

CSE 473: Artificial Intelligence Spring 2012

Reasoning about Uncertainty & Hidden Markov Models

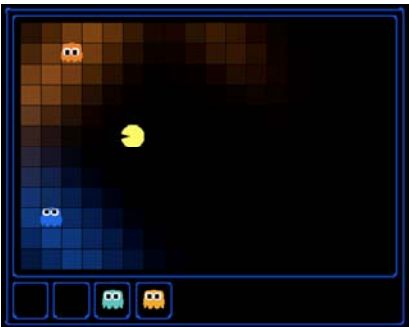
Daniel Weld

Many slides adapted from Dan Klein, Stuart Russell, Andrew Moore & Luke Zettlemoyer

Outline

- Probabilistic sequence models (and inference)
 - Bayesian Networks – Preview
 - Markov Chains
 - Hidden Markov Models
 - Exact Inference
 - Particle Filters

Going Hunting



4

Inference by Enumeration

- General case:

▪ Evidence variables: $E_1 \dots E_k = e_1 \dots e_k$	}	X_1, X_2, \dots, X_n <i>All variables</i>
▪ Query* variable: Q		
▪ Hidden variables: $H_1 \dots H_r$		
- We want: $P(Q|e_1 \dots e_k)$
- First, select the entries consistent with the evidence
- Second, sum out H to get joint of Query and evidence:

$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} \underbrace{P(Q, h_1 \dots h_r, e_1 \dots e_k)}_{X_1, X_2, \dots, X_n}$$
- Finally, normalize the remaining entries to conditionalize
- Obvious problems:
 - Worst-case time complexity $O(d^r)$
 - Space complexity $O(d^r)$ to store the joint distribution

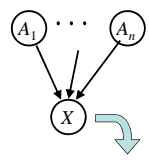
Bayes' Nets: Big Picture

- Two problems with using full joint distribution tables as our probabilistic models:
 - Unless there are only a few variables, the joint is WAY too big to represent explicitly
 - Hard to learn (estimate) anything empirically about more than a few variables at a time
- **Bayes' nets:** a technique for describing complex joint distributions (models) using simple, local distributions (conditional probabilities)
 - More properly called **graphical models**
 - We describe how variables locally interact
 - Local interactions chain together to give global, indirect interactions

Bayes' Net Semantics

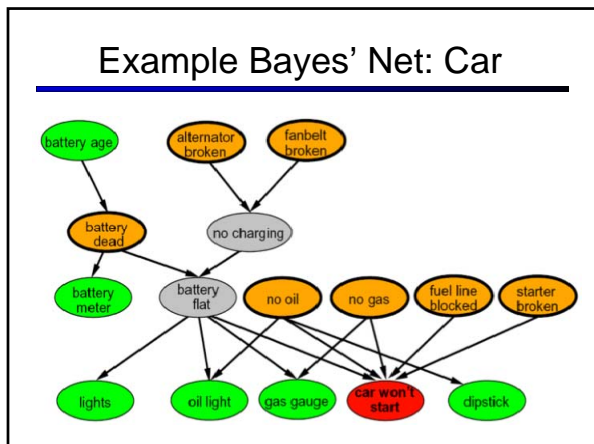
Formally:

- A set of **nodes**, one per variable X
- A **directed, acyclic graph**
- A **CPT for each node**
 - CPT = "Conditional Probability Table"
 - Collection of distributions over X , one for each combination of parents' values



$P(X|a_1 \dots a_n)$

A Bayes net = Topology (graph) + Local Conditional Probabilities



Hidden Markov Models

- Markov chains not so useful for most agents
 - Eventually you don't know anything anymore
 - Need observations to update your beliefs
- Hidden Markov models (HMMs)
 - Underlying Markov chain over states S
 - You observe outputs (effects) at each time step
 - POMDPs without actions (or rewards).
 - As a Bayes' net:

Hidden Markov Models

- Defines a joint probability distribution:

$$P(\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{E}_1, \dots, \mathbf{E}_n) = P(\mathbf{X}_1) P(\mathbf{E}_1 | \mathbf{X}_1) \prod_{i=2}^n P(\mathbf{X}_i | \mathbf{X}_{i-1}) P(\mathbf{E}_i | \mathbf{X}_i)$$

Example

- An HMM is defined by:
 - Initial distribution: $P(\mathbf{X}_1)$
 - Transitions: $P(\mathbf{X}_i | \mathbf{X}_{i-1})$
 - Emissions: $P(\mathbf{E}_i | \mathbf{X}_i)$

Ghostbusters HMM

- $P(X_i) = \text{uniform}$
- $P(X_i | X_{i-1}) = \text{usually move clockwise, but sometimes move in a random direction or stay in place}$
- $P(E_i | X_i) = \text{same sensor model as before: red means close, green means far away.}$

