

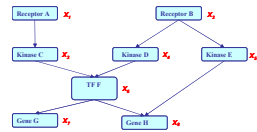
# Hidden Markov Model and Conditional Random Fields

## Probabilistic Graphical Models (10-708)

Lecture 12, Oct 29, 2007

Eric Xing

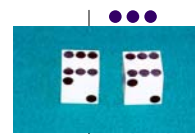
Reading: J-Chap. 12, and addition papers



1

- Feedbacks
  - Today
  - Office hour
- Recitation
- Project milestone
  - This Wednesday

# Hidden Markov Model revisit



- Transition probabilities between any two states

$$p(y_t^j = 1 | y_{t-1}^i = 1) = a_{i,j},$$

or  $p(y_t | y_{t-1}^i = 1) \sim \text{Multinomial}(a_{i,1}, a_{i,2}, \dots, a_{i,M}), \forall i \in I.$

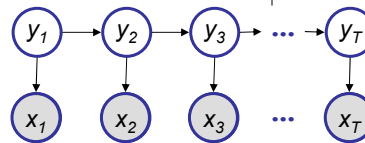
- Start probabilities

$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, \dots, \pi_M).$$

- Emission probabilities associated with each state

$$p(x_t | y_t^i = 1) \sim \text{Multinomial}(b_{i,1}, b_{i,2}, \dots, b_{i,K}), \forall i \in I.$$

or in general:  $p(x_t | y_t^i = 1) \sim f(\cdot | \theta_i), \forall i \in I.$



Eric Xing

3

# Applications of HMMs



- **Some early applications of HMMs**

- finance, but we never saw them
- speech recognition
- modelling ion channels

- In the mid-late 1980s HMMs entered genetics and molecular biology, and they are now firmly entrenched.

- **Some current applications of HMMs to biology**

- mapping chromosomes
- aligning biological sequences
- predicting sequence structure
- inferring evolutionary relationships
- finding genes in DNA sequence

Eric Xing

4



10

$$p(\bullet|y) = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

[illegible]

# The Forward Algorithm

- We want to calculate  $P(\mathbf{x})$ , the likelihood of  $\mathbf{x}$ , given the HMM

- Sum over all possible ways of generating  $\mathbf{x}$ :

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) = \sum_{y_1} \sum_{y_2} \cdots \sum_{y_N} \pi_{y_1} \prod_{t=2}^T a_{y_{t-1}, y_t} \prod_{t=1}^T p(x_t | y_t)$$

- To avoid summing over an exponential number of paths  $\mathbf{y}$ , define

$$\alpha(y_t^k = 1) = \alpha_t^k \stackrel{\text{def}}{=} P(x_1, \dots, x_t, y_t^k = 1) \quad (\text{the forward probability})$$

- The recursion:

$$\alpha_t^k = p(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

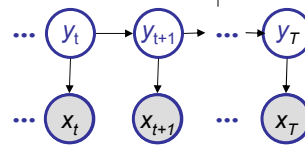
$$P(\mathbf{x}) = \sum_k \alpha_T^k$$

Eric Xing

7

# The Backward Algorithm

- We want to compute  $P(y_t^k = 1 | \mathbf{x})$ , the posterior probability distribution on the  $t^{\text{th}}$  position, given  $\mathbf{x}$



- We start by computing

$$P(y_t^k = 1, \mathbf{x}) = P(x_1, \dots, x_t, y_t^k = 1, x_{t+1}, \dots, x_T)$$

$P(\mathbf{x})$

$$= P(x_1, \dots, x_t, y_t^k = 1) P(x_{t+1}, \dots, x_T | x_1, \dots, x_t, y_t^k = 1)$$

$$= P(x_1 \dots x_t, y_t^k = 1) P(x_{t+1} \dots x_T | y_t^k = 1)$$



Forward,  $\alpha_t^k$

Backward,  $\beta_t^k = P(x_{t+1}, \dots, x_T | y_t^k = 1)$

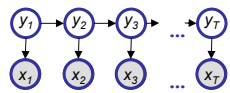
- The recursion:  $\beta_t^k = \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i$

