  $c(k)$  at iteration  $k$ . For clarity of analysis, we assume each state-action pair is visited only once, and individual policy is random  $\pi^i(A^{(i)}|s) = \frac{1}{2}$ . In the MMDP case, the transition matrix  $P^\pi$  is stable for the number of agents:

$$P^{\pi_1} = P^{\pi_k} = P^{\pi_1 \pi_2 \dots \pi_k} = \begin{bmatrix} 1 & 0 \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}. \quad (43)$$

For this reason,  $\rho_0 P^{\pi_1} P^{\pi_2} \dots P^{\pi_k}(s)$  does not change with the number of agents. As  $\varrho(s) = \sum_a \varrho(s, a) = \sum_a \frac{\sum_{s, a \in \mathcal{D}} \mathbf{1}[s=s, a=a]}{\sum_{s', a' \in \mathcal{D}} \mathbf{1}[s=s', a'=a']}$ , we can calculate  $\varrho(s)$  by counting state-action pairs in  $\mathcal{D}$  as follows

$$\varrho(s) = \frac{m}{2^{n+1}}. \quad (44)$$

The gradient of  $\varrho(s)$  is:

$$\varrho(s)' = \left(\frac{m}{2^{n+1}}\right)' = \frac{-m \cdot 2^{n+1} \ln 2}{(2^{n+1})^2} < 0. \quad (45)$$

Therefore,  $c(k)$  is monotonically increasing with the number of agents and  $C(\Pi^1) < C(\Pi^2) < \dots < C(\Pi^n)$ .  $\square$

#### B.4 Proof of Remark 2

**Remark 2.** For the optimization problem

$$\pi_{k+1} = \arg \max_{\pi} \mathbb{E}_{a \sim \pi(\cdot|\tau)}[Q^{\pi_k}(\tau, a)] \quad \text{s.t.} \quad D_{\text{KL}}(\pi \|\mu)[\tau] \leq \epsilon, \quad \sum_a \pi(a|\tau) = 1, \quad (46)$$

$$\text{the optimal policy is } \pi_{k+1}^*(a|\tau) = \frac{\mu(a|\tau) \exp(\frac{1}{\alpha} Q^{\pi_k}(\tau, a))}{\sum_{\tilde{a}} \mu(\tilde{a}|\tau) \exp(\frac{1}{\alpha} Q^{\pi_k}(\tau, \tilde{a}))}.$$

*Proof.* First, note the objective is a linear function of the decision variables  $\pi$ . All constraints are convex functions. Thus Equation 46 is a convex optimization problem. The Lagrangian equation is

$$\mathcal{L}(\pi, \alpha) = \mathbb{E}_{a \sim \pi}[Q^{\pi_k}(\tau, a)] + \alpha(\epsilon - D_{\text{KL}}(\pi \|\mu)[\tau]) + \lambda \left(1 - \sum_a \pi(a|\tau)\right), \quad (47)$$

where  $\alpha > 0$  denotes the Lagrangian coefficient. Differentiate  $\pi$  to get the following formula

$$\frac{\partial \mathcal{L}}{\partial \pi} = Q^{\pi_k}(\tau, a) - \alpha \left(1 + \log \left(\frac{\pi(a|\tau)}{\mu(a|\tau)}\right)\right) - \lambda. \quad (48)$$

Setting  $\frac{\partial \mathcal{L}}{\partial \pi}$  to zero, then:

$$\begin{aligned} Q^{\pi_k}(\tau, a) - \alpha \left(1 + \log \left(\frac{\pi(a|\tau)}{\mu(a|\tau)}\right)\right) - \lambda &= 0 \\ Q^{\pi_k}(\tau, a) &= \alpha \left(1 + \log \left(\frac{\pi(a|\tau)}{\mu(a|\tau)}\right)\right) + \lambda \\ \frac{Q^{\pi_k}(\tau, a)}{\alpha} - 1 - \frac{\lambda}{\alpha} &= \log \left(\frac{\pi(a|\tau)}{\mu(a|\tau)}\right) \\ \frac{\pi(a|\tau)}{\mu(a|\tau)} &= \exp \left(\frac{Q^{\pi_k}(\tau, a)}{\alpha} - 1 - \frac{\lambda}{\alpha}\right) \\ \pi(a|\tau) &= \mu(a|\tau) \exp \left(\frac{Q^{\pi_k}(\tau, a)}{\alpha} - 1 - \frac{\lambda}{\alpha}\right) \end{aligned} \quad (49)$$

