

Clustering

Clustering

- Process of grouping a set of examples
- Clustering generates a **partition** consisting of **cohesive groups or clusters** from given collection of examples



- The examples to be clustered are either **labelled** or **unlabelled**
 - Algorithms which cluster labelled examples:
 - Supervised clustering
 - Classification
 - Algorithms which cluster unlabelled examples:
 - Unsupervised clustering
 - Do not rely on predefined classes
 - **Learning by observation**, rather than learning by examples.

Clustering

- Clustering is a two step process
 - Step1: Partition the collection of examples (clustering)
 - Learning by observation (training phase)
 - Group the collection of examples into finite number of clusters such that the examples that are similar to one another within the same cluster and are dissimilar to examples in other clusters
 - Obtaining cluster labels
 - Unsupervised learning: Do not rely on predefined classes and class-labelled training examples
 - Step2: Assign cluster labels to examples
 - Testing phase

3

Categorization of Clustering Methods

- Partitioning methods:
 - These methods construct K partitions of the data, where each partition represents a cluster
 - Idea: Cluster the collection of examples based on the distance between examples
 - Results in spherical shaped cluster
 - 1. K -means algorithm
 - 2. K -medoids algorithm
 - 3. Gaussian mixture model
- Hierarchical methods:
 - These methods create a hierarchical decomposition of the collection of examples
 - Results in spherical shaped cluster
 - 1. Agglomerative approach (bottom-up approach)
 - 2. Divisive approach (top-down approach)

4

Categorization of Clustering Methods

- Density-based methods:
 - These methods cluster collection of examples based on the notion of density
 - General idea: To continue growing the given cluster as long as density (number of examples) in the neighbourhood exceeds some threshold
 - Example:
 - DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

5

Classical Portioning Methods

- Centroid-based technique:
 - Partition the collection of examples into K clusters based on the distance between examples
 - Cluster similarity is measured in regard to the sample mean of the examples within a cluster
 - Cluster centroid or center of gravity: Sample mean value of the examples within a cluster
 - Cluster center is used to represent the cluster
 - Example: K -means algorithm
- Representative object-based technique:
 - Actual example is considered to represent the cluster
 - One representative example per cluster
 - Example: K -medoids algorithm

6

***K*-Means Clustering Algorithm**

- Dividing the data into K groups or partitions
- **Given:** Training data, $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^d$ and K
- **Target:** Partition the set \mathcal{D} into K clusters (disjoint subsets), $\{\mathcal{D}_k\}_{k=1}^K$
 - Each of the clusters is associated with centers, $\boldsymbol{\mu}_k$, $k=1, 2, \dots, K$
 - Come up with the centers of clusters
 - Cluster center acts as a cluster representative
- Euclidean distance with center of a cluster can be used as a measure of dissimilarity

7

***K*-Means Clustering Algorithm**

1. Initialize the cluster center, $\boldsymbol{\mu}_k$, $k=1, 2, \dots, K$ using randomly selected K data points in \mathcal{D}
2. Assign each data point \mathbf{x}_n to cluster center k^*

$$k^* = \arg \min_k \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$
3. **Update** $\boldsymbol{\mu}_k$, $k=1, 2, \dots, K$: Re-compute $\boldsymbol{\mu}_k$ after assigning all the data points.

$$\hat{\boldsymbol{\mu}}_k = \frac{\sum_{\mathcal{D}_k} \mathbf{x}_n}{N_k}$$

N_k : Number of examples in cluster k
4. Repeat the steps 2 and 3 until the convergence

8

***K*-Means Clustering Algorithm**

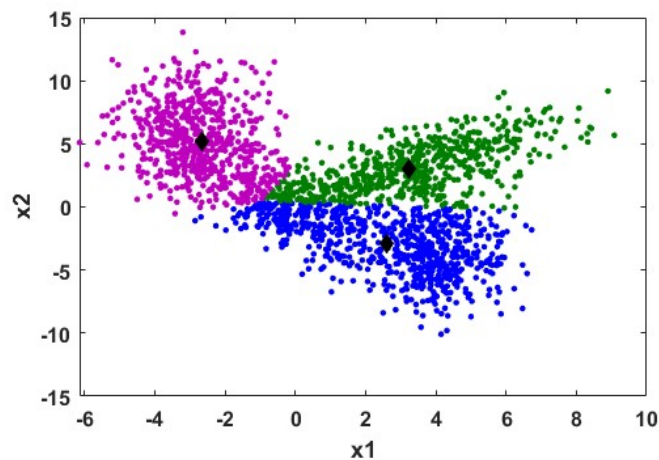
- **Convergence criteria:**
 - No change in the cluster assignment **OR**
 - The difference between the **distortion measure (J)** in the successive iteration falls below the threshold
 - **Distortion measure (J)** : Sum of the squares of the distance of each example to its assigned cluster center

