

MDP and RL

CS221: Section 4

Policy Iteration vs Value Iteration

SARSA vs Q-Learning

Function Approximation and Q-Learning

Policy Iteration vs Value Iteration

SARSA vs Q-Learning

Function Approximation and Q-Learning

Markov Decision Process (MDP)

$s \in \text{States}$

Actions(s): actions from state s

$T(s, a, s')$: probability next state is s' , given action a in state s

Note that $\sum_{s'} T(s, a, s') = 1$ for all $s \in \text{States}, a \in \text{Actions}(s)$

$R(s, a, s')$: reward for transition (s, a, s')

IsEnd(s): whether state marks end of game

Alternatively, $T(s, a, s) = 1$ and $R(s, a, s) = 0$

$0 \leq \gamma \leq 1$: discount factor



Markov Decision Process (MDP)

Given MDP

π : policy

take action $\pi(s)$ in state s

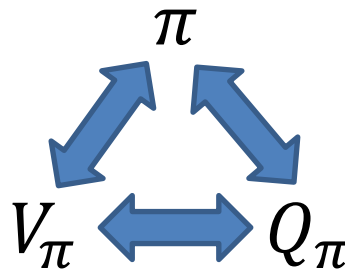
V_π : value function for policy π

$V_\pi(s)$ is expected utility of following π starting in s

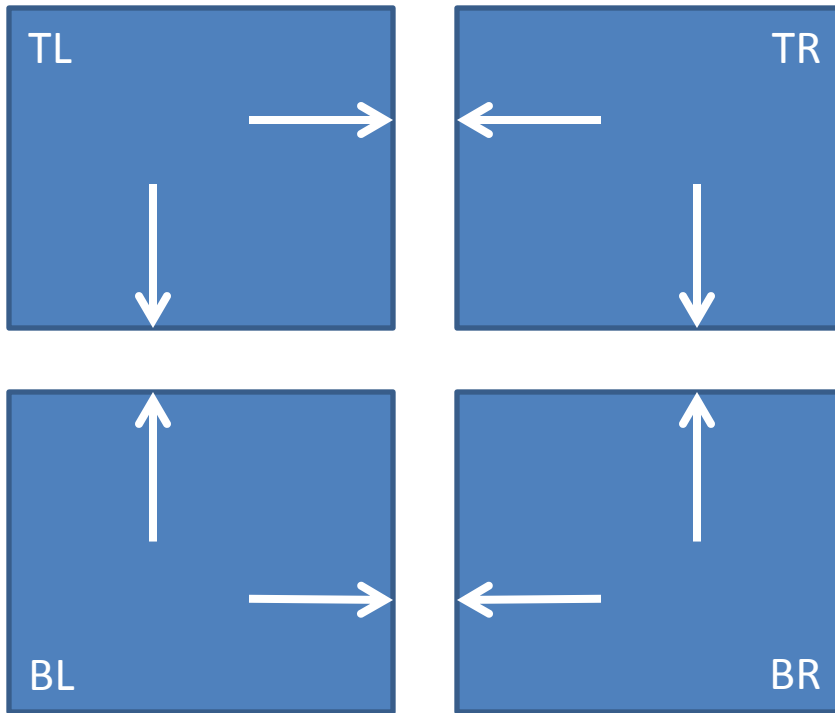
Q_π : Q – value of policy π

$Q_\pi(s, a)$ is expected utility of taking action a in s , then following π

Want to find optimal policy (that maximizes expected sum of discounted rewards)



Running Example



$States = \{TL, TR, NL, BR\}$

$EndState = BR$

$Actions(s) = \{H, V\}, \forall s$

$$T(s, a, s') = \begin{cases} s'(s, a), & \text{w.p. } 3/4 \\ s'(s, a^c), & \text{w.p. } 1/4 \end{cases}$$

except $T(BR, a, BR) = 1$

$Reward(s, a, s') = 4 \text{ if } s' = BR \text{ else } -1$

$\text{Simple policy } \pi(s) = H, \forall s$

Policy Iteration

Start with any π

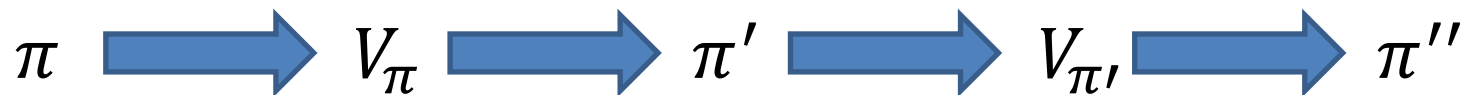
1. **Policy evaluation:** given π , compute V_π



2. **Policy improvement:** given V_π , compute improved policy π'



3. **Repeat** until policy converges/does not change



Policy Iteration

Start with any π

1. Policy evaluation: given π , compute V_π

a) Choose any V



Policy Iteration

Start with any π

1. Policy evaluation: given π , compute V_π

a) Choose any V

b) Update

$$V(s) := \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V(s')], \quad \forall s$$



Policy Iteration

Start with any π

1. Policy evaluation: given π , compute V_π

a) Choose any V

b) Update

$$V(s) := \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V(s')], \quad \forall s$$

or,

$$V := F_\pi V \quad (\text{compact})$$



Policy Iteration

Start with any π

1. Policy evaluation: given π , compute V_π

- a) Choose any V
- b) Update $V := F_\pi V$



Policy Iteration

Start with any π

1. Policy evaluation: given π , compute V_π

- a) Choose any V
- b) Update $V := F_\pi V$
- c) Repeat b) until convergence to obtain V_π

We know $V_\pi = F_\pi F_\pi \dots F_\pi V = F_\pi^\infty V$

Practically, update t_{PE} times to get reasonable approximation



Policy Evaluation Example

V_0

TL 0	TR 0
BL 0	BR 0

Simple policy $\pi(s) = H, \forall s$

Evaluate V_π

$V_0(s) = 0, \forall s$

Policy Evaluation Example

V_0		V_1	
TL	TR	TL	TR
0	0	-	-
BL	BR	BL	BR
0	0	-	-

Simple policy $\pi(s) = H, \forall s$

Evaluate V_π

$V_0(s) = 0, \forall s$

$$\begin{aligned}
 V_1(TL) &= \frac{3}{4} \times [-1 + V_0(TR)] + \frac{1}{4} \times [-1 + V_1(BL)] \\
 &= \frac{3}{4} \times [-1 + 0] + \frac{1}{4} \times [-1 + 0] \\
 &= -1
 \end{aligned}$$

Policy Evaluation Example

V_0		V_1	
TL	0	TL	-1
TR	0	TR	-
BL	0	BL	-
BR	0	BR	-

Simple policy $\pi(s) = H, \forall s$

Evaluate V_π

$V_0(s) = 0, \forall s$

$$\begin{aligned}
 V_1(TR) &= \frac{3}{4} \times [-1 + V_0(TL)] + \frac{1}{4} \times [4 + V_0(BR)] \\
 &= \frac{3}{4} \times [-1 + 0] + \frac{1}{4} \times [4 + 0] \\
 &= \frac{1}{4}
 \end{aligned}$$

Policy Evaluation Example

 V_0

TL 0	TR 0
BL 0	BR 0

 V_1

TL -1	TR $1/4$
BL $1/4$	BR 0

Policy Evaluation Example

V_0		V_1		V_2		...
TL	TR	TL	TR	TL	TR	
0	0	-1	$1/4$	$-1/8$	$-1/2$	
BL	BR	BL	BR	BL	BR	
0	0	$1/4$	0	$5/2$	0	

V_π obtained

Policy Iteration

Start with any π

1. Policy evaluation: given π , compute $V_\pi = F_\pi^\infty V$



Policy Iteration

Start with any π

1. **Policy evaluation:** given π , compute $V_\pi = F_\pi^\infty V$



2. **Policy improvement:** given V_π , compute improved policy π'



Policy Iteration

Start with any π

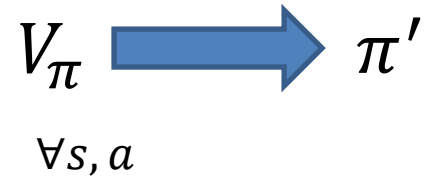
1. Policy evaluation: given π , compute $V_\pi = F_\pi^\infty V$



2. Policy improvement: given V_π , compute improved policy π'

a) Compute

$$Q_\pi(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_\pi(s')],$$



Policy Iteration

Start with any π

1. **Policy evaluation:** given π , compute $V_\pi = F_\pi^\infty V$



2. **Policy improvement:** given V_π , compute improved policy π'

a) Compute

$$Q_\pi(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_\pi(s')],$$



$\forall s, a$

b) Compute $\pi'(s) = \arg \max_{a \in \text{Actions}(s)} Q_\pi(s, a), \forall s$

Policy Iteration

Start with any π

1. **Policy evaluation:** given π , compute $V_\pi = F_\pi^\infty V$



2. **Policy improvement:** given V_π , compute improved policy π'

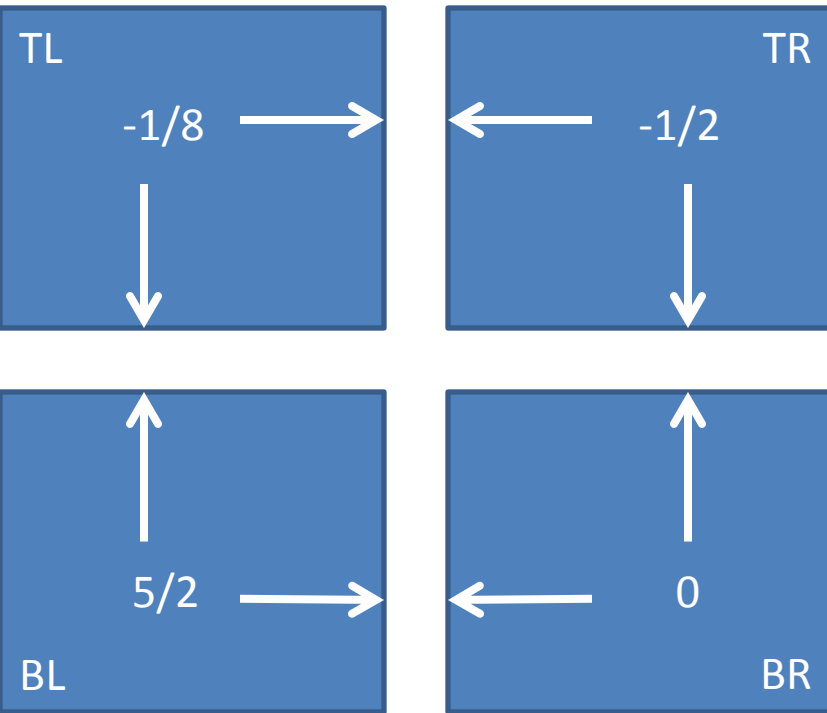


$$\pi'(s) = \arg \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_\pi(s')], \quad \forall s$$

