Markov Models and Hidden Markov Models

Robert Platt Northeastern University

Some images and slides are used from: 1. CS188 UC Berkeley 2. RN, AIMA

Markov Models

We have already seen that an MDP provides a useful framework for modeling stochastic control problems.

Markov Models: model any kind of temporally dynamic system.

Probability recap

- 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: $\forall r$

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

X and Y are conditionally independent given Z if and only if:

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

Probability again: Independence

Two random variables, *x* and *y*, are independent when:

$$\forall (x, y), P(x, y) = P(x)P(y) \iff x \perp y$$
$$x \not\perp y$$

The outcomes of two different coin flips are usually independent of each other

Probability again: Independence

If:
$$P(x, y) = P(x)P(y)$$

Then:
$$P(x) = P(x|y)$$

 $P(y) = P(y|x)$

Why?

Probability again: Independence

Two random variables, *x* and *y*, are independent when:

$$\forall (x, y), P(x, y) = P(x)P(y) \iff x \perp y$$
$$x \not\perp y$$

The outcomes of two different coin flips are usually independent of each other

	winter	!winter
snow	0.1	0.1
!snow	0.3	0.5

	winter	!winter
snow	0.1	0.1
!snow	0.3	0.5

Are snow and winter independent variables?

	winter	!winter
snow	0.1	0.1
!snow	0.3	0.5

Are snow and winter independent variables?

P(snow) = 0.2

P(winter) = 0.4

	winter	!winter
snow	0.1	0.1
!snow	0.3	0.5

Are snow and winter independent variables?

P(snow) = 0.2P(winter) = 0.4

What would the distribution look like if snow, winter were independent?

Conditional independence

Independence:

$$\forall (x, y), P(x, y) = P(x)P(y)$$
$$x \perp \!\!\!\perp y$$

Conditional independence:

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

Equivalent statements of conditional independence:

P(x|z) = P(x|z, y)

$$P(y|z) = P(y|z, x)$$

Conditional independence: example



P(toothache, catch | cavity) = P(toothache | cavity) P(catch | cavity)

Toothache and catch are conditionally independent given cavity – this is the "common cause" scenario covered in Bayes Nets...

Examples of conditional independence

What are the conditional independence relationships in the following?

- traffic, raining, late for work
- snow, cloudy, crash
- fire, smoke, alarm



Markov model can be used to model any sequential time process

- the weather
- traffic
- stock market
- news cycle

. . .



Since this is a Markov process, we assume transitions are Markov:

Process model: $P(X_t | X_{t-1}) = P(X_t | X_{t-1}, ..., X_1)$

Markov assumption: $X_t \perp X_{t-2} | X_{t-1}$



How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$

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



How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$

 $P(X_1, X_2, X_3, X_4) = \underbrace{P(X_1)P(X_2|X_1)P(X_3|X_2, X_1)P(X_4|X_3, X_2, X_1)}_{P(X_2, X_1)}$



How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$





How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$

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

Can we simplify this expression?



How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$

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



How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$

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



How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$

 $P(X_1, X_2, X_3, X_4) = P(X_1)P(X_2|X_1)P(X_3|X_2, X_1)P(X_4|X_3, X_2, X_1)$ $P(X_1, X_2, X_3, X_4) = P(X_1)P(X_2|X_1)P(X_3|X_2)P(X_4|X_3)$ T-1

In general: $P(X_1, X_2, \dots, X_T) = P(X_1) \prod_{t=1} P(X_{t+1}|X_t)$



How do we calculate: $P(X_1, X_2, X_3, X_4) = ?$ $P(X_1, X_2, X_3, X_4) = P(X_1)P(X_2|X_1)P(X \text{ Process model})$ $X_2, X_1)$ $P(X_1, X_2, X_3, X_4) = P(X_1)P(X_2|X_1)P(X_3|X_2)P(X_4|X_3)$ In general: $P(X_1, X_2, ..., X_T) = P(X_1) \prod P(X_{t+1}|X_t)$

Markov Processes: example

Two states: cloudy, sunny

X_{t-1}	X_t	X_t
sun	sun	0.8
sun	cloudy	0.2
cloudy	sun	0.3
cloudy	cloudy	0.7



 $t \in \{\text{mon, tues, weds, thurs, fri}\}$



But, suppose we want to predict the state at time T, given a prior distribution at time 1?

$$P(X_{2}) = \sum_{X_{1}} P(X_{1})P(X_{2}|X_{1})$$
$$P(X_{3}) = \sum_{X_{2}} P(X_{2})P(X_{3}|X_{2})$$
$$\vdots$$
$$P(X_{T}) = \sum_{X_{T-1}} P(X_{T-1})P(X_{T}|X_{T-1})$$

