

4/11/2024: Q-Learning

But first: What is the optimal policy if we know the full MDP?
we know $T(s, a, s')$ & Rewards (s, a, s')

$V_{\text{OPT}}(s)$ = maximum possible expected discounted sum of rewards for any policy, starting at state s
"optimal value"

$Q_{\text{OPT}}(s, a)$ = maximum possible expected discounted sum of rewards for any policy, starting at state s & forced to take action a
"Q-value"

$$V_{\text{OPT}}(s) = \begin{cases} 0 & \text{if } \text{IsEnd}(s) \\ \max_{a \in \text{Actions}(s)} Q_{\text{OPT}}(s, a) & \text{else} \end{cases}$$

play the best a play optimally after choosing a

$$Q_{\text{OPT}}(s, a) = \sum_{s' \in S} T(s, a, s') \left[\text{Reward}(s, a, s') + \gamma \cdot V_{\text{OPT}}(s') \right]$$

prob of going to s' Reward now (no discounting) Future reward (discounted)

If we know $Q_{\text{OPT}}(s, a)$ for all s & a ,

$$\text{Optimal policy } \pi^*(s) = \underset{a \in \text{Actions}(s)}{\text{argmax}} Q_{\text{OPT}}(s, a)$$

Lesson: If we can estimate $Q_{\text{OPT}}(s, a)$ well, we can deduce the optimal policy!

Now: Reinforcement Learning

- We believe the world works like an MDP
- But we don't know $T(s, a, s')$ & $Reward(s, a, s')$
- Can only learn about world by interacting with it.

RL training pseudocode:

For episode $\epsilon = 1, 2, 3, \dots$:

$S_t \leftarrow S_{start}$

For $t = 1, \dots$,

- Agent choose action $a_t = \pi_{act}(S_t)$
policy used to act during learning
- Agent receives:
 - Reward r_t
 - New state S_{t+1}
- Update agent's parameters (Act learning)
- If $IS_{End}(S_{t+1})$: break

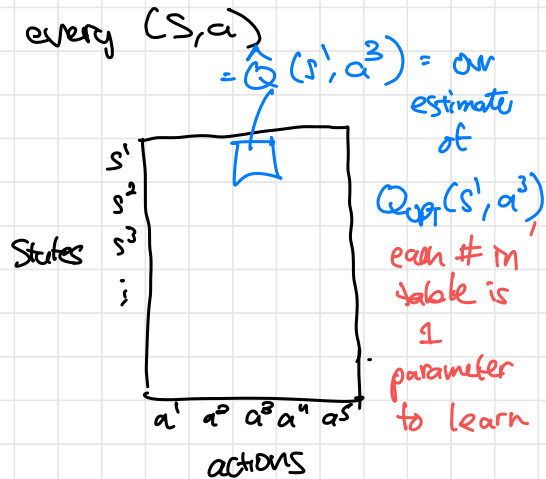
Learning algorithm: Q-Learning

Goal: Learn $Q_{opt}(s, a)$ for every (s, a)

Version 1: Tabular Q-Learning

Maintain a big table of all states \times all actions

Can initialize all to 0



How to learn $\hat{Q}(s,a)$?

We get 1 "training example" at a time consisting of (s, a, r, s')

Update Rule:

$$\hat{Q}(s,a) \leftarrow (1-\eta) \hat{Q}(s,a) + \eta (r + \gamma \hat{V}(s'))$$

Annotations:
- $(1-\eta)$: "learning rate" ($\hat{\approx} 0.1$)
- $\hat{Q}(s,a)$: "old value guess"
- r : "reward now"
- $\gamma \hat{V}(s')$: "anticipated future reward"
- The entire right-hand side is labeled "Estimate of total future reward".
- A blue line above the equation says "weighted average of old guess + new estimate of future reward".

$$\text{where } \hat{V}(s) = \begin{cases} \max_{a \in \text{Actions}(s)} \hat{Q}(s,a) & \text{if not ISEnd}(s) \\ 0 & \text{else} \end{cases}$$

How do we choose π_{act} ?

- Natural answer (wrong):

At each state s , choose $a \in \text{argmax}_{a \in \text{Actions}(s)} \hat{Q}(s,a)$
✓ optimal if $\hat{Q}(s,a) = Q_{opt}(s,a)$

✗ Bad idea early in training

Suppose we only visit s , take a , receive large reward.

$\Rightarrow \hat{Q}(s,a)$ would be large

\Rightarrow "Natural" policy would always take a in state s

All exploitation, no exploration

- Simple Fix: ϵ -Greedy

At each timestep, at state s :

- with probability $1-\epsilon$ choose $\text{argmax}_{a \in \text{Action}(s)} \hat{Q}(s,a)$
(Exploitation)

- with probability ϵ : choose random action in $\text{Actions}(s)$
(Exploration)

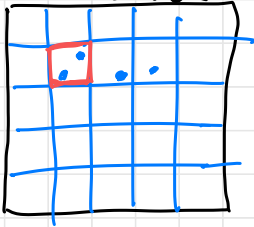
Usually $\epsilon = 0.1$ or 0.01

After training: Set $\epsilon = 0$

How to deal with large state spaces?

Option 1: Discretize state space

- E.g. state might be continuous (e.g. location of robot)
- Divide continuum into many buckets



Discretize area into
 5×5 grid = 25 states

Option 2: Version 2 of Q-learning = Q-learning with Linear Function Approximation

Idea: Q-learning is kind of like regression where
 $x = (s, a)$, $y = Q\text{-value}$

Not like supervised learning
b/c we don't know Q-values

Goal: Learn a linear model mapping (s, a) to Q-value

① Need feature function $\phi(s, a) \in \mathbb{R}^d$

② Learn parameter vector $w \in \mathbb{R}^d$
to predict $\hat{Q}(s, a) = w^T \phi(s, a)$

How to learn w ? Re-examine old Q-learning formula.

$$\begin{aligned} \hat{Q}(s, a) &\leftarrow (1-\eta) \hat{Q}(s, a) + \eta (r + \gamma \hat{V}(s')) \\ &= \hat{Q}(s, a) - \eta \left(\hat{Q}(s, a) - [r + \gamma \hat{V}(s')] \right) \end{aligned}$$

learning rate η (indicated by a red arrow pointing to the η term)

"gradient" (indicated by a red bracket under the term $\hat{Q}(s, a) - [r + \gamma \hat{V}(s')]$)

For Q-learning with Function Approx:
minimize squared error between

$\hat{Q}(s,a)$ and
"prediction"

$r + \hat{V}(s')$
"target"

$$\text{Loss}(w) = \frac{1}{2} \left(\hat{Q}(s,a) - [r + \hat{V}(s')] \right)^2$$

$$\nabla_w \text{Loss}(w) = \frac{1}{2} \cdot 2 \cdot (w^T \phi(s,a) - [r + \hat{V}(s')]) \cdot \phi(s,a)$$

so the gradient descent update is:-

$$w \leftarrow w - \eta \cdot \nabla_w \text{Loss}(w)$$

$$= w - \eta (w^T \phi(s,a) - [r + \hat{V}(s')]) \cdot \phi(s,a)$$