

Cluster Validation

Done By Venkatesh E
Priyadarshan

For cluster analysis, the question is how to evaluate the
“goodness” of the resulting clusters?

Then why do we want to evaluate them?

—>To avoid finding patterns in noise.

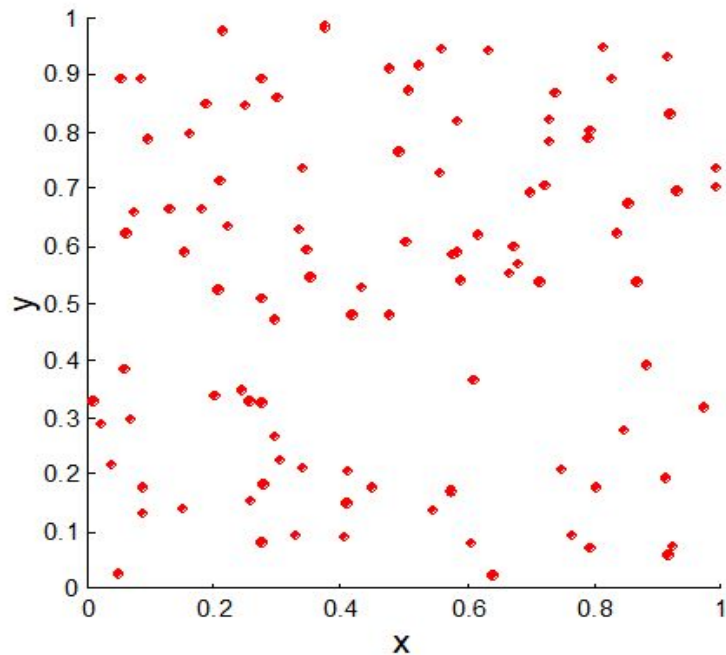
—>To compare clustering algorithms.

—>To compare two sets of clusters.

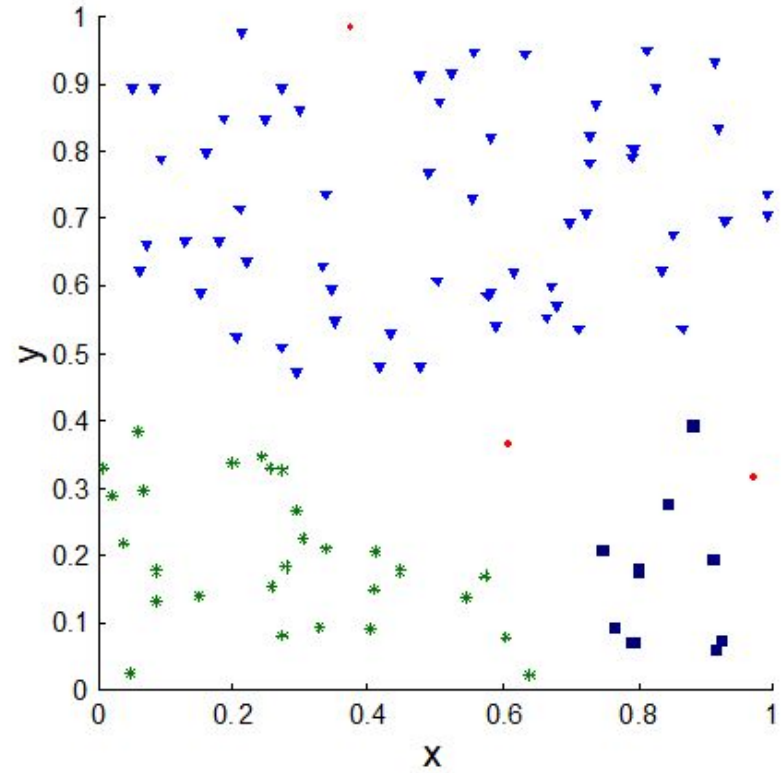
—>To compare two clusters.

