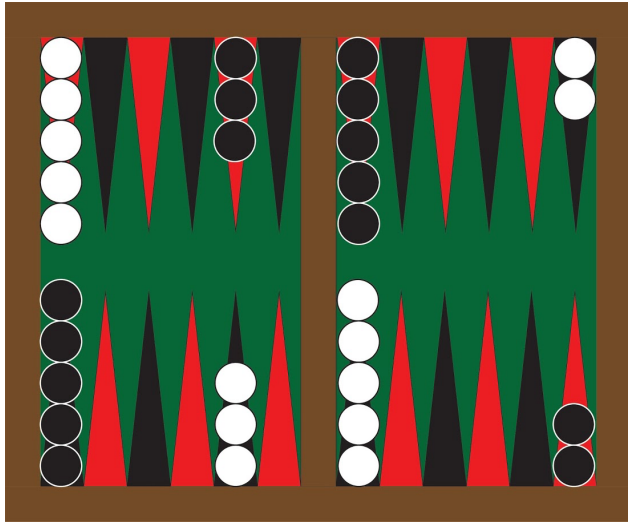


# Introduction and Basics of Markov Decision Process

**Wen Sun**

**CS 6789: Foundations of Reinforcement Learning**

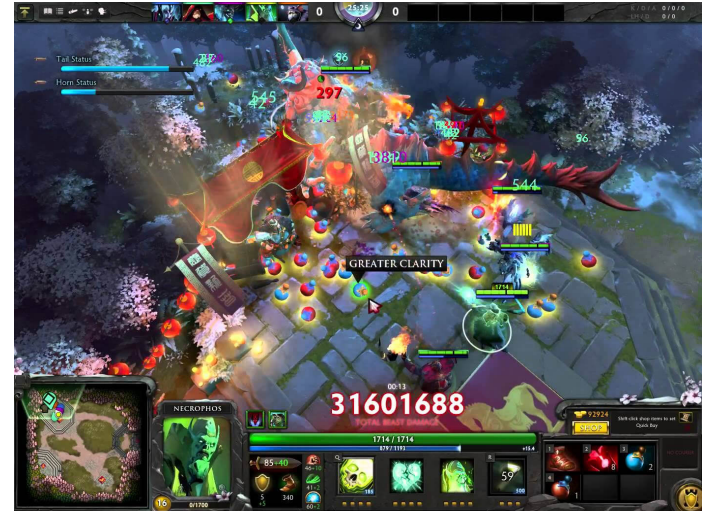
# Progress of RL in Practice



TD GAMMON [Tesauro 95]



[AlphaZero, Silver et.al, 17]

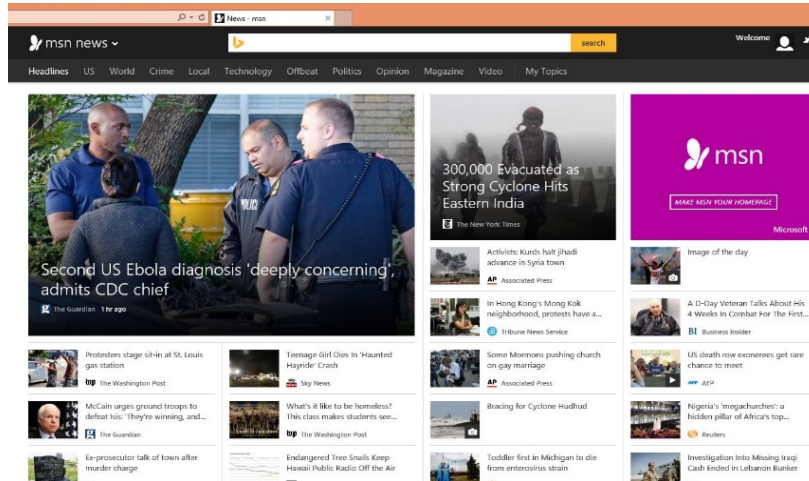


[OpenAI Five, 18]

