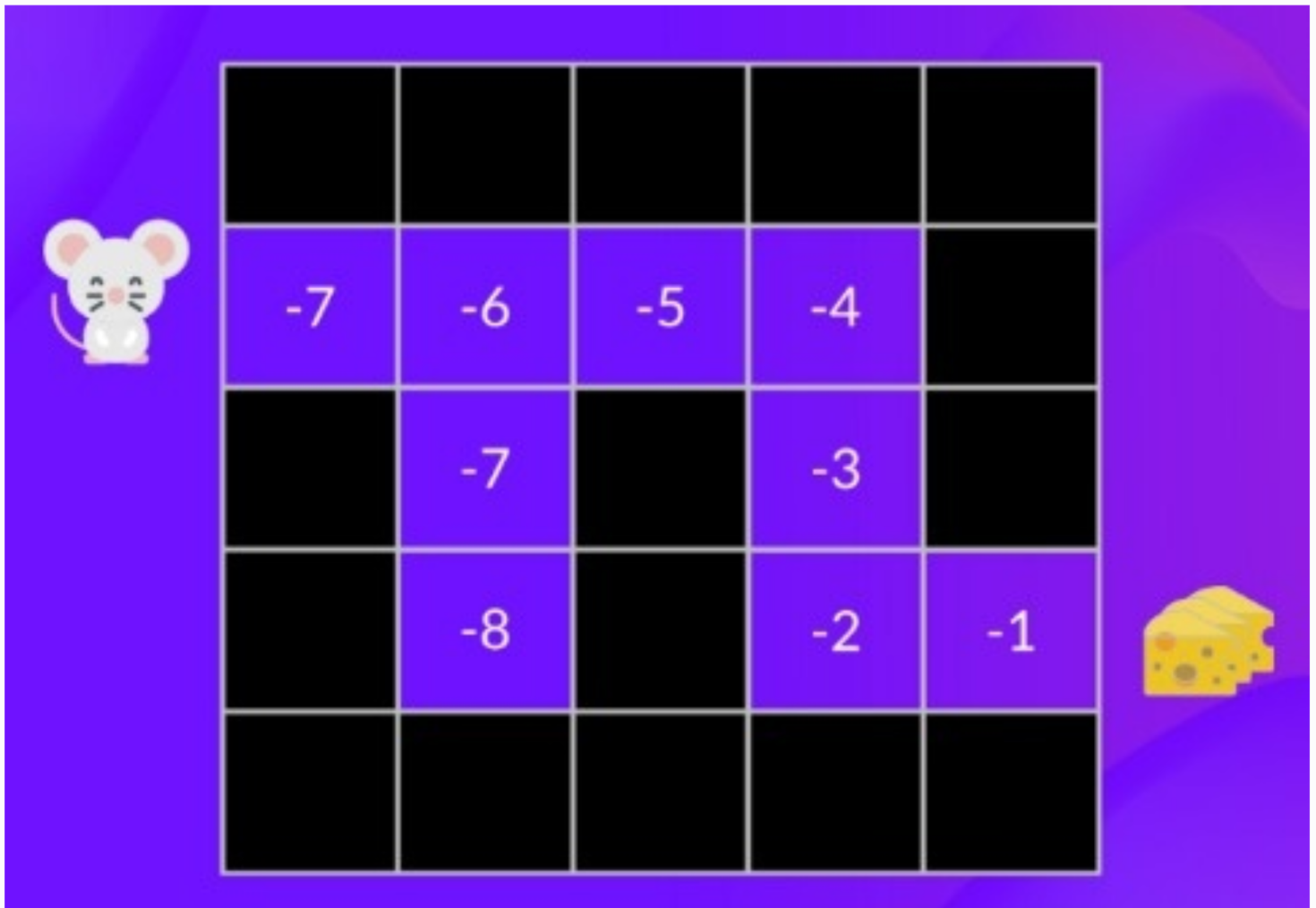


VALUE-BASED APPROACHES

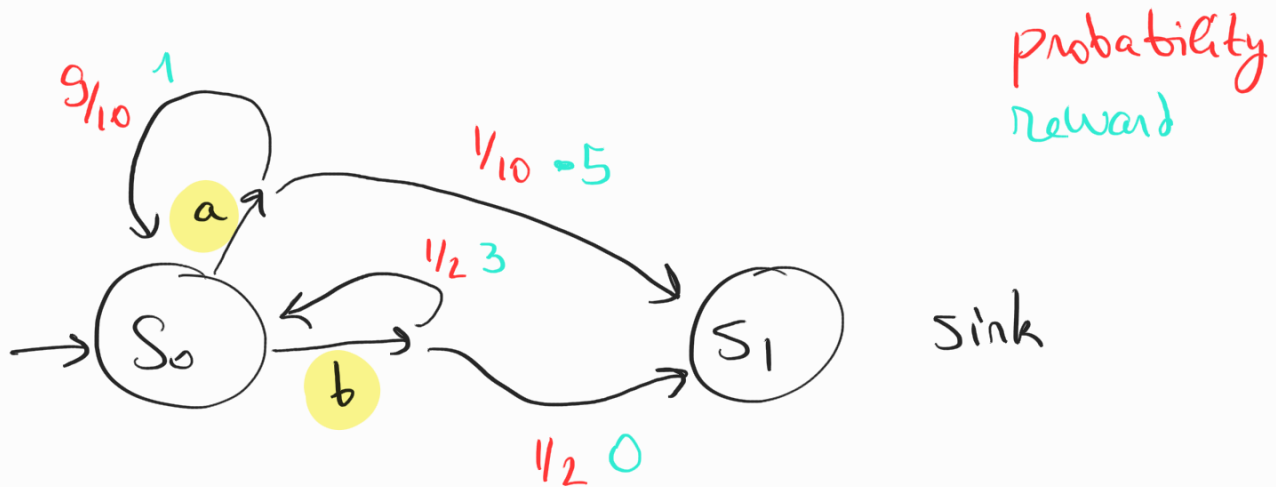
STATE VALUE FUNCTION

$$V_{\pi}: S \rightarrow \mathbb{R}$$

$$V_{\pi}(s) = \mathbb{E} \left[G \mid S(0) = s \wedge \pi \right]$$



BELLMAN EQUATIONS



$$\pi(s_0) = a$$

$$\begin{cases} V_{\pi}(s_0) = \frac{9}{10} (1 + \gamma V_{\pi}(s_0)) + \frac{1}{10} (-5 + \gamma V_{\pi}(s_1)) \\ V_{\pi}(s_1) = 0 \end{cases}$$

More generally:

$$V_{\pi}(s) = \mathbb{E}_{s', R} [R + \gamma V_{\pi}(s') \mid s, a]$$

optimal state value function:

$$V_* : S \rightarrow \mathbb{R}$$

$$V_*(s) = \sup_{\pi} V_{\pi}(s)$$

greedy policy

$$\pi(s) = \operatorname{argmax}_{a \in A} \mathbb{E} [R + \gamma V(s') \mid s, a]$$

STATE-ACTION VALUE FUNCTION

$$Q_{\pi} : S \times A \rightarrow \mathbb{R}$$

$$Q_{\pi}(s, a) = \mathbb{E} \left[G \mid S(0) = s \wedge A(0) = a \wedge \pi \right]$$

optimal state-action value function:

$$Q_{*} : S \times A \rightarrow \mathbb{R}$$

$$Q_{*}(s, a) = \sup_{\pi} Q_{\pi}(s, a)$$

greedy policy

$$\pi(s) = \underset{a \in A}{\operatorname{argmax}} \quad Q(s, a)$$

→ FIRST APPROACH : MONTE CARLO

we play the ϵ -greedy policy:

$$\pi(s) = \begin{cases} \operatorname{argmax}_{a \in A} Q(s, a) \\ \text{Uniform} \end{cases}$$

with probability $1-\epsilon$

with probability ϵ

REPEAT

sample a trajectory with ϵ -greedy policy