Policy Improvement Example

Given V_π , evaluate π' (improved policy)

Find $Q_\pi(s, a)$ for all states, choose best action

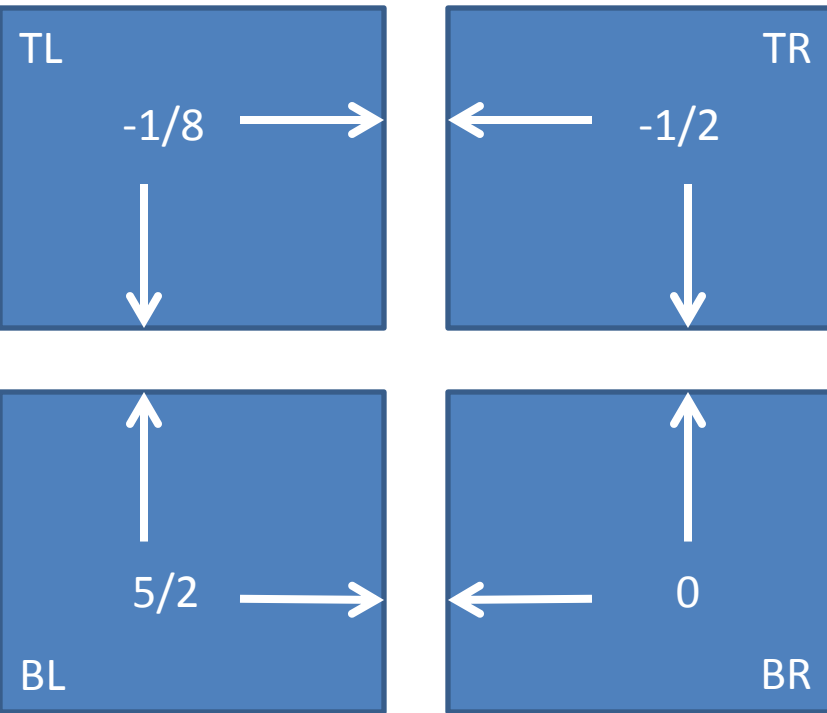


Policy Improvement Example

Given V_π , evaluate π' (improved policy)

Find $Q_\pi(s, a)$ for all states, choose best action

$$Q_\pi(TL, H) = \frac{3}{4} \times \left[-1 + \frac{-1}{2} \right] + \frac{1}{4} \times \left[-1 + \frac{5}{2} \right] = \frac{-3}{4}$$

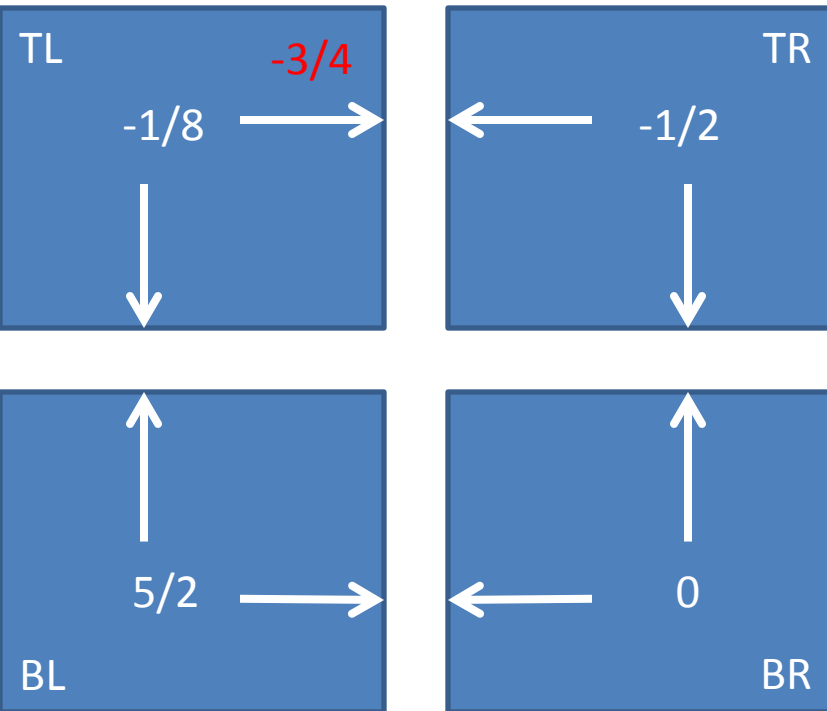


Policy Improvement Example

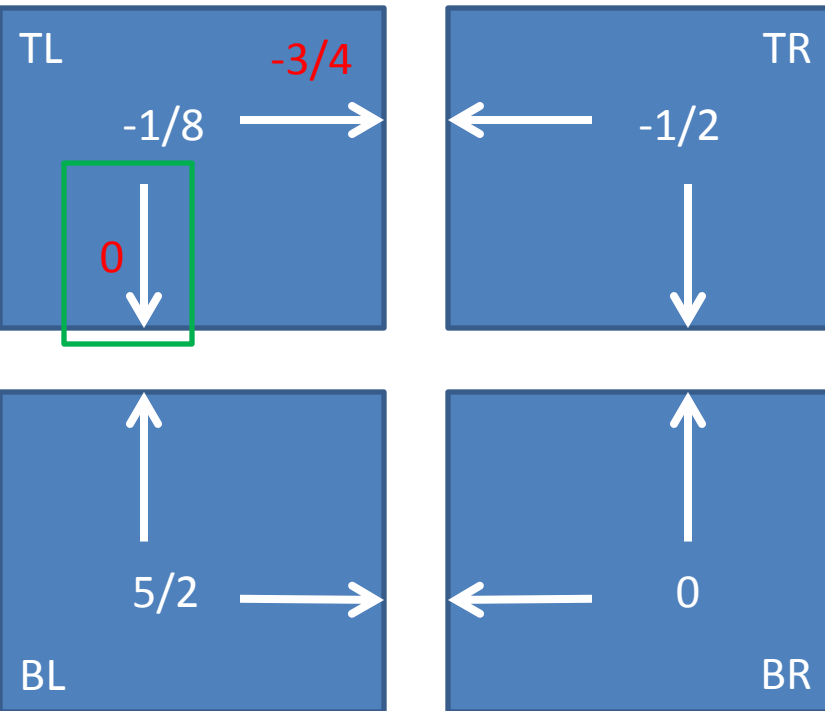
Given V_π , evaluate π' (improved policy)

Find $Q_\pi(s, a)$ for all states, choose best action

$$Q_\pi(TL, V) = \frac{1}{4} \times \left[-1 + \frac{-1}{2} \right] + \frac{3}{4} \times \left[-1 + \frac{5}{2} \right] = 0$$



Policy Improvement Example

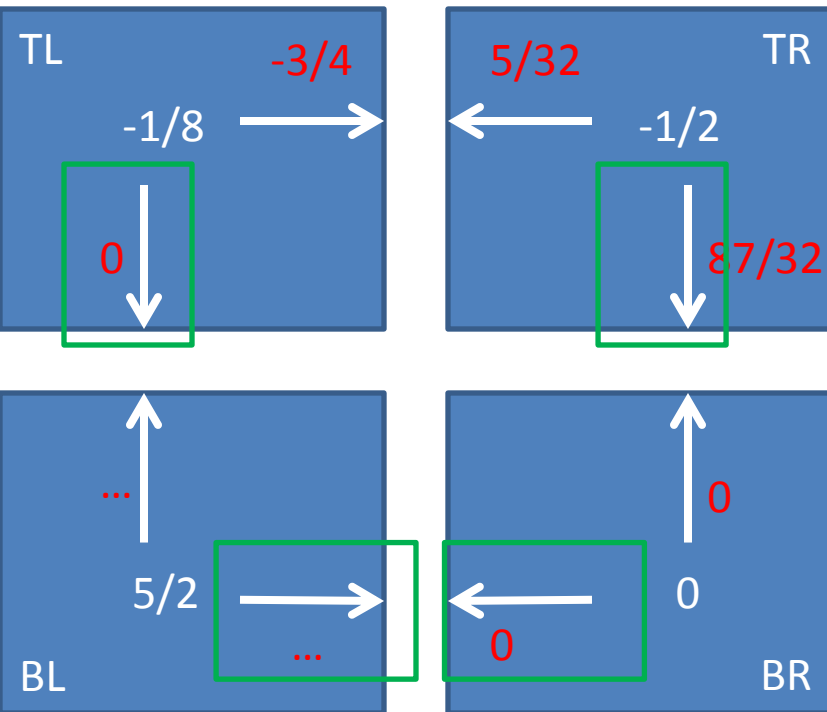


Given V_π , evaluate π' (improved policy)

Find $Q_\pi(s, a)$ for all states, choose best action

Optimal action
 $\pi'(TL) = V$

Policy Improvement Example



Given V_π , evaluate π' (improved policy)

Find $Q_\pi(s, a)$ for all states, choose best action

Optimal action

$$\pi'(TL) = V$$

$$\pi'(TR) = V$$

$$\pi'(BL) = H$$

Policy Iteration

Start with any π

1. **Policy evaluation:** given π , compute $V_\pi = F_\pi^\infty V$

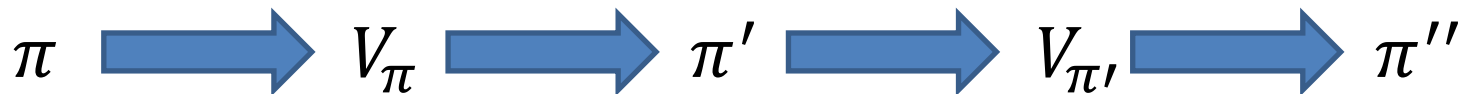


2. **Policy improvement:** given V_π , compute improved policy π'



$$\pi'(s) = \arg \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_\pi(s')], \quad \forall s$$

