

# Last Time

Trying Actions  
Adjusting  $\theta$  so that better  
actions are more likely

- What is Policy Gradient?
- What tricks are needed to make Policy Gradient work?
  - Log Derivative
  - Causality
  - Baseline Subtraction

# Map

# Map

Model  
Based

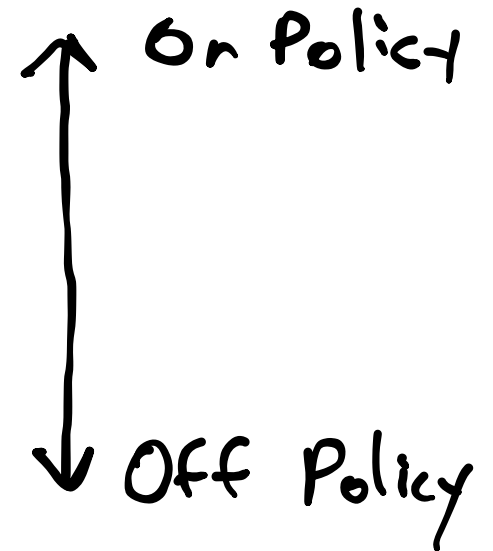
Model  
Free



# Map

Model  
Based

Model  
Free



# Map

Model  
Based

Model  
Free



MLMBTRL  
(learn  $T, R$ )

On Policy  
Off Policy

A vertical double-headed arrow on the right side of the diagram, indicating a spectrum or continuum between 'On Policy' and 'Off Policy' learning methods.

