



Markov Decision Processes (MDPs)

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 30th, 2007

©2005-2007 Carlos Guestrin

Thus far this semester

- Regression: $f: X \mapsto \mathbb{R}$
- Classification: $f: X \mapsto Y \in \{\text{hot, very hot, boiling hot}\}$
- Density estimation:

$$f: X \mapsto [0, 1]$$

$$\int_X f(x) dx = 1$$

Learning to act

- Reinforcement learning
- An agent
 - Makes sensor observations ^X
 - Must select action A (Y)
 - Receives rewards
 - positive for “good” states
 - negative for “bad” states



[Ng et al. '05]

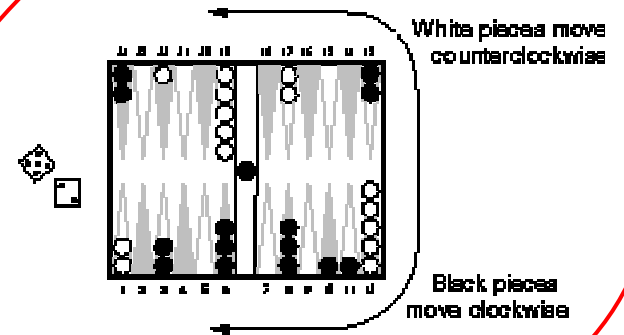
fully observable

X → true state of worlds

Learning to play backgammon

[Tesauro '95]

- Combines reinforcement learning with neural networks
- Played 300,000 games against itself
- Achieved grandmaster level!



Roadmap to learning about reinforcement learning

- When we learned about Bayes nets:

- First talked about formal framework:

- representation

- inference

- Then learning for BNs

- For reinforcement learning:

- Formal framework

- Markov decision processes

- Then learning

FreeCraft



peasant

footman

building

Real-time Strategy Game

Peasants collect resources and build

Footmen attack enemies

Buildings train peasants and footmen

States and actions

■ State space:

- Joint state \mathbf{x} of entire system

$$X = (P_1 = 0, P_2 = 7, E_c = 17, \dots)$$

■ Action space:

- Joint action $\mathbf{a} = \{a_1, \dots, a_n\}$ for all agents



enemy is "nature" or has a predefined strategy that doesn't or how you act. explicitly

States change over time

- Like an HMM, state changes over time
- Next state depends on current state and action selected

- e.g., action="build castle" likely to lead to a state where you have a castle

- Transition model:

- Dynamics of the entire system $P(\mathbf{x}'|\mathbf{x},\mathbf{a})$



$$P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \mathbf{a}^{(t-1)})$$

dynamics *you control*

Some states and actions are better than others

- Each state x is associated with a reward

- positive reward for successful attack
- negative for loss

- Reward function:

- Total reward $R(x)$

$R(x,a)$ ← depend on state & action

e.g., $R(x) = \begin{cases} 0 & \text{if playing} \\ 1 & \text{if won} \\ -1 & \text{if lost} \end{cases}$ at end

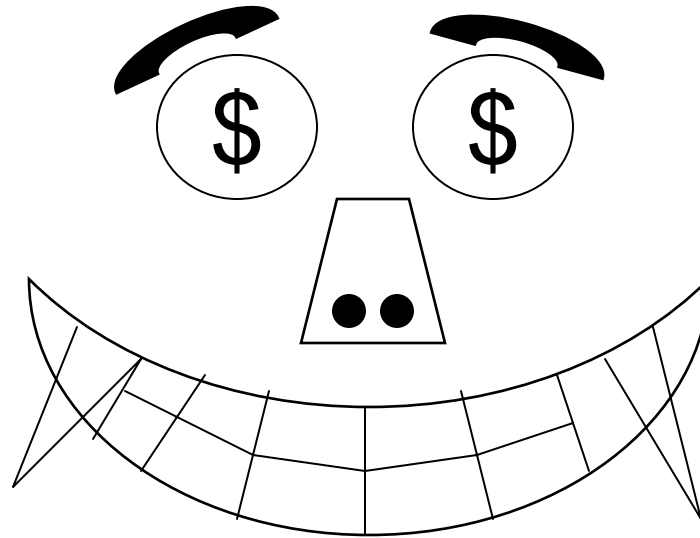


Discounted Rewards

An assistant professor gets paid, say, 20K per year.

How much, in total, will the A.P. earn in their life?

$$20 + 20 + 20 + 20 + 20 + \dots = \text{Infinity}$$



if inflation is
 $\frac{1}{8} = 1.03$

20k next
year
is worth only
¥20 this
year

What's wrong with this argument?

