

# Lets First Understand Clustering

## What is Clustering?

**Definition:**

- **Clustering is Unsupervised Machine Learning Technique (Unsupervised – Unlabeled data)** where we are trying grouping of **data points**. These groups are called **Clusters**. **Data points** that are in the **same group** should have **similar properties and/or features**, while **data points in different groups** should have **highly dissimilar properties and/or features**.
- **The aim is to segregate groups with similar traits and assign them into clusters.**

**Simple Use Case:**

Suppose, you are the head of a rental store and wish to understand preferences of your customers to scale up your business. It is definitely for you to look at details of each customer and devise a unique business strategy for each one of them? It is possible. But, what you can do is to cluster all of your customers into say 10 groups based on their purchasing habits and use a separate strategy for customers in each of these 10 groups. And this is what we call clustering.

## Why we use Clustering? or How its helps in Data Science??

- Clustering is used for things like **feature engineering** or **pattern discovery**.
- We can use clustering analysis to gain some valuable insights from our data by seeing what groups the data points falls.
- Statistical Analysis of Data.

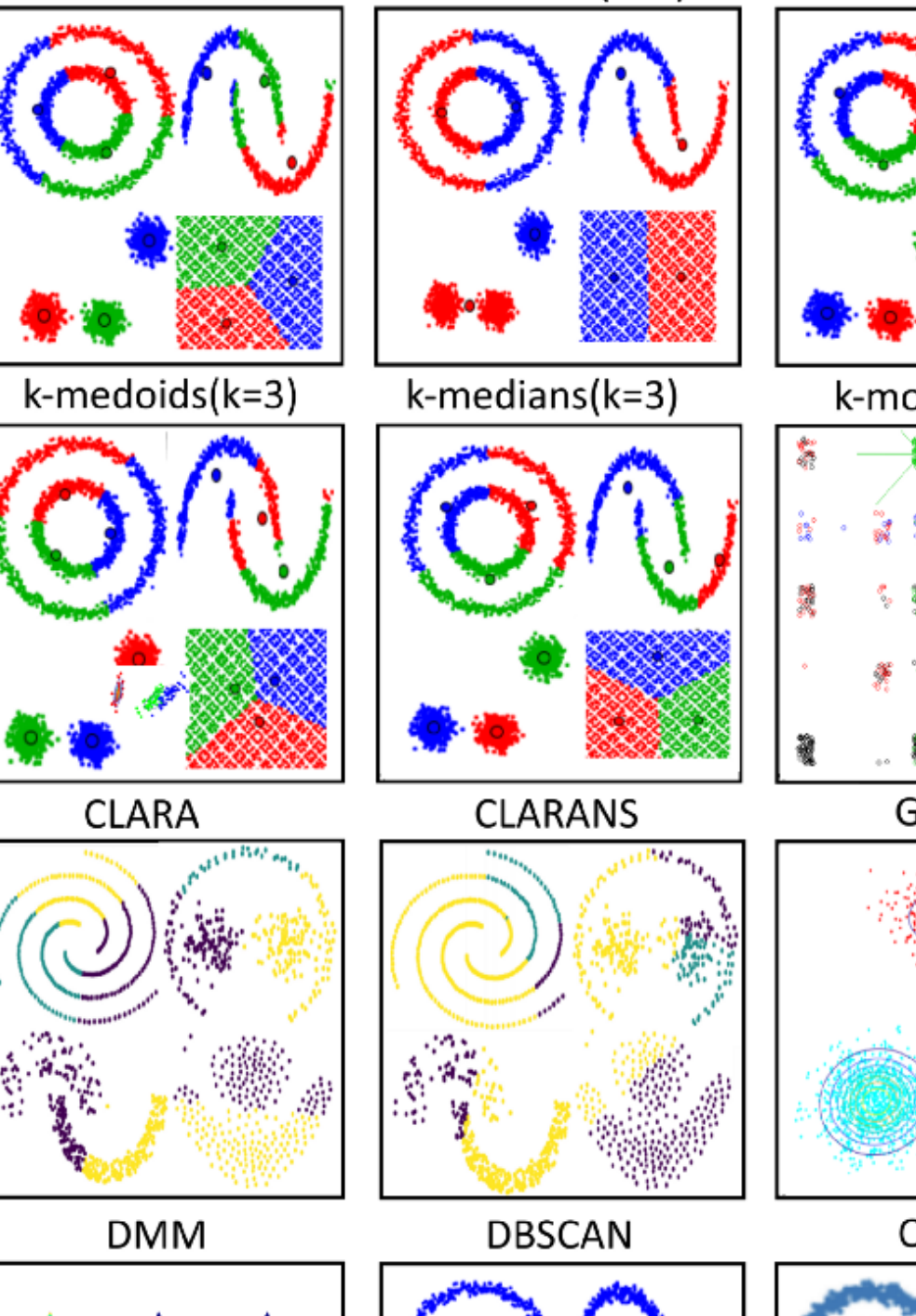
**Real-time Use-cases:**

- One of the most important application is related to image processing, detecting distinct kinds of pattern in image data. This can be very effective in biology research, distinguishing objects and identifying patterns.
- The personal data combined with shopping, location, interest, actions and a infinite number of indicators, can be analysed with this methodology, providing very important information and trends. Examples of this are the market research, marketing strategies, web analytics, and a lot of others.
- Other types of applications based on clustering algorithms are climatology, robotics, recommender systems, mathematical and statistical analysis, providing a broad spectrum of utilization.

**Possible Applications**

- Marketing: Finding groups of customers with similar behavior given a large database of customer data containing their properties and past buying records;
- Biology: Classification of plants and animals given their features;
- Libraries: Book ordering;
- Insurance: Identifying groups of motor insurance policy holders with a high average claim cost; identifying frauds;
- City-planning: Identifying groups of houses according to their house type, value and geographical location;
- Earthquake studies: Clustering observed earthquake epicenters to identify dangerous zones;
- WWW document classification: Clustering weblog data to discover groups of similar access patterns.

