

example: continued

$$V^{(k)}(A) = \max_{a=\{\rightarrow, \downarrow\}} \left(P(s' = B, v | s = A, a = \rightarrow) \cdot (V_{A \rightarrow B} + V^{(k-1)}(B)) \right)$$

$$+ P(s' = C, v | s = A, a = \rightarrow) \cdot (V_{A \rightarrow C} + V^{(k-1)}(C)),$$

$$P(s' = B, v | s = A, a = \downarrow)$$

$$\cdot (V_{A \rightarrow B} + V^{(k-1)}(B))$$

$$+ P(s' = C, v | s = A, a = \downarrow)$$

$$\cdot (V_{A \rightarrow C} + V^{(k-1)}(C)))$$

$$= \max(-20 \cdot 0.8 + 0 \cdot 0.2 V^{(k-1)}(B),$$

$$- 80 \cdot 0.2 + 0 \cdot 0.8 V^{(k-1)}(B))$$

$$V^{(k)}(B) = \max(19 \cdot 0.2 + 0 \cdot 0.8 V^{(k-1)}(A),$$

$$79 \cdot 0.8 + 0 \cdot 0.2 V^{(k-1)}(A))$$

| k | $V^{(k)}(A)$ | $V^{(k)}(B)$ |
|---|--------------|--------------|
| 1 | -20 · 0.8 | 79 · 0.8 |
| 3 | 39 · 0.8 | 88 · 0.4 |
| ⋮ | | |
| D | 51 | 90 |

→ here solved non-linear system → didn't know optimal policy

Bellman equations : for optimal value function V^* , i.e. the expected return of the optimal policy π^* , when starting in state s

$$V^*(s) = \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V^*(s')]$$

solve by value iteration algorithm, without knowing the optimal policy. ↪ about optimal policy is needed

- once we know V^* , we get the optimal policy by

policy extraction

$$\forall s: \pi^*(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V^*(s')]$$

Q - function or action-value function

- it is often better to consider value function and policy together

→ action-value function $Q(\pi(s,a))$

| | |
|---------------------------------|------------------------------------|
| state-value function $V_\pi(s)$ | action value function $Q_\pi(s,a)$ |
|---------------------------------|------------------------------------|

- expect return if we follow policy π and start in states

$$V_\pi(s) = \mathbb{E}_{\pi, t} \left[\sum_{t'=t+1}^{\infty} r^{t-t-1} R_{t'} | s_t = s \right]$$

$$= \sum_a \pi(a|s) Q_\pi(s,a)$$

- more freedom expected return if we start in state s , execute action a in the next move and only then follow policy π

$$Q_\pi(s,a) = \mathbb{E}_{\pi, t} \left[\sum_{t'=t+1}^{\infty} r^{t-t-1} R_{t'} | s_t = s, a_t = a \right]$$

can be transformed
into each other
 $a_t = a$

$$= \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi(s')]$$

- Optimal policy is $\pi^*(s) = \arg \max_a Q^*(s,a)$

$$\text{with } V^*(s) = \max_a Q^*(s,a)$$

$Q^*(s,a)$ where Q^* can be determined by the corresponding Bellman equations without knowing π^*

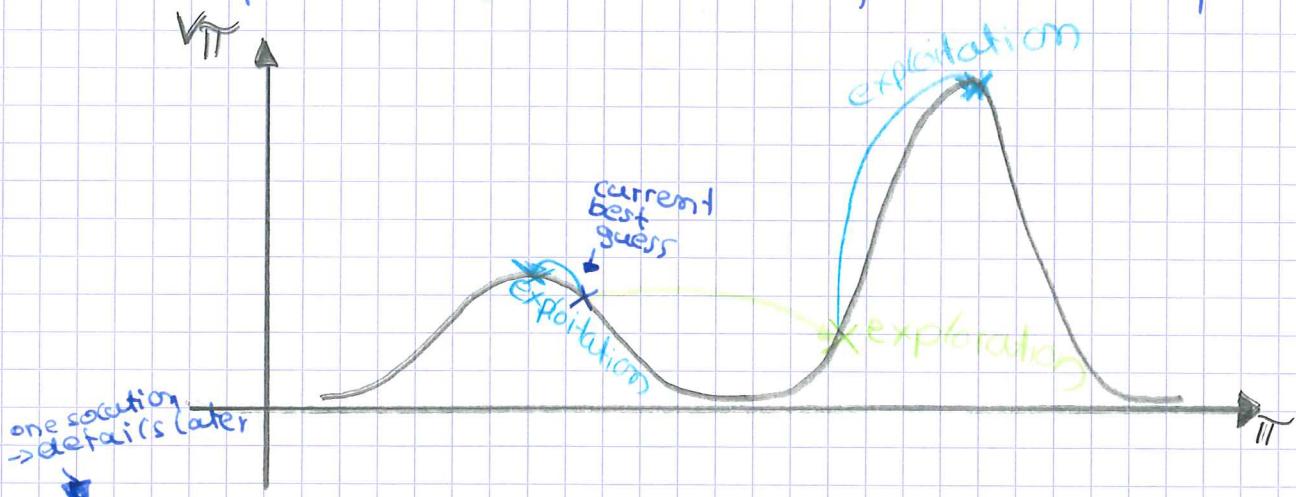
$$Q^*(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \max_a Q^*(s',a)]$$