Discounted Rewards

“A reward (payment) in the future is not worth quite as much as a reward now.”

- Because of chance of obliteration
- Because of inflation

Example:


Being promised \$10,000 next year is worth only 90% as much as receiving \$10,000 right now.

Assuming payment n years in future is worth only $(0.9)^n$ of payment now, what is the AP's Future Discounted Sum of Rewards ?

$$20 + \gamma 20 + \gamma^2 20 + \dots$$

$$= \frac{20}{1-\gamma}$$

Discount Factors



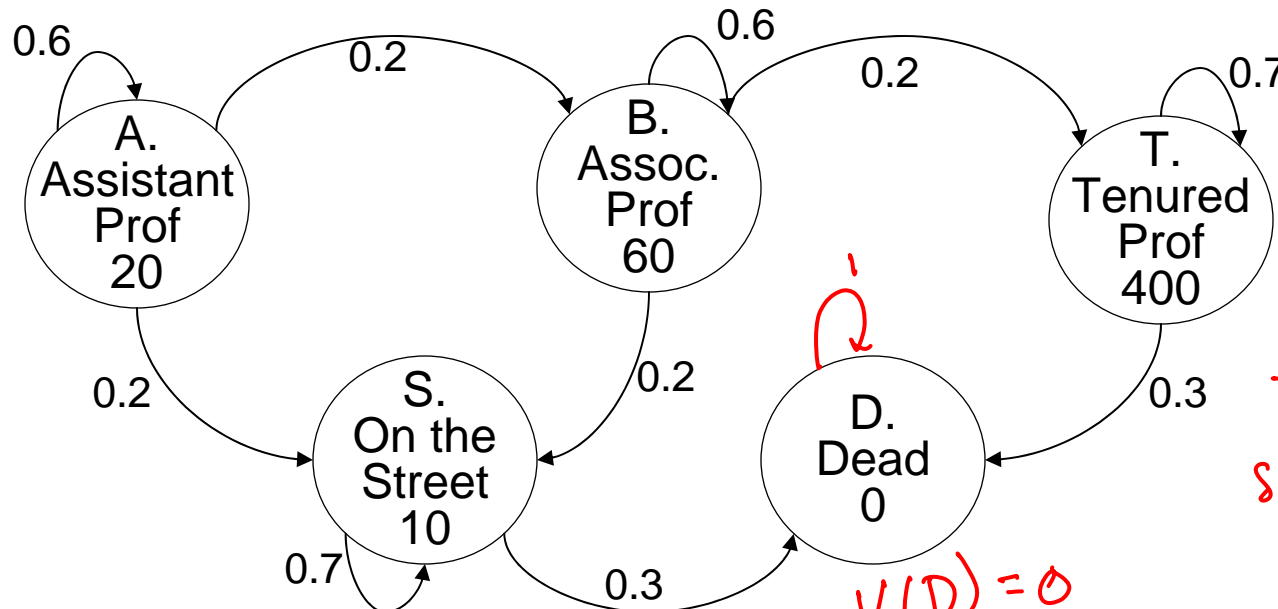
People in economics and probabilistic decision-making do this all the time.

The “Discounted sum of future rewards” using discount factor γ is

$$\begin{aligned} & (\text{reward now}) + \\ & \gamma (\text{reward in 1 time step}) + \\ & \gamma^2 (\text{reward in 2 time steps}) + \\ & \gamma^3 (\text{reward in 3 time steps}) + \\ & \quad : \\ & \quad : \quad (\text{infinite sum}) \end{aligned}$$

The Academic Life

Assume Discount Factor $\gamma = 0.9$



Define:

V_A = Expected discounted future rewards starting in state A

V_B = Expected discounted future rewards starting in state B

V_T = " " " " " " " T

V_S = " " " " " " " S

V_D = " " " " " " " D

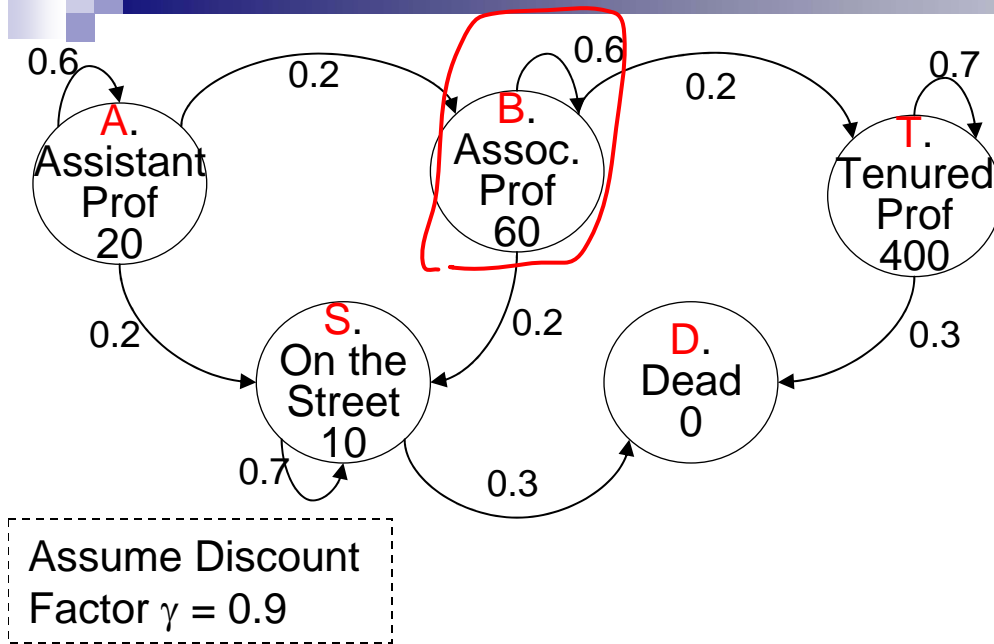
$$V(T) = 400 + 0.3\gamma V(D) + 0.7\gamma V(T)$$

$$V(T) = \frac{400}{1 - 0.7\gamma}$$

$V(D) = 0$

How do we compute V_A, V_B, V_T, V_S, V_D ?

Computing the Future Rewards of an Academic



$$\begin{aligned} V(B) = & 60 + \gamma \cdot 0.6 V(B) \\ & + \gamma \cdot 0.2 V(T) \\ & + \gamma \cdot 0.2 V(S) \end{aligned}$$

$$\begin{aligned} V(S) = & 10 + \gamma \cdot 0.7 V(S) \\ & + \gamma \cdot 0.3 V(D) \end{aligned}$$

Joint Decision Space

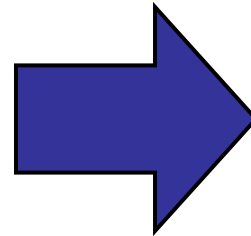
Markov Decision Process (MDP) Representation:

- State space:
 - Joint state \mathbf{x} of entire system
- Action space:
 - Joint action $\mathbf{a} = \{a_1, \dots, a_n\}$ for all agents
- Reward function:
 - Total reward $R(\mathbf{x}, \mathbf{a})$
 - sometimes reward can depend on action
- Transition model:
 - Dynamics of the entire system $P(\mathbf{x}' | \mathbf{x}, \mathbf{a})$



Policy

Policy: $\pi(\mathbf{x}) = \mathbf{a}$



At state \mathbf{x} ,
action \mathbf{a} for all
agents



$\pi(\mathbf{x}_0) =$ both peasants get wood



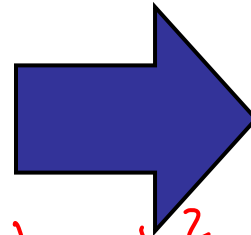
$\pi(\mathbf{x}_1) =$ one peasant builds
barrack, other gets gold



$\pi(\mathbf{x}_2) =$ peasants get gold,
footmen attack

Value of Policy

Value: $V_{\pi}(\mathbf{x})$



Expected long-term reward starting from \mathbf{x}

$$V(x_0) = E[R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) + \dots]$$

$$V_{\pi}(\mathbf{x}_0) = E[R(\mathbf{x}_0) + \gamma R(\mathbf{x}_1) + \gamma^2 R(\mathbf{x}_2) + \gamma^3 R(\mathbf{x}_3) + \gamma^4 R(\mathbf{x}_4) + \dots]$$

Future rewards discounted by $\gamma \in [0, 1)$

Start from \mathbf{x}_0



$R(\mathbf{x}_0)$

$\pi(\mathbf{x}_0)$



$R(\mathbf{x}_1)$

$\pi(\mathbf{x}_1)$



$R(\mathbf{x}_2)$

$\pi(\mathbf{x}_2)$



$R(\mathbf{x}_3)$

$\pi(\mathbf{x}_3)$



$R(\mathbf{x}_4)$

17

Computing the value of a policy

$$V_{\pi}(\mathbf{x}_0) = \mathbf{E}_{\pi}[R(\mathbf{x}_0) + \gamma R(\mathbf{x}_1) + \gamma^2 R(\mathbf{x}_2) + \gamma^3 R(\mathbf{x}_3) + \gamma^4 R(\mathbf{x}_4) + \dots]$$