## Clustering Visualization



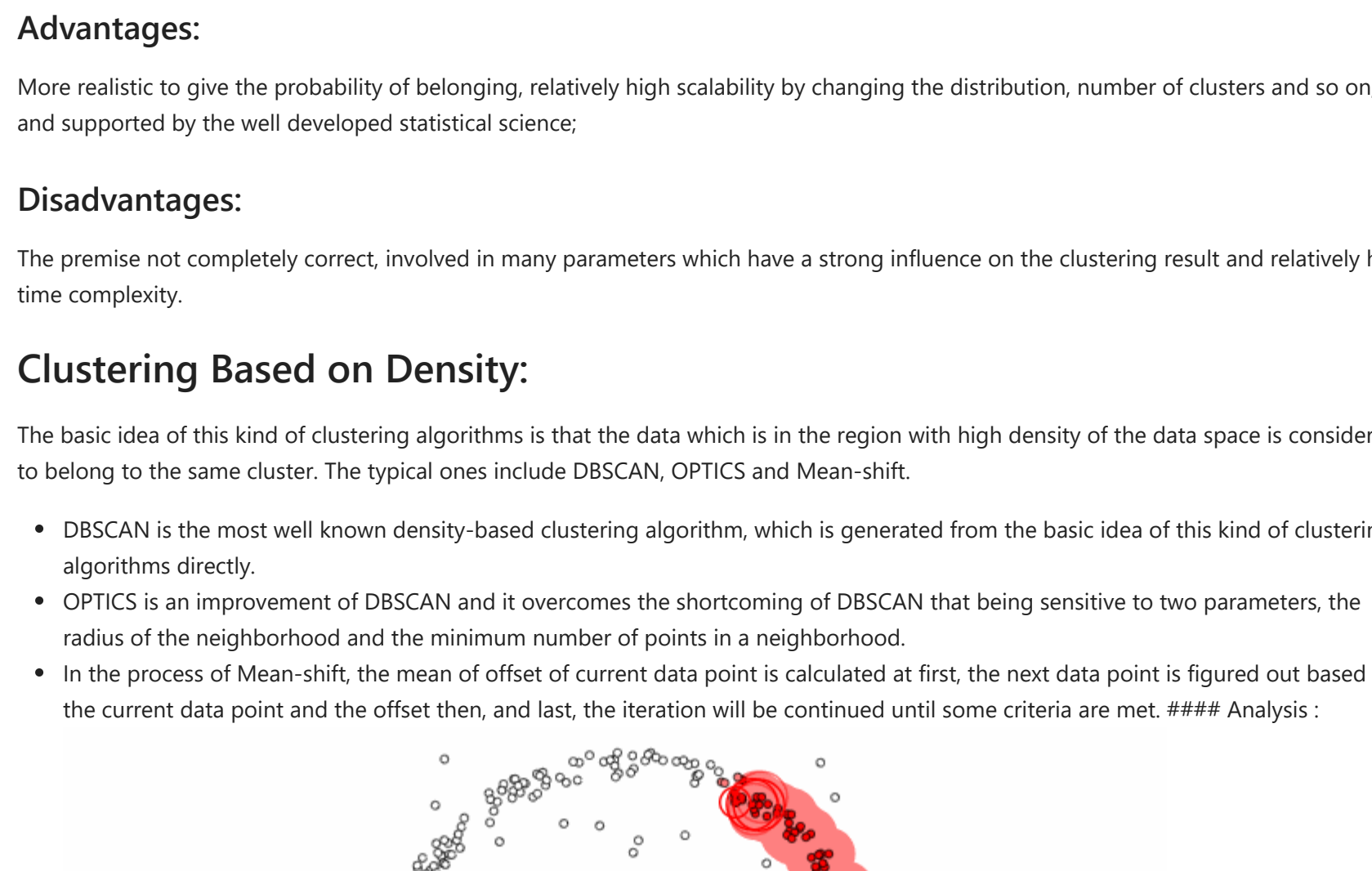
## Types of Clustering

**Based on the area of overlap**

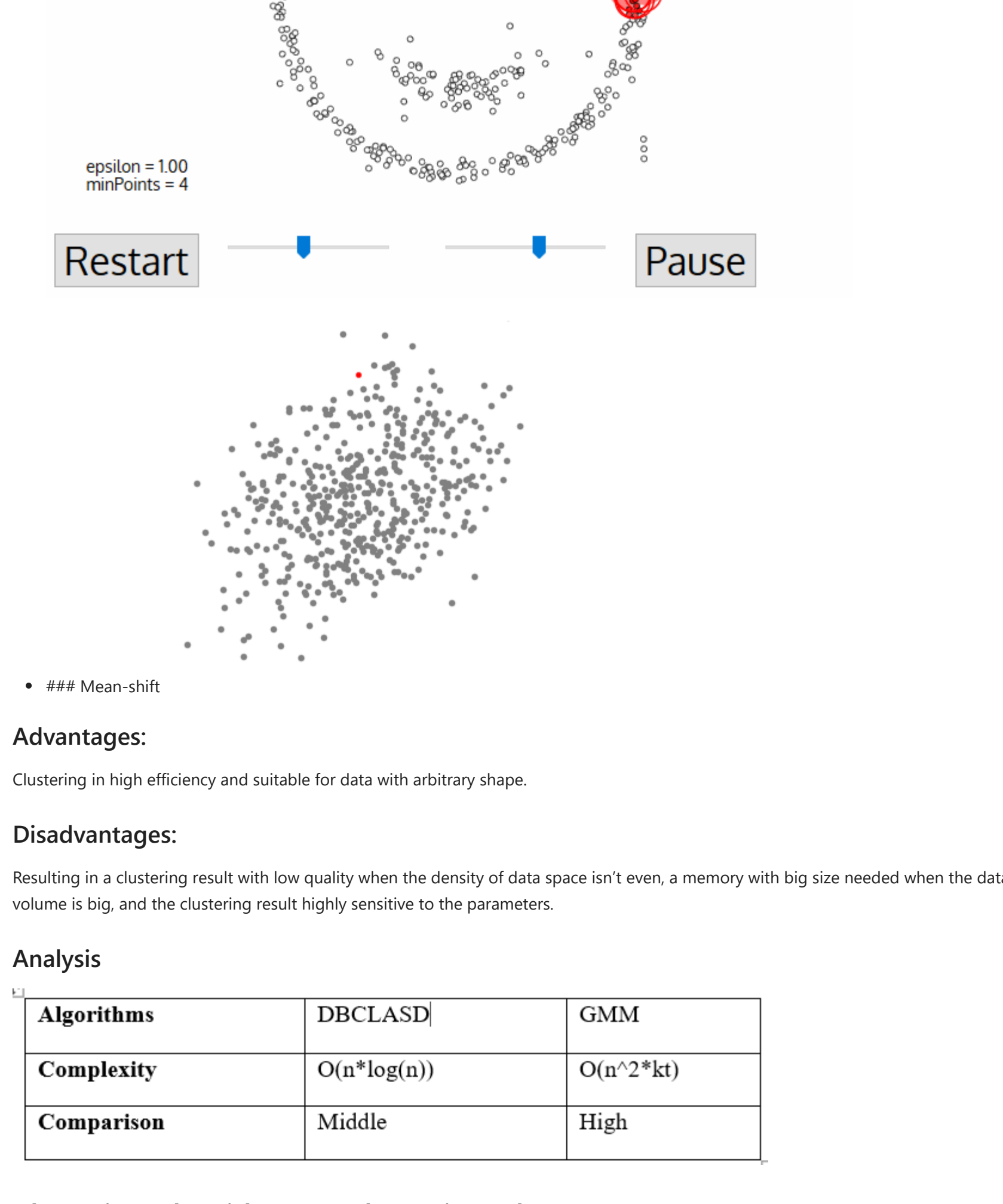
1. **Hard Clustering**
  - Clusters don't overlap=k-means, k-means++ Each data point either belongs to a cluster completely or not.
1. **Soft Clustering**
  - Clusters overlap=Fuzzy c-means, EM. A data object can exist in more than one cluster with a certain probability or degree of membership or
  - Instead of putting each data point into a separate cluster, a probability or likelihood of that data point will be in those clusters is assigned.

