

# Markov Decision Processes (MDPs)

Machine Learning – CSE446

Carlos Guestrin

University of Washington

June 3, 2013

©Carlos Guestrin 2005-2013

1



# Reinforcement Learning

training by feedback

©Carlos Guestrin 2005-2013

2

1

# Learning to act

- Reinforcement learning
- An agent
  - Makes sensor observations
  - Must select action
  - Receives rewards
    - positive for “good” states
    - negative for “bad” states



[Ng et al. '05]

©Carlos Guestrin 2005-2013

3

# Markov Decision Process (MDP) Representation

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

have Gold  
don't have it  
went Mining



©Carlos Guestrin 2005-2013

4

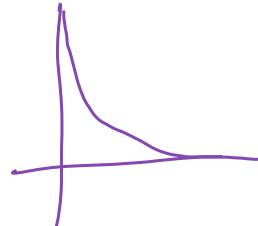
# Discount Factors

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

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

The “Discounted sum of future rewards” using discount factor  $\gamma$  is

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



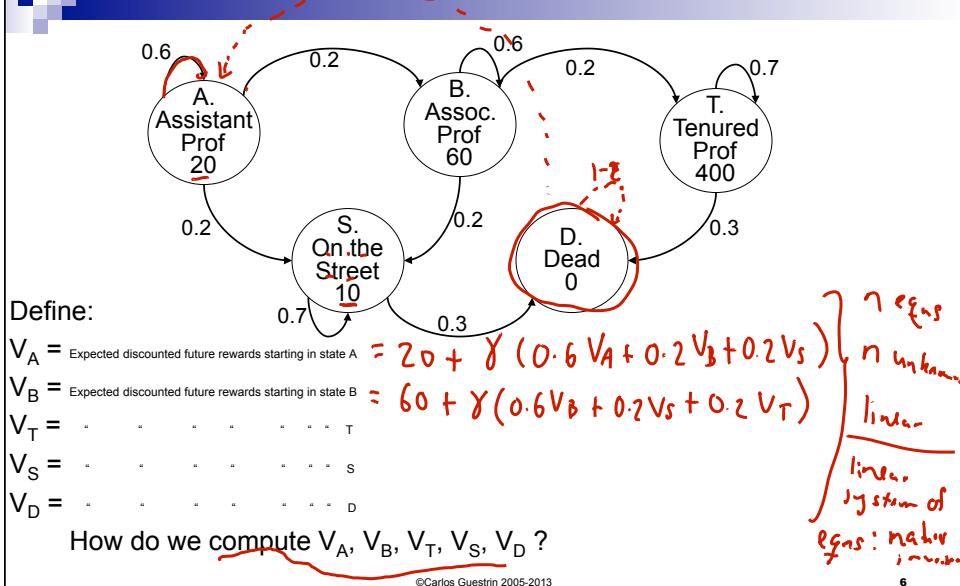
©Carlos Guestrin 2005-2013

5

## The Academic Life

$n$  States

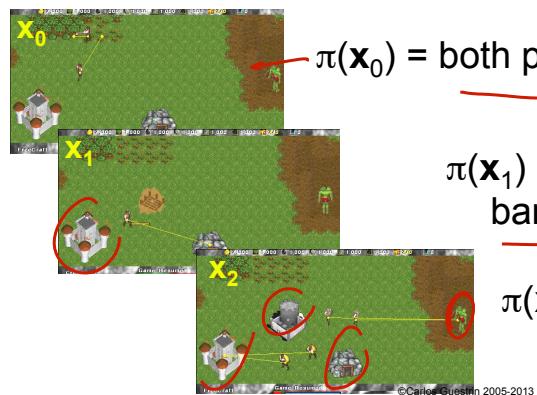
Assume Discount Factor  $\gamma = 0.9$



## Policy

Policy:  $\pi(x) = a$

At state  $x$ ,  
action  $a$  for all  
agents



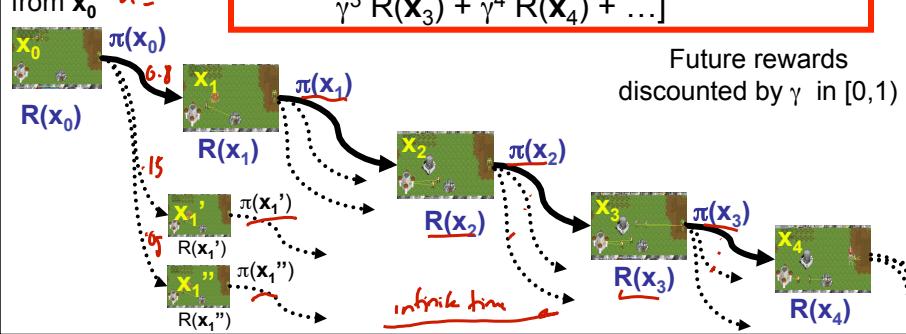
7

## Value of Policy

Value:  $V_\pi(x)$

Expected long-term reward starting from  $x$

$$V_\pi(x_0) = \mathbb{E}_\pi[R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) + \gamma^3 R(x_3) + \gamma^4 R(x_4) + \dots]$$



# Computing the value of a policy

$$V_\pi(x_0) = E_\pi[R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) + \gamma^3 R(x_3) + \gamma^4 R(x_4) + \dots]$$

- Discounted value of a state:

- value of starting from  $x_0$  and continuing with policy  $\pi$  from then on