Due to the second constraint in Equation 46, the policy is a probability distribution. Therefore, we adopt  $Z$  to normalize the result by moving the constant  $\mu(a|\tau) \exp(-1 - \frac{\lambda}{\alpha})$  to  $Z$ :

$$\pi_{k+1}^*(a|\tau) = \frac{1}{Z(\tau)} \mu(a|\tau) \exp \left(\frac{Q^{\pi_k}(\tau, a)}{\alpha}\right), \quad (50)$$

where  $Z(\tau) = \sum_{\tilde{a}} \mu(\tilde{a}|\tau) \exp(\frac{1}{\alpha} Q^{\pi_k}(\tau, \tilde{a}))$  is the normalizing partition function.  $\square$

#### B.5 Proof of Theorem 3

**Theorem 3.** Assuming the joint action-value function is linearly decomposed, we can decompose the multi-agent joint-policy under implicit constraint as follows

$$\pi = \arg \max_{\pi^1, \dots, \pi^n} \sum_i \mathbb{E}_{\tau^i, a^i \sim \mathcal{B}} \left[ \frac{1}{Z^i(\tau^i)} \log(\pi^i(a^i|\tau^i)) \exp \left(\frac{w^i(\tau) Q^i(\tau^i, a^i)}{\alpha}\right) \right], \quad (51)$$

where  $Z^i(\tau^i) = \sum_{\tilde{a}^i} \mu^i(\tilde{a}^i|\tau^i) \exp(\frac{1}{\alpha} w^i(\tau) Q^i(\tau^i, \tilde{a}^i))$  is the normalizing partition function.

*Proof.* Let  $\mathcal{J}_{\pi}$  denote the joint-policy loss. According to the assumption,  $\mathcal{J}_{\pi}$  is written:

$$\begin{aligned} \mathcal{J}_{\pi} &= \mathbb{E}_{\tau, a \sim \mathcal{B}} \left[ -\frac{1}{Z(\tau)} \log(\pi(a|\tau)) \exp \left(\frac{1}{\alpha} Q^{\pi}(\tau, a)\right) \right] \\ &= \mathbb{E}_{\tau, a^1, \dots, a^n \sim \mathcal{B}} \left[ -\frac{1}{Z(\tau)} \left( \sum_i \log(\pi^i(a^i|\tau^i)) \right) \exp \left(\frac{1}{\alpha} \left( \sum_i w^i(\tau) Q^i(\tau^i, a^i) + b(\tau) \right)\right) \right]. \end{aligned} \quad (52)$$

The loss function  $\mathcal{J}_\pi$  is equivalent to the following form by relocating the sum operator:

