


UNSUPERVISED

→ Don't have labels
→ it's all about how to group

2 types Unsupervised learning
grouping the data is called clustering
Reducing the dimensionality is another type

1. clustering the data
2. Hierarchical and Density based clustering
3. Gaussian Mixture model and cluster validation
4. principal Component analysis
5. Random projection and Independent Component analysis

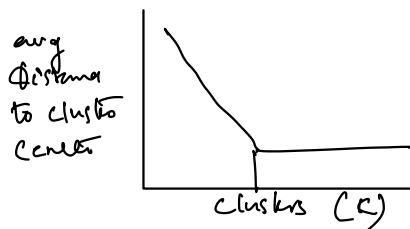
clustering - K MEANS algorithm

Recommending new items eg
Grouping the data by
closeness of the data

Choosing K value.

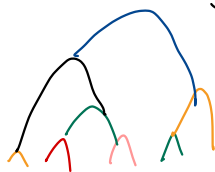
1. depends prior knowledge
2. Elbow method (no prior knowledge)

Elbow Method



Hierarchical density based clustering

DBSCAN - Density b



hierarchical clustering



Hierarchical clustering

Single link clustering

- * each point is a cluster
- * Smallest d/b/w the points
- * atleast 2 clusters

Complete link:-

from sklearn import datasets, cluster
from scipy.cluster.hierarchy import
dendrogram, Ward, Single

X =
linkage_matrix = Ward(X)
dendrogram(linkage_matrix)
plt.show()

from sklearn import AgglomerativeClustering
AgglomerativeClustering (linkage = 'complete')
ward = AgglomerativeClustering (n_clusters = 3)

word_predict = word_fit.predict(x)

from sklearn.metrics Adjust_rand_score

from sklearn.cluster import linkage

linkage_type = 'ward'

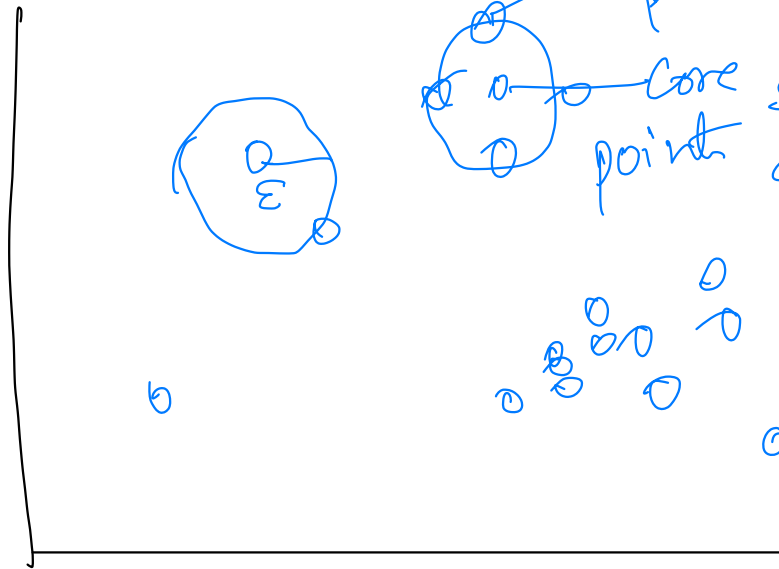
linkage_matrix = linkage(x, linkage_type)

from scipy.cluster.hierarchy import dendrogram

import matplotlib.pyplot as plt

plt.figure(figsize=(22,18))

dendrogram(linkage_matrix)



