

Recitation: HMM, GM and Learning Theory

Zeyu Jin

Outline

- Learning Theory
 - Uniform bound
 - $|H|$ and $VC(H)$
 - Insights
- GM
 - Factorized probability
 - D-separation
 - Inference
- HMM (recap)
 - Basic questions
 - Algorithms
 - Insight

Learning Theory

1. The question

- Want to know how good our classifier is

$$error_{true}(H) = ?$$

- However, H is trained on some data; **the randomness of data** makes this “?” a distribution. Let’s try

$$P(error_{true}(H) = p) = ?, \quad p \in [0,1]$$

- It is non-trivial

Learning Theory

1. The question

- With a family of models **H of certain complexity**, how many training **samples R** is needed in order to learn a model h with **reasonable training time** and **sufficient accuracy** on future data?
- We want answer

$$error_{true}(H(X^m)) = ?$$

Computationally efficient in polynomial time

Learning Theory

1. The question

- Distribution of error rate

$$\begin{aligned} &P(\text{error}_{\text{true}}(H) = p) \\ &= E[P(\text{error}_{\text{true}}(H(X)) = p \mid X = x)] \\ &= \iiint_X P(\text{error}_{\text{true}}(H(X)) = p \mid X = x) P_{\text{true}}(X = x) dx \end{aligned}$$

- Maybe we can try to get a uniform bound for this question

$$P(|\text{error}_{\text{true}}(H) - E_X[\text{error}_{\text{true}}(H)]| < \varepsilon) = ?$$

Learning Theory

1. The question

- Distribution of error rate

$$P(\text{error}_{\text{true}}(H) = p)$$

$$= E[P(\text{error}_{\text{true}}(H(X)) = p \mid X = x)]$$

$$= \iiint_X P(\text{error}_{\text{true}}(H(X)) = p \mid X = x) P_{\text{true}}(X = x) dx$$

- Maybe we can try to get a uniform bound for this question

$$P(|\text{error}_{\text{true}}(H) - E_X[\text{error}_{\text{true}}(H)]| < \varepsilon) = ?$$

- Still extremely hard. Maybe bound this probability

Learning Theory

1. Uniform Bound

- Bound the probability of the bounded error

$$P(|error_{true}(H) - E_X[error_{true}(H)]| < \varepsilon) > 1 - \delta$$

- Statisticians do have solution for this form!
- Three basic questions

- H is finite, $E_X[error_{true}(H)] = error_{train}(H)$ is 0 **PAC**
- H is finite, $E_X[error_{true}(H)] = error_{train}(H)$ is non-zero
- H is infinite

Learning Theory

2. Solutions

1) H is finite, $E_X[\text{error}_{\text{true}}(H)] = \text{error}_{\text{train}}(H)$ is 0

$$P(|\text{error}_{\text{true}}(H) - 0| < \varepsilon) \geq 1 - \delta \quad \varepsilon = \frac{\ln |H| + \ln(1/\delta)}{|X|}$$

2) H is finite, $E_X[\text{error}_{\text{true}}(H)] = \text{error}_{\text{train}}(H)$ is non-zero

$$P(|\text{error}_{\text{true}}(H) - E_X[\text{error}_{\text{true}}(H)]| < \varepsilon) > 1 - \delta$$

$$\varepsilon = \sqrt{\frac{|H| + \ln(1/\delta)}{2|X|}}$$

Learning Theory

2. Solutions

3) H is infinite

$$P(|error_{true}(H) - E_X[error_{true}(H)]| < \varepsilon) > 1 - \delta$$

$$\varepsilon = 8 \sqrt{\frac{VC(H) \left(\ln \frac{m}{VC(H)} + 1 \right) + \ln \frac{8}{\delta}}{2m}}$$

Learning Theory

3. Terms in this solutions

- For solution 1 and 2: $|H| = ?$
- For solution 3: $VC(H) = ?$

Learning Theory

3. Terms in this solutions

- **For solution 1 and 2: $|H| = ?$**
- For solution 3: $VC(H) = ?$

Instead of limiting the maximal depth of a decision tree, let's assume n binary attributes and binary class

What is $|H|$?

