

Markov Decision Processes (MDPs)

Machine Learning – CSE446
Carlos Guestrin
University of Washington

June 3, 2013
©Carlos Guestrin 2005-2013

18



Reinforcement Learning

training by feedback

©Carlos Guestrin 2005-2013

19

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

20

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 $P(\mathbf{x}' | \mathbf{x}, \mathbf{a})$



©Carlos Guestrin 2005-2013

21

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

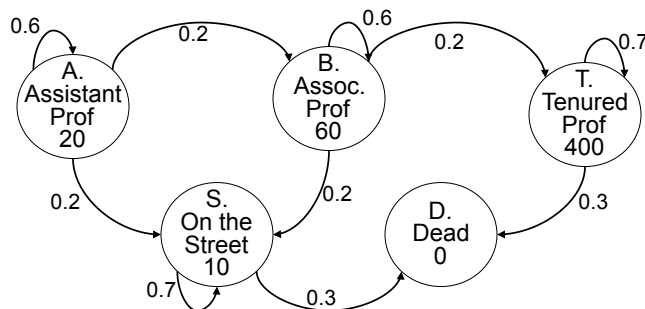
$$\begin{aligned}
 & (\text{reward now}) + \\
 & \gamma (\text{reward in 1 time step}) + \\
 & \gamma^2 (\text{reward in 2 time steps}) + \\
 & \gamma^3 (\text{reward in 3 time steps}) + \\
 & \vdots \\
 & \vdots \quad (\text{infinite sum})
 \end{aligned}$$

©Carlos Guestrin 2005-2013

22

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

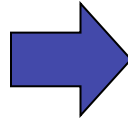
How do we compute V_A, V_B, V_T, V_S, V_D ?

©Carlos Guestrin 2005-2013

23

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



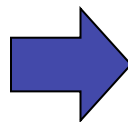
$\pi(\mathbf{x}_2)$ = peasants get gold,
footmen attack

© Carlos Guestrin 2005-2013

24

Value of Policy

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



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

Start
from \mathbf{x}_0



$R(\mathbf{x}_0)$

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



$R(\mathbf{x}_1)$

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



$R(\mathbf{x}_2)$

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



$R(\mathbf{x}_3)$

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



$R(\mathbf{x}_4)$

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

Future rewards
discounted by γ in $[0,1]$

\mathbf{x}_1'

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

$R(\mathbf{x}_1')$

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

$R(\mathbf{x}_1'')$

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

$R(\mathbf{x}_1''')$

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 \mathbf{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!

Simple approach for computing the value of a policy: Iteratively

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

- Can solve using a simple convergent iterative approach: (a.k.a. dynamic programming)

- Start with some guess V^0

- 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 $\|V_{t+1} - V_t\|_{\infty} < \epsilon$

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

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

©Carlos Guestrin 2005-2013

28

Unrolling the recursion

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

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

©Carlos Guestrin 2005-2013

29

Bellman equation

- 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

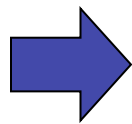
$$V^*(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V^*(\mathbf{x}')$$

©Carlos Guestrin 2005-2013

30

Optimal Long-term Plan

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



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

Optimal policy:

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

©Carlos Guestrin 2005-2013

31

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!!!

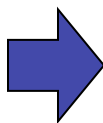
$$V_{\pi^*}(x) \geq V_{\pi}(x), \quad \forall x, \quad \forall \pi$$

©Carlos Guestrin 2005-2013

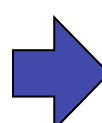
32

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

33

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
- 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 $\|V^* - V_{t+1}\|_\infty < \epsilon/(1-\gamma)$

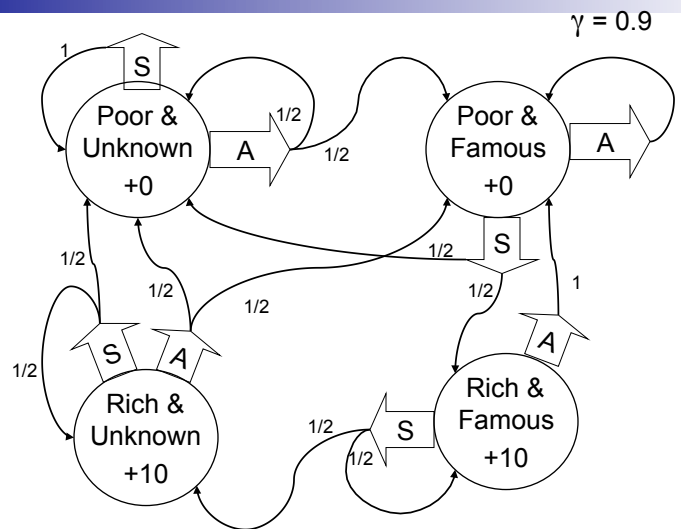
©Carlos Guestrin 2005-2013

34

A simple example

You run a startup company.

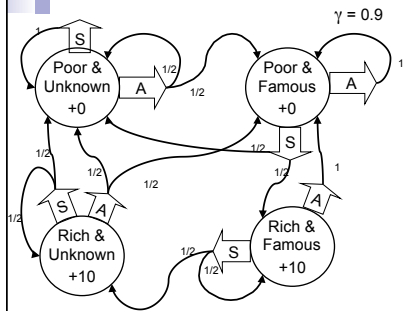
In every state you must choose between Saving money or Advertising.



©Carlos Guestrin 2005-2013

35

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



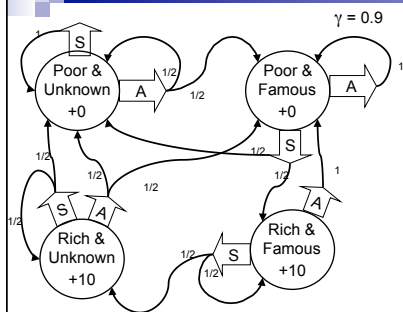
t	$V^t(\text{PU})$	$V^t(\text{PF})$	$V^t(\text{RU})$	$V^t(\text{RF})$
1				
2				
3				
4				
5				
6				

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

©Carlos Guestrin 2005-2013

36

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



t	$V^t(\text{PU})$	$V^t(\text{PF})$	$V^t(\text{RF})$	$V^t(\text{RU})$
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

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

©Carlos Guestrin 2005-2013

37

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

©Carlos Guestrin 2005-2013

38

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

39