SOME CLUSTERS WITH RANDOM DATAS

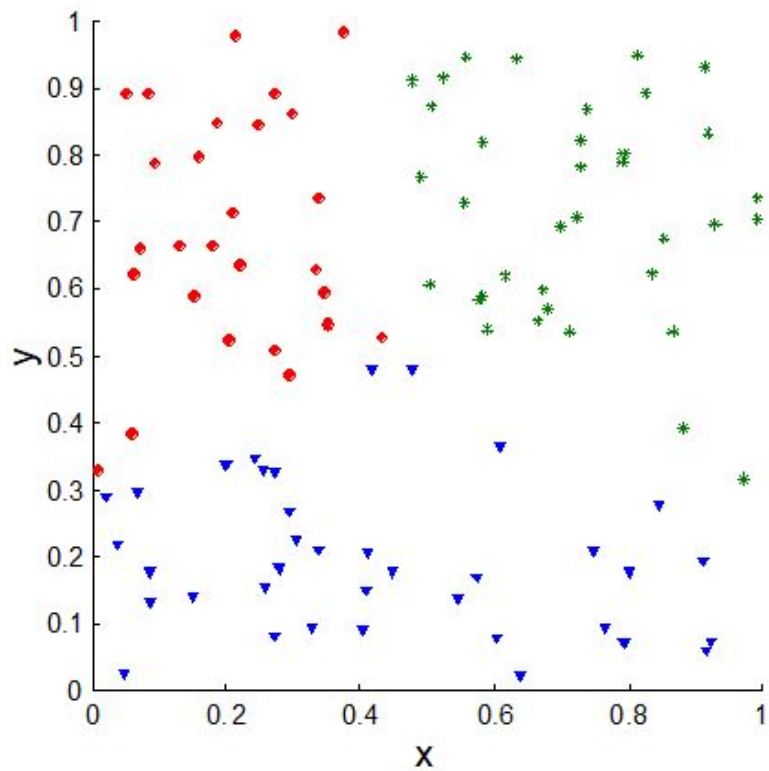
Random Points



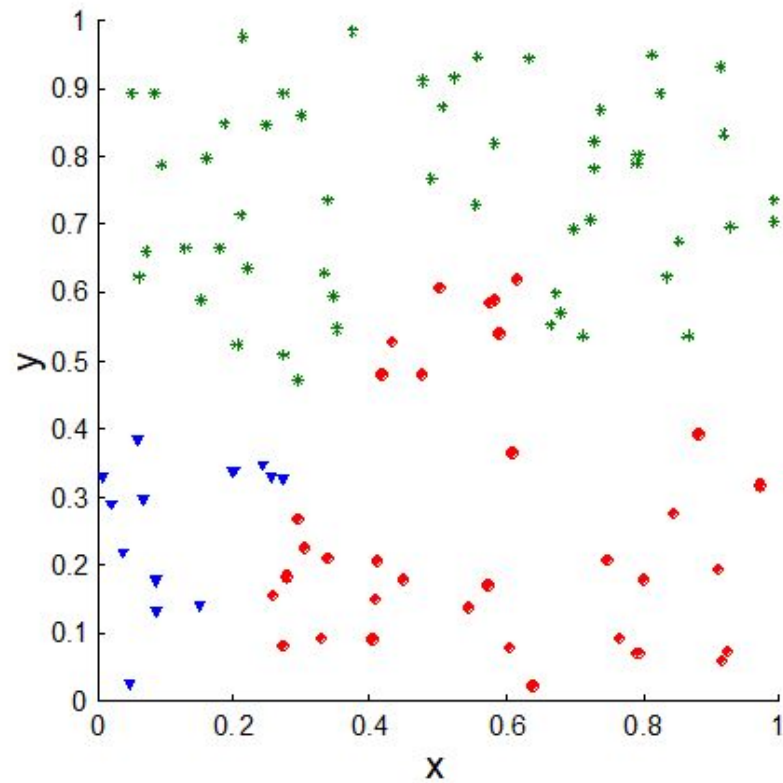
DBSCAN



K-Means



Complete Link



DIFFERENT ASPECTS OF CLUSTER VALIDATION

Aspect - 1

Determining the clustering tendency of a set of data, i.e., distinguishing whether non-random structure actually exists in the data.

Clustering tendency

Clustering tendency assessment **determines whether a given dataset contains meaningful clusters**

Aspect - 2

Comparing the results of a cluster analysis to externally known results, e.g., to externally given class labels.

