

Welcome

# COMP 136

## Day 19: Viterbi for HMMs

NOON

Agenda:

\* CP3 recap

\* Midterm recap

---- **midterm\_solutions.pdf on Piazza**

\* Q&A on Viterbi algorithm

logsumexp  
avoid nans

Upcoming Due Dates

- CP3 deadline extended to **Tue Apr 7 at NOON ET**
- Quiz3 released tomorrow Tue Apr 7, due **Thu Apr 9** 11:59pm ET
- Will be entirely online via gradescope
- Time limit: 20 minutes *~6*
- all multiple choice

HW4 released later today (due next week)

CP4 released in next day or so (will be pushed back a few days, stay tuned)

*next wed*

# CP3 recap

$$r_{nk} \triangleq \underbrace{P(z_n=k | x_n)} = \frac{P(x_n, z_n=k)}{P(x_n)}$$

- Autograd debugging tips? See @168
- Getting nans in your estep? See logsumexp tips in @183
- Getting nans in your entropy calculation? See @170

import scipy.special.logsumexp

$$r_{nk} \leftarrow \frac{\pi_k N(x_n | \mu_k, \sigma_k^2)}{\sum_l \pi_l N(x_n | \mu_l, \sigma_l^2)}$$

$$\log r_{nk} = a_{nk} - \log \sum_l e^{a_{nl}}$$

$$r_{nk} = e^{a_{nk} - \text{logsumexp}(a_n)}$$

$$\text{softmax}(a) = r_{nk} = 1 \quad \checkmark$$

$$a_{nk} \leftarrow \log \pi_k + \log N(x_n | \mu_k, \sigma_k^2)$$

$$r_{nk} = \frac{e^{a_{nk}}}{\sum_l e^{a_{nl}}}$$

$$r_{nk} = e^{[a_{nk} - \text{logsumexp}(a_n)]}$$

$\log P(x_n, z_n=k)$

$$\log \text{sumexp}(a_1 \dots a_k)$$

$$\log \left( \sum_k e^{a_k} \right)$$

$$k=3 \begin{cases} a_1 = -10000000 \\ a_2 = 1000 \\ a_3 = -1000 \end{cases}$$

$$\log \text{sumexp}(a_1 \dots a_k)$$

$$M = \max_k a_k$$

$$\sum_k e^{-1000}$$

