

AC For CDec - POMDP 大规模学习算法

<https://zhuanlan.zhihu.com/p/66571753>

● 摘要

作者在多智能体系统的子集设定内提出解决问题的算法，设定为城市，具有集中的 collective influence 的特性，即互相影响很突出。与之前方法不同的是，作者处理的问题是全局 reward，不显式分配给智能体，而是通过 AC 集中式算法处理信用分配，同时快速收敛低方差。我们同时拓展不同的 reward 使用信用分配机制。作者提出的方法在全局 reward 的设定下，具有良好性能。

● 介绍

CDec-POMDP 模型用于处理协作型 MDP 模型。作者提出了一种 CDec-POMDP 模型下的 AC 算法框架，该种设定的特点是智能体的影响可能微不足道，大量智能体存在的环境中实现该算法。

● CDec-POMDP 模型定义：

$$P(s_{t+1}, d_{t+1} | s_t, d_t, a_t) = P_g(d_{t+1} | d_t, \mathcal{T}(s_t, a_t)) \prod_{m=1}^M P_l(s_{t+1}^m | s_t^m, a_t^m, \mathcal{T}(s_t, a_t), d_t)$$

模型转化概率计算如上，全局的状态 d 的条件概率乘以局部转换条件概率。中心化训练和分散式执行算法有 BicNet、COMA 和 MADDPG。这些方法在智能体数量上就难以达到效果。作者提出的方法能解决这个问题，该方法是基于联合动作的，即将其他智能体的动作考虑进去。Q 函数的参数是与智能体数量无关的，因为它将多智能体之间的互动转化为一个智能体和它邻居智能体的互动，即一对多问题转换为一对一问题。该算法将独立于一个智能体考虑对于它自己最好的动作。该算法启发于 MFG 均场理论博弈，MFG 研究由个人决策的集合而产生的种群行为。MFG 通过定义两个微分方程模拟个人价值函数的后向动力学和智能体群体的合计价值分布。尽管后向动力学等价于 bellman 方程，MFG 的主要目标是进行基于模型的规划，并推断个体密度随时间的变化。该算法是 model-base 的，不同于 model-free，作者的目标是在 TD 框架下将 MFG 投入到可计算的 Q-learning。

● 方法

联合动作的维度与智能体的数量成比例。由于智能体战略性地执行动作，因此可能有竞争有合作，而同时评估带有联合动作的 Q 就变得不可行了，因为 Q 的估计复杂度很高。减少 Q 的复杂度计算的一个方法是将联合动作 Q 分解，因为交互可以分解为两个智能体之间的相互作用，所以对于智能体 i，该智能体的 Q 估计就是它与所有交互智能体的联合动作下的 Q 值的平均值。公式如下：

$$Q^j(s, a) = \frac{1}{N^j} \sum_{k \in \mathcal{N}(j)} Q^j(s, a^j, a^k), \quad (5)$$

其中 N 表示智能体 i 的邻居智能体的数量。该方法的分解能保证在成对交互中保留全局性。

1. Mean Field 估计

使用均场理论来估计公式 5 的成对 Q 函数。这里考虑离散动作空间，动作使用 one-hot 编码。对于智能体 j 的邻居智能体的动作，使用邻居智能体平均动作和一个小波动来表征。

$$a^k = \bar{a}^j + \delta a^{j,k}, \quad \text{where } \bar{a}^j = \frac{1}{N^j} \sum_k a^k, \quad (6)$$

