

11/2/2023 Finishing GMMs, Starting Dimensionality Reduction

Reminder: Naive Bayes

Dataset: $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

maximize \log likelihood of the data:

$$\begin{aligned} \sum_{i=1}^n \log P(x^{(i)}, y^{(i)}) &= \sum_{i=1}^n \log P(y^{(i)}) + \log P(x^{(i)} | y^{(i)}) \\ &= \sum_{i=1}^n \sum_{j=1}^k \underbrace{\mathbb{I}[y^{(i)}=j]}_{\substack{\text{inferred} \\ \text{probability}}} \left(\log P(Y=j) + \log P(x^{(i)} | Y=j) \right) \end{aligned}$$

Back to GMMs: Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$

Idea: Replace the $\mathbb{I}[z_i=j]$ quantity with R_{ij} = inferred probability that $z_i=j$

If: $R_{i1} = 0.8, R_{i2} = 0.2,$

maximize

$$0.8 \cdot [\log P(x^{(i)}, z_i=1)] + 0.2 [\log P(x^{(i)}, z_i=2)]$$

This is the objective called Expected Complete Log likelihood:

$$ECLL(\pi_{1:k}, N_{1:k}, \Sigma_{1:k})$$

$$= \sum_{i=1}^n \sum_{j=1}^k \underbrace{R_{ij}}_{\substack{\text{inferred} \\ \text{probability}}} \log P(x_i = x^{(i)}, z_i = j | \pi_{1:k}, N_{1:k}, \Sigma_{1:k})$$

Goal: maximize ECLL wrt.

$\pi_{1:k}, N_{1:k}, \Sigma_{1:k}$

Today: Compute optimal N_j , initial optimal μ_j, Σ_j

Plan: Take ∇_{N_j} ECU, set n to 0

$$= \sum_{i=1}^n \nabla_{N_j} R_{ij} \log P(x_i = x^{(i)}, z_i = j)$$

$$= \sum_{i=1}^n R_{ij} \nabla_{N_j} \left[\log P(z_i = j) + \log P(x_i = x^{(i)} | z_i = j) \right]$$

only depends on μ_j
not N_j

$$\frac{1}{(2\pi)^{d/2}} \frac{1}{\sqrt{\det(\Sigma_j)}} \cdot \exp\left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right)$$

constant

depends on Σ_j
not N_j

$$= \sum_{i=1}^n R_{ij} \nabla_{N_j} \left[-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right]$$

Fact: $\nabla_x x^T A x = 2Ax$

$$= \sum_{i=1}^n R_{ij} \cdot 2 \Sigma_j^{-1} \cdot (x^{(i)} - \mu_j) \cdot -1 = 0$$

$$= \sum_{i=1}^n R_{ij} \Sigma_j^{-1} (x^{(i)} - \mu_j) = 0$$

multiply both sides by Σ_j (on left)

$$= \sum_{i=1}^n R_{ij} x^{(i)} = \sum_{i=1}^n R_{ij} N_j = \left(\sum_{i=1}^n R_{ij} \right) N_j$$

$$N_j = \frac{\sum_{i=1}^n R_{ij} x^{(i)}}{\sum_{i=1}^n R_{ij}}$$

weighted average of $x^{(i)}$ weighted by probability that it's in cluster j

$P(\text{Example 1 in cluster } j)$
 $\leftarrow P(\text{Example 2 in cluster } j)$
 $\leftarrow \dots$

Expected # of points in cluster j

M -step formulas for π_j & Σ_j :

$$\pi_j = \frac{\sum_{i=1}^n R_{ij}}{n}$$

} Soft version of counting

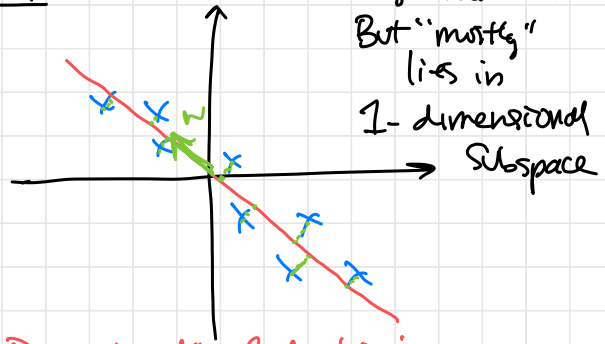
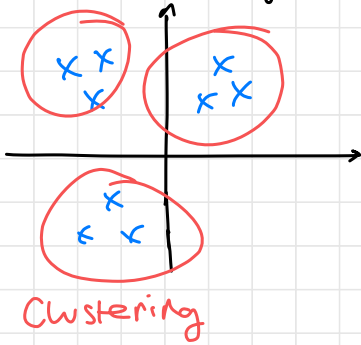
$\frac{\text{\# of points in cluster } j}{\text{total \# of points}}$