Learning Theory

3. Terms in this solutions

- For solution 1 and 2: $|H| = ?$
- **For solution 3: $VC(H) = ?$**
 - Find the **maximal N**
 - Where there **EXIST** N points in the problem's space
 - s.t. **ALL** element of the superset of these points (2^N)
 - can be picked out by S

Learning Theory

3. Terms in this solutions

- For solution 1 and 2: $|H| = ?$
- **For solution 3: $VC(H) = ?$**

What is the VC dimension of a 2D-circle?

$$\{(x, y) \mid x^2 + y^2 \leq R^2\}$$

Learning Theory

3. Terms in this solutions

- For solution 1 and 2: $|H| = ?$
- **For solution 3: $VC(H) = ?$**

What if the a circle plus a point?

$$\{(x, y) \mid x^2 + y^2 \leq R^2 \setminus (0, 0)\}$$

Learning Theory

4. Insights

– VC Dimension

- Find the **maximal N**
- Where there **EXIST** N points in the problem's space
- s.t. **ALL** element of the superset of these points (2^N)
- can be picked out by S

$S_N(H)$ = The number of elements of the superset of these N points can be picked out by H

$$\mathbb{P} \left(\sup_{A \in \mathcal{A}} |P_n(A) - P(A)| > \epsilon \right) \leq 8 (n+1)^d e^{-n\epsilon^2/32}.$$

$S_N(H)$ is not easy to obtain, but it can be shown that $s(\mathcal{A}, n) \leq (n+1)^d$.

Where d is VC dimension

Learning Theory

4. Insights

- Obtain error bound by simulation
 - Known: marginal distribution of data D , true model $p(Y|X)$

Repeat the following N times

- Draw m data from D for training; draw $k \gg m$ from D for test
- Draw y_i for each x_i from $p(Y|X)$
- Learn h based on your hypothesis space H
- Evaluate $\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h)$ on test data (or do it mathematically)

Then you will get a histogram of error which approximates $P(\text{error}_{\text{true}}(h))$. Solve

$$P(\text{error}_{\text{true}}(H) < \varepsilon) \geq 1 - \delta$$

Learning Theory

4. Insights

- Connection

$$P(\text{error}_{\text{true}}(H) < \varepsilon) \geq 1 - \delta$$

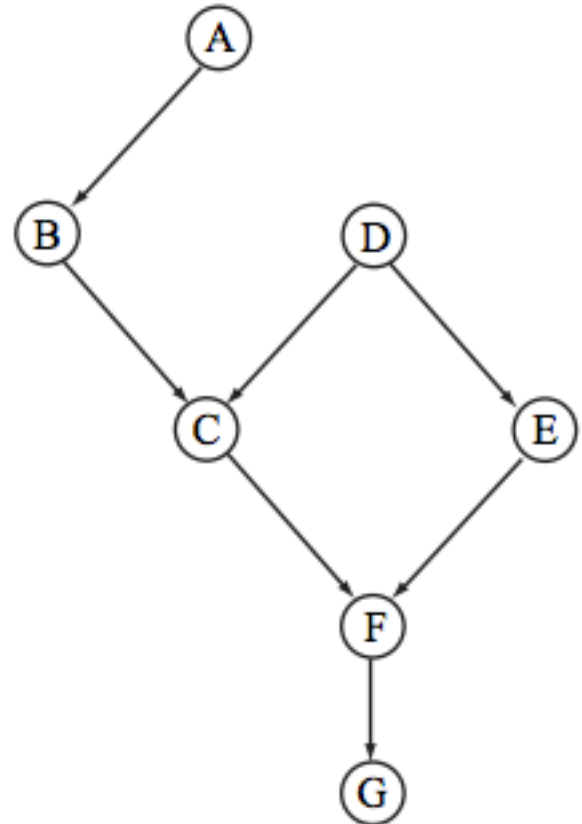
Confidence interval: with confidence

δ

we conclude that error of estimating $\text{error}_{\text{true}}$ is less than ε

Bayesian Networks

- Bayes Net \Leftrightarrow Factorized probability
 - Write Factorized probability