$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

z_{nk} is 1 if \mathbf{x}_n belongs to cluster k , otherwise 0

9

Illustration of *K*-Means Clustering



- Boundary between the cluster is linear
- **Hard clustering:** Each example must belong to exactly one group

10

K-Means Clustering Algorithm

- Dividing the data into K groups or partitions
- **Given:** Training data, $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^d$ and K
- **Target:** Partition the set \mathcal{D} into K clusters (disjoint subsets), $\{\mathcal{D}_k\}_{k=1}^K$
 - Each of the clusters is associated with centers, $\boldsymbol{\mu}_k$, $k=1, 2, \dots, K$
 - **Better representative for a cluster**
 - Come up with the **centers of clusters** and **covariance matrix**
 - Cluster center and covariance matrix act as **cluster representatives**
- **Mahalanobis distance** with cluster representatives can be used as a measure of dissimilarity

11

K-Means Clustering Algorithm

1. Initialize the cluster center, $\boldsymbol{\mu}_k$, $k=1, 2, \dots, K$ using randomly selected K data points in \mathcal{D}
2. Initialize the covariance matrix, $\boldsymbol{\Sigma}_k$, $k=1, 2, \dots, K$ using unit matrix
3. Assign each data point \mathbf{x}_n to cluster center k^*

$$k^* = \arg \min_k (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

4. **Update $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $k=1, 2, \dots, K$:** Re-compute $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ after assigning all the data points.

$$\hat{\boldsymbol{\mu}}_k = \frac{\sum_{\mathcal{D}_k} \mathbf{x}_n}{N_k} \quad \hat{\boldsymbol{\Sigma}}_k = \frac{\sum_{\mathcal{D}_k} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_k)^T}{N_k} \quad \begin{array}{l} N_k: \text{Number of} \\ \text{examples in cluster } k \end{array}$$

5. Repeat the steps 3 and 4 until the convergence

12

***K*-Means Clustering Algorithm**

- **Convergence criteria:**
 - No change in the cluster assignment **OR**
 - The difference between the **distortion measure** (J) in the successive iteration falls below the threshold

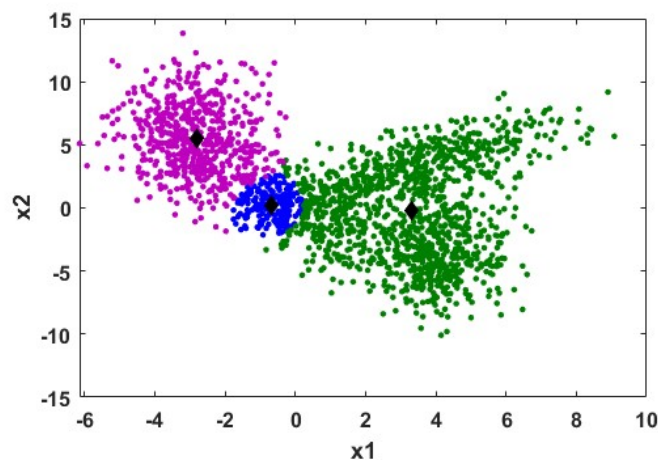
$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left[(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right]$$

