

Markov Decision Processes (MDPs)

Machine Learning – CSE546
Carlos Guestrin
University of Washington

December 2, 2014

©Carlos Guestrin 2005-2014

1



Reinforcement Learning

training by feedback

©Carlos Guestrin 2005-2014

2

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 2005-2014

3

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 2005-2014

4

Discount Factors

$$0 < \gamma < 1$$

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

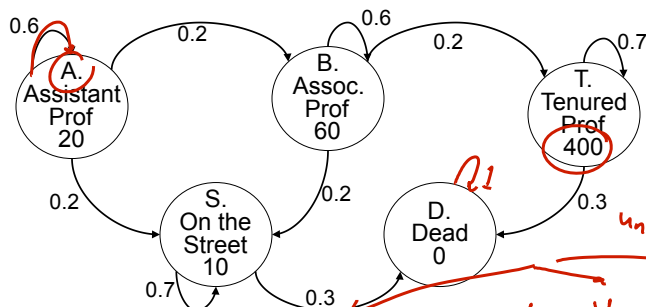
The “Discounted sum of future rewards” using discount factor γ is

$$\begin{aligned}
 & \text{(reward now)} + \gamma^2 R_3 + \dots \\
 & \gamma \text{ (reward in 1 time step)} + \dots \\
 & \gamma^2 \text{ (reward in 2 time steps)} + \dots \\
 & \gamma^3 \text{ (reward in 3 time steps)} + \dots \\
 & \vdots \\
 & \vdots \text{ (infinite sum)}
 \end{aligned}$$

Handwritten notes: $R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$
Rewards in future are worth exponentially less

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

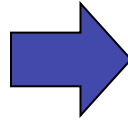
$$\begin{aligned}
 V_A &= 20 + \gamma (0.6 V_A + 0.2 V_B + 0.2 V_S) \\
 V_B &= 60 + \gamma (0.6 V_B + 0.2 V_T + 0.2 V_S)
 \end{aligned}$$

Handwritten notes: $V_D = 0$, $V_S = 10$, $V_T = 400$

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

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



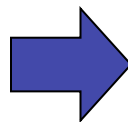
$\pi(\mathbf{x}_2) =$ peasants get gold,
footmen attack

© Carlos Guestrin 2005-2014

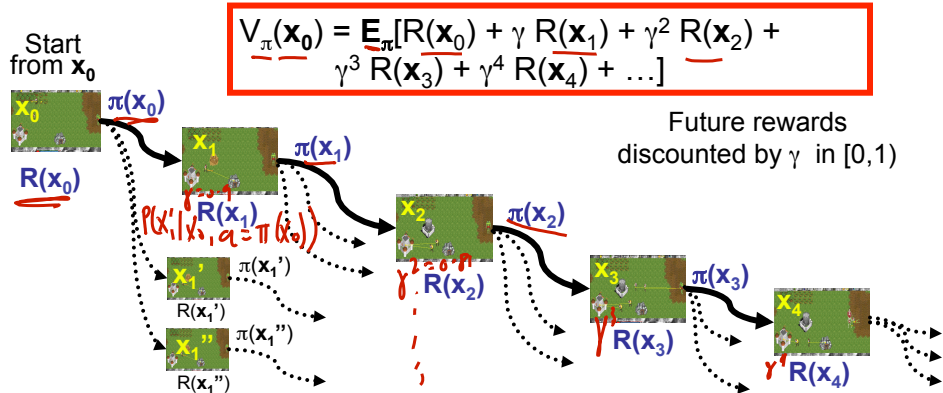
7

Value of Policy

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



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



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

$$\begin{aligned} V_{\pi}(x_0) &= E_{\pi}[R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) + \gamma^3 R(x_3) + \dots] \\ &= E_{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(x_t)\right] = E_{\pi}\left[R(x_0) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} R(x_t)\right] \\ &= R(x_0) + \gamma \sum_{x_1} P(x_1|x_0, a=\pi(x_0)) V_{\pi}(x_1) \end{aligned}$$

- A recursion!

