

Last Time

Policy Gradient

$\pi_{\theta}$

$$\theta \leftarrow \theta + \alpha \nabla U(\theta)$$

↑ estimate with 1 sim

Likelihood Ratio

Causality

Baseline subtraction

Model-Based

Learn TR  
solve



Model Free

Learn  $\pi$

Learn  $Q$

Today

Value Learning  
TD Learning

Incremental Mean Estimation

$$\hat{x}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$
$$= \frac{1}{m} \left( x^{(m)} + \sum_{i=1}^{m-1} x^{(i)} \right)$$

$$= \frac{1}{m} \left( x^{(m)} + (m-1) \hat{x}_{m-1} \right)$$

$$\rightarrow \hat{x}_m = \hat{x}_{m-1} + \frac{1}{m} (x^{(m)} - \hat{x}_{m-1})$$

$$Q(s,a) = Q(s,a) + \frac{1}{N(s,a)} (q - Q(s,a))$$

$$\hat{x} \leftarrow \hat{x} + \alpha (x - \hat{x})$$

↑ learning rate

$q = r + \gamma$  simulate!

## Q-learning (Conceptual)

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} T(s'|s,a) V(s')$$

$$= R(s,a) + \gamma \sum_{s'} T(s'|s,a) \max_{a'} Q(s',a')$$

$$Q(s,a) = E_{r,s'} [r + \gamma \max_{a'} Q(s',a')]$$

## Q-learning $(s,a,r,s')$

$$Q(s,a) \leftarrow 0 \quad \forall s,a$$

$$s \leftarrow s_0$$

$$\text{loop} \quad a \leftarrow \text{exploration\_policy}(Q)$$

$$s', r \leftarrow \text{step!}(\text{env}, a)$$

$$Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(s,a))$$

TD

e.g.  $\epsilon$ -greedy

book  
version

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

$$Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma Q(s',a') - Q(s,a))$$

Off-Policy

Q-learning

Fewer Samples

$\longleftrightarrow$

On-Policy

SARSA

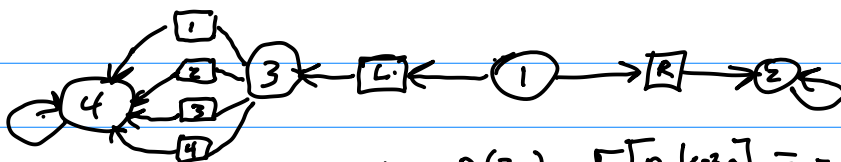
More stable

Policy Gradient

Learn only from  
actions that are  
in the policy they  
are executing

## Breakout Rooms

$$A(3) = \{1, 2, 3, 4\} \quad A(1) = \{L, R\}$$



↑ same reward  $R(3,a) = E[r | s=3, a] = -0.1$

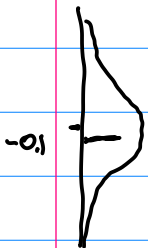
$$r \sim N(-0.1, 1)$$

After a few episodes

$$Q(3,1)$$

$$Q(3,2)$$

Will Q-learning work well?



Right side stays at 0

Max throws off estimate on Left

## Maximization Bias

$$\left. \begin{array}{l} Q_2(3,1) = -0.1 \quad Q_1(3,1) = 0.2 \\ \quad \quad \quad Q_1(3,2) = 0.1 \\ Q_2(3,3) = 0.1 \quad Q_1(3,3) = -0.3 \\ \quad \quad \quad Q_1(3,4) = -0.5 \end{array} \right\} \text{max} \rightarrow Q(1,L)$$

## Solution

Double Q-Learning

$Q_1 \quad Q_2$

$$Q_1(s,a) \leftarrow Q_1(s,a) + \alpha (R(s,a) + \gamma \underbrace{Q_2(s', \arg \max_{a'} Q_1(s', a'))}_{\text{max}} - Q_1(s,a))$$

S+B section 6.7

Credit Assignment

Problem: takes a long time to propagate backwards

Solution: Eligibility Traces (look at S+B)

Sarsa( $\lambda$ )

$Q(\lambda)$

TD( $\lambda$ )

$\lambda$ : decay param  
 $\in [0,1]$

$s, a, r, s', a'$

$$Q(s,a), N(s,a) \leftarrow 0 \quad \forall s,a$$

loop

$$a_{t+1} \leftarrow \text{explore}(Q)$$

$$N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$$

$$\delta \leftarrow r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

$$Q(s,a) \leftarrow Q(s,a) + \alpha \delta N(s,a) \quad \forall s,a$$

$$N(s,a) \leftarrow \gamma \lambda N(s,a) \quad \forall s,a$$

## Convergence

Q-learning converges to optimal Q-values w.p.1 S+B p.131

Sarsa converges to optimal Q-values w.p.1  
provided that  $\pi \rightarrow \text{greedy}(Q)$  S+B p.129

avg  
return  
for an  
episode



wall clock time

