

① Each Agent has observation space and continuous Action space.

Each Agent has 3 components

- ① Actor-network - Uses local observations for deterministic actions
- ② target actor-network - identical functionality for training stability
- ③ critic network - uses joint state action pair for estimate Q pairs

As critic learns the joint Q-value function over time, it sends appropriate Q-value approximations to the actor to help training.

② At each timestamp, agent stores following transition

$$(x, x', a_1, a_2, a_3, \dots, a_N, r_1, r_2, r_3, \dots, r_N)$$

we store joint state, next joint state and each of agent's received rewards.

③ Critic Updates

To update an agent's centralized critic, we use one-step lookahead TD error.

$$\begin{aligned} \mathcal{L}(O_i) &= \mathbb{E}_{x, a, r, x'} \left[ \left( Q_i^{\mu}(x, a_1, \dots, a_N) - y \right)^2 \right] \\ y &= r_i + \gamma Q_i^{\mu'}(x', a'_1, \dots, a'_N) \Big|_{a'_j = \mu'_j(O_j)} \quad \mu = \text{actor} \end{aligned}$$

④ Actor Updates

$$\nabla_{\theta_i} J(\mu_i) = \mathbb{E}_{x, a \sim D} \left[ \nabla_{\theta_i} \mu_i(a_i | O_i) \nabla_{a_i} Q_i^{\mu}(x, a_1, \dots, a_N) \Big|_{a_i = \mu_i(O_i)} \right]$$

We take gradient with respect to actor's parameters using central critic to guide us.

⑤ Policy Inference and Policy Ensembles

We can use probabilistic network & maximize the log probability of outputting another agent's observed action

$$\mathcal{L}(O_i^j) = - \mathbb{E}_{O_j, a_j} \left[ \log \hat{\mu}_i^j(a_j | O_j) + \lambda H(\hat{\mu}_i^j) \right]$$

Loss function for  $i^{\text{th}}$  agent estimating  $j^{\text{th}}$  agent's policy with an entropy regularizer  
Q value target (where  $\hat{y}$  is)

$$\hat{y} = r_i + \gamma Q_i^{\mu'}(x', \hat{\mu}_i^1(O_1), \dots, \hat{\mu}_i^N(O_N))$$

[We have removed the assumption that agent's knows each other policies]