公式如下：

第一项是邻居智能体的平均动作表征，第二项是小波动。平均动作可以将智能体 j 的邻居智能体动作分布显式表达。根据泰勒展开理论，成对 Q 函数如果二阶可微，则总 Q 函数可以泰勒展开：

$$\begin{aligned} Q^j(s, a) &= \frac{1}{N^j} \sum_k Q^j(s, a^j, a^k) \\ &= \frac{1}{N^j} \sum_k \left[Q^j(s, a^j, \bar{a}^j) + \nabla_{\bar{a}^j} Q^j(s, a^j, \bar{a}^j) \cdot \delta a^{j,k} \right. \\ &\quad \left. + \frac{1}{2} \delta a^{j,k} \cdot \nabla_{\bar{a}^j}^2 Q^j(s, a^j, \bar{a}^j) \cdot \delta a^{j,k} \right] \\ &= Q^j(s, a^j, \bar{a}^j) + \nabla_{\bar{a}^j} Q^j(s, a^j, \bar{a}^j) \cdot \left[\frac{1}{N^j} \sum_k \delta a^{j,k} \right] \\ &\quad + \frac{1}{2N^j} \sum_k \left[\delta a^{j,k} \cdot \nabla_{\bar{a}^j}^2 Q^j(s, a^j, \bar{a}^j) \cdot \delta a^{j,k} \right] \quad (7) \\ &= Q^j(s, a^j, \bar{a}^j) + \frac{1}{2N^j} \sum_k R_{s, a^j}^j(a^k) \approx Q^j(s, a^j, \bar{a}^j), \quad (8) \end{aligned}$$

$$\text{where } R_{s, a^j}^j(a^k) \triangleq \delta a^{j,k} \cdot \nabla_{\bar{a}^j}^2 Q^j(s, a^j, \bar{a}^j) \cdot \delta a^{j,k}, \quad \bar{a}^{j,k} = \bar{a}^j + \epsilon^{j,k} \delta a^{j,k}$$

其中小波动和为 0，即 $\sum_k \delta a^k = 0$ ，对于第二项，R 的输入变量 a^k 由其他智能体决定，因此是另一个分布，因此本质是一个随机变量，事实上可以进一步证明 R 的边界在 $[-2M, 2M]$ ，在成对 Q 函数是温和的条件下，即变化幅度不太大，R 是 M 阶可微的，因此 R 影响近乎不计接近 0。附录 B 具体展示了边界证明。简化假设所有智能体具有同质性和局部性，则 R 余项互相抵消。

像图 1 所示那样，一个智能体 j 的均场估计可以被简化为一个中心智能体即网格中某一点，和抽象出的虚拟邻居代表智能体，即网格中某一点附近的蓝色趋于邻居影响的交互。

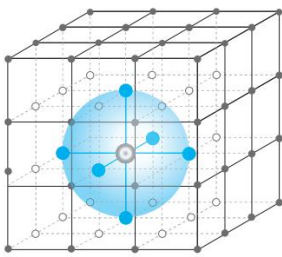


Figure 1: Mean field approximation. Each agent is represented as a node in the grid, which is only affected by the **mean** effect from its **neighbors** (the blue area). Many-agent interactions are effectively converted into two-agent interactions.

可以将交互简化成公式 8，同时均场 Q 函数的更新可以表达为：

$$Q_{t+1}^j(s, a^j, \bar{a}^j) = (1 - \alpha) Q_t^j(s, a^j, \bar{a}^j) + \alpha [r^j + \gamma v_t^j(s')], \quad (9)$$

式中均场状态的价值函数由公式 10 确定：

$$v_t^j(s') = \sum_{a^j} \pi_t^j(a^j | s', \bar{a}^j) \mathbb{E}_{\bar{a}^j(a^{-j}) \sim \pi_t^{-j}} [Q_t^j(s', a^j, \bar{a}^j)], \quad (10)$$

在博弈中，平均动作首先由前一个平均动作决定的参数化 Π 来计算的：

$$\bar{a}^j = \frac{1}{N^j} \sum_k a^k, \quad a^k \sim \pi_t^k(\cdot | s, \bar{a}^k), \quad (11)$$

策略更新依赖于平均动作 \mathbf{a} ，因此新的玻尔兹曼策略通过以下公式确认：

$$\pi_t^j(a^j | s, \bar{a}^j) = \frac{\exp(-\beta Q_t^j(s, a^j, \bar{a}^j))}{\sum_{a^{j'} \in \mathcal{A}^j} \exp(-\beta Q_t^j(s, a^{j'}, \bar{a}^j))}. \quad (12)$$