Eric Xing

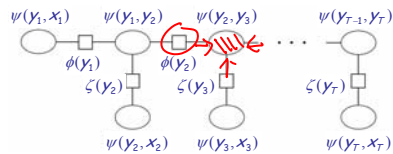
8

# The Junction tree algorithm

- A junction tree for the HMM



⇒



- Rightward pass

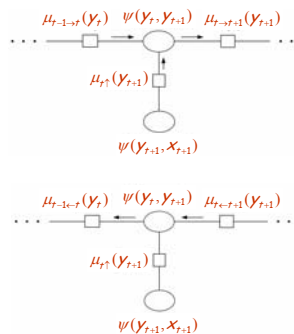
$$\begin{aligned}\mu_{t \rightarrow t+1}(y_{t+1}) &= \sum_{y_t} \psi(y_t, y_{t+1}) \mu_{t-1 \rightarrow t}(y_t) \mu_t^\uparrow(y_{t+1}) \\ &= \sum_{y_t} p(y_{t+1} | y_t) \mu_{t-1 \rightarrow t}(y_t) p(x_{t+1} | y_{t+1}) \\ &= p(x_{t+1} | y_{t+1}) \sum_{y_t} a_{y_t, y_{t+1}} \mu_{t-1 \rightarrow t}(y_t)\end{aligned}$$

- This is exactly the **forward algorithm**!

- Leftward pass ...

$$\begin{aligned}\mu_{t-1 \leftarrow t}(y_t) &= \sum_{y_{t+1}} \psi(y_t, y_{t+1}) \mu_{t \leftarrow t+1}(y_{t+1}) \mu_t^\uparrow(y_{t+1}) \\ &= \sum_{y_{t+1}} p(y_{t+1} | y_t) \mu_{t \leftarrow t+1}(y_{t+1}) p(x_{t+1} | y_{t+1})\end{aligned}$$

- This is exactly the **backward algorithm**!



Eric Xing

9

# Summary of the F-B algorithm

$$\alpha_t^k \stackrel{\text{def}}{=} \mu_{t-1 \rightarrow t}(k) = P(x_1, \dots, x_{t-1}, x_t, y_t^k = 1)$$

$$\alpha_t^k = p(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

$$\beta_t^k \stackrel{\text{def}}{=} \mu_{t \leftarrow t+1}(k) = P(x_{t+1}, \dots, x_T | y_t^k = 1)$$

$$\beta_t^k = \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i$$

$$\gamma_t^i \stackrel{\text{def}}{=} p(y_t^i = 1 | x_{1:T}) \propto \alpha_t^i \beta_t^i = \sum_j \xi_t^{i,j}$$

$$\xi_t^{i,j} \stackrel{\text{def}}{=} p(y_t^i = 1, y_{t+1}^j = 1, x_{1:T})$$

$$\propto \mu_{t-1 \rightarrow t}(y_t^i = 1) a_{i,j} \mu_{t \leftarrow t+1}(y_{t+1}^j = 1) p(x_{t+1} | y_{t+1}^j) p(y_{t+1} | y_t^i)$$

$$\xi_t^{i,j} = \alpha_{t-1}^i \beta_{t+1}^j a_{i,j} p(x_{t+1} | y_{t+1}^j = 1)$$

The matrix-vector form:

$$\begin{aligned}\alpha_t &\stackrel{\text{def}}{=} p(x_t | y_t = 1) \quad \alpha_t = (A^T \alpha_{t-1}) * B_t \\ \beta_t &\stackrel{\text{def}}{=} p(y_{t+1}^j = 1 | y_t = 1) \quad \beta_t = A(\beta_{t+1} * B_{t+1}) \\ \gamma_t &\stackrel{\text{def}}{=} p(y_t = 1 | x_{1:T}) \quad \gamma_t = \alpha_t (\beta_{t+1} * B_{t+1})^T * A \\ \gamma_t &= \alpha_t * \beta_t\end{aligned}$$