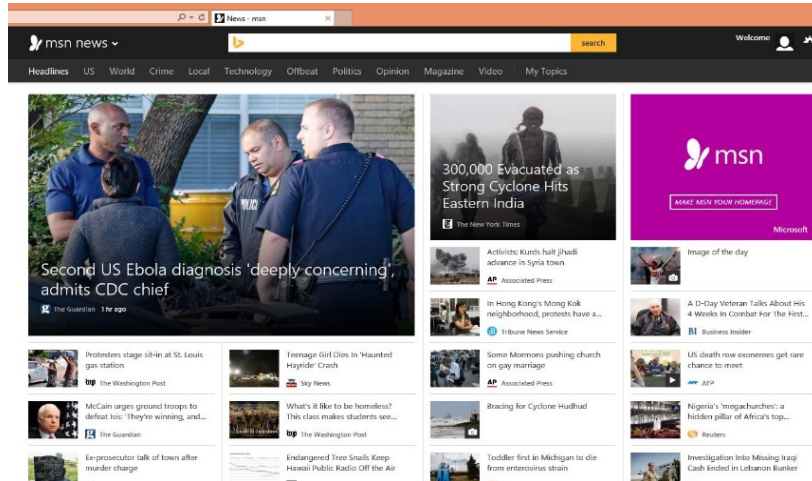
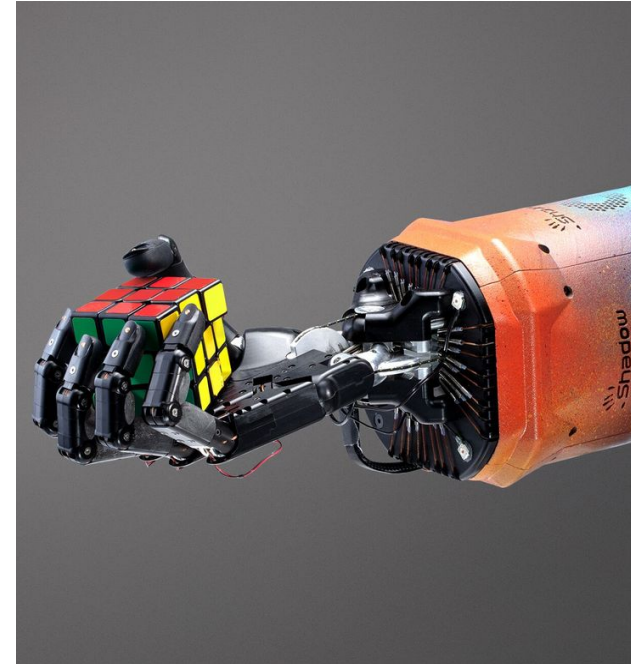
## RL in Real World:



# RL in Real World:



# RL in Real World:



# RL in Real World:

USER **this code is not working like i expect — how do i fix it?**

```
resultWorkerErr := make(chan error)
defer close(resultWorkerErr)
go func() {
    defer cancel()
    resultWorkerErr <- b.resultWorker(ctx)
}()

err := b.worker(ctx)
cancel()
if err == nil {
    return <-resultWorkerErr
}
return multierror.Append(err, <-resultWorkerErr)
```

CHATGPT It's difficult to say what's wrong with the code without more context. Can you provide more information about what the code is supposed to do and what isn't working as expected? Also, is this the entire code or just a part of it?

USER **just part of it — the error never surfaces. I think there is something wrong with the channel**

CHATGPT It's difficult to say without more information about what the code is supposed to do and what's happening when it's executed. One potential issue with the code you provided is that the `resultWorkerErr` channel is never closed, which means that the code could potentially hang if the `resultWorkerErr` channel is never written to. This

Training Language models  
using RL, e.g., chatGPT

# **This course focuses on RL Theory**

When and Why RL works!  
(Convergence, sample / computation complexity, etc)

## **Four main themes we will cover in this course:**

1. Fundamentals (MDPs and Optimal planning)
2. Exploration (sample complexity)
3. Policy Gradient (global convergence)
4. Learning Partially observable models



# Logistics

Four (HW0-HW3) assignments (total 55%), Course Project (40%), Reading (5%)

(HW0 10%, HW1-3 15% each)

HW0 is out today and due in one week

# Prerequisites (HW0)

## Deep understanding of Machine Learning, Optimization, Statistics

ML: sample complexity analysis for supervised learning (PAC)

Opt: Convex (linear) optimization, e.g., gradient descent for convex functions

Stats: basics of concentration (e.g., Hoeffding's), tricks such as union bound

# Prerequisites (HW0)

## Deep understanding of Machine Learning, Optimization, Statistics

ML: sample complexity analysis for supervised learning (PAC)

Opt: Convex (linear) optimization, e.g., gradient decent for convex functions

Stats: basics of concentration (e.g., Hoeffding's), tricks such as union bound

**Check out HW0 asap!**

## Course projects (40%)

- Team work: size 3
- Midterm report (5%), Final presentation (15%), and Final report (20%)
- Basics: **survey** of a set of similar RL theory papers. Reproduce analysis and provide a coherent story
- Advanced: **identify** extensions of existing RL papers, **formulate** theory questions, and **provide** proofs

# Course Notes:

# Reinforcement Learning Theory & Algorithms

- Book website: <https://rltheorybook.github.io/>
- Many lectures will correspond to chapters in Version 2.
- Reading assignment (5%) is from this book and additional notes
- Please let us know if you find typos/errors in the book!  
We appreciate it!

