

## Q-learning with function approximation

Given  $w^- \rightarrow$  update to new  $w$

$$L(w) = \frac{1}{2} \underbrace{\left( r_{k+1} + \gamma \max_a Q(s_{k+1}, a; w^-) - Q(s_k, a_k; w) \right)}_{\text{target}}^2$$

minimize  $L(w) :$   $w \leftarrow w^- - \alpha \nabla_w L(w^-)$

$$\begin{aligned} \nabla_w L(w) &= - \underbrace{\left( r_{k+1} + \gamma \max_a Q(s_{k+1}, a; w^-) - Q(s_k, a_k; w) \right)}_{\delta_k} \nabla_w Q(s_k, a_k; w) \\ &= - \delta_k \nabla_w Q(s_k, a_k; w) \end{aligned}$$

update rule:  $w \leftarrow w^- + \alpha \delta_k \nabla_w Q(s_k, a_k; w^-)$

$\nwarrow \delta_k(s_k, a_k, r_{k+1}, s_{k+1}; w^-)$

$\rightarrow$  doesn't work due to correlated updates with global dependence on  $w$

Mnih 2015: experience replay  $\rightarrow$  update for a collection of historical states.

$$L(w) = E_{(s,a,r,s') \sim \mathcal{D}} \left[ \frac{1}{2} \left( \underbrace{r + \gamma \max_a Q(s', a; w^-)}_{\text{target}} - Q(s, a; w) \right)^2 \right]$$

history  $\nearrow$

$$= \sum_{(s,a,r,s') \in \mathcal{D}} \frac{1}{N} \left[ ( \text{---} )^2 \right]$$

$$\nabla_w L(w) = - E_{(s,a,r,s') \sim \mathcal{D}} \left[ \delta(s, a, r, s'; w^-, w) \nabla_w Q(s, a; w) \right]$$

$$w \leftarrow w^- - \alpha \nabla_w L(w^-)$$

samples in history  $\mathcal{D}$  were not generated with current  $Q$ , but with old  $Q$  policies  $\rightarrow$  need off-policy Q-learning