3. **Repeat** until policy converges/does not change



Policy Iteration

Start with any π

1. **Policy evaluation:** given π , compute $V_\pi = F_\pi^\infty V$

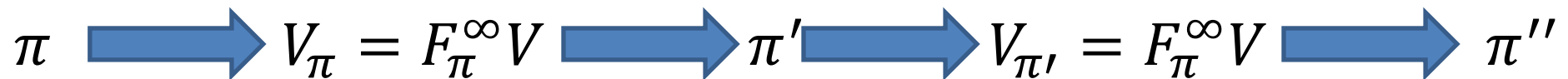


2. **Policy improvement:** given V_π , compute improved policy π'



$$\pi'(s) = \arg \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_\pi(s')], \quad \forall s$$

3. **Repeat** until policy converges/does not change



Policy Iteration

Start with any π

1. **Policy evaluation:** given π , compute $V_\pi = F_\pi^\infty V$

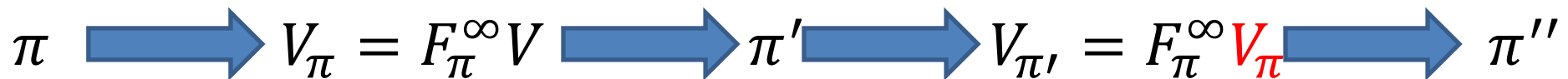


2. **Policy improvement:** given V_π , compute improved policy π'



$$\pi'(s) = \arg \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_\pi(s')], \quad \forall s$$

3. **Repeat** until policy converges/does not change



Warm start!

Value Iteration

Start with any V

1. Compute Q-values: given V , compute Q

$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \quad \forall s, a$$

Value Iteration

Start with any V

1. Compute Q-values: given V , compute Q

$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \quad \forall s, a$$

2. Compute V' : given Q , compute V'

$$V'(s) = \max_{a \in \text{Actions}(s)} Q(s, a), \quad \forall s$$

Value Iteration

Start with any V

1. Compute Q-values: given V , compute Q

$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \quad \forall s, a$$

2. Compute V' : given Q , compute V'

$$V'(s) = \max_{a \in \text{Actions}(s)} Q(s, a), \quad \forall s$$

3. Repeat until value function converges

Value Iteration

Start with any V

1. **Compute V'** : given V , compute V'

$$V'(s) = \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \quad \forall s$$

2. **Repeat** until value function converges

Value Iteration

Start with any V

1. Compute V' : given V , compute π' , then compute V'

a) Compute $\pi'(s) = \arg \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \forall s$

b) Compute $V' = F_{\pi'} V$

2. Repeat until value function converges

Value Iteration

Start with any V

1. Compute V' : given V , compute π' , then compute V'

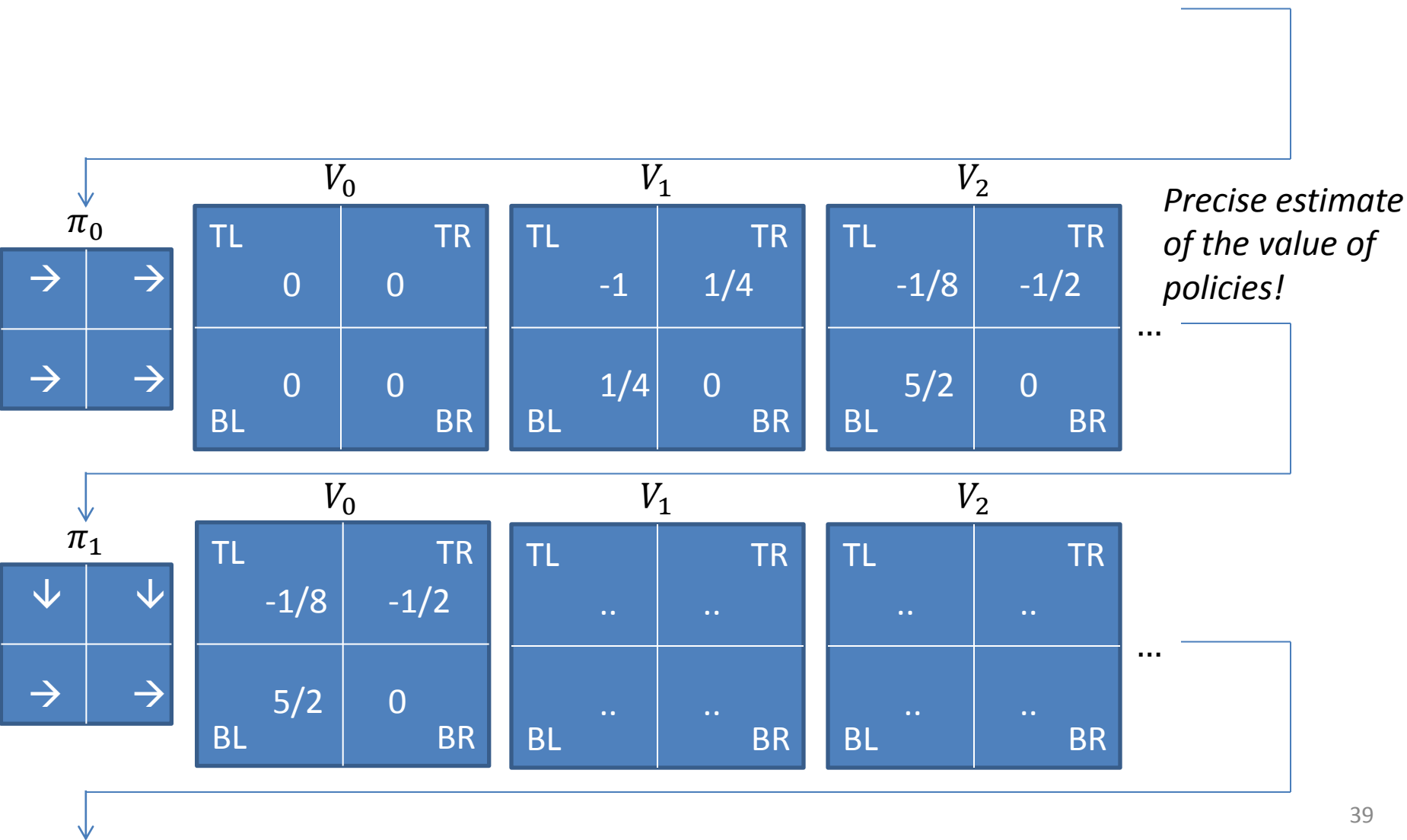
a) Compute $\pi'(s) = \arg \max_{a \in \text{Actions}(s)} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \forall s$

b) Compute $V' = F_{\pi'} V$

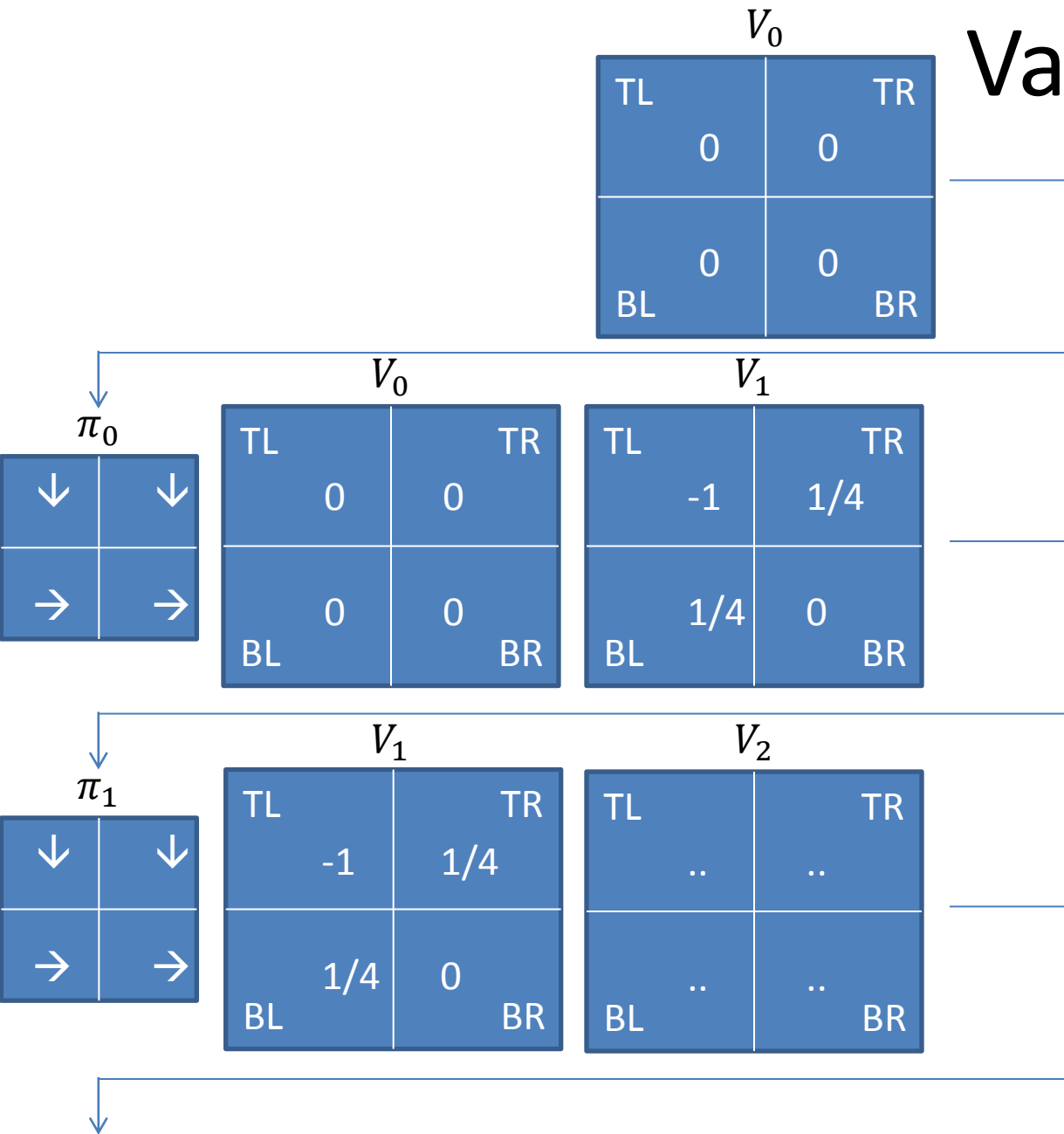
2. Repeat until value function converges