$$e^M (e^{a_1-M} + e^{a_2-M} + \dots + e^{a_k-M})$$

$$\rightarrow \log \left( e^{-1000} + e^{-1000} + e^{-1000} \right)$$

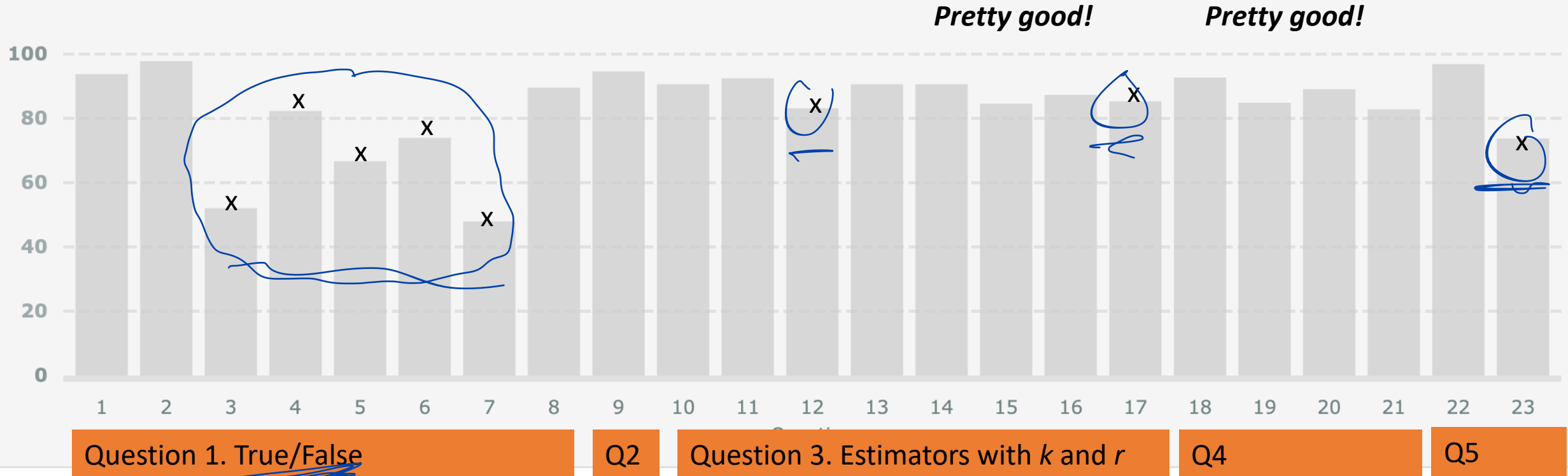
$$-1000 + \log(3)$$

$$\rightarrow \log \left( e^{-1000} (e^0 + e^0 + e^0) \right) = \log(e^{-1000}) + \log(e^0 + e^0 + e^0)$$

# Midterm

Percentiles:

- 10<sup>th</sup> : 0.77
- 50<sup>th</sup> : 0.87
- 90<sup>th</sup> : 0.93



midterm solutions.pdf

Piazza  
→ Resources

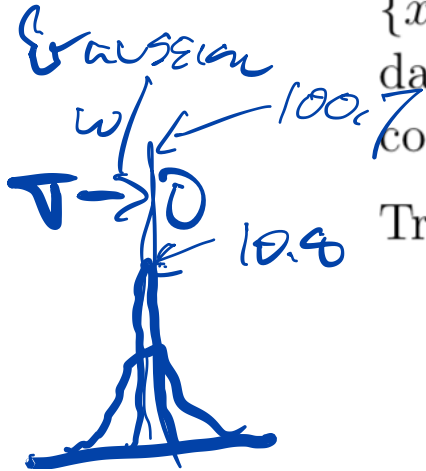


1c

$p(x)$  is a PDF  
for real  $x$

PDF  $\geq 0$   
Density

(c) You compute the log marginal likelihood  $\log p(\mathbf{x})$  of a dataset  $\mathbf{x} = \{x_n\}_{n=1}^N$ . Each observation is a real number:  $x_n \in \mathbb{R}$ . For a specific dataset, you compute  $\log p(\mathbf{x}) = 2.345$ . There must be a bug in your code; this value is impossible for a valid probabilistic model.



True  False

PDFs of real-valued variables can be any value  $> 0$ .  
 $e^{2.345}$  is a valid PDF value.

(d) Consider a random variable  $S$  with  $K$  possible values.  $\{2, 4, 8, 16, \dots, 2^K\}$ .

$$p(x) = e^{2.345} > 0$$

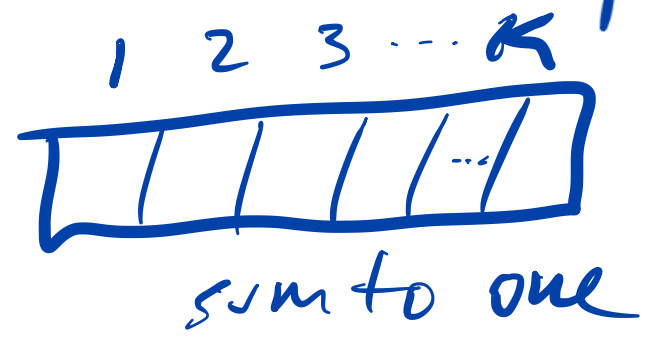
1d

$|\Omega| = K$   
 $\{a, b, c, d, \dots, K\}$

(d) Consider a random variable  $S$  with  $K$  possible values,  $\{2, 4, 8, 16, \dots, 2^K\}$ . We wish to define the marginal distribution  $p(S)$ . We can define this distribution in terms of a parameter vector that requires exactly  $K - 1$  distinct scalar values ( $K - 1$  degrees of freedom).