Eric Xing

10

## Posterior decoding



- We can now calculate

$$P(y_t^k = 1 | \mathbf{x}) = \frac{P(y_t^k = 1, \mathbf{x})}{P(\mathbf{x})} = \frac{\alpha_t^k \beta_t^k}{P(\mathbf{x})}$$

- Then, we can ask

- What is the most likely state at position  $t$  of sequence  $\mathbf{x}$ :

$$k_t^* = \arg \max_k P(y_t^k = 1 | \mathbf{x})$$

- Note that this is an MPA of a **single** hidden state, what if we want a MPA of a whole hidden state sequence?

- Posterior Decoding:  $\{y_t^{k_t^*} = 1 : t = 1 \dots T\}$

- This is different from MPA of a **whole sequence** of hidden states

- This can be understood as **bit error rate** vs. **word error rate**

Example:  
MPA of  $\mathbf{x}$ ?  
MPA of  $(\mathbf{x}, \mathbf{y})$ ?

$x$	$y$	$P(x, y)$
0	0	0.35
0	1	0.05
1	0	0.3
1	1	0.3

Eric Xing

11

## Viterbi decoding

- GIVEN  $\mathbf{x} = x_1, \dots, x_T$ , we want to find  $\mathbf{y} = y_1, \dots, y_T$ , such that  $P(\mathbf{y} | \mathbf{x})$  is maximized:

- Let  $y^* = \arg \max_{\mathbf{y}} P(\mathbf{y} | \mathbf{x}) = \arg \max_{\mathbf{y}} P(\mathbf{y}, \mathbf{x})$

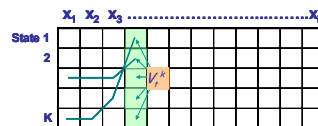
$$V_t^k = \max_{\{y_1, \dots, y_{t-1}\}} P(x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t, y_t^k = 1)$$

= Probability of most likely **sequence of states** ending at state  $y_t = k$

- The recursion:

$$V_t^k = p(x_t | y_t^k = 1) \max_i a_{i,k} V_{t-1}^i$$

- Underflows are a significant problem



$$p(x_1, \dots, x_t, y_1, \dots, y_t) = \pi_{y_1} a_{y_1, y_2} \dots a_{y_{t-1}, y_t} b_{y_t, x_t}$$

- These numbers become extremely small – **underflow**
- Solution: Take the logs of all values:  $V_t^k = \log p(x_t | y_t^k = 1) + \max_i (\log(a_{i,k}) + V_{t-1}^i)$

Eric Xing

12

## The Viterbi Algorithm – derivation



- Define the viterbi probability:

$$\begin{aligned}
 V_{t+1}^k &= \max_{\{y_1, \dots, y_t\}} P(x_1, \dots, x_t, y_1, \dots, y_t, x_{t+1}, y_{t+1}^k = 1) \\
 &= \max_{\{y_1, \dots, y_t\}} P(x_{t+1}, y_{t+1}^k = 1 | x_1, \dots, x_t, y_1, \dots, y_t) P(x_1, \dots, x_t, y_1, \dots, y_t) \\
 &= \max_{\{y_1, \dots, y_t\}} P(x_{t+1}, y_{t+1}^k = 1 | y_t) P(x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t, y_t) \\
 &= \max_i P(x_{t+1}, y_{t+1}^k = 1 | y_t^i = 1) \max_{\{y_1, \dots, y_{t-1}\}} P(x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t, y_t^i = 1) \\
 &= \max_i P(x_{t+1}, y_{t+1}^k = 1) a_{i,k} V_t^i \\
 &= P(x_{t+1}, y_{t+1}^k = 1) \max_i a_{i,k} V_t^i
 \end{aligned}$$

Eric Xing

13

## Computational Complexity and implementation details



- What is the running time, and space required, for Forward, and Backward?

$$\begin{aligned}
 \alpha_t^k &= p(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k} \\
 \beta_t^k &= \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i \\
 V_t^k &= p(x_t | y_t^k = 1) \max_i a_{i,k} V_{t-1}^i
 \end{aligned}$$

Time:  $O(K^2M)$ ;

Space:  $O(KM)$ .

