#### CS 188: Artificial Intelligence Course Summary



Instructors: John Canny and Oliver Grillmeyer --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

#### Announcements

- Final Exam
  - Thursday 5/15 from 3:00-6:00pm
  - See Exam Logistics on CS 188 website
- Please fill out course evaluations
- Oliver's RRR office hours will be remote
  - Use CS 188 zoom link
  - **5/6, 5/8 2:30-4:00**
  - **5/13 3:30-5:00**

# CS 188: Artificial Intelligence

#### **Review: Markov Decision Processes and RL**



University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

# MDPs/RL vs Human Thought



System II reasoning (human)

# System I reasoning (human)

# **Markov Decision Processes**

- An MDP is defined by:
  - A set of states  $s \in S$
  - A set of actions  $a \in A$
  - A transition function T(s, a, s')
    - Probability that a from s leads to s', i.e., P(s' | s, a)
    - Also called the model or the dynamics
  - A reward function R(s, a, s')
    - Sometimes just R(s) or R(s')
  - A start state
  - Maybe a terminal state

#### MDPs are non-deterministic search problems

- One way to solve them is with expectimax search
- We'll have a new tool soon



# Grid World Example



S	а	s'	R
(1,1)	north		

T(s, a, s'):

- T((1,1), north, (2,1)) = 0.8
- T((1,1), north, (1,2)) = 0.1
- T((1,1), north, (1,1)) = 0.1

# Grid World Example



S	а	s'	R
(1,1)	north	(2,1)	-0.1

**R(s, a, s'):** R((1,1), north, (2,1)) = -0.1

# What is Markov about MDPs?

- "Markov" generally means that given the present state, the future and the past are independent
- For Markov decision processes, "Markov" means action outcomes depend only on the current state

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \dots S_0 = s_0)$$

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

 This is just like search, where the successor function could only depend on the current state (not the history)



Andrey Markov (1856-1922)

# Policies

- In deterministic single-agent search problems, we wanted an optimal plan, or sequence of actions, from start to a goal
- For MDPs, we want an optimal policy  $\pi^*: S \rightarrow A$ 
  - A policy  $\pi$  gives an action for each state
  - An optimal policy is one that maximizes expected utility if followed
  - An explicit policy defines a reflex agent
- Expectimax didn't compute entire policies
  - It computed the action from a single state only



Optimal policy when R(s, a, s') = -0.03 for all non-terminals s

# Discounting

- It's reasonable to maximize the sum of rewards
- It's also reasonable to prefer rewards now to rewards later
- One solution: values of rewards decay exponentially



# Solving MDPs



# **Optimal Quantities**

- The value (utility) of a state s:
  - V\*(s) = expected utility starting in s and acting optimally
- The value (utility) of a q-state (s,a):
  - Q\*(s,a) = expected utility starting out having taken action a from state s and (thereafter) acting optimally
- The optimal policy:
  - $\pi^*(s)$  = optimal action from state s



# Snapshot of Demo – Gridworld V Values

○ ○ ○ Gridworld Display				
	0.64 )	0.74 →	0.85 )	1.00
	•		• 0.57	-1.00
	• 0.49	∢ 0.43	• 0.48	∢ 0.28
	VALUES	SAFTER 1	LOO ITERA	ATIONS

# Snapshot of Demo – Gridworld Q Values



# Value Iteration

- Start with V<sub>0</sub>(s) = 0: no time steps left means an expected reward sum of zero
- Given vector of V<sub>k</sub>(s) values, do one ply of expectimax from each state:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$



Repeat until convergence

0.0	Gridworl	d Display		
	<b>^</b>	<b>^</b>		
0.00	0.00	0.00	0.00	
		<b>^</b>		
0.00		0.00	0.00	
	<b>^</b>	<b>^</b>	<b>^</b>	
0.00	0.00	0.00	0.00	
VALUES AFTER O TTERATIONS				

0	0	Gridworl	d Display	-	
	▲ 0.00	• 0.00	0.00 )	1.00	
	•		∢ 0.00	-1.00	
	•	• 0.00	• 0.00	0.00	
	VALUES AFTER 1 TTERATIONS				

k=2

0 0	Gridworl	d Display	al and "Restored al advances they "New as you manual way"
• 0.00	0.00 >	0.72 ▸	1.00
•		• 0.00	-1.00
• 0.00	• 0.00	• 0.00	0.00
VALUES AFTER 2 ITERATIONS			

k=3

0	0	Gridworl	d Display	
	0.00 )	0.52 →	0.78 )	1.00
	• 0.00		• 0.43	-1.00
	• 0.00	• 0.00	• 0.00	0.00
	VALUES AFTER 3 ITERATIONS			

k=4

0 0	Gridworl	d Display		
0.37 ▸	0.66 )	0.83 )	1.00	
•		• 0.51	-1.00	
• 0.00	0.00 →	• 0.31	∢ 0.00	
VALUES AFTER 4 ITERATIONS				

k=5

0.0	0	Gridworl	d Display		
	0.51 )	0.72 )	0.84 )	1.00	
	• 0.27		• 0.55	-1.00	
	•	0.22 →	• 0.37	∢ 0.13	
	VALUES AFTER 5 ITERATIONS				

k=6

Gridworld Display					
	0.59)	0.73 )	0.85 )	1.00	
	• 0.41		• 0.57	-1.00	
	• 0.21	0.31 →	▲ 0.43	∢ 0.19	
	VALUES AFTER 6 ITERATIONS				

00	Gridworl	d Display	
0.62)	0.74 →	0.85 →	1.00
• 0.50		• 0.57	-1.00
▲ 0.34	0.36 →	▲ 0.45	◀ 0.24
VALU	ES AFTER	7 ITERA	FIONS

k=8

0 0	Gridworl	d Display	
0.63 )	0.74 →	0.85 )	1.00
• 0.53		• 0.57	-1.00
• 0.42	0.39 )	• 0.46	∢ 0.26
VALUES AFTER 8 ITERATIONS			

k=9

Gridworld Display				
ſ	0.64 )	0.74 ▸	0.85 )	1.00
	• 0.55		• 0.57	-1.00
	• 0.46	0.40 →	• 0.47	∢ 0.27
VALUES AFTER 9 ITERATIONS				

Gridworld Display				
	0.64 )	0.74 ▸	0.85 )	1.00
	• 0.56		• 0.57	-1.00
	• 0.48	∢ 0.41	• 0.47	◀ 0.27
VALUES AFTER 10 ITERATIONS				

0 0	C Cridworld Display			
	0.64 )	0.74 →	0.85 →	1.00
	• 0.56		• 0.57	-1.00
	•	◀ 0.42	• 0.47	∢ 0.27
VALUES AFTER 11 ITERATIONS				

O O Gridworld Display			
0.64	▶ 0.74 ▶	0.85 )	1.00
• 0.57		• 0.57	-1.00
• 0.49	◀ 0.42	• 0.47	◀ 0.28
VALUES AFTER 12 ITERATIONS			

O O Gridworld Display			
0.64 )	0.74 →	0.85 )	1.00
•		•	
0.57		0.57	-1.00
<b>^</b>		<b>^</b>	
0.49	∢ 0.43	0.48	∢ 0.28
VALUES AFTER 100 ITERATIONS			

# The Bellman Equations



# The Bellman Equations

 Definition of "optimal utility" via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values

$$V^{*}(s) = \max_{a} Q^{*}(s, a)$$

$$Q^{*}(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{*}(s') \right]$$

$$Thes V^{*}(s) = \max_{a} \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{*}(s') + \gamma V^{*}(s') \right]$$



# Value Iteration

Bellman equations characterize the optimal values:

$$V^{*}(s) = \max_{a} \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{*}(s') \right]$$

• Value iteration computes them:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

- Value iteration is just a fixed point solution method
  - ... though the V<sub>k</sub> vectors are also interpretable as time-limited values



# **Policy Evaluation**

- How do we calculate the V's for a fixed policy  $\pi$ ?
- Idea 1: Turn recursive Bellman equations into updates (like value iteration)

$$V_0^{\pi}(s) = 0$$
  
$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^{\pi}(s')$$

- Efficiency: O(S<sup>2</sup>) per iteration
- Idea 2: Without the maxes, the Bellman equations are just a linear system
  - Solve with Matlab (or your favorite linear system solver)



# **Computing Actions from Values**

- Let's imagine we have the optimal values V\*(s)
- How should we act?
  - It's not obvious!
- We need to do a mini-expectimax (one step)



$$\pi^{*}(s) = \arg\max_{a} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{*}(s')]$$

• This is called **policy extraction**, since it gets the policy implied by the values

# **Computing Actions from Q-Values**

- Let's imagine we have the optimal q-values:
- How should we act?
  - Completely trivial to decide!

$$\pi^*(s) = \arg\max_a Q^*(s,a)$$



Important lesson: actions are easier to select from q-values than values!

# **Policy Iteration**

- Alternative approach for optimal values:
  - Step 1: Policy evaluation: calculate utilities for some fixed policy (not optimal utilities!) until convergence
  - Step 2: Policy improvement: update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
  - Repeat steps until policy converges

#### This is policy iteration

- It's still optimal!
- Can converge (much) faster under some conditions
# **Example: Policy Iteration**



#### Improved Policy using Q-Values



#### Improve again – Optimal!



Q-Values for the above policies

# **Reinforcement Learning**

- Still assume a Markov decision process (MDP):
  - A set of states  $s \in S$
  - A set of actions (per state) A
  - A model T(s,a,s')
  - A reward function R(s,a,s')
- Still looking for a policy π(s)



- New twist: don't know T or R
  - I.e. we don't know which states are good or what the actions do
  - Must try out actions and states to learn
    - Q1: How to learn from things tried? (today, Passive Reinforcement Learning)
    - Q2: What to decide to try? (Thursday, Active Reinforcement Learning)

## **Classical Reinforcement Learning Diagram**



- Basic idea:
  - Must (learn to) act so as to maximize expected rewards
  - All learning is based on observed samples of outcomes!

