CSE 573: Artificial Intelligence

Hanna Hajishirzi Hidden Markov Models

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



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
$$P(X_1, X_2, ... X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2)...$$

= $\prod_{i=1}^n P(X_i|X_1, ..., X_{i-1})$

- **X,** Y independent if and only if: $\forall x, y : P(x, y) = P(x)P(y)$
- lacksquare X and Y are conditionally independent given Z if and only if: $X \perp\!\!\!\perp Y \mid Z$

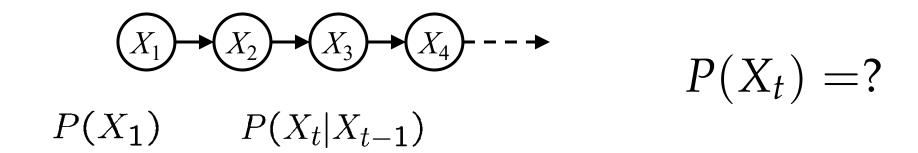
$$\forall x, y, z : P(x, y|z) = P(x|z)P(y|z)$$

Reasoning over Time or Space

- Often, we want to reason about a sequence of observations
 - Speech recognition
 - Robot localization
 - User attention
 - Medical monitoring
- Need to introduce time (or space) into our models

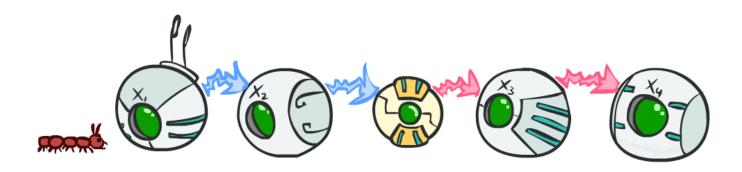
Markov Models

o Value of X at a given time is called the state



- o Parameters: called transition probabilities or dynamics, specify how the state evolves over time (also, initial state probabilities)
- o Stationarity assumption: transition probabilities the same at all times
- o Same as MDP transition model, but no choice of action
- o 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



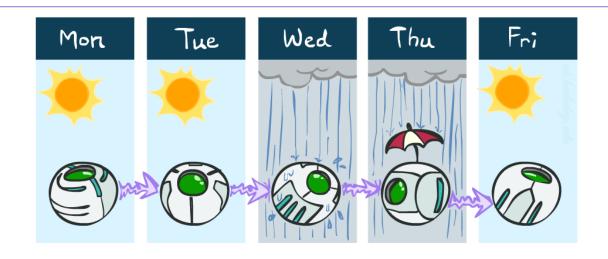
Basic conditional independence:

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

Example Markov Chain: Weather

States: X = {rain, sun}

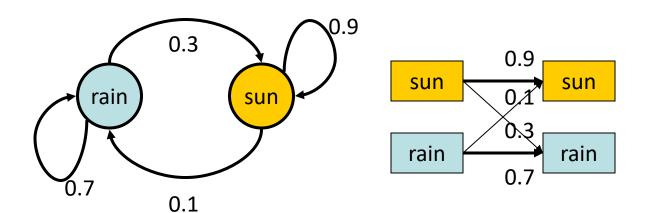
Initial distribution: 1.0 sun



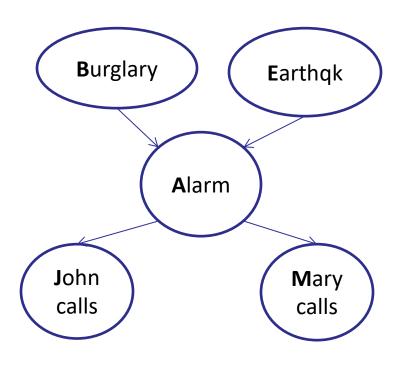
• CPT $P(X_t | X_{t-1})$:

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



Bayes Nets -- Independence



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

Markov Models (Markov Chains)

$$X_1$$
 X_2 X_3 X_4 X_N

- A Markov model defines
 - o a joint probability distribution:

$$P(X_1, X_2, X_3, X_4) =$$

More generally:

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

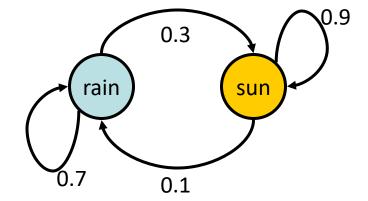
$$P(X_1,\ldots,X_n)=P(X_1)\prod_{t=2}^N P(X_t|X_{t-1})$$
 • Why? • Chain Rule,

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

Indep. Assumption?

