

L13

□ recap : Unsupervised learning

□ Reinforcement Learning (RL)

□ Model free RL

□ Policy gradients

Supervised learning	Unsupervised learning
<ul style="list-style-type: none">— Regression— Classification	<ul style="list-style-type: none">— data visualization— Dimensionality reduction

— Data is static .

— learning is offline .

— Prediction do not affect data samples

Reinforcement Learning (RL)

— Dataset is dynamically changing

— learning is online

— Explicit modeling of feedback loops



↓ Robotics

[RL \longleftrightarrow Control theory
Dynamic System, Optimal nonlinear Control]

Setting

- Environment
- Agent
- Actions at time step
- Actions \rightarrow rewards
- Goal: Maximum rewards.

$$S_{t+1} = f(S_t, a_t)$$

↳ state transition

function / dynamical system

$$\text{Reward} = r(S_t, a_t)$$

Goal:
Agent observes

$$\tau = s, a_0, s_1, a_1, \dots, s_{T-1}$$

↳ Trajectory / rollout 轨迹
↳ Horizon

$$r(S_t, a_t) \text{ for } t=0 \dots (T-1)$$

Predict a strategy/policy for for deciding the next step action

Policy: $a_t = \pi(z_t)$
 Function z_t \nearrow previous timestamps

Apply ML to solve this prediction problems!

→ Choose representation for π (policy)

→ Choose a loss function

$$R(\pi) = \sum_{t=0}^{T-1} -r(s_t, a_t)$$

subject to $s_{t+1} = f(s_t, a_t)$

$a_t = \pi(z_t)$, so is given π

try to learn policy π

→ Optimize this function

linear models, logistic, neural nets

Deep RL

Need to figure out $\frac{\partial R}{\partial \pi}$

Idea: Assume that π is not deterministic, instead it predicts a prob distribution of actions

$\hookrightarrow a_t$ random $\hookrightarrow R$ is a random variable
 s_t random

Modified Problem $\min E_{\pi(s)} R(\pi) = E - \sum_{t=0}^{T-1} r(a_t, s_t)$
 s.t. $s_{t+1} = f(s_t, a_t)$

Training Samples: \rightarrow multiple trajectories/rollouts
 $a_t = \pi(z_t)$

Supervised $\xrightarrow{\text{only 2 per rollout (x)}}$ RL $\xrightarrow{\text{Unsupervised (x)}}$

- Q-learning
- Value iterations
- Actor Critic

Policy gradients:

assume π is a linear policy parametered by θ

$$\frac{\partial E_{\pi(\theta)} [R(z)]}{\partial \theta} = ?$$

log-derivative

$$\frac{\partial}{\partial \theta} \log \pi(z) = \frac{1}{\pi(z)} \cdot \frac{\partial}{\partial \theta} \pi(z)$$

$$\text{i.e. } \frac{\partial \pi(z)}{\partial \theta} = \pi(z) \cdot \frac{\partial}{\partial \theta} \log \pi(z)$$

$$\begin{aligned} \frac{\partial E_{\pi(\theta)} [R(z)]}{\partial \theta} &= \frac{\partial}{\partial \theta} \sum_z \pi(z) R(z) \\ &= \sum_z \pi(z) \frac{\partial}{\partial \theta} R(z) \\ &= \sum_z \underbrace{\pi(z)} \underbrace{R(z) \frac{\partial}{\partial \theta} \log \pi(z)} \end{aligned}$$

$$\nabla E_{\pi(\theta)} [R(z)] = E_{\pi(\theta)} \left[R(z) \frac{\partial}{\partial \theta} \log \pi(z) \right]$$

Gradient of expected rewards

= Expected value of Reward \times log derivative

\therefore Just like in S&D,

- sample different roll outs
- compute approximation to the gradient of expected rewards.
- update θ in that direction

Repeat

REINFORCE

(根据 $\Delta \log \pi(\tau)$)

1) Sample rollout $\tau = (s_0, a_0, \dots, s_T)$

2) Compute $k(\tau) = \frac{1}{T} \sum_{t=0}^{T-1} r(s_t, a_t)$

reward 函数

越大越好

3) $\theta \leftarrow \theta - \eta k(\tau) \frac{\partial \log \pi(\tau)}{\partial \theta}$

$k(\tau) \downarrow, \uparrow$

↳ R appears but never its gradient

↳ Useful when reward are discrete

↳ Environment f do not appear, only the rollout $\tau \dots k(\tau)$

"Model free reinforcement learning"

↳ "Model-based-RL"

Derivative free optimizer

Gradient descent $\theta \leftarrow \theta - \eta \nabla k(\theta)$

Random search

- 1. pick a random vector v
- 2. compute step size (either negative or positive) that minimize $R(\theta + \eta v)$
- 3. $\theta \leftarrow \theta + \eta v$