z_{nk} is 1 if \mathbf{x}_n belongs to cluster k , otherwise 0

- **Hard clustering:** Each example must belong to exactly one group

13

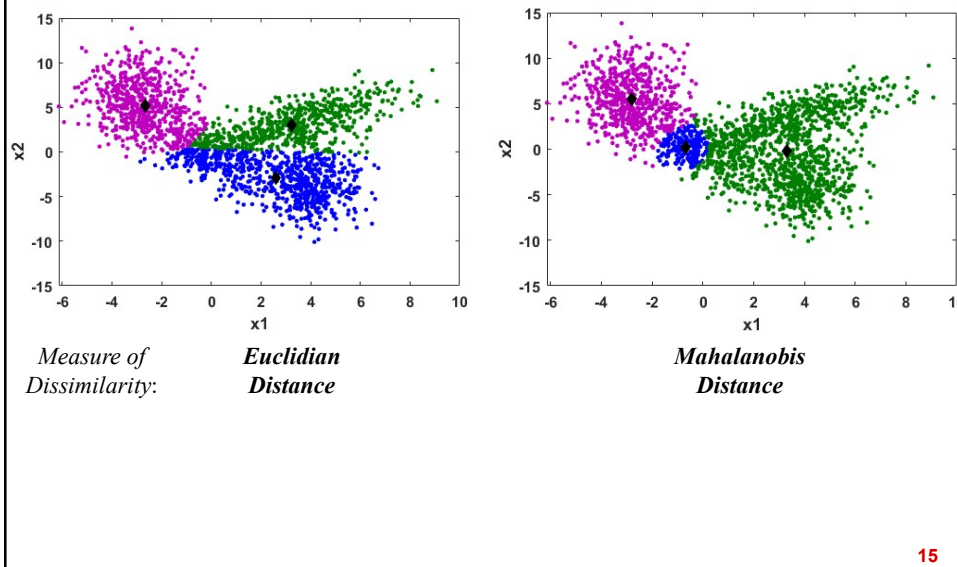
Illustration of *K*-Means Clustering



- Boundary between the cluster is quadratic
- **Hard clustering:** Each example must belong to exactly one group

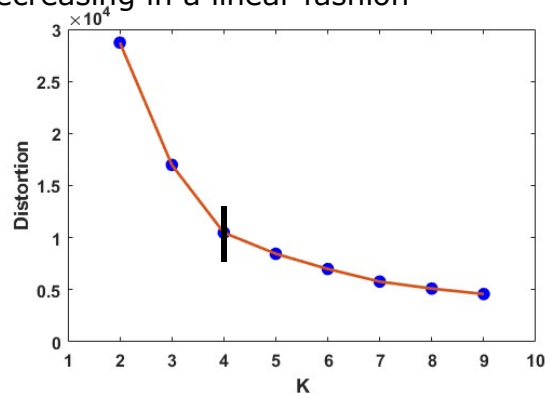
14

Illustration of K -Means Clustering

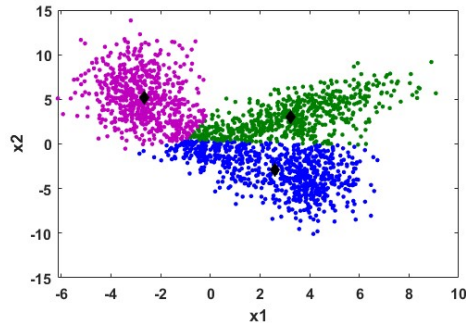


Elbow Method to Choose K

- Determine the **distortion measure** for different values of K
- Plot the **K vs Distortion**
- **Optimal number of clusters**: Select the value of K at the "elbow" i.e. the point after which the distortion start decreasing in a linear fashion



Soft Clustering



- **Soft clustering**: Each example belong to each group with some probability
 - Fuzzyness at the boundary of the clusters
- Gaussian mixture model (GMM) is one of the soft clustering techniques
- GMM can be seen as similar to K-means clustering
- Each cluster is represented as Gaussian density

17

Gaussian Mixture Model (GMM)

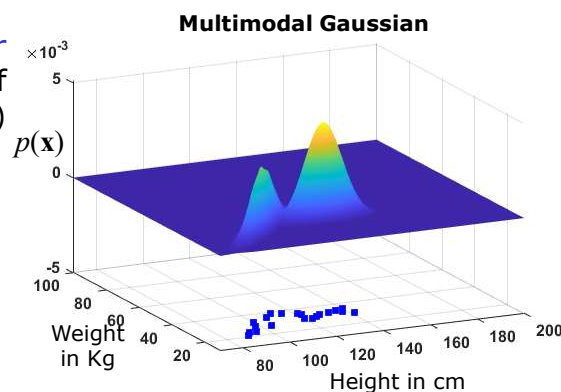
- **Given**: Training data having N samples

$$\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \dots, \mathbf{x}_N\}, \quad \mathbf{x}_n \in \mathbb{R}^d$$

- GMM is a **linear superposition** of multiple **Gaussian components**:

$$p(\mathbf{x}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right)$$



18

Gaussian Mixture Model (GMM)

- GMM is a linear superposition of multiple Gaussians:

$$p(\mathbf{x}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- For a d -dimensional feature vector representation of data, the parameters of GMM are
 - Mixture coefficients, w_k , $k = 1, 2, \dots, K$
 - Mixture weight or Strength of each clusters (or mixtures or modes)
 - Property: $\sum_{k=1}^K w_k = 1$
 - d -dimensional mean vector, $\boldsymbol{\mu}_k$, $k = 1, 2, \dots, K$
 - $d \times d$ size covariance matrices, $\boldsymbol{\Sigma}_k$, $k = 1, 2, \dots, K$
- Training process objective:
 - Partition the data into K groups
 - To estimate the parameters of the each cluster in GMM