Policy Iteration

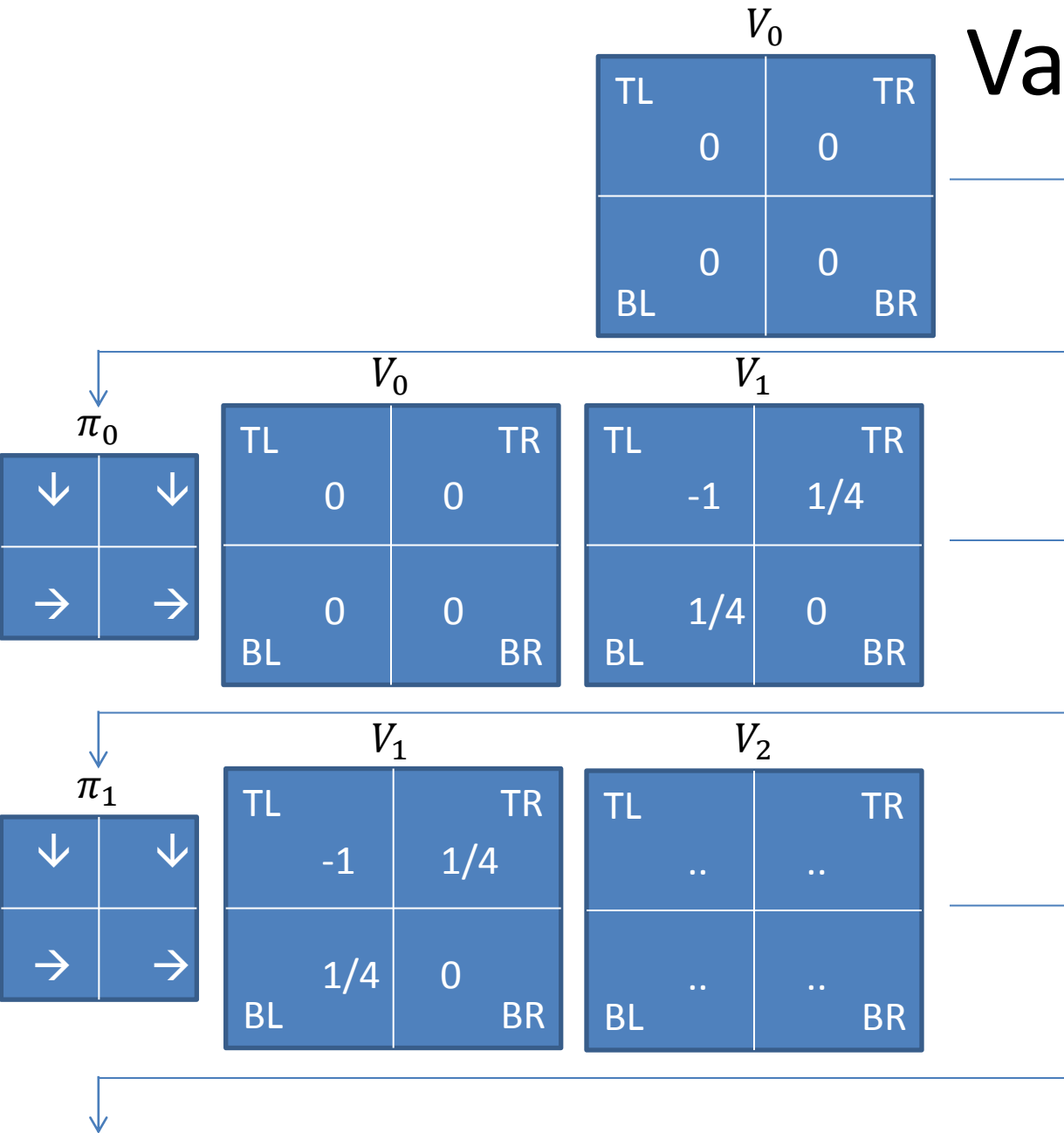


Value Iteration



Crude estimate of the value of policies!

Value Iteration



Crude estimate of the value of policies!

Exact numbers on slide may not be correct

Policy Iteration vs Value Iteration

Policy Iteration

$$\pi \longrightarrow V_{\pi} = F_{\pi}^{\infty} V \longrightarrow \pi' \longrightarrow V_{\pi'} = F_{\pi'}^{\infty} V$$

$$V \longrightarrow \pi \longrightarrow V' = F_{\pi} V \longrightarrow \pi' \longrightarrow V'' = F_{\pi'} V'$$

Value Iteration

Policy Iteration vs Value Iteration

Policy Iteration

Precise Policy Evaluation + Policy Improvement

$$\pi \longrightarrow V_{\pi} = F_{\pi}^{\infty} V \longrightarrow \pi' \longrightarrow V_{\pi'} = F_{\pi'}^{\infty} V$$

$$V \longrightarrow \pi \longrightarrow V' = F_{\pi} V \longrightarrow \pi' \longrightarrow V'' = F_{\pi'} V'$$

Crude Policy Evaluation + Policy Improvement

Value Iteration

Policy Iteration vs Value Iteration

SARSA vs Q-Learning

Function Approximation and Q-Learning

Reinforcement Learning (RL)

$s \in \text{States}$

Actions(s): actions from state s

$T(s, a, s')$: probability next state is s' , given action a in state s

Note that $\sum_{s'} T(s, a, s') = 1$ for all $s \in \text{States}, a \in \text{Actions}(s)$

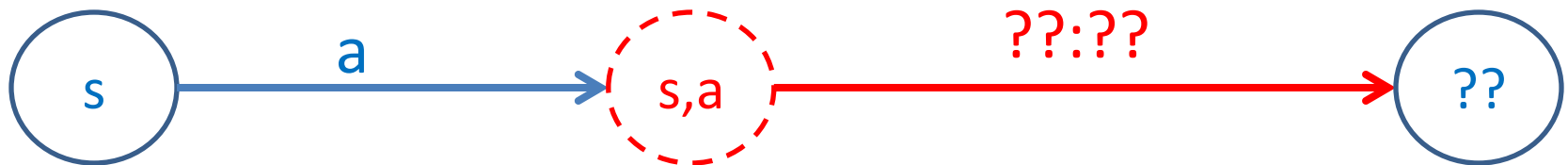
UNKNOWN

$R(s, a, s')$: reward for transition (s, a, s')

IsEnd(s): whether state marks end of game

Alternatively, $T(s, a, s) = 1$ and $R(s, a, s) = 0$

$0 \leq \gamma \leq 1$: discount factor



Reinforcement Learning (RL)

Algorithm	Estimate	Comments
Model-based Monte Carlo	$T(s, a, s'), R(s, a, s')$	Find Q_{opt}, π_{opt} afterward
Model-free Monte Carlo	Q_π using utility u	No optimal policy (on-policy)
SARSA	Q_π using utility r, Q_π	No optimal policy (on-policy)
Q-Learning	Q_{opt} using utility r, Q_{opt}	Optimal policy (off-policy)

Reinforcement Learning (RL)

Algorithm	Estimate	Comments
Model-based Monte Carlo	$T(s, a, s'), R(s, a, s')$	Find Q_{opt}, π_{opt} afterward
Model-free Monte Carlo	Q_π using utility u	No optimal policy (on-policy)
SARSA	Q_π using utility r, Q_π	No optimal policy (on-policy)
Q-Learning	Q_{opt} using utility r, Q_{opt}	Optimal policy (off-policy)

Starting in state s_0 , apply policy π
 s_0

s_0 ;

Reinforcement Learning (RL)

Algorithm	Estimate	Comments
Model-based Monte Carlo	$T(s, a, s'), R(s, a, s')$	Find Q_{opt}, π_{opt} afterward
Model-free Monte Carlo	Q_π using utility u	No optimal policy (on-policy)
SARSA	Q_π using utility r, Q_π	No optimal policy (on-policy)
Q-Learning	Q_{opt} using utility r, Q_{opt}	Optimal policy (off-policy)

Starting in state s_0 , apply policy π

$s_0; \mathbf{a}_1$

s_0

$\mathbf{a}_1 \sim \pi(s_0)$

Reinforcement Learning (RL)

Algorithm	Estimate	Comments
Model-based Monte Carlo	$T(s, a, s'), R(s, a, s')$	Find Q_{opt}, π_{opt} afterward
Model-free Monte Carlo	Q_π using utility u	No optimal policy (on-policy)
SARSA	Q_π using utility r, Q_π	No optimal policy (on-policy)
Q-Learning	Q_{opt} using utility r, Q_{opt}	Optimal policy (off-policy)

Starting in state s_0 , apply policy π

$s_0; \mathbf{a}_1, r_1, s_1$

s_0

$\mathbf{a}_1 \sim \pi(s_0)$

$s_1 \sim T(s_0, a_1, \cdot)$

$r_1 \sim P_R(s_0, a_1, s_1)$

Reinforcement Learning (RL)

Algorithm	Estimate	Comments
Model-based Monte Carlo	$T(s, a, s'), R(s, a, s')$	Find Q_{opt}, π_{opt} afterward
Model-free Monte Carlo	Q_π using utility u	No optimal policy (on-policy)
SARSA	Q_π using utility r, Q_π	No optimal policy (on-policy)
Q-Learning	Q_{opt} using utility r, Q_{opt}	Optimal policy (off-policy)

Starting in state s_0 , apply policy π

$s_0; \mathbf{a}_1, r_1, s_1; \mathbf{a}_2, r_2, s_2; \mathbf{a}_3, r_3, s_3$

s_0

$\mathbf{a}_1 \sim \pi(s_0)$

$s_1 \sim T(s_0, a_1, \cdot)$

$r_1 \sim P_R(s_0, a_1, s_1)$

$\mathbf{a}_2 \sim \pi(s_1)$

$s_2 \sim T(s_1, a_2, \cdot)$

$r_2 \sim P_R(s_1, a_2, s_2)$

\vdots

SARSA vs Q-Learning

SARSA: Update $\hat{Q}_\pi(s, a)$ based on (s, a, r, s', a')

$$\hat{Q}_\pi(s, a) \leftarrow (1 - \eta)\hat{Q}_\pi(s, a) + \eta[r + \gamma\hat{Q}_\pi(s', a')]$$

Helps evaluate policy π by sampling transitions and rewards from policy π
 \Rightarrow **On-policy** algorithm

SARSA evaluates policy π . Policy Improvement comes next. Iterate.