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

$P(X_i)$

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

$P(X_i | X_{i-1}, 2 >)$

P(red 3)	P(orange 3)	P(yellow 3)	P(green 3)
0.05	0.15	0.5	0.3

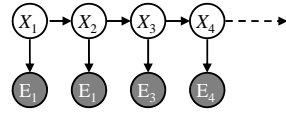
$P(E_i | X_i)$

HMM Computations

- Given
 - joint $P(X_{1:n}, E_{1:n})$
 - evidence $E_{1:n} = e_{1:n}$
- Inference problems include:
 - Filtering, find $P(X_t | e_{1:t})$ for some t
 - Smoothing, find $P(X_t | e_{1:n})$ for some t

HMM Computations

- Given
 - joint $P(X_{1:n}, E_{1:n})$
 - evidence $E_{1:n} = e_{1:n}$



```

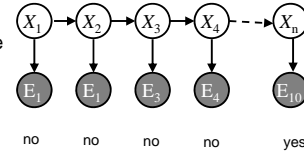
    graph LR
      X1((X1)) --> X2((X2))
      X2 --> X3((X3))
      X3 --> X4((X4))
      X4 -.-> Xn[...]
      E1((E1)) --> X1
      E2((E2)) --> X2
      E3((E3)) --> X3
      E4((E4)) --> X4
    
```

- Inference problems include:
 - **Filtering**, find $P(X_t/e_{1:n})$ for all t
 - **Smoothing**, find $P(X_t/e_{1:n})$ for all t
 - **Most probable explanation**, find

$$x^*_{1:n} = \operatorname{argmax}_{x_{1:n}} P(x_{1:n}/e_{1:n})$$

HMM Computations

- **State = number**
 - T++ add sum of 2 dice
- Observation:
 - Yes/no
 - Is sum 75?



```

    graph LR
      X1((X1)) --> X2((X2))
      X2 --> X3((X3))
      X3 --> X4((X4))
      X4 -.-> Xn[...]
      E1((E1)) --> X1
      E2((E2)) --> X2
      E3((E3)) --> X3
      E4((E4)) --> X4
      E10((E10)) --> Xn
    
```

no no no no yes

- **Most probable explanation**, find

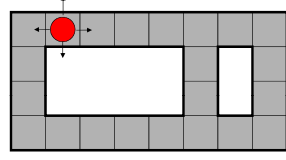
$$x^*_{1:n} = \operatorname{argmax}_{x_{1:n}} P(x_{1:n}/e_{1:n})$$

Filtering / Monitoring

- Filtering, or monitoring, is the task of tracking the distribution $B(X)$ (the belief state) over time
- We start with $B(X)$ 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

Example: Robot Localization

Example from Michael Pfeiffer

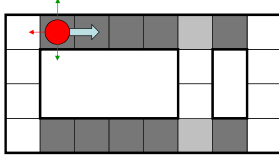


Prob 0 1

t=0

Sensor model: never more than 1 mistake
Motion model: may not execute action with small prob.

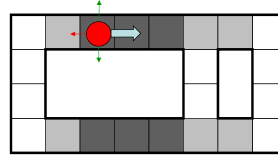
Example: Robot Localization



Prob 0 1

t=1

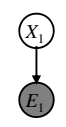
Example: Robot Localization



Prob 0 1


t=2

Inference Recap: Simple Cases



$P(X_1|e_1)$

That's my rule!

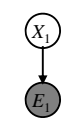


$$P(x_1|e_1) = P(x_1, e_1)/P(e_1)$$

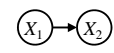
$$\propto_{X_1} P(x_1, e_1)$$

$$= P(x_1)P(e_1|x_1)$$

Inference Recap: Simple Cases



$P(X_1|e_1)$



$P(X_2)$

$$P(x_1|e_1) = P(x_1, e_1)/P(e_1)$$

$$\propto_{X_1} P(x_1, e_1)$$

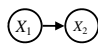
$$= P(x_1)P(e_1|x_1)$$

$$P(x_2) = \sum_{x_1} P(x_1, x_2)$$

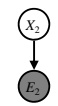
$$= \sum_{x_1} P(x_1)P(x_2|x_1)$$

Online Belief Updates

- Every time step, we start with current $P(X | \text{evidence})$
- We update for time:



$$P(x_t|e_{1:t-1}) = \sum_{x_{t-1}} P(x_{t-1}|e_{1:t-1}) \cdot P(x_t|x_{t-1})$$
- We update for evidence:

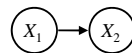


$$P(x_t|e_{1:t}) \propto_X P(x_t|e_{1:t-1}) \cdot P(e_t|x_t)$$

Passage of Time

- Assume we have current belief $P(X | \text{evidence to date})$

$$B(X_t) = P(X_t|e_{1:t})$$
- Then, after one time step passes:

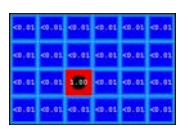


$$P(X_{t+1}|e_{1:t}) = \sum_{x_t} P(X_{t+1}|x_t)P(x_t|e_{1:t})$$
- Or, compactly:

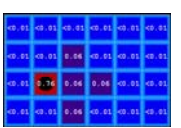
$$B'(X') = \sum_x P(X'|x)B(x)$$
- Basic idea: beliefs get "pushed" through the transitions
 - With the "B" notation, we have to be careful about what time step t the belief is about, and what evidence it includes

Example: Passage of Time


Without observations, uncertainty "accumulates"



T = 1



T = 2



T = 5

$$B'(X') = \sum_x P(X'|x)B(x)$$

Transition model: ghosts usually go clockwise

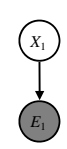
Observations

- Assume we have current belief $P(X | \text{previous evidence})$:

$$B'(X_{t+1}) = P(X_{t+1}|e_{1:t})$$
- Then:

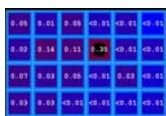
$$P(X_{t+1}|e_{1:t+1}) \propto P(e_{t+1}|X_{t+1})P(X_{t+1}|e_{1:t})$$
- Or:

$$B(X_{t+1}) \propto P(e_t|X)B'(X_{t+1})$$
- Basic idea: beliefs reweighted by likelihood of evidence
- Unlike passage of time, we have to renormalize

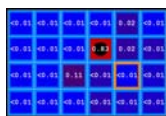


Example: Observation

- As we get observations, beliefs get reweighted, uncertainty “decreases”



Before observation



After observation

$$B(X) \propto P(e|X)B'(X)$$

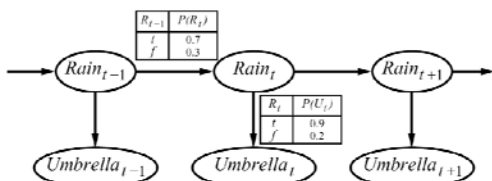
The Forward Algorithm

- We want to know: $B_t(X) = P(X_t|e_{1:t})$
- We can derive the following updates

$$\begin{aligned} P(x_t|e_{1:t}) &\propto_X P(x_t, e_{1:t}) \\ &= \sum_{x_{t-1}} P(x_{t-1}, x_t, e_{1:t}) \\ &= \sum_{x_{t-1}} P(x_{t-1}, e_{1:t-1})P(x_t|x_{t-1})P(e_t|x_t) \\ &= P(e_t|x_t) \sum_{x_{t-1}} P(x_t|x_{t-1})P(x_{t-1}, e_{1:t-1}) \end{aligned}$$

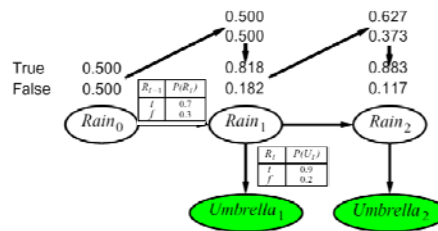
- To get $B_t(X)$ compute each entry and normalize

Example



- An HMM is defined by:
 - Initial distribution: $P(X_1)$
 - Transitions: $P(X_t|X_{t-1})$
 - Emissions: $P(E_t|X_t)$

Forward Algorithm



Example Pac-man



Summary: Filtering

- Filtering is the inference process of finding a distribution over X_T given e_1 through e_T : $P(X_T | e_{1:T})$
- We first compute $P(X_1 | e_1)$: $P(x_1|e_1) \propto P(x_1) \cdot P(e_1|x_1)$
- For each t from 2 to T , we have $P(X_{t-1} | e_{1:t-1})$
- Elapse time:** compute $P(X_t | e_{1:t-1})$

$$P(x_t|e_{1:t-1}) = \sum_{x_{t-1}} P(x_{t-1}|e_{1:t-1}) \cdot P(x_t|x_{t-1})$$

- Observe:** compute $P(X_t | e_{1:t-1}, e_t) = P(X_t | e_{1:t})$

$$P(x_t|e_{1:t}) \propto P(x_t|e_{1:t-1}) \cdot P(e_t|x_t)$$

Recap: Reasoning Over Time

- Stationary Markov models
 - Diagram: $X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow X_4 \rightarrow \dots$ with transitions to 'rain' (0.7) and 'sun' (0.5) states.
 - Probabilities: $P(X_1)$, $P(X|X_{-1})$, $P(E|X)$
- Hidden Markov models
 - Diagram: $X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow X_4 \rightarrow \dots$ with observations E_1, E_2, E_3, E_4 .
 - Table:

X	E	P
rain	umbrella	0.9
rain	no umbrella	0.1
sun	umbrella	0.2
sun	no umbrella	0.8

Particle Filtering

- Sometimes $|X|$ is too big for exact inference
 - $|X|$ may be too big to even store $B(X)$
 - E.g. when X is continuous
 - $|X|^2$ may be too big to do updates
- Solution: approximate inference
 - Track samples of X , not all values
 - Samples are called **particles**
 - Time per step is linear in the number of samples
 - But: number needed may be large
 - In memory: list of particles, not states
- How robot localization works in practice
 - Diagram: 3x3 grid showing particles (red dots) moving from a single dot to multiple dots.

