# CS221 Problem Workout

Week 8

Stanford University

# Introduction

#### Jenn Grannen



General OH: Mondays HW OH: Fridays 9:30-11:00 Bytes + Zoom 9:30-11:00 Huang Basement

#### Jeremy Kim



 HW OH: Wednesdays
 3:30-5:30 Huang

 HW OH: Sunday
 2:00-3:00 Online

# Outline

# HMM Review

# • Bayesian networks: Learning

- Maximum likelihood
- Smoothing
- EM Algorithm
- Problem discussion

# Hidden Markov Model (HMM) Review

**Defining HMMs** 

### **Defining HMMs - Ice Cream Example**

Famous problem by Jason Eisner (2002) where you want to predict if a day was COLD or HOT (your hidden states) based on records of the # of ice creams (your known evidence) Eisner ate that day.



- $S = \{s_1...s_N\}$ , N states (2 states: cold or hot)
- $A = a_{11}...a_{ij}...a_{NN}$ , transition probabilities (e.g. cold  $\rightarrow$  hot?)
- $B = b_i(o_t)$ , emission probabilities (e.g. 3 ice creams  $\rightarrow$  hot?)
- $\pi = \{\pi_1 ... \pi_N\}$ , initial probabilities (e.g. start  $\rightarrow$  hot?)

## **Defining HMMs - Ice Cream Example**

Famous problem by Jason Eisner (2002) where you want to predict if a day was COLD or HOT (your hidden states) based on records of the # of ice creams (your known evidence) Eisner ate that day.



HMM problem to motivate the Forward Algorithm:

 Given HMM λ (like above), what is the probability P(O|λ) of a specific observation sequence O (evidence e.g. 3 1 3)?

# Hidden Markov Model (HMM) Review

**The Forward Algorithm** 



HMM problem to motivate the Forward Algorithm:

 Given HMM λ (like above), what is the probability P(O|λ) of a specific observation sequence O (evidence e.g. 3 1 3)?

First, consider an easier problem: suppose our states are not hidden (we just have a "Markov model") and we have Q = (hot hot cold). What is the probability (aka likelihood) of O = 3, 1, 3?



First, consider an easier problem: suppose our states are not hidden (we just have a "Markov model") and we have Q = (hot hot cold).

What is the probability (aka likelihood) of O = 3, 1, 3?

$$P(O|Q) = \prod_t^T P(o_t|q_t)$$
  
 $P(3 \ 1 \ 3| ext{hot cold}) = P(3| ext{hot})P(1| ext{hot})P(3| ext{cold})$ 



Simplification: probability O = 3, 1, 3 given Q = (hot hot cold)?  $P(3 \ 1 \ 3|hot hot cold) = P(3|hot)P(1|hot)P(3|cold)$ 





Simplification: probability O = 3, 1, 3 given Q = (hot hot cold)?  $P(3 \ 1 \ 3|hot hot cold) = P(3|hot)P(1|hot)P(3|cold)$ 

Back to the original problem: we don't know the actual weather sequence – it's a HIDDEN Markov model!

What is the probability of 3 1 3 given the HMM?



Back to the original problem: we don't know the actual weather sequence – it's a HIDDEN Markov model!

What is the probability of 3 1 3 given the HMM?

Brute Force: Sum over all possible weather sequences:  $P(3 \ 1 \ 3, \text{ cold cold cold})? P(3 \ 1 \ 3, \text{ hot cold cold})? P(3 \ 1 \ 3, \text{ hot hot cold})? etc...?$ Then add them all together...



What is the probability of  $3\ 1\ 3$  given the HMM?

Brute Force: Sum over all possible weather sequences:  $P(3 \ 1 \ 3, \text{ cold cold })? P(3 \ 1 \ 3, \text{ hot cold cold})? etc...?$ 

P(O, Q) is the joint probability:

$$P(O, Q) = P(O|Q)P(Q) = \prod_{t}^{T} P(o_t|q_t) \prod_{t}^{T} P(q_t|q_{t-1})$$



P(O, Q) is the joint probability:

$$P(O, Q) = P(O|Q)P(Q) = \prod_{t}^{T} P(o_t|q_t) \prod_{t}^{T} P(q_t|q_{t-1})$$

Example: joint probability of  $O = 3 \ 1 \ 3$  and Q = hot hot cold P(O, Q) = P(3|hot)P(1|hot)P(3|cold)P(hot|start)P(hot|hot)P(cold|hot)



Example: joint probability of  $O = 3 \ 1 \ 3$  and Q = hot hot cold P(O, Q) = P(3|hot)P(1|hot)P(3|cold)P(hot|start)P(hot|hot)P(cold|hot)





What is the probability of 3 1 3 given the HMM?