SARSA vs Q-Learning

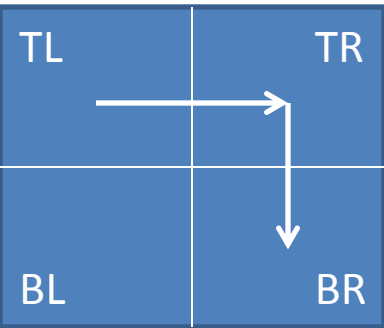
SARSA: Update $\hat{Q}_\pi(s, a)$ based on (s, a, r, s', a')

$$\hat{Q}_\pi(s, a) \leftarrow \hat{Q}_\pi(s, a) - \eta [\hat{Q}_\pi(s, a) - (r + \gamma \hat{Q}_\pi(s', a'))]$$

Helps evaluate policy π by sampling transitions and rewards from policy π
 \Rightarrow **On-policy** algorithm

SARSA evaluates policy π . Policy Improvement comes next. Iterate.

SARSA Example

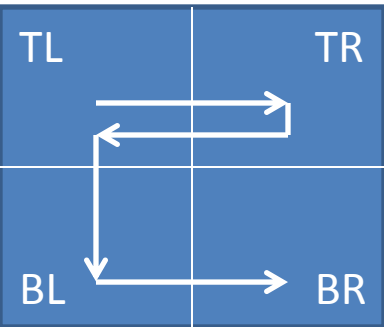


Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

SARSA Example



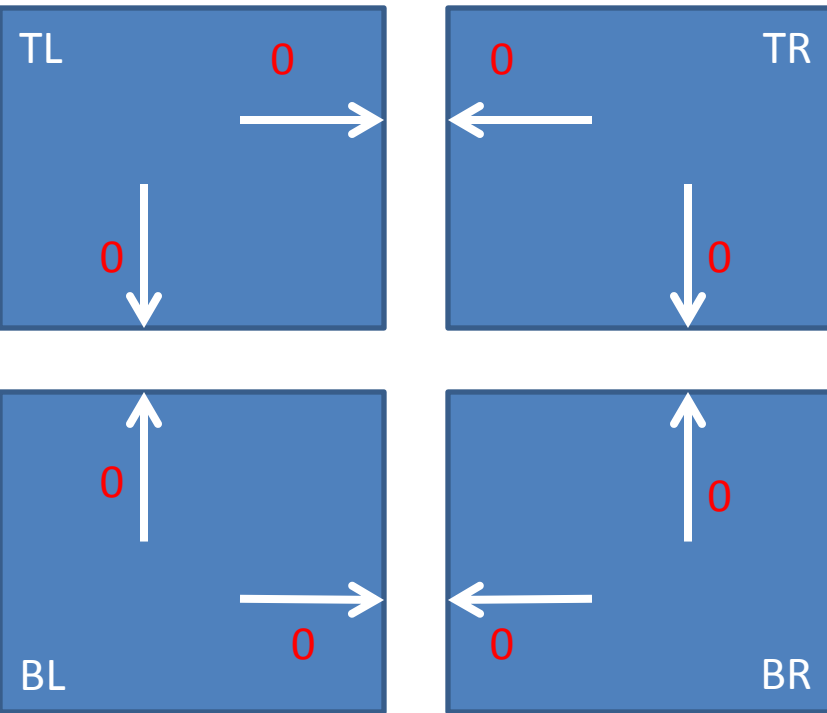
Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

SARSA Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

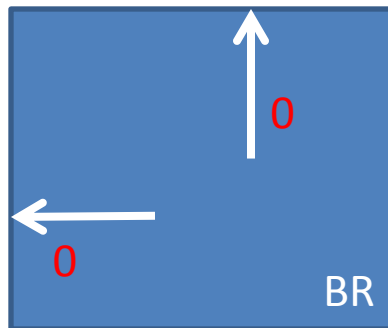
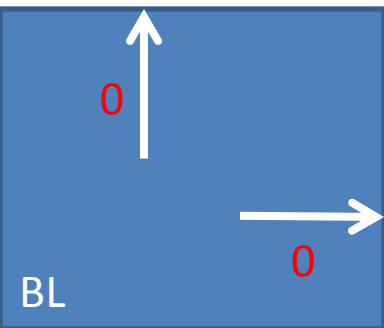
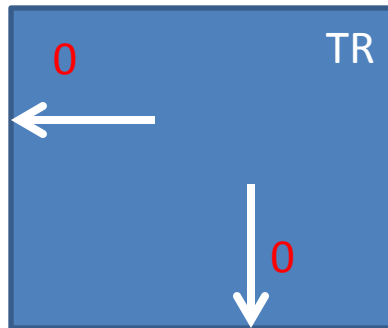
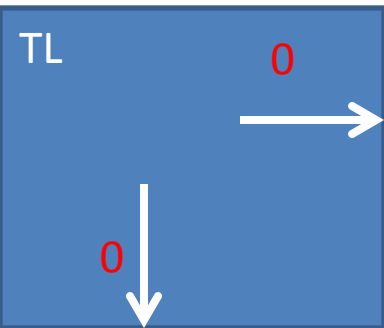
1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

SARSA

Set $Q_{\pi}(s, a) = 0$

SARSA Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. **TL**; **H**, -1, **TR**; **H**, 4, **BR**

2. **TL**; **H**, -1, **TR**; **H**, -1, **TL**; **H**, -1, **BL**; **H**, 4, **BR**

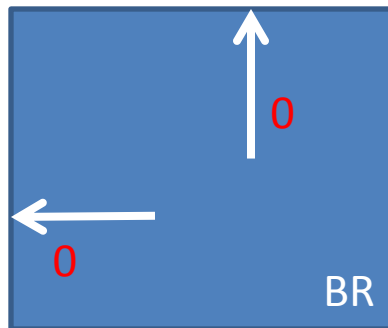
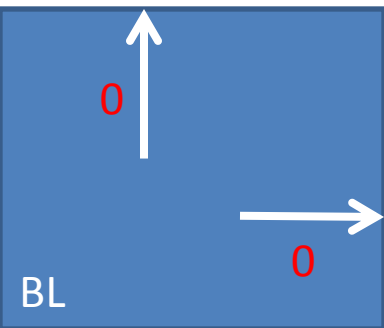
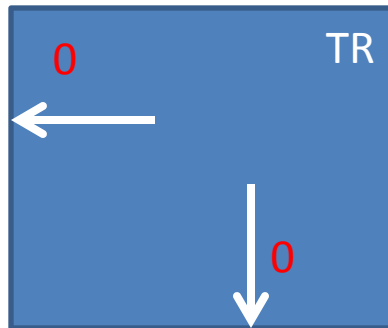
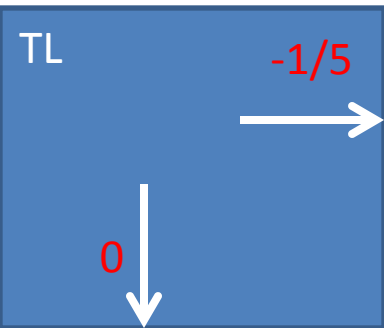
SARSA

Set $Q_{\pi}(s, a) = 0$

Consider **TL**, **H**, -1, **TR**, **H**
 $s \quad a \quad r \quad s' \quad a'$

$$\begin{aligned} Q_{\pi}(TL, H) &:= \frac{4}{5} \times Q_{\pi}(TL, H) + \frac{1}{5} \times [-1 + Q_{\pi}(TR, H)] \\ &= \frac{4}{5} \times 0 + \frac{1}{5} \times [-1 + 0] = \frac{-1}{5} \end{aligned}$$

SARSA Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

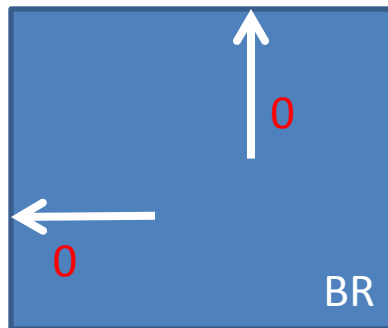
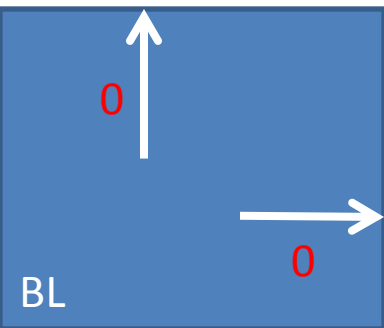
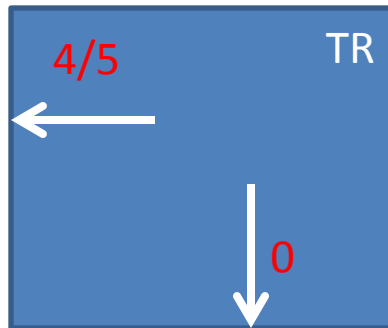
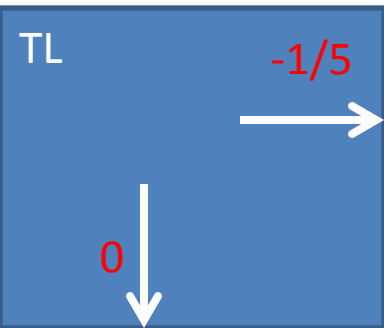
SARSA