19

Expectation-Maximization (EM) for GMMs

- Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters
 - Initialize the mean vectors $\boldsymbol{\mu}_k$, covariance matrices $\boldsymbol{\Sigma}_k$ and mixing coefficients w_k , and evaluate the initial value of the log likelihood
 - E-step**: Evaluate the responsibilities $\gamma_k(\mathbf{x})$ using the current parameter values – Assign the data points to each cluster

20

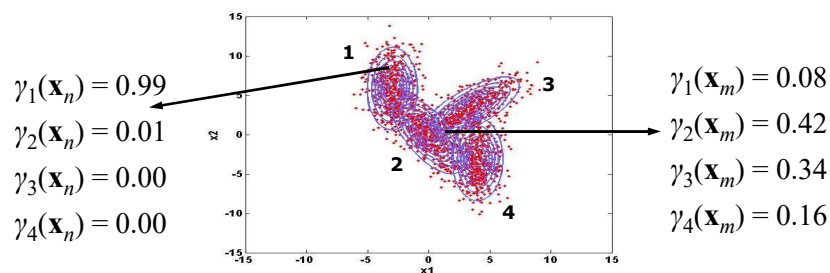
EM Method – Responsibility Term

- A quantity that plays an important role is the **responsibility term**, $\gamma_k(\mathbf{x})$

- It is given by

$$\gamma_k(\mathbf{x}) = \frac{w_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

- w_k : **mixture coefficient** or **prior probability of cluster k** ,
- $\gamma_k(\mathbf{x})$ gives the **posterior probability of the cluster k** for the observation \mathbf{x}



21

Expectation-Maximization (EM) for GMMs

- Given a Gaussian mixture model, the goal is to **maximize the likelihood function with respect to the parameters**

1. Initialize the **mean vectors $\boldsymbol{\mu}_k$** , **covariance matrices $\boldsymbol{\Sigma}_k$** and **mixing coefficients w_k** , and evaluate the initial value of the log likelihood

2. **E-step**: Evaluate the responsibilities $\gamma_k(\mathbf{x})$ using the current parameter values

3. **M-step**: Re-estimate the parameters $\boldsymbol{\mu}_k^{new}$, $\boldsymbol{\Sigma}_k^{new}$ and w_k^{new} using the current responsibilities

$$\boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(\mathbf{x}_n) \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$w_k^{new} = \frac{N_k}{N}$$

$$N_k = \sum_{n=1}^N \gamma_k(\mathbf{x}_n)$$

- N_k : Effective number of points assigned to the cluster k

22

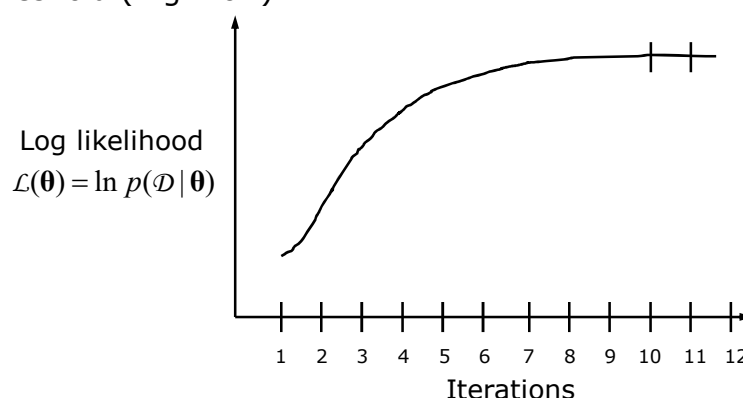
Expectation-Maximization (EM) for GMMs

- Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters
 1. Initialize the mean vectors μ_k , covariance matrices Σ_k and mixing coefficients w_k , and evaluate the initial value of the log likelihood
 2. **E-step**: Evaluate the responsibilities $\gamma_k(\mathbf{x})$ using the current parameter values
 3. **M-step**: Re-estimate the parameters μ_k^{new} , Σ_k^{new} and w_k^{new} using the current responsibilities
 4. Evaluate the log likelihood and check for convergence of the log likelihood
 - If the convergence criterion is not satisfied return to step 2