Brute Force: Sum over all possible weather sequences:

 $P(3 \ 1 \ 3, \text{ cold cold cold}) + P(3 \ 1 \ 3, \text{ hot cold cold}) + P(3 \ 1 \ 3, \text{ hot hot cold}) + ...$ 

This is a  $N^T$  operation with N states and T observations!

Not efficient for more complex problems!



What is the probability of 3 1 3 given the HMM?

 $P(3 \ 1 \ 3, \text{ cold cold } cold) + P(3 \ 1 \ 3, \text{ hot cold cold}) + P(3 \ 1 \ 3, \text{ hot hot cold}) + ...$ 

This is a  $N^T$  operation with N states and T observations!

Forward Algorithm does this in  $O(N^2T)$  via dynamic programming!







Formally, for each cell  $\alpha_t(j)$  in our lattice structure, we compute

$$\alpha_t(j) = \sum_i^N \alpha_{t-1}(i) a_{ij} b_j(o_t)$$

and the probability of sequence 3 1 3 is at the end

$$\mathsf{P}(O|\lambda) = \sum_{i}^{N} lpha_{\mathcal{T}}(i)$$

# Hidden Markov Model (HMM) Review

**Relating back to Lecture** 

## The Problem of Filtering





Problem of Filtering: what is the distribution of a hidden state  $H_i$  based on the observations aka evidence (*E* in lecture) so far? Check your understanding: what is the distribution of  $q_2$  in the ice cream example given observations *O*:  $o_1 = 3$  and  $o_2 = 1$ ?

## The Problem of Filtering





Problem of Filtering: what is the distribution of  $q_2$  in the ice cream example given observations O:  $o_1 = 3$  and  $o_2 = 1$ ?

$$P(q_2 = \mathsf{H} \mid o_1, o_2) = rac{0.0404}{0.0404 + 0.069}$$
,  $P(q_2 = \mathsf{C} \mid o_1, o_2) = rac{0.069}{0.0404 + 0.069}$ 





Problem of Smoothing: what is the distribution of a hidden state  $H_i$  based ALL observations aka evidence from start to end? Forward Algorithm is not enough! What if hypothetically a later transition is 0?



For Smoothing, need Forward AND Backward passes!

- Forward: compute  $\alpha_t(i)$  or F from lecture.
- Backward: compute  $\beta_t(i)$  or B from lecture.
- Define S = FB, that is for each cell in the lattice, multiply the forward and backward results together.

What happens now if there is a 0 along the backward pass?



- Forward: compute  $\alpha_t(i)$  or F from lecture.
- Backward: compute  $\beta_t(i)$  or B from lecture.
- Define S = FB

Suppose  $\beta_2(1) = 0.03$ ,  $\beta_2(2) = 0.02$  (made up numbers). What is the distribution of  $q_2$  given all observations *O*?



Suppose  $\beta_2(1) = 0.03$ ,  $\beta_2(2) = 0.02$  (made up numbers). What is the distribution of  $q_2$  given all observations O?

$$P(q_2 = H | O) = \frac{0.0404 * 0.02}{0.0404 * 0.02 + 0.069 * 0.03}$$
$$P(q_2 = C | O) = \frac{0.069 * 0.03}{0.0404 * 0.02 + 0.069 * 0.03}$$

Check back on the lecture slides to make sure you see the parallel!

# Hidden Markov Model (HMM) Review

**Particle Filtering** 

## **Motivation for Particle Filtering**



For T observations and N possible states (i.e. |domain| = N), the Forward-Backward Algorithm is  $O(2 * N^2 T) \rightarrow O(N^2 T)$ .

