

Machine Learning

10-701/15-781, Fall 2006

Hidden Markov Model

Eric Xing

Lecture 15, November 2, 2006



Reading: Chap. 13, C.B book



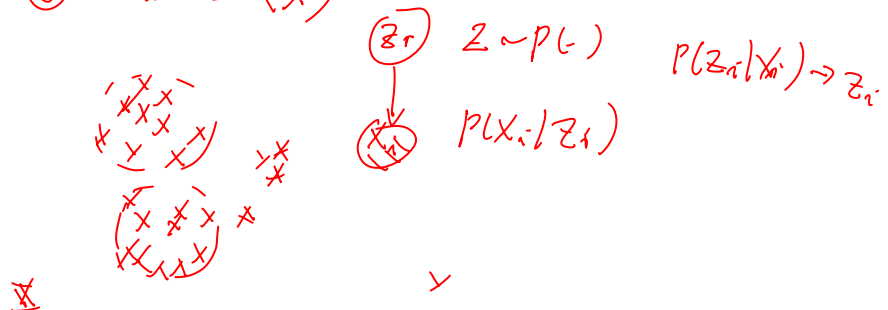
Eric Xing

1

① EM: MIX

② HMM

③ $x := \langle x \rangle$



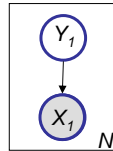
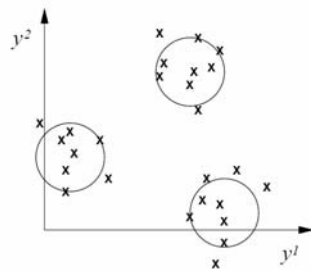
Eric Xing

2

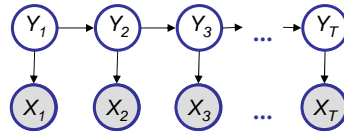
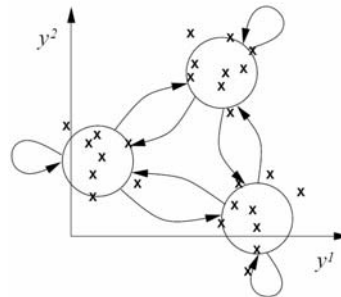
Hidden Markov Model: from static to dynamic mixture models



Static mixture



Dynamic mixture



Eric Xing

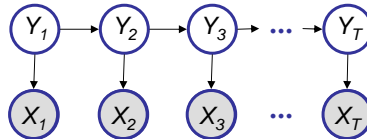
3

Hidden Markov Models



The underlying source:
genomic entities,
dice,

The sequence:
Play NT,
sequence of rolls,



Eric Xing

4

Example: The Dishonest Casino



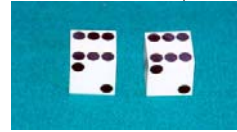
A casino has two dice:

- Fair die
 $P(1) = P(2) = P(3) = P(5) = P(6) = 1/6$
- Loaded die
 $P(1) = P(2) = P(3) = P(5) = 1/10$
 $P(6) = 1/2$

Casino player switches back-&-forth
between fair and loaded die once every
20 turns

Game:

1. You bet \$1
2. You roll (always with a fair die)
3. Casino player rolls (maybe with fair die,
maybe with loaded die)
4. Highest number wins \$2



Eric Xing

5

Puzzles Regarding the Dishonest Casino



GIVEN: A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

QUESTION

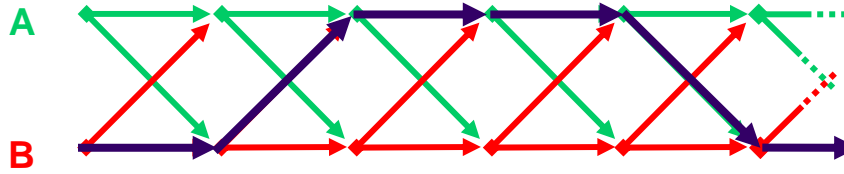
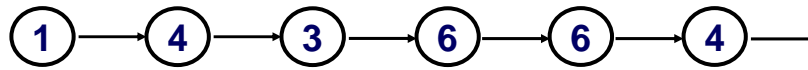
- How likely is this sequence, given our model of how the casino works?
 - This is the **EVALUATION** problem in HMMs
- What portion of the sequence was generated with the fair die, and what portion with the loaded die?
 - This is the **DECODING** question in HMMs
- How “loaded” is the loaded die? How “fair” is the fair die? How often does the casino player change from fair to loaded, and back?
 - This is the **LEARNING** question in HMMs