Aspect - 3

Evaluating how well the results of a cluster analysis fit the data *without* reference to external information. We should use only the data given.

Aspect - 3

Evaluating how well the results of a cluster analysis fit the data *without* reference to external information. We should use only the data given.

Aspect - 4

Comparing the results of two different sets of cluster analyses to determine which is better.

Aspect - 5

Determining the 'correct' number of clusters.

For 2, 3, and 4, we can further distinguish whether we want to evaluate the entire clustering or just individual clusters.

FRAMEWORK FOR CLUSTER VALIDATION

NEED A FRAMEWORK TO INTERPRET ANY MEASURES

For example,

If our measure of evaluation has the value, 10,
is that good, fair, or poor?

STATISTICS PROVIDE A FRAMEWORK FOR CLUSTER VALIDATION

- The more “atypical” a clustering result is, the more likely it represents valid structure in the data
- Can compare the values of an index that result from random data or clusterings to those of a clustering result.
 - If the value of the index is unlikely, then the cluster results are valid
- These approaches are more complicated and harder to understand.

FOR COMPARING CLUSTERS FRAMEWORK IS LESS NECESSARY

However, there is the question of whether the difference between two
index values is significant

MEASURE OF CLUSTER VALIDITY

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.
 - **External Index:** Used to measure the extent to which cluster labels match externally supplied class labels.
 - **Internal Index:** Used to measure the goodness of a clustering structure *without* respect to external information.
 - **Relative Index:** Used to compare two different clusterings or clusters.
 - Often an external or internal index is used for this function, e.g., SSE or entropy
 - Sometimes these are referred to as **criteria** instead of **indices**
 - However, sometimes criterion is the general strategy and index is the numerical measure that implements the criterion.
-

External Validation

Algorithm 21.4: Algorithm for matching partitions and clusters

MatchPartitionCluster ($P, C, match$):

```
1 foreach  $p \in P$  do
2    $match(p) \leftarrow \emptyset$ 
3   foreach  $c \in C$  do
4      $overlap(p, c) \leftarrow \frac{|p \cap c|}{|p|}$ 
5   while  $overlap \neq \emptyset$  do
6      $(p_{max}, c_{max}) \leftarrow GetMaxOverlap(overlap)$ 
7      $match(p_{max}) \leftarrow c_{max}$ 
8      $overlap \leftarrow overlap - \{overlap(p_{max}, *), overlap(*, c_{max})\}$ 
```

CORRELATION MEASURES

Is a statistical technique which determines how one variables moves/changes in relation with the other variable.

HUBERTS TAU STATISTICS

$$\Gamma = \frac{1}{m} \sum_{i=1}^{n-1} \sum_{j=i+1}^n X_P(i,j) X_C(i,j)$$

NORMALIZED TAU STATISTICS

$$\hat{f} = \frac{\frac{1}{m} \sum_{i=1}^{n-1} \sum_{j=i+1}^n (X_P(i,j) - \mu_P)(X_C(i,j) - \mu_C)}{\sigma_P \sigma_C}$$

where μ_P and μ_C are the means and σ_P and σ_C are the variances of the matrices X_C and X_P .

MEASURING CLUSTER VALIDITY VIA CORRELATIONS

TWO MATRICES

Proximity Matrix

“Incidence” Matrix

- One row and one column for each data point
- An entry is 1 if the associated pair of points belong to the same cluster
- An entry is 0 if the associated pair of points belongs to different clusters

COMPUTE THE CORRELATION B/W 2 MATRICES

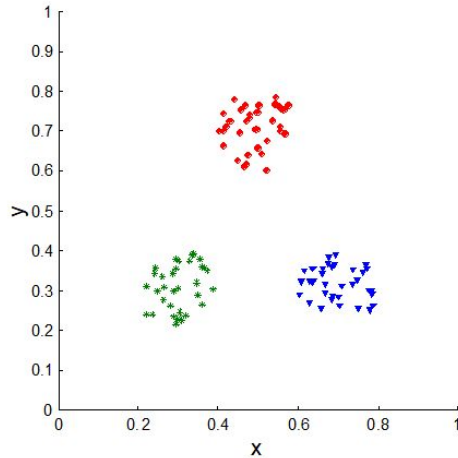
Since the matrices are symmetric, only the correlation between

$n(n-1) / 2$ entries needs to be calculated.