$$\Sigma_j = \frac{\sum_{i=1}^n R_{ij} (x^{(i)} - N_j) (x^{(i)} - N_j)^T}{\sum_{i=1}^n R_{ij}}$$

Note: Formula for covariance matrix of dataset is $\frac{1}{n} \sum_{i=1}^n (x^{(i)} - N) (x^{(i)} - N)^T$

\rightarrow This is covariance but weighted by probability of each point being in cluster j

Dimensionality Reduction



Dimensionality Reduction:
Find a low-dimensional subspace
that preserves most of the info in
our dataset

Method: Principal Component Analysis (PCA)

Common use case: high dim data \rightarrow 2D for visualization

Starting point: Try to find best 1-D projection

Key assumption: Data has mean 0

$$\text{i.e. } \frac{1}{n} \sum_{i=1}^n x^{(i)} = 0$$

Ensure by computing mean, subtract it from every example

What is our parameter?

we have 1 parameter vector $w \in \mathbb{R}^d$
that defines the 1-D subspace we will project onto

Force $\|w\| = 1$

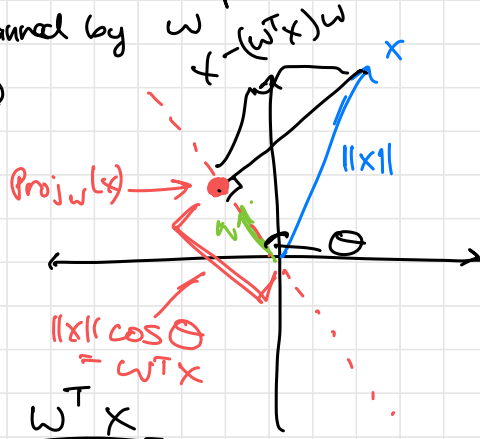
What is a good loss function for w ?

"Reconstruction Error": How well can we reconstruct $\{x^{(1)}, \dots, x^{(n)}\}$ based on only the projections of $x^{(1)}, \dots, x^{(n)}$ onto subspace spanned by w

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

minimize

project onto $\text{Span}(w)$



minimize

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

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

$$\Rightarrow \|x\| \cos \theta = w^T x$$

$$\Rightarrow \text{Proj}_w(x) = \underbrace{(w^T x)}_{\text{right distance}} \underbrace{w}_{\text{right direction}}$$

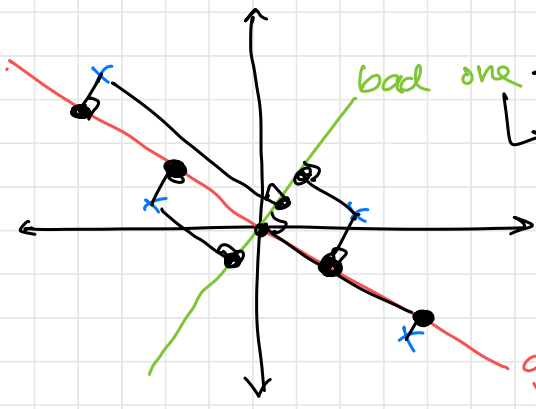
Equivalent view: minimizing reconstruction error is same as maximizing variance of points after projection

Pythagorean Theorem: $(w^T x)^2 + \text{ReconError for example } i = \|x\|^2$

want to maximize want to minimize Fixed

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

Note: $\frac{1}{n} \sum_{i=1}^n (w^T x^{(i)})^2 = \text{Variance of } w^T x$
(since $E[x] = 0$)



large distance =
high reconstruction error X
less spread out after projection
= low variance X

good subspace → small distance
to line = small
reconstruction
error ✓

↓
resulting points
spread out =
high variance ✓