### Model-Free Reinforcement Learning



# **Policy Evaluation: Problem Setting**

#### Simplified task: policy evaluation

- Input: a fixed policy  $\pi(s)$
- You don't know the transitions T(s,a,s')
- You don't know the rewards R(s,a,s')
- Goal: learn the state values

#### In this case:

- Learner is "along for the ride"
- No choice about what actions to take
- Just execute the policy and learn from experience
- This is NOT offline planning! You actually take actions in the world.



# **Problems with Direct Policy Evaluation**

#### What's good about direct evaluation?

 It eventually computes the correct average values, using just sample transitions

#### What bad about it?

- It wastes information about state connections
- Each state must be learned separately
- So, it takes a long time to learn

#### **Output Values**



If B and E both go to C under this policy, how can their values be different?

## **Temporal Difference Learning**

- Big idea: learn from every experience!
  - Update V(s) each time we experience a transition (s, a, s', r)
  - Likely outcomes s' will contribute updates more often
- Temporal difference learning of values
  - Policy still fixed, still doing evaluation!
  - Move values toward value of whatever successor occurs: running average

Sample of V(s):  $sample = R(s, \pi(s), s') + \gamma V^{\pi}(s')$ Update to V(s):  $V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + (\alpha)sample$ Same update:  $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(sample - V^{\pi}(s))$ 



#### Problems with TD Value Learning

- TD value leaning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages
- However, if we want to turn values into an optimal policy, we're sunk:

$$\pi(s) = \arg\max_{a} Q(s, a)$$
$$Q(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V(s') \right]$$

- Idea: learn Q-values, not values
- Makes action selection model-free too!



### **Recall: Q-Value Iteration**

- Value iteration: find successive (depth-limited) values
  - Start with  $V_0(s) = 0$ , which we know is right
  - Given V<sub>k</sub>, calculate the depth k+1 values for all states:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

- But Q-values are more useful, so compute them instead
  - Start with  $Q_0(s,a) = 0$ , which we know is right
  - Given Q<sub>k</sub>, calculate the depth k+1 q-values for all q-states:

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

# Q-Learning

Q-Learning: sample-based Q-value iteration

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

- Learn Q(s,a) values as you go
  - Receive a sample (s,a,s',r)
  - Consider your old estimate: Q(s, a)
  - Consider your new sample estimate:

 $sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$ 

Incorporate the new estimate into a running average:

 $Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha) [sample]$ 

[Demo: Q-learning – gridworld (L10D2)] [Demo: Q-learning – crawler (L10D3)]

# **Q-Learning Properties**

- Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!
- This is called off-policy learning
- Caveats:
  - You have to explore enough
  - You have to eventually make the learning rate small enough
  - ... but not decrease it too quickly
  - Basically, in the limit, it doesn't matter how you select actions (!)



## Exploration - ε-greedy

#### Several schemes for forcing exploration

- Simplest: random actions (ε-greedy)
  - Every time step, flip a coin
  - With (small) probability  $\varepsilon$ , act randomly
  - With (large) probability 1- $\varepsilon$ , act on current policy



# **Exploration Functions**

#### When to explore?

- Random actions: explore a fixed amount
- Better idea: explore areas whose badness is not (yet) established, eventually stop exploring

Exploration function

 Takes a value estimate u and a visit count n, and returns an optimistic utility, e.g.



f(u,n) = u + k/n

**Regular Q-Update:**  $Q(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} Q(s', a')$ 

• Note: this propagates the "bonus" back to states that lead to unknown states as well! Modified Q-Update:  $Q(s,a) \leftarrow_{\alpha} R(s,a,s') + \gamma \max_{a'} f(Q(s',a'), N(s',a'))$ 

[Demo: exploration – Q-learning – crawler – exploration function (L11D4)]

# Q-Learning with Experience Replay

#### Problem:

 Need to repeat same (s,a,s',r) transitions in environment many times to propagate values

#### Solution:

- Collect transitions in a memory buffer and "replay" them to update Q values
  - Uses memory of transitions only, no need to repeat them in environment
- Evidence of such experience replay in the brain



# Q-Learning with Experience Replay

#### • At each step:

- Receive a sample transition (s,a,s',r)
- Add (s,a,s',r) to replay buffer
- Repeat N times:
  - Randomly pick transition (s,a,s',r) from replay buffer
  - Make sample based on (s,a,s',r):  $sample = R(s,a,s') + \gamma \max_{a'} Q(s',a')$
  - Update Q based on picked sample:  $Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha) [sample]$



### **Feature-Based Representations**

- Solution: describe a state using a vector of features (properties) f<sub>1</sub>, f<sub>2</sub>, ...
  - Features are functions from states to real numbers (often in [0,1]) that capture important properties of the state
  - Example features:
    - Distance to closest ghost
    - Distance to closest dot
    - Number of ghosts
    - 1 / (dist to dot)<sup>2</sup>
    - Is Pacman in a tunnel? (0/1)
    - ..... etc.
    - Is it the exact state on this slide?
  - Can also describe a q-state (s, a) with features (e.g. action moves closer to food)



## **Linear Value Functions**

Using a feature representation f<sub>1</sub>, f<sub>2</sub>, ... we can write a q function (or value function) for any state using a few weights w<sub>1</sub>, w<sub>2</sub>, ... :

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \ldots + w_n f_n(s)$$

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a)$$

- Advantage: our experience is summed up in a few powerful numbers w<sub>1</sub>, w<sub>2</sub>, ...
- Disadvantage: states may share features but actually be very different in value!



the same value if we don't include ghost positions as a feature:

# Approximate Q-Learning

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a)$$

Q-learning with linear Q-functions:

transition = (s, a, r, s')difference =  $\left[r + \gamma \max_{a'} Q(s', a')\right] - Q(s, a)$   $Q(s, a) \leftarrow Q(s, a) + \alpha$  [difference] Exact Q's  $w_i \leftarrow w_i + \alpha$  [difference]  $f_i(s, a)$  Approximate Q's



- Intuitive interpretation:
  - Adjust weights of active features
  - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features
- Formal justification: online least squares, gradient descent

# **Policy Search**

- Problem: often the feature-based policies that work well (win games, maximize utilities) aren't the ones that approximate V / Q best
  - Q-learning's priority: get Q-values close (modeling)
  - Action selection priority: get ordering of Q-values right (prediction)
  - We'll see this distinction between modeling and prediction again later in the course
- Solution: learn policies  $\pi$  that maximize rewards, not the Q values that predict them
- Policy search: start with an ok solution (e.g. Q-learning) then fine-tune by hill climbing on feature weights

# **Policy Search**

- Simplest policy search:
  - Start with an initial linear value function or Q-function
  - Nudge each feature weight up and down and see if your policy is better than before
- Problems:
  - How do we tell the policy got better?
  - Need to run many sample episodes!
  - If there are a lot of features, this can be impractical
- Better methods exploit lookahead structure, sample wisely, change multiple parameters...
  - Policy Gradient, Proximal Policy Optimization (PPO) are examples

### The Story So Far: MDPs and RL

#### **Known MDP: Offline Solution**

Goal	Technique
Compute V*, Q*, π*	Value / policy iteration
Evaluate a fixed policy $\boldsymbol{\pi}$	Policy evaluation

#### Unknown MDP: Model-Based

Goal	Technique
Compute V*, Q*, π*	VI/PI on approx. MDP
Evaluate a fixed policy $\pi$	PE on approx. MDP

#### Unknown MDP: Model-Free

Goal	Technique
Compute V*, Q*, π*	Q-learning
Evaluate a fixed policy $\pi$	Value Learning



[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

# **Random Variables**

- A random variable is some aspect of the world about which we (may) have uncertainty
  - R = Is it raining?
  - T = Is it hot or cold?
  - D = How long will it take to drive to work?
  - L = Where is the ghost?
- We denote random variables with capital letters
- Like variables in a CSP, random variables have domains
  - R in {true, false} (often write as {+r, -r})
  - T in {hot, cold}
  - D in [0, ∞)
  - L in possible locations, maybe {(0,0), (0,1), ...}



# **Probability Distributions**

- Associate a probability with each value of that random variable
  - Temperature:

• Weather:











W	Р
sun	0.6
rain	0.1
fog	0.3
meteor	0.0

# **Probabilistic Models**

- A probabilistic model is a joint distribution over a set of random variables
- Probabilistic models:
  - (Random) variables with domains
  - Assignments are called *outcomes*
  - Joint distributions: say whether assignments (outcomes) are likely
  - *Normalized:* sum to 1.0
  - Ideally: only certain variables directly interact

Т	W	Р
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

Distribution over T,W



# **Marginal Distributions**

P(x,y)

 $\mathbf{D}(\mathbf{T})$ 

- Marginal distributions are sub-tables which eliminate random variables
- Marginalization (summing out): Combine collapsed rows by adding

$$P(T, W)$$

$$P(T, W)$$

$$P(t) = \sum_{w} P(t, w)$$

$$\frac{T \quad W \quad P}{hot \quad sun \quad 0.4}$$

$$P(w) = \sum_{t} P(t, w)$$

$$P(W)$$

 $x_2$ 

hidden (unobserved) variables

## **Conditional Probabilities**

- A simple relation between joint and conditional probabilities
  - In fact, this is taken as the *definition* of a conditional probability

evidence

 $P(a|b) = \frac{P(a,b)}{P(b)}$ query = (proportion of b where a holds)