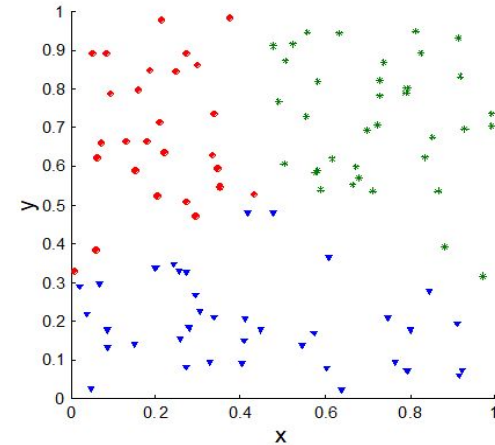
HIGHER CORRELATION MEANS POINTS ARE CLOSE TO EACH OTHER IN
SAME CLUSTER.

NOT A GOOD MEASURE FOR SOME DENSITY OR CONTIGUITY BASED
CLUSTERS.

CORRELATION OF INCIDENCE AND PROXIMITY MATRICES FOR K-MEANS CLUSTERING OF THE FOLLOWING 2 DATA SETS



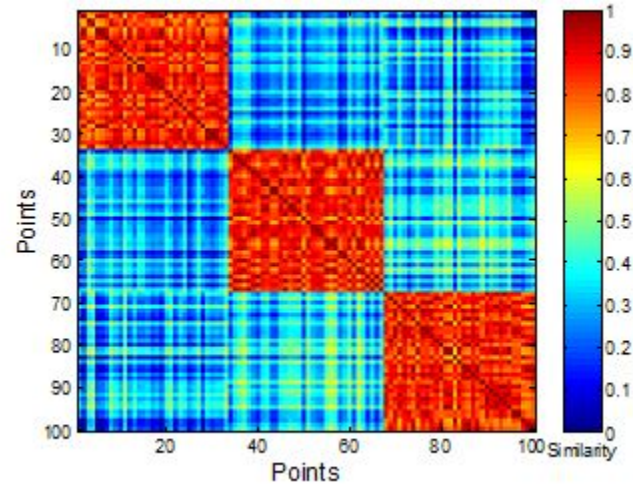
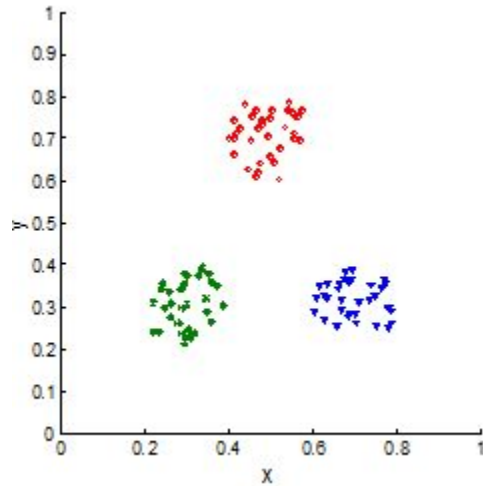
Corr=-0.9235



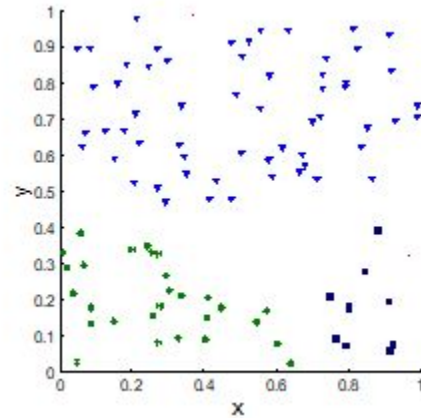
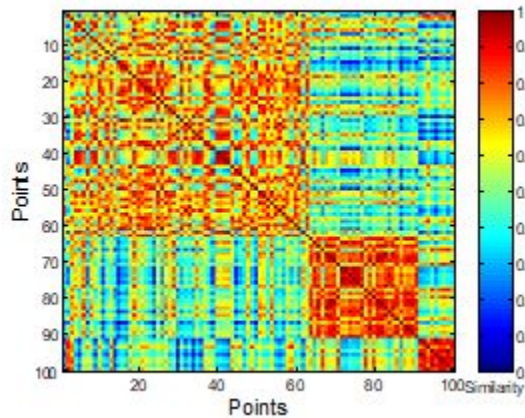
Corr=-0.5810

USING SIMILARITY MATRIX FOR CLUSTER VALIDATION

Order the similarity matrix with respect to cluster labels and inspect visually.