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

- A recursion!

$$\begin{aligned} V_\pi(x_0) &= E_\pi\left[\sum_{t=0}^{\infty} \gamma^t R(x_t)\right] = E_\pi\left[R(x_0) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} R(x_t)\right] \\ &= R(x_0) + \gamma E_\pi\left[R(x_1) + \sum_{t=2}^{\infty} \gamma^{t-1} R(x_t)\right] \\ &= R(x_0) + \gamma \underbrace{E_\pi\left[R(x_1) + \sum_{t=2}^{\infty} \gamma^{t-1} R(x_t)\right]}_{V_\pi(x_1)} \quad \left. \begin{array}{l} \text{n linear} \\ \text{equations} \end{array} \right\} \quad \left. \begin{array}{l} \text{matrix} \\ \text{inversion} \end{array} \right\} \\ &= R(x_0) + \gamma \sum_{x'} P(x' | x_0, \pi(x_0)) V_\pi(x') \quad \left. \begin{array}{l} \text{n unknowns} \\ \text{is your friend!} \end{array} \right\} \end{aligned}$$

©Carlos Guestrin 2005-2013

9

## Simple approach for computing the value of a policy: Iteratively

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

✓ from previous slide  
recursion

- Can solve using a simple convergent iterative approach:

(a.k.a. dynamic programming)  $V^*(x) = R(x)$

- Start with some guess  $V^0$

✓ estimate from previous iteration

- Iteratively say:

$$V_\pi^{t+1}(x) \leftarrow R(x) + \gamma \sum_{x'} P(x' | x, a = \pi(x)) V_\pi^t(x')$$

- Stop when  $\|V_\pi^{t+1} - V_\pi^t\|_\infty < \varepsilon$

$$\text{■ means that } \|V_\pi^{t+1} - V_\pi^t\|_\infty < \varepsilon / (1-\gamma)$$

↑ estimate is close to true value

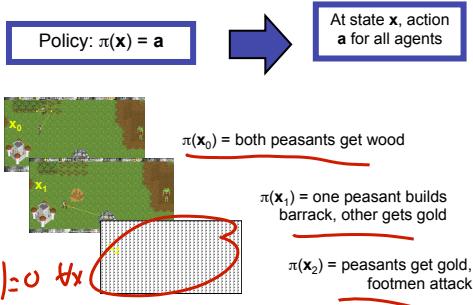
©Carlos Guestrin 2005-2013

10

## 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!
- $$V(x) = \max_a R(x, a)$$



©Carlos Guestrin 2005-2013

11

## Unrolling the recursion

n states

- Choose actions that lead to best value in the long run
  - Optimal value policy achieves optimal value  $V^*$