### Normalization Trick



# Inference by Enumeration

- General case:
  - Evidence variables:
  - Query\* variable:
  - Hidden variables:
  - Step 1: Select the entries consistent with the evidence



 Step 2: Sum out H to get joint of Query and evidence We want:

\* Works fine with multiple query variables, too

 $P(Q|e_1\ldots e_k)$ 

Step 3: Normalize



 $Z = \sum_{q} P(Q, e_1 \cdots e_k)$  $P(Q|e_1 \cdots e_k) = \frac{1}{Z} P(Q, e_1 \cdots e_k)$ 





 $P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(\underbrace{Q, h_1 \dots h_r, e_1 \dots e_k})$  $X_1, X_2, \ldots, X_r$ 

0.15

## Inference with Bayes' Rule

Example: Diagnostic probability from causal probability:

$$P(\text{cause}|\text{effect}) = \frac{P(\text{effect}|\text{cause})P(\text{cause})}{P(\text{effect})}$$

• Example:

M: meningitis, S: stiff neck

$$\begin{array}{c} P(+m) = 0.0001 \\ P(+s|+m) = 0.8 \\ P(+s|-m) = 0.01 \end{array} \begin{array}{c} \mbox{Example givens} \end{array}$$

 $P(+m|+s) = \frac{P(+s|+m)P(+m)}{P(+s)} = \frac{P(+s|+m)P(+m)}{P(+s|+m)P(+m) + P(+s|-m)P(-m)} = \frac{0.8 \times 0.0001}{0.8 \times 0.0001 + 0.01 \times 0.999}$ 

 $P(+m | +s) \approx 0.008$ 

### CS 188: Artificial Intelligence

### Bayes' Nets



[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

## Independence

• Two variables are *independent* if:

$$\forall x, y : P(x, y) = P(x)P(y)$$

- This says that their joint distribution *factors* into a product two simpler distributions
- Another form:

$$\forall x, y : P(x \mid y) = P(x)$$

- We write:  $X \perp \!\!\!\perp Y$
- Independence is a simplifying modeling assumption
  - *Empirical* joint distributions: at best "close" to independent
  - What could we assume for {Weather, Traffic, Cavity, Toothache}?



# **Conditional Independence**

- Unconditional (absolute) independence very rare between variables in the same system (why?)
- Conditional independence is our most basic and robust form of knowledge about uncertain environments.
- X is conditionally independent of Y given Z written if and only if:  $X \! \perp \!\!\!\perp Y | Z$

$$\forall x, y, z : P(x, y | z) = P(x | z)P(y | z)$$
  
or, equivalently, if and only if

$$\forall x, y, z : P(x \mid y, z) = P(x \mid z)$$
  
or

$$\forall x, y, z : P(y \mid x, z) = P(y \mid z)$$

# Conditional Independence and the Chain Rule



# Conditional Independence and the Chain Rule

Chain rule:  $P(X_1, X_2, ..., X_n) = P(X_1)P(X_2 | X_1)P(X_3 | X_1, X_2)...$ Partial Dependency graph: Independences from the Graph:





# Example: Coin Flips



#### No interactions between variables: absolute independence
#### Example: Alarm Network



#### Example: Alarm Network



+e

+e

-e

-e

-b

-b

-b

+a

-a

+a

-a

0.71

0.001

0.999

$$P(+b, -e, +a, -j, +m) =$$

#### Example: Alarm Network



#### **Conditional Independence**

X and Y are independent if

$$\forall x, y \ P(x, y) = P(x)P(y) \ \neg \neg \neg \Rightarrow \ X \bot\!\!\!\bot Y$$

X and Y are conditionally independent given Z

$$\forall x, y, z \ P(x, y|z) = P(x|z)P(y|z) \dashrightarrow X \perp Y|Z$$

(Conditional) independence is a property of a distribution

Example:

 $A larm \bot Fire | Smoke$ 



#### **D-separation: Overview**

#### D-separation:

 a condition / algorithm for answering conditional independence queries from just studying the graph

#### How:

- Study independence properties for triples
- Analyze complex cases as composition of triples

#### **Recap of Triples**



#### **D-Separation**

- A path is active if each (overlapping) triple is active:
  - Note: e.g. for a path A B C D E, the triples are:
  - A B C, B C D, C D E

Note: all it takes to block a path is a single inactive segment

- Are X and Y "D-separated" given evidence variables {Z}?
  - Consider all (undirected) paths from X to Y
  - If none of the paths are active, then X and Y are D-separated given {Z}
  - On the other hand, if there is at least one active path, then X and Y are not D-separated given {Z}
- Independence and D-separation:

X and Y are guaranteed conditionally independent given {Z} IF AND ONLY IF X and Y are d-separated given {Z}

 $\rightarrow$  just need to check the graph



## Example

Red = Nodes are conditionally independent given the evidence Blue = Nodes are d-separated given the evidence

 $R \perp\!\!\!\!\perp B$ YesYes $R \perp\!\!\!\!\perp B | T$ No?? $R \perp\!\!\!\!\!\perp B | T'$ No??



# Inference

- Inference: calculating some useful quantity from a joint probability distribution
- Examples:
  - Posterior probability

 $P(Q|E_1 = e_1, \dots E_k = e_k)$ 

- Most likely explanation:
  - $\operatorname{argmax}_q P(Q = q | E_1 = e_1 \ldots)$



# Inference by Enumeration

- General case:
  - Evidence variables:
  - Query\* variable:
  - Hidden variables:
- $E_{1} \dots E_{k} = e_{1} \dots e_{k}$  Q  $H_{1} \dots H_{r}$   $X_{1}, X_{2}, \dots X_{n}$ All variables

 $P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(\underbrace{Q, h_1 \dots h_r, e_1 \dots e_k}_{X_1, X_2, \dots X_n})$ 

We want:

\* Works fine with multiple query variables, too

$$P(Q|e_1\ldots e_k)$$

 Step 1: Select the entries consistent with the evidence

-3

-1

5

 $\odot$ 

Pa

0.05

0.25

0.2

0.01

0.07

0.15



Step 3: Normalize



 $Z = \sum_{q} P(Q, e_1 \cdots e_k)$  $P(Q|e_1 \cdots e_k) = \frac{1}{Z} P(Q, e_1 \cdots e_k)$ 

#### Inference by Enumeration in Bayes' Net

- Given unlimited time, inference in BNs is easy
- Reminder of inference by enumeration by example:

 $P(B \mid +j,+m) \propto_B P(B,+j,+m)$ 

e,a

$$=\sum_{e,a} P(B,e,a,+j,+m)$$



= P(B)P(+e)P(+a|B,+e)P(+j|+a)P(+m|+a) + P(B)P(+e)P(-a|B,+e)P(+j|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)P(+m|-a)PP(B)P(-e)P(+a|B,-e)P(+j|+a)P(+m|+a) + P(B)P(-e)P(-a|B,-e)P(+j|-a)P(+m|-a)P(+m|-a)P(-a|B,-e)P(+j|-a)P(-a|B,-e)P(+j|-a)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,-e)P(-a|B,

### **Factor Summary**

- In general, when we write  $P(Y_1 \dots Y_N | X_1 \dots X_M)$ 
  - It is a "factor," a multi-dimensional array
  - Its values are  $P(y_1 \dots y_N | x_1 \dots x_M)$
  - Any assigned (=lower-case) X or Y is a dimension selected from the array



# Inference by Enumeration: Procedural Outline

- Track objects called factors
- Initial factors are local CPTs (one per node)

P(R)			
+r	0.1		
-r	0.9		



P(L T)				
	+t	+	0.3	
	+t	-	0.7	
	-t	+	0.1	
	-t	-1	0.9	

Any known values are selected

• E.g. if we know

, the initial factors are

0.9

$$L = +\ell$$





#### **Operation 1: Join Factors**

- First basic operation: joining factors
- Combining factors:
  - Just like a database join
  - Get all factors over the joining variable
  - Build a new factor over the union of the variables involved
- Example: Join on R





 $\forall r, t : P(r, t) = P(r) \cdot P(t|r)$ 

## Marginalizing Early! (aka VE)



#### **General Variable Elimination**

• Query: 
$$P(Q|E_1 = e_1, \dots E_k = e_k)$$

- Start with initial factors:
  - Local CPTs (but instantiated by evidence)
- While there are still hidden variables (not Q or evidence):
  - Pick a hidden variable H<sub>i</sub>
  - Join all factors mentioning H<sub>i</sub>
  - Eliminate (sum out) H<sub>i</sub>
- Join all remaining factors and normalize







# Sampling

- Sampling is a lot like repeated simulation
  - Predicting the weather, basketball games, ...
- Basic idea
  - Draw N samples from a sampling distribution S
  - Compute an approximate posterior probability
  - Show this converges to the true probability P

- Why sample?
  - Learning: get samples from a distribution you don't know
  - Inference: getting a sample is faster than computing the right answer (e.g. with variable elimination)



# Sampling

#### Sampling from given distribution

- Step 1: Get sample *u* from uniform distribution over [0, 1)
  - E.g. random() in python
- Step 2: Convert this sample *u* into an outcome for the given distribution by having each target outcome associated with a sub-interval of [0,1) with sub-interval size equal to probability of the outcome

#### Example



 $\begin{array}{l} 0 \leq u < 0.6, \rightarrow C = red \\ 0.6 \leq u < 0.7, \rightarrow C = green \\ 0.7 \leq u < 1, \rightarrow C = blue \end{array}$ 

- If random() returns u = 0.83, then our sample is C = blue
- E.g, after sampling 8 times:



# **Prior Sampling**

- For i = 1, 2, ..., n
  - Sample x<sub>i</sub> from P(X<sub>i</sub> | Parents(X<sub>i</sub>))
- Return (x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>)



## **Prior Sampling**



# Example

- We'll get a bunch of samples from the BN:
  - +c, -s, +r, +w
  - +c, +s, +r, +w
  - -c, +s, +r, -w
  - +c, -s, +r, +w
  - -c, -s, -r, +w
- If we want to know P(W)
  - We have counts <+w:4, -w:1>
  - Normalize to get P(W) = <+w:0.8, -w:0.2>
  - This will get closer to the true distribution with more samples
  - Can estimate anything else, too
  - What about P(C | +w)? P(C | +r, +w)? P(C | -r, -w)?
  - Fast: can use fewer samples if less time (what's the drawback?)