Eric Xing

6

A Stochastic Generative Model

- Observed sequence:



- Hidden sequence (a parse or segmentation):



Eric Xing

7

Definition (of HMM)

- Observation space

Alphabetic set: $\mathcal{C} = \{c_1, c_2, \dots, c_K\}$
Euclidean space: \mathbb{R}^d

- Index set of hidden states

$\mathcal{I} = \{1, 2, \dots, M\}$

- Transition probabilities between any two states

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

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

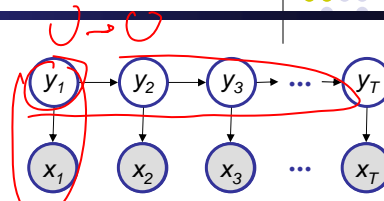
- Start probabilities

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

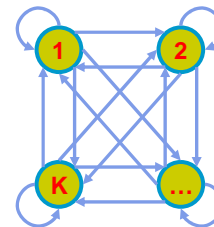
- Emission probabilities associated with each state

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

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



Graphical model



State automata

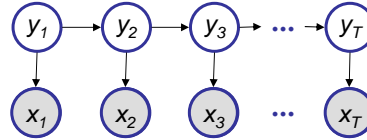
Eric Xing

8

Probability of a Parse

$$p(y_i) \quad p(y_i | y_{i-1})$$

- Given a sequence $x = x_1, \dots, x_T$ and a parse $y = y_1, \dots, y_T$,
- To find how likely is the parse:
(given our HMM and the sequence)



$$\begin{aligned} p(x, y) &= p(x_1, \dots, x_T, y_1, \dots, y_T) \quad (\text{Joint probability}) \\ &= p(y_1) p(x_1 | y_1) p(y_2 | y_1) p(x_2 | y_2) \dots p(y_T | y_{T-1}) p(x_T | y_T) \\ &= p(y_1) p(y_2 | y_1) \dots p(y_T | y_{T-1}) \times p(x_1 | y_1) p(x_2 | y_2) \dots p(x_T | y_T) \\ &= p(y_1, \dots, y_T) p(x_1, \dots, x_T | y_1, \dots, y_T) \end{aligned}$$

$$\text{Let } \pi_{y_1} \stackrel{\text{def}}{=} \prod_{i=1}^M [\pi_i]^{y_1^i}, \quad a_{y_i, y_{i+1}} \stackrel{\text{def}}{=} \prod_{j=1}^M [a_{ij}]^{y_i^j y_{i+1}^j}, \quad \text{and } b_{y_i, x_i} \stackrel{\text{def}}{=} \prod_{k=1}^K [b_k]^{y_i^j x_i^k},$$

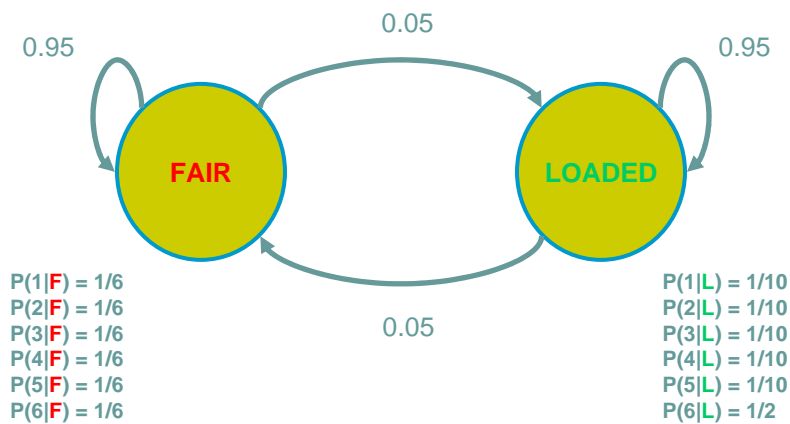
$$= \pi_{y_1} a_{y_1, y_2} \dots a_{y_{T-1}, y_T} b_{y_1, x_1} \dots b_{y_T, x_T}$$