True  False

Yes, by the sum to one requirement, if we know  $K-1$  values we can determine the probab. of the  $K$ th outcome.



$K = 2$   
 $[0.6 \mid 0.4] = 1$

1e

- (e) When doing regression, if the inverse of  $\Phi^T \Phi$  does not exist, this means there is no possible weight vector  $w$  that maximizes the likelihood.

True

False

There are infinitely many  $w$  that maximize lik.

~~How many lines go through this point?~~

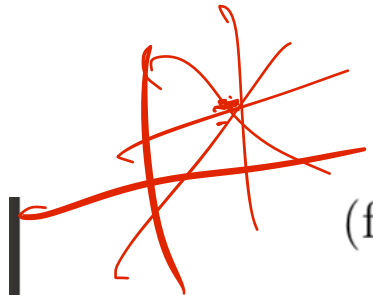
$$y = mx + b$$

1f

$$M = 1$$

$N = 2$  w/ same  $\{x_n, t_n\}$

$$\Phi = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



(f) When doing regression, if the number of examples  $N$  is much larger than the number of features  $M$ , there is always a unique weight vector  $w$  that maximizes the likelihood.

True

False

We still need  $\Phi^T \Phi$  to be full rank (thus invertible)

Counterexample:

$N$  points w/ duplicates; like  $(x_1, t_1)$  repeated  $N$  times

(g) When doing regression, if we place a Normal prior on the precision  $\beta$ ,

188

W Normal prior  
Normal lik  
 $N(y|\beta^T x, \beta)$

→ Normal post  
w/  $y, x \sim N$

(g) When doing regression, if we place a Normal prior on the precision  $\beta$ , N times  
the posterior on  $\beta$  is also Normal.

True

False

Normal prior on  $\beta > 0$   
doesn't make sense.

Not conjugate.

(h) We could create a generalized linear model for binary outcomes  $y_n$

3c

$$X = \sum_{n=1}^N x_n$$

$$k(x_n | \theta)$$

likelihood

$$\theta \in \Theta$$

$$r(\theta | \lambda)$$

 $\lambda$  hyperparam

(c) Define the evidence, the marginal probability of the observed data (integrating away  $\theta$ ), in terms of the functions  $r$  and  $k$ .

log evidence

$$\log p(x) = \int p(x, \theta) d\theta \quad \text{sum rule}$$

$$= \int \left[ \prod_{n=1}^N k(x_n | \theta) \right] r(\theta | \lambda) d\theta$$

Common mistakes:

- $\log$

-

# 3g: Note what consistency actually means

(g) List one advantage and one disadvantage of maximum likelihood estimation.

PRO: Consistent as  $N \rightarrow \infty$   
• Possible to do for any probabilistic model  
• Easier than alternatives (faster, simpler)

CON: Tend to overfit when data small

Consistency

If data drawn from  $K(x_n|\theta)$   
- then as  $N \rightarrow \infty$

$\theta_{ML} \Rightarrow \theta_{true}$

if model wrong  
as  $N \rightarrow \infty$   
 ~~$\theta_{ML}$  will be the true par~~



3h

(h) Describe two reasonable ways to select hyperparameter  $\lambda$  from a grid of possible values  $\lambda_1, \dots, \lambda_G$ . For each one, specify the performance metric and how you would use the  $N$  training examples.

(1) Pick value with largest evidence

$$\max_{\lambda_g \in \{\lambda_1, \dots, \lambda_G\}}$$

$$\log P(x | \lambda_g)$$

see expression from 3c

whole dataset

(2) Pick value w/ best heldout likelihood on validation set, using MAP estimate

Divide data into

80% train:  $x^{tr}$   
20% valid:  $x^{va}$

$$\max_{\lambda_g \in \{\lambda_1, \dots, \lambda_G\}}$$

valid

$$\sum_{n=1}^N \log$$

$$k(x_{n2}^{va} | \theta_{MAP}(x^{tr}, \lambda_g))$$

why can't we use MLE?