- Discounted value of a state:

- value of starting from x_0 and continuing with policy π from then on

$$\begin{aligned} V_{\pi}(x_0) &= E_{\pi}[R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) + \gamma^3 R(x_3) + \dots] \\ &= E_{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(x_t)\right] \end{aligned}$$

- A recursion!

$$\begin{aligned} V_{\pi}(x_0) &= E_{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(x_t)\right] = E_{\pi}\left[R(x_0) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} R(x_t)\right] \\ &= E_{\pi}[R(x_0)] + \gamma E_{\pi}\left[\sum_{t=1}^{\infty} \gamma^{t-1} R(x_t)\right] \\ V_{\pi}(x_0) &= R(x_0) + \gamma E_{x_1}[V_{\pi}(x_1)] = R(x_0) + \gamma \sum_{x_1} P(x_1 | x_0, a = \pi(x_0)) V_{\pi}(x_1) \end{aligned}$$

Computing the value of a policy 1 – the matrix inversion approach

n possible states

$$V_{\pi}(x) = R(x) + \gamma \sum_{x'} P(x' | x, a = \pi(x)) V_{\pi}(x')$$

- Solve by simple matrix inversion:

$$V_{\pi} = R + \gamma P_{\pi} V_{\pi}$$

$$(I - \gamma P_{\pi}) V_{\pi} = R$$

$$V_{\pi} = (I - \gamma P_{\pi})^{-1} R$$

simple matrix inversion

Handwritten definitions and matrix structures:

- $V_{\pi} = \begin{pmatrix} V_{\pi}(x_1) \\ \vdots \\ V_{\pi}(x_n) \end{pmatrix}$ (size $n \times 1$)
- $R = \begin{pmatrix} R(x_1) \\ \vdots \\ R(x_n) \end{pmatrix}$ (size $n \times 1$)
- $P_{\pi} = \begin{pmatrix} P(x_1 | x_1, a = \pi(x_1)) & \dots & P(x_n | x_1, a = \pi(x_1)) \\ \vdots & \ddots & \vdots \\ P(x_1 | x_n, a = \pi(x_n)) & \dots & P(x_n | x_n, a = \pi(x_n)) \end{pmatrix}$ (size $n \times n$)
- Annotations:
 - $x' = x^{(t)}$ (next time step)
 - $x = x^{(t-1)} = i$ (current time step)
 - Matrix element: $P(x^{(t)} | x^{(t-1)} = i, a = \pi(x^{(t-1)} = i))$

Computing the value of a policy 2 – iteratively

$$V_{\pi} = R + \gamma P_{\pi} V_{\pi}$$

$$V_{\pi}(x) = R(x) + \gamma \sum_{x'} P(x' | x, a = \pi(x)) V_{\pi}(x')$$

- If you have 1000,000 states, inverting a 1000,000x1000,000 matrix is hard!
- Can solve using a simple convergent iterative approach: (a.k.a. dynamic programming)
 - Start with some guess V_0
 - Iteratively say:
 - $V_{t+1} = R + \gamma P_{\pi} V_t$
 - Stop when $\|V_{t+1} - V_t\|_{\infty} \leq \epsilon$
 - means that $\|V_{\pi} - V_{t+1}\|_{\infty} \leq \epsilon/(1-\gamma)$

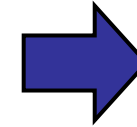
good guess: $V_0 = R$

converges to $V_{\pi} = (I - \gamma P_{\pi})^{-1} R$
because $R + \gamma P_{\pi} V_{\pi}$
is a contraction

But we want to learn a Policy

- So far, told you how good a policy is...
- But how can we choose the best policy???
- Suppose there was only one time step:
 - world is about to end!!!
 - select action that maximizes reward!

Policy: $\pi(\mathbf{x}) = \mathbf{a}$



At state \mathbf{x} , action \mathbf{a} for all agents



$\pi(\mathbf{x}_0)$ = both peasants get wood



$\pi(\mathbf{x}_1)$ = one peasant builds barrack, other gets gold



$\pi(\mathbf{x}_2)$ = peasants get gold, footmen attack

$$\underset{\substack{\text{one} \\ \text{timestep} \\ \text{to go}}}{V(\mathbf{x}_0)} = \underset{\mathbf{a}}{\operatorname{arg\,max}} R(\mathbf{x}_0) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{a}, \mathbf{x}_0) R(\mathbf{x}')$$

Another recursion!