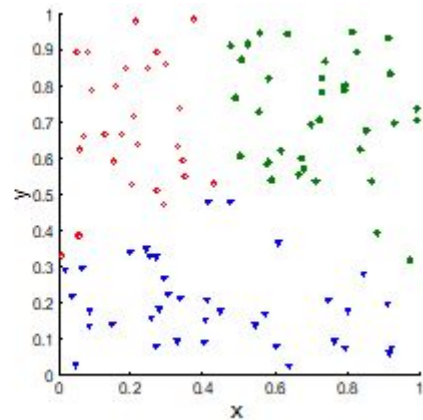
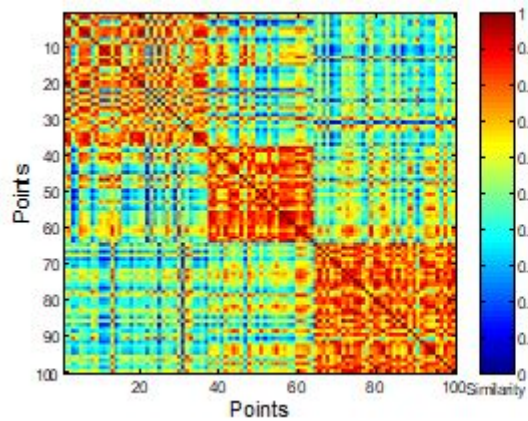


Clusters in random data are not so crisp



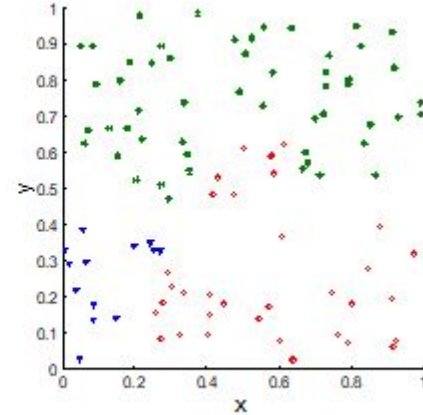
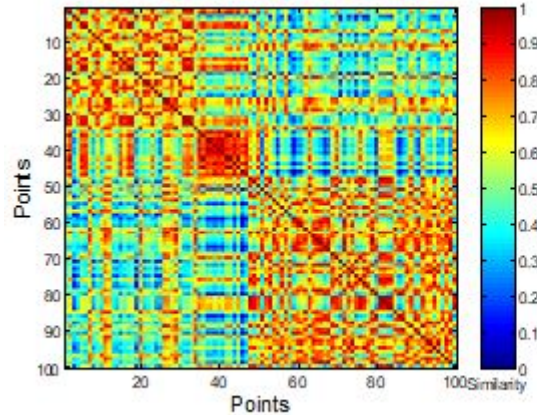
DBSCAN