 $P(x_1) = 1$ Suppose is it sunny on mon... $P(x_2) = P(x_2|x_1)P(x_1)$ Prob sunny tues = 0.8 $P(x_3) = P(x_3|x_2)P(x_2) + P(x_3|\bar{x}_2)P(\bar{x}_2)$ Prob sunny weds = 0.64 + 0.06 = 0.7 $P(x_4) = P(x_4|x_3)P(x_3) + P(x_4|\bar{x}_3)P(\bar{x}_3)$ Prob sunny thurs = 0.56 + 0.09 = 0.65 $P(x_5) = P(x_5|x_4)P(x_4) + P(x_5|\bar{x}_4)P(\bar{x}_4)$ Prob sunny fri $= 0.52 \pm 0.105 = 0.625$ $P(x_{\infty}) = 0.6$

Suppose is it cloudy on mon...

$$P(x_1) = 0$$

Prob sunny tues

$$P(x_2) = P(x_2|x_1)P(x_1) + P(x_2|\bar{x}_1)P(\bar{x}_1)$$

= 0 + 0.3 = 0.3

Prob sunny weds
$$P(x_3) = P(x_3|x_2)P(x_2) + P(x_3|\bar{x}_2)P(\bar{x}_2)$$
$$= 0.24 + 0.21 = 0.45$$

Prob sunny thurs
$$P(x_4) = P(x_4|x_3)P(x_3) + P(x_4|\bar{x}_3)P(\bar{x}_3)$$
$$= 0.36 + 0.165 = 0.53$$

Prob sunny fri

$$P(x_5) = P(x_5|x_4)P(x_4) + P(x_5|\bar{x}_4)P(\bar{x}_4)$$

= 0.424 + 0.141 = 0.565
$$P(x_{\infty}) = 0.6$$

Suppose is it cloudy on mon...

$$P(x_1) = 0$$

Prob sunny tues

$$P(x_2) = P(x_2|x_1)P(x_1) + P(x_2|\bar{x}_1)P(\bar{x}_1)$$

= 0+0.3 = 0.3

 $x_5|\bar{x}_4)P(\bar{x}_4)$

5

Prob sunny weds
$$P(x_3) = P(x_3|x_2)P(x_2) + P(x_3|\bar{x}_2)P(\bar{x}_2)$$

 $= 0.24 + 0.21 = 0.45$

Prob sunny thurs
$$P(x_4) = P(x_4|x_3)P(x_3) + P(x_4|\bar{x}_3)P(\bar{x}_3)$$
$$= 0.36 + 0.165 = 0.53$$

Prob sunny

Converge to same distribution regardless of starting point – called the "stationary distribution"

$$P(x_{\infty}) = 0.6$$

An aside: the stationary distribution

How might you calculate the stationary distribution?

Let:
$$p_t = \begin{pmatrix} p(sun) \\ p(cloudy) \end{pmatrix}$$
 $T = \begin{pmatrix} 0.8 & 0.3 \\ 0.2 & 0.7 \end{pmatrix}$

Then:
$$p_{t+1} = Tp_t$$

 $p_{t+n} = T^n p_t$

Stationary distribution is the value for p such that: p = Tp

An aside: the stationary distribution

How might you calculate the stationary distribution?

Let:
$$p_t = \begin{pmatrix} p(sun) \\ p(cloudy) \end{pmatrix}$$
 $T = \begin{pmatrix} 0.8 & 0.3 \\ 0.2 & 0.7 \end{pmatrix}$

Then:
$$p_{t+1} = Tp_t$$

 $p_{t+n} = T^n p_t$

Stationary distribution is the value for p such that: p = Tp

How calculate p that satisfies this eqn?

Hidden Markov Models (HMMs)

Hidden Markov Models:

- extension of the Markov model
- state is assumed to be "hidden"

Hidden Markov Models (HMMs)



Examples:

- speech to text; tracking in computer vision' robot localization

Hidden Markov Models (HMMs)



<u>Sensor Markov Assumption:</u> the current observation depends only on current state:

$$P(E_t | X_t, X_{t-1}, \dots, X_1) = P(E_t | X_t)$$
$$E_t \perp X_{t-1} | X_t$$

HMM example



HMM Filtering

Given a prior distribution, $P(X_1)$, and a series of observations, E_1, \ldots, E_T , calculate the posterior distribution: $P(X_t|E_1, \ldots, E_T)$

Two steps:



HMM Filtering

Given a prior distribution, $P(X_1)$, and a series of observations, E_1, \ldots, E_T , calculate the posterior distribution: $P(X_t|E_1, \ldots, E_T)$

Two steps:



HMM Filtering



Process update



This is just forward simulation of the Markov Model

Process update: example

$$B'(X_{t+1}) = \sum_{X_t} P(X_{t+1}|X_t, e_{1:t}) B(X_t)$$



If we only do the process update, then we typically lose information over time – when might this not be true?