# Map

Model  
Based

Model  
Free

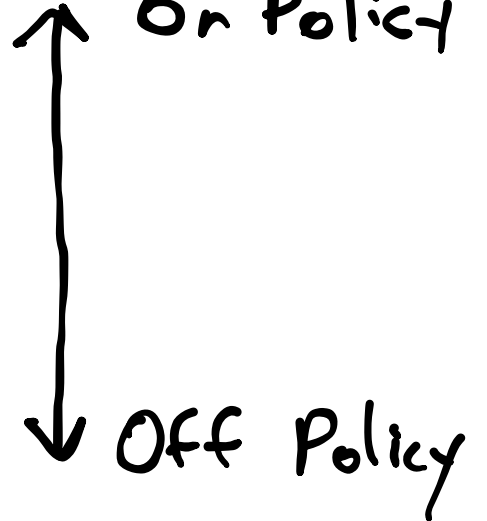


learn Q

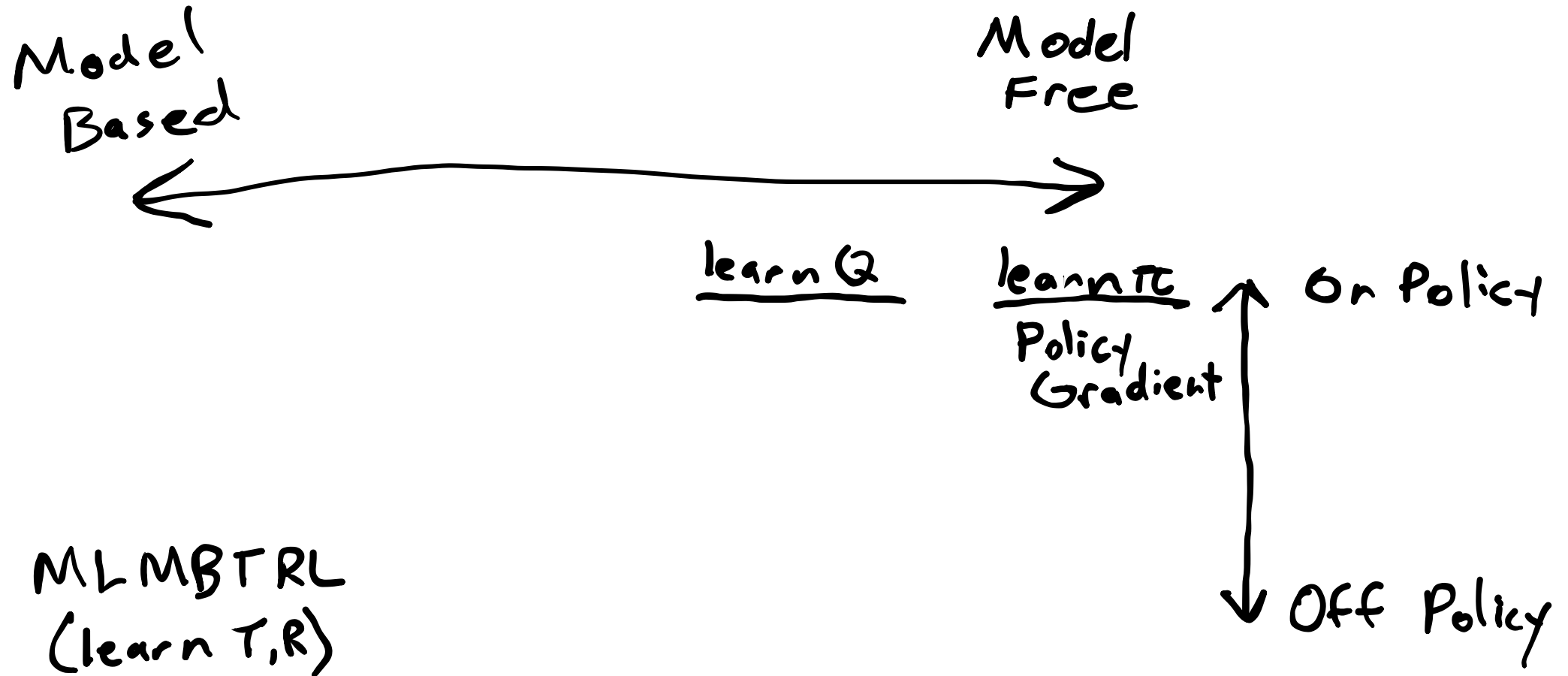
learn  $\pi$

On Policy

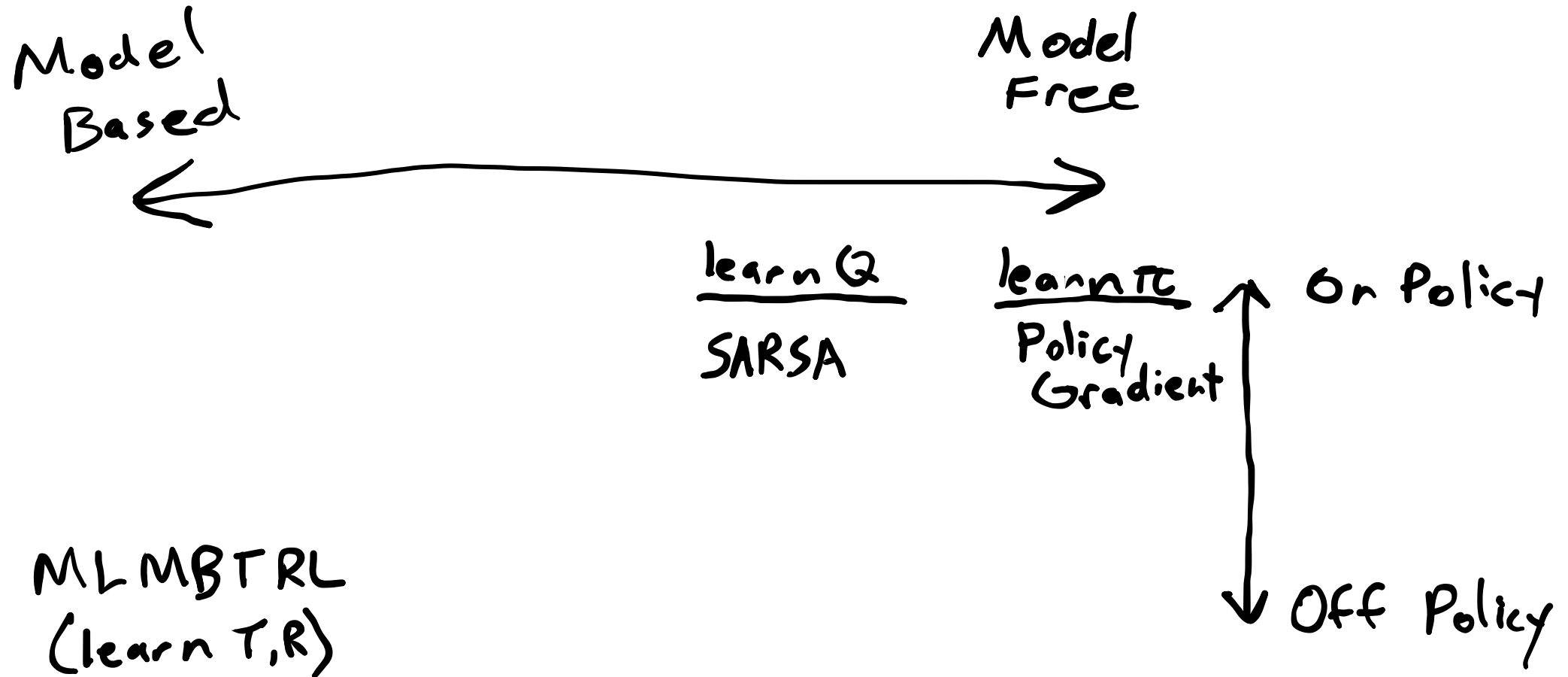
MLMBTRL  
(learn T, R)