- Useful implementation technique to avoid underflows
  - Viterbi: sum of logs
  - Forward/Backward: rescaling at each position by multiplying by a constant

Eric Xing

14

# Learning HMM



- **Supervised learning:** estimation when the “right answer” is known
  - **Examples:**
    - GIVEN:** a genomic region  $x = x_1 \dots x_{1,000,000}$  where we have good (experimental) annotations of the CpG islands
    - GIVEN:** the casino player allows us to observe him one evening, as he changes dice and produces 10,000 rolls
- **Unsupervised learning:** estimation when the “right answer” is unknown
  - **Examples:**
    - GIVEN:** the porcupine genome; we don't know how frequent are the CpG islands there, neither do we know their composition
    - GIVEN:** 10,000 rolls of the casino player, but we don't see when he changes dice
- **QUESTION:** Update the parameters  $\theta$  of the model to maximize  $P(x|\theta)$  -  
-- Maximal likelihood (ML) estimation

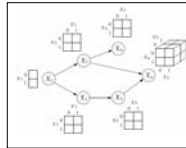
Eric Xing

15

# Learning HMM: two scenarios



- Supervised learning: if only we knew the true state path then ML parameter estimation would be trivial
  - E.g., recall that for complete observed tabular BN:
 



$$\theta_{ijk}^{ML} = \frac{n_{ijk}}{\sum_{i,j,k} n_{ijk}}$$

$$a_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^T y_{n,t-1}^i y_{n,t}^j}{\sum_n \sum_{t=2}^T y_{n,t-1}^i}$$

$$b_{jk}^{ML} = \frac{\#(i \rightarrow k)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=1}^T y_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^T y_{n,t}^i}$$
- What if  $y$  is continuous? We can treat  $\{(x_{n,t}, y_{n,t}) : t=1:T, n=1:N\}$  as  $N \times T$  observations of, e.g., a GLIM, and apply learning rules for GLIM ...
- Unsupervised learning: when the true state path is unknown, we can fill in the missing values using inference recursions.
  - The Baum Welch algorithm (i.e., EM)
    - Guaranteed to increase the log likelihood of the model after each iteration
    - Converges to local optimum, depending on initial conditions

Eric Xing

16



# The Baum Welch algorithm

- The complete log likelihood

$$\mathcal{L}_c(\theta; \mathbf{x}, \mathbf{y}) = \log p(\mathbf{x}, \mathbf{y}) = \log \prod_n \left( p(y_{n,1}) \prod_{t=2}^T p(y_{n,t} | y_{n,t-1}) \prod_{t=1}^T p(x_{n,t} | x_{n,t}) \right)$$

- The expected complete log likelihood

$$\langle \mathcal{L}_c(\theta; \mathbf{x}, \mathbf{y}) \rangle = \sum_n \left( \langle y_{n,1}^i \rangle_{p(y_{n,1}^i | \mathbf{x}_n)} \log \pi_i \right) + \sum_n \sum_{t=2}^T \left( \langle y_{n,t-1}^i y_{n,t}^j \rangle_{p(y_{n,t-1}^i, y_{n,t}^j | \mathbf{x}_n)} \log a_{i,j} \right) + \sum_n \sum_{t=1}^T \left( \langle x_{n,t}^k y_{n,t}^i \rangle_{p(y_{n,t}^i | \mathbf{x}_n)} \log b_{i,k} \right)$$

- EM

- The E step

$$\gamma_{n,t}^i = \langle y_{n,t}^i \rangle = p(y_{n,t}^i = 1 | \mathbf{x}_n)$$

$$\xi_{n,t}^{i,j} = \langle y_{n,t-1}^i y_{n,t}^j \rangle = p(y_{n,t-1}^i = 1, y_{n,t}^j = 1 | \mathbf{x}_n)$$