This can still be slow if N is large! Or consider if the domain is based on a continuous function, e.g. instead of just hot or cold, we have to consider a spectrum of floating point temperatures [0, 100].



### Big idea of Particle Filtering: introduce sampling!

1. First, we propose assignments aka particles to each hidden state by sampling from the transition probabilities.

Example: proposing a value for  $q_1$  involves sampling from P(H|start) = 0.8 and P(C|start) = 0.2, i.e. we have an 80% chance to pick hot, 20% chance to pick cold.



2 Second, we weight each assignment by the emission probabilities.

Example: suppose we have 3 particles of  $q_1 = H$ ,  $q_1 = H$ ,  $q_1 = C$ , and we have the observation  $o_1 = 1$ . Then the weights of our particles are P(1|H) = 0.2, P(1|H) = 0.2, P(1|C) = 0.5 respectively.



- 1. First, we propose assignments aka particles to each hidden state by sampling from the transition probabilities.
- 2. Second, we weight each assignment by the emission probabilities.
- 3. Third, we resample new assignments from the particles based on the weight distributions.

3 Third, we resample new assignments from the particles based on the weight distributions.

Example: suppose we have 3 particles of  $q_1 = H$ ,  $q_1 = H$ ,  $q_1 = C$ , and we have the observation  $o_1 = 1$ .

Then the weights of our particles are P(1|H) = 0.2, P(1|H) = 0.2, P(1|H) = 0.2, P(1|C) = 0.5 respectively.

Now to resample, we have the distribution:

• 
$$P(q_1 \to H) = \frac{0.2}{0.2 + 0.2 + 0.5}$$

• 
$$P(q_1 \to H) = \frac{0.2}{0.2 + 0.2 + 0.5}$$

• 
$$P(q_1 \to C) = \frac{0.5}{0.2 + 0.2 + 0.5}$$

Notice how even though our initial proposal had a higher chance to pick  $q_1 = H$ , we now have a higher chance to get  $q_1 = C$ ! The resampling takes into account the observations!



Suppose after all that, we have new assignments for our 3 particles:  $q_1 = C$ ,  $q_1 = C$ ,  $q_1 = H$ ... And repeated the propose process for  $q_2$  to get:  $(q_1, q_2) = (C, H)$ ;  $(q_1, q_2) = (C, C)$ ;  $(q_1, q_2) = (H, C)$  with  $o_2 = 3$ ...



And repeated the propose process for  $q_2$  to get:  $(q_1, q_2) = (C, H)$ ;  $(q_1, q_2) = (C, C)$ ;  $(q_1, q_2) = (H, C)$  with  $o_2 = 3...$ 

The weight process then assigns the particles:

- $(q_1, q_2) = (C, H)$ : P(3|H) = 0.4
- $(q_1, q_2) = (C, C)$ : P(3|C) = 0.1
- $(q_1, q_2) = (H, C)$ : P(3|C) = 0.1

The weight process then assigns the particles:

- $(q_1, q_2) = (C, H)$ : P(3|H) = 0.4
- $(q_1, q_2) = (C, C)$ : P(3|C) = 0.1
- $(q_1, q_2) = (H, C)$ : P(3|C) = 0.1

And the resample process then samples from the above 3 options, that is:

- $(q_1, q_2) = (C, H)$  has a 4/6 chance of being picked.
- The other two each have a 1/6 chance of being picked.

And so a possible resampling result might yield the particles:  $(q_1, q_2) = (C, H), (C, H), \text{ and } (C, C).$ And you'd repeat the process with  $q_3...$ 

# Outline

HMM Review

# • Bayesian networks: Learning

- Maximum likelihood
- Smoothing
- EM Algorithm
- Problem discussion

# Bayesian networks: Learning

Given local probability distributions, i.e. P(x | parents(x))

Find conditional P(Q | E=e)

Inference

Given observations / samples

Find the local distributions, i.e. P(x | Parents(x))

Learning

# Example

Variables:

- Genre  $G \in \{ drama, comedy \}$
- Rating  $R \in \{1, 2, 3, 4, 5\}$

$$\bigcirc G \longrightarrow R \qquad \mathbb{P}(G = g, R = r) = p_G(g)p_R(r \mid g)$$

