

Q-Learning

main RL method

- Return after time step t is the discounted total reward of the rest of the game

$$G_t = \sum_{t'=t+1}^T r^{t'-t-1} R_{t'}$$

- Value function: expected return in state s when we're playing according to policy π : $v_\pi(s)$

- Action value function: expected return in states, when we execute action a and then use policy π for the rest of the game: $q_\pi(s, a)$
 $= Q_\pi(s, a)$

$$v_\pi(s) = \mathbb{E}_{a \sim \pi(a|s)} [Q_\pi(s, a)]$$

- Bellmann equations for q -function:

$$Q_\pi(s, a) = \mathbb{E}_{s', R \sim p(s', R | s, a)} [R + \gamma v_\pi(s')]$$

Specifically, for the optimal policy $\pi^*(s) := \arg \max_a Q^*(s, a)$

$$Q^*(s, a) = \mathbb{E}_{s', R \sim p(s', R | s, a)} [R + \gamma \max_a Q^*(s', a)]$$

- for the Q -function of suboptimal policy π , we can get a new policy by

$$\pi'(s) = \arg \max_a Q_\pi(s, a)$$

one can prove that $\pi'(s)$ is never worse than $\pi(s) \Rightarrow$ possibility of „off-policy learning“:

games played with a sub-optimal policy π can still be used as training data for an improved policy π'

\Rightarrow iterate this over long training sessions $\Rightarrow \pi' \rightarrow \pi^*$

- goal: learn $\hat{\pi} \approx \pi^*$ and $\hat{Q} \approx Q^*$ without learning a full world model $p(s', R | s, a)$, because this is very difficult

„model-free RL“

- trick: approximate expectation in the Bellmann equations by a single instance

\Rightarrow temporal difference learning

$$\begin{aligned} Q^*(s, a) &= Q_\pi(s, a) + [Q^*(s', a) - Q_\pi(s, a)] \\ &= Q_\pi(s, a) + \underbrace{[\mathbb{E}_p [R + \gamma V^*(s')] - Q_\pi(s, a)]}_{\text{Expect over } p, \text{ the "world"} \quad \text{approximate: Expectation } \rightarrow \text{single instance}} \\ &\quad \text{V}^*(s') \rightarrow \text{our current } v_\pi(s') \end{aligned}$$

$$Q^*(s, a) \approx Q_\pi(s, a) + [R + \gamma v_\pi(s') - Q_\pi(s, a)]$$

\Rightarrow Learning rule: add a fraction of the correction (as in stochastic gradient descent)

$$Q'(s,a) \leftarrow Q_\pi(s,a) + \alpha [R + \gamma v_\pi(s') - Q_\pi(s,a)]$$

α : learning rate, hyperparameter

algorithm:

① collect training data with the current policy π

② for all moves in TS: perform above update

Q-Learning Algorithm

① initialisation $Q^{(0)}(s,a)$ (arbitrary, but good guess speeds up convergence)

important: $Q^{(0)}(s=\text{terminal state}) = 0$

② outer loop: $t = 1, \dots, E$ (play many games)

③ inner loop: $t = 1, \dots, T_t$ (time steps of game t)

④ decide about next move according to

- ϵ -greedy policy $a_{t,t} \sim \begin{cases} \underset{a}{\arg \max} Q^{(t-1)}(s_{t,t}, a) & \text{current state} \\ \text{uniform } (a) & \text{with prob } 1-\epsilon \\ \text{random} & \text{with prob } \epsilon \end{cases}$

- softmax policy $a_{t,t} \sim \text{softmax}_g(Q^{(t-1)})$ with $\underbrace{\text{hyperparam}}_g \beta$

⑤ for all $t=1, \dots, T_t$: update Q-function

$$Q^{(t)}(s_{t,t}, a_{t,t}) \leftarrow Q^{(t-1)}(s_{t,t}, a_{t,t}) + \alpha [R_{t,t+1} + \gamma V^{(t-1)}(s_{t,t+1}) - Q^{(t-1)}(s_{t,t}, a_{t,t})]$$

⑥ final policy: $\pi^*(s) = \underset{a}{\arg \max} Q^{(t)}(s, a)$

two major variants: what's the current guess for $v^{(t-1)}(s_{t,t+1})$

- Q-Learning: $v^{(t-1)}(s) = \max_a Q^{(t-1)}(s, a)$ (use more we should have done)

- SARSA algorithm: $v^{(t-1)}(s_{t,t+1}) = Q^{(t-1)}(s_{t,t+1}, a_{t,t+1})$ (use more we actually executed at time $t+1$)

more accurate variant for updates: use actual rewards of multiple time steps:

n-step Q-learning

$$Q^{(t)}(s_{\tau,t}, a_{\tau,t}) \leftarrow Q^{(t-1)}(s_{\tau,t}, a_{\tau,t}) + \alpha \left[\sum_{t'=t+1}^{t+n} r^{t'-t-1} R_{t'} + \gamma^n v^{(t-1)}(s_{\tau,t+n}) - Q^{(t-1)}(s_{\tau,t}, a_{\tau,t}) \right]$$

in practice: maintain an experience buffer = collection of moves from all games so far (possibly pruned to keep size manageable)