- Marginal probability: $p(x) = \sum_y p(x, y) = \sum_{y_1} \sum_{y_2} \dots \sum_{y_N} \pi_{y_1} \prod_{t=2}^T a_{y_{t-1}, y_t} \prod_{t=1}^T p(x_t | y_t)$
- Posterior probability: $p(y | x) = p(x, y) / p(x)$

Eric Xing

9

The Dishonest Casino Model



Eric Xing

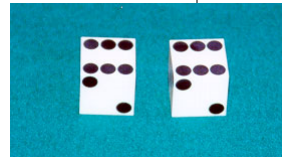
10

Example: the Dishonest Casino



- Let the sequence of rolls be:

- $x = 1, 2, 1, 5, 6, 2, 1, 6, 2, 4$



- Then, what is the likelihood of

- $y = \text{Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair?}$

(say initial probs $a_{0\text{Fair}} = 1/2$, $a_{0\text{Loaded}} = 1/2$)

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

$$\frac{1}{2} \times P(1 | \text{Fair}) P(\text{Fair} | \text{Fair}) P(2 | \text{Fair}) P(\text{Fair} | \text{Fair}) \dots P(4 | \text{Fair}) =$$

$$\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = .00000000521158647211 = 5.21 \times 10^{-9}$$

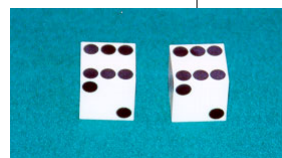
Eric Xing

11

Example: the Dishonest Casino



- So, the likelihood the die is fair in all this run is just 5.21×10^{-9}



- OK, but what is the likelihood of

- $\pi = \text{Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded?}$

$$\frac{1}{2} \times P(1 | \text{Loaded}) P(\text{Loaded} | \text{Loaded}) \dots P(4 | \text{Loaded}) =$$

$$\frac{1}{2} \times (1/10)^8 \times (1/2)^2 (0.95)^9 = .00000000078781176215 = 0.79 \times 10^{-9}$$

- Therefore, it is after all 6.59 times more likely that the die is fair all the way, than that it is loaded all the way

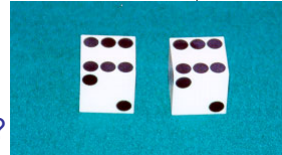
Eric Xing

12

Example: the Dishonest Casino



- Let the sequence of rolls be:
 - $x = 1, 6, 6, 5, 6, 2, 6, 6, 3, 6$
- Now, what is the likelihood $\pi = F, F, \dots, F$?
 - $\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = 0.5 \times 10^{-9}$, same as before
- What is the likelihood $y = L, L, \dots, L$?



$$\frac{1}{2} \times (1/10)^4 \times (1/2)^6 (0.95)^9 = .00000049238235134735 = 5 \times 10^{-7}$$

- So, it is 100 times more likely the die is loaded

Eric Xing

13

Three Main Questions on HMMs



1. Evaluation

GIVEN an HMM M and a sequence x ,
 FIND Prob ($x | M$)
 ALGO. Forward

2. Decoding

GIVEN an HMM M and a sequence x ,
 FIND the sequence y of states that maximizes, e.g., $P(y | x, M)$,
 or the most probable subsequence of states
 ALGO. Viterbi, Forward-backward

3. Learning

GIVEN an HMM M , with unspecified transition/emission probs.,
 and a sequence x ,
 FIND parameters $\theta = (\pi_i, a_{ij}, \eta_{ik})$ that maximize $P(x | \theta)$
 ALGO. Baum-Welch (EM)