 $\mathcal{D}_{\mathsf{train}} = \{(\mathsf{d},4), (\mathsf{d},4), (\mathsf{d},5), (\mathsf{c},1), (\mathsf{c},5)\}$ 

Parameters:  $\theta = (p_G, p_R)$ 

Example borrowed from lecture slides

# Outline

- Bayesian networks: Learning
  - Maximum likelihood
  - Smoothing
  - EM Algorithm
- Problem discussion

# Maximum likelihood

Input: training examples  $\mathcal{D}_{\text{train}}$  of full assignments

```
Output: parameters \theta = \{p_d : d \in D\}
```



Slide borrowed from lecture slides

# Outline

- Bayesian networks: Learning
  - Maximum likelihood
  - Smoothing
  - EM Algorithm
- Problem discussion

# Smoothing

Why?

• What if count is 0? Should P be 0?

How to solve?

• Initialize all counts with a non-zero constant  $\boldsymbol{\lambda}$ 

Observations

- Larger  $\lambda$  -> more uniform distributions, less influenced by data
- Smaller  $\lambda$  -> more influenced by data
- Infinite data -> effect of  $\lambda$  vanishes

# Final algorithm

Input: training examples  $\mathcal{D}_{\text{train}}$  of full assignments

```
Output: parameters \theta = \{p_d : d \in D\}
```



Slide modified from lecture slides

# Outline

# • Bayesian networks: Learning

- Maximum likelihood
- Smoothing
- EM Algorithm
- Problem discussion

# EM Algorithm

Variables: H is hidden, E = e is observed Example:

$$G$$
  
 $H = G$   $E = (R_1, R_2)$   $e = (1, 2)$   
 $\theta = (p_G, p_R)$ 

Maximum marginal likelihood objective:

$$\begin{split} & \max_{\theta} \prod_{e \in \mathcal{D}_{\text{train}}} \mathbb{P}(E = e; \theta) \\ & = \max_{\theta} \prod_{e \in \mathcal{D}_{\text{train}}} \sum_{h} \mathbb{P}(H = h, E = e; \theta) \end{split}$$

Slide borrowed from lecture slides

# EM Algorithm

Initialize  $\boldsymbol{\theta}$  randomly

Until convergence:

## # E Step

for each e in Data:

for each h:

 $q(h; e) = P(H = h | E = e; \theta) \dots$  inference

# Update table from {e, count(e)} to {(h,e), (q(h; e) x count(e)}

# Now no variables are hidden

## # M step

update(θ) using Table {(h,e), (q(h; e) x count(e)} ... MLE

# Summary

• Given data learn the parameters of bayesian net

MLE  $p \propto count(x_i;$  $parents(x_i))$  **Smoothing**  $p \propto count(x_i;$ parents(x\_i)) +  $\lambda$  EM Data is *incomplete E step:* compute counts *M step:* MLE

# Outline

- HMM Review
- Bayesian networks: Learning
  - Maximum likelihood
  - Smoothing
  - EM Algorithm
- Problem discussion

# Problem: P2, Winter 2021 Exam 2



# How does the bayesian net look?



# Learning using EM algorithm

- Data = {H, H, T}
- $\lambda_0$  and  $\lambda$  are given. To find:  $p_{\chi}$  and  $p_{\gamma}$
- Why do we need EM?
  - $\circ$  What is not observed?
  - C<sub>i</sub> is not observed
- How do we use EM?
  - Compute  $q(c_i)$  using  $p'_x$  and  $p'_y$
  - Use ML to update  $p'_{\chi}$  and  $p'_{\gamma}$

# Given q's compute update

	T=1	T=2	T=3
Х	0.1	0.5	0.3
Y	0.9	0.5	0.7

Data = {H, H, T}

Compute:

C <sub>i</sub>	O <sub>i</sub>	Count
?	?	?

# Given q's compute update

	T=1	T=2	T=3
X	0.1	0.5	0.3
Y	0.9	0.5	0.7

C <sub>i</sub>	O <sub>i</sub>	Count
x	Н	0.1
X	н	0.5
X	Т	0.3

$$P(H \mid X) = p'_{X}$$
  
= (0.1 + 0.5) / (0.1 + 0.5 + 0.3) ... MLE

# Thank You