Inputs

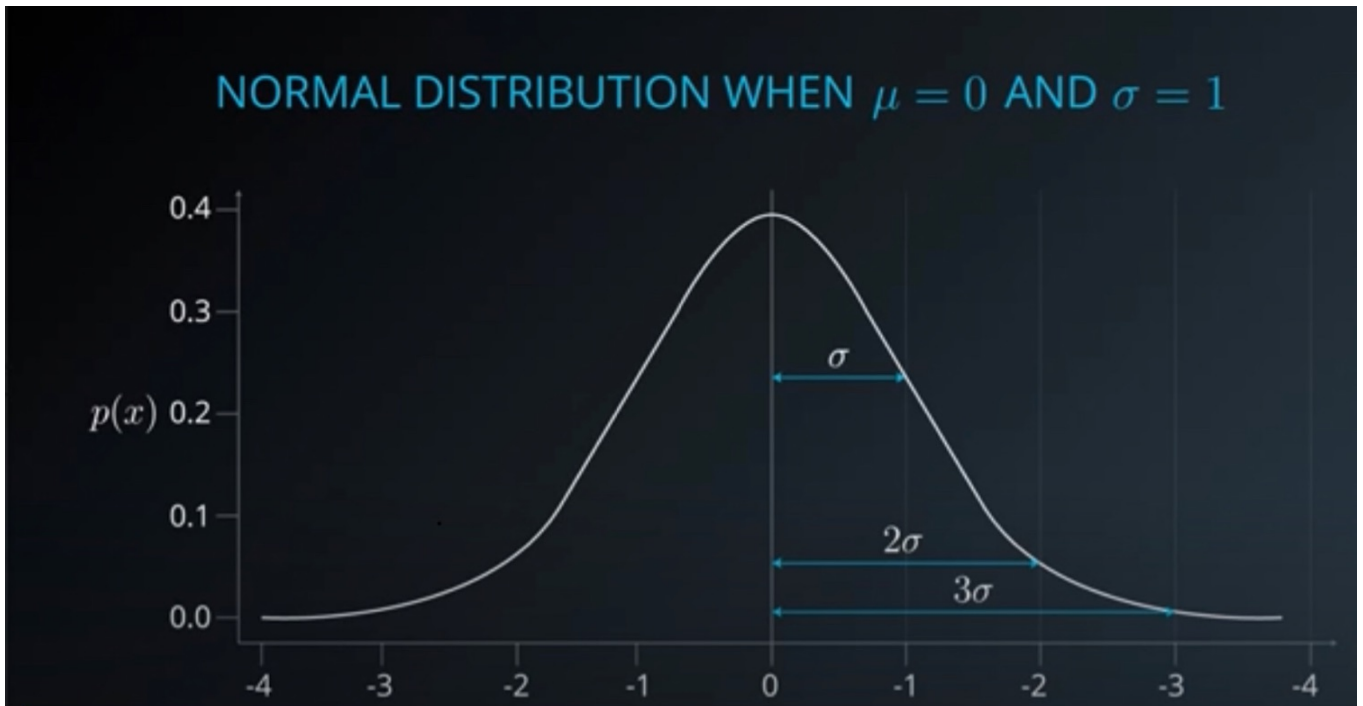
Epsilon -

search distance
around point

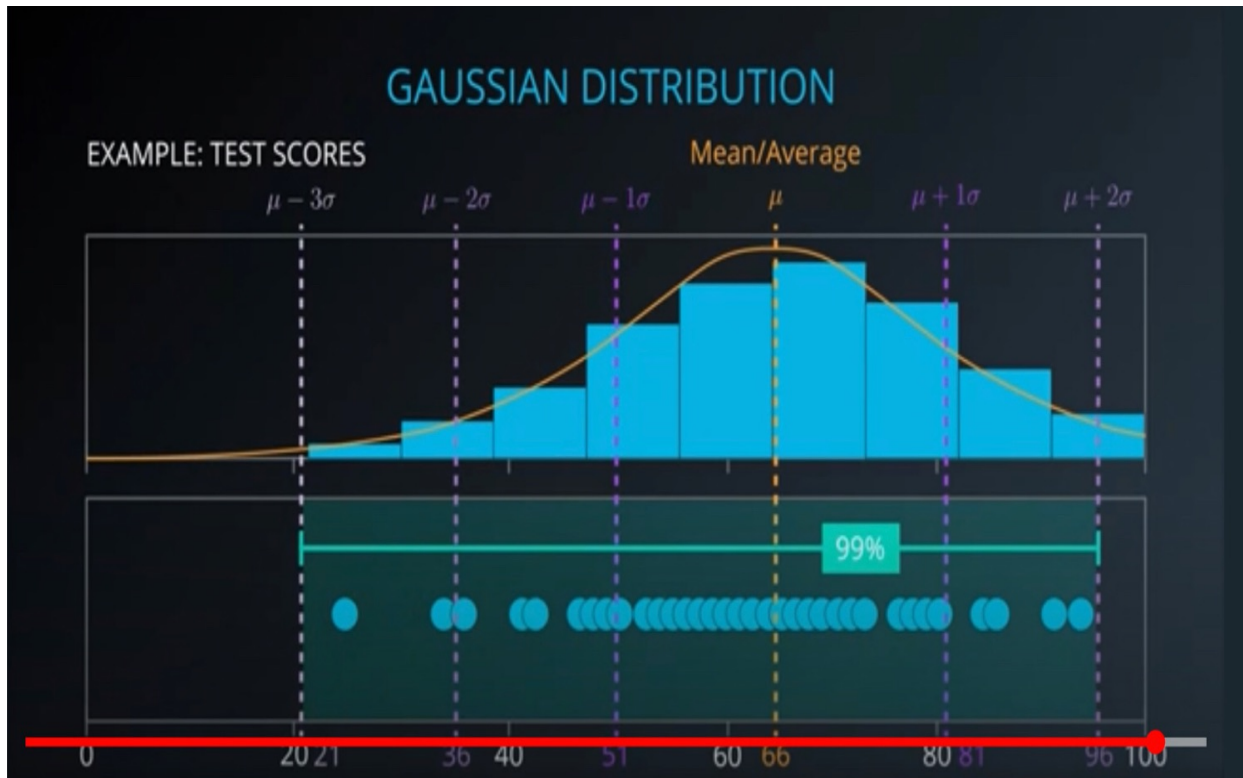
Minpts = 5

Minimum no
of required to form
density cluster

GAUSSIAN MIXTURE MODEL & CLUSTERING



Plotting gaussian distribution

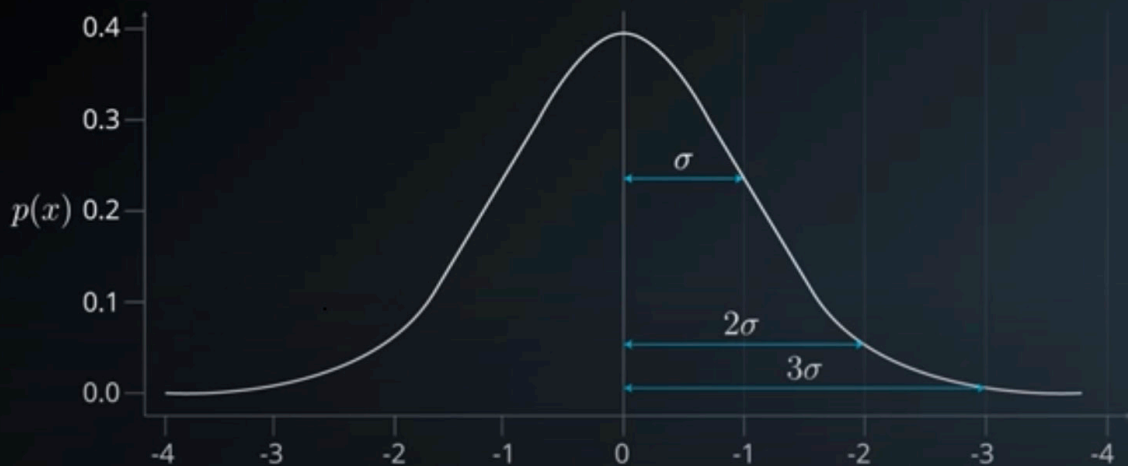


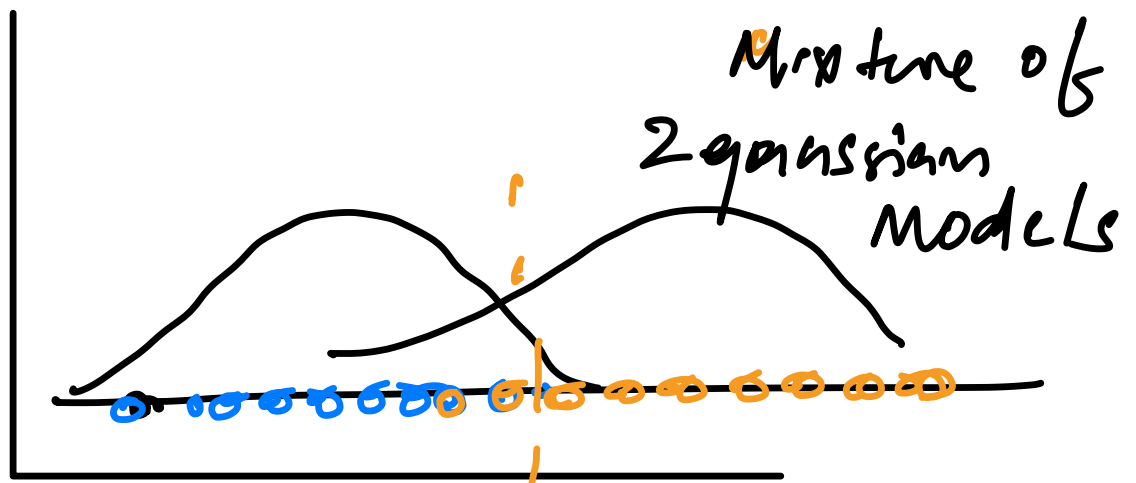
$\mu = \text{mean}$
 $\sigma = \text{std}$

$\mu - 3\sigma$
 \downarrow
 $\mu + 3\sigma$
 $\text{---} 99\%$

$\mu - 1\sigma$ — $\mu + 1\sigma$ — 68%
 $\mu - 2\sigma$ — $\mu + 2\sigma$ — 95%

NORMAL DISTRIBUTION WHEN $\mu = 0$ AND $\sigma = 1$