compute return G_t

$$V(s) \leftarrow V(s) + \alpha (G - V(s))$$

$$\text{NEW} = \text{OLD} + \alpha [\text{CURRENT} - \text{OLD}]$$

→ SECOND APPROACH: TEMPORAL DIFFERENCE

Key idea: bootstrapping

REPEAT

sample a step (s, a, r, s') with ϵ -greedy policy

$$V(s) \leftarrow V(s) + \alpha [r + \gamma V(s') - V(s)]$$

$$\text{NEW} = \text{OLD} + \alpha [\text{CURRENT} - \text{OLD}]$$

Key difference:

- Monte Carlo requires a full trajectory before making an update
BUT smaller variance
- Temporal Difference updates at each step

Q-LEARNING

- Value-based : $Q : S \times A \rightarrow \mathbb{R}$
- Temporal Difference
- Off policy : uses one policy for *acting* and another one for *updating*

REPEAT

sample a step (s, a, r, s') with ϵ -greedy policy

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a) \right]$$

$$\text{NEW} = \text{OLD} + \alpha \left[\text{CURRENT} - \text{OLD} \right]$$

SARSA

- Value-based : $Q : S \times A \rightarrow \mathbb{R}$
- Temporal Difference
- on-policy : uses one policy for *acting* and for *updating*

REPEAT

Sample a step (S, a, r, S') with ϵ -greedy policy

Sample a step (S', a', r', S'') with ϵ -greedy policy

$$Q(S, a) \leftarrow Q(S, a) + \alpha [r + \gamma Q(S', a') - Q(S, a)]$$

$$\text{NEW} = \text{OLD} + \alpha [\text{CURRENT} - \text{OLD}]$$

EXPERIENCE REPLAY

running new samples each time is :