## A Comprehensive Survey of Clustering Algorithms.

Link : <https://link.springer.com/article/10.1007/s40745-015-0040-1>



**Various Clustering Algorithms:**



## Lets Discuss different type of method used for clustering

### Clustering Algorithm Based on Distribution

The basic idea is that the data, generated from the same distribution, belongs to the same cluster if there exists several distributions in the original data. The typical algorithms are DBCLASD and GMM. The core idea of DBCLASD, a dynamic incremental algorithm, is that if the distance between a cluster and its nearest data point satisfies the distribution of expected distance which is generated from the existing data points of that cluster, the nearest data point should belong to this cluster. The core idea of GMM is that GMM consists of several Gaussian distributions from which the original data is generated and the data, obeying the same independent Gaussian distribution, is considered to belong to the same cluster.

**Analysis**

Algorithms	DBCLASD	GMM
Complexity	$O(n^2 \log(n))$	$O(n^2 \cdot kt)$
Comparison	Middle	High

**Advantages:**

More realistic to give the probability of belonging, relatively high scalability by changing the distribution, number of clusters and so on, and supported by the well developed statistical science;

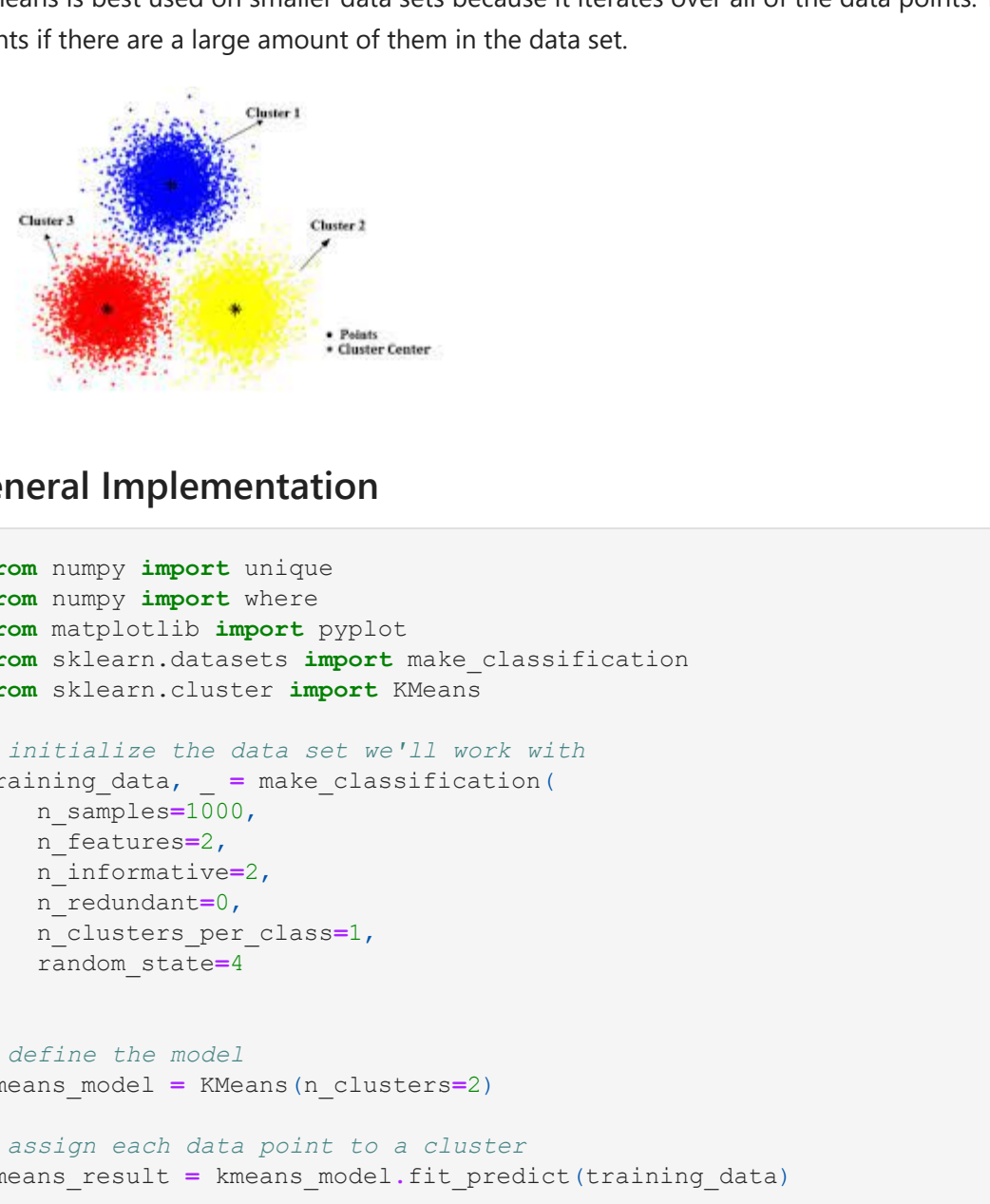
**Disadvantages:**

The premise not completely correct, involved in many parameters which have a strong influence on the clustering result and relatively high time complexity.