Eric Xing

14

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.**

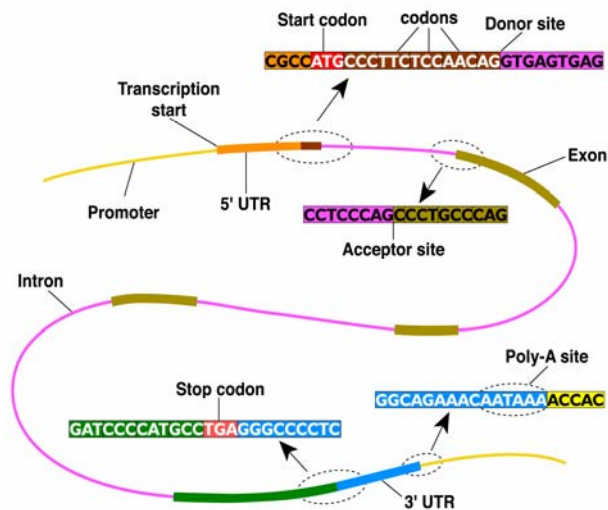
- **Some current applications of HMMs to biology**

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

Eric Xing

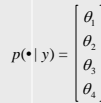
15

Typical structure of a gene



Eric Xing

16

[illegible]

- **Evaluation:** What is the probability of the observed sequence? **Forward**
- **Decoding:** What is the probability that the state of the 3rd roll is loaded, given the observed sequence? **Forward-Backward**
- **Decoding:** What is the most likely die sequence? **Viterbi**
- **Learning:** Under what parameterization are the observed sequences most probable? **Baum-Welch (EM)**

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

19

The Forward Algorithm – derivation

- Compute the forward probability:

$$\begin{aligned} \alpha_t^k &= P(x_1, \dots, x_{t-1}, x_t, y_t^k = 1) \\ &= \sum_{y_{t-1}} P(x_1, \dots, x_{t-1}, x_t, y_{t-1}, y_t^k = 1) \\ &= \sum_{y_{t-1}} P(x_1, \dots, x_{t-1}, y_{t-1}) P(y_t^k = 1 | y_{t-1}, x_1, \dots, x_{t-1}) P(x_t | y_t^k = 1, x_1, \dots, x_{t-1}, y_{t-1}) \\ &= \sum_{y_{t-1}} P(x_1, \dots, x_{t-1}, y_{t-1}) P(y_t^k = 1 | y_{t-1}) P(x_t | y_t^k = 1) \\ &= P(x_t | y_t^k = 1) \sum_i P(x_1, \dots, x_{t-1}, y_{t-1}^i = 1) P(y_t^k = 1 | y_{t-1}^i = 1) \\ &= P(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k} \end{aligned}$$

$$\text{Chain rule: } P(A, B, C) = P(A)P(B|C)P(C|A, B)$$

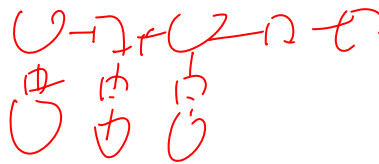
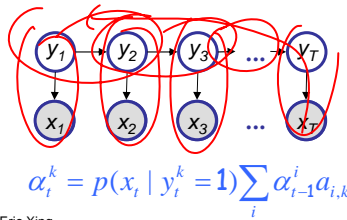
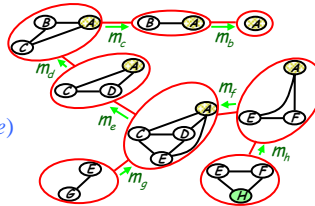
Eric Xing

20

Recall the Elimination and Message Passing Algorithm

- Elimination \equiv message passing on a **clique tree**

$$m_e(a, c, d) = \sum_e p(e | c, d) m_g(e) m_f(a, e)$$



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

Eric Xing

21

The Forward Algorithm

- We can compute α_t^k for all k, t , using dynamic programming!