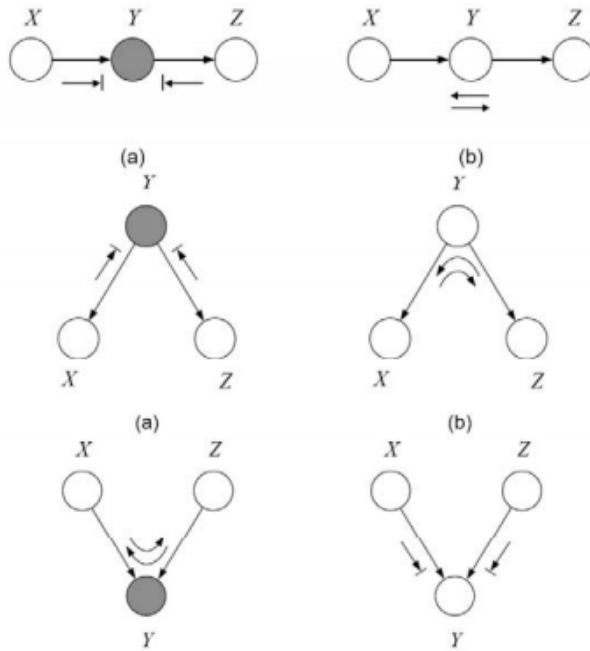


Bayesian Networks

- Bayes Net \Leftrightarrow Factorized probability
 - What's the Bayes net for
 - Naïve Bayes?
 - Full Bayes?
 - k-th order Markov Model
 - Hidden Markov model

Bayesian Networks

- Understand dependency in BN – D-separation



X and Y are D-separated by Z

If all the path from X to Y are blocked

Maybe we lost one case in class

Bayesian Networks

- Understand dependency in BN – D-separation

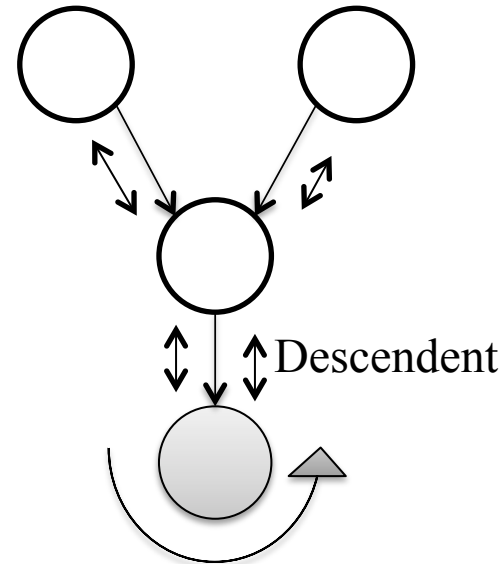
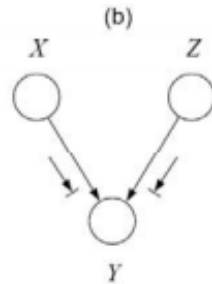
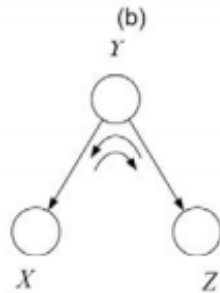
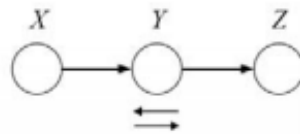
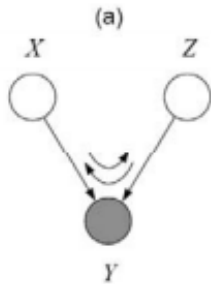
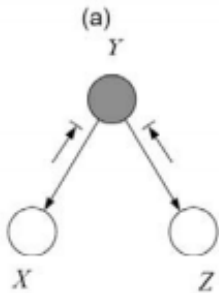
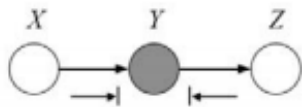
Original Definition

- (a) the arrows on the path meet either head-to-tail or tail-to-tail at the node, and the node is in the set C , or
- (b) the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in the set C .

