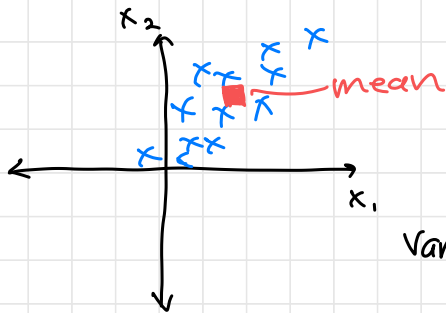


10/31/2023: Gaussian Mixture Models (GMM)

1. How to think of clusters w/ non-spherical shape?
2. What is a GMM?
3. Inference - Assign datapoint to cluster?
4. Learning - Decide shapes of clusters & assign points to clusters at same time



$$\begin{aligned}\text{Mean } \mu &= E[X] \\ &= \frac{\sum_{i=1}^n x^{(i)}}{n}\end{aligned}$$

$$\begin{aligned}\text{Variance of } x_1 &= E[(x_1 - \mu_1)^2] \\ &= \frac{1}{n} \sum_{i=1}^n (x_1^{(i)} - \mu_1)^2\end{aligned}$$

Covariance between x_1 & x_2

$$= E[(x_1 - \mu_1)(x_2 - \mu_2)] = \frac{1}{n} \sum_{i=1}^n (x_1^{(i)} - \mu_1)(x_2^{(i)} - \mu_2)$$

Correlation between x_1 & x_2 =
$$\frac{\text{Covariance}(x_1, x_2)}{\sqrt{\text{Var}(x_1) \text{Var}(x_2)}}$$

Covariance > 0

\Rightarrow

positively correlated,

Covariance < 0

\Leftarrow

negatively correlated

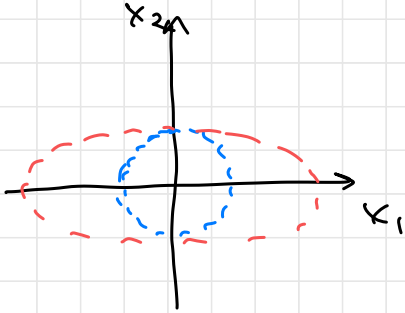
Covariance Matrix
$$\Sigma = \begin{pmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) \end{pmatrix}$$

matrix that summarizes

- variance in every dimension

- covariance between every two dimensions

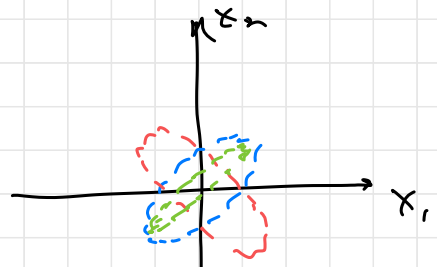
} captures shape of distribution



$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}$$

large variance for x_1
small variance for x_2



$$\Sigma = \begin{pmatrix} 1 & 1/2 \\ 1/2 & 1 \end{pmatrix}$$

Correlation is 0.5

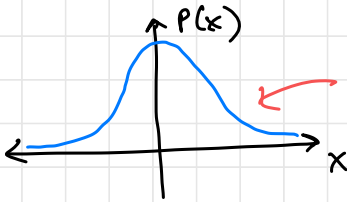
$$\Sigma = \begin{pmatrix} 1 & -1/2 \\ -1/2 & 1 \end{pmatrix}$$

Correlation is 0.9

$$\Sigma = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix}$$

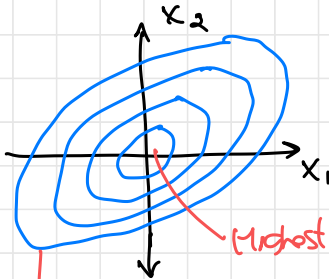
Multivariate Gaussians

In 1 Dimension: Univariate Gaussian distribution



Determined by mean = μ
variance = σ^2

In d-dimensions: Multivariate Gaussian distribution



Determined by mean $\mu \in \mathbb{R}^d$
and covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$

points with same probability density

constants

Formula for probability density:

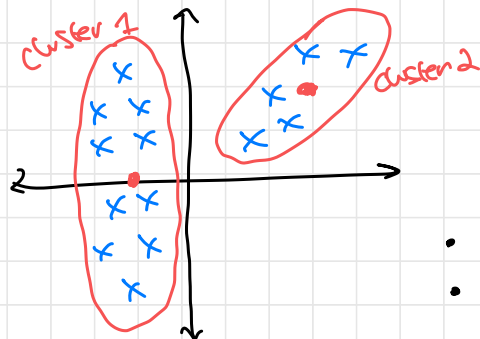
$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2}} \cdot \frac{1}{\sqrt{\det(\Sigma)}} \cdot \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

Compare with 1-D case:

$$\frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\sqrt{\sigma^2}} \cdot \exp\left(-\frac{1}{2} \cdot \frac{(x-N)^2}{\sigma^2}\right)$$

What is a GMM?