Example Markov Chain: Weather

o Initial distribution: 1.0 sun



• What is the probability distribution after one step?

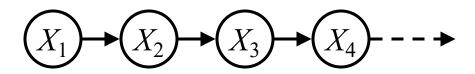
$$P(X_2 = sun) = \sum_{x_1} P(x_1, X_2 = sun) = \sum_{x_1} P(X_2 = sun | x_1) P(x_1)$$

$$P(X_2 = \text{sun}) = P(X_2 = \text{sun}|X_1 = \text{sun})P(X_1 = \text{sun}) + P(X_2 = \text{sun}|X_1 = \text{rain})P(X_1 = \text{rain})$$

$$0.9 \cdot 1.0 + 0.3 \cdot 0.0 = 0.9$$

Mini-Forward Algorithm

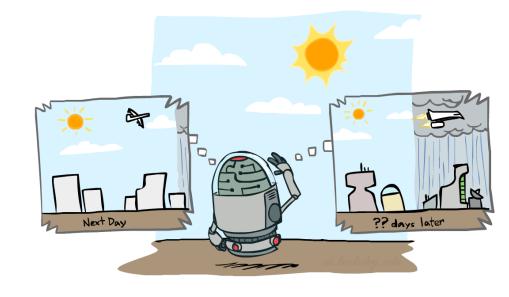
Question: What's P(X) on some day t?



$$P(x_1) = known$$

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

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



Example Run of Mini-Forward Algorithm

From initial observation of sun

From initial observation of rain

• From yet another initial distribution $P(X_1)$:

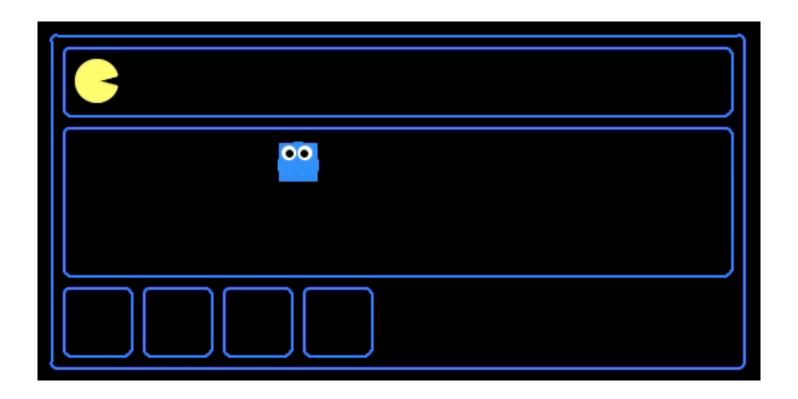
$$\left\langle \begin{array}{c} p \\ 1-p \\ P(X_1) \end{array} \right\rangle \qquad \cdots \qquad \left\langle \begin{array}{c} 0.75 \\ 0.25 \\ P(X_{\infty}) \end{array} \right\rangle$$

13

[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



Stationary Distributions

o For most chains:

- o 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?

$$X_1$$
 X_2 X_3 X_4 X_4

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

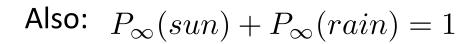
$$P_{\infty}(rain) = P(rain|sun)P_{\infty}(sun) + P(rain|rain)P_{\infty}(rain)$$

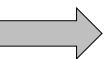
$$P_{\infty}(sun) = 0.9P_{\infty}(sun) + 0.3P_{\infty}(rain)$$

