

Week - 5

~~GMM~~ Multivariate Gaussian.

$$g(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

GMM

$$\lambda = \sum_{k=1}^K w_k g_k(x, \mu_k, \Sigma_k)$$
$$\sum_{k=1}^K w_k = 1$$

How to train a GMM

EM - Expectation Maximization

↓

Statistics

↓

Fitting GMM to given data

1) Iterative Procedure

In each Iteration

→ compute statistics

Wts

→ Mean

→ σ

Objective! maximizing representation of data of x by λ

It. therefore λ_1, λ_2

$$P(x/\lambda_1) \neq P(x/\lambda_2)$$

$$P(x/\lambda_2) > P(x/\lambda_1)$$

As λ_2 is better to λ_1

Convergence

$$P(x/\lambda_2) - P(x/\lambda_1) \leq \text{Threshold.}$$

To begin with $x = \{x_1, \dots, x_m\}$
+ Initial Model λ - GMM

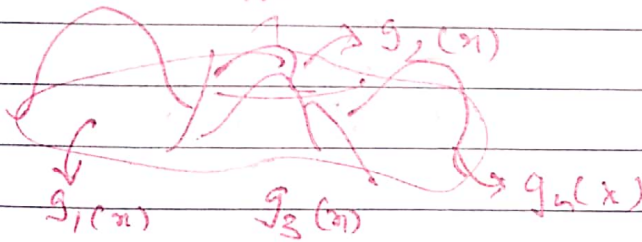
Ways to Generate Initial Model \rightarrow Random
 \rightarrow k-means clustering algo

"k" \rightarrow M data.
partition \downarrow non-overlapping
k partition. } Random.

k-means

cluster datapoints into "k" non-overlapping data.

$$\lambda = \sum_{k=1}^K w_k g_k(x_i, \mu_k, \Sigma_k)$$



$$\lambda = \sum_{k=1}^K w_k g_k(x, \mu_k, \Sigma_k)$$

EA

$$x = \{x_1, \dots, x_m\}$$

$P(k/x_i, \lambda) = x_i$ belonging to k^{th} mixture model in λ model.

$$= \frac{w_k g_k(x_i | \mu_k, \Sigma_k)}{\sum_{m=1}^K w_m g_m(x_i | \mu_m, \Sigma_m)}$$

$$\bar{w}_k = \frac{1}{M} \sum_{i=1}^M P(k/x_i, \lambda)$$

$$\bar{\mu}_k = \frac{\sum_{i=1}^M P(k/x_i, \lambda) \cdot x_i}{\sum_{i=1}^M P(k/x_i, \lambda)}$$

$$\bar{\Sigma}_k = \frac{\sum_{i=1}^M P(k/x_i, \lambda) (x_i - \bar{\mu}_k)(x_i - \bar{\mu}_k)^T}{\sum_{i=1}^M P(k/x_i, \lambda)}$$

$$\bar{\lambda} = \sum_{k=1}^K \bar{w}_k g_k(x_i, \bar{\mu}_k, \bar{\Sigma}_k)$$

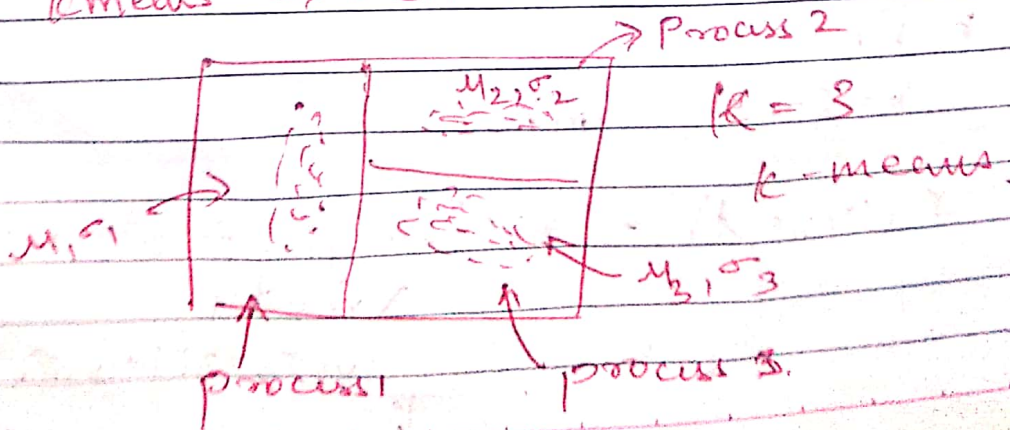
$$P(x/\bar{\lambda}) = P(x, \bar{\lambda}) \geq \text{Threshold}$$