⇒ update step ⑥ uses a random batch from the experience buffer instead of just the previous game

⇒ better, especially when the batch is sampled with non-uniform probabilities: prefer instances where the current guess

$Q^{(t-1)}$ is very bad

$$\left[(G_{\tau,t} - Q^{(t-1)}(s_{\tau,t}, a_{\tau,t}))^2 \gg 0 \right]$$

↑
actual return ↑
current guess

Problem: If implemented naively, $\hat{Q}(s,a)$ is a huge matrix of size $|S| \times |A|$ (# of states x # of actions)

- matrix cannot be stored in memory

- we cannot collect enough training data to fit

all matrix elements



analogously to problems of multidimensional histograms

How to do Q-Learning in practice?

① easy possibility: group states into k clusters where action values are similar ⇒ $\hat{Q}(c_t, a)$ has only size $k \times |A|$

with handcrafted or learned clustering (good actions should be equal for all members of a cluster)

② reduce it to a regression problem: $\hat{Y} = f(x)$

two possibilities:

I: Y is a vector of length $|A|$, holding the expected return for every action possible in state $X = s$

II: Y is a scalar, giving the expected return of the move $X = (s, a)$

define response ($Y[a]$ for possibility I or Y for possibility II)

- MC value estimation:

$$Y_{\tau,t} = \sum_{t'=t+1}^{\infty} \gamma^{t'-t-1} R_{\tau,t'} \quad (\text{use actual observed rewards})$$

↑
Monte Carlo

- TD Q-Learning : $Y_{t,t} = R_{t,t+1} + \gamma \max_a Q^{(t-1)}(s_{t,t+1}, a)$

γ
temporal difference

- n-step TD Q-Learning : $Y_{t,t} = \sum_{t'=t+1}^{t+n} \gamma^{t'-t-1} R_{t,t'} + \gamma^n \max_a Q^{(t-1)}(s_{t,t+n}, a)$

- SARSA : $Y_{t,t} = R_{t,t+1} + \gamma Q^{(t-1)}(s_{t,t+1}, a_{t,t+1})$

• define X : I : hand crafted feature detector $X = \psi(s)$

II : use $X = s$ and have the regression learn optimal features (especially effective when $f(x)$ is a neural network)

• Simplest possibility : Linear regression

$$Q(s,a) = \psi(s) \cdot \beta_a \quad (\text{one weight vector for each action } a)$$

training by batch gradient descent:

① sample batch of training instances where action a was performed

$$\textcircled{2} \text{ update } \beta_a \leftarrow \beta_a + \frac{\alpha}{N_{\text{batch}}} \sum_{t,t' \in \text{batch}} \underbrace{\psi(s_{t,t'})^T (Y_{t,t} - \psi(s_{t,t'}) \cdot \beta_a)}_{\text{derivative of squared loss}}$$

- only works when features $\psi(s)$ are very expressive

=> use standard feature extraction methods (PCA, ICA, LLE) in combination with hand-designed preprocessing

=> exploit symmetries of the game [Bomberman: mirroring of S or rotation of s by $90^\circ, 180^\circ, 270^\circ$ are also legal states]

• more powerful : regression forest

- can train one forest for each action a or a single forest for outputting a vector of action values

=> can learn non-linear mappings from state representation $\psi(s) \rightarrow Y$

=> still needs good features $\psi(s)$, but less so than linear regression

• Deep Q-learning : typically use $X = S$ (two feature extraction) and learn optimal features in the first layers of a NN
(entire network determines values as before)

- example: ATARI game agents by Deep Mind : $X = \text{pixels of the current video frame} \stackrel{?}{=} \text{same as human players}$

- process this with a "convolutional neural network"
- advantage: learned features are superior to hand-crafted ones
- disadvantage: find a good network architecture and train to convergence can take quite long \Rightarrow often does not finish by final project deadline

Reward shaping

- auxiliary rewards during training when official rewards are very sparse
 - idea: define intermediate goals for agent
- e.g. chess: taking an opponent's piece is usually good, but not always (may be a trap or sacrifice)

\Rightarrow reward shaping must be done with care, otherwise the agent may converge to a bad policy

theory: modify the rewards by a potential function

$$R'_{\tau, t+1} = R_{\tau, t+1} + \underbrace{F(S_{\tau, t}, S_{\tau, t+1})}_{\text{reward shaping}}$$

$$F(S_{\tau, t}, S_{\tau, t+1}) = r \Phi(S_{\tau, t+1}) - \Phi(S_{\tau, t})$$

$\Phi(s)$ must not depend on the action

optimal theoretical choice: $\Phi(s) = v^*(s)$, but:

- we do not know $v^*(s)$

- is not very good during exploration



Bei dem optimalen Pfad kann jede random action zum Tod führen
↳ viele Fehlversuche, dann gestrichelten Pfad bevorzugen