- Two time steps: address tradeoff
 - good reward now
 - better reward in the future

$$V(x_0) = \max_a R(x_0) + \gamma \sum_{x'} P(x'|x_0, a) V(x')$$

going to Vegas & spending
all my savings
 $V.$

saving for retirement

include a bunch of max's

Unrolling the recursion

- Choose actions that lead to best value in the long run

↙ optimal □ Optimal value policy achieves optimal value V^*

$$\underline{V^*(x_0)} = \max_{a_0} R(x_0, a_0) + \gamma E_{a_0} [\underbrace{\max_{a_1} R(x_1, a_1) + \gamma^2 E_{a_1} [\max_{a_2} R(x_2, a_2) + \dots]}_{V^*(x_1)}]$$

$$V^*(x_0) = \max_{a_0} R(x_0, a_0) + \gamma \sum_{x'} P(x' | x_0, a) V^*(x')$$

Bellman equation [Bellman 60]

- Evaluating policy π :

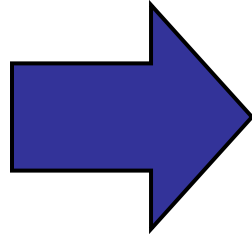
$$V_{\pi}(x) = R(x) + \gamma \sum_{x'} P(x' | x, a = \pi(x)) V_{\pi}(x')$$

- Computing the optimal value V^* - Bellman equation

$$V^*(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V^*(\mathbf{x}')$$

Optimal Long-term Plan

Optimal value
function $V^*(\mathbf{x})$



Optimal Policy: $\pi^*(\mathbf{x})$

$$Q^*(\mathbf{x}, \mathbf{a}) = R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V^*(\mathbf{x}')$$

Optimal policy:

$$\pi^*(\mathbf{x}) = \arg \max_{\mathbf{a}} Q^*(\mathbf{x}, \mathbf{a})$$

$$\pi^*(x) = \arg \max_a R(x,a) + \gamma \sum_{x'} P(x'|x,a) V^*(x')$$

Interesting fact – Unique value

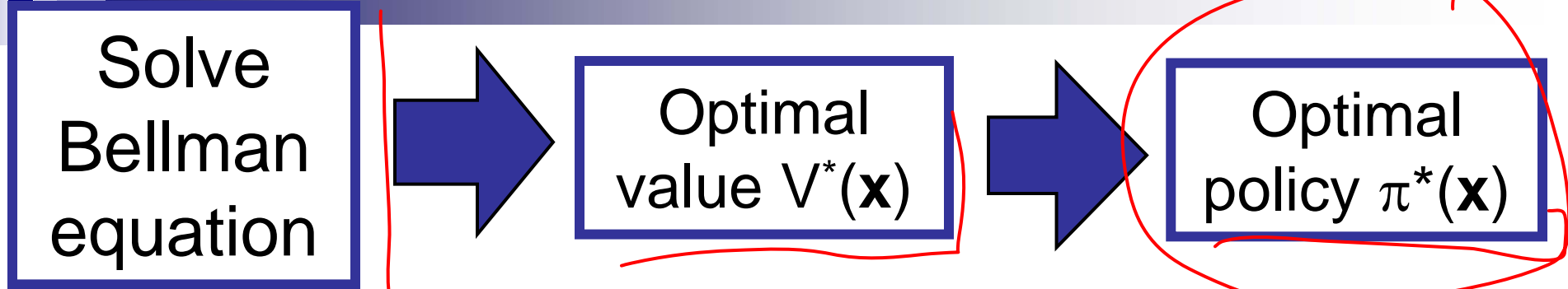
$$\underline{V^*(\mathbf{x})} = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V^*(\mathbf{x}')$$

- *Slightly surprising fact.* There is only one V^* that solves Bellman equation!
 - there may be many optimal policies that achieve V^*
- *Surprising fact.* optimal policies are good everywhere!!!

$$V^*(x) \doteq \underline{V_{\pi^*}(x)} \geq V_{\pi}(x), \quad \forall x, \quad \forall \pi$$

no worse *everywhere* *any other policy*

Solving an MDP



$$V^*(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V^*(\mathbf{x}')$$

Bellman equation is non-linear!!!

Many algorithms solve the Bellman equations:

- Policy iteration [Howard '60, Bellman '57]
- Value iteration [Bellman '57]
- Linear programming [Manne '60]
- ...