# Map

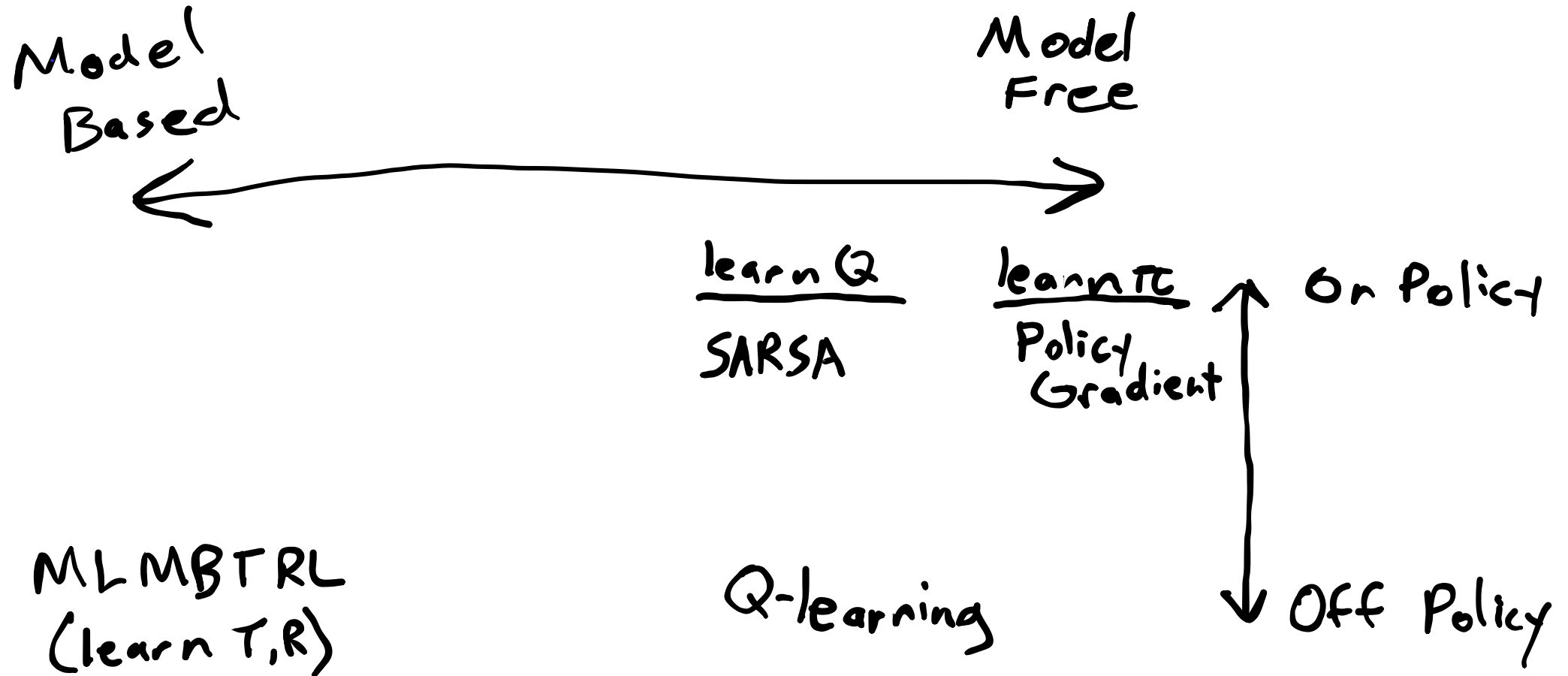


# Map





# Map



# Today

# Today

- Basic On- and Off-Policy **value based** model free RL algorithms

# Today

- Basic On- and Off-Policy **value based** model free RL algorithms
- Tricks for tabular value based RL algorithms

# Today

- Basic On- and Off-Policy **value based** model free RL algorithms
- Tricks for tabular value based RL algorithms
- What makes it hard for some approaches to use Off-Policy data?

# Review: Why learn Q?

$$V(s)$$

$$R(s,a)$$

$$Q(s,a)$$

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s,a)$$



# Incremental Mean Estimation



# Incremental Mean Estimation

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





# Incremental Mean Estimation

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



# Incremental Mean Estimation

$$\begin{aligned}\hat{x}_m &= \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ &= \frac{1}{m} \left( x^{(m)} + \underbrace{\sum_{i=1}^{m-1} x^{(i)}}_{(m-1)\hat{x}_{m-1}} \right) \\ &= \frac{1}{m} \left( x^{(m)} + (m-1)\hat{x}_{m-1} \right)\end{aligned}$$



# Incremental Mean Estimation

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



# Incremental Mean Estimation

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

```
function simulate!( $\pi$ ::MonteCarloTreeSearch, s, d= $\pi$ .d)
    if d ≤ 0
        return  $\pi$ .U(s)
    end
     $\mathcal{P}$ , N, Q, c =  $\pi$ . $\mathcal{P}$ ,  $\pi$ .N,  $\pi$ .Q,  $\pi$ .c
     $\mathcal{A}$ , TR,  $\gamma$  =  $\mathcal{P}$ . $\mathcal{A}$ ,  $\mathcal{P}$ .TR,  $\mathcal{P}$ . $\gamma$ 
    if !haskey(N, (s, first( $\mathcal{A}$ )))
        for a in  $\mathcal{A}$ 
            N[(s,a)] = 0
            Q[(s,a)] = 0.0
        end
        return  $\pi$ .U(s)
    end
    a = explore( $\pi$ , s)
    s', r = TR(s,a)
    q = r +  $\gamma$ *simulate!( $\pi$ , s', d-1)
    N[(s,a)] += 1
    Q[(s,a)] += (q-Q[(s,a)])/N[(s,a)]
    return q
end
```



# Incremental Mean Estimation

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

```
function simulate!( $\pi$ ::MonteCarloTreeSearch, s, d= $\pi$ .d)
    if d  $\leq$  0
        return  $\pi$ .U(s)
    end
     $\mathcal{P}$ , N, Q, c =  $\pi$ . $\mathcal{P}$ ,  $\pi$ .N,  $\pi$ .Q,  $\pi$ .c
     $\mathcal{A}$ , TR,  $\gamma$  =  $\mathcal{P}$ . $\mathcal{A}$ ,  $\mathcal{P}$ .TR,  $\mathcal{P}$ . $\gamma$ 
    if !haskey(N, (s, first( $\mathcal{A}$ )))
        for a in  $\mathcal{A}$ 
            N[(s,a)] = 0
            Q[(s,a)] = 0.0
        end
        return  $\pi$ .U(s)
    end
    a = explore( $\pi$ , s)
    s', r = TR(s,a)
    q = r +  $\gamma$ *simulate!( $\pi$ , s', d-1)
    Q[(s,a)] += (q-Q[(s,a)))/N[(s,a)]
end
```