Model-free Reinforcement Learning

we can't know
needed
knowledge of prob.

- Methods so far require knowledge of the environment behavior $p(s', r | s, a)$, i.e. a model of the environment. In general, such a model is unknown.
- two approaches
 1. model-based: first learn $p(s', r | s, a)$, then derive value functions and policies via Bellman equations
 => true understanding of the world, but very difficult
 (similar generative models in classification)
 2. learn values and policies directly, without explicit construction of $p(s', r | s, a)$ (world knowledge is implicit)
 => simpler, but lacks explicit "insight" \rightarrow no idea what it is doing
 (similar to discriminative modeling in classification)
- model free approaches must balance two goals
 - learn something about the environment
 \rightarrow exploration
 - learn useful actions \rightarrow trade-off
 \rightarrow exploitation

- exploitation: local search to improve the current policy
- exploration: global search to find an entirely new policy



perform the best known action most time, but a random action every once in a while

Theory : policy is a valid exploration/exploitation trade-off if after an infinite number of moves, every state is visited and -ii- action a is tried infinitely often.

goal: solve Bellman equations for V^* or Q^* without knowing transition prob $p(s', r | s, a)$

Monte-Carlo (MC) value estimation

- value = expected return \rightarrow estimate by empirical mean
- simplification: assume that every game $T = 1, 2, \dots$ ends after a finite number of moves T_T
 each game = one episode
 each move = experience event
- define expected return G_{T+} in step t of episode T

$$G_{T+} = \sum_{t'=t+1}^{T_T} \gamma^{t'-t-1} R_{t'+1}$$

↑
return in episode T
↑
immediate rewards after step $t-1$ in episode T

- run many episode with some exploration/exploitation policy π and compute empirical mean

$$\hat{V}_{\pi}(s) = \frac{1}{N_s} \sum_{\substack{\tau, t: s_{\tau,t}=s \\ \text{total number of visits to state } s}} G_{\tau+}$$

↑
start in state s and some the experience

- the mean can be computed incrementally

⑥ init $\hat{V}(s) = 0$, $N_s = 0$

- ⑦ repeat for all τ, t

if $s_t = s$:

$$\hat{V}(s) \leftarrow \hat{V}(s) + \frac{1}{N_s} (G_{\tau+} - \hat{V}(s))$$

↑
old guess
↑
new guess
↑
correction
learning rate

looks like gradient descent

generalized to arbitrary learning rates α

$$\hat{V}_{\tau+} \leftarrow \hat{V}_{\tau+} + \alpha (G_{\tau+} - \hat{V}_{\tau+})$$

α : hyper-parameter

- if α const \rightarrow old experience will be forgotten
 → make α smaller over time
 ↳ better in practise
 → stabilize the system

Temporal Difference (TD) value estimation

- don't calculate G_{π^*} exactly, but approximate with current guess of $\hat{V}(s_t)$

$$G_{\pi^*} = \sum_{t'=t+1}^{T^*} \gamma^{t'-t-1} R_{t'} = R_{t+1} + \gamma \sum_{t'=t+2}^{T^*} \gamma^{t'-t-2} R_{t'}$$

$\hat{V}(s_{t+1}) \approx \hat{V}(s_{t+1})$

$$\approx R_{t+1} + \gamma \hat{V}(s_{t+1})$$

\Rightarrow update rule:

$$\hat{V}(s_{t+1}) \leftarrow \hat{V}(s_{t+1}) + \alpha (R_{t+1} + \gamma \hat{V}(s_{t+2}) - \hat{V}(s_{t+1}))$$