Goal: given data, produce clusters with custom shapes



How? For each cluster, we will learn:

- π_1, π_2 : size of each cluster
 - μ_1, μ_2 : center of each cluster
 - Σ_1, Σ_2 : covariance of each cluster
- } params of GMM model

For this dataset, we want to learn:

$$\pi_1 = \frac{2}{3}, \quad \pi_2 = \frac{1}{3}$$

$$\mu_1 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 9 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}$$

Probabilistic Story of GMMs:

For each $i=1, \dots, n$: # of examples in dataset

① Randomly sample cluster Z_i
where $P(Z_i = j) = \pi_j$

② Randomly sample example X_i
from multivariate Gaussian with mean μ_{Z_i} covariance Σ_{Z_i}

We only observe that $X_i = x^{(i)}$ in the dataset
random variable X_i observed value $x^{(i)}$

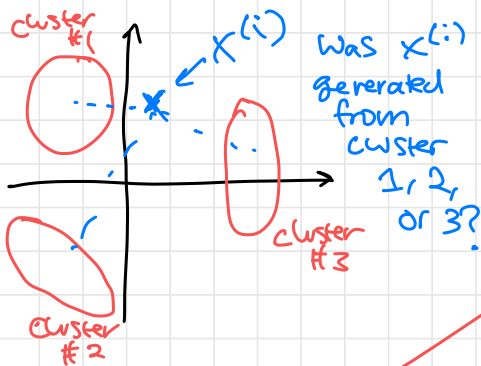
We never directly observe Z_i ("latent variable")

Inference - Inferring a probability distribution of a latent random variable conditioned on observed random variables

For GMMs: Given:

- Observed value $x^{(i)}$ for X_i
- Known parameters $\pi_{1:k}, \mu_{1:k}, \Sigma_{1:k}$

Compute $P(Z_i | X_i = x^{(i)}; \pi_{1:k}, \mu_{1:k}, \Sigma_{1:k})$



$$P(Z_i = j) = \pi_j$$

How? Bayes Rule

$$P(Z_i = j | X_i = x^{(i)}) = \frac{P(Z_i = j) P(X_i = x^{(i)} | Z_i = j)}{\sum_{b=1}^K P(Z_i = b) P(X_i = x^{(i)} | Z_i = b)}$$

$P(X_i = x^{(i)} | Z_i = j)$
= multivariate Gaussian pdf
with mean μ_j ,
covariance Σ_j

Result: for each $x^{(i)}$, we get distribution

$$\begin{aligned} P(Z_i = 1 | X_i = x^{(i)}) &= 0.6 \\ P(Z_i = 2 | X_i = x^{(i)}) &= 0.1 \\ P(Z_i = 3 | X_i = x^{(i)}) &= 0.3 \end{aligned}$$

produce a "hard" assignment
by choosing most likely option

assign to cluster #1

"soft assignment" / "soft clustering"
b/c it's all probabilistic

"hard assignment"

Finally: How do we learn $\mu_{i:k}$, $N_{i:k}$, $\Sigma_{i:k}$?

Algorithm: Expectation-Maximization (EM)

Very general method whenever you have:

- Latent variables
- Unknown parameters

Strategy: Alternate between updating each one

① E-step: Infer latent variable distribution using current guess of parameters

↕ Alternate ↕

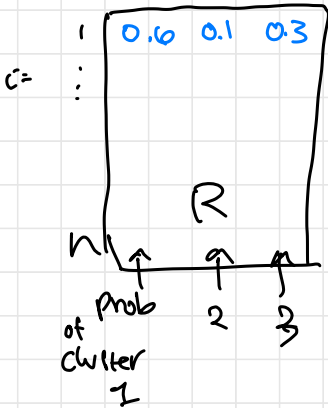
↗ make assignments based on ← means centroids

② M-step: Choose parameters that best fit the data based on inferred distribution of latent variables

↗ Choose new k-means centroids based on newest assignment

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

$$R_{ij} = P(Z_i = j \mid X_i = x^{(i)}; \text{current guess of parameters})$$



Uses the inference procedure from before

M-step

Takes in :

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

} less info than in supervised learning

Can't do MLE w/o actual values of Z_i 's

But something similar:

We can maximize Expected Complete Log-likelihood (ECLL)

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

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \log P(X_i = x_i^{(i)}, Z_i = j; \pi, \mu, \Sigma)$$

Prob of example i in cluster j

complete log likelihood (includes X_i & Z_i)

expectation