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

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

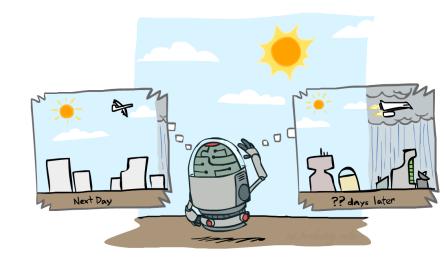
$$P_{\infty}(rain) = 1/3P_{\infty}(sun)$$





$$P_{\infty}(sun) = 3/4$$

$$P_{\infty}(rain) = 1/4$$



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

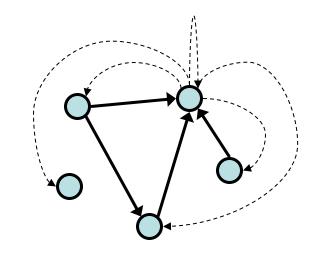
Application of Stationary Distribution: Web Link Analysis

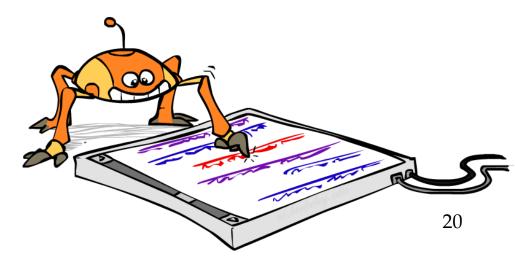
PageRank over a web graph

- o Each web page is a possible value of a state
- o Initial distribution: uniform over pages
- o Transitions:
 - With prob. c, uniform jump to a random page (dotted lines, not all shown)
 - o With prob. 1-c, follow a random outlink (solid lines)

Stationary distribution

- Will spend more time on highly reachable pages
- o E.g. many ways to get to the Acrobat Reader download page
- o 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)

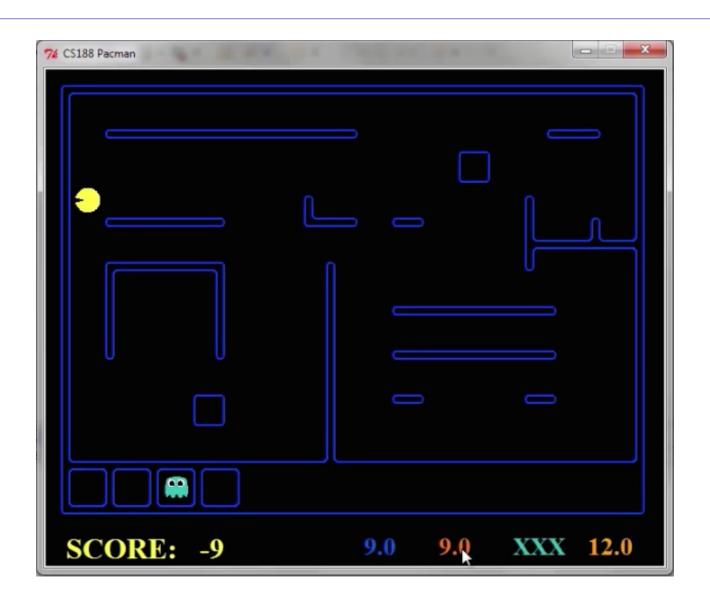




Hidden Markov Models

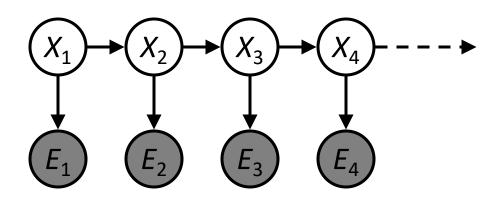


Pacman – Sonar



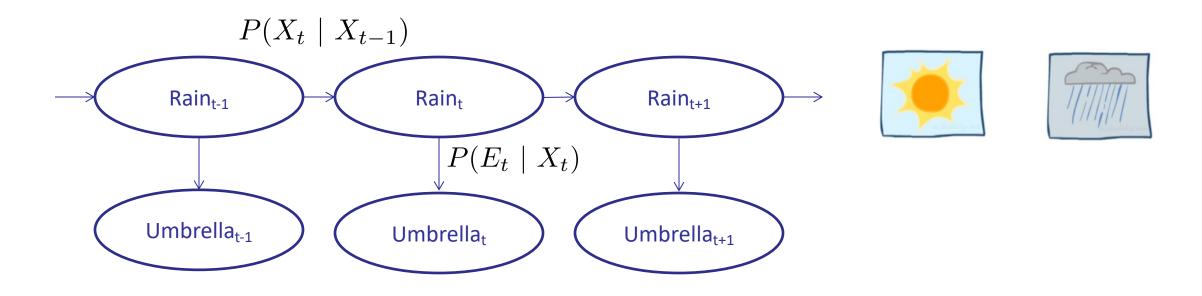
Hidden Markov Models

- 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:

o Initial distribution: $P(X_1)$

o Transitions: $P(X_t \mid X_{t-1})$

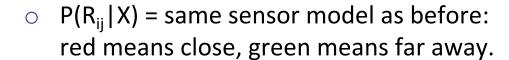
o Emissions: $P(E_t \mid X_t)$

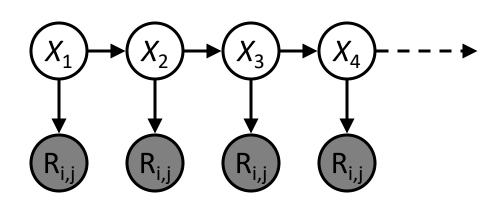
R _{t-1}	R _t	$P(R_t R_{t-1})$
+r	+r	0.7
+r	-r	0.3
-r	+r	0.3
-r	-r	0.7

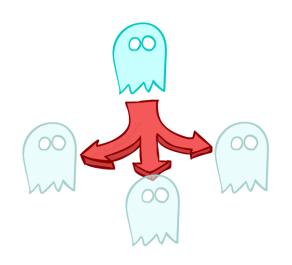
R _t	U _t	$P(U_t R_t)$
+r	+u	0.9
+r	-u	0.1
-r	+u	0.2
-r	-u	0.8

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

 $P(X_1)$

1/6	16	1/2
0	1/6	0
0	0	0

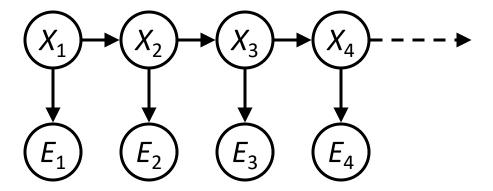
$$P(X | X' = <1,2>)$$

Video of Demo Ghostbusters – Circular Dynamics -- HMM



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:

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

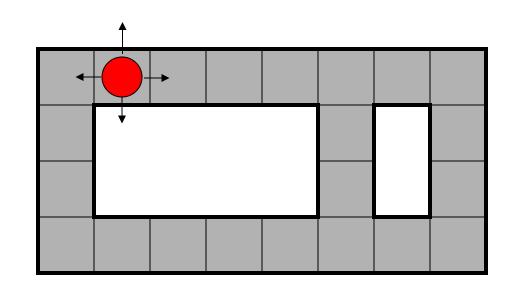
Machine translation HMMs:

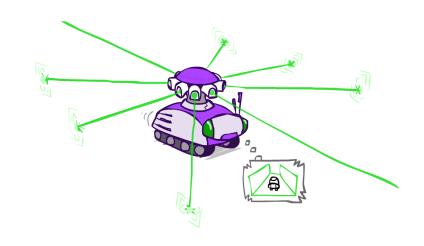
- Observations are words (tens of thousands)
- States are translation options

Filtering / Monitoring

- Filtering, or monitoring, is the task of tracking the distribution $B_t(X) = P_t(X_t \mid e_1, ..., e_t)$ (the belief state) over time
- \circ We start with $B_1(X)$ in an initial setting, usually uniform
- \circ 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