通过公式 11 和公式 12，对于所有智能体的平均动作及其策略就可以确认。虽然直观上没法说明智能体 j 的交互能够稳定下来，但是通过实验发现，上式通过几次迭代后收敛。

与纳什 \mathbf{v} 函数不同的是，作者使用公式 10 的均场 \mathbf{v} 函数代替纳什 \mathbf{v} 函数，同时我们定义均场操作算子：

$$\mathcal{H}^{\text{MF}} Q(s, \mathbf{a}) = \mathbb{E}_{s' \sim p} [r(s, \mathbf{a}) + \gamma \mathbf{v}^{\text{MF}}(s')]. \quad (13)$$

事实上，我们可以证明该算子定义了一种压缩映射，那就是应用算子等价于迭代地更新 Q ，在适当条件下，均场 Q 会收敛到纳什 Q 。

2. 实现

在 MF - Q 中，loss function 定义如下：

$$\mathcal{L}(\phi^j) = (y^j - Q_{\phi^j}(s, a^j, \bar{a}^j))^2, \quad y^j = r^j + \gamma \mathbf{v}_{\phi_-^j}^{\text{MF}}(s')$$

其中 y 由目标网络确定，而 loss function 的微分如下：

$$\nabla_{\phi^j} \mathcal{L}(\phi^j) = (y^j - Q_{\phi^j}(s, a^j, \bar{a}^j)) \nabla_{\phi^j} Q_{\phi^j}(s, a^j, \bar{a}^j), \quad (14)$$

Mean Field Q-learning 算法如下：

Algorithm 1 Mean Field Q-learning (MF-Q)

Initialise Q_{ϕ^j} , $Q_{\phi_-^j}$, and \bar{a}^j for all $j \in \{1, \dots, N\}$
while training not finished **do**
 for $m = 1, \dots, M$ **do**
 For each agent j , sample action a^j from Q_{ϕ^j} by Eq. (12), with the current mean action \bar{a}^j and the exploration rate β
 For each agent j , compute the new mean action \bar{a}^j by Eq. (11)
 Take the joint action $\mathbf{a} = [a^1, \dots, a^N]$ and observe the reward $\mathbf{r} = [r^1, \dots, r^N]$ and the next state s'
 Store $\langle s, \mathbf{a}, \mathbf{r}, s', \bar{\mathbf{a}} \rangle$ in replay buffer \mathcal{D} , where $\bar{\mathbf{a}} = [\bar{a}^1, \dots, \bar{a}^N]$
 for $j = 1$ to N **do**
 Sample a minibatch of K experiences $\langle s, \mathbf{a}, \mathbf{r}, s', \bar{\mathbf{a}} \rangle$ from \mathcal{D}
 Sample action a_-^j from $Q_{\phi_-^j}$ with $\bar{a}_-^j \leftarrow \bar{a}^j$
 Set $y^j = r^j + \gamma \mathbf{v}_{\phi_-^j}^{\text{MF}}(s')$ by Eq. (10)
 Update the Q -network by minimizing the loss $\mathcal{L}(\phi^j) = \frac{1}{K} \sum (y^j - Q_{\phi^j}(s^j, a^j, \bar{a}^j))^2$
 Update the parameters of the target network for each agent j with learning rate τ :

$$\phi_-^j \leftarrow \tau \phi^j + (1 - \tau) \phi_-^j$$

而 MF - AC 算法的 actor 又可以通过以下公式训练：

$$\nabla_{\theta^j} \mathcal{J}(\theta^j) \approx \nabla_{\theta^j} \log \pi_{\theta^j}(s) Q_{\phi^j}(s, a^j, \bar{a}^j) \Big|_{a=\pi_{\theta^j}(s)}.$$

而 critic 的训练 loss function 可以通过公式 14 来确定。

Mean Field AC 算法如下:

Algorithm 2 Mean Field Actor-Critic (MF-AC)

