# ② Deep Q-learning DQN

Published in 2013 - 2015 Natine

Approach: Q-learning $\longrightarrow$ what if q-values are represented by a neural network?

Bellman equations: $Q^*$ is the only solution to that equation

$$Q^*(s,a) = \mathbb{E}_{R,s'}\left[ R + \gamma \max_{a' \in A} Q^*(s',a') \right]$$

$$\curvearrowright \Delta(s,a)$$

target model $\theta'$        prediction model $\theta$

$$\mathcal{L}(\theta) = \frac{1}{2}\left( \mathbb{E}_{\substack{(s,a,r,s') \\ \sim \sigma_\theta}} \left[ R + \gamma \max_{a' \in A} Q_{\theta'}(s',a') - Q_\theta(s,a) \right] \right)^2$$

two occurrences of $\theta$

$\left[\begin{array}{l} \text{if } \mathcal{L}(\theta)=0 \text{ then } Q_\theta \text{ satisfies Bellman equation} \\ \quad \Rightarrow Q_\theta = Q_* \end{array}\right.$

Key idea: use two networks !

$\quad\hookrightarrow$ similarly to off·policy learning

$\quad\hookrightarrow$ similarly to the maximisation bias

$$\mathcal{L}(\theta) = \frac{1}{2}\left(\mathop{\mathbb{E}}_{\substack{(s,a,r,s') \\ \sim \sigma_\theta}}\left[\underbrace{r + \gamma \max_{a'\in A} Q_{\theta'}(s',a')}_{\text{fixed}} - Q_\theta(s,a)\right]\right)^2$$

Implementation details:

- (prioritised) experience replay

- reward clipping

DQN algorithm

initialise two models : prediction model $\theta$

target model $\theta'$

$\theta \rightsquigarrow \sigma_\theta(s) = \arg\max_{a\in A} q_\theta(s,a)$

Iterate:

constant: batch size

Batch sampling (B)

Simulate the environment using $\varepsilon$· grady from $6_\Theta$ giving full trajectories

$\rightarrow$ we break them down into B steps :

$$(s, a, r, s')$$

Second step : Update the networks

\* at every iteration, update the prediction model :

Compute $\widehat{\mathcal{L}}(\Theta)$

$$\mathcal{L}(\Theta) = \frac{1}{2}\left( \underset{(s,a,r,s') \sim 6_\Theta}{\mathbb{E}} \left[ r + \gamma \underset{a' \in A}{\max} Q_{\Theta'}(s', a') - Q_\Theta(s,a) \right] \right)^2$$

$$\underbrace{\quad}_{fixed}$$

$$\widehat{\mathcal{L}}(\Theta) = \frac{1}{2} \text{ average over steps from the batch :}$$

$$(s,a,r,s') : \quad r + \gamma \underset{a' \in A}{\max} Q_{\Theta'}(s', a') - Q_\Theta(s,a)$$

if you have a neural network / ML model parametrised by $\Theta$, you have functions for

computing $\nabla_\Theta$ and for updating :

optimisation step

$$\widehat{\mathcal{L}}(\Theta) \xrightarrow{\text{gradient}} \nabla_\Theta \widehat{\mathcal{L}}(\Theta) \longrightarrow \Theta_{new}$$

* every N iterations, update the target model :

$$\theta' \leftarrow \theta \longleftarrow \text{parameters of the prediction model}$$

⌐ parameters of the target model

$$N \simeq 10$$