Example from Michael Pfeiffer

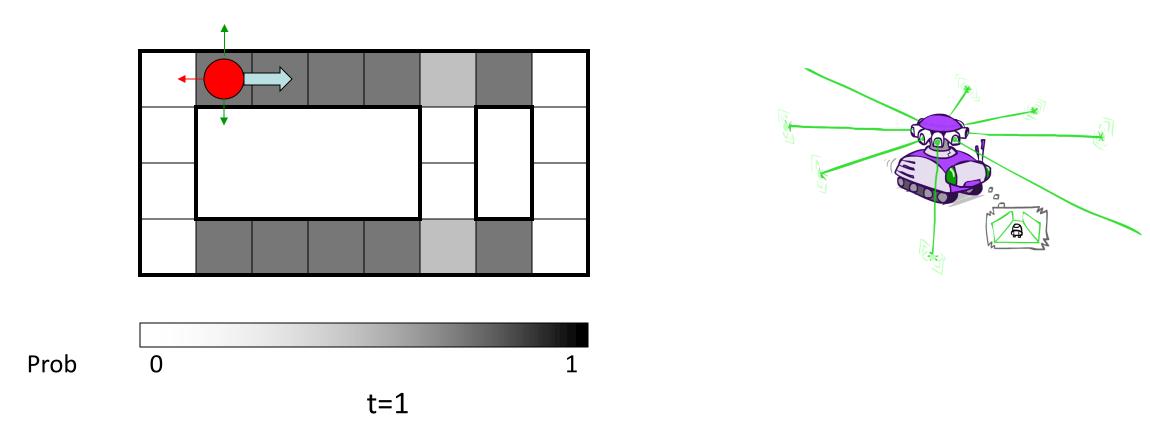




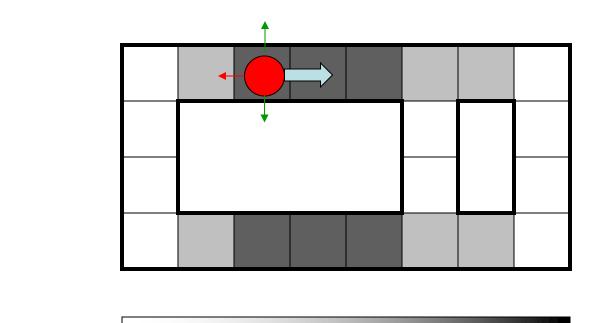


Sensor model: can read in which directions there is a wall, never more than 1 mistake

Motion model: may not execute action with small prob.

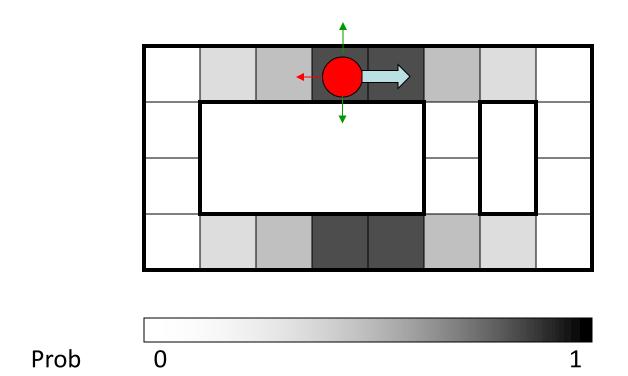


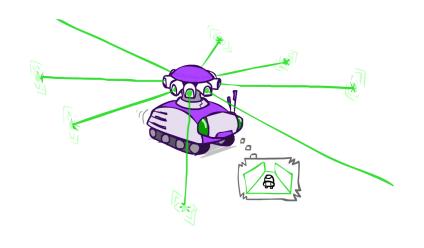
Lighter grey: was possible to get the reading, but less likely b/c required 1 mistake



