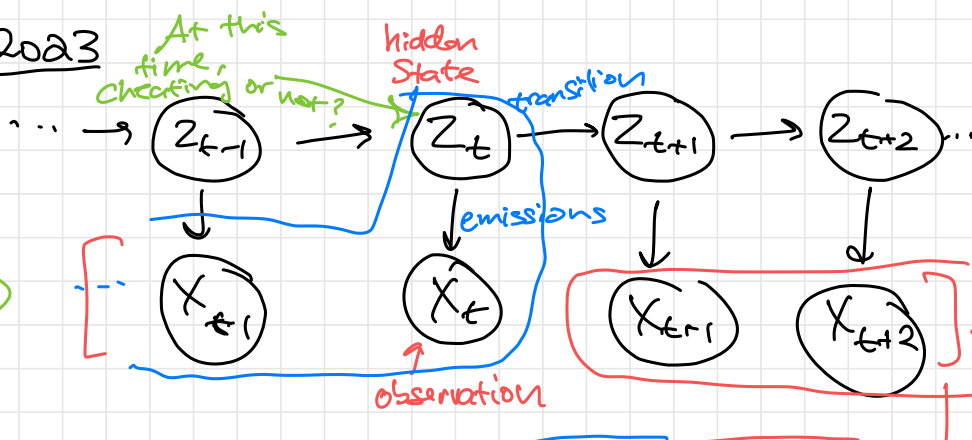


3/30/2023



$$p(z_t | x_{1:T}) \propto p(z_t, x_{1:T}) = \underbrace{p(x_{1:t}, z_t)}_{\alpha_t(z_t) \text{ "forward"}}$$

$$\underbrace{p(x_{t+1:T} | z_t)}_{\beta_t(z_t) \text{ "backward"}}$$

"proportional to"
 Compute for every value of z_t
 then normalize by the sum

Known

α -recursion: Suppose we know $\alpha_{t-1}(i) = p(x_{1:t-1}, z_{t-1} = i)$ for every i

What is $\alpha_t(j)$ for all j marginalize out z_{t-1}

$$\alpha_t(j) = \sum_{i=1}^k p(x_{1:t-1}, z_{t-1} = i) p(z_t = j | z_{t-1} = i) p(x_t | z_t = j)$$

$$= \sum_{i=1}^k \alpha_{t-1}(i) \cdot A_{ij} \cdot p(x_t | z_t = j)$$

Base case: $\alpha_1(j) = p(x_1, z_1 = j) = \underbrace{p(z_1 = j)}_{\text{prior}} \underbrace{p(x_1 | z_1 = j)}_{\text{emission}}$

For $\beta_t(j)$: Compute it based on $\beta_{t+1}(i)$ for every i

All together: To infer $p(Z_t | X_{1:T})$

① Compute α_t 's and β_t 's recursively

② Compute $p(Z_t = j | X_{1:T}) = \frac{\alpha_t(j) \beta_t(j)}{\sum_{i=1}^K \alpha_t(i) \beta_t(i)}$

Learning HMM parameters

Again EM

How to compute good parameters?

Idea of EM: If we had complete data, our lives are easy

| | | | | | |
|----|-----|-----|-----|-----|-----|
| z: | 2 | 1 | 3 | 1 | 2 |
| x: | 1.6 | 6.1 | 2.2 | 9.2 | 1.4 |

Sequence 1

$$p(z_t) = \begin{cases} 1/2 & \text{if } z_t = 1 \\ 1/2 & \text{if } z_t = 2 \\ 0 & \text{if } z_t = 3 \end{cases}$$

$$p(x_t | z_t = 1)$$

= Normal with

$$N = 8$$

$$\sigma^2 = 6$$

| | | | | |
|----|-----|-----|-----|-----|
| z: | 1 | 3 | 2 | 1 |
| x: | 7.4 | 2.1 | 1.2 | 9.7 |

Sequence 2

$$p(x_t | z_t = 2) =$$

Normal with

$$N \text{ of } 1.4$$

$$\sigma^2 = 0.1$$

$$p(z_t | z_{t-1} = 1) = \begin{cases} 0 & \text{if } z_t = 1 \\ 1/3 & \text{if } z_t = 2 \\ 2/3 & \text{if } z_t = 3 \end{cases}$$

Now: Suppose we don't know z_t 's

E-step creates fictitious data (w/ smart guessing)

M-step estimates params on fictitious data

$p(z_t)$:

E-step: Compute $p(z_1 | X_{1:T})$

"pseudocounts"

Suppose you get:

Sequence 1:

$$[0.2, 0.7, 0.1]$$

Sequence 2:

$$[0.6, 0.3, 0.1]$$

Pretend our data looks like

10 copies of Sequence 1 where

- 2 have $z_1 = 1$
- 7 have $z_1 = 2$
- 1 has $z_1 = 3$

10 copies of Sequence 2 where

- 6 have $z_1 = 1$
- 3 have $z_1 = 2$
- 1 has $z_1 = 3$

M-step (Estimate parameters)

Just treat this fictitious data as real
↳ Count things to estimate parameters

$$P(z_1 = 1) = \frac{8}{20}$$

$$P(z_1 = 2) = \frac{10}{20}$$

$$P(z_1 = 3) = \frac{2}{20}$$

For emissions:

E-step: Infer $P(z_t | x_{1:T})$ for every t

Suppose for some t , $x_t = 1.7$

$$P(z_t | x_{1:T}) = [.7, .1, .2]$$

Then we have

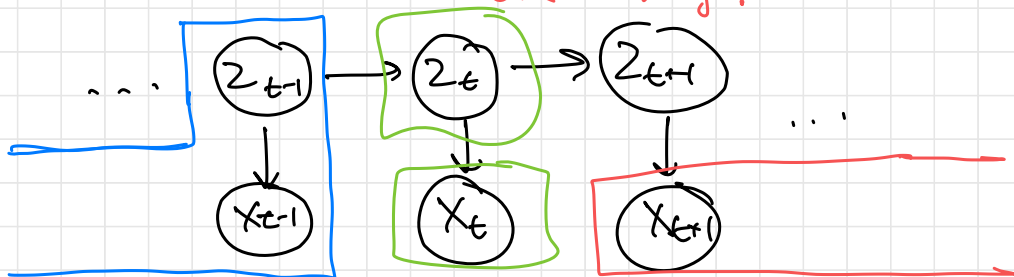
- .7 "counts" of $(z_1 = 1, x_t = 1.7)$
- .1 "counts" of $(z_1 = 2, x_t = 1.7)$
- .2 "counts" of $(z_1 = 3, x_t = 1.7)$

Transitions:

E-step: We want pseudocounts of how many times
state $i \rightarrow$ state j

$$p(z_{t-1}, z_t | X_{1:T})$$

Based on observations, which pairs (z_{t-1}, z_t)
are likely?



$$p(z_{t-1}, z_t | X_{1:T}) = p(x_{1:t-1}, z_{t-1}) \cdot p(x_{t+1:T} | z_t) \cdot$$

$\alpha_{t-1}(z_{t-1})$ $\beta_t(z_t)$

$$\cdot p(z_t | z_{t-1}) p(x_t | z_t)$$

then normalize over all pairs of (z_{t-1}, z_t)

M-step: Use these pseudocounts to estimate
transition probabilities

Announcements

- HW3 starter code
- Section: Neural network optimizers

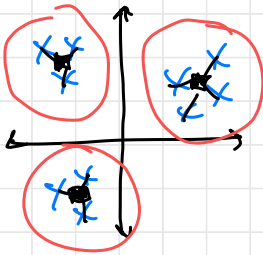
↑
Variants of SGD

$$P(z_{t+1}, x_{1:T}) = P(\underbrace{z_t, x_{1:t}}_A, \underbrace{x_{t+1:T}}_B)$$

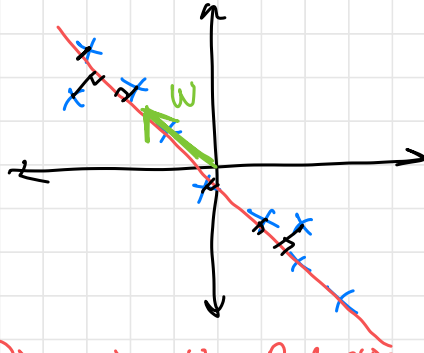
$$P(\underbrace{z_t, x_{1:t}}_A) \cdot P(\underbrace{x_{t+1:T}}_B \mid \underbrace{z_t, x_{1:t}}_A)$$

$$= P(z_t, x_{1:t}) \cdot P(x_{t+1:T} \mid z_t)$$

Dimensionality Reduction



(Clustering)



points in \mathbb{R}^2
but a single
line captures
all of information

Dimensionality Reduction
Given $\{x^{(1)}, \dots, x^{(n)}\} \in \mathbb{R}^d$
Find a lower-dimensional subspace
that preserves most of the information

Method: Principal Component Analysis (PCA)

For now: want to find a good 1-D projection

Key assumption: Data has mean 0

↳ In practice: compute mean of the data & subtract it away

Parameter $w \in \mathbb{R}^d$ that defines the 1-D subspace
 $\|w\| = 1$

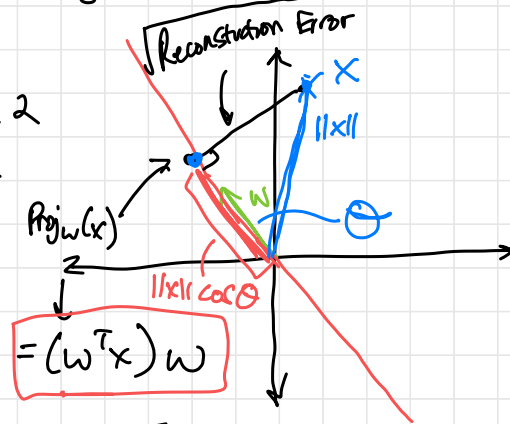
What loss function describes a "good" choice of w ?

Reconstruction Error

$$\sum_{i=1}^n \|x^{(i)} - \text{Proj}_w(x^{(i)})\|^2$$

projection onto w

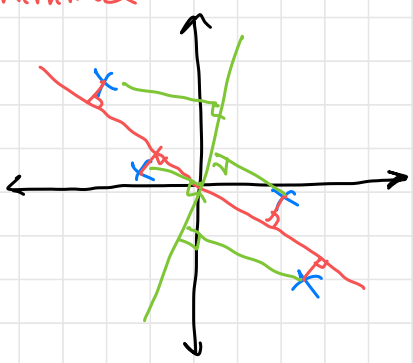
$$= \sum_{i=1}^n \|x^{(i)} - (w^T x^{(i)}) \cdot w\|^2$$



PCA chooses w to minimize this loss function

$$\cos \theta = \frac{w^T x}{\|x\| \|w\|}$$

$$\|x\| \cos \theta = \frac{w^T x}{\|w\|} = w^T x$$



Small recon. error = good
 Large recon error = bad

Equivalent view:

Maximize variance of points after projection

By Pythagorean Theorem:

$$(w^T x)^2 + \text{ReconError} = \|x\|^2$$

↑ maximize this \Leftrightarrow ↑ want to minimize ↑ Fixed

Maximize $\sum_{i=1}^n (w^T x^{(i)})^2$

$\frac{1}{n} \sum_{i=1}^n (w^T x^{(i)})^2$ is just variance of $w^T x$