

4/2/2024: Learning GMMs

Given data, how to learn $\pi_{1:k}$, $\mu_{1:k}$, and $\Sigma_{1:k}$?

Algorithm: Expectation-Maximization (EM)

Used whenever you have:

- ① Latent variables (that are unknown) z_1, \dots, z_n
- ② Unknown parameters (need to learn) $\pi_{1:k}, \mu_{1:k}, \Sigma_{1:k}$

Strategy: Alternate between guessing each one

① E-step ("Expectation"):

Infer latent variable distribution using current guess of parameters

In K-Means

\approx make assignments based on current centroids

② M-step ("Maximization"):

Choose parameters that fit the data best based on inferred distribution of latent variables

\approx choose new centroids based on new assignment

E-step For each $i=1, \dots, n$, we infer:

$$R_{ij} = p(z_i = j | x_i = x^{(i)}; \text{current guess of } \pi_{1:k}, \mu_{1:k}, \Sigma_{1:k})$$

\uparrow
Run inference using Bayes Rule

1	0.6	0.1	0.3
...			
n			

$\uparrow \uparrow \uparrow$
prob of being from cluster j

M-step

Inputs are:

- Actual values of X_i 's
- Inferred distribution of all Z_i 's

Goal: Find the best π, N, Σ

Imagine: we actually had true value of Z_i for each i . Then we could apply Maximum Likelihood Estimation:

$$\begin{aligned} \text{Maximize } & \sum_{i=1}^n \log P(Z_i, X_i) \\ &= \sum_{i=1}^n \log P(Z_i) + \log P(X_i | Z_i) \\ &= \sum_{i=1}^n \sum_{j=1}^K \mathbb{1}[Z_i=j] \left[\log P(Z_i=j) + \log P(X_i | Z_i=j) \right] \end{aligned}$$

Actual M-step: Approximate the MLE objective using inferred R_{ij} 's

Expected Complete Log-Likelihood (ECLL)

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

$$\sum_{i=1}^n \sum_{j=1}^K R_{ij} \cdot \left[\log P(Z_i=j) + \log P(X_i=x^{(i)} | Z_i=j) \right]$$

how much to weight the possibility that $Z_i=j$

"complete log likelihood" includes x 's and Z 's

Expected value with R_{ij} 's as weights

Now: Derive formula for N_j to maximize ELL, discuss formulas for Π_j & Σ_j

Plan: Take ∇_{N_j} ELL, set to 0, solve for N_j

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

Depends only on Π_j
 ∇ will be 0

For each j , this depends on that cluster's N & Σ
 \Rightarrow Ignore all terms except for the N_j we care about

$$= \sum_{i=1}^n R_{ij} \nabla_{N_j} \log P(X_i=x^{(i)} | Z_i=j)$$

Plug in multivariate gaussian pdf formula

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

constant
 doesn't depend on N_j

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

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

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

multiply by Σ_j on left

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

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

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

★

$P(x_1 \text{ is in cluster } j)$
 $P(x_2 \text{ is in cluster } j)$
 \vdots
 $+ P(x_n \text{ is in cluster } j)$

Expected # of points in cluster j

Weighted average of $x^{(i)}$
 weighted by the
 probability that each
 $x^{(i)}$ is in cluster j

M-step formulas for π_j & Σ_j

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

"Soft" version of

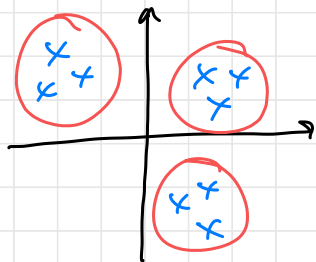
$\frac{\text{\# of points in cluster } j}{\text{total \# points in data}}$

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

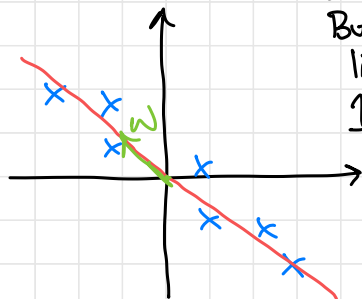
Formula for covariance matrix is $\frac{1}{n} \sum_{i=1}^n (x^{(i)} - N)(x^{(i)} - N)^T$

"Weighted average" version of covariance matrix formula

Dimensionality Reduction



Clustering



Data in \mathbb{R}^2
But "mostly" lies on
1-dimensional
Subspace

Dimensionality Reduction:
Find a low-dimensional subspace
that preserves most information
in the dataset

Method: Principal Components Analysis (PCA)

Commonly: map high-dim data \rightarrow 2 dimensions

Starting Point: Try to find best 1-D subspace

Key Assumption: Data has mean 0

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

Enforce by computing mean of data,
then subtract it from each example

What is our parameter to learn?

Learn a single parameter vector $w \in \mathbb{R}^d$
that defines the 1-D subspace to project onto

original dimension of data

We'll focus only on w where $\|w\| = 1$
since we only care about direction

"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 defined by W

goal:
make this
"lost info" as
little as possible