### Clustering Based on Density:

The basic idea of this kind of clustering algorithms is that the data which is in the region with high density of the data space is considered to belong to the same cluster. The typical ones include DBSCAN, OPTICS and Mean-shift.

- DBSCAN is the most well known density-based clustering algorithm, which is generated from the basic idea of this kind of clustering algorithms directly.
- OPTICS is an improvement of DBSCAN and it overcomes the shortcoming of DBSCAN that being sensitive to two parameters, the radius of the neighborhood and the minimum number of points in a neighborhood.
- In the process of Mean-shift, the mean of offset of current data point is calculated at first, the next data point is figured out based on the current data point and the offset then, and last, the iteration will be continued until some criteria are met. ##### Analysis :



**Advantages:**

Clustering in high efficiency and suitable for data with arbitrary shape.

**Disadvantages:**

Resulting in a clustering result with low quality when the density of data space isn't even, a memory with big size needed when the data volume is big, and the clustering result highly sensitive to the parameters.

**Analysis**

Algorithms	DBCLASD	GMM
Complexity	$O(n^2 \log(n))$	$O(n^2 \cdot kt)$
Comparison	Middle	High

### Clustering Algorithm Based on Hierarchy

Hierarchical clustering as the name suggests is an algorithm that builds hierarchy of clusters. This algorithm starts with all the data points assigned to a cluster of their own. Then two nearest clusters are merged into the same cluster. In the end, this algorithm terminates when there is only a single cluster left.