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

$$V^*(x_0) = \max_{a_0} R(x_0, a_0) + \gamma E_{a_0} [V^*(x_1)]$$

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

↑  
not linear

*n unknowns  
 $V(x)$*

*n equations  
but not linear,  
matrix inversion  
is not your friend!!*

©Carlos Guestrin 2005-2013

12

## Bellman equation

- Evaluating policy  $\pi$ :

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

- Computing the optimal value  $V^*$  - Bellman equation

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

if I know  $V^*$ :

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

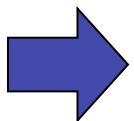
↖ act optimally

©Carlos Guestrin 2005-2013

13

## Optimal Long-term Plan

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



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

## Optimal policy:

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

©Carlos Guestrin 2005-2013

14

## Interesting fact – Unique value

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

- *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!!!

$$V_{\pi^*}(\mathbf{x}) \geq V_{\pi}(\mathbf{x}), \quad \forall \mathbf{x}, \quad \forall \pi$$

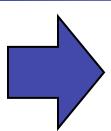
↗  
 if  $\pi^*$  is  
 an optimal  
 policy

©Carlos Guestrin 2005-2013

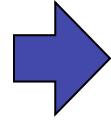
15

## Solving an MDP

Solve  
Bellman  
equation



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



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

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

**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]
- ...

©Carlos Guestrin 2005-2013

16

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

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

- Start with some guess  $V^0(x) = \max_a R(x, a)$
  - Iteratively say:
    - $V^{t+1}(x) \leftarrow \max_a R(x, a) + \gamma \sum_{x'} P(x' | x, a) V^t(x')$
  - Stop when  $\|V^{t+1} - V^t\|_\infty < \epsilon$ 
    - means that  $\|V - V^{t+1}\|_\infty < \epsilon/(1-\gamma)$
- ↓ value from previous iteration  
 monotonically increasing if  
 $\forall x \exists a : R(x, a) \geq \epsilon$