Representation: Particles

- Our representation of $P(X)$ is now a list of N particles (samples)
 - Generally, $N \ll |X|$
 - Storing map from X to counts would defeat the point
- $P(x)$ approximated by number of particles with value x
 - So, many x will have $P(x) = 0!$
 - More particles, more accuracy
- For now, all particles have a weight of 1
 - Diagram: 3x3 grid with red dots.
 - Particles: (3,3), (2,3), (3,3), (3,3), (3,2), (3,3), (3,2), (2,1), (3,3), (3,3), (2,1)

Particle Filtering: Elapse Time

- Each particle is moved by sampling its next position from the transition model
 - $x' = \text{sample}(P(X'|x))$
 - This is like prior sampling – samples' frequencies reflect the transition probs
 - Here, most samples move clockwise, but some move in another direction or stay in place
- This captures the passage of time
 - If we have enough samples, close to the exact values before and after (consistent)
 - Diagram: 3x3 grid showing particles moving from one state to another.

Particle Filtering: Observe

- How handle noisy observations?
 - Diagram: 3x3 grid with red dots and a red box highlighting a specific state.
- Suppose sensor gives red reading?
 - Diagram: 3x3 grid with red dots and a red box highlighting a specific state.

Particle Filtering: Observe

Slightly trickier:

- Don't do **rejection sampling** (why not?)
- We don't sample the observation, we fix it
- Instead: **downweight samples** based on the evidence (form of likelihood weighting)

$$w(x) = P(e|x)$$

$$B(X) \propto P(e|X)B'(X)$$

- Note: as before, probabilities **don't sum to one**, since most have been downweighted (in fact they sum to an approximation of $P(e)$)
- Diagram: 3x3 grid showing particles being downweighted (smaller dots) and a red box highlighting a specific state.

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

Old Particles:

- (3,3) w=0.1
- (2,1) w=0.9
- (2,1) w=0.9
- (3,1) w=0.4
- (3,2) w=0.3
- (2,2) w=0.4
- (1,1) w=0.4
- (3,1) w=0.4
- (2,1) w=0.9
- (3,2) w=0.3

New Particles:

- (2,1) w=1
- (2,1) w=1
- (2,1) w=1
- (3,2) w=1
- (2,2) w=1
- (2,1) w=1
- (1,1) w=1
- (3,1) w=1
- (2,1) w=1
- (1,1) w=1

Particle Filtering Summary

- Represent current belief $P(X | \text{evidence to date})$ as set of n samples (actual assignments $X=x$)
- For each new observation e :
 - Sample transition, once for each current particle x

$$x' = \text{sample}(P(X'|x))$$
 - For each new sample x' , compute importance weights for the new evidence e :

$$w(x') = P(e|x')$$
 - Finally, normalize by resampling the importance weights to create N new particles

Robot Localization

- In robot localization:
 - We know the map, but not the robot's position
 - Observations may be vectors of range finder readings
 - State space and readings are typically continuous (works basically like a very fine grid) and so we cannot store $B(X)$
 - Particle filtering is a main technique

Robot Localization

QuickTime™ and a GIF decompressor are needed to see this picture.

Which Algorithm?

Particle filter, uniform initial beliefs, 25 particles

SCORE: 0

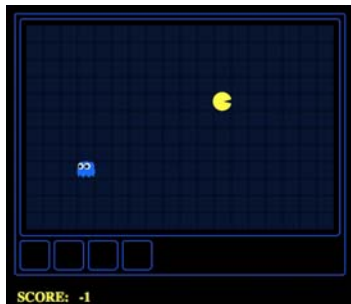
Which Algorithm?

Particle filter, uniform initial beliefs, 300 particles

SCORE: 0

Which Algorithm?

Exact filter, uniform initial beliefs



P4: Ghostbusters

- **Plot:** Pacman's grandfather, Grandpac, learned to hunt ghosts for sport.
- He was blinded by his power, but could hear the ghosts' banging and clanging.
- **Transition Model:** All ghosts move randomly, but are sometimes biased
- **Emission Model:** Pacman knows a "noisy" distance to each ghost

Noisy distance prob
True distance = 8