Clusters in random data are not so crisp

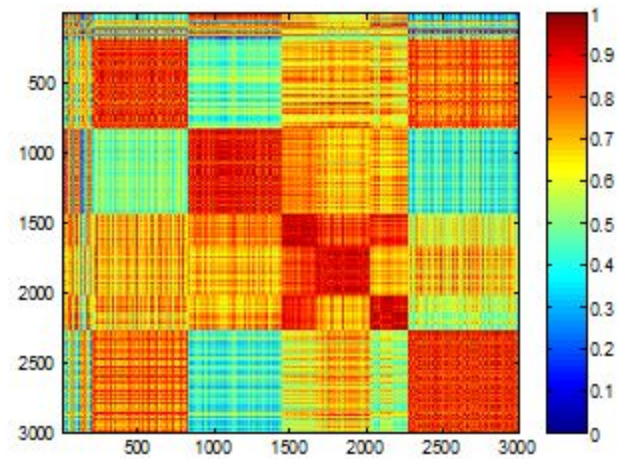
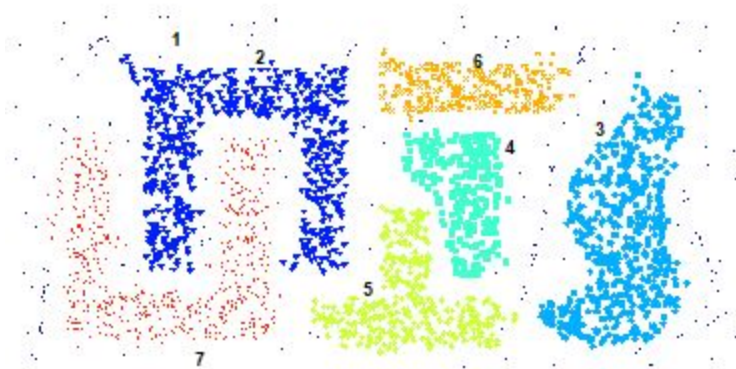


K-MEANS

Clusters in random data are not so crisp



COMPLETE LINK

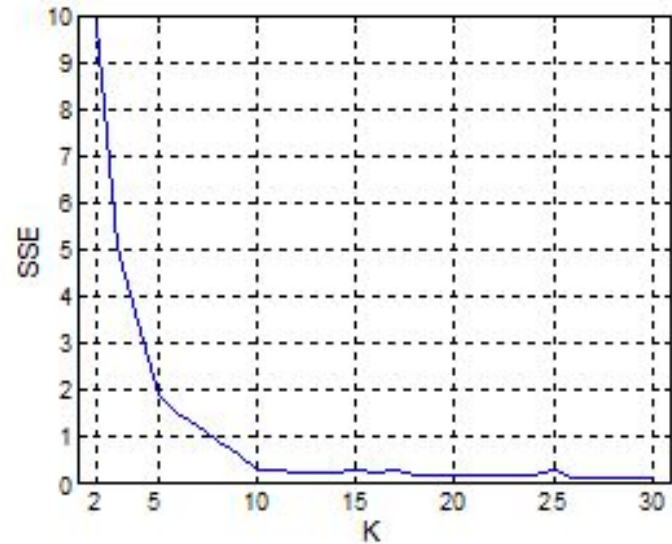
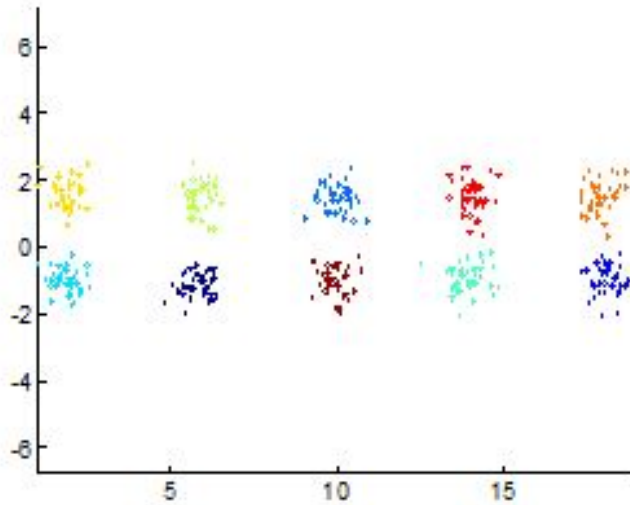