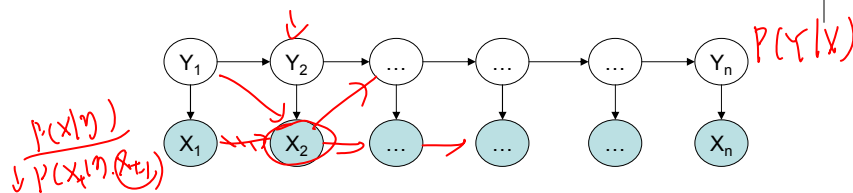
- The M step ("symbolically" identical to MLE)

$$\pi_i^{ML} = \frac{\sum_n \gamma_{n,1}^i}{N} \quad a_{ij}^{ML} = \frac{\sum_n \sum_{t=2}^T \xi_{n,t}^{i,j}}{\sum_n \sum_{t=1}^{T-1} \gamma_{n,t}^i} \quad b_{ik}^{ML} = \frac{\sum_n \sum_{t=1}^T \gamma_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^T \gamma_{n,t}^i}$$

Eric Xing

17

# Shortcomings of Hidden Markov Model



- HMM models capture dependencies between each state and **only** its corresponding observation

- NLP example: In a sentence segmentation task, each segmental state may depend not just on a single word (and the adjacent segmental stages), but also on the (non-local) features of the whole line such as line length, indentation, amount of white space, etc.

- Mismatch between learning objective function and prediction objective function

- HMM learns a joint distribution of states and observations  $P(\mathbf{Y}, \mathbf{X})$ , but in a prediction task, we need the conditional probability  $P(\mathbf{Y}|\mathbf{X})$

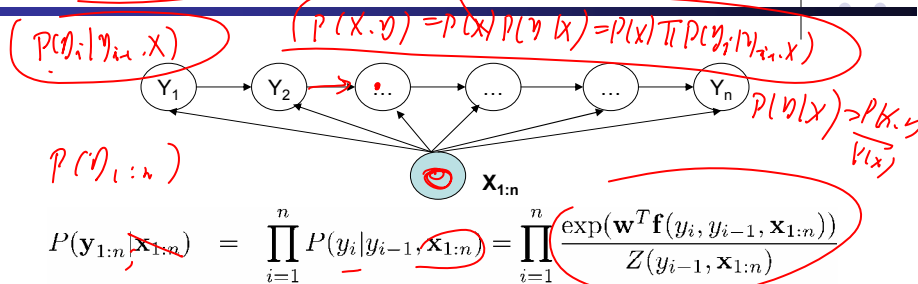
$$P(\mathbf{x}, \mathbf{y}) \rightarrow P(\mathbf{x}) \rightarrow P(\mathbf{y}|\mathbf{x}) = \frac{P(\mathbf{x}, \mathbf{y})}{P(\mathbf{x})}$$

Eric Xing

18

## Solution:

### Maximum Entropy Markov Model (MEMM)

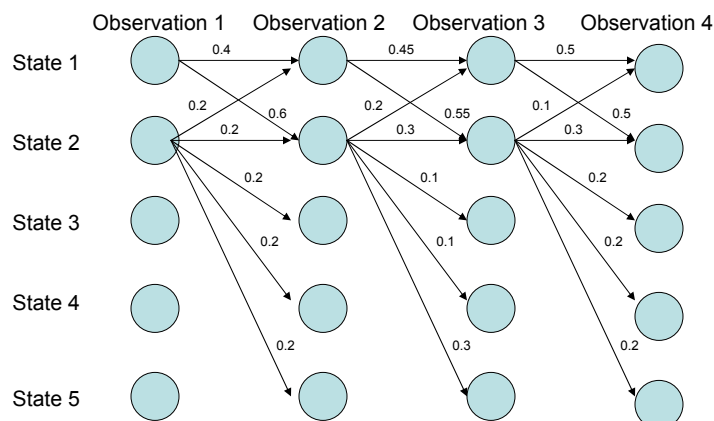


- Models dependence between each state and the full observation sequence explicitly
  - More expressive than HMMs
- Discriminative model
  - Completely ignores modeling  $P(\mathbf{X})$ : saves modeling effort
  - Learning objective function consistent with predictive function:  $P(\mathbf{Y}|\mathbf{X})$