# Incremental Mean Estimation

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

```
function simulate!( $\pi$ ::MonteCarloTreeSearch, s, d= $\pi$ .d)
    if d  $\leq$  0
        return  $\pi$ .U(s)
    end
     $\mathcal{P}$ , N, Q, c =  $\pi$ . $\mathcal{P}$ ,  $\pi$ .N,  $\pi$ .Q,  $\pi$ .c
     $\mathcal{A}$ , TR,  $\gamma$  =  $\mathcal{P}$ . $\mathcal{A}$ ,  $\mathcal{P}$ .TR,  $\mathcal{P}$ . $\gamma$ 
    if !haskey(N, (s, first( $\mathcal{A}$ )))
        for a in  $\mathcal{A}$ 
            N[(s,a)] = 0
            Q[(s,a)] = 0.0
        end
    end
    return  $\pi$ .U(s)
end
a = explore( $\pi$ , s)
s', r = TR(s,a)
q = r +  $\gamma$ *simulate!( $\pi$ , s', d-1)
Q[(s,a)] += (q-Q[(s,a)])/N[(s,a)]
end
```

loop

$$\hat{x} \leftarrow \hat{x} + \overset{\text{learning rate}}{\alpha} (\underbrace{x - \hat{x}})$$

# Incremental Mean Estimation

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

```
function simulate!( $\pi$ ::MonteCarloTreeSearch, s, d= $\pi$ .d)
    if d  $\leq$  0
        return  $\pi$ .U(s)
    end
     $\mathcal{P}$ , N, Q, c =  $\pi$ . $\mathcal{P}$ ,  $\pi$ .N,  $\pi$ .Q,  $\pi$ .c
     $\mathcal{A}$ , TR,  $\gamma$  =  $\mathcal{P}$ . $\mathcal{A}$ ,  $\mathcal{P}$ .TR,  $\mathcal{P}$ . $\gamma$ 
    if !haskey(N, (s, first( $\mathcal{A}$ )))
        for a in  $\mathcal{A}$ 
            N[(s,a)] = 0
            Q[(s,a)] = 0.0
        end
    end
    return  $\pi$ .U(s)
end
a = explore( $\pi$ , s)
s', r = TR(s,a)
q = r +  $\gamma$ *simulate!( $\pi$ , s', d-1)
Q[(s,a)] += (q-Q[(s,a)])/N[(s,a)]
end
```

loop

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

"Temporal Difference  
(TD) Error"

# Q Learning


$$\hat{q}(r, s')$$



# Q Learning

Want:  $Q(s, a) \leftarrow Q(s, a) + \alpha (\hat{q} - Q(s, a))$



$\hat{q}(r, s')$

# Q Learning

Want:  $Q(s, a) \leftarrow Q(s, a) + \alpha (\hat{q}(r, s') - Q(s, a))$



$\hat{q}(r, s')$

# Q Learning

Want:  $Q(s, a) \leftarrow Q(s, a) + \alpha (\hat{q}(r, s') - Q(s, a))$

$$Q(s, a) = R(s, a) + \gamma E[V(s')]$$



# Q Learning

Want:  $Q(s, a) \leftarrow Q(s, a) + \alpha (\hat{q}(r, s') - Q(s, a))$

$$\begin{aligned} Q(s, a) &= R(s, a) + \gamma E[V(s')] \\ &= R(s, a) + \gamma E \left[ \max_{a'} Q(s', a') \right] \end{aligned}$$

  
 $\hat{q}(r, s')$

# Q Learning

Want:  $Q(s, a) \leftarrow Q(s, a) + \alpha (\hat{q}(r, s') - Q(s, a))$

$$\begin{aligned} Q(s, a) &= R(s, a) + \gamma E[V(s')] \\ &= R(s, a) + \gamma E \left[ \max_{a'} Q(s', a') \right] \\ &= E \left[ r + \gamma \max_{a'} Q(s', a') \right] \end{aligned}$$

$\hat{q}(r, s')$

# Q learning

TD

# Q learning

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

$$s \leftarrow s_0$$

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\operatorname{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\operatorname{env})$$