Value iteration (a.k.a. dynamic programming) – the simplest of all

$$V^*(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V^*(\mathbf{x}')$$

- Start with some guess V_0 $= R$
- Iteratively say:
 - $V_{t+1}(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V_t(\mathbf{x}')$

- Stop when $\|V_{t+1} - V_t\|_{\infty} \leq \varepsilon$
 - means that $\|V^* - V_{t+1}\|_{\infty} \leq \varepsilon / (1 - \gamma)$

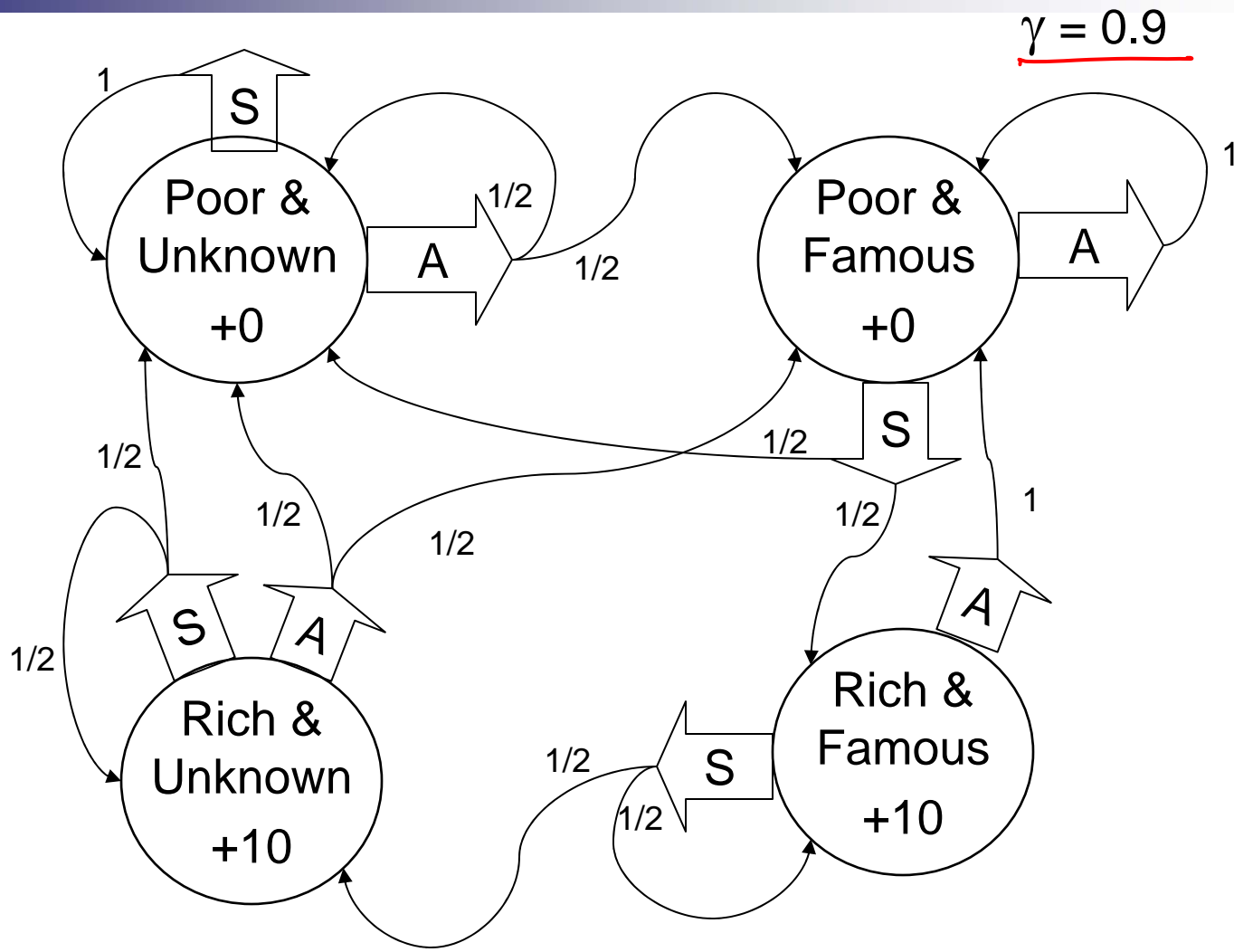
$\Rightarrow \Pi_{t+1} = \text{Greedy Policy } V_{t+1}$

$$\|V^* - V_{\Pi_{t+1}}\|_{\infty} \leq \frac{\varepsilon \gamma}{(1 - \gamma)^2}$$