$N \times x_n$   
 $k(\cdot)$   
 $F(\cdot)$

$$\sum_n \log k(x_n | \theta)$$

$\theta$  does not depend on  $x$

private chat your answer



# 5b: Lagrange

$$\max_{w \in \mathbb{R}^D} \sum_{d=1}^D w_d c_d, \quad \text{subject to: } \sum_d w_d^2 = 1$$

*w is unit vector*

(b) Solve for  $w$  (a length- $D$  vector) and  $\lambda \neq 0$ , via the system of  $D+1$  equations:  $\frac{d}{dw_1} \mathcal{L}(w, \lambda) = 0, \dots, \frac{d}{dw_D} \mathcal{L}(w, \lambda) = 0$ , and  $\nabla_{\lambda} \mathcal{L}(w, \lambda) = 0$

$$\mathcal{L}(w, \lambda) = \sum_d w_d c_d + \lambda \left(1 - \sum_d w_d^2\right)$$

first deriv = 0  
does not alone  
make maximizer

Sanity check!

$D=1 \quad C=1$

$w = \begin{matrix} +1 \\ -1 \end{matrix} \quad \begin{matrix} +1 \cdot 1 = +1 \\ -1 \cdot 1 = -1 \end{matrix}$

## System of D+1 eqs

$$\frac{\partial}{\partial w_1} \mathcal{L} = 0 \rightarrow \frac{\partial}{\partial w_1} [w^T c + \lambda (1 - \sum_d w_d^2)] = 0$$

*take deriv add  $\lambda 2w_1$*

$$c_1 - \lambda 2w_1 = 0 \rightarrow c_1 = \lambda 2w_1 \quad (1)$$

$$\vdots$$

$$c_D = \lambda 2w_D \quad (D)$$

$$\frac{\partial}{\partial \lambda} \mathcal{L} = 0 \rightarrow \frac{\partial}{\partial \lambda} [w^T c + \lambda (1 - \sum_d w_d^2)] = 0$$

$$1 - \sum_d w_d^2 = 0 \quad (D+1)$$

Now take (1)-(D) & add squares of both sides

$$c_1^2 = \lambda^2 4 w_1^2$$

$$\vdots$$

$$c_D^2 = \lambda^2 4 w_D^2$$

$$\sum_d c_d^2 = \lambda^2 4 \left( \sum_d w_d^2 \right)$$

Sub in  $\sum_d w_d^2 = 1$  from (D+1) we get

$$\sum_d c_d^2 = 4 \lambda^2 \quad \lambda \neq 0$$

$$\lambda = \pm \frac{\sqrt{c_1^2 + c_2^2 + \dots + c_D^2}}{2}$$

Now plug  $\lambda^*$  back into (1)...(D) to solve for  $w_1^* \dots w_D^*$

$$w_1 = \frac{c_1}{2\lambda} = \frac{c_1}{\sqrt{c_1^2 + \dots + c_D^2}}$$

$$w_2 = \frac{c_2}{\sqrt{c_1^2 + \dots + c_D^2}}$$

$$\vdots$$

$$w_D = \frac{c_D}{\sqrt{c_1^2 + \dots + c_D^2}}$$

is  $w^*$  a unit vector?

That's it!