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

$$s \leftarrow s'$$

*$\epsilon$ -greedy exploration*

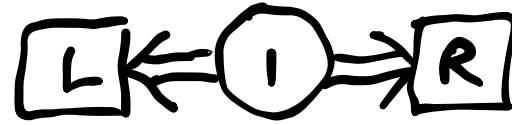
*aggressive*

**TD**

# Illustrative Problem



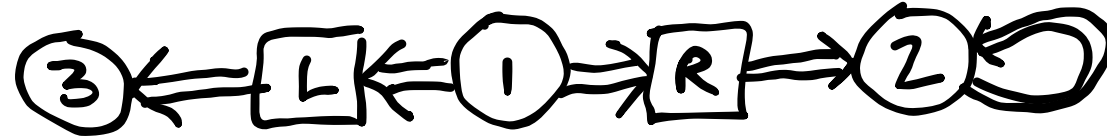
# Illustrative Problem



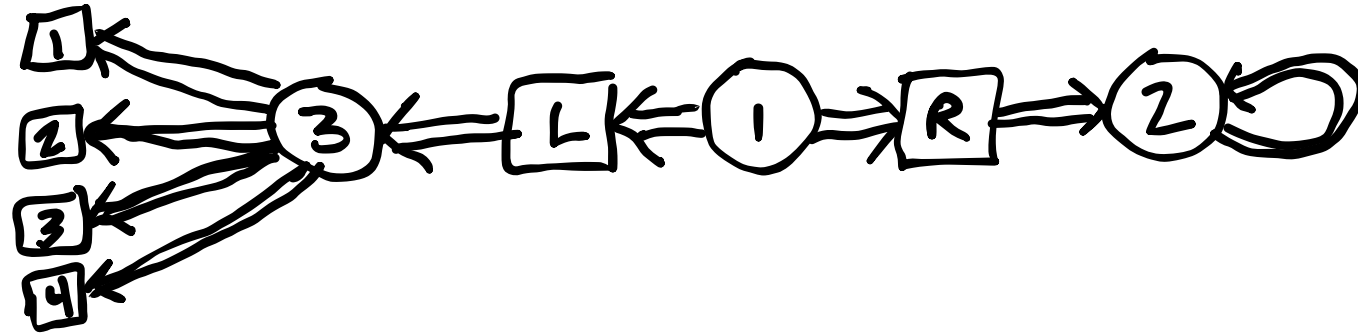
# Illustrative Problem



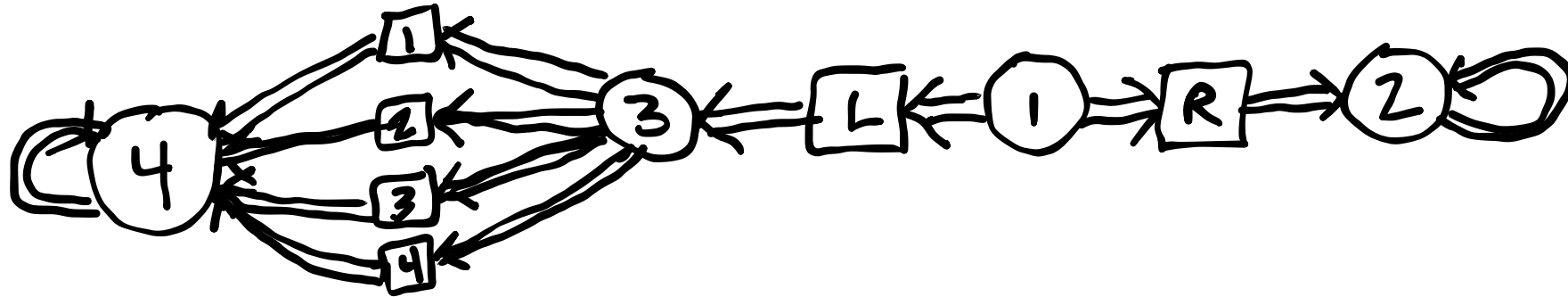
# Illustrative Problem



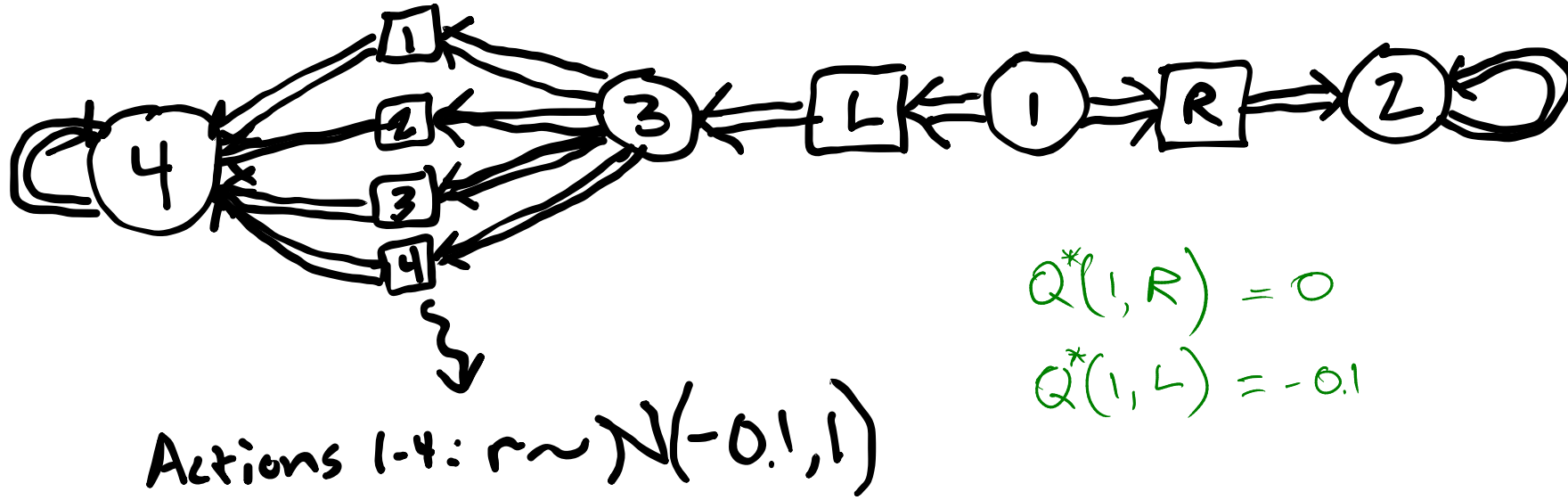
# Illustrative Problem



# Illustrative Problem

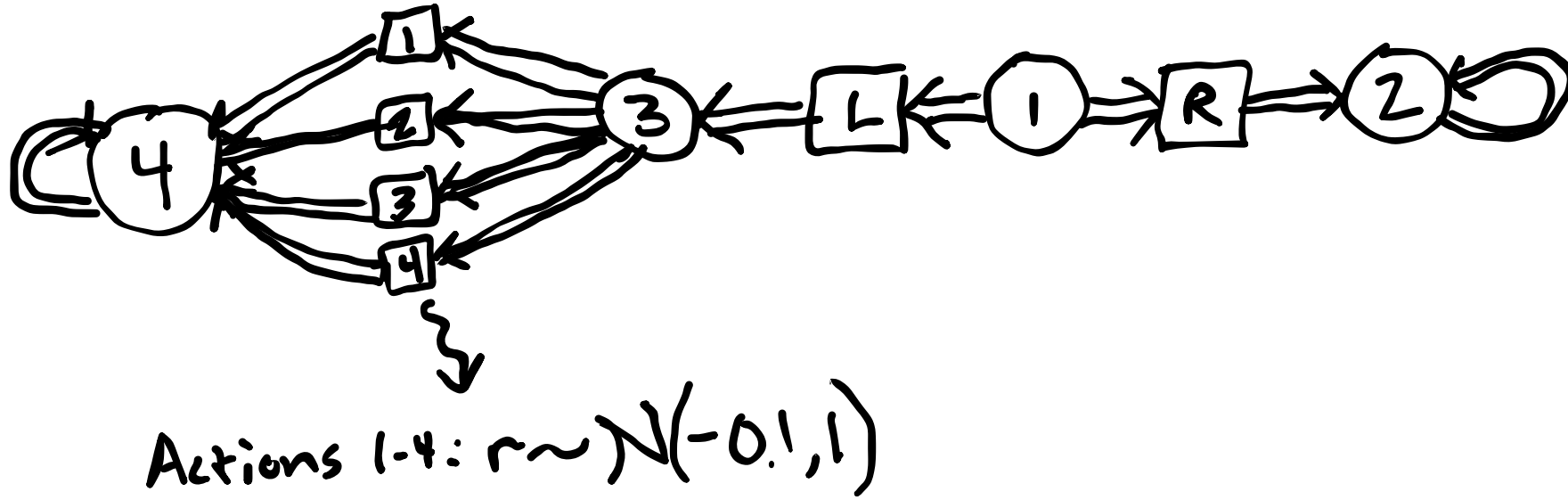


# Illustrative Problem



$$Q^*(1, R) = 0$$
$$Q^*(1, L) = -0.1$$