$$\begin{aligned}
\mathcal{J}_\pi &= \sum_i \mathbb{E}_{\tau, a^1, \dots, a^n \sim \mathcal{B}} \left[ -\frac{1}{Z(\tau)} \log(\pi^i(a^i | \tau^i)) \exp \left( \frac{\sum_i w^i(\tau) Q^i(\tau^i, a^i) + b(\tau)}{\alpha} \right) \right] \\
&= \sum_i \mathbb{E}_{\tau, a^1, \dots, a^n \sim \mathcal{B}} \left[ -\frac{1}{Z(\tau)} \log(\pi^i(a^i | \tau^i)) \exp \left( \frac{w^i(\tau) Q^i(\tau^i, a^i)}{\alpha} \right) \right. \\
&\quad \left. \exp \left( \frac{\sum_{j \neq i} w^j(\tau) Q^j(\tau^j, a^j) + b(\tau)}{\alpha} \right) \right] \\
&= \sum_i \mathbb{E}_{\tau, a^i \sim \mathcal{B}} \mathbb{E}_{a^{j \neq i} \sim \mathcal{B}} \left[ -\frac{1}{Z(\tau)} \log(\pi^i(a^i | \tau^i)) \exp \left( \frac{w^i(\tau) Q^i(\tau^i, a^i)}{\alpha} \right) \right. \\
&\quad \left. \exp \left( \frac{\sum_{j \neq i} w^j(\tau) Q^j(\tau^j, a^j) + b(\tau)}{\alpha} \right) \right] \\
&= \sum_i \mathbb{E}_{\tau, a^i \sim \mathcal{B}} \left[ -\frac{1}{Z^i(\tau^i)} \log(\pi^i(a^i | \tau^i)) \exp \left( \frac{w^i(\tau) Q^i(\tau^i, a^i)}{\alpha} \right) \right], \\
\\
Z^i(\tau^i) &= \frac{\sum_{\tilde{a}^i} \sum_{\tilde{a}^{j \neq i}} \mu(\tilde{a} | \tau) \exp \left( \frac{1}{\alpha} w^i(\tau) Q^i(\tau^i, \tilde{a}^i) \right) \exp \left( \frac{1}{\alpha} (\sum_{j \neq i} w^j(\tau) Q^j(\tau^j, \tilde{a}^j) + b(\tau)) \right)}{\mathbb{E}_{\tilde{a}^{j \neq i} \sim \mathcal{B}} \left[ \exp \left( \frac{1}{\alpha} (\sum_{j \neq i} w^j(\tau) Q^j(\tau^j, \tilde{a}^j) + b(\tau)) \right) \right]} \\
&= \frac{\sum_{\tilde{a}^i} \sum_{\tilde{a}^{j \neq i}} \mu^i(\tilde{a}^i | \tau^i) \mu^{j \neq i}(\tilde{a}^j | \tau^j) \exp \left( \frac{1}{\alpha} w^i(\tau) Q^i(\tau^i, \tilde{a}^i) \right)}{\sum_{\tilde{a}^{j \neq i}} \mu^{j \neq i}(\tilde{a}^j | \tau^j) \exp \left( \frac{1}{\alpha} (\sum_{j \neq i} w^j(\tau) Q^j(\tau^j, \tilde{a}^j) + b(\tau)) \right)} \\
&\quad \exp \left( \frac{1}{\alpha} \left( \sum_{j \neq i} w^j(\tau) Q^j(\tau^j, \tilde{a}^j) + b(\tau) \right) \right) \\
&= \sum_{\tilde{a}^i} \mu^i(\tilde{a}^i | \tau^i) \exp \left( \frac{1}{\alpha} w^i(\tau) Q^i(\tau^i, \tilde{a}^i) \right).
\end{aligned} \tag{53}$$

(54)  $\square$

## C Additional Results

### C.1 Additional Ablation Results in StarCraft II

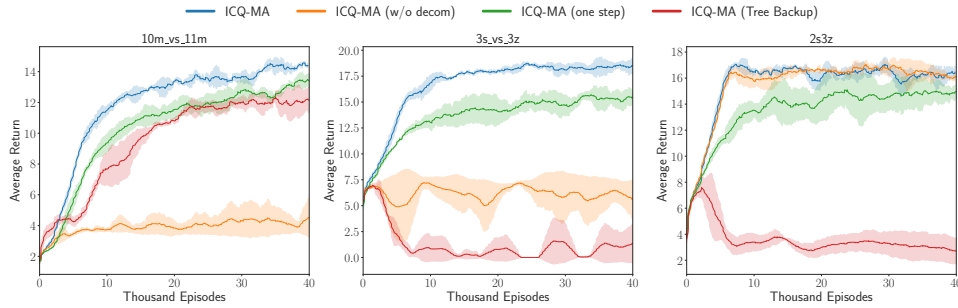


Figure 6: Module ablation study in additional StarCraft II environments.