Observation update



Observation update

 $B(X_{t+1}) = \eta P(e_{t+1}|X_{t+1})B'(X_{t+1})$

0.05	0.01	0.05	<0.01	<0.01	<0.01
0.02	0.14	0.11	0.35	<0.01	<0.01
0.07	0.03	0.05	<0.01	0.03	<0.01
0.03	0.03	<0.01	<0.01	<0.01	<0.01

Before observation

<0.01	<0.01	<0.01	<0.01	0.02	<0.01
<0.01	<0.01	<0.01	0.83	0.02	<0.01
<0.01	<0.01	0.11	<0.01	<0.01	<0.01
<0.01	<0.01	<0.01	<0.01	<0.01	<0.01

After observation

Observations enable the system to gain information

- a single observation may not determine system state exactly
- but, the more observations, the better













Prob

0

1



Prob



Prob

0

1



Prob





w_t	P(w_t)
sun	0.5
cloudy	0.5



?

?

sun

cloudy

0.5

0.5

sun

cloudy



w_t	P(w_t)	w_t	P(w_t)
sun	0.5	sun	0.55
cloudy	0.5	cloudy	0.45

X {t-1}	X t	X t
sun	sun	0.8
sun	cloudy	0.2
cloudy	sun	0.3
cloudy	cloudy	0.7

X_t	$P(g_t X_t)$
sun	0.7
cloudy	0.4



w_t	P(w_t)
sun	0.5
cloudy	0.5

w_t	P(w_t)	
sun	0.55	
cloudy	0.45	

w_t	P(w_t)
sun	?
cloudy	?

X {t-1}	X t	X t
sun	sun	0.8
sun	cloudy	0.2
cloudy	sun	0.3
cloudy	cloudy	0.7

X_t	$P(g_t X_t)$
sun	0.7
cloudy	0.4



w_t	P(w_t)	
sun	0.5	
cloudy	0.5	

w_t	P(w_t)
sun	0.55
cloudy	0.45

w_t	P(w_t)
sun	0.68
cloudy 0.31	



w_t	P(w_t)
sun	0.68
cloudy	0.31



w_t	P(w_t)
sun	0.68
cloudy	0.31





Particle Filtering



Representation: Particles

P(x) approximated by number of particles with	Particles
value x	(3,3) (2,3)
So, many x may have P(x) = 0!	(3,3) (3,2)
 More particles, more accuracy 	(3,3) (3,2) (1,2)
	(1,2) (3,3)
For now, all particles have a waight of 1	(3,3) (2,3)

Our representation of P(X) is now a list of N particles (samples)

For now, all particles have a weight of 1

Generally, N << |X|

Storing map from X to counts would defeat the point



Particle Filtering: Elapse Time

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

 $x' = \operatorname{sample}(P(X'|x))$

- 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)



Particles:

(3,3)(2,3)

(3.3)

(3.2)

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

(2.3)

Particles:

(3,2) (2,3)

(3,2)

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

(3,2)(2,2)

Particle Filtering: Observe

ignely	circitic				
Don't	sample	observ	vation,	fix	it

Slightly trickier

 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 downweighted (in fact they now sum to (N times) an approximation of P(e))



Particles:

(3,2) w=.9 (2,3) w=.2 (3,2) w=.9 (3,1) w=.4 (3,3) w=.4

(3.2) w = .9

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

(3.2) w = .9

(2.2) w = .4



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







Recap: Particle Filtering

Particles: track samples of states rather than an explicit distribution



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





Particle Filter Localization (Sonar)



Particle Filter Localization (Laser)



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

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: G₁^a = (3,3) G₁^b = (5,3)
- Elapse time: Sample a successor for each particle
 - Example successor: G₂^a = (2,3) G₂^b = (6,3)
- Observe: Weight each <u>entire</u> sample by the likelihood of the evidence conditioned on the sample
 - Likelihood: $P(\mathbf{E_1^a} | \mathbf{G_1^a}) * P(\mathbf{E_1^b} | \mathbf{G_1^b})$
- Resample: Select prior samples (tuples of values) in proportion to their likelihood