CSE 573: Artificial Intelligence

Hanna Hajishirzi Hidden Markov Models

slides adapted from Dan Klein, Pieter Abbeel ai.berkeley.edu And Dan Weld, Luke Zettelmoyer



Announcements

PS3: Due today
PS4 -> Released
HW2 -> Released

Probability Summary

Conditional probability

$$P(x|y) = \frac{P(x,y)}{P(y)}$$

- Product rule P(x,y) = P(x|y)P(y)
- Chain rule $P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2)\dots$ $= \prod_{i=1}^n P(X_i|X_1, \dots, X_{i-1})$
- X, Y independent if and only if: $\forall x, y : P(x, y) = P(x)P(y)$
- X and Y are conditionally independent given Z if and only if: $X \perp \!\!\!\perp Y | Z$ $\forall x, y, z : P(x, y|z) = P(x|z)P(y|z)$

Recap: Bayes' Net Representation

- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node
 - A collection of distributions over X, one for each combination of parents' values $P(X|a_1 \dots a_n)$
- Bayes' nets implicitly encode joint distributions
 - As a product of local conditional distributions
 - To see what probability a BN gives to a full assignment, multiply all the relevant and $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | parents(X_i))$





Quiz: Bayes' Rule

Ρ

0.1

0.9

0.7

0.3



• What is $P(W \mid dry)$?

Quiz: Bayes' Rule

Ρ

0.1

0.9

0.7

0.3



• What is $P(W \mid dry)$?

 $P(sun|dry) \sim P(dry|sun)P(sun) = .9^*.8 = .72$ $P(rain|dry) \sim P(dry|rain)P(rain) = .3^*.2 = .06$ P(sun|dry)=12/13P(rain|dry)=1/13

Ghostbusters, Revisited

• Let's say we have two distributions:

- **Prior distribution** over ghost location: P(G)
 - Let's say this is uniform
- Sensor reading model: $P(R \mid G)$
 - Given: we know what our sensors do
 - \circ R = reading color measured at (1,1)
 - E.g. P(R = yellow | G=(1,1)) = 0.1
- We can calculate the posterior distribution P(G | r) over ghost locations given a reading using Bayes' rule: $P(g|r) \propto P(r|g)P(g)$

0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11



[Demo: Ghostbuster – with probability (L12D2)]

7

Reasoning over Time or Space

• Often, we want to reason about a sequence of observations

- \circ Speech recognition
- Robot localization
- o User attention
- o Medical monitoring
- Need to introduce time (or space) into our models

Markov Models

• Value of X at a given time is called the **state**



- Parameters: called transition probabilities or dynamics, specify how the state evolves over time (also, initial state probabilities)
- Stationarity assumption: transition probabilities the same at all times
- Same as MDP transition model, but no choice of action
- A (growable) BN: We can always use generic BN reasoning on it if we truncate the chain at a fixed length

Markov Assumption: Conditional Independence



- Basic conditional independence:
 - Past and future independent given the present
 - Each time step only depends on the previous
 - This is called the (first order) Markov property

Example Markov Chain: Weather

• States: X = {rain, sun}

Initial distribution: 1.0 sun



■ CPT P(X_t | X_{t-1}):

X _{t-1}	X _t	P(X _t X _{t-1})
sun	sun	0.9
sun	rain	0.1
rain	sun	0.3
rain	rain	0.7

Two new ways of representing the same CPT



11

Bayes Nets -- Independence



Bayes Net Chain Rule

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | parents(X_i))$$
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1 \dots x_{i-1})$$

Markov Models (Markov Chains)

$$(X_1) \rightarrow (X_2) \rightarrow (X_3) \rightarrow (X_4) \cdots \rightarrow (X_N)$$

- A Markov model defines
 - a joint probability distribution:
 - $P(X_1, X_2, X_3, X_4) =$
- More generally:

 $P(X_1, X_2, \dots, X_T) = P(X_1)P(X_2|X_1)P(X_3|X_2)\dots P(X_T|X_{T-1})$

$$P(X_1, \dots, X_n) = P(X_1) \prod_{t=2}^{N} P(X_t | X_{t-1})$$
 Why?

 Chain Rule, Indep. Assumption?

