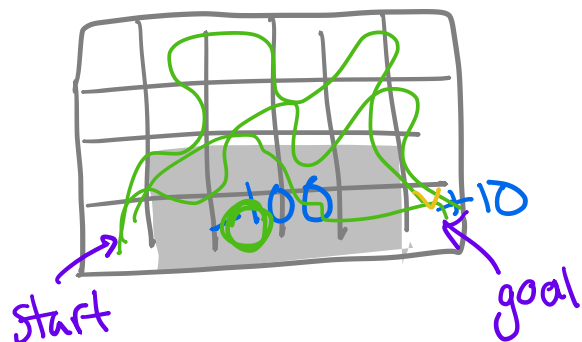


Goal find $\pi^*(a|s)$

How? learn $Q^*(s,a) \rightarrow \pi(a|s) = \operatorname{argmax}_a Q(s,a)$

① Q-Learning

② Deep Q-Networks



deterministic $P(s'|s,a)$

ϵ -greedy policy

$$a = \begin{cases} \operatorname{argmax}_a Q(s,a) & 1-\epsilon \\ \text{random} & \epsilon \end{cases}$$

Q-Learning Algorithm

initialize Q-values $Q(s,a) = 0 \quad \forall (s,a)$

for num episodes

initialize state s

while not done

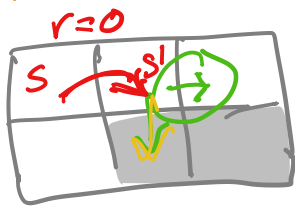
select action $a = \begin{cases} \operatorname{argmax}_a Q(s,a) & 1-\epsilon \\ \text{random} & \epsilon \end{cases}$

take action a and observe reward r and next state s'

update Q-value

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[\underbrace{r + \gamma \max_{a'} Q(s',a')}_{\text{target}} - \underbrace{Q(s,a)}_{\text{curr. est}} \right]$$

$$Q^*(s,a) = r + \gamma \max_{a'} Q^*(s',a')$$



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

error (TD-error, RPE)

↳ drives learning

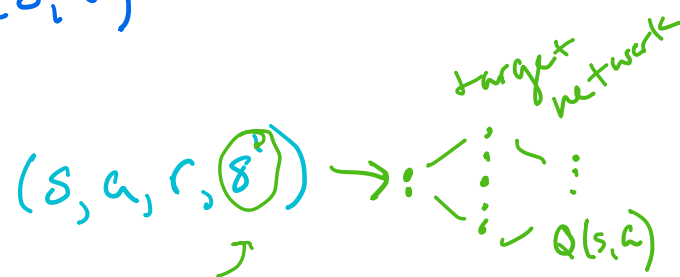
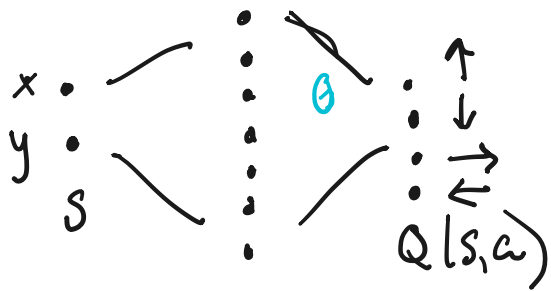
$$\operatorname{argmax}_{a'} P(s'|s',a')$$

Problem w/ Q-Learning : TABULAR

Deep RL : use neural networks for function approximation

Deep Q-Networks :

- use NN to approximate $Q(s,a)$



$$\mathcal{L}(\theta) = \left\langle \left(\underbrace{r + \gamma \max_{a'} Q(s', a'; \theta^-)}_{\text{target}} - \underbrace{Q(s, a; \theta)}_{\text{curr. est.}} \right)^2 \right\rangle_{(s,a,r,s') \sim \text{Unif. (Mem)}}$$

① correlated inputs \rightarrow experience replay

② high variance gradient \rightarrow target network

periodically set $\theta^- = \theta$

$$\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} Q(s, a; \theta)$$

DQN Algorithm

initialize Q-network with random weights θ

initialize target network with weights $\theta^- = \theta$

for num episodes

initializing state s

while not done:

select action $a = \begin{cases} \arg\max_a Q(s, a; \theta) & 1-\epsilon \\ \text{random} & \epsilon \end{cases}$

take action a and observe reward r and next state s'

store experience (s, a, r, s') in memory
 sample random batch of (s, a, r, s') from memory
 calculate target

$$y = \begin{cases} r + \gamma \max_{a'} Q(s', a'; \theta^-) & \text{if } s \text{ not terminal} \\ r & \text{if } s \text{ terminal} \end{cases}$$

 perform gradient descent step on $(y - Q(s, a; \theta))^2$ wrt. θ
 every C steps update target network $\theta^- = \theta$

Main (Deep) RL Algorithm Classes

① Value-based (e.g. DQN)

$s \rightarrow \text{network} \rightarrow Q(s, a)$

② Policy-based (e.g. Policy Gradient) $J(\theta) = \left\langle \sum_{t=0}^T \gamma^t r_t \right\rangle_{\pi_\theta}$

$s \rightarrow \text{network} \rightarrow \pi(a|s)$

$$\nabla_{\theta} J(\theta) = \left\langle \sum_{t=0}^T \nabla_{\theta} \pi(a_t | s_t) \hat{G}_t \right\rangle_{\pi_\theta}$$

$\hat{G}_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$

③ Actor-Critic

$s \rightarrow \text{Actor} \rightarrow \pi(a|s)$

$s \rightarrow \text{Critic} \rightarrow Q(s, a)$

④ Combinations!