DBSCAN

INTERNAL MEASURES : SSE

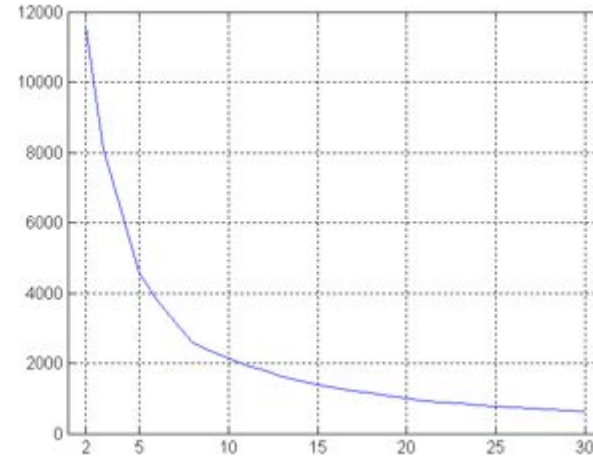
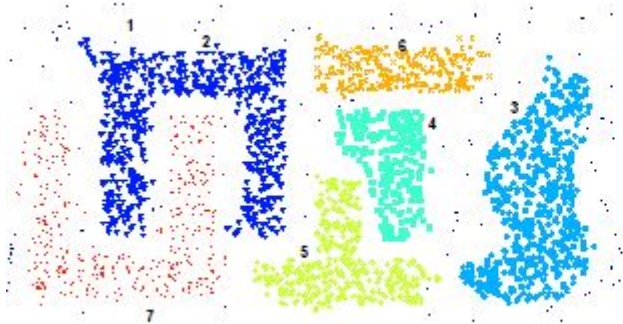
- Clusters in more complicated figures aren't well separated.
- Internal Index: Used to measure the goodness of a clustering structure without respect to external information.
- SSE is good for comparing two clusterings or two clusters (average SSE).
- Can also be used to estimate the number of clusters

Number of clusters can be determined



INTERNAL MEASURES : SSE

SSE curve for a more complicated data set.



SSE Cluster found using K-Means

INTERNAL MEASURE : COHESION & SEPARATION

Cluster Cohesion : Measures

how closely related are objects in
a cluster

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

Cluster Separation : Measure

how distinct or well-separated a
cluster is from other clusters

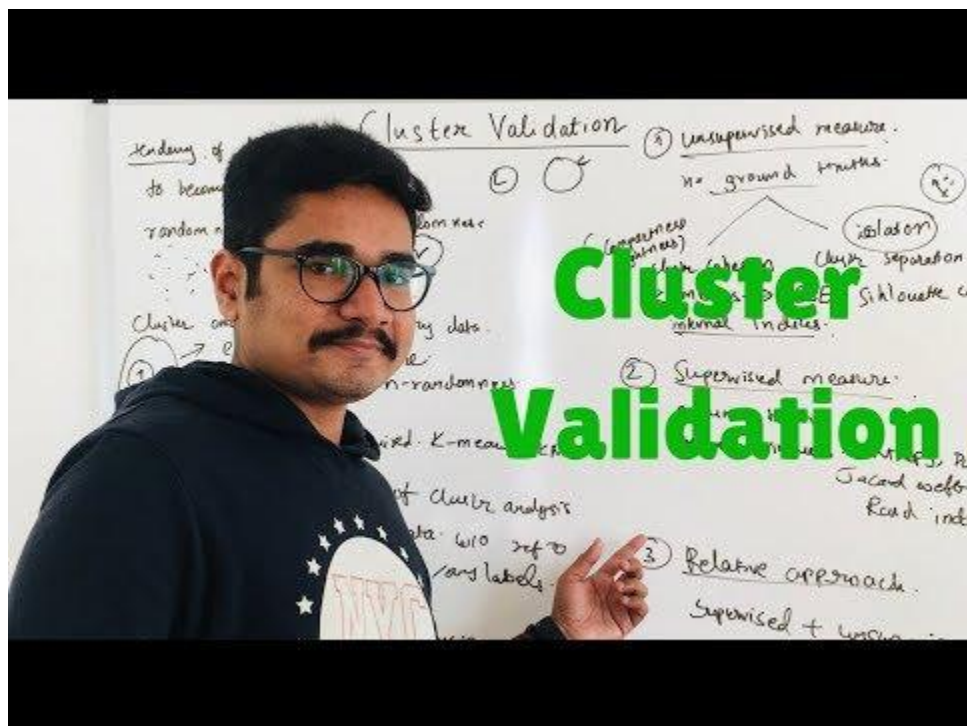
$$BSS = \sum_i |C_i| (m - m_i)^2$$

COMMENT MADE BY JAIN & DUBE IN BOOK (ALGORITHMS OF CLUSTERING DATA)

“The validation of clustering structures is the most difficult and frustrating part of cluster analysis.

Without a strong effort in this direction, cluster analysis will remain a black art accessible only to those true believers who have experience and great courage.”

FOR MORE DETAILS :



THANK YOU