# **Reflections on Prior Sampling**

#### Pros:

- Much simpler than enumeration or variable elimination: We only ever need samples (scalar values) for each variable, not probabilities.
- Therefore we only ever need k rows from one CPT table to sample a variable X that has k possible values. No potentially exponential increase with number of variables.
- Therefore it doesn't matter as much how sparse the graph is: we still only need k rows from each CPT table, regardless of how many rows (how many parents) it has.

Cons:

- So far we can't deal with evidence.
- We don't get exact values, and its expensive to get accurate estimates of small probabilities. E.g. estimating a 0.001 probability with 1% relative error requires around 10<sup>7</sup> samples.

## **Rejection Sampling**

#### Let's say we want P(C)

- No point keeping all samples around
- Just tally counts of C as we go

#### Let's say we want P(C | +s)

- Same thing: tally C outcomes, but ignore (reject) samples which don't have S=+s
- This is called rejection sampling
- It is also consistent for conditional probabilities (i.e., correct in the limit)



+C, -S, +r, +W +C, +S, +r, +W -C, +S, +r, -W +C, -S, +r, +W -C, -S, -r, +W

## **Rejection Sampling**

- Input: evidence instantiation
- For i = 1, 2, ..., n
  - Sample x<sub>i</sub> from P(X<sub>i</sub> | Parents(X<sub>i</sub>))
  - If x<sub>i</sub> not consistent with evidence
    - Reject: return no sample is generated in this cycle

Return (x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>)



## **Reflections on Rejection Sampling**

Pros:

- Inherits all the pros of prior sampling.
- Now we can deal with evidence.

Cons:

- Dealing with evidence can be very costly. Our rejection rate is 1- the marginal probability of the evidence, and getting N good samples requires taking N/p samples overall, where p is the marginal probability of the evidence.
- We still don't get exact values, and its still expensive to get accurate estimates of small probabilities. E.g. estimating a 0.001 probability with 1% relative error requires around 10<sup>7</sup> samples, multiplied by 1/marginal probability of the evidence.

#### Likelihood Weighting



# Likelihood Weighting



#### Likelihood Weighting Estimates

• We use the weights to estimate probabilities:

$$p(x+) \approx \frac{\sum_{s \in \{s \mid x+(s)\}} w_s}{\sum_s w_s}$$

Where  $w_s$  is the weight of a sample *s*, and x + (s) is true if X = x + in sample *s*.

# Likelihood Weighting

- Likelihood weighting is good
  - We have taken evidence into account as we generate the sample
  - E.g. here, W's value will get picked based on the evidence values of S, R
  - More of our samples will reflect the state of the world suggested by the evidence

- Likelihood weighting doesn't solve all our problems
  - Evidence influences the choice of downstream variables, but not upstream ones (C isn't more likely to get a value matching the evidence)
- We would like to consider evidence when we sample every variable (leads to Gibbs sampling)



# **Gibbs Sampling**

- Procedure: keep track of a full instantiation x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>. Start with an arbitrary instantiation consistent with the evidence. Sample one variable at a time, conditioned on all the rest, but keep evidence fixed. Keep repeating this for a long time.
- Property: in the limit of repeating this infinitely many times the resulting samples come from the correct distribution (i.e. conditioned on evidence).
- *Rationale*: both upstream and downstream variables condition on evidence.
- In contrast: likelihood weighting only conditions on upstream evidence, and hence weights obtained in likelihood weighting can sometimes be very small. Sum of weights over all samples is indicative of how many "effective" samples were obtained, so we want high weight.

# Gibbs Sampling: Conditioning Variables

 No! Remember that the node equations for children depend on their other parents. The complete set is called the Markov Blanket of the node, and includes the other parents of the sampled node's children (orange):



# Gibbs Sampling: Node probability

- To sample, we first compute a factor that includes all node CPTs that depend on A.
- Then we normalize it to get a conditional probability for A.
- Finally we sample to get a new value for A.



# Gibbs Sampling Example with Evidence

Step 1: Fix evidence



- Step 2: Initialize other variables
  - Randomly



Steps 3: Repeat

■ R = +r

- Choose a non-evidence variable X
- Resample X from P( X | all other variables)



Eventually: Compute probabilities with counts over many samples, e.g. p(S|+r), p(W|+r), p(C, W|+r)

# **Reflections on Gibbs Sampling**

Pros:

- Similar to other sampling methods: Only needs k rows from each CPT table for a variable X with k values.
- It also doesn't matter how sparse the graph is, in fact Gibbs sampling typically converges faster on denser graphs, because it *mixes* faster.
- Samples are unweighted, and come from the exact posterior probability conditioned on the evidence (eventually). So estimates can be fairly fast (modulo mixing).

Cons:

- There is a "warm-up" period for the sampler to reach the final distribution.
- Because samples are correlated, need more of them to get estimates at a given accuracy compared to other sampling methods.
- Both of the above depend on "mixing time," for which smaller is better.

There is much theory and many techniques to improve Gibbs sampling.

#### **Decision Networks**



#### **Decision Networks**


#### **Decision Networks**

#### New node types:

- Chance nodes (circular or oval, just like BNs)
- Actions (rectangles, like actions in MDPs)

 Utility node (diamond, like rewards in MDPs)



- For each state B and action A, we compute expected utility by averaging over C, conditioned on B.
- Then we take the maximum over A for each B.
- Finally we can average over B to get the MEU.



В	Α	С	U(A,C)
+b	+a	+C	5
+b	+a	-C	3.2
+b	-a	+c	-1
+b	-a	-C	-4
•••	•••	•••	

	В	А	U(A,B)
	+b	+a	4.1
	+b	-a	-2.3
	-b	+a	0.4
	-b	-a	-1.5

- For each state B and action A, we compute expected utility by averaging over C, conditioned on B.
- Then we take the maximum over A for each B.
- Finally we can average over B to get the MEU.







	В	MEU(B)	
	+b	4.1	
	-b	0.4	

Choose best A | B

- For each state B and action A, we compute expected utility by averaging over C, conditioned on B.
- Then we take the maximum over A for each B.
- Finally we can average over B to get the MEU.





- For each state B and action A, we compute expected utility by averaging over C, conditioned on B.
- Then we take the maximum over A for each B.
- Finally we can average over B to get the MEU.





### Value of Information

 Value of Information is the difference in MEU between networks with different action conditioning (information).



## Value of Information/VPI

 Value of Information is the difference in MEU between networks with different action conditioning (information).



# **VPI** Properties

VPI(E'|e) is value of knowing E' given evidence e.
Nonnegative:

 $\forall E', e : \mathsf{VPI}(E'|e) \ge 0$ 

- Nonadditive (think of observing  $E_j$  twice)  $VPI(E_j, E_k | e) \neq VPI(E_j | e) + VPI(E_k | e)$
- Order-independent

 $VPI(E_j, E_k|e) = VPI(E_j|e) + VPI(E_k|e, E_j)$  $= VPI(E_k|e) + VPI(E_j|e, E_k)$ 







### POMDPs

#### MDPs have:

- States S
- Actions A
- Transition function P(s'|s,a) (or T(s,a,s'))
- Rewards R(s,a,s')

#### POMDPs add:

- Observations O
- Observation function P(o|s) (or O(s,o))
- POMDPs are MDPs over belief states b (distributions over S)
- We'll be able to say more in a few lectures



#### CS 188: Artificial Intelligence

#### Hidden Markov Models



#### University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

#### Hidden Markov Models

Markov chains OK for games, weak for real robots

- Need observations to update your beliefs
- Hidden Markov models (HMMs)
  - Underlying Markov chain over states X
  - You observe outputs (effects) at each time step





#### Example: Weather HMM







#### An HMM is defined by:

- Initial distribution:
- Transitions:
- Emissions:

 $P(X_1)$   $P(X_t \mid X_{t-1})$   $P(E_t \mid X_t)$ 

Transitions			
R <sub>t-1</sub>	R <sub>t</sub>	$P(R_{t}   R_{t\text{-}1})$	R <sub>t</sub>
+r	+r	0.7	+r
+r	-r	0.3	+r
-r	+r	0.3	-r
-r	-r	0.7	-r

_		
Lm		Inne
	155	

R <sub>t</sub>	Ut	P(U <sub>t</sub>   R <sub>t</sub> )
+r	+u	0.9
+r	-u	0.1
-r	+u	0.2
-r	-u	0.8

#### **Conditional Independence**

- HMMs have two important independence properties:
  - Markov hidden process: future depends on past via the present
  - Current observation independent of all else given current state



- Does this mean that evidence variables are guaranteed to be independent?
  - No, they are correlated by the hidden state

#### Inference: Base Cases



**Passage of Time:** 

#### **Observation:**



#### Two Steps: Passage of Time + Observation



#### **Particle Filtering**



# **Particle Filtering**

- Filtering: approximate solution
- Sometimes |X| is too big to use exact inference
  - |X| may be too big to even store B(X)
  - E.g. X is continuous
- Solution: approximate inference
  - Track samples of X, not all values
  - Samples are called particles
  - Typically, there are multiple samples per time step
  - Particles do not interact with each other, and computing time per step is linear in the number of samples
  - But: number needed may be large
  - In memory: list of particles, not states
- This is how robot localization works in practice
- Particle is just new name for sample