Bishop 8.2.2

Bayesian Networks

- Understand dependency in BN – D-separation

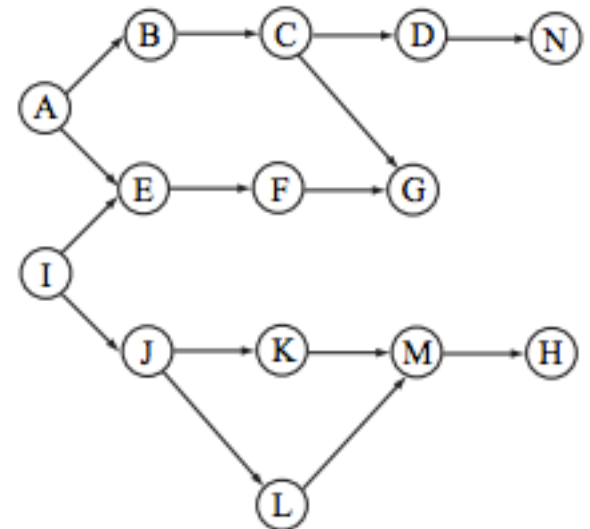


Bayesian Networks

- Understand dependency in BN – D-separation

Exam problem TRUE/FALSE

- (a) $P(D, H) = P(D)P(H)$
- (b) $P(A, I) = P(A)P(I)$
- (c) $P(A, I|G) = P(A|G)P(I|G)$
- (d) $P(J, G|F) = P(J|F)P(G|F)$
- (e) $P(J, M|K, L) = P(J|K, L)P(M|K, L)$
- (f) $P(E, C|A, G) = P(E|A, G)P(C|A, G)$
- (g) $P(E, C|A) = P(E|A)P(C|A)$



Bayesian Networks

- Understand dependency in BN – D-separation

Exam problem TRUE/FALSE

Key...

(a) $P(D, H) = P(D)P(H)$

(b) $P(A, I) = P(A)P(I)$

(c) $P(A, I|G) = P(A|G)P(I|G)$

(d) $P(J, G|F) = P(J|F)P(G|F)$

(e) $P(J, M|K, L) = P(J|K, L)P(M|K, L)$

(f) $P(E, C|A, G) = P(E|A, G)P(C|A, G)$

(g) $P(E, C|A) = P(E|A)P(C|A)$

- a) Yes, blocked by on E on one path and G on another path
- b) Yes, blocked by E
- c) No, G is a descendent of E
- d) No, the path JIEABCG is unblocked
- e) Yes, blocked on both paths
- f) No, path EFGC unblocked
- g) Yes, EABC blocked by A, and EFGC blocked by G

Bayesian Networks

- Inference
 - What is inference?
the process of computing answers to queries about the distribution P defined by given BN
 - Likelihood
 - Conditional probability (we will see one example after this slide)
 - Most probable assignment (most likely states sequence in HMM)
 - Methods?
 - Variable elimination, belief propagation (do exact calculation)
 - Gibbs sampling (simulation)

Bayesian Networks

- Inference 1: variable elimination

$$P(G = T|A = T) = ?$$

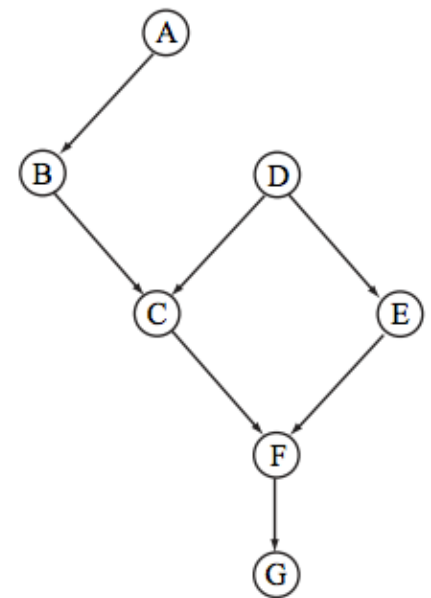
$$\begin{aligned} P(G = T|A = T) &= \frac{P(G = T, A = T)}{P(A = T)} = \sum_{BCDEF} P(A = T, B, C, D, E, F, G = T) \\ &= \sum_{BCDEF} P(B|A = T)P(D)P(C|B, D)P(E|C, D)P(F|C, E)P(G = T|F) \\ &= \sum_B P(B|A = T) \sum_D P(D) \sum_C P(C|B, D) \sum_E P(E|D) \sum_F P(F|C, E)P(G = T|F) \end{aligned}$$

$$f_{F,G}(c, e, G = T) = \sum_F P(F|C, E)P(G = T|F)$$

$$f_{C,D}(c, d, G = T) = \sum_E P(E|d)f_{F,G}(c, E, G = T)$$

$$f_{B,D}(b, d, G = T) = \sum_c P(f_c(c, B, D))f_{C,D}(c, d, G = T)$$

...