Eric Xing

19

## MEMM: the Label bias problem



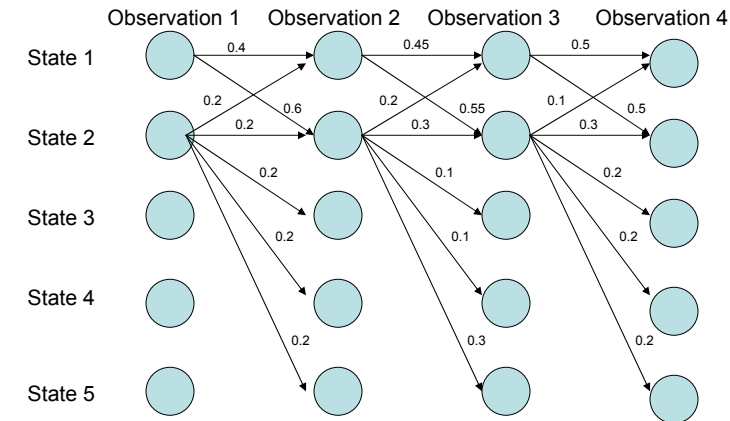
What the local transition probabilities say:

- State 1 almost always prefers to go to state 2
- State 2 almost always prefer to stay in state 2

Eric Xing

20

## MEMM: the Label bias problem



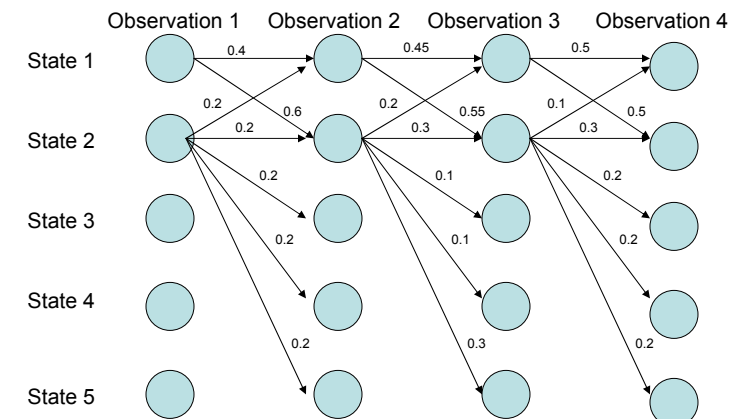
Probability of path 1-> 1-> 1-> 1:

- $0.4 \times 0.45 \times 0.5 = 0.09$

Eric Xing

21

## MEMM: the Label bias problem



Probability of path 2->2->2->2 :

- $0.2 \times 0.3 \times 0.3 = 0.018$

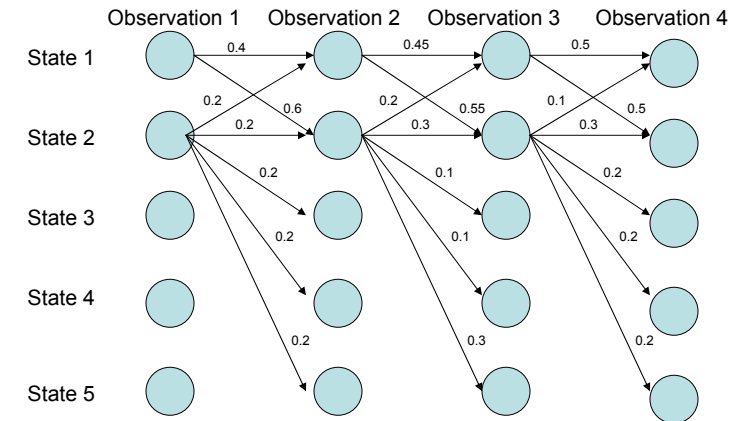
Other paths:

1-> 1-> 1-> 1: 0.09

Eric Xing

22

## MEMM: the Label bias problem



Probability of path 1->2->1->2:

- $0.6 \times 0.2 \times 0.5 = 0.06$

Other paths:

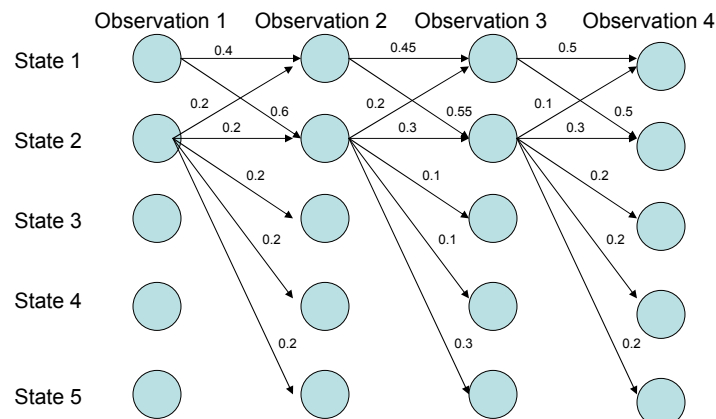
1->1->1->1: 0.09

2->2->2->2: 0.018

Eric Xing

23

## MEMM: the Label bias problem



Probability of path 1->1->2->2:

- $0.4 \times 0.55 \times 0.3 = 0.066$

Other paths:

1->1->1->1: 0.09

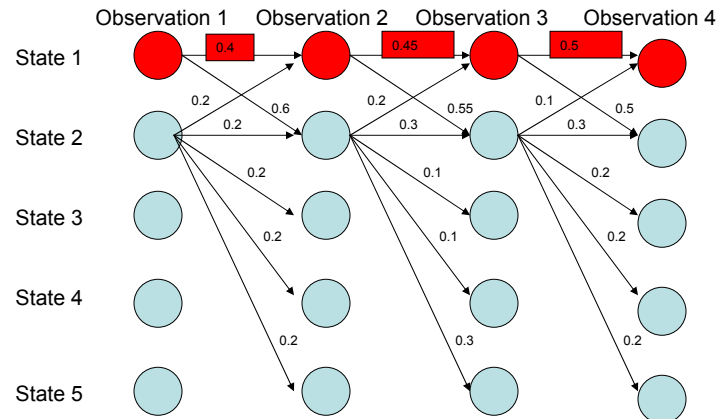
2->2->2->2: 0.018

1->2->1->2: 0.06

Eric Xing

24

## MEMM: the Label bias problem



**Most Likely Path:** 1 → 1 → 1 → 1

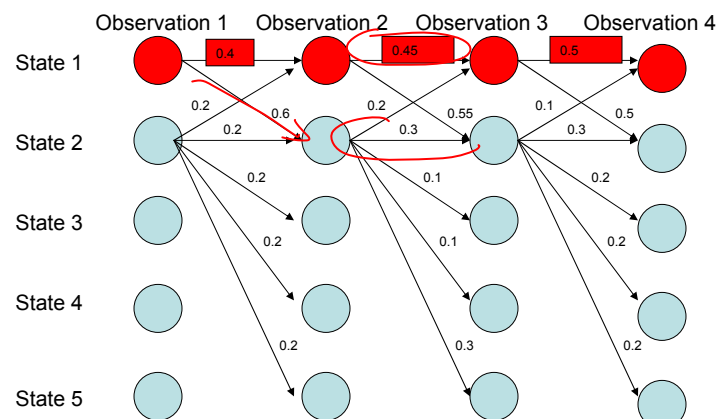
- Although **locally** it seems state 1 wants to go to state 2 and state 2 wants to remain in state 2.