0.0	0.1	0.0
0.0	0.0	0.2
0.0	0.2	0.5





 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model



 Each particle is moved by sampling its next position from the transition model

- This is like prior sampling samples' frequencies reflect the transition probabilities
- Here, most samples move clockwise, but some move in another direction or stay in place
- This captures the passage of time
  - If enough samples, close to exact values before and after (consistent)



### Particle Filtering: Observe

#### Slightly trickier:

- Don't sample observation, fix it
- Similar to likelihood weighting, downweight samples based on the evidence

w(x) = P(e|x) $B(X) \propto P(e|X)B'(X)$ 

 As before, the probabilities don't sum to one, since all have been down-weighted (in fact they now sum to (N times) an approximation of P(e))



# Recall: Sampling from a Set

#### Sampling from given distribution

- Step 1: Get sample *u* from uniform distribution over [0, 1)
  - E.g. random() in python
- Step 2: Convert this sample *u* into an outcome for the given distribution by having each target outcome associated with a sub-interval of [0,1) with sub-interval size equal to probability of the outcome

#### Example

 $\begin{array}{l} 0 \leq u < 0.6, \rightarrow C = red \\ 0.6 \leq u < 0.7, \rightarrow C = green \\ 0.7 \leq u < 1, \rightarrow C = blue \end{array}$ 

- If random() returns u = 0.83, then our sample is C = blue
- E.g, after sampling 8 times:



### Particle Filtering: Resample

- Rather than tracking weighted samples, we resample
- N times, we choose from our weighted sample distribution (i.e. draw with replacement)
- This is equivalent to renormalizing the distribution
- Now the update is complete for this time step, continue with the next one



(3,2)

(2,2)(3,2)

(2,3)

(3,3)

(3,2)

(1,3)(2,3)(3,2)(3,2) 3

•

Ο

3

### **Recap: Particle Filtering**

#### Particles: track samples of states rather than an explicit distribution



[Demos: ghostbusters particle filtering (L15D3,4,5)]

#### **Dynamic Bayes Nets**



# Dynamic Bayes Nets (DBNs)

- We want to track multiple variables over time, using multiple sources of evidence
- Idea: Repeat a fixed Bayes net structure at each time
- Variables from time t can condition on those from t-1



Dynamic Bayes nets are a generalization of HMMs



#### Exact Inference in DBNs

- Variable elimination applies to dynamic Bayes nets
- Procedure: "unroll" the network for T time steps, then eliminate variables until  $P(X_T | e_{1:T})$  is computed



 Online belief updates. Eliminate all variables from the previous time step; store factors for current time only

#### **DBN** Particle Filters

- A particle is a complete sample for a time step
- Initialize: Generate prior samples for the t=1 Bayes net
  - Example particle:  $G_1^a = (3,3) G_1^b = (5,3)$
- Elapse time: Sample a successor for each particle
  - Example successor:  $G_2^a = (2,3) G_2^b = (6,3)$
- Observe: Weight each <u>entire</u> sample by the likelihood of the evidence conditioned on the sample
   Likelihood: P(E<sub>1</sub><sup>a</sup> | G<sub>1</sub><sup>a</sup>) \* P(E<sub>1</sub><sup>b</sup> | G<sub>1</sub><sup>b</sup>)
- **Resample:** Select prior samples (tuples of values) in proportion to their likelihood
# CS 188: Artificial Intelligence

Search



### University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley (ai.berkeley.edu).]

### State Space Graphs vs. Search Trees



Each NODE in in the search tree is an entire PATH in the state space graph.

We construct both on demand – and we construct as little as possible.



## Depth-First Search



## Depth-First Search

Strategy: expand a deepest node first

Implementation: Fringe is a LIFO stack





# Depth-First Search (DFS) Properties

- What nodes DFS expand?
  - Some left prefix of the tree.
  - Could process the whole tree!
  - If m is finite, takes time O(b<sup>m</sup>)
- How much space does the fringe take?
  - Only has siblings on path to root, so O(bm)
- Is it complete?
  - m could be infinite, so only if we prevent that
- Is it optimal?
  - No, it finds the "leftmost" solution, regardless of depth or cost



### **Breadth-First Search**



### **Breadth-First Search**

Strategy: expand a shallowest node first

Implementation: Fringe is a FIFO queue





# Breadth-First Search (BFS) Properties

### What nodes does BFS expand?

- Processes all nodes above shallowest solution
- Let depth of shallowest solution be s
- Search takes time O(b<sup>s</sup>)
- How much space does the fringe take?
  - Has roughly the last tier, so O(b<sup>s</sup>)
- Is it complete?
  - s must be finite if a solution exists, so yes!
- Is it optimal?
  - Only if costs are all 1 (more on costs later)



### **Uniform Cost Search**



## **Uniform Cost Search**

Strategy: expand a cheapest node first:

*Fringe is a priority queue (priority: cumulative cost)* 





# Uniform Cost Search (UCS) Properties

- What nodes does UCS expand?
  - Processes all nodes with cost less than cheapest solution!
  - If that solution costs  $C^*$  and arcs cost at least  $\epsilon$  , then the "effective depth" is roughly  $C^*\!/\!\epsilon$
  - Takes time O(b<sup>C\*/ε</sup>) (exponential in effective depth)
- How much space does the fringe take?
  - Has roughly the last tier, so O(b<sup>C\*/ε</sup>)
- Is it complete?
  - Assuming best solution has a finite cost and minimum arc cost is positive, yes!
- Is it optimal?
  - Yes! (Proof next lecture via A\*)



## CS 188: Artificial Intelligence

### Informed Search



Spring 2025

University of California, Berkeley

# Greedy Search



# Greedy Search

- Strategy: expand a node that you think is closest to a goal state
  - Heuristic: estimate of distance to nearest goal for each state

- A common case:
  - Best-first takes you straight to the (wrong) goal

Worst-case: like a badly-guided DFS





[Demo: contours greedy empty (L3D1)] [Demo: contours greedy pacman small maze (L3D4)]

## A\* Search



# Combining UCS and Greedy

- Uniform-cost orders by path cost, or backward cost g(n)
- Greedy orders by goal proximity, or *forward cost* h(n)



A\* Search orders by the sum: f(n) = g(n) + h(n)

## Properties of A\*



## **Admissible Heuristics**

• A heuristic *h* is *admissible* (optimistic) if:

 $0 \leq h(n) \leq h^*(n)$ 

where  $h^*(n)$  is the true cost to a nearest goal

Examples:





 Coming up with admissible heuristics is most of what's involved in using A\* in practice.

### Comparison



Greedy

### **Uniform Cost**



# A\*: Summary

- A\* uses both backward costs and (estimates of) forward costs
- A\* is optimal with admissible / consistent heuristics
- Heuristic design is key: often use relaxed problems



## CS 188: Artificial Intelligence

### **Constraint Satisfaction Problems**





John Canny, Oliver Grillmeyer University of California, Berkeley

# Example: Map Coloring

- Variables: WA, NT, Q, NSW, V, SA, T
- Domains:  $D = \{red, green, blue\}$
- Constraints: adjacent regions must have different colors

Implicit: WA  $\neq$  NT

Explicit: (WA, NT) ∈ {(red, green), (red, blue), ...}
Solutions are assignments satisfying all constraints, e.g.:

{WA=red, NT=green, Q=red, NSW=green, V=red, SA=blue, T=green}





### **Constraint Graphs**



# Varieties of Constraints

#### Varieties of Constraints

Unary constraints involve a single variable (equivalent to reducing domains), e.g.:

### $SA \neq green$

Binary constraints involve pairs of variables, e.g.:

 $SA \neq WA$ 

- Higher-order constraints involve 3 or more variables:
   e.g., cryptarithmetic column constraints
- Preferences (soft constraints):
  - E.g., red is better than green
  - Often representable by a cost for each variable assignment
  - Gives constrained optimization problems
  - (We'll ignore these until we get to Bayes' nets)



# **Backtracking Search**

Backtracking search is the basic uninformed algorithm for solving CSPs

#### Idea 1: One variable at a time

- Variable assignments are commutative, so fix ordering
- I.e., [WA = red then NT = green] same as [NT = green then WA = red]
- Only need to consider assignments to a single variable at each step

### Idea 2: Check constraints as you go

- I.e. consider only values which do not conflict with previous assignments
- Might have to do some computation to check the constraints
- "Incremental goal test"
- Depth-first search with these two improvements is called *backtracking search* (not the best name)
- Can solve n-queens for n ≈ 25



# Backtracking Example



# **Improving Backtracking**

- General-purpose ideas give huge gains in speed
- Ordering:
  - Which variable should be assigned next?
  - In what order should its values be tried?
- Filtering: Can we detect inevitable failure early?
- Structure: Can we exploit the problem structure?



# Filtering: Forward Checking

- Filtering: Keep track of domains for unassigned variables and cross off bad options
- Forward checking: Cross off values that violate a constraint when added to the existing assignment



### Video of Demo Coloring – Backtracking with Forward Checking

### **Filtering: Forward Checking**

- Filtering: Keep track of domains for unassigned variables and cross off bad options
- Forward checking: Cross off values that violate a constraint when added to the existing assignment



# CS 188: Artificial Intelligence

**Constraint Satisfaction Problems II** 

Instructors: John Canny and Oliver Grillmeyer

University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

# Consistency of A Single Arc

■ An arc X → Y is consistent iff for every x in the tail there is some y in the head which could be assigned without violating a constraint



• Forward checking: Enforcing consistency of arcs pointing to each new assignment

# Arc Consistency of an Entire CSP

• A simple form of propagation makes sure all arcs are consistent:



- Important: If X loses a value, neighbors of X need to be rechecked!
- Arc consistency detects failure earlier than forward checking
- Can be run as a preprocessor or after each assignment
- What's the downside of enforcing arc consistency?