Set $Q_{\pi}(s, a) = 0$

Consider $TR, H, 4, BR, -$
 $s \quad a \quad r \quad s' \quad a'$

$$\begin{aligned} Q_{\pi}(TR, H) &:= \frac{4}{5} \times Q_{\pi}(TR, H) + \frac{1}{5} \times [4 + Q_{\pi}(BR, -)] \\ &= \frac{4}{5} \times 0 + \frac{1}{5} \times [4 + 0] = \frac{4}{5} \end{aligned}$$

SARSA Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

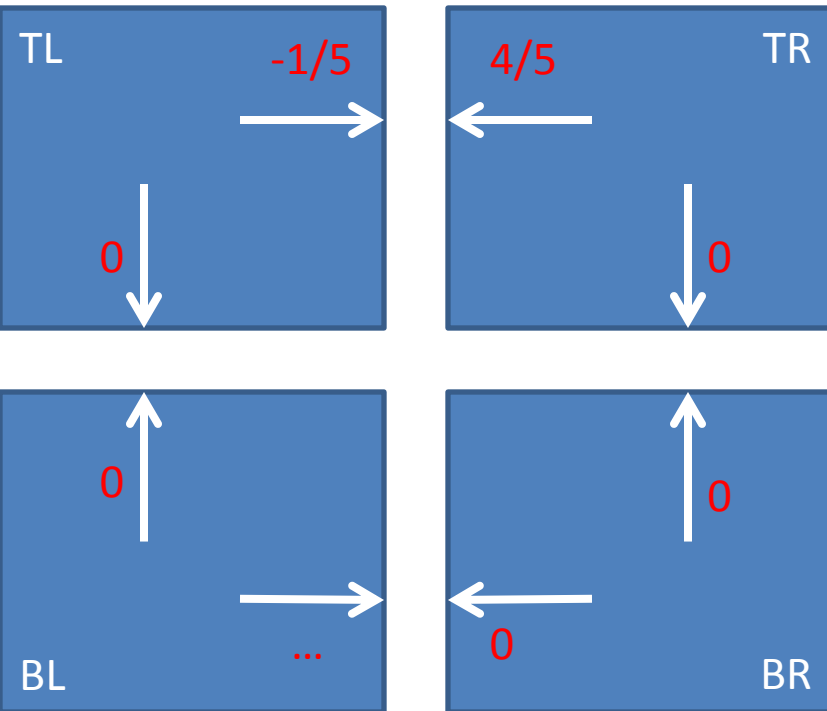
SARSA

Set $Q_\pi(s, a) = 0$

Consider $TL, H, -1, TR, H$
 $s \quad a \quad r \quad s' \quad a'$

$Q_\pi(TL, H) := \dots$

SARSA Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. *TL; H, -1, TR; H, 4, BR*

2. *TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR*

SARSA

Set $Q_\pi(s, a) = 0$

Consider $TL, H, -1, TR, H$
 $s \quad a \quad r \quad s' \quad a'$

$Q_\pi(TL, H) := \dots$

Like Policy Evaluation – given enough iterations, get precise estimates

SARSA vs Q-Learning

SARSA: Update $\hat{Q}_\pi(s, a)$ based on (s, a, r, s', a')

$$\hat{Q}_\pi(s, a) \leftarrow (1 - \eta)\hat{Q}_\pi(s, a) + \eta[r + \gamma\hat{Q}_\pi(s', a')]$$

Helps evaluate policy π by sampling transitions and rewards from policy π
 \Rightarrow **On-policy** algorithm

SARSA evaluates policy π . Policy Improvement comes next. Iterate.

Q-Learning: Update $\hat{Q}_{opt}(s, a)$ based on (s, a, r, s')

$$\hat{Q}_{opt}(s, a) \leftarrow (1 - \eta)\hat{Q}_{opt}(s, a) + \eta[r + \gamma\hat{V}_{opt}(s')]$$

Helps evaluate optimal policy by sampling transitions and rewards from policy π
 \Rightarrow **Off-policy** algorithm

Use (randomized) policy to explore π to explore all (s, a) , find greedy policy!

SARSA vs Q-Learning

SARSA: Update $\hat{Q}_\pi(s, a)$ based on (s, a, r, s', a')

$$\hat{Q}_\pi(s, a) \leftarrow (1 - \eta)\hat{Q}_\pi(s, a) + \eta[r + \gamma\hat{Q}_\pi(s', a')]$$

Helps evaluate policy π by sampling transitions and rewards from policy π
 \Rightarrow **On-policy** algorithm

SARSA evaluates policy π . Policy Improvement comes next. Iterate.

Q-Learning: Update $\hat{Q}_{opt}(s, a)$ based on (s, a, r, s')

$$\hat{Q}_{opt}(s, a) \leftarrow \hat{Q}_{opt}(s, a) - \eta[\hat{Q}_{opt}(s, a) - (r + \gamma\hat{V}_{opt}(s'))]$$

Helps evaluate optimal policy by sampling transitions and rewards from policy π
 \Rightarrow **Off-policy** algorithm

Use (randomized) policy to explore π to explore all (s, a) , find greedy policy!

SARSA vs Q-Learning

SARSA: Update $\hat{Q}_\pi(s, a)$ based on (s, a, r, s', a')

$$\hat{Q}_\pi(s, a) \leftarrow (1 - \eta)\hat{Q}_\pi(s, a) + \eta[r + \gamma\hat{Q}_\pi(s', a')]$$

Helps evaluate policy π by sampling transitions and rewards from policy π
 \Rightarrow **On-policy** algorithm

SARSA evaluates policy π . Policy Improvement comes next. Iterate.

Q-Learning: Update $\hat{Q}_{opt}(s, a)$ based on (s, a, r, s')

$$\hat{Q}_{opt}(s, a) \leftarrow (1 - \eta)\hat{Q}_{opt}(s, a) + \eta \left[r + \gamma \max_{a' \in A(s')} \hat{Q}_{opt}(s', a') \right]$$

Helps evaluate optimal policy by sampling transitions and rewards from policy π
 \Rightarrow **Off-policy** algorithm

Use (randomized) policy to explore π to explore all (s, a) , find greedy policy!

SARSA vs Q-Learning

SARSA: Update $\hat{Q}_{\pi}(s, a)$ based on (s, a, r, s', a')

$$\hat{Q}_{\pi}(s, a) \leftarrow (1 - \eta)\hat{Q}_{\pi}(s, a) + \eta[r + \gamma\hat{Q}_{\pi}(s', a')]$$

Helps evaluate policy π by sampling transitions and rewards from policy π
 \Rightarrow **On-policy** algorithm

SARSA evaluates policy π . Policy Improvement comes next. Iterate.

Q-Learning: Update $\hat{Q}_{opt}(s, a)$ based on (s, a, r, s')

$$\hat{Q}_{opt}(s, a) \leftarrow (1 - \eta)\hat{Q}_{opt}(s, a) + \eta\left[r + \gamma \max_{a' \in A(s')} \hat{Q}_{opt}(s', a')\right]$$

Helps evaluate optimal policy by sampling transitions and rewards from policy π
 \Rightarrow **Off-policy** algorithm

Use (randomized) policy to explore π to explore all (s, a) , find greedy policy!

SARSA vs Q-Learning

SARSA: Update $\hat{Q}_{\pi}(s, a)$ based on (s, a, r, s', a')

$$\hat{Q}_{\pi}(s, a) \leftarrow (1 - \eta)\hat{Q}_{\pi}(s, a) + \eta[r + \gamma\hat{Q}_{\pi}(s', a')]$$

Helps evaluate policy π by sampling transitions and rewards from policy π
 \Rightarrow **On-policy** algorithm

SARSA evaluates policy π . Policy Improvement comes next. Iterate.

Q-Learning: Update $\hat{Q}_{opt}(s, a)$ based on (s, a, r, s')

$$\hat{Q}_{opt}(s, a) \leftarrow (1 - \eta)\hat{Q}_{opt}(s, a) + \eta\left[r + \gamma \max_{a' \in A(s')} \hat{Q}_{opt}(s', a')\right]$$

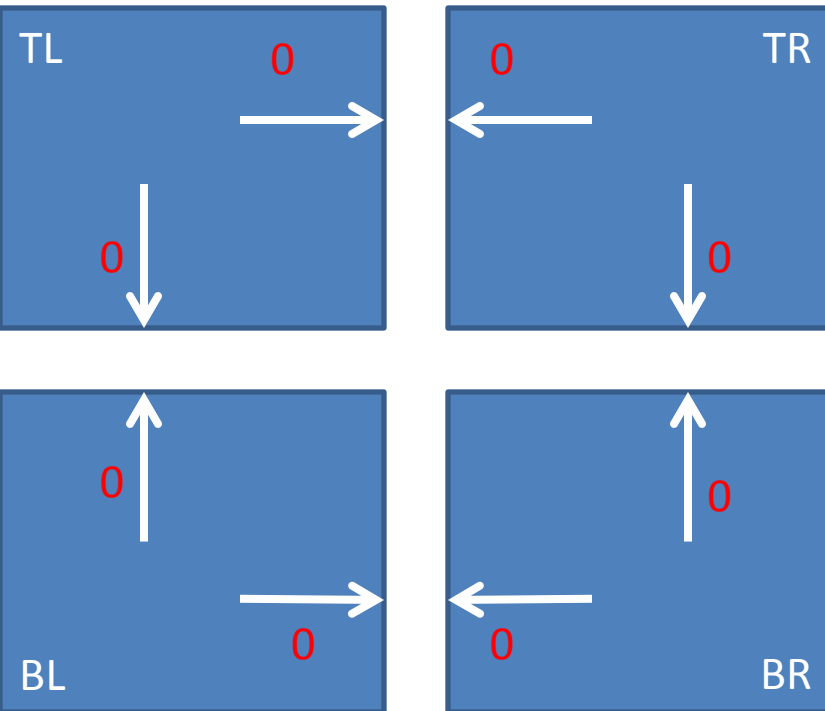
Helps evaluate optimal policy by sampling transitions and rewards from policy π
 \Rightarrow **Off-policy** algorithm

Example from notes:

Use (randomized) policy to explore π to explore all (s, a) , find greedy policy!