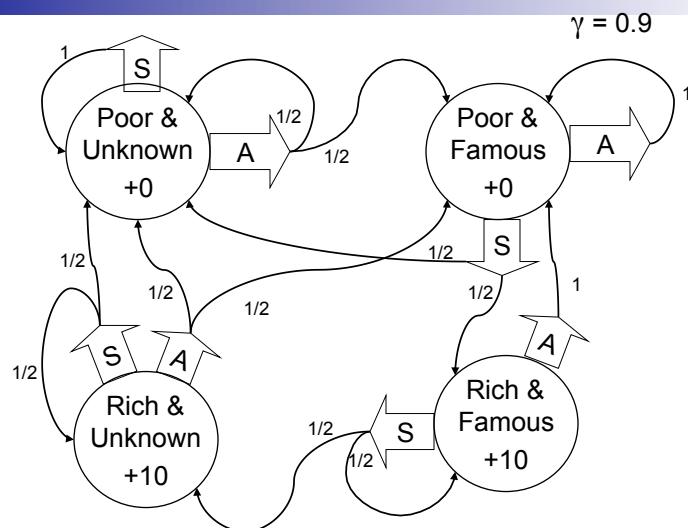
*Converges to  $V^*$*   
 $\nabla V^0$

©Carlos Guestrin 2005-2013

17

## A simple example

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

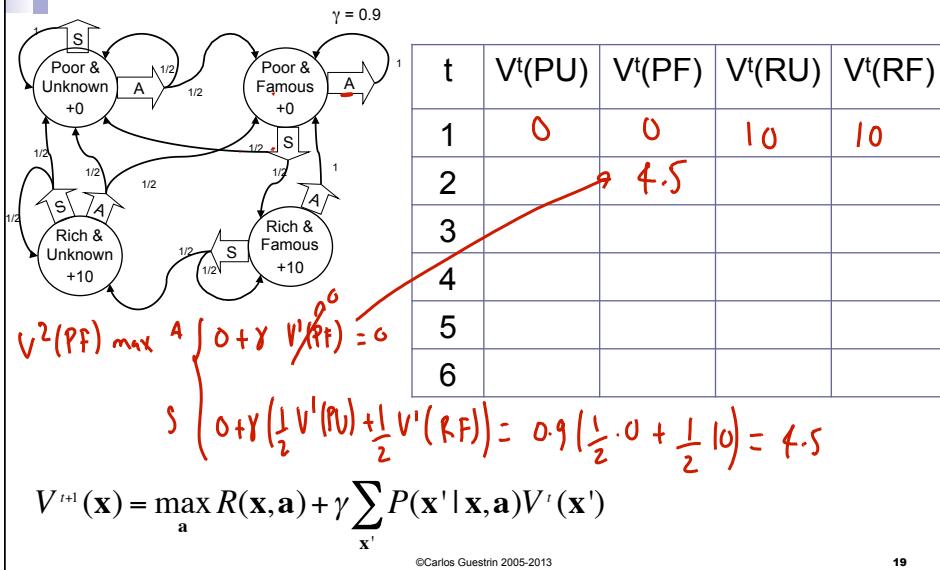


©Carlos Guestrin 2005-2013

18

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

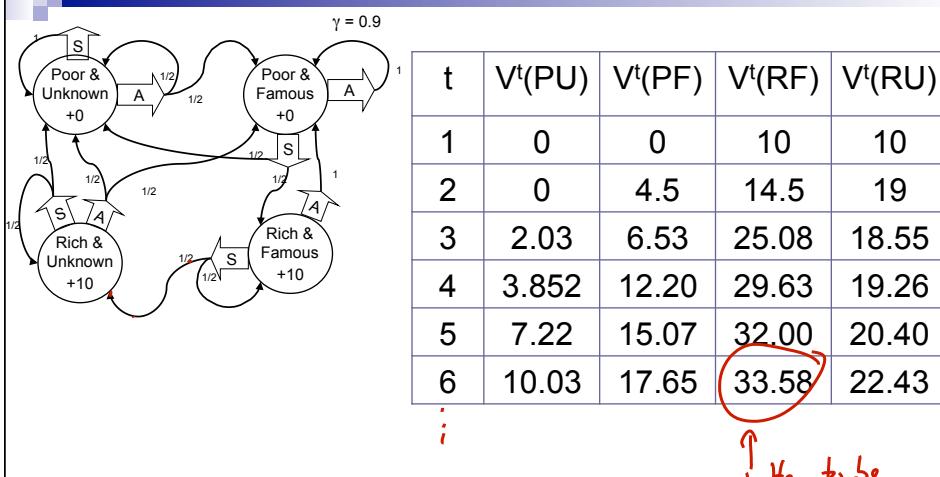
$$R(x, a) = R(x)$$



©Carlos Guestrin 2005-2013

19

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



©Carlos Guestrin 2005-2013

20

## What you need to know

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

©Carlos Guestrin 2005-2013

21

## Acknowledgment

- This lecture contains some material from Andrew Moore's excellent collection of ML tutorials:
  - <http://www.cs.cmu.edu/~awm/tutorials>

©Carlos Guestrin 2005-2013

22

Don't know  $R$ ,  
Don't know  $P$ ,  
Don't know  $\pi$

# Reinforcement Learning

Machine Learning – CSE446

Carlos Guestrin

University of Washington

June 5, 2013

©Carlos Guestrin 2005-2013

23

## 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.

©Carlos Guestrin 2005-2013

## Formalizing the (online) reinforcement learning problem

- Given a set of states  $\underline{\mathbf{X}}$  and actions  $\underline{\mathbf{A}}$ 
  - in some versions of the problem size of  $\underline{\mathbf{X}}$  and  $\underline{\mathbf{A}}$  unknown
- Interact with world at each time step  $t$ :
  - world gives state  $\underline{x}_t$  and reward  $\underline{r}_t$        $\underline{x}_t, \underline{r}_t, \underline{a}_t \rightarrow \underline{x}_{t+1}$
  - you give next action  $\underline{a}_t$
- **Goal:** (quickly) learn policy that (approximately) maximizes long-term expected discounted reward

©Carlos Gueñin 2005-2013