# Outline

1. Definition of infinite horizon discounted MDPs

2. Bellman Optimality

3. State-action distribution

# Supervised Learning

# Supervised Learning

Given i.i.d examples at training:



( ,cat )



( ,cat )






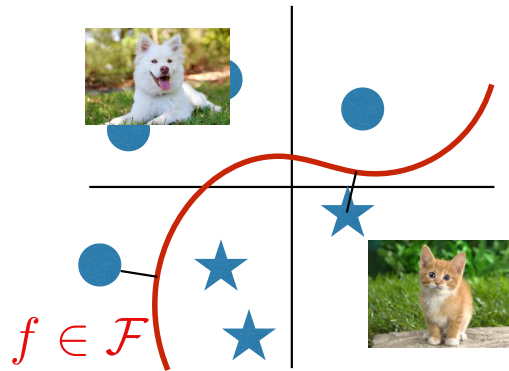
( ,dog )



# Supervised Learning

Given i.i.d examples at training:

(  ,cat ) (  ,cat ) (  ,dog )

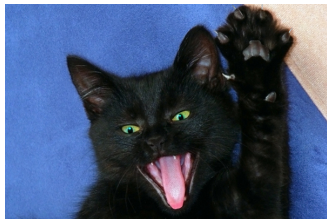


# Supervised Learning

Given i.i.d examples at training:



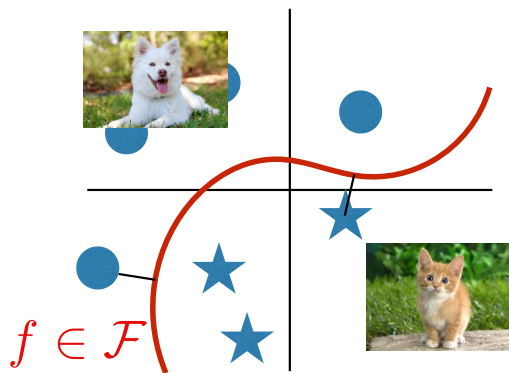
,cat



,cat



,dog



Passive:

Prediction



Data Distribution

AgentLinear

Selected Actions:

RIGHT

SPEED

Active:

Decisions



Data Distribution

AgentLinear

Selected Actions:

RIGHT

SPEED

Active:

Decisions



Data Distribution

AgentLinear

Selected Actions:

RIGHT

SPEED

Active:

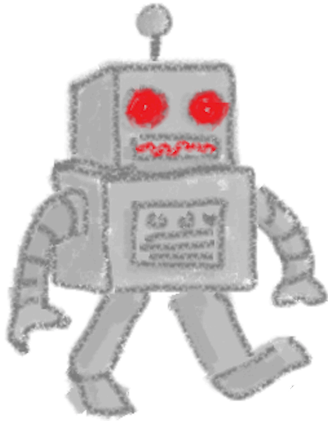
Decisions



Data Distribution

# Markov Decision Process

Learning  
Agent

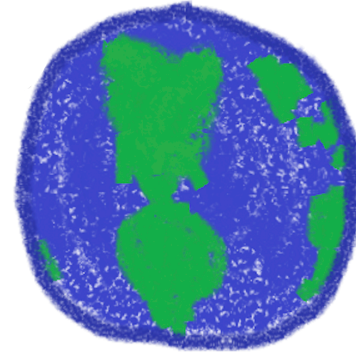


$$a \sim \pi(s)$$

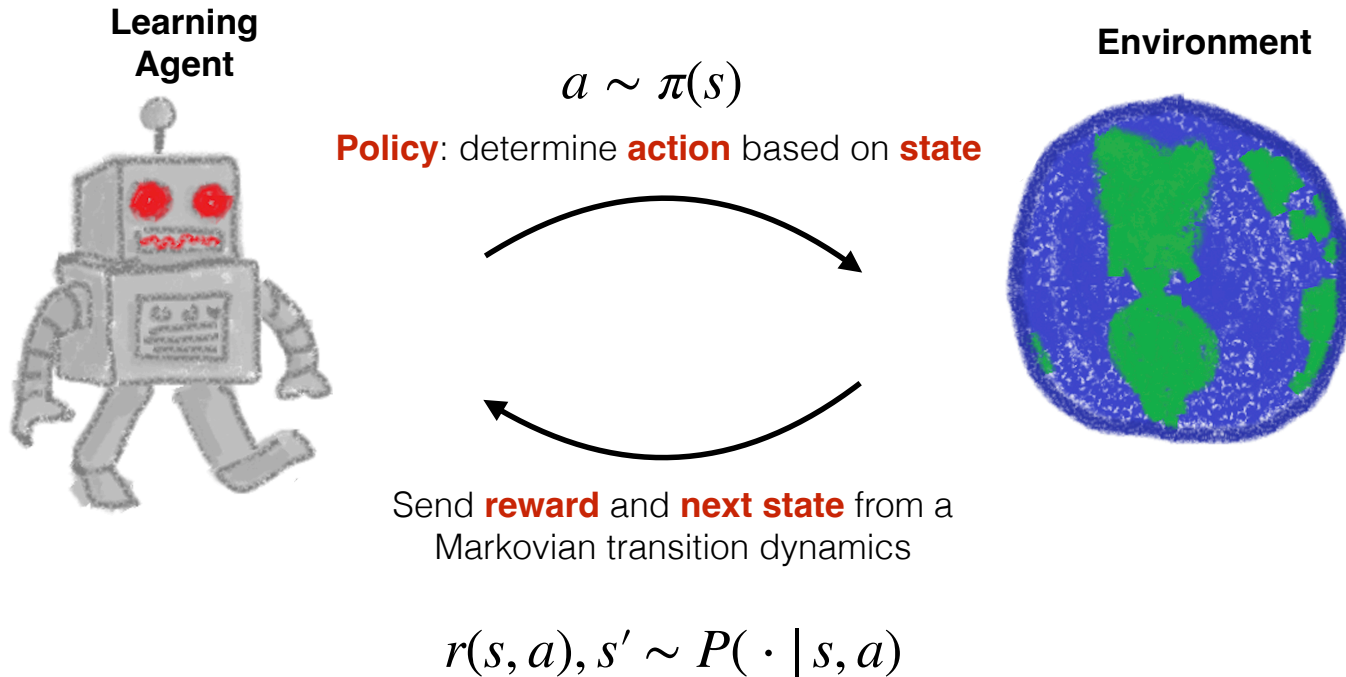
**Policy:** determine **action** based on **state**