- n -step TD update: only use current guess after n steps of the sum

$$G_{\pi^*} = R_{t+1} + R_{t+2} + \dots + \gamma^{n-1} R_{t+n-1} + \gamma^n$$

$$\cdot \sum_{t'=t+n}^{t'+n-1} \gamma^{t'-t-n} R_{t'}$$

$$G_{\pi^*} \approx \sum_{t'=t+1}^{t'+n-1} \gamma^{t'-t-1} R_{t'} + \gamma^n \hat{V}(s_{t+n})$$

update rule:

$$\hat{V}(s_{t+1}) \leftarrow \hat{V}(s_{t+1}) + \alpha \left(\sum_{t'=t+1}^{t+n-1} \gamma^{t'-t-1} R_{t'} + \gamma^n \hat{V}(s_{t+n}) - \hat{V}(s_{t+1}) \right)$$

n : hyperparameter (typical:

$$n = 3, \dots, 10$$

\rightarrow reduces the variance of our estimates for G_{π^*} , each $R_{t'}$ is used n times

In general: TD updates have lower variance than MC updates, but corrections are slightly biased due to the use of $\hat{V}(s_{t+n})$
bias-variance trade-off

Q-learning: learn Q^* instead of V^*

* best for the game • we can always improve the current policy by updating $\pi(s)$: $\hat{\pi}(s) \leftarrow \operatorname{argmax}_a \hat{Q}^{\hat{\pi}}(s, a)$

Theorem: if everything is discrete (π, Q are stored as tables), update will never make π worse.

Q-learning algorithm

① choose initial guess for $Q(s, a)$ (good guess → outer)

② for episodes $t = 1, 2, \dots$

a) create episode $s_{t+1}, a_{t+1}, r_{t+1}, s_{t+2}, a_{t+2}, r_{t+2}, \dots$ termination

using an exploration/exploitation policy (→ outer)

b) for $t = 1, \dots, T_t - 1$

$$\hat{Q}(s_{t+1}, a_{t+1}) \leftarrow \hat{Q}(s_{t+1}, a_{t+1})$$

$$+ \alpha [r_{t+1} + \gamma \hat{V}(s_{t+2}) - \hat{Q}(s_{t+1}, a_{t+1})]$$

correction

two variants: 1. α -learning:

$$\hat{V}(s_{t+1}) = \max_{a'} \hat{Q}(s_{t+1}, a')$$

2. SARSA

$$\hat{V}(s_{t+1}) = \hat{Q}(s_{t+1}, a_{t+1})$$

③ $\forall s \quad \hat{\pi}(s) = \operatorname{argmax}_a \hat{Q}(s, a)$

variants: - update Q after every move, not after entire episode

⇒ rest of the episode will use the updated policy already

- n-step updates (equal to n-step TD)

corrections:

$$\sum_{t'=t+1}^{t+n-1} \gamma^{t'-t-1} r_{t+1} + \gamma^n \hat{V}(s_{t+n}) - \hat{Q}(s_{t+1}, a_{t+1})$$

- exploration / exploitation policy:

- ϵ -greedy policy:

$$\pi(a|s) = \begin{cases} 1 & [a = a'] \text{ with } a' = \underset{a}{\operatorname{argmax}} Q(s, a) \text{ with prob } 1-\epsilon \\ \text{uniform}(a) & \text{with prob } \epsilon \end{cases}$$

ϵ : hyperparameter, anneal over time

$$\epsilon = \frac{\epsilon_0}{t(\tau, t)} \leftarrow \text{increases with } \tau \text{ and } t$$

- softmax policy at temperature ρ

$$\pi(a|s) \sim \frac{e^{Q(s, a) / \rho}}{\sum_{a'} e^{Q(s, a') / \rho}} \quad \begin{array}{l} \text{if } \rho \rightarrow 0 \\ a' = \underset{a}{\operatorname{argmax}} Q(s, a) \\ \rightarrow \text{greedy policy} \end{array}$$

$$\quad \quad \quad \begin{array}{l} \text{if } \rho \rightarrow \infty \\ a \sim \text{uniform}(a') \end{array}$$

$$\text{anneal } \rho \text{ over time } \rho = \frac{\rho_0}{t(\tau, t)} \leftarrow \text{increasing with } \tau \text{ and } t$$

- good initial guess:

• required that $Q(s, a) = 0$ if s is a terminal state

- other Q are arbitrary ≥ 0 , but good choice speed-up convergence

- "reward shaping": if rewards are very sparse ($r_a \approx 0$ most of the time), convergence is slow

⇒ use hand-crafted auxiliary rewards to guide algorithm towards good solutions, but without changing the optimal $\hat{\pi} \approx \pi^*$:

$\overset{(a)}{Q}(s, a) = \phi(s)$ ("potential"): does not prefer any action a
 $\phi(s) \approx V^*(s)$ start close to true values