### C.2 Additional Results in MMDP

Due to the space limits, we put the complete results in MMDP in Figure 7. BCQ gradually diverges as the number of agents increases, while ICQ has accurate estimates.

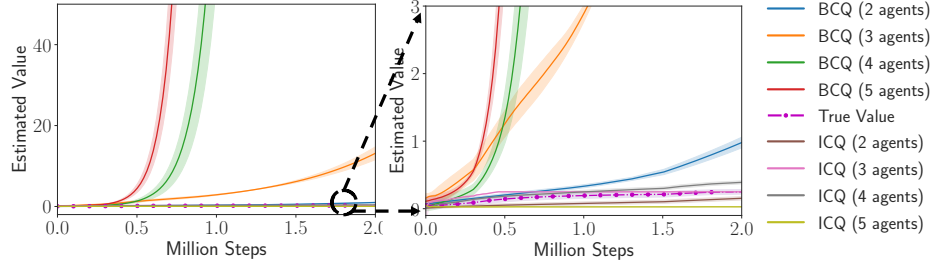


Figure 7: Additional results in MMDP.

### C.3 Additional Results in D4RL

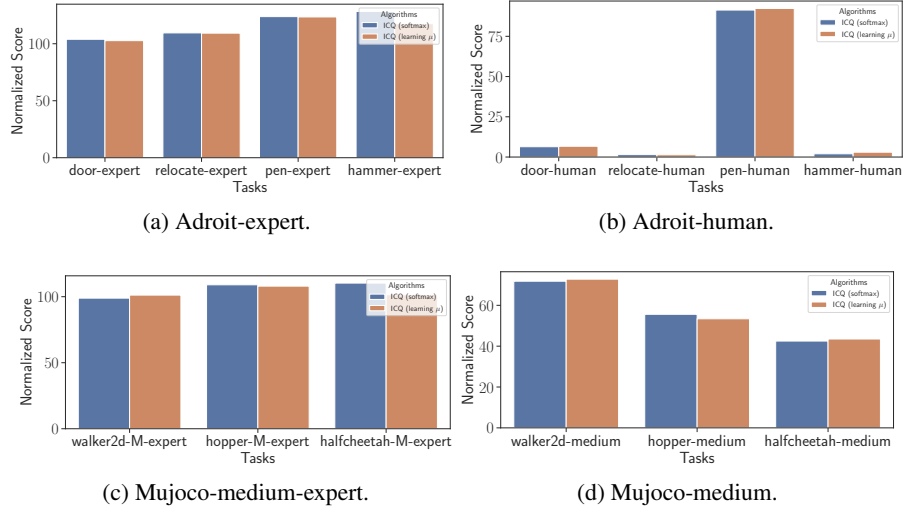


Figure 8: The performance on D4RL tasks with different implementation of ICQ.

### C.4 Ablation Study

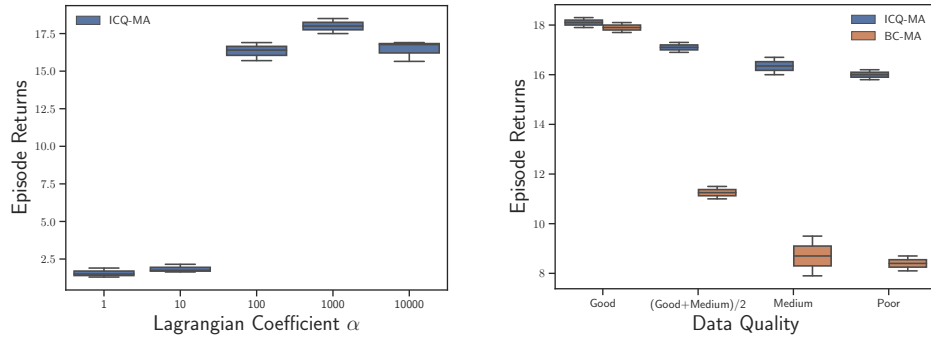


Figure 9: Ablation study on MMM map.