23

Expectation-Maximization (EM) for GMMs

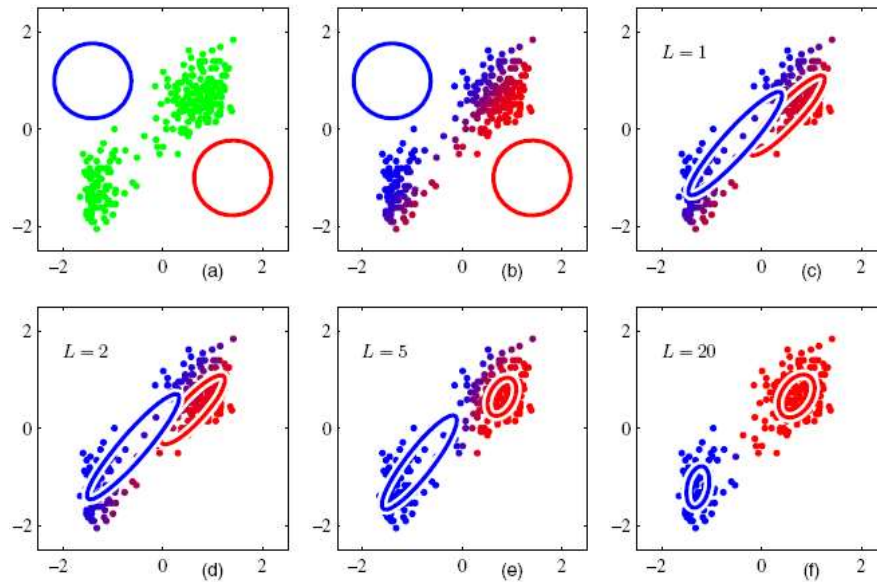
- Convergence criterion**: Difference between log likelihoods of successive iterations fall below a threshold (E.g. 10^{-3})



$$\mathcal{L}(\theta) = \ln p(\mathcal{D} | \theta) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) \right)$$

24

Illustration of Parameter Estimation

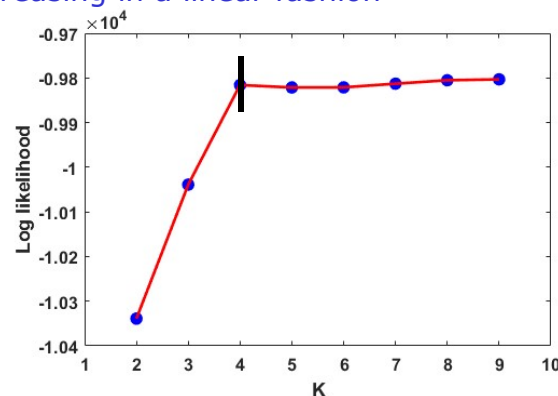


C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.

25

Elbow Method to Choose K

- Determine the **total data log likelihood** for different values of K
- Plot the K vs **log likelihood**
- **Optimal number of clusters**: Select the value of K at the "elbow" i.e. the point after which the log likelihood start **increasing in a linear fashion**



26

***K*-Medoid Clustering Algorithms**

- Related to *K*-means clustering
- The *K*-means algorithm is sensitive to outliers because an example with extremely large value may substantially distort the distribution of data
- Solution: One of the data points is chosen as representative of cluster, instead of mean value of the cluster
- It replaces the means of cluster with modes
- Partitioning around medoids
- A medoid of a finite dataset: The data point from the set, whose average dissimilarity (distance) to all the points is minimal
 - The most centrally located point in the set

27

***K*-Medoid Clustering Algorithm**

- Given: Training data, $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^d$ and K
1. Initialize the medoid, $\hat{\mathbf{x}}_k$, $k=1, 2, \dots, K$ using randomly selected K data points in \mathcal{D}
 2. Assign each data point \mathbf{x}_n to the closest medoid

$$k^* = \arg \min_k \|\mathbf{x}_n - \hat{\mathbf{x}}_k\|^2$$
 3. Update medoids $\hat{\mathbf{x}}_k$, $k=1, 2, \dots, K$
 - For each data point \mathbf{x}_n assigned to a cluster k compute the average dissimilarity (distance) of \mathbf{x}_n to all the data points assigned to cluster k

$$\text{Average dissimilarity for } \mathbf{x}_n = \frac{\sum_{\mathbf{x}_m \in \mathcal{D}_k} \|\mathbf{x}_n - \mathbf{x}_m\|^2}{N_k}$$

N_k : Number of examples in cluster k
 - Select the example with minimum average dissimilarity as medoid
 4. Repeat the steps 2 and 3 until the convergence

28

K-Medoid Clustering Algorithm

- Convergence criteria:
 - No change in the cluster assignment **OR**
 - The difference between the distortion measure (absolute-error) (J) in the successive iteration falls below the threshold
 - Distortion measure (J) : Sum of the squares of the distance of each example to its corresponding reference point (medoid)

$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \hat{\mathbf{x}}_k\|^2$$

z_{nk} is 1 if \mathbf{x}_n belongs to cluster k , otherwise 0

- Optimal number of clusters (k) can obtained using elbow method

29

Evaluation of Clustering: Purity Score

- Lets us assume that class index for each example is given
- Purity score: Purity is a measure of the extent to which clusters contain a single class
 - For each cluster, count the number of data points from the most common class in said cluster
 - Take the sum over all clusters and divide by the total number of data points
- Let M be the number of classes, $C_1, C_2, \dots, C_m, \dots, C_M$
- Let K be the number of clusters, $k = 1, 2, \dots, K$
- Let N be the number of data points