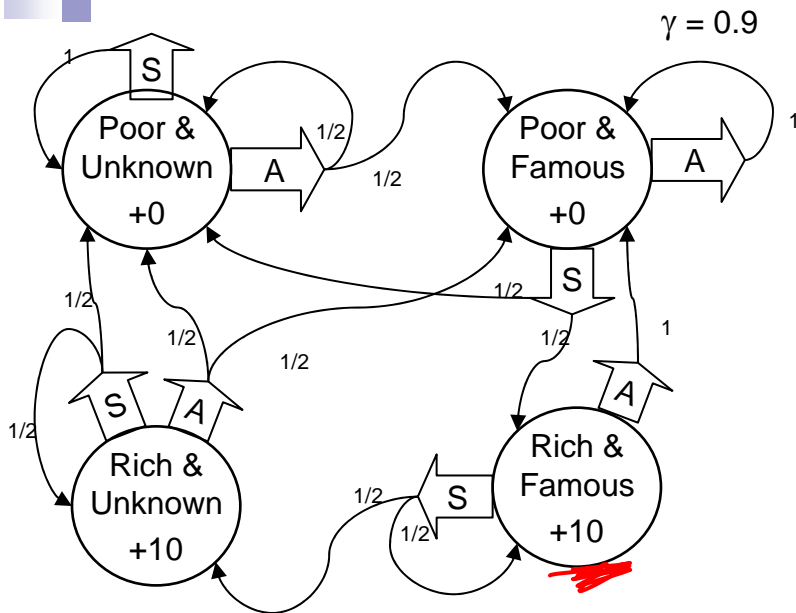
converges to V^*
can get Π^*

A simple example

You run a startup company.
In every state you must choose between Saving money or Advertising.



Let's compute $V_t(x)$ for our example



| t | $V_t(\text{PU})$ | $V_t(\text{PF})$ | $V_t(\text{RU})$ | <u>$V_t(\text{RF})$</u> |
|---|------------------|------------------|------------------|------------------------------------|
| 1 | 0 | 0 | 10 | 10 |
| 2 | | | | 19 |
| 3 | | | | |
| 4 | | | | |
| 5 | | | | |
| 6 | | | | |

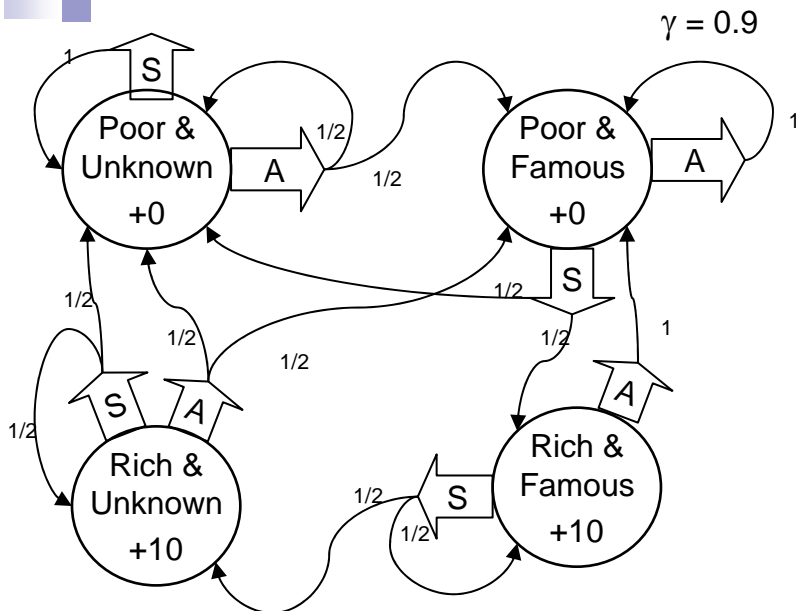
$$V_{t+1}(\text{RF}) = 10 + \gamma$$

$a=A \begin{cases} 1 \cdot V_t(\text{PF}) \\ 0.5 V_t(\text{RF}) + 0.5 V_t(\text{RU}) \end{cases}$

$\max_{a=S}$

$$V_{t+1}(\mathbf{x}) = \max_a R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V_t(\mathbf{x}')$$

Let's compute $V_t(\mathbf{x})$ for our example



| t | $V_t(\text{PU})$ | $V_t(\text{PF})$ | $V_t(\text{RU})$ | $V_t(\text{RF})$ |
|---|------------------|------------------|------------------|------------------|
| 1 | 0 | 0 | 10 | 10 |
| 2 | 0 | 4.5 | 14.5 | 19 |
| 3 | 2.03 | 6.53 | 25.08 | 18.55 |
| 4 | 3.852 | 12.20 | 29.63 | 19.26 |
| 5 | 7.22 | 15.07 | 32.00 | 20.40 |
| 6 | 10.03 | 17.65 | 33.58 | 22.43 |

$$V_{t+1}(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V_t(\mathbf{x}')$$