Simple approach for computing the value of a policy: Iteratively

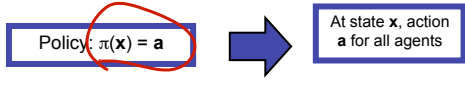
$$V_{\pi}(x) = R(x) + \gamma \sum_{x'} P(x' | x, a = \pi(x)) V_{\pi}(x')$$

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

- Start with some guess V^0 ← value for each state ← e.g. ϕ or $V^0(x) = R(x)$
- Iteratively say:
 - $V_{\pi}^{t+1}(x) \leftarrow R(x) + \gamma \sum_{x'} P(x' | x, a = \pi(x)) V_{\pi}^t(x')$ *iterate*
- 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

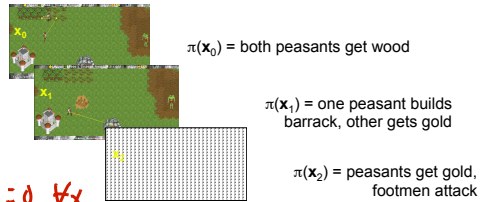
- So far, told you how good a policy is...
- But how can we choose the best policy???



- Suppose there was only one time step:
 - world is about to end!!!
 - select action that maximizes reward!

$$V^{t+1}(x) = 0 \quad \forall x$$

$$V(x) = \max_a R(x, a)$$



Unrolling the recursion

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

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

$$V^*(x_0) = \max_a R(x, a) + \gamma \sum_{x_1} P(x_1 | x_0, a) V^*(x_1)$$

↑
for each x_0

n - states
 k - actions
 n unknowns
 → Value of each state
 n - equations
 ⇒ hard because nonlinear

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

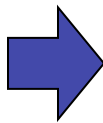
$$V^*(\mathbf{x}) = \max_{\mathbf{a}} R(\mathbf{x}, \mathbf{a}) + \gamma \sum_{\mathbf{x}'} P(\mathbf{x}' | \mathbf{x}, \mathbf{a}) V^*(\mathbf{x}')$$

©Carlos Guestrin 2005-2014

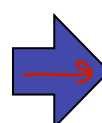
13

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 2005-2014

15

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

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

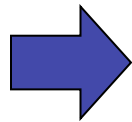
- Start with some guess V^0 ← ϕ
- Iteratively say:
 - $V^{t+1}(x) \leftarrow \max_a R(x, a) + \gamma \sum_{x'} P(x' | x, a) V^t(x')$ iterate
- Stop when $\|V_{t+1} - V_t\|_\infty < \epsilon$
 - means that $\|V^* - V_{t+1}\|_\infty < \epsilon / (1 - \gamma)$

©Carlos Guestrin 2005-2014

16

Optimal Long-term Plan

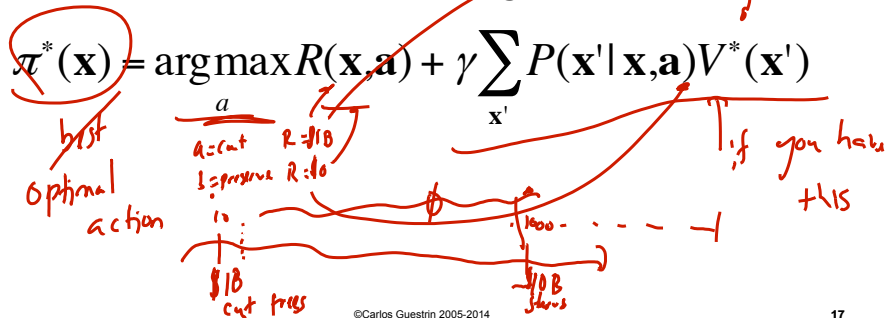
Optimal value function $V^*(\mathbf{x})$



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

Optimal policy:

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



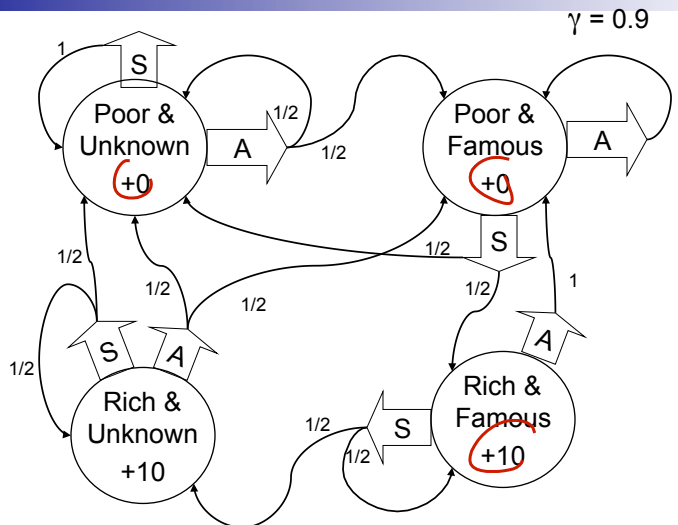
©Carlos Guestrin 2005-2014

17

A simple example

You run a startup company.