<http://web.stanford.edu/class/cs221/lectures/index.html#include=mdp2.js&slideIndex=46>

Q-Learning Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. **TL**; **H**, -1, **TR**; **H**, 4, **BR**

2. **TL**; **H**, -1, **TR**; **H**, -1, **TL**; **H**, -1, **BL**; **H**, 4, **BR**

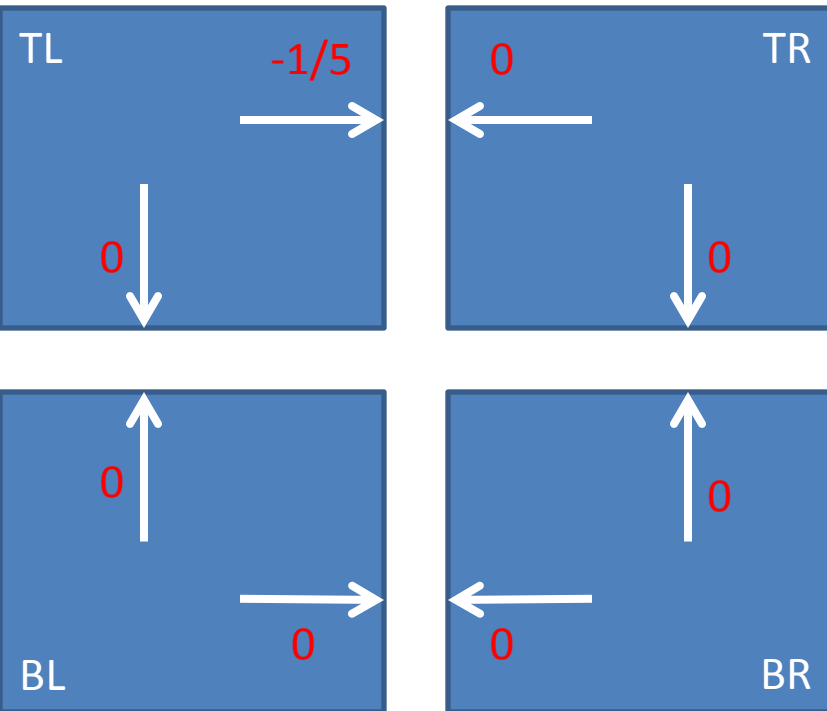
Q - learning

Set $Q_{opt}(s, a) = 0$

*Consider **TL**, **H**, -1, **TR***
 $s \quad a \quad r \quad s'$

$$Q_{opt}(TL, H) := \frac{4}{5} \times Q_{opt}(TL, H) + \frac{1}{5} \times \left[-1 + \max_{a \in A(s)} Q_{opt}(TR, a) \right]$$

Q-Learning Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

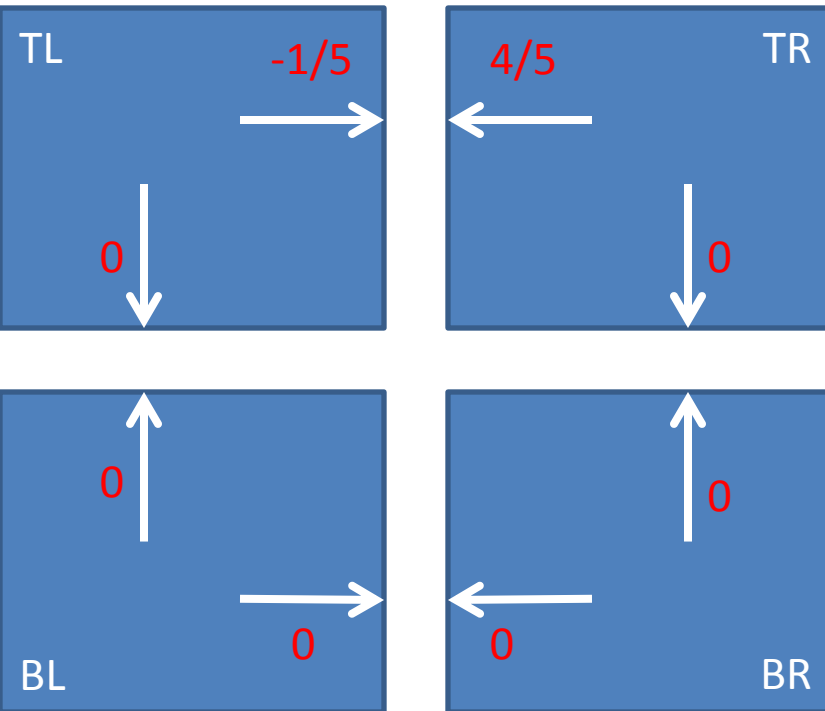
Q-learning

Set $Q_{opt}(s, a) = 0$

Consider $TL, H, -1, TR$
 $s \quad a \quad r \quad s'$

$$\begin{aligned}
 Q_{opt}(TL, H) &:= \frac{4}{5} \times Q_{opt}(TL, H) + \frac{1}{5} \times \left[-1 + \max_{a \in A(s)} Q_{opt}(TR, a) \right] \\
 &= \frac{4}{5} \times 0 + \frac{1}{5} \times \left[-1 + 0 \right] = \frac{-1}{5}
 \end{aligned}$$

Q-Learning Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

Q - learning

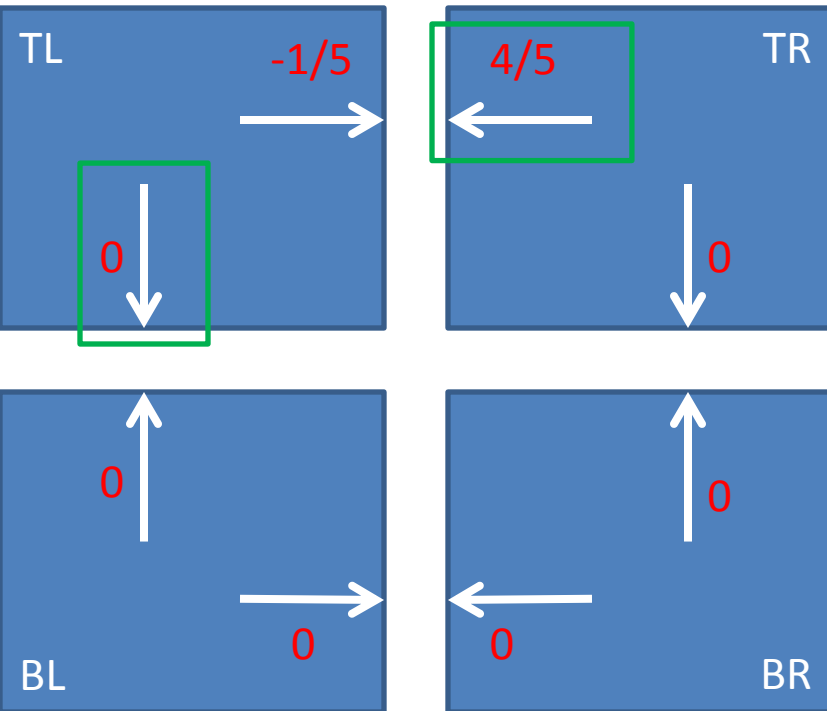
Set $Q_{opt}(s, a) = 0$

Consider $TR, H, 4, BR$

$s \quad a \quad r \quad s'$

$$\begin{aligned}
 Q_{opt}(TR, H) &:= \frac{4}{5} \times Q_{opt}(TR, H) + \frac{1}{5} \times \left[4 + \max_{a \in A(s)} Q_{opt}(BR, a) \right] \\
 &= \frac{4}{5} \times 0 + \frac{1}{5} \times [4 + 0] = \frac{4}{5}
 \end{aligned}$$

Q-Learning Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

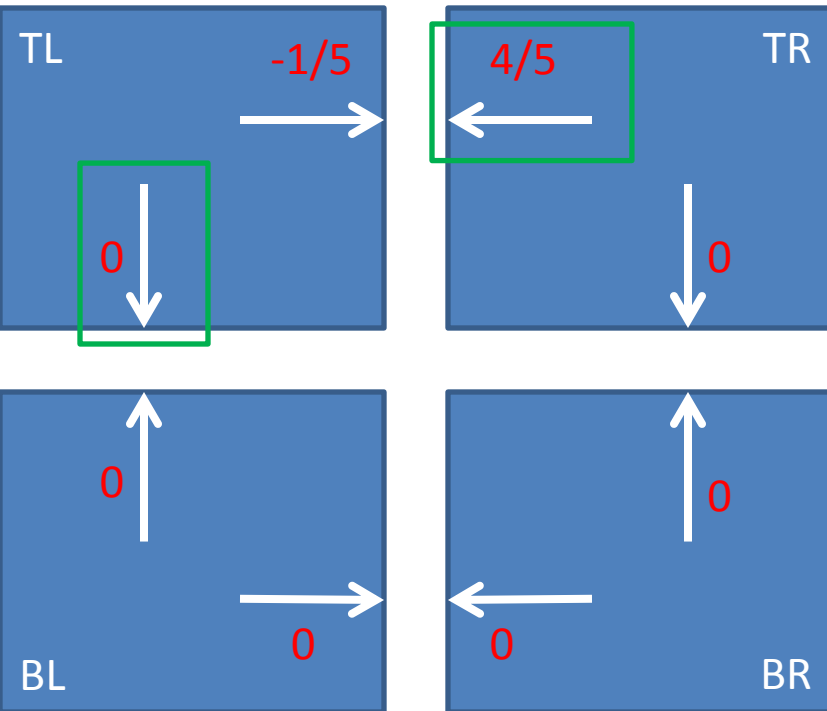
2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

Q – learning

Set $Q_{opt}(s, a) = 0$

Optimal policy different from π
 \Rightarrow off – policy algorithm

Q-Learning Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. *TL; H, -1, TR; H, 4, BR*

2. *TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR*

Q – learning

Set $Q_{opt}(s, a) = 0$

Optimal policy different from π
 \Rightarrow off – policy algorithm

Problems: $Q(s, V)$ never explored

Need to use exploration/randomized policy

Policy Iteration vs Value Iteration

SARSA vs Q-Learning

Function Approximation and Q-Learning

Function Approximation

For optimal policy, compute $Q_{opt}(s, a)$ for all states $s \in S, a \in A(s)$

Our current setup does not generalize to unseen states!