- ↳ wasteful
- ↳ biased by the current strategy
- ↳ not giving the maximum amount of information

EXPERIENCE REPLAY

the buffer has a fixed size

A buffer stores steps (s, a, r, s')

Two actions:

- (1) sample a trajectory and add each step to the buffer (independently!)
- (2) sample from the buffer to update

PRIORITISED EXPERIENCE REPLAY

- when we add (s, a, r, s') we compute the bias

$$B = r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a)$$

to obtain a distribution we apply a softmax:

$$\exp(B) / \sum_{B'} \exp(B')$$

- we sample from the buffer with this distribution

DOUBLE Q-LEARNING

Issue:

$$Q(s, a) = Q(s, a) + \alpha \left(r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a) \right)$$

→ tends to overshoot

THE MAXIMISATION BIAS

X_1, \dots, X_n random variables

Goal: estimate $\max_i \mathbb{E}[X_i]$

single estimator:

for each i : Sample X_i a number of times

→ choose $\operatorname{argmax}_i \hat{\mathbb{E}}[X_i] = *$

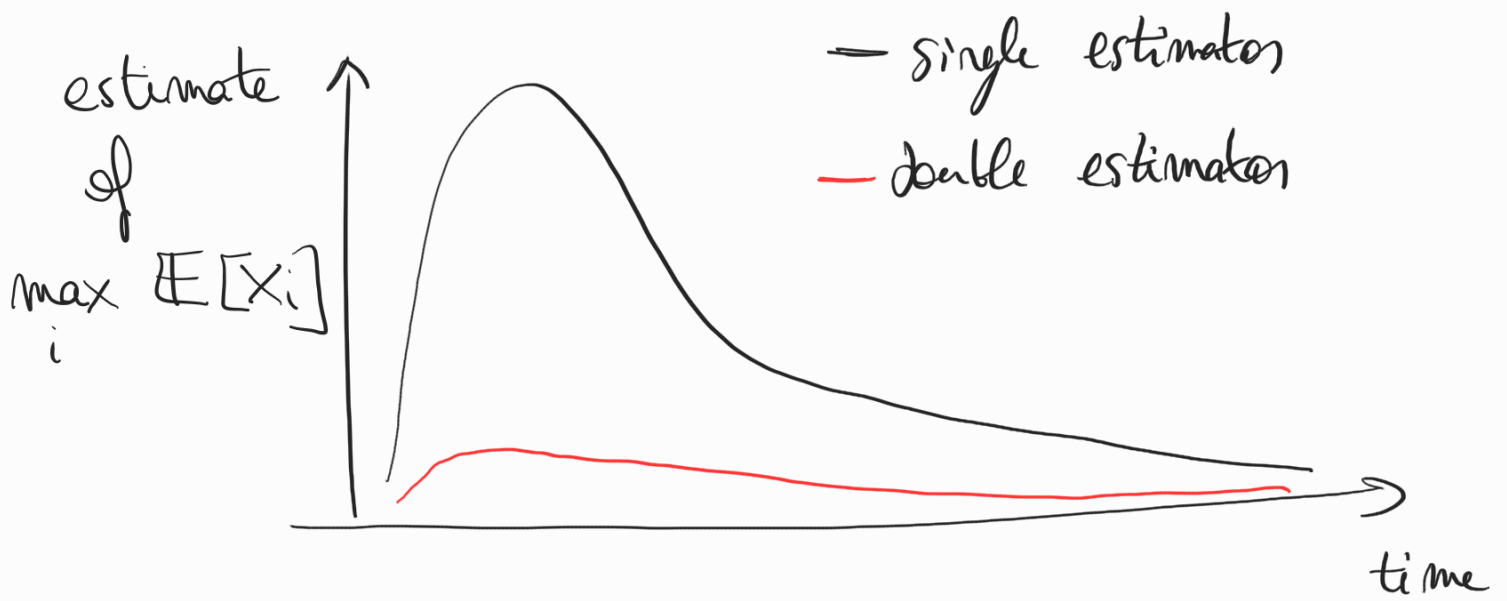
→ evaluate by $\hat{\mathbb{E}}[X_*]$

double estimator:

for each i : Sample X_i a number of times, twice

→ choose $\operatorname{argmax}_i \hat{\mathbb{E}}[X_i] = *$ with first set

→ evaluate by $\hat{\mathbb{E}}[X_*]$ with second set



DOUBLE Q-LEARNING ALGORITHM

with probability $1/2$:

- $a' = \operatorname{argmax}_{a'} Q_2(s', a')$
- $Q_1(s, a) \leftarrow Q_1(s, a) + \alpha (r + \gamma Q_1(s', a') - Q_1(s, a))$

else:

symmetrically