- One common inference problem:
 - Compute marginals $P(X_t)$ for all time steps t

Example Markov Chain: Weather

• Initial distribution: 1.0 sun



• What is the probability distribution after one step? $P(X_2 = sun) = \sum_{x_1} P(x_1, X_2 = sun) = \sum_{x_1} P(X_2 = sun | x_1) P(x_1)$

 $P(X_2 = \text{sun}) = P(X_2 = \text{sun}|X_1 = \text{sun})P(X_1 = \text{sun}) + P(X_2 = \text{sun}|X_1 = \text{rain})P(X_1 = \text{rain}) + 0.9 \cdot 1.0 + 0.3 \cdot 0.0 = 0.9$

Mini-Forward Algorithm

• Question: What's P(X) on some day t?

$$(X_1) \rightarrow (X_2) \rightarrow (X_3) \rightarrow (X_4) - - - \rightarrow$$

$$P(x_1) = known$$

$$P(x_t) = \sum_{x_{t-1}} P(x_{t-1}, x_t)$$

=
$$\sum_{x_{t-1}} P(x_t \mid x_{t-1}) P(x_{t-1})$$

Forward simulation



Example Run of Mini-Forward Algorithm

From initial observation of sun

$$\begin{pmatrix} 1.0 \\ 0.0 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix} \begin{pmatrix} 0.84 \\ 0.16 \end{pmatrix} \begin{pmatrix} 0.804 \\ 0.196 \end{pmatrix} \longrightarrow \begin{pmatrix} 0.75 \\ 0.25 \end{pmatrix}$$

$$P(X_1) P(X_2) P(X_3) P(X_4) P(X_{\omega})$$

• From initial observation of rain $\begin{cases}
0.0 \\
1.0 \\
P(X_1)
\end{cases}
\begin{cases}
0.3 \\
0.7 \\
P(X_2)
\end{cases}
\begin{cases}
0.48 \\
0.52 \\
P(X_3)
\end{cases}
\begin{cases}
0.588 \\
0.412 \\
P(X_4)
\end{cases}
\longrightarrow
\begin{cases}
0.75 \\
0.25 \\
P(X_{\infty})
\end{cases}$

• From yet another initial distribution $P(X_1)$: $\begin{pmatrix} p \\ 1-p \\ P(X_1) \end{pmatrix}$... $P(X_m)$

16 [Demo: L13D1,2,3]

Pac-man Markov Chain

Pac-man knows the ghost's initial position, but gets no observations!



Video of Demo Ghostbusters Circular Dynamics

ile Edit Naviga	ate Search Project Run Window Help	
	** • • • • • • • • • • • • • • • • • •	E Pydev for team
Console <terminated></terminated>	83 > basic	
<		

Stationary Distributions

• For most chains:

- Influence of the initial distribution gets less and less over time.
- The distribution we end up in is independent of the initial distribution

Stationary distribution:

- The distribution we end up with is called the stationary distribution P_∞ of the chain
- It satisfies

$$P_{\infty}(X) = P_{\infty+1}(X) = \sum_{x} P(X|x)P_{\infty}(x)$$



Example: Stationary Distributions

 \circ Question: What's P(X) at time t = infinity?

$$X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow X_4 \rightarrow - - \rightarrow$$

 $P_{\infty}(sun) = P(sun|sun)P_{\infty}(sun) + P(sun|rain)P_{\infty}(rain)$ $P_{\infty}(rain) = P(rain|sun)P_{\infty}(sun) + P(rain|rain)P_{\infty}(rain)$

 $P_{\infty}(sun) = 0.9P_{\infty}(sun) + 0.3P_{\infty}(rain)$ $P_{\infty}(rain) = 0.1P_{\infty}(sun) + 0.7P_{\infty}(rain)$

 $P_{\infty}(sun) = 3P_{\infty}(rain)$ $P_{\infty}(rain) = 1/3P_{\infty}(sun)$

Also: $P_{\infty}(sun) + P_{\infty}(rain) = 1$





X _{t-1}	X _t	$P(X_t X_{t-1})$
sun	sun	0.9
sun	rain	0.1
rain	sun	0.3
rain	rain	0.7

Application of Stationary Distribution: Web Link Analysis

• PageRank over a web graph

- Each web page is a possible value of a state
- Initial distribution: uniform over pages
- Transitions:

With prob. c, uniform jump to a random page (dotted lines, not all shown)
With prob. 1-c, follow a random outlink (solid lines)

• Stationary distribution

- Will spend more time on highly reachable pages
- E.g. many ways to get to the Acrobat Reader download page
- Google 1.0 returned the set of pages containing all your keywords in decreasing rank, now all search engines use link analysis along with many other factors (rank actually getting less important over time)



Hidden Markov Models



Pacman – Sonar

74 CS188 Pacman	
SCORE: -9	9.0 9.0 XXX 12.0

Hidden Markov Models

• Markov chains not so useful for most agents

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:

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

R _{t-1}	R_{t}	$P(R_{t} R_{t})$	R _t	l
+r	+r	0.7	+r	+
+r	-r	0.3	+r	-1
-r	+r	0.3	-r	+
-r	-r	0.7	-r	-1

R_{t}	U_{t}	$P(U_t R_t)$
+r	+u	0.9
+r	-u	0.1
-r	+u	0.2
-r	-u	0.8

Example: Ghostbusters HMM

- \circ P(X₁) = uniform
- P(X|X') = usually move clockwise, but sometimes move in a random direction or stay in place



1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

P(X₁)

P(R_{ij} | X) = same sensor model as before:
 red means close, green means far away.