What if

- Number of states $|S|$ is huge?
 - Extreme: what if continuous state? Infinitely many of them...
- Number of actions $|A(s)|$ is huge?
- Majority of neighboring states should have similar Q-values

Function Approximation

IDEA:

Express $Q_{opt}(s, a)$ as a weighted combination of its features

$$Q_{opt}(s, a) \approx w \cdot \Phi(s, a) = w_1\phi_1(s, a) + w_2\phi_2(s, a) + \dots + w_n\phi_n(s, a)$$

$\Phi(s, a)$ = Feature vector

w = Weights

Function Approximation

IDEA:

Express $Q_{opt}(s, a)$ as a weighted combination of its features

$$Q_{opt}(s, a) \approx w \cdot \Phi(s, a) = Q_{opt}(s, a; w)$$

Q-Learning with function approximation: Update $\hat{Q}_{opt}(s, a; w)$ based on (s, a, r, s')

$$w \leftarrow w - \eta \left[\hat{Q}_{opt}(s, a; w) - \left(r + \gamma \hat{V}_{opt}(s'; w) \right) \right] \Phi(s, a)$$

Amounts to SGD on:

$$\min_w \sum_{(s, a, r, s')} \left(\hat{Q}_{opt}(s, a; w) - \left(r + \gamma \hat{V}_{opt}(s'; w) \right) \right)^2$$

Function Approximation

IDEA:

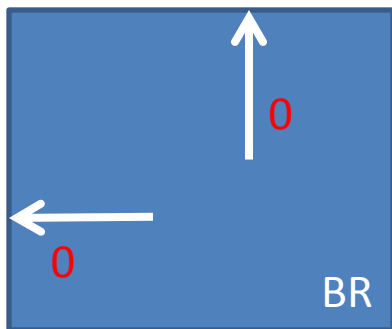
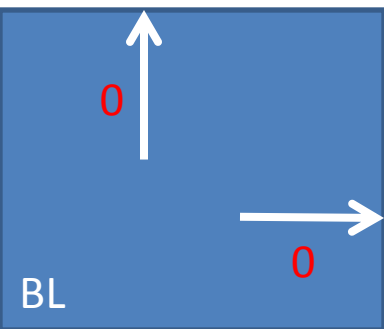
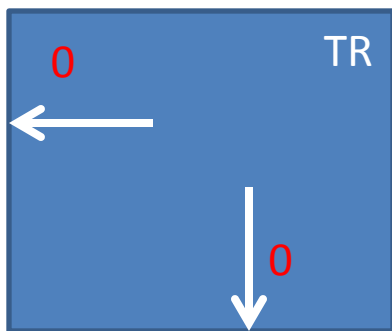
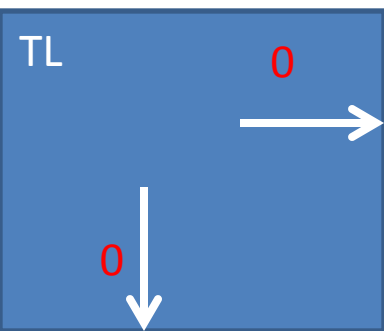
Express $Q_{opt}(s, a)$ as a weighted combination of its features

$$Q_{opt}(s, a) \approx w \cdot \Phi(s, a) = Q_{opt}(s, a; w)$$

Q-Learning with function approximation: Update $\hat{Q}_{opt}(s, a; w)$ based on (s, a, r, s')

$$w \leftarrow w - \eta \left[\hat{Q}_{opt}(s, a; w) - \left(r + \gamma \max_{a' \in A(s')} \hat{Q}_{opt}(s', a'; w) \right) \right] \Phi(s, a)$$

Q-Learning with Function Approximation Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

Q – learning with function approximation

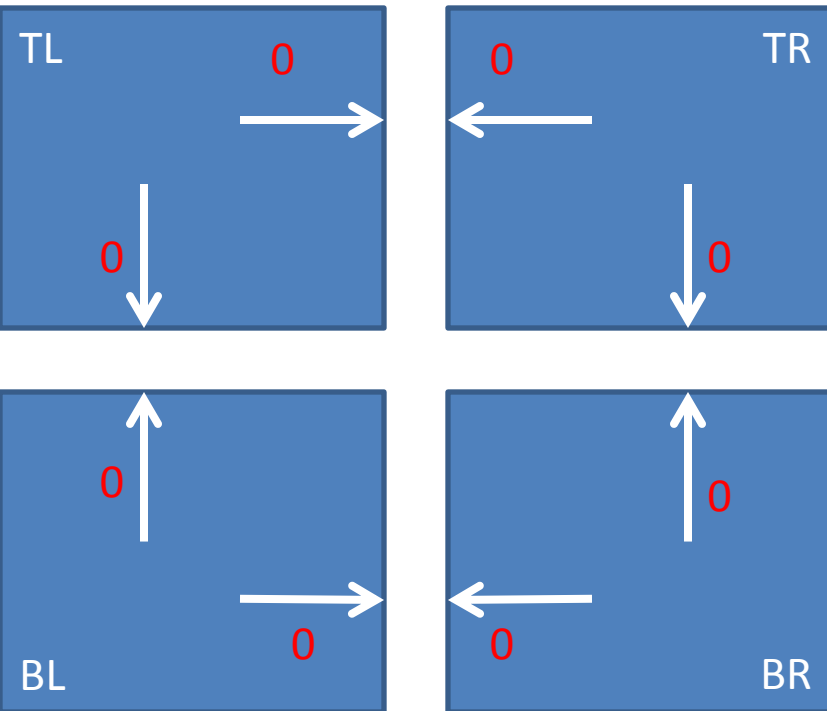
$\Phi(s, a) = [1, \mathbf{1} = \{s = * L\}, \mathbf{1} = \{a = H\},$

$\mathbf{1} = \{s = * L\} \mathbf{1} = \{a = H\}]^T$

$w = [w_1, w_2, w_3, w_4]^T$

Set $w = 0 \Rightarrow Q_{opt}(s, a; w) = 0$

Q-Learning with Function Approximation Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

Q – learning with function approximation

$\Phi(s, a) = [1, \mathbf{1} = \{s = * L\}, \mathbf{1} = \{a = H\},$

$\mathbf{1} = \{s = * L\} \mathbf{1} = \{a = H\}]^T$

$w = [w_1, w_2, w_3, w_4]^T$

Set $w = 0 \Rightarrow Q_{opt}(s, a; w) = 0$

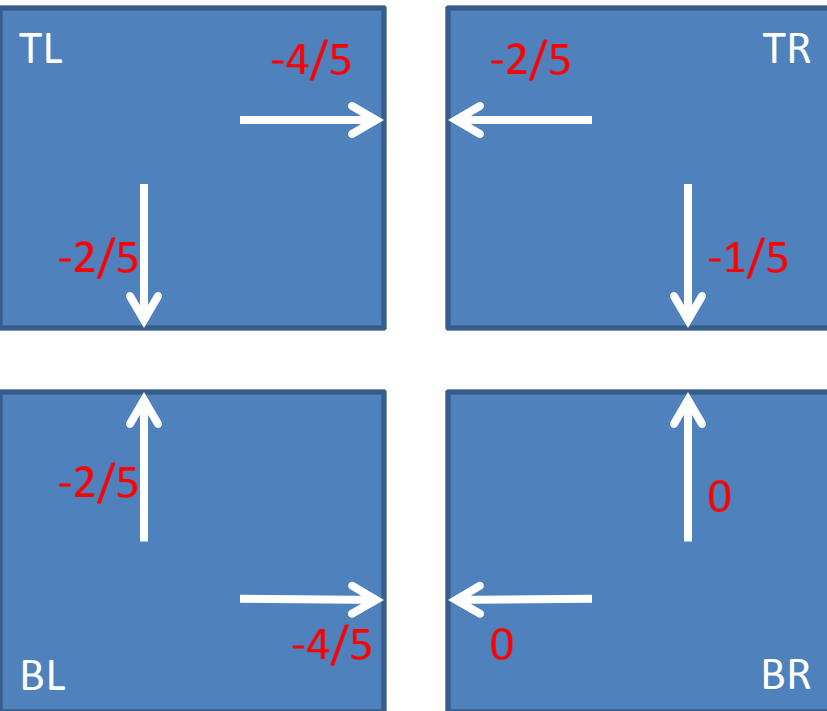
Consider $TL, H, -1, TR$

$s \quad a \quad r \quad s'$

$$w := [0 \ 0 \ 0 \ 0]^T - \frac{1}{5} \left[Q_{opt}(TL, H; w) - \left(-1 + \max_{a' \in A(TR)} Q_{opt}(TR, a'; w) \right) \right] \Phi(TL, H)$$

$$= [0 \ 0 \ 0 \ 0]^T - \frac{1}{5} \left[\begin{array}{cc} 0 & - \left(-1 + \max_{a' \in A(TR)} 0 \right) \end{array} \right] [1 \ 1 \ 1 \ 1]^T$$

Q-Learning with Function Approximation Example



Fixed policy $\pi(s) = H, \forall s$

Generate sample paths

1. $TL; H, -1, TR; H, 4, BR$

2. $TL; H, -1, TR; H, -1, TL; H, -1, BL; H, 4, BR$

Q – learning with function approximation

$\Phi(s, a) = [1, \mathbf{1} = \{s = * L\}, \mathbf{1} = \{a = H\},$

$\mathbf{1} = \{s = * L\} \mathbf{1} = \{a = H\}]^T$

$w = [w_1, w_2, w_3, w_4]^T$

Set $w = 0 \Rightarrow Q_{opt}(s, a; w) = 0$

Consider $TL, H, -1, TR$
 $s \quad a \quad r \quad s'$

$$w := [0 \ 0 \ 0 \ 0]^T - \frac{1}{5} \left[Q_{opt}(TL, H; w) - \left(-1 + \max_{a' \in A(TR)} Q_{opt}(TR, a'; w) \right) \right] \Phi(TL, H)$$

$$= \left[\frac{-1}{5} \ \frac{-1}{5} \ \frac{-1}{5} \ \frac{-1}{5} \right]^T$$