$$P(x/\bar{\lambda}) = \frac{1}{M} \sum_{i=1}^M P(x_i/\bar{\lambda})$$

↑
 $\max_k P(x_i/\bar{\lambda})$

Clustering by GMM

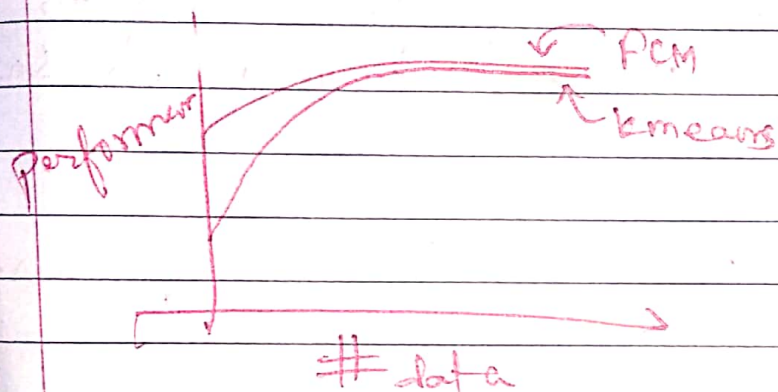
kmeans \rightarrow FCM \rightarrow GMM



FCM Advantages over k-means.

mean vector in FCM is more stable as it gives more data for mean vector calculation.

GMM : takes care spread (σ) of data
So better fit than FCM.



FCM performs better
in case of less
data.

when we go for GMM we have μ , σ^2 (variance)

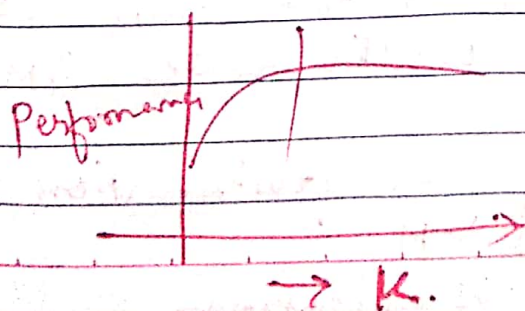
μ - First order statistics (FOS)

σ^2 - Second order (SOS)

as GMM has SOS it is better representative of model

data requirement Grows in exponential order
SOS requires much more data than FOS.

How to choose 'k' in GMM



Cluster purity:-

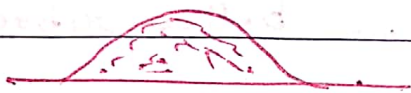
~~Recall~~ Depending value of k that many clusters are created; we also have reference which point belongs to which cluster

ideally if there are 3 clusters, when n but $k=3$, data points belonging to particular cluster (reference) falls to that cluster then cluster purity is said to 100%

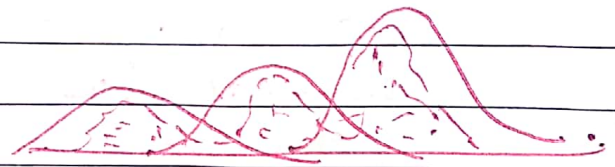
performance is measured in terms of cluster purity.

~~Q. How to set~~

Why Gaussians are used so extensively.



a have Gaussian
A fit Gaussian



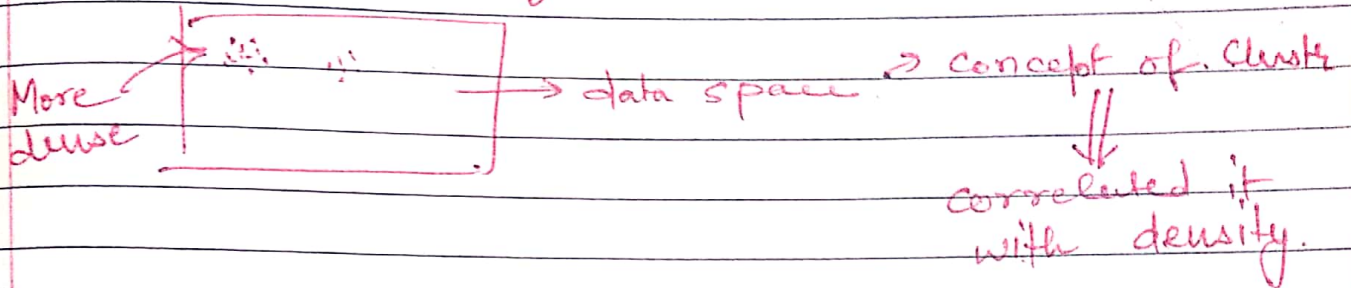
A mixture of Gaussian
model data can be
- fit. ($K=3$)

So a non-Gaussian distribution can be
model it with multivariate GMM.