## D Experimental Details

### D.1 Implementation details of ICQ

We provide the two implementation options of our methods regards whether learning  $\mu$  to calculate  $\rho$ .

**Learning an auxiliary behavior model  $\hat{\mu}$ .** We first consider to learn the behavior policy  $\hat{\mu}$  using conditional variational auto-encoder as BCQ. Next, we will sample actions  $n$  times ( $n = 100$  in our experiment) from  $\hat{\mu}$  to calculate  $Z(\tau)$  on each value update:

$$\rho(\tau, a) = \frac{\exp(\frac{Q(\tau, a)}{\alpha})}{Z(\tau)} \approx \frac{\exp(\frac{Q(\tau, a)}{\alpha})}{\mathbb{E}_{\tilde{a} \sim \hat{\mu}} \exp(\frac{Q(\tau, \tilde{a})}{\alpha})}. \quad (55)$$

If  $\hat{\mu} \approx \mu$ , this method is favored as it provides an accurate approximation. However, since it may introduce unseen pairs sampled from the learned behavior model, it is against the principle of our analysis. Nevertheless, we believe it is still a better choice compared with BCQ. If there is any unseen pair  $\tau, \tilde{a}$  with large extrapolation error sampled from  $\hat{\mu}$ , e.g.  $Q_B(\tau, \tilde{a}) \gg Q_M(\tau, \tilde{a})$ , we will have  $\hat{\rho}(\tau, a) < \rho(\tau, a)$ , which means the unsafe estimation is truncated and the resulting target  $Q$ -value tends to be conservative.

**Approximate with softmax operation over a mini-batch.** We have the following measure to approximately calculate  $\rho$  without  $\mu$ , which reduces the computational complexity:

$$\rho(\tau, a) = \frac{\exp(\frac{Q(\tau, a)}{\alpha})}{Z(\tau)} \approx \frac{\exp(\frac{Q(\tau, a)}{\alpha})}{\sum_{(\tau', a') \sim \text{mini-batch}} \exp(\frac{Q(\tau', a')}{\alpha})}, \quad (56)$$

where  $Z^i(\tau^i)$  is approximated by softmax operation over mini-batch samples. The benefit of the softmax operation is that it does not include any unseen pairs, which is consistent with our theoretical analysis. However, the price is that the softmax operation ignores the difference of states over a mini-batch, which introduces an additional bias. However, we find it does not harm the performance a lot in practice. There are also some previous works using softmax to deal with the partition function, such as AWAC [29] and VMPO [45], which has been confirmed to promote performance improvement.

Considering the concise form of the softmax operation, we prefer the the second version in the multi-agent tasks. We conduct ablation studies of these two measures on D4RL to demonstrate their superior performance (see Figure 8).

### D.2 Baselines Details

**BCQ-MA** is trained by minimizing the following loss:

$$\mathcal{J}_Q^{\text{BCQ}}(\phi, \psi) = \mathbb{E}_{\tau \sim \mathcal{B}, a \sim \mu} \left[ \left( r(\tau, a) + \gamma \max_{\tilde{a}^{[j]}} Q^\pi(\tau', \tilde{a}^{[j]}; \phi', \psi') - Q^\pi(\tau, a; \phi, \psi) \right)^2 \right], \quad (57)$$