P(X|X'=<1,2>)

Video of Demo Ghostbusters – Circular Dynamics -- HMM

	ghostbusters (beliefs dynamic, center)	
	2 ghostbusters (beliefs dynamic, circle)	
	🤗 3 ghostousters (beliefs dynamic, basic)	
	4 pacman sonar.py (no beliefs)	
	🧉 5 pacman sonar.py	
	🧬 6 ghostbusters (beliefs dynamic, circle, particles) (tons)	
	7 ghostbusters (beliefs dynamic, circle, particles)	
	8 ghostbusters (beliefs dynamic, circle, particles, some)	
	🥭" 9 ghostbusters (beliefs dynamic, circle, no noise)	
	🧉 1st class pacman	
	Run As	
	Run Configurations	
	Organize Favorites	
Console 🔀	enter	
Console 🛛	enter	
Console 🕅 ninated> ce	enter	
Console 🕅 ninated> ce	enter	
Console 🕅 ninated> ce	enter	
ionsole 🛛	enter	

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 tend to correlated by the hidden state]

Real HMM Examples

• Robot tracking:

- Observations are range readings (continuous)
- States are positions on a map (continuous)

• Speech recognition HMMs:

- Observations are acoustic signals (continuous valued)
- States are specific positions in specific words (so, tens of thousands)

• Machine translation HMMs:

- Observations are words (tens of thousands)
- States are translation options

Filtering / Monitoring

- Filtering, or monitoring, is the task of tracking the distribution $B_t(X) = P_t(X_t | e_1, ..., e_t)$ (the belief state) over time
- We start with $B_1(X)$ in an initial setting, usually uniform
- As time passes, or we get observations, we update B(X)
- The Kalman filter was invented in the 60's and first implemented as a method of trajectory estimation for the Apollo program

























Inference: Find State Given Evidence

 $\,\circ\,\,$ We are given evidence at each time and want to know

$$B_t(X) = P(X_t | e_{1:t})$$

 $_{\odot}$ Idea: start with P(X₁) and derive B_t in terms of B_{t-1}

 $_{\rm O}$ equivalently, derive B_{t+1} in terms of B_t

Background: Probabilistic Inference

- Probabilistic inference: compute a desired probability from other known probabilities (e.g. conditional from joint)
- We generally compute conditional probabilities
 - P(on time | no reported accidents) = 0.90
 - These represent the agent's *beliefs* given the evidence

• Probabilities change with new evidence:

- P(on time | no accidents, 5 a.m.) = 0.95
- P(on time | no accidents, 5 a.m., raining) = 0.80
- Observing new evidence causes *beliefs to be updated*



 \circ P(W)?

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

 \circ P(W)?

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

 \circ P(W)?

P(sun)=.3+.1+.1+.15=.65

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

 \circ P(W)?

P(sun)=.3+.1+.1+.15=.65 P(rain)=1-.65=.35

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

- General case:
 - Evidence variables:
 - Query* variable:
 - Hidden variables:
- $\left. \begin{array}{c} E_1 \dots E_k = e_1 \dots e_k \\ Q \\ H_1 \dots H_r \end{array} \right\} \begin{array}{c} X_1, X_2, \dots X_n \\ \text{All variables} \end{array}$
- 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



Step 3: Normalize





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

 \circ P(W | winter)?

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

 \circ P(W | winter)?

P(sun|winter)~.1+.15=.25

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

 \circ P(W | winter)?

P(rain|winter)~.05+.2=.25

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

\circ P(W | winter)?

P(sun|winter)~.25 P(rain|winter)~.25 P(sun|winter)=.5 P(rain|winter)=.5

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

0	P(W	winter,	hot)?
	`	· /	- /

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

0	P(W	winter,	hot)?
	`	· /	- /

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

• $P(W \mid winter, hot)$?

P(sun|winter,hot)~.1 P(rain|winter,hot)~.05

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

• $P(W \mid winter, hot)$?

P(sun|winter,hot)~.1 P(rain|winter,hot)~.05 P(sun|winter,hot)=2/3 P(rain|winter,hot)=1/3

S	Т	W	Р
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

Obvious problems:

- Worst-case time complexity O(dⁿ)
- Space complexity O(dⁿ) to store the joint distribution

Next Topic

• Inference in HMMs