*Remember: Delete from the tail!* 

# **Ordering: Minimum Remaining Values**

- Variable Ordering: Minimum remaining values (MRV):
  - Choose the variable with the fewest legal left values in its domain



- Why min rather than max?
- Also called "most constrained variable"
- "Fail-fast" ordering



# Ordering: Least Constraining Value

### Value Ordering: Least Constraining Value

- Given a choice of variable, choose the *least* constraining value
- I.e., the one that rules out the fewest values in the remaining variables
- Note that it may take some computation to determine this! (E.g., rerunning filtering)
- Why least rather than most?
- Combining these ordering ideas makes 1000 queens feasible





### **Nearly Tree-Structured CSPs**



Conditioning: instantiate a variable, prune its neighbors' domains

Cutset conditioning: instantiate (in all ways) a set of variables such that the remaining constraint graph is a tree

Cutset size c gives runtime O( (d<sup>c</sup>) (n-c) d<sup>2</sup> ), very fast for small c

## **Cutset Conditioning**


# Summary: CSPs

#### CSPs are a special kind of search problem:

States are partial assignments Goal test defined by constraints

Basic solution: backtracking search

Speed-ups:

Ordering Filtering Structure



Iterative min-conflicts is often effective in practice

# Hill Climbing

### Simple, general idea:

Start wherever

Repeat: move to the best neighboring state If no neighbors better than current, quit

### What's bad about this approach?

Complete? Optimal?

What's good about it?



# Hill Climbing Diagram



## Beam Search

Like greedy hillclimbing search, but keep K states at all times:





Greedy Search

**Beam Search** 

Variables: beam size, encourage diversity? The best choice in MANY practical settings Complete? Optimal? Why do we still need optimal methods?

### CS 188: Artificial Intelligence Game Trees: Adversarial Search



[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley (ai.berkeley.edu). [Updated slides from: Stuart Russell and Dawn Song]

# Zero-Sum Games



- Zero-Sum Games
  - Agents have opposite utilities (values on outcomes)
  - Pure competition:
    - One *maximizes*, the other *minimizes*



- General-Sum Games
  - Agents have *independent* utilities (values on outcomes)
  - Cooperation, indifference, competition, shifting alliances, and more are all possible
- Team Games
  - Common payoff for all team members

## **Adversarial Game Trees**



# **Minimax Values**



V(s) =known

## Minimax Example



# **Minimax Pruning**



The order of generation matters: more pruning is possible if good moves come first

# **Evaluation Functions**

Evaluation functions score non-terminals in depth-limited search



White slightly better

Black winning

- Ideal function: returns the actual minimax value of the position
- In practice: typically weighted linear sum of features:

$$Eval(s) = w_1 f_1(s) + w_2 f_2(s) + \ldots + w_n f_n(s)$$

- E.g.  $f_1(s) = (num white queens num black queens), etc.$
- Or a more complex nonlinear function (e.g., NN) trained by self-play RL

# **Depth Matters**

- Evaluation functions are always imperfect
- The deeper in the tree the evaluation function is buried, the less the quality of the evaluation function matters
- An important example of the tradeoff between complexity of features and complexity of computation





# Summary

- Games are decision problems with multiple agents
  - Huge variety of issues and phenomena depending on details of interactions and payoffs
- For zero-sum games, optimal decisions defined by minimax
  - Implementable as a depth-first traversal of the game tree
  - Time complexity O(b<sup>m</sup>), space complexity O(bm)
- Alpha-beta pruning
  - Preserves optimal choice at the root
  - Alpha/beta values keep track of best obtainable values from any max/min nodes on path from root to current node
  - Time complexity drops to  $O(b^{m/2})$  with ideal node ordering
- Exact solution is impossible even for "small" games like chess



[These slides were created by Dan Klein, Pieter Abbeel for CS188 Intro to AI at UC Berkeley (ai.berkeley.edu).]

## Worst-Case vs. Average Case



Idea: Uncertain outcomes controlled by chance, not an adversary!

## Expectimax Pseudocode





v = (1/2) (8) + (1/3) (24) + (1/6) (-12) = 10

# **Reminder:** Probabilities

- A random variable represents an event whose outcome is unknown
- A probability distribution is an assignment of weights to outcomes
- Example: Traffic on freeway
  - Random variable: T = whether there's traffic
  - Outcomes: T in {none, light, heavy}
  - Distribution: P(T=none) = 0.25, P(T=light) = 0.50, P(T=heavy) = 0.25
- Some laws of probability (more later):
  - Probabilities are always non-negative
  - Probabilities over all possible outcomes sum to one
- As we get more evidence, probabilities may change:
  - P(T=heavy) = 0.25, P(T=heavy | Hour=8am) = 0.60
  - We'll talk about methods for reasoning and updating probabilities later



# **Reminder: Expectations**

- The expected value of a function of a random variable is the average, weighted by the probability distribution over outcomes
- Example: How long to get to the airport?





# Mixed Layer Types

- E.g. Backgammon
- Expectiminimax
  - Environment is an extra "random agent" player that moves after each min/max agent
  - Each node computes the appropriate combination of its children



# **Multi-Agent Utilities**

What if the game is not zero-sum, or has multiple players?

**1,6,**6

7,1,2

<mark>6,1,</mark>2

7,2,1

<mark>5,1</mark>,7

1,5,2

<mark>5,2</mark>,5

7,7,1

- Generalization of minimax:
  - Terminals have utility tuples
  - Node values are also utility tuples
  - Each player maximizes its own component
  - Can give rise to cooperation and competition dynamically...



# Maximum Expected Utility

- Why should we average utilities? Why not minimax?
- Principle of maximum expected utility:
  - A rational agent should chose the action that maximizes its expected utility, given its knowledge
- Questions:
  - Where do utilities come from?
  - How do we know such utilities even exist?
  - How do we know that averaging even makes sense?
  - What if our behavior (preferences) can't be described by utilities?



# Preferences

- An agent must have preferences among:
  - Prizes: *A*, *B*, etc.
  - Lotteries: situations with uncertain prizes

L = [p, A; (1-p), B]



#### Notation:

- Preference:  $A \succ B$
- Indifference:  $A \sim B$



# **Rational Preferences**

#### The Axioms of Rationality

Orderability  $(A \succ B) \lor (B \succ A) \lor (A \sim B)$ Transitivity  $(A \succ B) \land (B \succ C) \Rightarrow (A \succ C)$ Continuity  $A \succ B \succ C \Rightarrow \exists p \ [p, A; \ 1-p, C] \sim B$ Substitutability  $A \sim B \Rightarrow [p, A; 1-p, C] \sim [p, B; 1-p, C]$ Monotonicity  $A \succ B \Rightarrow$  $(p \ge q \Leftrightarrow [p, A; 1-p, B] \succeq [q, A; 1-q, B])$ 



Theorem: Rational preferences imply behavior describable as maximization of expected utility

# CS 188: Artificial Intelligence Naïve Bayes



#### Instructors: John Canny and Oliver Grillmeyer --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

# **Model-Based Classification**

### Model-based approach

- Build a model (e.g. Bayes' net) where both the output label and input features are random variables
- Instantiate any observed features
- Query for the distribution of the label conditioned on the features

### Challenges

- What structure should the BN have?
- How should we learn its parameters?



# Inference for Naïve Bayes

- Goal: compute posterior distribution over label variable Y
  - Step 1: get joint probability of label and evidence for each label

$$P(Y, f_1 \dots f_n) = \begin{bmatrix} P(y_1, f_1 \dots f_n) \\ P(y_2, f_1 \dots f_n) \\ \vdots \\ P(y_k, f_1 \dots f_n) \end{bmatrix} \longrightarrow \begin{bmatrix} P(y_1) \prod_i P(f_i | y_1) \\ P(y_2) \prod_i P(f_i | y_2) \\ \vdots \\ P(y_k) \prod_i P(f_i | y_k) \end{bmatrix}$$

 $P(Y|f_1 \dots f_n)$ 

- Step 2: sum to get probability of evidence
- Step 3: normalize by dividing Step 1 by Step 2

# Naïve Bayes for Text

#### Bag-of-words Naïve Bayes:

- Features: W<sub>i</sub> is the word at position i
- As before: predict label conditioned on feature variables (spam vs. ham)
- As before: assume features are conditionally independent given label
- New: each W<sub>i</sub> is identically distributed

Word at position *i, not i<sup>th</sup> word in the dictionary!* 

• Generative model: 
$$P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i | Y)$$

- "Tied" distributions and bag-of-words
  - Usually, each variable gets its own conditional probability distribution P(F|Y)
  - In a bag-of-words model
    - Each position is identically distributed
    - All positions share the same conditional probs P(W|Y)
    - Why make this assumption?
  - Called "bag-of-words" because model is insensitive to word order or reordering

# **Example: Spam Filtering**

• Model: 
$$P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i | Y)$$

What are the parameters?



### P(W|spam)

the	:	0.0156
to	:	0.0153
and	:	0.0115
of	:	0.0095
you	:	0.0093
a	:	0.0086
with:		0.0080
from:		0.0075
•••		

#### P(W|ham)

the :	0.0210
to :	0.0133
of :	0.0119
2002:	0.0110
with:	0.0108
from:	0.0107
and :	0.0105
a :	0.0100
•••	

Where do these tables come from?