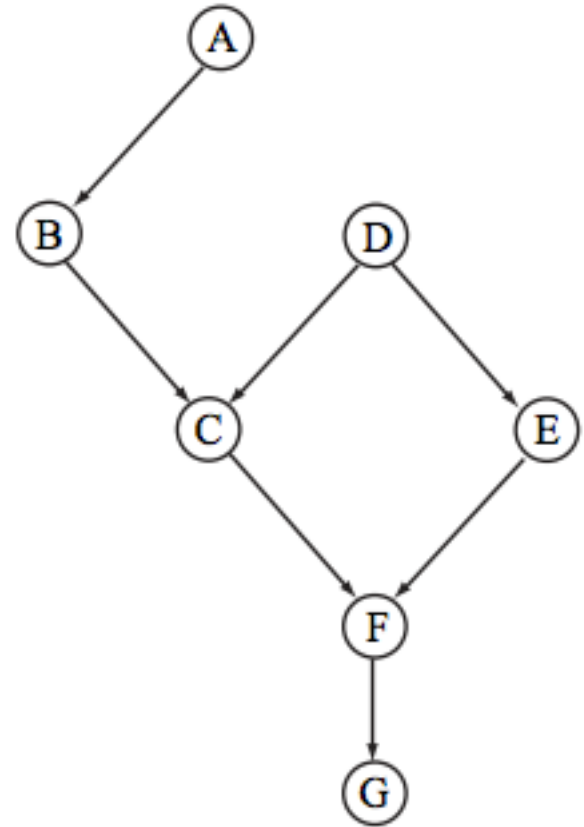


Binary for all RV

Bayesian Networks

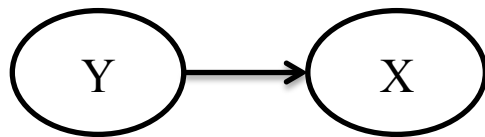
- Inference 2: sampling
 - Naïve sampling:
 - (A,B,C,D,E,F,G) each time from $p(A,B,C,D,E,F,G)$
 - Calculate $P(G|A=T)$ by counting
 - Have problem with rare event
 - Weighted sampling
 - If $P(A=T)$ is rare, just set $A=T$ and sample $(A=T,B,C,D,E,F,G)$
 - When calculating $P(G|A=T)$, the number of $(A=T,b,c,d,e,f,g)$ is weighted by $p(A=T)$

$$P(G = T | A = T) = ?$$



HMM (recap)

- Static vs. time series

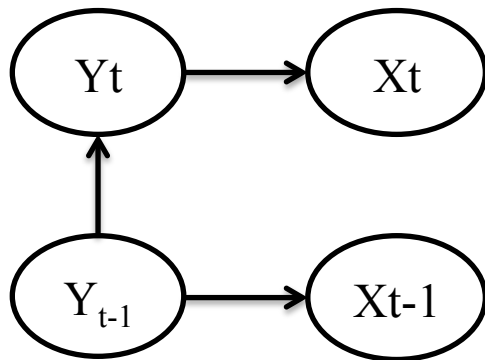


Discriminative

$$P(Y|X)$$

Generative

$$P(X, Y)$$



$P(Y_t, Y_{t-1} \dots | X_t X_{t-1} \dots)$ $P(X_t, Y_t, X_{t-1}, Y_{t-1} \dots)$
Conditional random field

Conditioning on no variable, all Xs and Ys are correlated

HMM (recap)

- Basic questions

1. Parameters

2. Factorization

3. Inference

4. Learning

- K: number of states

- M: number of observations

– Initial state: $P(y_1)$

– Transition: $P(y_t | y_{t-1})$

– Emission: $P(x_t | y_t)$

#par. shorthand

K-1 $\pi_i = P(y_1=i)$

$K*(K-1)$ $a_{ij} = P(y_{t+1}=j | y_t=i)$

$K*(M-1)$ $b_{ik} = P(x_t=k | y_t=i)$

HMM (recap)

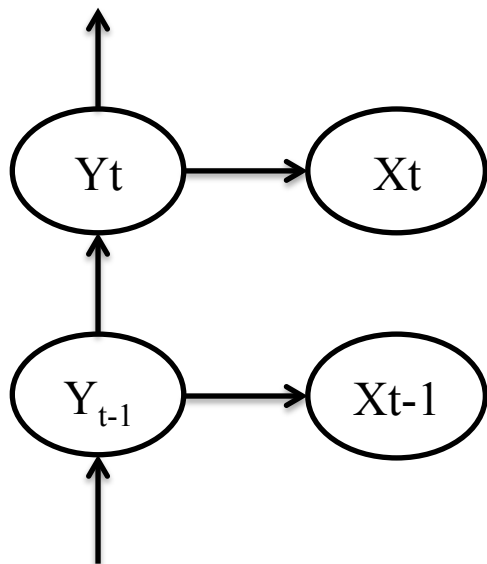
- Basic questions

1. Parameters
2. **Factorization**
3. Inference
4. Learning

HMM is Generative

Complete likelihood based on given parameters is

$$\begin{aligned} & P(x_1, \dots, x_T, y_1, \dots, y_T) \\ &= P(y_1) p(x_1 | y_1) p(y_2 | y_1) \dots P(y_T | y_{T-1}) P(x_T | y_T) \\ &= P(y_1) \prod_t P(y_t | y_{t-1}) P(x_t | y_t) \end{aligned}$$