$$\tilde{a}^{[j]} = a^{[j]} + \xi(\tau, a^{[j]})$$

where  $Q^\pi(\tau, a) = w^i(\tau)Q^i(\tau^i, a^i) + b(\tau)$  and  $\xi(\tau, a^{[j]})$  denotes the perturbation model, which is decomposed as  $\xi^i(\tau^i, a^{i, [j]})$ . If  $\frac{a^{i, [j]} \sim G^i(\tau^i; \psi^i)}{\max_{\{a^{i, [j]} \sim G^i(\tau^i; \psi^i)\}_{j=1}^m}} \leq \zeta$  in agent  $i$ ,  $a^{i, [j]}$  is considered an unfamiliar action and  $\xi^i(\tau^i, a^{i, [j]})$  will mask  $a^{i, [j]}$  in maximizing  $Q^i$ -value operation.

**CQL-MA** is trained by minimizing the following loss:

$$\mathcal{J}_Q^{\text{CQL}}(\phi, \psi) = \alpha^{\text{CQL}} \mathbb{E}_{\tau \sim \mathcal{B}} \left[ \sum_i \log \sum_{a^i} \exp(w^i(\tau)Q^i(\tau^i, a^i) + b(\tau)) - \mathbb{E}_{a \sim \mu(a|\tau)} [Q^\pi(\tau, a)] \right]$$

$$+ \frac{1}{2} \mathbb{E}_{\tau \sim \mathcal{B}, a \sim \mu(a|\tau)} \left[ (y^{\text{CQL}} - Q^\pi(\tau, a))^2 \right]$$

$$\mathcal{J}_\pi^{\text{CQL}}(\theta) = \sum_i \mathbb{E}_{\tau^i, a^i \sim \mathcal{B}} \left[ -\log(\pi^i(a^i | \tau^i; \theta_i)) Q^i(\tau^i, a^i) \right], \quad (58)$$

where we adopt the decomposed policy gradient to train  $\pi$ , and  $y^{\text{CQL}}$  is calculated based on  $n$ -step off-policy estimation (e.g., Tree Backup algorithm). Besides,  $w^i(\tau) = w^i(\tau; \psi)$ ,  $b(\tau) = b(\tau; \psi)$  and  $Q^\pi(\tau, a) = Q^\pi(\tau, a; \phi, \psi)$ .

**BC-MA** only optimize  $\pi$  by minimizing the following loss:

$$\mathcal{J}_\pi^{\text{BC}}(\theta) = \sum_i \mathbb{E}_{\tau^i, a^i \sim \mathcal{B}} [-\log(\pi^i(a^i | \tau^i; \theta_i))]. \quad (59)$$

## E Multi-Agent Offline Dataset Based on StarCraft II

We divide maps in StarCraft II into three classifications based on difficulty (see Table 2). We divide behavior policies into three classifications based on the episode returns (see Table 3).

Table 2: Classification of maps in the dataset.

Difficulties	Maps
Easy	MMM, 2s_vs_3z, 3s_vs_3z, 3s5z, 2s3z, so_many_baneling
Hard	10m_vs_11m, 2c_vs_64zg
Super Hard	MMM2, 27m_vs_30m

Table 3: Classification of behavior policies in the dataset.

Level	Episode Returns
Good	15 ~ 20
Medium	10 ~ 15
Poor	0 ~ 10

### E.1 Hyper-parameters

Hyper-parameters in multi-agent tasks are respectively presented in Table 4. Please refer to our official code for the hyper-parameter in single-agent tasks.

Table 4: Multi-agent hyper-parameters sheet

Hyper-parameter	Value
Shared	
Policy network learning rate	$5 \times 10^{-4}$
Value network learning rate	$10^{-4}$
Optimizer	Adam
Discount factor $\gamma$	0.99
Parameters update rate $d$	600
Gradient clipping	20
Mixer network dimension	32
RNN hidden dimension	64
Activation function	ReLU
Batch size	16
Replay buffer size	$1.2 \times 10^4$
Others	
Lagrangian coefficient $\alpha$	1000 or 100
$\lambda$	0.8
$\alpha^{\text{CQL}}$	2.0
$\zeta$	0.3