# Spam Example

Word	P(w spam)	P(w ham)	Tot Spam	Tot Ham
(prior)	0.33333	0.66666	-1.1	-0.4

# Important Concepts

- Data: labeled instances (e.g. emails marked spam/ham)
  - Training set
  - Held out set
  - Test set
- Features: attribute-value pairs which characterize each x
- Experimentation cycle
  - Learn parameters (e.g. model probabilities) on training set
  - (Tune hyperparameters on held-out set)
  - Compute accuracy of test set
  - Very important: never "peek" at the test set!
- Evaluation (many metrics possible, e.g. accuracy)
  - Accuracy: fraction of instances predicted correctly
- Overfitting and generalization
  - Want a classifier which does well on test data
  - Overfitting: fitting the training data very closely, but not generalizing well
  - We'll investigate overfitting and generalization formally in a few lectures



# **Confusion Matrix**



# **Performance Metrics**

### Accuracy

- Number of correct predictions from the entire data set: (TP + TN) / (TP + FP + TN + FN)
- Precision
  - Number of correct positive predictions from total positive predictions: TP / (TP + FP)
- Recall
  - Number of correct positive predictions from the actual positive samples: TP / (TP + FN)
- F-score
  - Harmonic mean of Precision and Recall: 2TP / (2TP + FP + FN)

# Overfitting



# **Parameter Estimation**

- Estimating the distribution of a random variable
- Elicitation: ask a human (why is this hard?)
- Empirically: use training data (learning!)
  - E.g.: for each outcome x, look at the *empirical rate* of that value:

$$P_{\mathsf{ML}}(x) = \frac{\mathsf{count}(x)}{\mathsf{total samples}}$$



b

 $P_{\rm MI}(r) = 2/3$ 

This is the estimate that maximizes the *likelihood of the data*

$$L(x,\theta) = \prod_{i} P_{\theta}(x_i)$$

# Laplace Smoothing

#### Laplace's estimate:

 Pretend you saw every outcome once more than you actually did



$$P_{LAP}(x) = \frac{c(x) + 1}{\sum_{x} [c(x) + 1]}$$
$$= \frac{c(x) + 1}{N + |X|}$$

$$P_{ML}(X) =$$

$$P_{LAP}(X) =$$

 Can derive this estimate with Dirichlet priors (see cs281a)

# Summary

- Bayes rule lets us do diagnostic queries with causal probabilities
- The naïve Bayes assumption takes all features to be independent given the class label
- We can build classifiers out of a naïve Bayes model using training data
- Smoothing estimates is important in real systems
CS 188: Artificial Intelligence Perceptrons and Logistic Regression



Instructors: John Canny & Oliver Grillmeyer —- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

#### **Feature Vectors**



## Linear Classifiers

- Inputs are feature values
- Each feature has a weight
- Sum is the activation



activation<sub>w</sub>(x) = 
$$\sum_{i} w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
  - Positive, output +1
  - Negative, output -1



# Weights

- Binary case: compare features to a weight vector
- Learning: figure out the weight vector from examples



# **Binary Decision Rule**

- In the space of feature vectors
  - Examples are points
  - Any weight vector is a hyperplane
  - One side corresponds to Y=+1
  - Other corresponds to Y=-1





w

BIAS	:	-3
free	:	4
money	•	2
• • •		

## Learning: Binary Perceptron

- Start with weights = 0
- For each training instance:
  - Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \ge 0\\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

- If correct (i.e., y=y\*), no change!
- If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y\* is -1.

$$w = w + y^* \cdot f$$



## **Multiclass Decision Rule**

- If we have multiple classes:
  - A weight vector for each class:

 $w_y$ 

Score (activation) of a class y:

 $w_y \cdot f(x)$ 

Prediction highest score wins

$$y = \arg \max_{y} w_{y} \cdot f(x)$$





Binary = multiclass where the negative class has weight zero

## Learning: Multiclass Perceptron

- Start with all weights = 0
- Pick up training examples one by one
- Predict with current weights

$$y = \arg \max_y w_y \cdot f(x)$$

- If correct, no change!
- If wrong: lower score of wrong answer, raise score of right answer

$$w_y = w_y - f(x)$$
$$w_{y^*} = w_{y^*} + f(x)$$



#### Example: Multiclass Perceptron

"win the vote" "win the election" "win the game"

```
Sample A
          wS . fA = 1; wP . fA = 0; wT . fA = 0
BIAS : 1
win : 1
game : 0
                        wS . fB = -2; wP . fB = 3; wT . fB = 0
              Sample B
vote : 1
              BIAS : 1
the : 1
              win : 1
              game : 0
                                       wS . fC = -2; wP . fC = 3; wT . fC = 0
                             Sample C
              vote : 0
                             BIAS : 1
              the : 1
                             win : 1
                             game : 1
                             vote : 0
                             the : 1
```



# **Properties of Perceptrons**

- Separability: true if some parameters get the training set perfectly correct
- Convergence: if the training is separable, perceptron will eventually converge (binary case)
- Mistake Bound: the maximum number of mistakes (binary case) related to the margin or degree of separability

mistakes 
$$< \frac{k}{\delta^2}$$





Non-Separable



#### Non-Separable Case: Probabilistic Decision



#### How to get probabilistic decisions?

- Perceptron scoring:  $z = w \cdot f(x)$
- If  $z = w \cdot f(x)$  very positive  $\rightarrow$  want probability going to 1
- If  $z = w \cdot f(x)$  very negative  $\rightarrow$  want probability going to 0

Sigmoid function

$$\phi(z) = \frac{1}{1 + e^{-z}}$$



#### Best w?

Maximum likelihood estimation:

$$\max_{w} \quad ll(w) = \max_{w} \quad \sum_{i} \log P(y^{(i)} | x^{(i)}; w)$$
 with:

$$P(y^{(i)} = +1 | x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$
$$P(y^{(i)} = -1 | x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

= Logistic Regression

## **Multiclass Logistic Regression**

- Recall Perceptron:
  - A weight vector for each class:
  - Score (activation) of a class y:
  - Prediction highest score wins

$$w_y \cdot f(x)$$
  
 $y = \arg \max_y w_y \cdot f(x)$ 

f(m)

 $w_y$ 



How to make the scores into probabilities?

$$z_{1}, z_{2}, z_{3} \rightarrow \underbrace{\frac{e^{z_{1}}}{e^{z_{1}} + e^{z_{2}} + e^{z_{3}}}, \frac{e^{z_{2}}}{e^{z_{1}} + e^{z_{2}} + e^{z_{3}}}, \frac{e^{z_{3}}}{e^{z_{1}} + e^{z_{2}} + e^{z_{3}}}, \frac{e^{z_{3}}$$

#### Best w?

Maximum likelihood estimation:

with:

$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)} | x^{(i)}; w)$$

$$P(y^{(i)}|x^{(i)};w) = \frac{e^{w_{y^{(i)}} \cdot f(x^{(i)})}}{\sum_{y} e^{w_{y} \cdot f(x^{(i)})}}$$

#### = Multi-Class Logistic Regression

## CS 188: Artificial Intelligence

## **Optimization and Neural Nets**



Instructors: John Canny and Oliver Grillmeyer --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

## **1-D Optimization**



- Then step in best direction
- Or, evaluate derivative:

$$\frac{\partial g(w_0)}{\partial w} = \lim_{h \to 0} \frac{g(w_0 + h) - g(w_0 - h)}{2h}$$

Tells which direction to step into

## **Gradient Ascent**

- Perform update in uphill direction for each coordinate
- The steeper the slope (i.e. the higher the derivative) the bigger the step for that coordinate
- E.g., consider:  $g(w_1, w_2)$ 
  - Updates:

$$w_1 \leftarrow w_1 + \alpha * \frac{\partial g}{\partial w_1}(w_1, w_2)$$
$$w_2 \leftarrow w_2 + \alpha * \frac{\partial g}{\partial w_2}(w_1, w_2)$$

Updates in vector notation:

$$w \leftarrow w + \alpha * \nabla_w g(w)$$
  
with:  $\nabla_w g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_1}(w) \\ \frac{\partial g}{\partial w_2}(w) \end{bmatrix}$ 

= gradient

## **Optimization Procedure: Gradient Ascent**

• init 
$$w$$
  
• for iter = 1, 2, ...  
 $w \leftarrow w + \alpha * \nabla g(w)$ 

- *α*: learning rate --- hyperparameter that needs to be chosen carefully
- How? Try multiple choices
  - Crude rule of thumb: update changes w about 0.1 1 %

#### Batch Gradient Ascent on the Log Likelihood Objective

$$\max_{w} ll(w) = \max_{w} \sum_{i} \log P(y^{(i)} | x^{(i)}; w)$$

$$g(w)$$

• init 
$$W$$
  
• for iter = 1, 2, ...  
 $w \leftarrow w + \alpha * \sum_{i} \nabla \log P(y^{(i)} | x^{(i)}; w)$ 

#### Stochastic Gradient Ascent on the Log Likelihood Objective

$$\max_{w} ll(w) = \max_{w} \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

**Observation:** once gradient on one training example has been computed, might as well incorporate before computing next one

• init w• for iter = 1, 2, ... • pick random j  $w \leftarrow w + \alpha * \nabla \log P(y^{(j)} | x^{(j)}; w)$ 

#### Mini-Batch Gradient Ascent on the Log Likelihood Objective

$$\max_{w} ll(w) = \max_{w} \sum_{i} \log P(y^{(i)} | x^{(i)}; w)$$

**Observation:** gradient over small set of training examples (=mini-batch) can be computed in parallel, might as well do that instead of a single one

• init w• for iter = 1, 2, ... • pick random subset of training examples J  $w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)} | x^{(j)}; w)$ 

#### **Multi-class Logistic Regression**

#### = special case of neural network



## Deep Neural Network = Also learn the features!



 $z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$ 

g = nonlinear activation function

#### **Common Activation Functions**

Sigmoid Function



g'(z) = g(z)(1 - g(z))