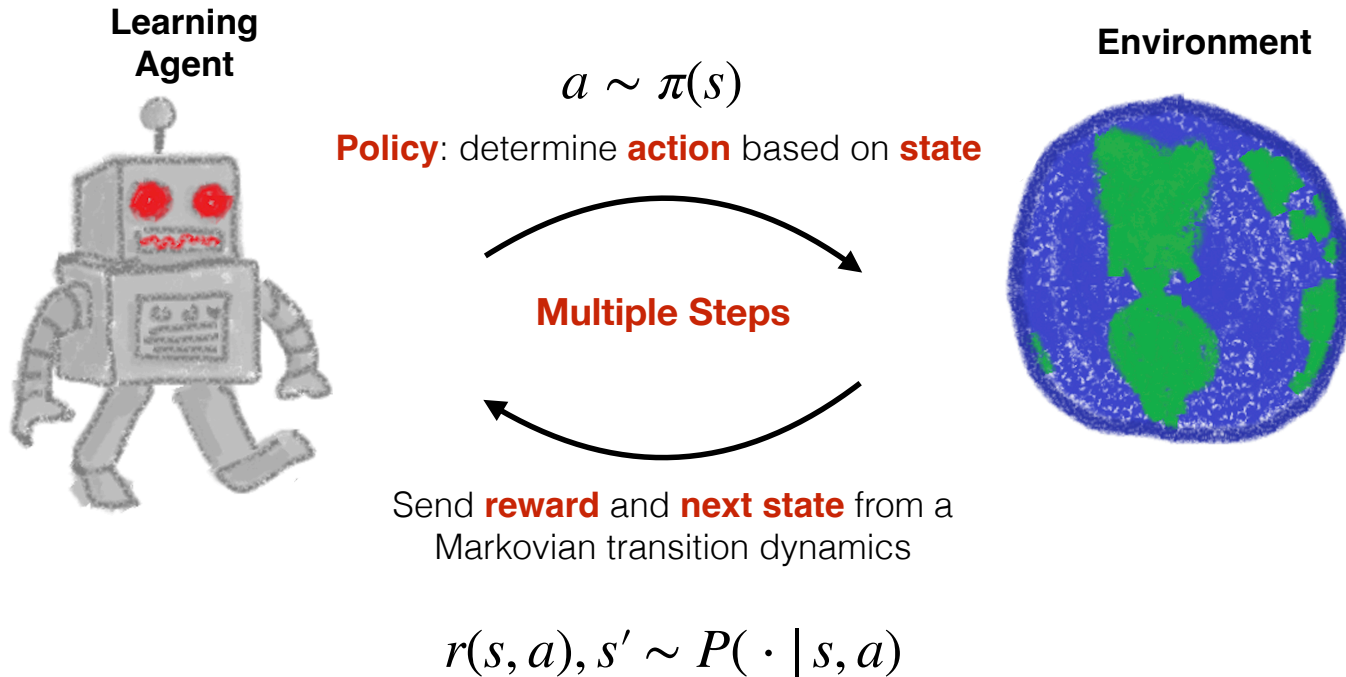
Environment



# Markov Decision Process

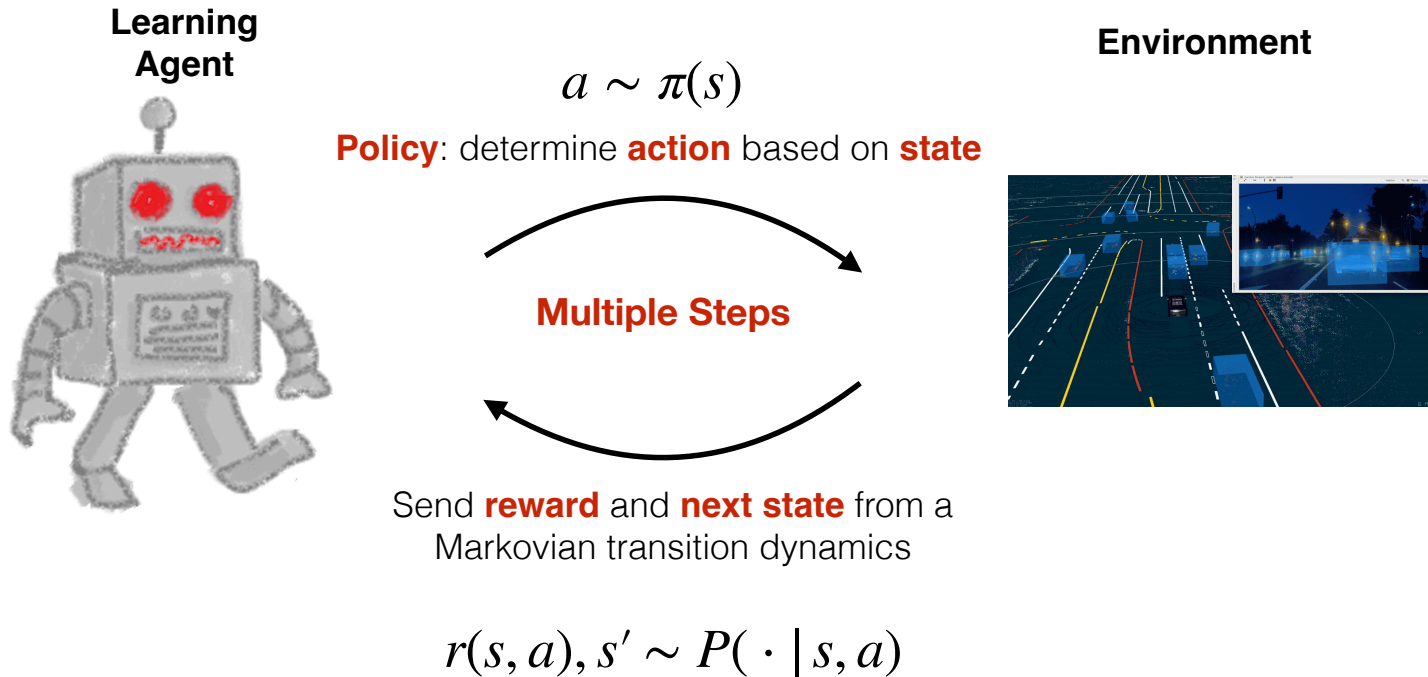


# Markov Decision Process

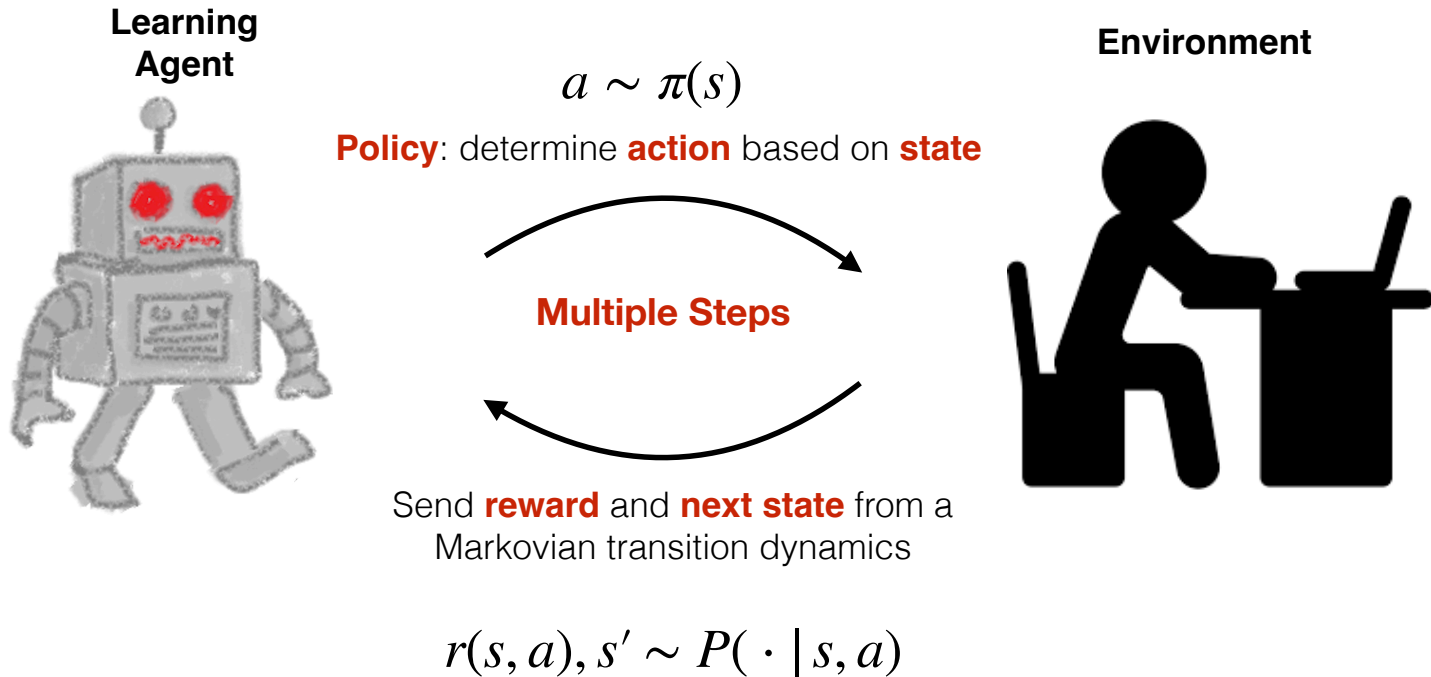




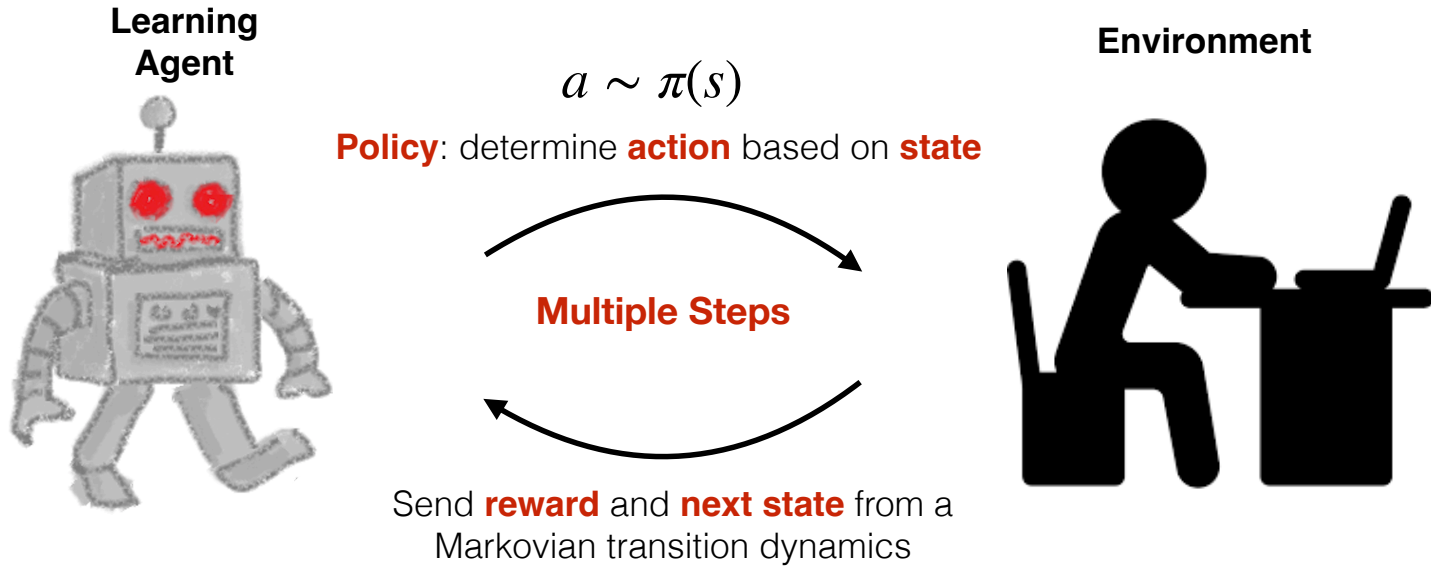
# Markov Decision Process



# Markov Decision Process



# Markov Decision Process





$$r(s, a), s' \sim P(\cdot | s, a)$$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r_0, s_1 \sim P(s_0, a_0), a_1 \sim \pi(s_1), r_1 \dots$$