30

Evaluation of Clustering: Purity Score

- For each cluster k ,
 - Count the number of data points from each class
 - Consider the number of data points of most common class

$$\max_m |N_k \cap C_m|$$

$|N_k \cap C_m|$ is the number of data points in k^{th} cluster belonging to class m

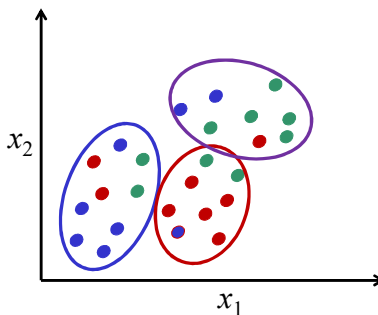
- Take the sum over all clusters, k
- Divide by the total number of data points (N)

$$\text{Purity Score} = \frac{1}{N} \sum_{k=1}^K \max_m |N_k \cap C_m|$$

31

Illustration of Computing Purity Score

- Number of data points, $N = 25$
- number of classes, $M = 3$
- number of clusters, $K = 3$
- *Cluster 1*: Number of examples of **Blue Class** are more, i.e. **5**
- *Cluster 2*: Number of examples of **Red Class** are more, i.e. **5**
- *Cluster 3*: Number of examples of **Green Class** are more, i.e. **5**
- **Purity score: $(5+5+5)/25 = 0.60$**



32

Hierarchical Clustering

Hierarchical Clustering Algorithms

- These methods create a **hierarchical decomposition** of the collection of examples
- Produce nested sequence of data partitions
- These sequence can be depicted using a tree structure
- Hierarchical clustering method works by grouping data points into a **tree of clusters**
- Hierarchical algorithms are either **agglomerative** or **divisive**
 - This classification of hierarchical clustering is depending on whether the hierarchical decomposition is formed in a
 - **Bottom-up (merging)** OR
 - **Top-down (splitting)** fashion

Agglomerative Hierarchical Clustering

- Bottom-up approach
- This strategy starts by placing each example in its own cluster (atomic clusters) and then merges these atomic clusters into larger and larger clusters
- Starts with N clusters where each example is a cluster
- At each successive step (level), the most similar pair of clusters are merged
 - The measure of closeness (intercluster similarity) is considered to decide which two clusters are merged
 - At each level, number of clusters is reduces by one
- The process continues till all the examples are in a single cluster or until certain termination conditions are satisfied
 - Termination condition could be
 - Number of clusters
 - Intercluster similarity between each pair of cluster is within a certain threshold

35

Agglomerative Hierarchical Clustering

- ***Once two examples are placed in the same cluster at a level, they remain in same cluster at all subsequent levels***
- Example: AGglomerative NESTing (AGNES)
- Most hierarchical clustering methods belong to tis category
- They differ only in their definition of intercluster similarity
- Intercluster similarity is to identifying two closest cluster for merging
- When there is one example in a cluster, two closest clusters are found by computing minimum Euclidian distance between two clusters
- However, there is no unique way to find the two closest clusters when there are more than one data points in each clusters

36

Agglomerative Hierarchical Clustering

- Different **intercluster similarities** to find similarity between the clusters having more than one examples:
 - **Minimum distance** between any two examples from two clusters C_i and C_j

$$d_{\min}(C_i, C_j) = \min_{\mathbf{x} \in C_i, \mathbf{x}' \in C_j} \|\mathbf{x} - \mathbf{x}'\|$$

- **Distance between the centers** of two clusters C_i and C_j

$$d_{\text{mean}}(C_i, C_j) = \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|$$

- Where $\boldsymbol{\mu}_i$ is the center of C_i and $\boldsymbol{\mu}_j$ is the center of C_j

- **Average distance** of a point in one cluster (C_i) to a point in another cluster (C_j)

$$d_{\text{avg}}(C_i, C_j) = \frac{1}{N_i N_j} \sum_{\mathbf{x} \in C_i} \sum_{\mathbf{x}' \in C_j} \|\mathbf{x} - \mathbf{x}'\|$$

- Where N_i and N_j are the number of examples in clusters C_i and C_j respectively