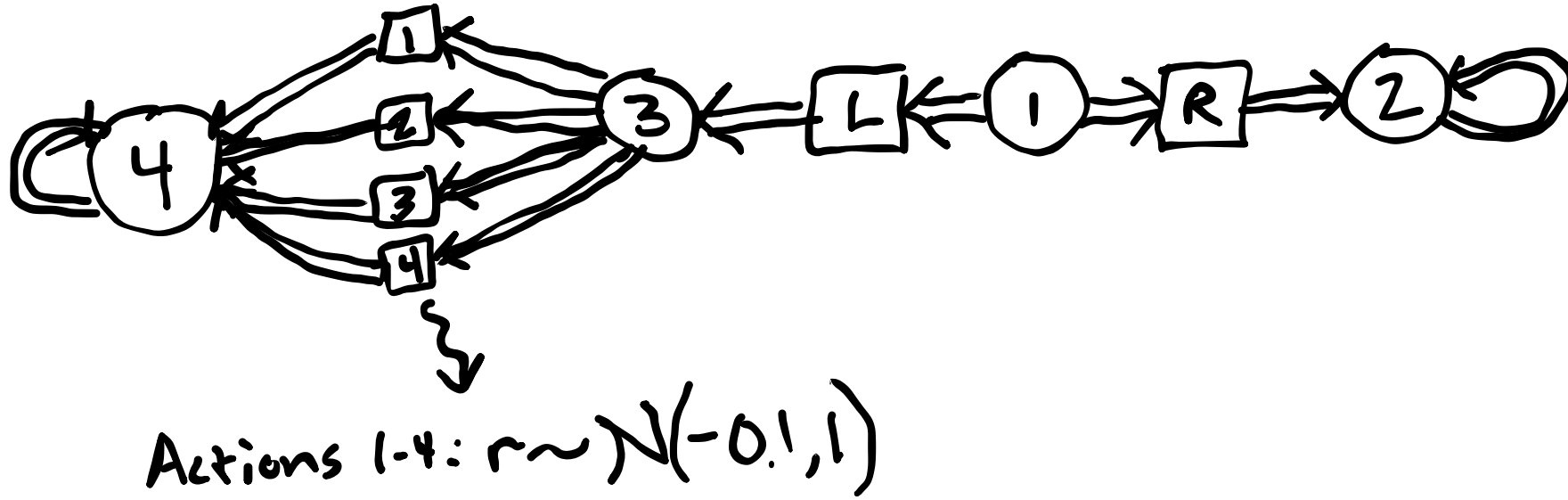
# Illustrative Problem



1. After a few episodes, what is  $Q(3, a)$  for  $a$  in 1-4?

$Q(3, 1)$

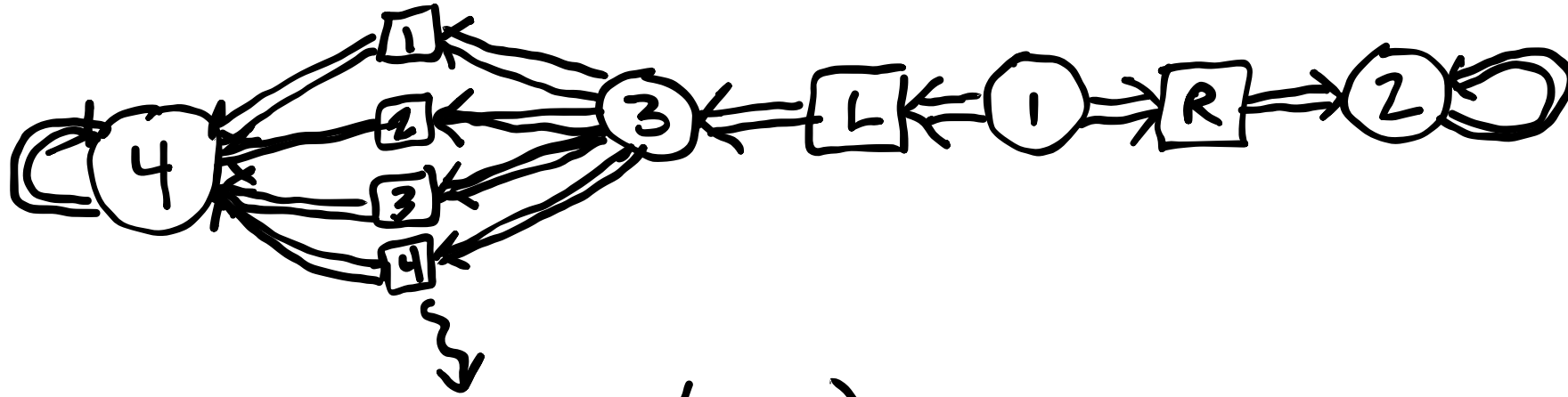
# Illustrative Problem



1. After a few episodes, what is  $Q(3, a)$  for  $a$  in 1-4?
2. After a few episodes, what is  $Q(1, L)$ ?



# Illustrative Problem



Actions 1-4:  $r \sim \mathcal{N}(\underline{-0.1}, 1)$

1. After a few episodes, what is  $Q(3, a)$  for  $a$  in 1-4?
2. After a few episodes, what is  $Q(1, L)$ ?
3. Why is this a problem and what are some possible solutions?

$$\left. \begin{aligned} Q(3, 4) &= -0.5 \\ Q(3, 1) &= 0.1 \\ Q(3, 2) &= -0.3 \end{aligned} \right\}$$

$$Q(1, L) = 0.1$$

# Big Problem: Maximization Bias

# Big Problem: Maximization Bias

Even if all  $Q(s', a')$  unbiased,  $\max_{a'} Q(s', a')$  is biased!

# Big Problem: Maximization Bias

Even if all  $Q(s', a')$  unbiased,  $\max_{a'} Q(s', a')$  is biased!

Solution: Double Q Learning

# Big Problem: Maximization Bias

Even if all  $Q(s', a')$  unbiased,  $\max_{a'} Q(s', a')$  is biased!

Solution: Double Q Learning      $Q_1, Q_2$

# Big Problem: Maximization Bias

Even if all  $Q(s', a')$  unbiased,  $\max_{a'} Q(s', a')$  is biased!