	<b>Learn from Experience</b>	<b>Generalize</b>	<b>Interactive</b>	<b>Exploration</b>	<b>Credit assignment</b>
<b>Supervised Learning</b>					
<b>Reinforcement Learning</b>					

Table content based on slides from Emma Brunskill

	<b>Learn from Experience</b>	<b>Generalize</b>	<b>Interactive</b>	<b>Exploration</b>	<b>Credit assignment</b>
<b>Supervised Learning</b>					
<b>Reinforcement Learning</b>					

	<b>Learn from Experience</b>	<b>Generalize</b>	<b>Interactive</b>	<b>Exploration</b>	<b>Credit assignment</b>
<b>Supervised Learning</b>	✓	✓			
<b>Reinforcement Learning</b>	✓	✓			

	<b>Learn from Experience</b>	<b>Generalize</b>	<b>Interactive</b>	<b>Exploration</b>	<b>Credit assignment</b>
<b>Supervised Learning</b>	✓	✓			
<b>Reinforcement Learning</b>	✓	✓	✓		

	<b>Learn from Experience</b>	<b>Generalize</b>	<b>Interactive</b>	<b>Exploration</b>	<b>Credit assignment</b>
<b>Supervised Learning</b>	✓	✓			
<b>Reinforcement Learning</b>	✓	✓	✓	✓	



	<b>Learn from Experience</b>	<b>Generalize</b>	<b>Interactive</b>	<b>Exploration</b>	<b>Credit assignment</b>
<b>Supervised Learning</b>	✓	✓			
<b>Reinforcement Learning</b>	✓	✓	✓	✓	✓

# Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

*state space*  
*Action set*

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$s' \sim P(\cdot | s, a)$$

$$r(s, a)$$

$$s_0 \sim \mu_0$$

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

# Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

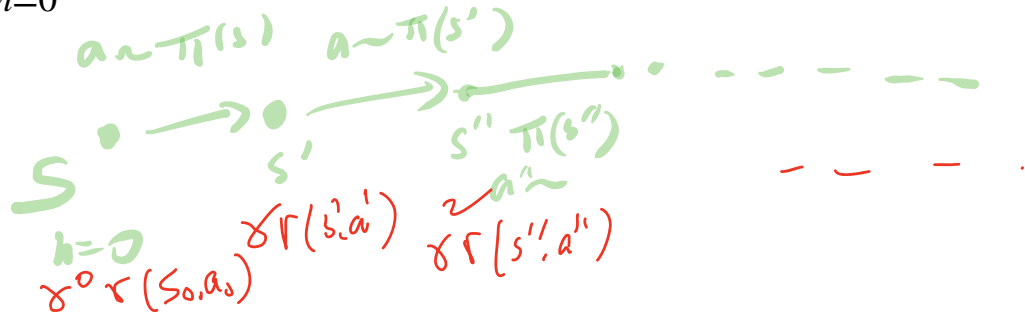
# Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

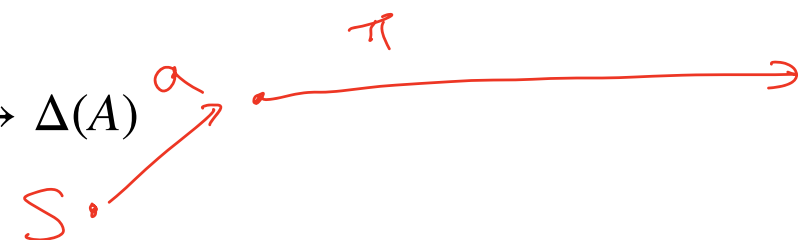
Value function  $V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$



# Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

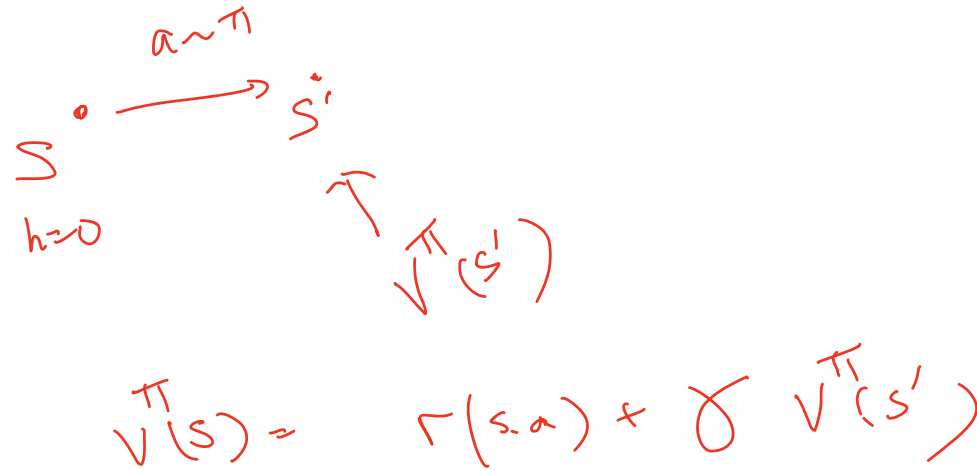
$$\text{Policy } \pi : S \mapsto \Delta(A)$$