Initialize Q_{ϕ^j} , $Q_{\phi_-^j}$, π_{θ^j} , $\pi_{\theta_-^j}$, and \bar{a}^j for all $j \in \{1, \dots, N\}$

while training not finished **do**

For each agent j , sample action $a^j = \pi_{\theta^j}(s)$; compute the new mean action $\bar{a} = [\bar{a}^1, \dots, \bar{a}^N]$

Take the joint action $\mathbf{a} = [a^1, \dots, a^N]$ and observe the reward $\mathbf{r} = [r^1, \dots, r^N]$ and the next state s'

Store $\langle s, \mathbf{a}, \mathbf{r}, s', \bar{\mathbf{a}} \rangle$ in replay buffer \mathcal{D}

for $j = 1$ to N **do**

Sample a minibatch of K experiences $\langle s, \mathbf{a}, \mathbf{r}, s', \bar{\mathbf{a}} \rangle$ from \mathcal{D}

Set $y^j = r^j + \gamma v_{\phi_-^j}^{\text{MF}}(s')$ by Eq. (10)

Update the critic by minimizing the loss $\mathcal{L}(\phi^j) = \frac{1}{K} \sum (y^j - Q_{\phi^j}(s, a^j, \bar{a}^j))^2$

Update the actor using the sampled policy gradient:

$$\nabla_{\theta^j} \mathcal{J}(\theta^j) \approx \frac{1}{K} \sum \nabla_{\theta^j} \log \pi_{\theta^j}(s') Q_{\phi_-^j}(s', a_-^j, \bar{a}_-^j) \Big|_{a_-^j = \pi_{\theta_-^j}(s')}$$

Update the parameters of the target networks for each agent j with learning rates τ_ϕ and τ_θ :

$$\begin{aligned} \phi_-^j &\leftarrow \tau_\phi \phi^j + (1 - \tau_\phi) \phi_-^j \\ \theta_-^j &\leftarrow \tau_\theta \theta^j + (1 - \tau_\theta) \theta_-^j \end{aligned}$$

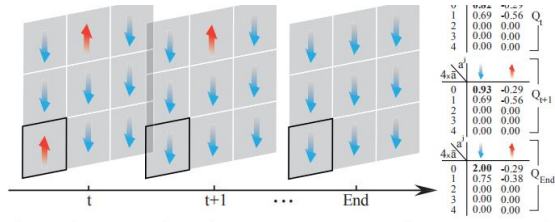


Figure 2: MF- Q iterations on a 3×3 stateless toy example. The goal is to coordinate the agents to an agreed direction. Each agent has two choices of actions: *up* \uparrow or *down* \downarrow . The reward of each agent's staying in the same direction as its $[0, 1, 2, 3, 4]$ neighbors are $[-2.0, -1.0, 0.0, 1.0, 2.0]$, respectively. The neighbors are specified by the four directions on the grid with cyclic structure on all directions, *e.g.* the first row and the third row are adjacent. The reward for the highlighted agent j on the bottom left at time $t+1$ is 2.0, as all neighboring agents stay down in the same time. We listed the Q -tables for agent j at three time steps where \bar{a}^j is the percentage of neighboring ups. Following Eq. 9, we have $Q_{t+1}^j(\uparrow, \bar{a}^j = 0) = Q_t^j(\uparrow, \bar{a}^j = 0) + \alpha[r^j - Q_t^j(\uparrow, \bar{a}^j = 0)] = 0.82 + 0.1 \times (2.0 - 0.82) = 0.93$. The rightmost plot shows the convergent scenario where the Q -value of staying down is 2.0, which is the largest reward in the environment.

图二为 MF - Q 的迭代过程示意图。

3. 收敛性证明

**

● 实验

1. 环境
2. 与非联合动作强化学习方法进行比较

3. 其他智能体学习策略的影响
4. 使用一组策略训练策略的影响

- 总结

作者方法模拟了 **MA** 环境中的交互动力学。**MF Q-learning** 学习每一个智能体从其他智能体得到的均场效应的最佳响应。作者提出将一对多问题转化为一对一，即该智能体与邻居智能体的问题。文中也证明了该算法收敛到 **nash Q** 值。