

# Learning in HMMs and Bayes Nets

Ron Parr  
CPS 271

With some content courtesy of Lise Getoor

## The Usual Caveats

- This is only a very brief overview
- A proper treatment takes:
  - More time than we have
  - More patience than you have this late in the semester ☺

## First: HMMs

- The easy case:
  - Suppose you have a complete record of
    - Underlying states
    - Observations
- What would you do???

## Estimate Probabilities

- Estimate probabilities from relative frequencies
- $P(s_j | s_i) = \#(s_j \text{ to } s_i \text{ transitions}) / \#(s_i \text{ occurrences})$
- $P(o_i | s_j) = \#(o_i \text{ observations in } s_j) / \#(s_j \text{ occurrences})$
- Why is this valid???

  - Markov assumption
  - Relative frequency is max likelihood solution

## What if Underlying States not Known?

- All we see is trajectories of observations
- True states are hidden variables
- EM to the rescue!
- E step:
  - Use forward-backward/variable elimination to estimate  $P(s_i, s_j)$  for all pairs, for all trajectories
- M step:
  - Compute,  $P(s_i | s_j)$ ,  $P(o_i | s_j)$  given results of E-step
  - Note: For soft EM, this means adding fractional state transitions

## Getting the pairwise probabilities

- Forward and backward steps are already doing most of the work for this:

$$\begin{aligned}
 P(s_i | e_{1..e_0}) &= \frac{P(e_i | S_i, e_{1..e_0})P(S_i | e_{1..e_0})}{P(e_i | e_{1..e_0})} & P(S_i | e_{1..e_0}) &= \alpha P(e_{1..e_{i-1}} | S_{1..i-1}, e_{1..e_0})P(S_i | e_{1..e_0}) \\
 &= \alpha P(e_i | S_i, e_{1..e_0})P(S_i | e_{1..e_0}) & &= \alpha P(e_{1..e_{i-1}} | S_{1..i-1})P(S_i | e_{1..e_0}) \\
 &= \alpha P(e_i | S_i)P(S_i | e_{1..e_0}) & P(e_{1..e_{i-1}} | S_{1..i-1}) &= \sum_{s_{1..i-1}} P(e_{1..e_{i-1}} | S_{1..i-1}, s_{1..i-1})P(S_{1..i-1} | s_{1..i-1}) \\
 &= \alpha P(e_i | S_i) \sum_{s_{i-1}} P(S_i | S_{i-1})P(S_{i-1} | e_{1..e_0}) & &= \sum_{s_{i-1}} P(e_{1..e_{i-1}} | S_{1..i-1}, s_{i-1})P(S_{i-1} | S_i) \\
 & & &= \sum_{s_{i-1}} P(e_{1..e_{i-1}} | S_{1..i-1})P(e_{i-1} | s_{i-1}, S_{i-1})P(S_{i-1} | S_i)
 \end{aligned}$$

- Also expressed as  $p(x_i, x_j)_t \propto f(x_i)_t P(x_j | x_i)_t b(x_j)_t$
- Smoothed estimate is the normalized product of the forward and backward estimates

## EM for HMMs

- Converges to local optimum
- Most commonly used method for HMM learning
- Used for speech recognition, target detection, tracking, etc.
- How to pick number of states?
  - Just guess?
  - Just pick a big number?
  - When does it matter?
  - Can cross validation help?

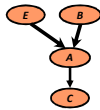
## Learning Bayesian Networks: Known Structure Case

- Given a network structure  $G$
- Learn parameters for network

### Goal

- Construct a network that is “closest” to probability that generated the data

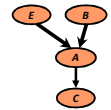
## Learning Parameters for a Bayesian Network



- Training data have the form:

$$D = \begin{bmatrix} E[1] & B[1] & A[1] & C[1] \\ \vdots & \vdots & \vdots & \vdots \\ E[M] & B[M] & A[M] & C[M] \end{bmatrix}$$

## Learning Parameters for a Bayesian Network

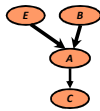


- Since we assume i.i.d. samples, likelihood function is

$$\mathcal{L}(\Theta; D) = \prod_m P(E[m], B[m], A[m], C[m]; \Theta)$$

↑  
Network parameters (CPTs)

## Learning Parameters for a Bayesian Network

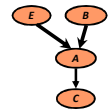


- By definition of Bayes net, we get

$$\begin{aligned} \mathcal{L}(\Theta; D) &= \prod_m P(E[m], B[m], A[m], C[m]; \Theta) \\ &= \prod_m P(E[m]; \Theta) P(B[m]; \Theta) P(A[m] | B[m], E[m]; \Theta) P(C[m] | A[m]; \Theta) \end{aligned}$$

$$\begin{bmatrix} E[1] & B[1] & A[1] & C[1] \\ \vdots & \vdots & \vdots & \vdots \\ E[M] & B[M] & A[M] & C[M] \end{bmatrix}$$

## Learning Parameters for a Bayesian Network



- Rewriting terms, we get

$$\begin{aligned} \mathcal{L}(\Theta; D) &= \prod_m P(E[m], B[m], A[m], C[m]; \Theta) \\ &= \prod_m P(E[m]; \Theta) \prod_m P(B[m]; \Theta) \prod_m P(A[m] | B[m], E[m]; \Theta) \prod_m P(C[m] | A[m]; \Theta) \end{aligned}$$

$$\begin{bmatrix} E[1] & B[1] & A[1] & C[1] \\ \vdots & \vdots & \vdots & \vdots \\ E[M] & B[M] & A[M] & C[M] \end{bmatrix}$$

## General Bayesian Networks

Generalizing for any Bayesian network:

$$\begin{aligned}
 \mathcal{L}(\Theta : D) &= \prod_m \mathcal{P}(x_1[m], \dots, x_n[m] : \Theta) && \text{i.i.d. samples} \\
 &= \prod_m \prod_i \mathcal{P}(x_i[m] | \rho_{a_i}[m] : \Theta_i) && \text{Network factorization} \\
 &= \prod_i \prod_m \mathcal{P}(x_i[m] | \rho_{a_i}[m] : \Theta_i) \\
 &= \prod_i \mathcal{L}_i(\Theta_i : D)
 \end{aligned}$$

- The likelihood **decomposes** according to the structure of the network.
- Bottom line: Optimize each CPT individually

## Learning BNs with Missing Data

- Some variable values may be missing
- Assume no systematic pattern to omissions
- And the trick is...
- You guessed it, EM!
- E step computes distribution over missing vars
- M step does “fully observable” Bayes net learning (using results from E step)

## Learning w/unknown Structure

- Great application: Learning structure of biological regulatory networks
- Recall that we typically want to maximize:

$$\frac{P(D|\theta)P(\theta)}{P(D)}$$

- Our M step previously maximized  $P(D|\theta)$
- $\theta$  now includes the space of models
- Whis is this more complicated now?
  - Closed form solution? (kind of)
  - Need for regularization

## Need for Regularization

- Fully connected model will almost always have higher likelihood if  $P(\theta)$  is uniform!
- Typically introduce a structural prior to penalize structures that are complex (e.g. high parent count)
- This complicates search:
  - No closed form solution for best BN
  - Usually do some kind of local search

$$\frac{P(D|\theta)P(\theta)}{P(D)}$$

## Conclusions

- Learning w/known structure, full observability is easy
- Learning with partial observability is trickier
  - EM is our friend
  - We like EM
- Learning with unknown structure is quite tricky