• **why?**

Eric Xing

25

## MEMM: the Label bias problem



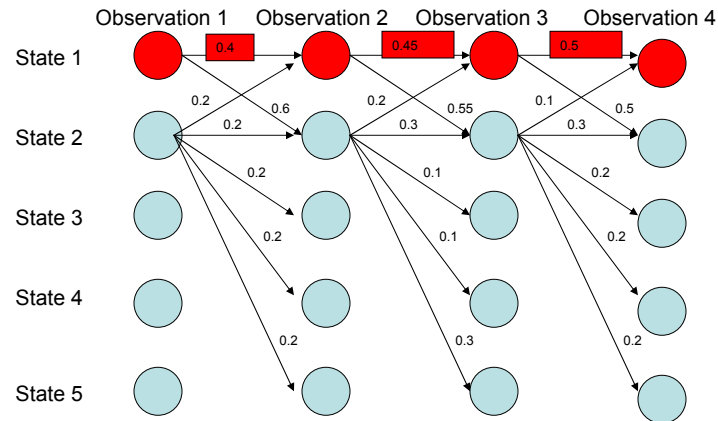
**Most Likely Path:** 1 → 1 → 1 → 1

- State 1 has only two transitions but state 2 has 5:
- Average transition probability from state 2 is lower

Eric Xing

26

## MEMM: the Label bias problem



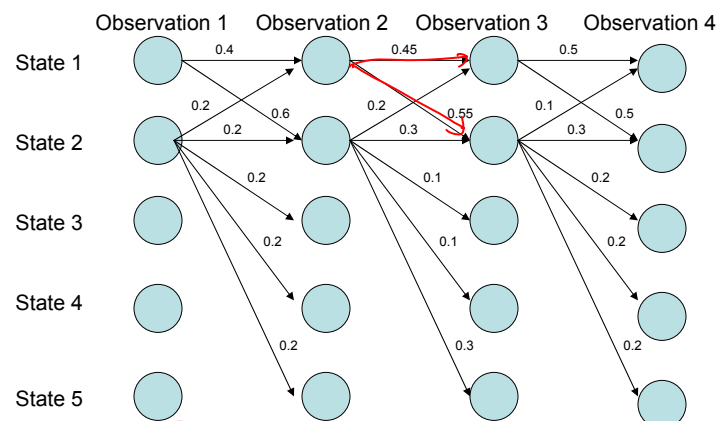
### Label bias problem in MEMM:

- Preference of states with lower number of transitions over others

Eric Xing

27

## Solution: Do not normalize probabilities locally

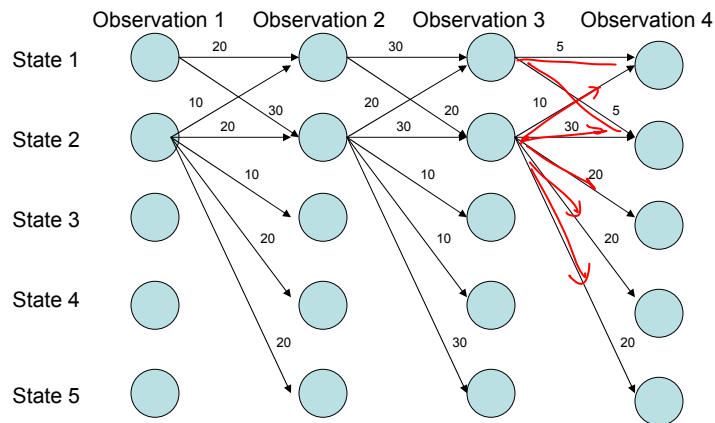


From local probabilities ...

Eric Xing

28

## Solution: Do not normalize probabilities locally



From local probabilities to local potentials

- States with lower transitions do not have an unfair advantage!

Eric Xing

29

- Office hour

1 - 2 pm Mon

- Mid-term project milestone due today

- Next week

4:30 - 5:30 this Friday

Eric Xing

30