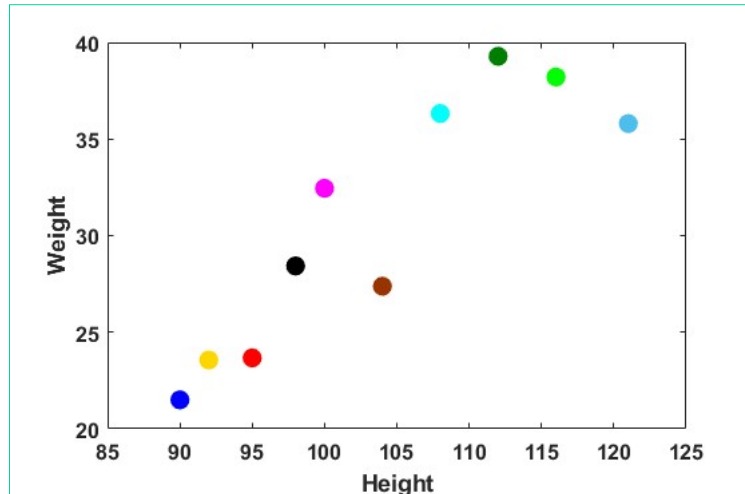
37

Agglomerative Hierarchical Clustering

- **Given**: Training data, $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N, \mathbf{x}_n \in \mathbb{R}^d$
- **Target**: Partition the data
- **Step1**: N clusters where each example is a cluster
- **Step2**: Compute **intercluster similarity** between each pair of clusters
- **Step3**: Choose a pair of clusters that are **most similar** (**minimum intercluster distance**) and merge them
- **Step4**: Repeat **Step2** and **Step3** until **all the examples are in a single cluster** or until **certain termination conditions are satisfied**