Hyperbolic Tangent



 $g'(z) = 1 - g(z)^2$ 

Rectified Linear Unit (ReLU)



 $g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$ 

## Deep Neural Network: Also Learn the Features!

Training the deep neural network is just like logistic regression:

$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)} | x^{(i)}; w)$$

- Much larger weight vector to learn
- Keep training (adjust weights with gradient ascent) until we meet our performance criteria or validation set performance starts decreasing

# Summary of Key Ideas

 $\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)} | x^{(i)}; w)$ 

- Optimize probability of label given input
- Continuous optimization
  - Gradient ascent:
    - Compute steepest uphill direction = gradient (= just vector of partial derivatives)
    - Take step in the gradient direction
    - Repeat (until held-out data accuracy starts to drop = "early stopping")

#### Deep neural nets

- Last layer = still logistic regression
- Now also many more layers before this last layer
  - = computing the features
  - $\rightarrow$  the features are learned rather than hand-designed
- Universal function approximation theorem
  - If neural net is large enough
  - Then neural net can represent any continuous mapping from input to output with arbitrary accuracy
  - But remember: need to avoid overfitting / memorizing the training data → early stopping!
- Automatic differentiation gives the derivatives efficiently (how? = outside of scope of 188)

## CS 188: Artificial Intelligence

## Neural Nets (wrap-up) and Decision Trees



Instructors: John Canny and Oliver Grillmeyer --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

## **Object Detection**



#### Features and Generalization





Image



## Performance

# ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

# Inductive Learning (Science)

- Simplest form: learn a function from examples
  - A target function: *g*
  - Examples: input-output pairs (x, g(x))
  - E.g. x is an email and g(x) is spam / ham
  - E.g. *x* is a house and *g*(*x*) is its selling price

#### Problem:

- Given a hypothesis space *H*
- Given a training set of examples X<sub>i</sub>
- Find a hypothesis h(x) such that  $h \sim g$
- Includes:
  - Classification (outputs = class labels)
  - Regression (outputs = real numbers)
- How do perceptron and naïve Bayes fit in? (*H*, *h*, *g*, etc.)



## Inductive Learning

• Curve fitting (regression, function approximation):



- Consistency vs. simplicity
- Ockham's razor

## **Decision Trees**



#### Features

- Features, aka attributes
  - Sometimes: TYPE=French
  - Sometimes:  $f_{\text{TYPE=French}}(x) = 1$

Example	Attributes									Target	
Litempre	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait
$X_1$	T	F	F	Т	Some	\$\$\$	F	T	French	0–10	Т
$X_2$	T	F	F	Т	Full	\$	F	F	Thai	30–60	F
$X_3$	F	T	F	F	Some	\$	F	F	Burger	0–10	Т
$X_4$	T	F	Т	Т	Full	\$	F	F	Thai	10–30	Т
$X_5$	T	F	Т	F	Full	\$\$\$	F	T	French	>60	F
$X_6$	F	T	F	Т	Some	\$\$	Т	T	Italian	0–10	Т
$X_7$	F	T	F	F	None	\$	Т	F	Burger	0–10	F
$X_8$	F	F	F	Т	Some	\$\$	Т	T	Thai	0–10	Т
$X_9$	F	T	Т	F	Full	\$	Т	F	Burger	>60	F
$X_{10}$	T	T	T	Т	Full	\$\$\$	F	T	Italian	10–30	F
$X_{11}$	F	F	F	F	None	\$	F	F	Thai	0–10	F
$X_{12}$	T	T	Т	Т	Full	\$	F	F	Burger	30–60	Т

### **Decision Trees**

- Compact representation of a function:
  - Truth table
  - Conditional probability table
  - Regression values
- True function
  - Realizable: in *H*


# Choosing an Attribute

 Idea: a good attribute splits the examples into subsets that are (ideally) "all positive" or "all negative"



 So: we need a measure of how "good" a split is, even if the results aren't perfectly separated out

# Entropy

- General answer: if prior is  $\langle p_1, ..., p_n \rangle$ :
  - Information is the expected code length

$$H(\langle p_1, \dots, p_n \rangle) = E_p \log_2 1/p_i$$
$$= \sum_{i=1}^n -p_i \log_2 p_i$$

- Also called the entropy of the distribution
  - More uniform = higher entropy
  - More values = higher entropy
  - More peaked = lower entropy
  - Rare values almost "don't count"



# Information Gain

- Back to decision trees!
- For each split, compare entropy before and after
  - Difference is the information gain
  - Problem: there's more than one distribution after split!





## Example: Learned Tree

Decision tree learned from these 12 examples:



- Substantially simpler than "true" tree
  - A more complex hypothesis isn't justified by data
- Also: it's reasonable, but wrong

### Example: Miles Per Gallon

mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	low	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europe
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low	medium	low	medium	70to74	asia
bad	4	low	medium	low	low	70to74	asia
bad	8	high	high	high	low	75to78	america
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
bad	8	high	high	high	low	70to74	america
good	8	high	medium	high	high	79to83	america
bad	8	high	high	high	low	75to78	america
good	4	low	low	low	low	79to83	america
bad	6	medium	medium	medium	high	75to78	america
good	4	medium	low	low	low	79to83	america
good	4	low	low	medium	high	79to83	america
bad	8	high	high	high	low	70to74	america
good	4	low	medium	low	medium	75to78	europe
bad	5	medium	medium	medium	medium	75to78	europe

40 Examples

# Find the First Split

 Look at information gain for each attribute

Note that each attribute is correlated with the target!

What do we split on?



### Second Level





# **Reminder: Overfitting**

### Overfitting:

- When you stop modeling the patterns in the training data (which generalize)
- And start modeling the noise (which doesn't)

### We had this before:

- Naïve Bayes: needed to smooth
- Perceptron: early stopping



## Pruning example



# Regularization

- MaxP<sub>CHANCE</sub> is a regularization parameter
- Generally, set it using held-out data (as usual)



## Two Ways of Controlling Overfitting

### Limit the hypothesis space

- E.g. limit the max depth of trees
- Easier to analyze

### Regularize the hypothesis selection

- E.g. chance cutoff
- Skip most of the hypotheses unless data is clear
- Usually done in practice

## Large Language Model Transformers

0



# Large Language Models

- Feature engineering
  - Text tokenization
  - Word embeddings
- Deep neural networks
  - Autoregressive models
  - Self-attention mechanisms
  - Transformer architecture
- Multi-class classification

### Supervised learning

- Self-supervised learning
- Instruction tuning
- Reinforcement learning
  - ... from human feedback (RLHF)
- Policy search
  - Policy gradient methods
- Beam search

### **Text Tokenization**

#### GPT-3.5 & GPT-4 GPT-3 (Legacy)

Many words map to one token, but some don't: indivisible.

Unicode characters like emojis may be split into many tokens containing the underlying bytes: 000000

Sequences of characters commonly found next to each other may be grouped together: 1234567890

Text Token IDs



https://platform.openai.com/tokenizer

### **Text Tokenization**

#### GPT-3.5 & GPT-4 GPT-3 (Legacy)

[8607, 4339, 2472, 311, 832, 4037, 11, 719, 1063, 1541, 956, 25, 3687, 23936, 382, 35020, 5885, 1093, 100166, 1253, 387, 6859, 1139, 1690, 11460, 8649, 279, 16940, 5943, 25, 11410, 97, 248, 9468, 237, 122, 271, 1542, 45045, 315, 5885, 17037, 1766, 1828, 311, 1855, 1023, 1253, 387, 41141, 3871, 25, 220, 4513, 10961, 16474, 15]

Text Token IDs



https://platform.openai.com/tokenizer

# Word Embeddings



# What do word embeddings look like?

Features learned in language models:



ig.ft.com/generative-ai

# **Autoregressive Models**



## **Self-Attention Mechanisms**



Instead of conditioning on all input tokens equally...



 Pay more attention to relevant tokens!

### **Self-Attention Mechanisms**



ig.ft.com/generative-ai





## **Multi-Headed Attention**



## **Multi-Headed Attention**

### Head 6: previous word



## **Multi-Headed Attention**

### Head 4: pronoun references

[CLS]	[CLS]	[CLS]	[CLS]	[CLS]	[CLS]
the	the	the	the	the	the
girl	girl	girl	girl	girl	girl
and	and	and	and	and	and
the	the	the	the	the	the
boy	boy	boy	boy	boy	boy
walked	walked	walked	walked	walked	walked
home	home	home	home	home	home
		. //	/ / ·	. ///	
[SEP]	[SEP]	[SEP]	[SEP]	[SEP]	[SEP]
she ${\color{red}{\leftarrow}}$	she	she	she	she	she
took	took	took	took	took	took
his	his	his	his	his	his
hand	hand	hand	hand	hand	hand
in	in	in	in	in	in
hers	hers	hers	hers	hers	hers
[SEP]	[SEP]	[SEP]	[SEP]	[SEP]	[SEP]

#### https://github.com/jessevig/bertviz

## **Transformer Architecture**



## **Transformer Architecture**





# **Pre-Training and Fine-Tuning**

### Pre-Train: train a large model with a lot of data on a selfsupervised task

- Predict next word / patch of image
- Predict missing word / patch of image
- Predict if two images are related (contrastive learning)
- 2 Fine-Tune: continue training the same model on task you care about

# Instruction Tuning

Task 1 = predict next word

(learns to mimic human-written text)

- Query: "What is population of Berkeley?"
- Human-like completion: "This question always fascinated me!"
- Task 2 = generate **helpful** text
  - Query: "What is population of Berkeley?"
  - Helpful completion: "It is 117,145 as of 2021 census."
- Fine-tune on collected examples of helpful human conversations
- Also can use Reinforcement Learning

### Beam Search

