

# Q-learning

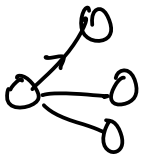
optimal control

$$\arg \min \sum_{t=0}^T J_t(x_t, u_t)$$

s.t.  $x_{t+1} = f(x_t, u_t)$   
unknown

LQR we assumed system dynamics to be known and linear

Dijkstra or A-star we assumed dynamics to be deterministic and known

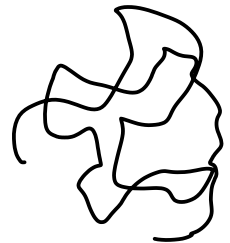


Learning function from data

$$\begin{pmatrix} x_{t+1} \\ y_{t+1} \\ \theta_{t+1} \end{pmatrix} = \begin{pmatrix} x_t + v \cos \theta_t dt \\ y_t + v \sin \theta_t dt \\ \theta_t + \omega_t dt \end{pmatrix}$$

Data

$$\mathcal{D} = \begin{cases} x_t, u_t \rightarrow x_{t+1} \\ x_1, u_1 \rightarrow x_2 \\ x_2, u_2 \rightarrow x_3 \\ \vdots \\ x_{100}, u_{100} \rightarrow x_{101} \end{cases}$$



Solution

learn the function  $f = Ax_t + Bu_t$   
To be learnt

$$\min_{A, B} \sum_{x_t, u_t, x_{t+1} \in \mathcal{D}} \|x_{t+1} - (Ax_t + Bu_t)\|_2^2$$

$f = NN \rightarrow$  LQR, Dijkstra, A-star

Model-based learning  
system dynamics

Model-free learning

In LQR

Cost to go (Value function)

$$V_T(x_T) = J_T(x_T)$$

$$V_t^*(x_t) = \sum_{k=t}^T J_k(x_k, \pi^*(x_k))$$

Optimal cost to go

$$V_t^*(x_t) = \min_{u_t} J_t(x_t, u_t) + V_{t+1}^*(x_{t+1})$$

$f(x_t, u_t)$   
LQR  
Dijkstra

Bellman equation

The main idea is called dynamic programming

$$f_i = f_{i-1} + f_{i-2}$$

Recursion inefficient

```
def fib(n):  
    if n == 1:  
        return 1  
    else:  
        return fib(n-1) + fib(n-2)
```

DP efficient

```
fibs = [1, 1]  
for i in range(n):  
    f[i+2] = f[i] + f[i+1]
```

Q-function = Action <sup>max</sup> Value function

$$V_t^*(x_t) = \sum_{k=t}^T J_k(x_k, \pi^*(x_k))$$

cost = min  
rewards = max

$$Q_t^*(x_t, u_t) = J_t(x_t, u_t) + \underbrace{\sum_{k=t+1}^T J_k(x_k, \pi^*(x_k))}_{\text{cost function}}$$

$$V_t^*(x_t) = \min_{u_t} Q_t^*(x_t, u_t)$$

$$Q_t^*(x_t, u_t) = \underbrace{J_t(x_t, u_t)}_{\text{cost function}} + V_{t+1}^*(x_{t+1})$$

policy function  $\rightarrow \pi^*(x_t) = \arg \min_{u_t} Q_t^*(x_t, u_t)$  ✓

$$= \arg \min_{u_t} \underbrace{J_t(x_t, u_t)}_{\text{functional form of the cost}} + V_{t+1}^*(x_{t+1})$$

X

$$V_t^*(x_t) = \min_{u_t} J_t(x_t, u_t) + V_{t+1}^*(x_{t+1})$$

$$Q_t^*(x_t, u_t) = J_t(x_t, u_t) + V_{t+1}^*(x_{t+1})$$

$$Q_t^*(x_t, u_t) = \underbrace{J_t(x_t, u_t)}_{\text{need not be functional}} + \min_{u_{t+1}} Q(x_{t+1}, u_{t+1})$$

$$J_t(x_t, u_t) = \|x_t - x_g\|_2^2$$

Q-learning is about learning

$$j_t \sim \underbrace{J_t(x_t, u_t)}_{\text{cost function from data}}$$

The Q-function directly from data

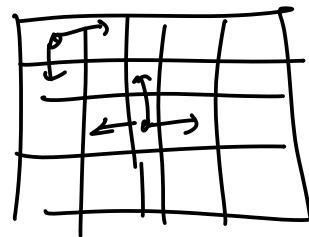
$$D \begin{cases} x_1, u_1 \rightarrow x_2, j_1 \\ x_2, u_2 \rightarrow x_3, j_2 \\ x_3, u_3 \rightarrow x_4, j_3 \end{cases}$$

$$\begin{aligned} j_1 &\sim J_1(x_1, u_1) \\ j_2 &\sim J_2(x_2, u_2) \end{aligned}$$

How to represent Q function

① Tabular Q function

$$Q \begin{bmatrix} \overbrace{0, 0, 0}^{x_t} & \underbrace{0, 1, 2, 3}_{u_t} \end{bmatrix} \rightarrow Q^*(x_t, u_t)$$



② Q-function as a NN

$$Q^*(x_t, u_t; W) = W_2^T \sigma \left( W_1 \begin{bmatrix} x_t \\ u_t \end{bmatrix} \right)$$

optimal action  $\rightarrow \pi^*(x_t) = \min_{u_t} Q(x_t, u_t)$

Q-learning algorithm

① Exploration (Random actions)

Initialize  $Q(x_t, u_t) = 0.2$  | optimistic encourages exploration  
for tabular Use bellman equations to update

$t \rightarrow 1 \dots 100$   
99% of actions to be random  $\rightarrow$  1% of actions  
 $\epsilon$ -greedy policy

② Exploitation/learn Q function

Neural Network

$$Q(x_t, u_t; W) = 0.2$$

Replay buffer + shuffle

Loss function  
gradient on this

$$L(W) = \left\| \underbrace{Q_t^*(x_t, u_t; W)}_{\text{gradient on this}} - \left( \delta_t + \min_{u_{t+1}} \underbrace{Q(x_{t+1}, u_{t+1}; W)}_{\text{freeze}} \right) \right\|_2^2$$

assume fixed  $x$

$$+ \min_{u_{t+1}} \underbrace{Q(x_{t+1}, u_{t+1}; W)}_{\text{freeze}}$$

target network

$$\boxed{W_{i+1} = W_i - \alpha \nabla_W L(W)}$$

gradient descent