$w$  should be the unit vector that points in same direction as  $c$

$$\lambda = \pm \frac{\sqrt{c_1^2 + \dots + c_D^2}}{2}$$

$$w_1 = \frac{\pm c_1}{\sqrt{c_1^2 + \dots + c_D^2}}$$

$$w_D = \frac{c_D}{\sqrt{c_1^2 + \dots + c_D^2}}$$

# Q&A Viterbi

① Task: find "best" state seq

$$z_1, z_2, \dots, z_T = \underset{z_{1:T}}{\operatorname{argmax}} P(z_{1:T} | x_{1:T}, \theta)$$

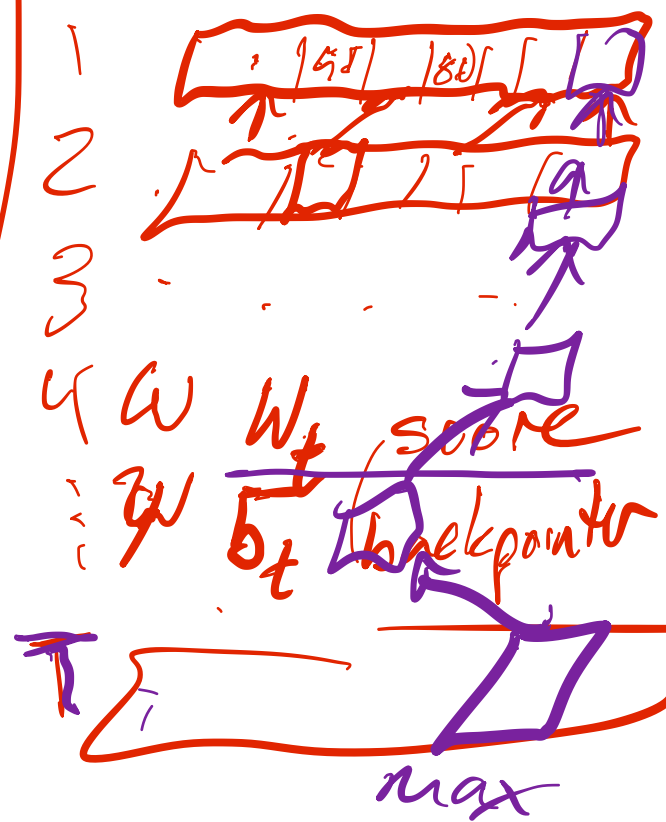
note

$[\hat{z}_1, \dots, \hat{z}_T]$  with  $x_{1:T}$

compared to

$\hat{z}_1, \dots, \hat{z}_T, \hat{z}_{T+1}$  with  $x_{1:T+1}$

② solve w/ dyn prog



# Big Picture

$X_{1:T}$

(1) Estimate  
HMM params  
max likelihood  
(E-M)

$$\theta = \{\pi, A, \mu, \Sigma\}$$

(2) given  $X_{1:T}, \theta$   
what is  $\hat{z}_{1:T}$ ? (Viterbi)

(3) given  $X_{1:T}, \theta$ , what is  $\mathcal{P}(X_{1:T} / \theta)$

# EM

alternate

$$q(z/s) \leftarrow P(z_{1:T} / x_{1:T}, \theta)$$

$\pi, A, \mu, \sigma \leftarrow \text{updated}$

After E step

$$q(z/s) = P(z_{1:T} / x_{1:T}, \theta)$$

$$\alpha = \log P(x_{1:T} / \theta)$$

$$= \mathbb{E}_q \left[ \frac{\log p(x/\theta) + \log p(z/x, \theta) - \log p(z/x, \theta)}{\log p(x/\theta)} \right]$$

optimizing

$$\mathcal{L}(x, s, \pi, A, \mu, \sigma)$$

$$= \mathbb{E}_q \left[ \log p(x, z) - \log q(z) \right]$$

EM algorithm actually computes  $\log p(x|\theta)$

Correct

EM algo

forward algo

for  $t$

$s \leftarrow E_{\text{step}}$

$\theta \leftarrow M_{\text{step}}$

eval

$$\log p(x_{1:T}|\theta)$$

---

Forward algorithm

$$\alpha_{t,k}$$

$$P(z_t = k | x_{1:t}, \theta) =$$

joint  $x_{1:t}, z_t = k$

$$\frac{\quad}{x_{1:t}}$$