HMM (recap)

- Basic questions

1. Parameters

$$P(y_t | X_{1:t}) = ?$$

2. Factorization

$$P(y_t | X_{1:T}) = ?$$

- 3. Inference**

$$\operatorname{argmax}_y P(y_{1:T} | X_{1:T}) = ?$$

4. Learning

Before that, we have the following tools

1. Forward probability

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

2. Backward probability

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

HMM (recap)

- Basic questions

1. Parameters

$$\underline{P(y_t | X_{1:t}) = ?}$$

2. Factorization

$$P(y_t | X_{1:T}) = ?$$

3. **Inference**

$$\operatorname{argmax}_y P(y_{1:T} | X_{1:T}) = ?$$

4. Learning

$$p(y_t = i | X_{1:t}) = \frac{p(x_1, \dots, x_t, y_t = i)}{p(x_1, \dots, x_t)} = \frac{\alpha_t^i}{\boxed{p(x_1, \dots, x_t)} = \sum_{i=1}^k \alpha_T^i}$$

HMM (recap)

- Basic questions

1. Parameters

$$P(y_t | X_{1:t}) = ?$$

2. Factorization

$$\underline{P(y_t | X_{1:T})} = ?$$

- 3. Inference**

$$\operatorname{argmax}_y P(y_{1:T} | X_{1:T}) = ?$$

4. Learning

$$\begin{aligned} p(y_t = i | x_1, \dots, x_T) &= \frac{p(y_t = i, x_1, \dots, x_T)}{p(x_1, \dots, x_T)} \\ &= \frac{p(y_t = i, x_1, \dots, x_t) p(x_{t+1}, \dots, x_T | y_t = i, x_1, \dots, x_t)}{p(x_1, \dots, x_T)} \\ &= \frac{\alpha_t^i \beta_t^i}{\boxed{p(x_1, \dots, x_T)}} = \sum_{i=1}^k \alpha_T^i \end{aligned}$$

?

HMM (recap)

- Basic questions

1. Parameters

$$P(y_t | X_{1:t}) = ?$$

2. Factorization

$$P(y_t | X_{1:T}) = ?$$

3. Inference

$$\underline{\operatorname{argmax}_y P(y_{1:T} | X_{1:T}) = ?}$$

4. Learning

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

Viterbi

HMM (recap)

- Basic questions

- Parameters
- Factorization
- Inference
- Learning**

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

E-step

$$\begin{aligned} \gamma_{n,t}^i &= \langle \gamma_{n,t}^i \rangle = p(\gamma_{n,t}^i = 1 | \mathbf{x}_n) \\ \xi_{n,t}^{i,j} &= \langle \gamma_{n,t-1}^i \gamma_{n,t}^j \rangle = p(\gamma_{n,t-1}^i = 1, \gamma_{n,t}^j = 1 | \mathbf{x}_n) \end{aligned}$$

M-step

$$\begin{aligned} \pi_i^{ML} &= \frac{\sum_n \gamma_{n,1}^i}{N} & a_{ij}^{ML} &= \frac{\sum_n \sum_{t=2}^T \xi_{n,t}^{i,j}}{\sum_n \sum_{t=1}^{T-1} \gamma_{n,t}^i} & b_{ik}^{ML} &= \frac{\sum_n \sum_{t=1}^T \gamma_{n,t}^i \mathbf{x}_{n,t}^k}{\sum_n \sum_{t=1}^{T-1} \gamma_{n,t}^i} \end{aligned}$$

HMM (recap)

- Computational complexity
 - Forward: K states, N time points $\Rightarrow O(K^2N)$
 - Backward: $O(K^2N)$
 - $P(y_t|X_{1:t})$: Forward + sum of forward = $O(K^2N)$
 - $P(y_t|X_{1:T})$: Forward + backward + sum of forward = $O(K^2N)$
 - Viterbi: forward + $O(K^2N)$ updates of V = $O(K^2N)$

HMM (recap)

- Other mutation of HMM
 - IO-HMM
 - Kalman Filter
 - MEMM
 - Spectral HMM (algorithm)