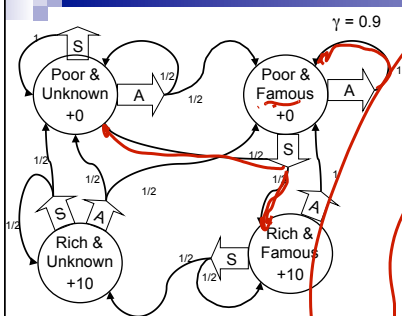
In every state you must choose between Saving money or Advertising.



©Carlos Guestrin 2005-2014

18

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	0	0	10	10
2		4.5		
3				
4				
5				
6				

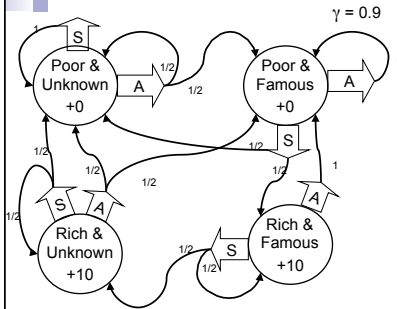
$$V^2(\text{PF}) = \max_a \left\{ \begin{array}{l} A \left\{ 0 + 0.9 V^1(\text{PF}) \right\} \\ S \left\{ 0 + 0.9 (0.5 V^1(\text{PU}) + 0.5 V^1(\text{RF})) \right\} \end{array} \right.$$

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

©Carlos Guestrin 2005-2014

19

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	0	0	10	10
2	0	4.5	14.5	19
3	2.03	9.46	17.44	25.08
4	5.17	13.61	20.17	29.13
5	8.45	16.91	22.88	32.19
6	11.41	19.62	25.43	34.78
∞	31.59	38.60	44.02	54.02

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

©Carlos Guestrin 2005-2014

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_{π}
- Optimal value function and optimal policy
 - Bellman equation
- Solving Bellman equation
 - with value iteration, policy iteration and linear programming

©Carlos Guestrin 2005-2014

21

Acknowledgment

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

Reinforcement Learning

*Rewards
&
probs are unknown*

Machine Learning – CSE546
Carlos Guestrin
University of Washington

December 4, 2014

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 2005-2014

Formalizing the (online) reinforcement learning problem

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

©Carlos Guestrin 2005-2014

The “Credit Assignment” Problem

I'm in state <u>43</u> ,	reward = <u>0</u> ,	action = <u>2</u>
“ “ “ <u>39</u> ,	“ = <u>0</u> ,	“ = <u>4</u>
“ “ “ <u>22</u> ,	“ = <u>0</u> ,	“ = <u>1</u>
“ “ “ <u>21</u> ,	“ = <u>0</u> ,	“ = <u>1</u>
“ “ “ <u>21</u> ,	“ = <u>0</u> ,	“ = <u>1</u>
“ “ “ <u>13</u> ,	“ = <u>0</u> ,	“ = <u>2</u>
“ “ “ <u>54</u> ,	“ = <u>0</u> ,	“ = <u>2</u>
“ “ “ <u>26</u> ,	“ = <u>100</u> ,	

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.

26

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

©Carlos Guestrin 2005-2014

Rmax – A model-based approach

©Carlos Guestrin 2005-2014

29

Given a dataset – learn model

Given data, learn (MDP) Representation:

- Dataset: $x_t, a_t \rightarrow r_t, x_{t+1}$

- Learn reward function:

□ $R(x,a)$ $R(x_t, a_t) = r_t$



- Learn transition model:

□ $P(x'|x,a)$

$\begin{matrix} \nearrow & \nearrow & \nearrow \\ \text{next} & \text{current} & \text{state} \end{matrix}$

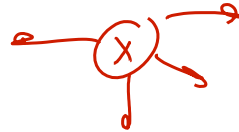
MLE

$$\frac{\text{Count}(x_{t+1}=x', x_t=x, A_t=a)}{\text{Count}(x_t=x, A_t=a)}$$

©Carlos Guestrin 2005-2014

Planning with insufficient information

- Model-based approach:
 - estimate $R(\mathbf{x}, \mathbf{a})$ & $P(\mathbf{x}'|\mathbf{x}, \mathbf{a})$
 - obtain policy by value or policy iteration, or linear programming
 - No credit assignment problem!
 - learning model, planning algorithm takes care of "assigning" credit
- What do you plug in when you don't have enough information about a state?
 - don't reward at a particular state
 - plug in 0?
 - plug in smallest reward (R_{\min})? ← *pessimistic*
 - plug in largest reward (R_{\max})? ← *optimistic*
 - don't know a particular transition probability?



