

plain_text (0.96)

the k-mean algorithm. The iterative formula for updating the cluster center is shown below [43].

isolate_formula (0.87)

$$a_k = \frac{\sum_{i=1}^n z_{ik} x_i}{\sum_{i=1}^n z_{ik}} \quad (1)$$

isolate_formula (0.89)

$$z_{ik} = \begin{cases} 1 & \text{if } \|x_i - a_k\|^2 = \min_{l \leq k \leq c} \|x_i - a_l\|^2 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

plain_text (0.96)

where $\|x_i - a_k\|$ is the Euclidean distance between x_i and the cluster center a_k .

title (0.91)

3.2 Reinforcement learning

plain_text (0.98)

When faced with a problem, different actions are adopted to influence the outcome by perceiving the environment and maximizing the benefits of learning by interaction, which is also referred to as reinforcement learning. Current reinforcement learning algorithms utilize neural networks to extract high-dimensional characteristics from observed data and represent their strategies or value functions as function approximators.

The agent receives state s_t at each time step t , and selects an action a_t from a collection of actions \mathcal{A} according to the policy model π , where π is a map of states to actions a_t the next state s_{t+1} and scalar reward r_t as return. The whole process runs continuously to maximize rewards $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$, in which γ is the discount factor ($0 < \gamma < 1$) and r_t is the reward at step t . In value-based model-free reinforcement learning methods, the action-value function is represented by a neural network and other function approximators.

title (0.91)

3.2.1 A2C and A3C

plain_text (0.98)

A2C (Advantage Actor-critic) combines policy-based and value-based approaches, which maintain a policy $\pi(a_t | s_t; \theta)$ and an estimate of the value function $V(s_t; \theta)$. Policies $\pi(a_t | s_t; \theta)$ and $V(s_t; \theta)$ are mainly defined as the policy network and critic network, respectively. θ can be optimized by the gradient ascent method as follows [44].

isolate_formula (0.94)

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (r_t^n + V_\pi(s_{t+1}^n) - V_\pi(s_t^n)) \nabla \log p_\theta(a_t^n | s_t^n) \quad (3)$$

plain_text (0.98)

A2C can learn from only one work with the environment, which differs from A3C (Asynchronous Advantage Actor-critic). A3C operates asynchronously with a global network and multiple branch networks learning from multiple works. Each branch network can be understood as a replica of the global network. Asynchronous training updates

plain_text (0.96)

multiple branch network parameters to uninterruptedly update the global network parameters. The updates performed can be seen below [45].

isolate_formula (0.91)

$$\nabla_{\theta'} \log \pi(a_t | s_t; \theta') A(s_t, a_t; \theta, \theta_v) + \beta \nabla_{\theta'} H(\pi(s_t; \theta_v)) \quad (4)$$

plain_text (0.97)

$\beta \nabla_{\theta'} H(\pi(s_t; \theta'))$ is the entropy, which can increase exploration of the environment and prevent premature convergence leading to suboptimal results. $A(s_t, a_t; \theta, \theta_v)$ is the advantage function, which is given as follows [45]

isolate_formula (0.95)

$$A(s_t, a_t; \theta, \theta_v) = \sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta_v) - V(s_t; \theta_v) \quad (5)$$

plain_text (0.96)

The parameter θ_v of the value network is then updated by the gradient descent of

isolate_formula (0.95)

$$\nabla_{\theta_v} \left[\left(\sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta_v) - V(s_t; \theta_v) \right) \right] \quad (6)$$

title (0.91)

3.2.2 Parallel Advantage Actor-Critic (PAAC)

plain_text (0.97)

A general framework for deep reinforcement learning proposed by Clemente [42], and multiple participants can be trained synchronously on one machine.

Gradients in PAAC are calculated using mini batch experiences. Multiple environment instances are trained in parallel and may explore different states at any given time. The advantage is reducing the relevance of the state encountered and contributing to stable training [45].

The updated formulas for policy gradients ∇_{θ}^π and value gradients ∇_{θ}^V are shown below [42]

isolate_formula (0.94)

$$\nabla_{\theta}^\pi \approx \frac{1}{n_e \cdot t_{\max}} \sum_{e=1}^{n_e} \sum_{t=1}^{t_{\max}} \left(Q^{(t_{\max}-t+1)}(s_{e,t}, a_{e,t}; \theta, \theta_v) - V(s_{e,t}; \theta_v) \right) \nabla_{\theta} \log \pi(a_{e,t} | s_{e,t}; \theta) + \beta \nabla_{\theta} H(\pi(s_{e,t}; \theta)) \quad (7)$$

isolate_formula (0.94)

$$\nabla_{\theta_v}^V \approx \nabla_{\theta_v} \frac{1}{n_e \cdot t_{\max}} \sum_{e=1}^{n_e} \sum_{t=1}^{t_{\max}} \left(Q^{(t_{\max}-t+1)}(s_{e,t}, a_{e,t}; \theta, \theta_v) - V(s_{e,t}; \theta_v) \right) \quad (8)$$

title (0.95)

3.3 Parallel Advantage Actor-Critic K-means (PAAC-K)

plain_text (0.96)

The PAAC-K model consists of a parallel advantage actor-critic, K-means, reward functions, and stepwise training