Initialization:

$$\alpha_1^k = P(x_1 | y_1^k = 1) \pi_k$$

Iteration:

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

Termination:

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

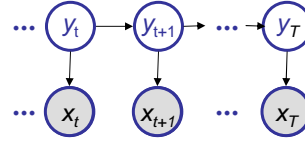
$$\begin{aligned} \alpha_1^k &= P(x_1, y_1^k = 1) \\ &= P(x_1 | y_1^k = 1) P(y_1^k = 1) \\ &= P(x_1 | y_1^k = 1) \pi_k \end{aligned}$$

Eric Xing

22

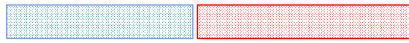
The Backward Algorithm

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



- We start by computing

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



Forward, α_t^k

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

- The recursion:

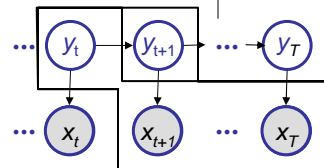
$$\beta_t^k = \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i$$

Eric Xing

23

The Backward Algorithm – derivation

- Define the backward probability:



$$\begin{aligned} \beta_t^k &= P(x_{t+1}, \dots, x_T | y_t^k = 1) \\ &= \sum_{y_{t+1}^i} P(x_{t+1}, \dots, x_T, y_{t+1}^i | y_t^k = 1) \\ &= \sum_i P(y_{t+1}^i = 1 | y_t^k = 1) p(x_{t+1} | y_{t+1}^i = 1, y_t^k = 1) P(x_{t+2}, \dots, x_T | x_{t+1}, y_{t+1}^i = 1, y_t^k = 1) \\ &= \sum_i P(y_{t+1}^i = 1 | y_t^k = 1) p(x_{t+1} | y_{t+1}^i = 1) P(x_{t+2}, \dots, x_T | y_{t+1}^i = 1) \\ &= \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i \end{aligned}$$

Chain rule: $P(A, B, C | \alpha) = P(A, \alpha) P(B | C, \alpha) P(C | A, B, \alpha)$

Eric Xing

24

The Backward Algorithm



- We can compute β_t^k for all k, t , using dynamic programming!

Initialization:

$$\beta_T^k = 1, \forall k$$

Iteration:

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

Termination:

$$P(\mathbf{x}) = \sum_k \alpha_1^k \beta_1^k$$

Eric Xing

25

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 to 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 states

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

Example:
MPA of X ?
MPA of (X, Y) ?

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

Eric Xing

26

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:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) = \operatorname{argmax}_{\pi} P(\mathbf{y}, \mathbf{x})$$

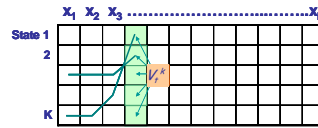
- Let

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

= 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

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

- 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

27

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

28

The Viterbi Algorithm



- Input: $\mathbf{x} = x_1, \dots, x_T$

Initialization:

$$V_1^k = P(x_1 | y_1^k = 1) \pi_k$$

Iteration:

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

$$\text{Ptr}(k, t) = \arg \max_i a_{i,k} V_{t-1}^i$$

Termination:

$$P(\mathbf{x}, \mathbf{y}^*) = \max_k V_T^k$$

TraceBack:

$$y_T^* = \arg \max_k V_T^k$$

$$y_{t-1}^* = \text{Ptr}(y_t^*, t)$$

Eric Xing

29

Computational Complexity and implementation details



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

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

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

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

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

30

Learning HMM: two scenarios

- **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 θ of the model to maximize $P(x|\theta)$ --- Maximal likelihood (ML) estimation

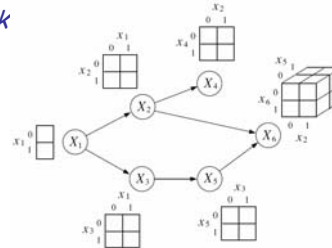
Eric Xing

31

Recall MLE for observed BN

- Assume each CPD is represented as a table (multinomial) where

$$\theta_{ijk} \stackrel{\text{def}}{=} p(X_i = j \mid X_{\pi_i} = k)$$



- The log-likelihood is

$$\ell(\theta; \mathcal{D}) = \log \prod_{i,j,k} \theta_{ijk}^{n_{ijk}} = \sum_{i,j,k} n_{ijk} \log \theta_{ijk}$$

- Using a Lagrange multiplier to enforce so $\sum_j \theta_{ijk} = 1$ we get

$$\theta_{ijk}^{ML} = \frac{n_{ijk}}{\sum_{j'} n_{ij'k}}$$

Eric Xing

32

Supervised ML estimation

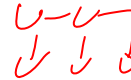


- Given $x = x_1 \dots x_N$ for which the true state path $y = y_1 \dots y_N$ is known,