©Carlos Guestrin 2005-2014

A surprisingly simple approach for model based RL – The Rmax algorithm [Brafman & Tenenbholz]

- **Optimism in the face of uncertainty!!!!**
 - heuristic shown to be useful long before theory was done (e.g., Kaelbling '90)
- If you don't know reward for a particular state-action pair, set it to R_{\max} !!!
- If you don't know the transition probabilities $P(\mathbf{x}'|\mathbf{x}, \mathbf{a})$ from some some state action pair \mathbf{x}, \mathbf{a} assume you go to a **magic, fairytale** new state \mathbf{x}_0 !!!
 - $R(\mathbf{x}_0, \mathbf{a}) = R_{\max}$
 - $P(\mathbf{x}_0|\mathbf{x}_0, \mathbf{a}) = 1$

©Carlos Guestrin 2005-2014

The Rmax algorithm

Initialization:

- Add state x_0 to MDP
- $R(x,a) = R_{\max}, \forall x,a$
- $P(x_0|x,a) = 1, \forall x,a$
- all states (except for x_0) are unknown

Repeat

- obtain policy for current MDP and Execute policy
- for any visited state-action pair, set reward function to appropriate value $R(x_t, a_t) = r_t$
- if visited some state-action pair x,a enough times to estimate $P(x'|x,a)$
 - update transition probs. $P(x'|x,a)$ for x,a using MLE
 - recompute policy

©Carlos Guestrin 2005-2014

Visit enough times to estimate $P(x'|x,a)$?

How many times are enough?

- use Chernoff Bound!

Chernoff Bound:

- X_1, \dots, X_n are i.i.d. Bernoulli trials with prob. θ
- $P(|1/n \sum_i X_i - \theta| > \epsilon) \leq \exp\{-2n\epsilon^2\}$

©Carlos Guestrin 2005-2014

Putting it all together

- **Theorem:** With prob. at least $1-\delta$, Rmax will reach a ϵ -optimal policy in time polynomial in: num. states, num. actions, T , $1/\epsilon$, $1/\delta$
 - Every T steps:
 - achieve near optimal reward (great!) or
 - visit an unknown state-action pair ! num. states and actions is finite, so can't take too long before all states are known

©Carlos Guestrin 2005-2014

What you need to know about RL...

- Neither supervised, nor unsupervised learning
- Try to learn to act in the world, as we travel states and get rewards
- Model-based & Model-free approaches
- Rmax, a model based approach:
 - Learn model of rewards and transitions
 - Address exploration-exploitation tradeoff
 - Simple algorithm, great in practice

©Carlos Guestrin 2005-2014

39

Closing....

©Carlos Guestrin 2005-2014

40

What you have learned this quarter

- Learning is function approximation
- Point estimation
- Regression
- LASSO
- Subgradient
- Stochastic gradient descent
- Coordinate descent
- Discriminative v. Generative learning
- Naïve Bayes
- Logistic regression
- Bias-Variance tradeoff
- Decision trees
- Cross validation
- Boosting
- Instance-based learning
- Perceptron
- SVMs
- Kernel trick
- PAC learning
- Bayes nets
 - representation, parameter and structure learning
- K-means
- EM
- Mixtures of Gaussians
- Dimensionality reduction, PCA
- MDPs
- Reinforcement learning



©Carlos Guestrin 2005-2014

BIG PICTURE

- Improving the performance at some task though experience!!! 😊
 - before you start any learning task, remember the fundamental questions:

What is the learning problem?

From what experience?

What model?

What loss function are you optimizing?

With what optimization algorithm?

Which learning algorithm?

With what guarantees?

How will you evaluate it?

©Carlos Guestrin 2005-2014

You have done a lot!!!

- And (hopefully) learned a lot!!!
 - Implemented
 - LASSO
 - LR
 - Perceptron
 - Clustering
 - ...
 - Answered hard questions and proved many interesting results
 - Completed (I am sure) an amazing ML project
 - And did excellently on the final!
- Now you are ready for one of the most sought-after careers in industry today!!! 😊

Thank You for the
Hard Work!!!

©Carlos Guestrin 2005-2014