$$\text{Value function } V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$\text{Q function } Q^\pi(s, a) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

# Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$



# Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

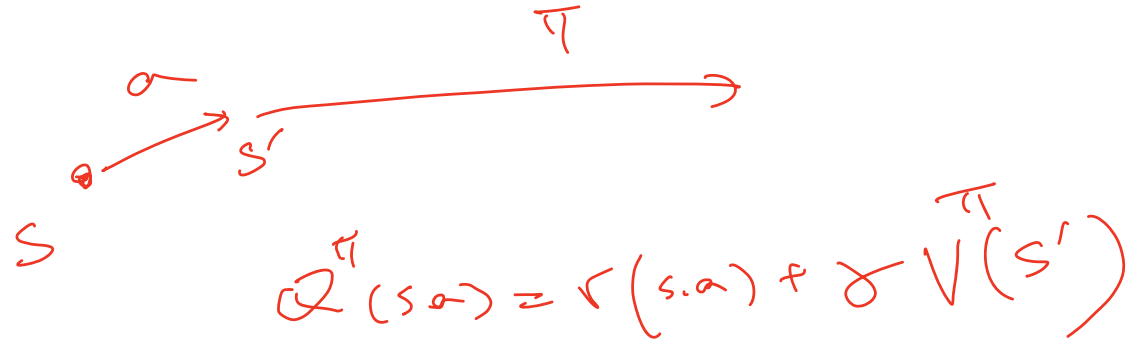
$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s') \right]$$

# Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s') \right]$$

$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$



$Q^\pi(s, a) = r(s, a) + \gamma V^\pi(s')$



# Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s') \right]$$

$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s')$$

# Outline



1. Definition of infinite horizon discounted MDPs

2. Bellman Optimality

3. State-action distribution

# Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^* : S \mapsto A, \text{ s.t.}, V^{\pi^*}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

# Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^* : S \mapsto A, \text{ s.t.}, V^{\pi^*}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

We denote  $V^* := V^{\pi^*}$ ,  $Q^* := Q^{\pi^*}$

# Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

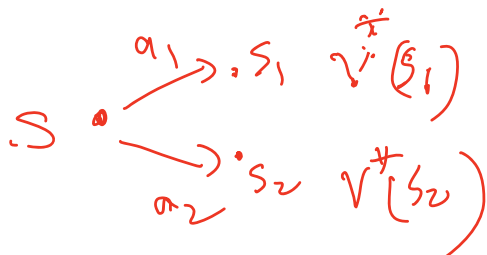
$$\pi^* : S \mapsto A, \text{ s.t., } V^{\pi^*}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

$$\text{We denote } V^* := V^{\pi^*}, Q^* := Q^{\pi^*}$$

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$



# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

$Q^*(s, a)$

Denote  $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$ , we will prove  $V^{\hat{\pi}}(s) = V^*(s), \forall s$

---

# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote  $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$ , we will prove  $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s')$$

↑  
Bellman-Eqn



# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote  $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$ , we will prove  $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s')$$

$$\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s')$$

$$Q^*(s, a) \quad \hat{\pi}(s) = \arg \max_a Q^*(s, a)$$

← repeat

Started

# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote  $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$ , we will prove  $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s')$$

$$\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s')$$

$$= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right]$$

Bellman - Eqn  
 $\hat{\pi}(s')$

# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote  $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$ , we will prove  $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \end{aligned}$$

# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote  $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$ , we will prove  $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s')$$

$$\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s')$$

$$= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right]$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right]$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[ r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right]$$

$\Delta$   
second-row

# Proof of Bellman Optimality

## Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote  $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$ , we will prove  $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s')$$

$$\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s')$$

$$= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right]$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right]$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[ r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right]$$

$$\leq \mathbb{E} \left[ r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots \right] = V^{\hat{\pi}}(s)$$

$V^*(s) \leq V^{\hat{\pi}}(s)$   
 $V^*(s) \geq V^{\hat{\pi}}(s)$   
 $\therefore \pi^*$  is optimal