1. Initialize gaussian distribution
2. Soft data clustering
3. Re estimate gaussian Maximization
4. Evaluate log likelihood to check the Con

Initialize gaussian distribution

1. Apply K-MEANS to the dataset identify the clusters
pick Random Mean and Std Deviation
2. Soft clustering find the membership of the clusters
use the probability density function of the normal distribution function to identify the membership of the cluster
3. Re estimate the parameters for the gaussians
find the new mean that comes from the weighted avg of the points and do find the membership of the point
apply this process
4. Evaluate log likelihood

From sklearn-mixture import Gaussian_Mixture

Gaussian-mixture (n-components = 5)

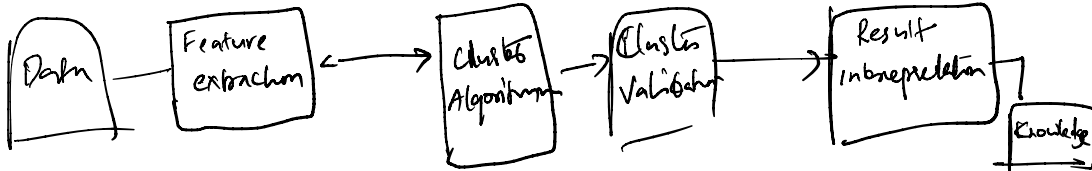
GMM Clustering

Advantages

1. soft clustering
2. cluster shape flexibility

Disadvantages

1. Sensitive to initialization values
2. possible to converge local optimum



Adjusted Rand index — External indices

Ranges from -1 to 1

Silhouette coefficient | Internal indices

Ranges from -1 to 1

Silhouette coefficient helps to identify the best values for
for suitable k value silhouette coefficient
will be higher.