

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

# Reinforcement Learning

Machine Learning – CSE446

Carlos Guestrin

University of Washington

June 5, 2013

©Carlos Guestrin 2005-2013

1

## 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. ↗  $P(77|3, a=2)$

**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  $X$  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  $r_t$   $\underline{x_t}, \underline{r_t}, \underline{a_t} \rightarrow \underline{x_{t+1}}$
  - you give next action  $a_t$
- **Goal:** (quickly) learn policy that (approximately) maximizes long-term expected discounted reward

©Carlos Guestrin 2005-2013

## The “Credit Assignment” Problem

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

©Carlos Guestrin 2005-2013

## Two main reinforcement learning approaches

- Model-based approaches:
  - explore environment, then 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

©Carlos Guestrin 2005-2013

# Rmax – A model-based approach

©Carlos Guestrin 2005-2013

7

## 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:

□  $P(x'|x, a)$

counts like BNs  
next state  
current state  
MLE

count(x', x, a) = 1  
count(x, a) — action = 1  
current state = 32

43 → 32 → 43 → 31 ...

1 2 2 1 ...



©Carlos Guestrin 2005-2013

## 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 ( $R_{\min}$ )?
    - plug in largest reward ( $R_{\max}$ )?
  - don't know a particular transition probability?

©Carlos Guestrin 2005-2013

## Some challenges in model-based RL 2: Exploration-Exploitation tradeoff

- A state may be very hard to reach
  - waste a lot of time trying to learn rewards and transitions for this state
  - after a much effort, state may be useless
- A strong advantage of a model-based approach:
  - you know which states estimate for rewards and transitions are bad
  - can (try) to plan to reach these states
  - have a good estimate of how long it takes to get there

©Carlos Guestrin 2005-2013

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

## ■ 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 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 2005-2013

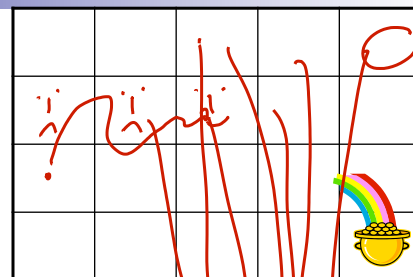
## Understanding $R_{\max}$

future rewards are not subject  

$$V^g(x_0) = R_{\max} + \gamma R_{\max} + \gamma^2 R_{\max} \dots$$

$$= R_{\max} / (1 - \gamma)$$

- With  $R_{\max}$  you either:
  - **explore** – visit a state-action pair you don't know much about
    - because it seems to have lots of potential
  - **exploit** – spend all your time on known states
    - even if unknown states were amazingly good, it's not worth it
- Note: you never know if you are exploring or exploiting!!!



You have to know  $R_{\max}$ , not true in life



©Carlos Guestrin 2005-2013

# Implicit Exploration-Exploitation Lemma

## ■ Lemma: every $T$ time steps, either:

- **Exploits**: achieves near-optimal reward for these  $T$ -steps, or
- **Explores**: with high probability, the agent visits an unknown state-action pair
  - learns a little about an unknown state
- $T$  is related to mixing time of Markov chain defined by MDP
  - time it takes to (approximately) forget where you started

©Carlos Guestrin 2005-2013

## 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
- 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 2005-2013

only works in finite state spaces

Lemma either

near-optimal

OR

learn something new

Plan with optimism until you learn something new then replan

$$R(x_t, a_t) = r_t$$

## Visit enough times to estimate $P(\mathbf{x}'|\mathbf{x},\mathbf{a})$ ?

- How many times are enough?

- use Chernoff Bound!

- **Chernoff Bound:**

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

©Carlos Guestrin 2005-2013

## Putting it all together

■ **Theorem:** With prob. at least  $1-\delta$ ,  $R_{\max}$  will reach a  $\epsilon$ -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 2005-2013



## 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 through optimism*
  - Learn model of rewards and transitions
  - Address exploration-exploitation tradeoff
  - Simple algorithm, great in practice

# Closing....

# 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

*Handwritten signature*



©Carlos Guestrin 2005-2013

# 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 2005-2013

# 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!

## Thank You for the Hard Work!!!

©Carlos Gue~~ra~~tin 2005-2013