**Advantages:**

Suitable for the data set with arbitrary shape and attribute of arbitrary type, the hierarchical relationship among clusters easily detected, and relatively high scalability in general;

**Disadvantages:**

Relatively high in time complexity in general, the number of clusters needed to be preset.

**Analysis**

Algorithms	BIRCH	CURE	ROCK	Chameleon
Complexity	$O(n)$	$O(s^2 \cdot s)$	$O(n^2 \cdot \log(n))$	$O(n^2)$
Comparison	Low	Low	High	High

### Clustering Algorithm Based on Fuzzy Theory

The basic idea of this kind of clustering algorithms is that the discrete value of belonging label,  $\{0, 1\}$ , is changed into the continuous interval  $[0, 1]$ , in order to describe the belonging relationship among objects more reasonably. Typical algorithms of this kind of clustering include FCM, FCS and MM. The core idea of FCM is to get membership of each data point to every cluster by optimizing the object function, FCS, different from the traditional fuzzy clustering algorithms, takes the multidimensional hypersphere as the prototype of each cluster, so as to cluster with the distance function based on the hypersphere. MM, based on the Mountain Function, is used to find the center of cluster

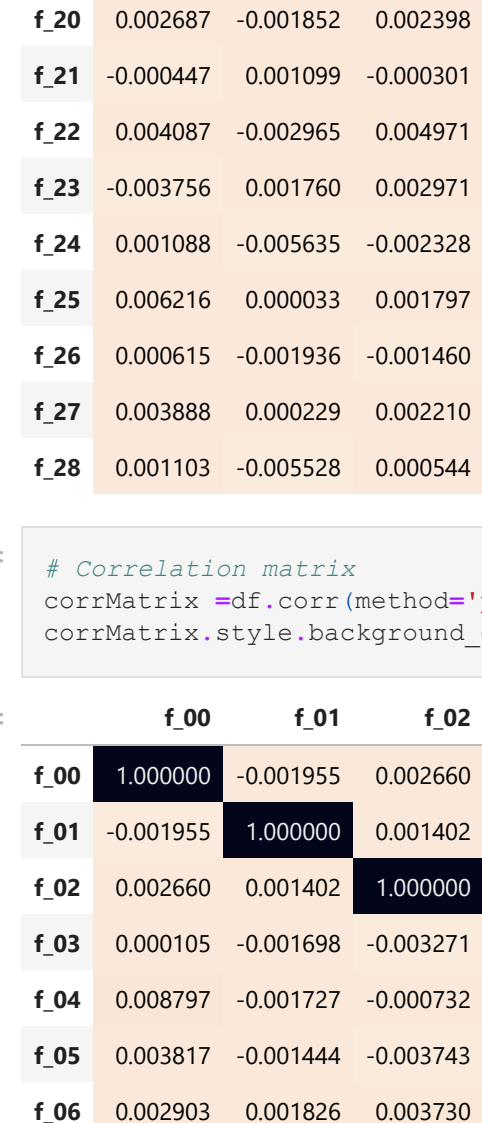
**Advantages:**

More realistic to give the probability of belonging, relatively high accuracy of clustering;

**Disadvantages:**

Relatively low scalability in general, easily drawn into local optimal, the clustering result sensitive to the initial parameter values, and the number of clusters needed to be preset.

Algorithms	FCM	FCS	MM
Complexity	$O(n)$	Kernel	$O(v^{1/2} \cdot n)$
Comparison	Low	High	Middle



## Now Lets Jumps to Implementation Section

### Famous Clustering Algorithms:

1. **K-means clustering algorithm**
2. **DBSCAN clustering algorithm**
3. **Gaussian Mixture Model algorithm**
4. **BIRCH algorithm**
5. **Affinity Propagation clustering algorithm**
6. **Mean-Shift clustering algorithm**
7. **OPTICS algorithm**
8. **Agglomerative Hierarchy clustering algorithm**

### K-means clustering algorithm

K-means clustering is the most commonly used clustering algorithm. It's a centroid-based algorithm and the simplest unsupervised learning algorithm.

This algorithm tries to minimize the variance of data points within a cluster. It's also how most people are introduced to unsupervised machine learning.

K-means is best used on smaller data sets because it iterates over all of the data points. That means it'll take more time to classify data points if there are a large amount of them in the data set.