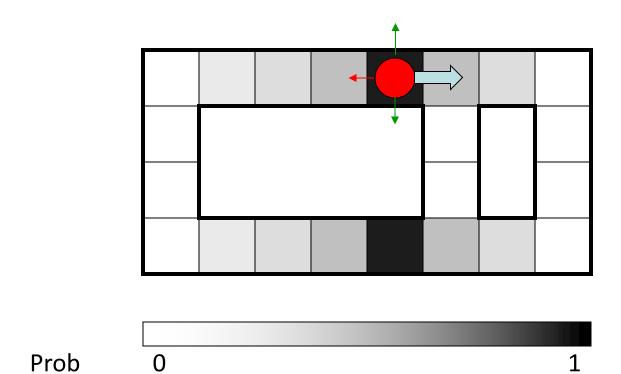


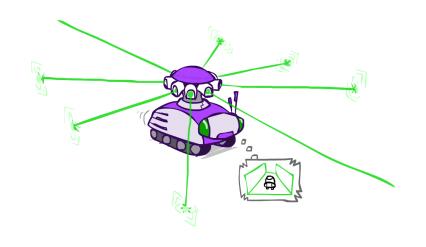
Prob 0 1



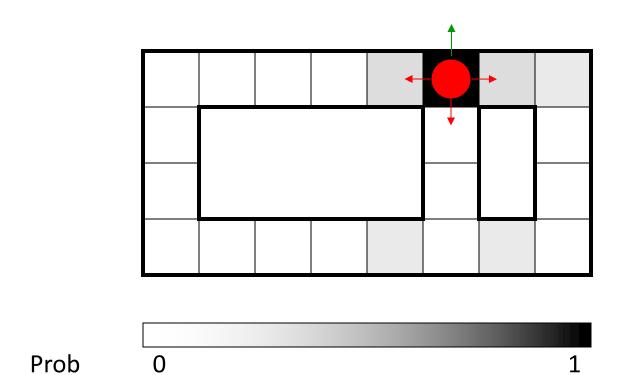


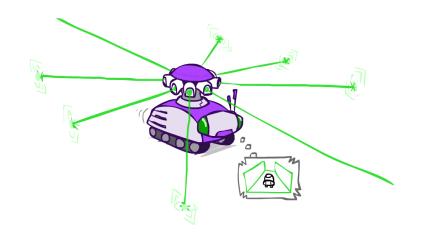
t=3





t=4





t=5

Inference: Find State Given Evidence

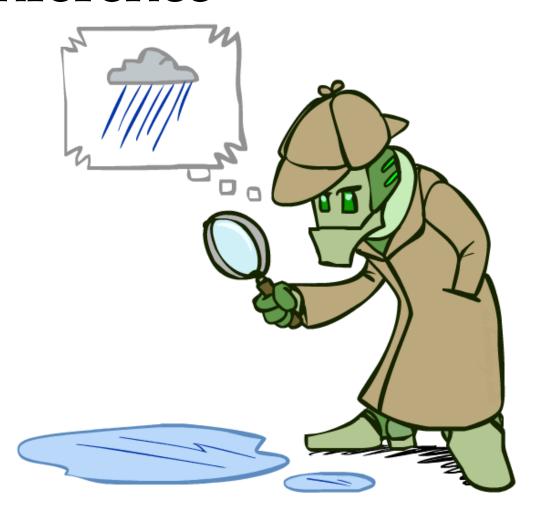
We are given evidence at each time and want to know

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

- Idea: start with P(X₁) and derive B_t in terms of B_{t-1}
 - 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
 - o P(on time ∣ no reported accidents) = 0.90
 - o These represent the agent's beliefs given the evider
- Probabilities change with new evidence:
 - o P(on time | no accidents, 5 a.m.) = 0.95
 - o P(on time | no accidents, 5 a.m., raining) = 0.80
 - o Observing new evidence causes beliefs to be updated



Inference by Enumeration

o 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

Inference by Enumeration

o 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

Inference by Enumeration

 \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

o 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

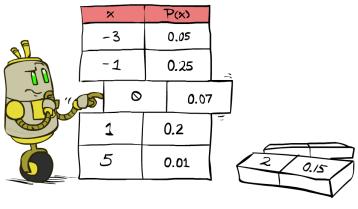
General case:

o Evidence variables: $E_1 \dots E_k = e_1 \dots e_k$ o Query* variable: Qo Hidden variables: $H_1 \dots H_r$ $X_1, X_2, \dots X_n$ We want:

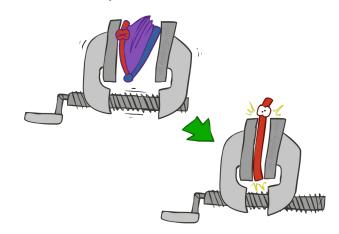
* Works fine with multiple query variables, too

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

 Step 1: Select the entries consistent with the evidence



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



$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(Q, h_1 \dots h_r, e_1 \dots e_k)$$

$$X_1, X_2, \dots X_n$$

Step 3: Normalize

$$\times \frac{1}{Z}$$

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

$$P(Q|e_1 \cdots e_k) = \frac{1}{Z} P(Q, e_1 \overset{44}{\cdots} e_k)$$

o 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

○ 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

○ 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

○ 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

○ 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 | 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

o P(W ∣ 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 | 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