Solution: Double Q Learning  $Q_1, Q_2$

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

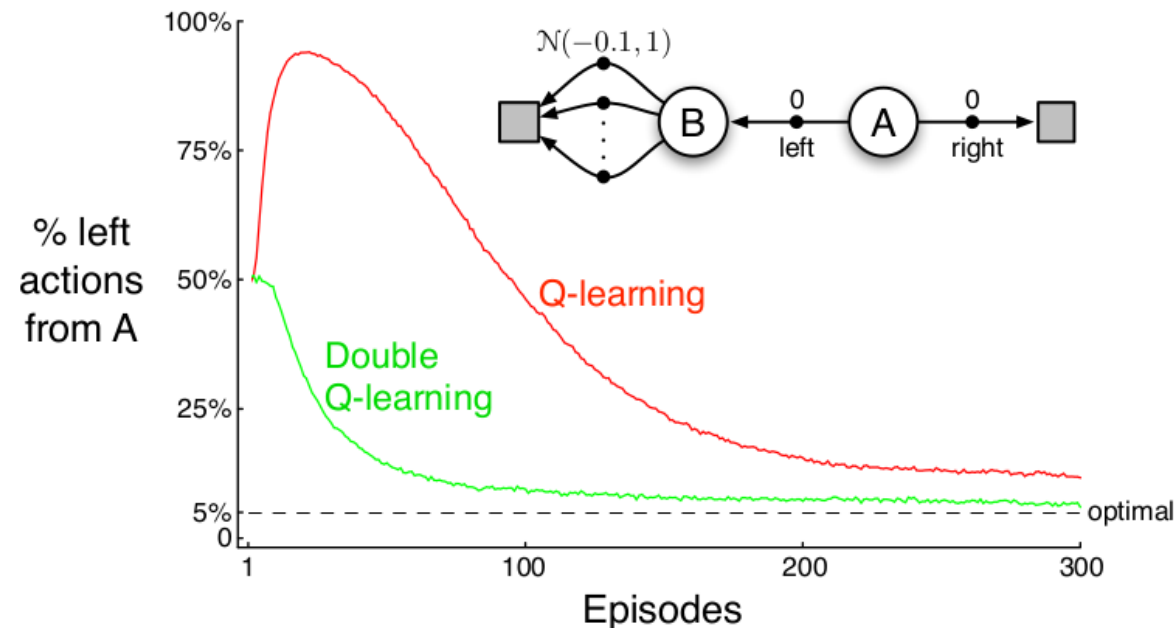
$r + \gamma \max_{a'} Q(s', a')$

# Big Problem: Maximization Bias

Even if all  $Q(s', a')$  unbiased,  $\max_{a'} Q(s', a')$  is biased!

Solution: Double Q Learning  $Q_1, Q_2$

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



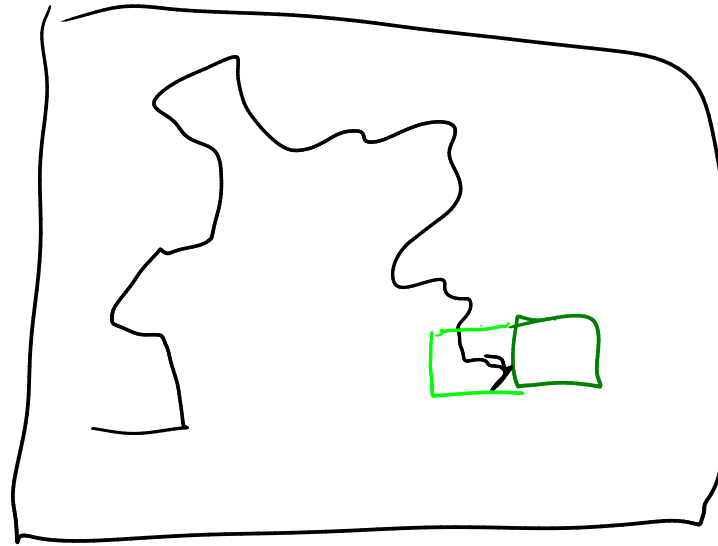
# SARSA

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

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



# Eligibility Traces



# SARSA- $\lambda$

# SARSA- $\lambda$

Games

## Half-Life at 20: why it is the most important shooter ever made

From its opening scenes, Valve's pioneering sci-fi horror game reinvented storytelling and universe building - what made it such a terrifying success?



📺 'It taught a whole generation of big-budget game developers how to tell stories' ... the Half-Life box art.  
Illustration: Valve

# SARSA- $\lambda$

Games

## Half-Life at 20: why it is the most important shooter ever made

From its opening scenes, Valve's pioneering sci-fi horror game reinvented storytelling and universe building - what made it such a terrifying success?



It taught a whole generation of big-budget game developers how to tell stories' ... the Half-Life box art. Illustration: Valve

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

initialize  $s, a, r, s'$

loop

$$a' \leftarrow \operatorname{argmax} Q(s', a) \text{ w.p. } 1 - \epsilon, \quad \text{rand}(A) \text{ o.w.}$$

$$N(s, a) \leftarrow N(s, a) + 1$$

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

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

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

$$s \leftarrow s', \quad a \leftarrow a'$$

$$r \leftarrow \text{act}!(\text{env}, a)$$

$$s' \leftarrow \text{observe}(\text{env})$$

$$\lambda \in [0, 1)$$

# Convergence

# Convergence

- Q learning converges to optimal Q-values w.p. 1 (Sutton and Barto, p. 131)

# Convergence

- Q learning converges to optimal Q-values w.p. 1  
(Sutton and Barto, p. 131)
- SARSA converges to optimal Q-values w.p. 1 ***provided that***  
 $\pi \rightarrow \text{greedy}$   
(Sutton and Barto, p. 129)

# On vs Off-Policy



# On vs Off-Policy

On Policy

# On vs Off-Policy

On Policy

Off Policy

# On vs Off-Policy

On Policy

Off Policy

SARSA:

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

# On vs Off-Policy

## On Policy

SARSA:

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

## Off Policy

Q-learning:

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

# On vs Off-Policy

## On Policy

SARSA:

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

## Off Policy

Q-learning:

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

Will eligibility traces work with Q-learning?

# On vs Off-Policy

## On Policy

SARSA:

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

## Off Policy

Q-learning:

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

Will eligibility traces work with Q-learning?

Not easily

# On vs Off-Policy

## On Policy

SARSA:

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

## Off Policy

Q-learning:

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

Will eligibility traces work with Q-learning?

Not easily

Policy Gradient:

$$\theta \leftarrow \theta + \alpha \sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau)$$

# Today

Double Q

- Basic On- and Off-Policy **value based** model free RL algorithms
- Tricks for tabular value based RL algorithms
- What makes it hard for some approaches to use Off-Policy data?

SARSA- $\lambda$