Policy iteration – Another approach for computing π^*

- Start with some guess for a policy π_0
- Iteratively say:

- evaluate policy:

$$V_t(\mathbf{x}) = R(\mathbf{x}, \mathbf{a} = \pi_t(\mathbf{x})) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a} = \pi_t(\mathbf{x})) V_t(\mathbf{x}')$$

- improve policy:

$$V_t = (I - \gamma P_{\pi_t})^{-1} R$$
$$\pi_{t+1}(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V_t(\mathbf{x}')$$

- Stop when
 - policy stops changing
 - usually happens in about 10 iterations
 - or $\|V_{t+1} - V_t\|_{\infty} \leq \varepsilon$
 - means that $\|V^* - V_{t+1}\|_{\infty} \leq \varepsilon / (1 - \gamma)$

open problem:
PI converges in
polynomial time?

Policy Iteration & Value Iteration: Which is best ???

It depends.

Lots of actions? Choose **Policy Iteration**

Already got a fair policy? **Policy Iteration**

Few actions, acyclic? **Value Iteration**

Best of Both Worlds:

Modified Policy Iteration [Puterman]

...a simple mix of value iteration and policy iteration

3rd Approach

Linear Programming

LP Solution to MDP

[Manne '60]

Value computed by linear programming:

$$\begin{aligned} &\text{minimize: } \sum_{\mathbf{x}} V(\mathbf{x}) \\ &\text{subject to: } \begin{cases} V(\mathbf{x}) \geq R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V(\mathbf{x}') \\ \forall \mathbf{x}, \mathbf{a} \end{cases} \end{aligned}$$

- One variable $V(\mathbf{x})$ for each state
- One constraint for each state \mathbf{x} and action \mathbf{a}
- Polynomial time solution

What you need to know



- What's a Markov decision process
 - state, actions, transitions, rewards
 - a policy
 - value function for a policy
 - computing V_π
- Optimal value function and optimal policy
 - Bellman equation
- Solving Bellman equation
 - with value iteration, policy iteration and linear programming

Acknowledgment



- This lecture contains some material from Andrew Moore's excellent collection of ML tutorials:

- <http://www.cs.cmu.edu/~awm/tutorials>



Reinforcement Learning

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 30th, 2007

©2005-2007 Carlos Guestrin

The Reinforcement Learning task



World: You are in state 34.
Your immediate reward is 3. You have possible 3 actions.

Robot: I'll take action 2.

World: You are in state 77.
Your immediate reward is -7. You have possible 2 actions.

Robot: I'll take action 1.

World: You're in state 34 (again).
Your immediate reward is 3. You have possible 3 actions.

Formalizing the (online) reinforcement learning problem

- Given a set of states \mathbf{X} and actions \mathbf{A}
 - in some versions of the problem size of \mathbf{X} and \mathbf{A} unknown
- Interact with world at each time step t :
 - world gives state \mathbf{x}_t and reward r_t
 - you give next action \mathbf{a}_t
- **Goal:** (quickly) learn policy that (approximately) maximizes long-term expected discounted reward

The “Credit Assignment” Problem



| | | |
|------------------|-------------|------------|
| I'm in state 43, | reward = 0, | action = 2 |
| “ “ “ 39, | “ = 0, | “ = 4 |
| “ “ “ 22, | “ = 0, | “ = 1 |
| “ “ “ 21, | “ = 0, | “ = 1 |
| “ “ “ 21, | “ = 0, | “ = 1 |
| “ “ “ 13, | “ = 0, | “ = 2 |
| “ “ “ 54, | “ = 0, | “ = 2 |
| “ “ “ 26, | “ = 100, | |

Yippee! I got to a state with a big reward! But which of my actions along the way actually helped me get there??

This is the **Credit Assignment** problem.

Exploration-Exploitation tradeoff

- You have visited part of the state space and found a reward of 100
 - is this the best I can hope for???
- **Exploitation:** should I stick with what I know and find a good policy w.r.t. this knowledge?
 - at the risk of missing out on some large reward somewhere
- **Exploration:** should I look for a region with more reward?
 - at the risk of wasting my time or collecting a lot of negative reward

Two main reinforcement learning approaches

■ Model-based approaches:

- explore environment → learn model ($P(\mathbf{x}'|\mathbf{x},\mathbf{a})$ and $R(\mathbf{x},\mathbf{a})$) (almost) everywhere
- use model to plan policy, MDP-style
- approach leads to strongest theoretical results
- works quite well in practice when state space is manageable

■ Model-free approach:

- don't learn a model → learn value function or policy directly
- leads to weaker theoretical results
- often works well when state space is large