~~Thus~~ So theoretically all distribution can be
represented by GMM.
⇒ Data can be Gaussianisation

"Density" Based Clustering

↓
literal meaning

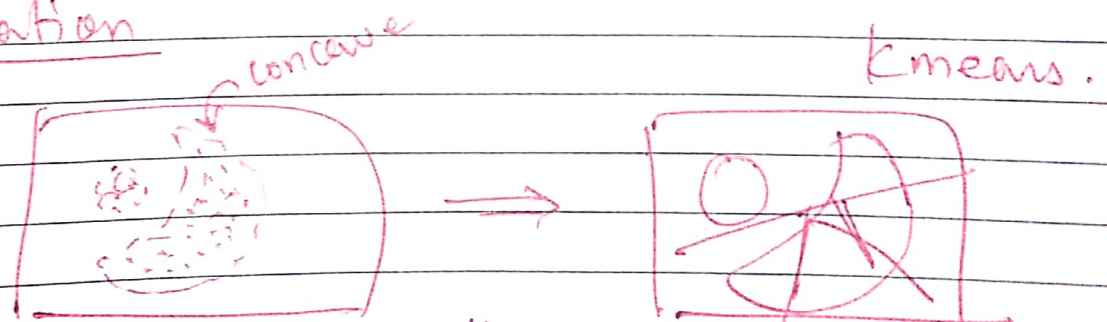


More data points in neighbours tells us about presence of a ~~data~~ cluster

Denseness is also attribute to come up with a type of clustering (or exploited)

DBSCAN (Density based spatial clustering Application & noise)

Motivation



Can't distribute a concave surface
In Kmeans, FCM, GMM there is no means to reject outliers

in Density based we could argue that it is not dense enough to be a cluster so an outlier.

- Prone to outlier: outliers are assigned to cluster, after that they pull mean vectors towards them which makes their cluster prone to far point than near points.

DBSCAN can

DBSCAN Density based can find out any arbitrary shape cluster without being affected by noise.

Density Based Spatial Clustering of Application 4 Noise

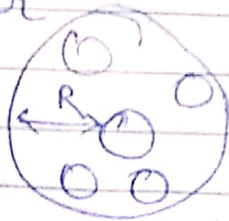
R: Radius of Neighbour hood

M: Min number of neighbour

Intuition/Idea.

If a point is in a cluster of it should be near to lots of other points in that cluster

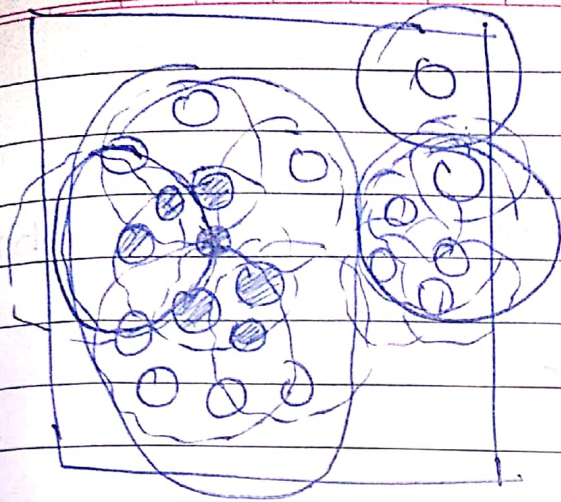
R



If a taking any point if in a given Radius 'R' it have more than M neighbours we call it Dense area or ~~area~~ cluster.

3 types of Data point

- Core ~~border~~
- Border
- Outlier



$$R=2$$

$$M=6$$

Core : If there are M points or more in a given radius.

Border If there are points but less than M it is border

Outlier : if points in neighbourhood is zero

Soft Computing

Core points surrounded by border points generate arbitrary shape.

Disadvantage:

→ Need to select 2 parameter ' R ' & ' M '

Advantage.

- Arbitrary shaped cluster
- Robust to Classifier
- Doesn't require specification of cluster