- Define:

$$A_{ij} = \# \text{ times state transition } i \rightarrow j \text{ occurs in } y$$

$$B_{ik} = \# \text{ times state } i \text{ in } y \text{ emits } k \text{ in } x$$



- We can show that the maximum likelihood parameters θ are:

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

$$b_{ik}^{ML} = \frac{\#(i \rightarrow k)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=1}^T y_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^T y_{n,t}^i} = \frac{B_{ik}}{\sum_k B_{ik}}$$

Handwritten notes: $y \rightarrow ()$, $y, y \rightarrow x$, $(y_t) (x_t)$, $P(\eta_{t+1}, \eta_t | x)$, $(P(\eta | x))$

- 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 Gaussian, and apply learning rules for Gaussian ...

(Homework!)

Eric Xing

33

Supervised ML estimation, ctd.



- Intuition:

- When we know the underlying states, the best estimate of θ is the average frequency of transitions & emissions that occur in the training data

- Drawback:

- Given little data, there may be **overfitting**:
 - $P(x|\theta)$ is maximized, but θ is unreasonable
 - 0 probabilities – VERY BAD**

- Example:

- Given 10 casino rolls, we observe

$x = 2, 1, 5, 6, 1, 2, 3, 6, 2, 3$
 $y = F, F, F, F, F, F, F, F, F, F$

- Then:

$a_{FF} = 1; \quad a_{FL} = 0$
 $b_{F1} = b_{F3} = .2;$
 $b_{F2} = .3; b_{F4} = 0; b_{F5} = b_{F6} = .1$

Eric Xing

34

Pseudocounts



- Solution for small training sets:
 - Add pseudocounts
 - A_{ij} = # times state transition $i \rightarrow j$ occurs in y + R_{ij}
 - B_{ik} = # times state i in y emits k in x + S_{ik}
 - R_{ij}, S_{ij} are pseudocounts representing our prior belief
 - Total pseudocounts: $R_i = \sum_j R_{ij}$, $S_i = \sum_k S_{ik}$,
 - --- "strength" of prior belief,
 - --- total number of imaginary instances in the prior
- Larger total pseudocounts \Rightarrow strong prior belief
- Small total pseudocounts: just to avoid 0 probabilities --- smoothing

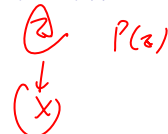
Eric Xing

35

Unsupervised ML estimation



- Given $x = x_1 \dots x_N$ for which the true state path $y = y_1 \dots y_N$ is unknown,



EXPECTATION MAXIMIZATION

0. Starting with our best guess of a model M , parameters θ .
1. Estimate A_{ij}, B_{ik} in the training data
 - How? $A_{ij} = \sum_{n,t} \langle y_{n,t-1}^i y_{n,t}^j \rangle$ $B_{ik} = \sum_{n,t} \langle y_{n,t}^i x_{n,t}^k \rangle$, How? (homework)
2. Update θ according to A_{ij}, B_{ik}
 - Now a "supervised learning" problem
3. Repeat 1 & 2, until convergence

This is called the Baum-Welch Algorithm

We can get to a provably more (or equally) likely parameter set θ each iteration

Eric Xing

36

The Baum Welch algorithm



- The complete log ~~likelihood~~

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

- The expected complete log likelihood

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

- 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)

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

Eric Xing

37

The Baum-Welch algorithm -- comments



Time Complexity:

$$\# \text{ iterations} \times O(K^2N)$$

- Guaranteed to increase the log likelihood of the model
- Not guaranteed to find globally best parameters
- Converges to local optimum, depending on initial conditions
- Too many parameters / too large model: Overt-fitting

Eric Xing

38



GLIM

$$p(\quad) = \exp \left[\sum w_i x_i \right]$$

$$g^y = \exp \ln b \cdot y$$

$$\ln = \ln b + \ln y$$

$$\ln \exp - \frac{(x-m)^2}{\sigma^2}$$

$$\left(\frac{(x-m)^2}{\sigma^2} \right) \rightarrow (x) \alpha^2$$

Eric Xing

39



Eric Xing

40