38

Agglomerative Hierarchical Clustering

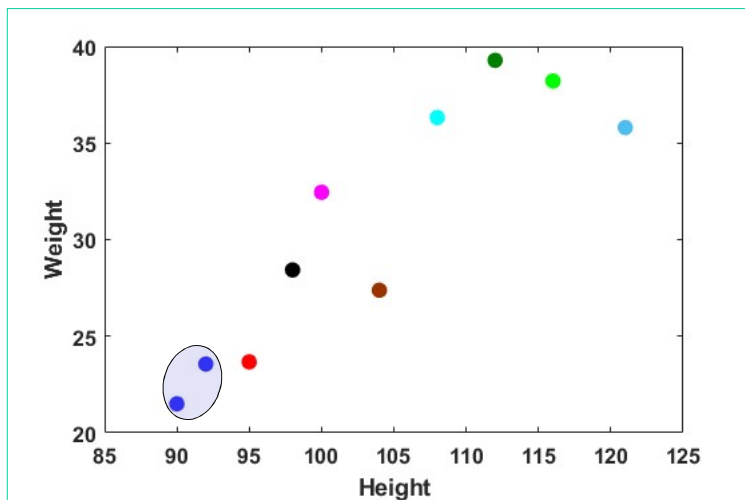


Level: 0

Number of clusters: 10

39

Agglomerative Hierarchical Clustering

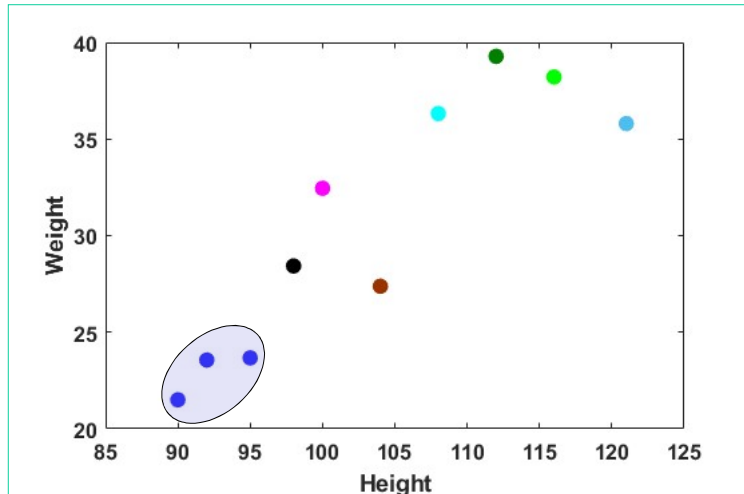


Level: 1

Number of clusters: 9

40

Agglomerative Hierarchical Clustering

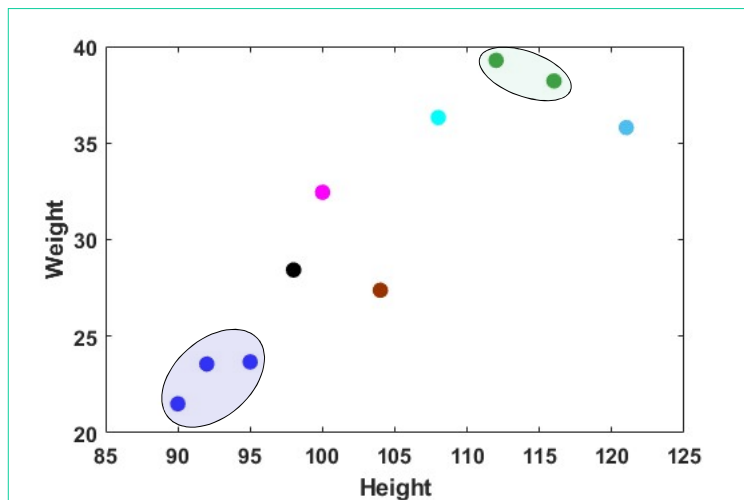


Level: 2

Number of clusters: 8

41

Agglomerative Hierarchical Clustering

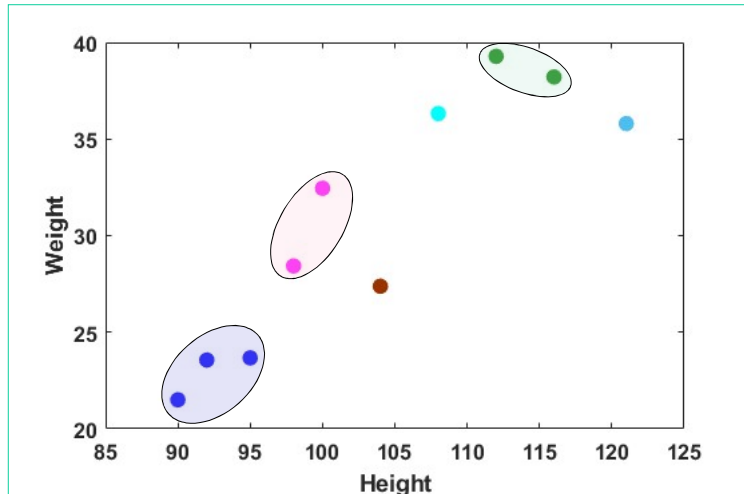


Level: 3

Number of clusters: 7

42

Agglomerative Hierarchical Clustering

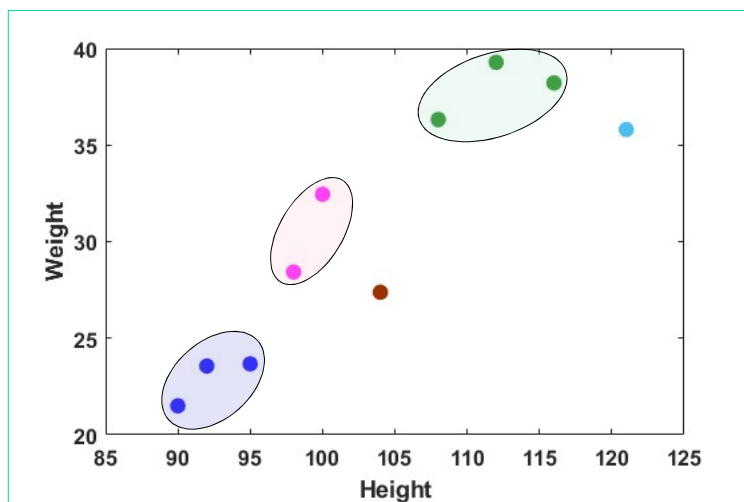


Level: 4

Number of clusters: 6

43

Agglomerative Hierarchical Clustering

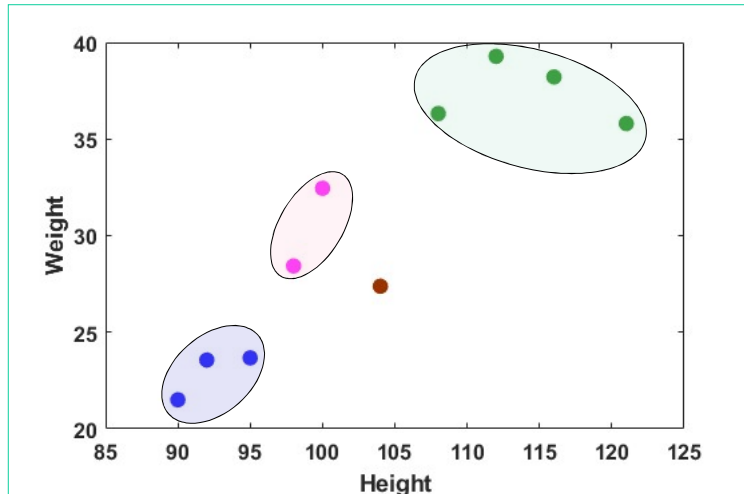


Level: 5

Number of clusters: 5

44

Agglomerative Hierarchical Clustering

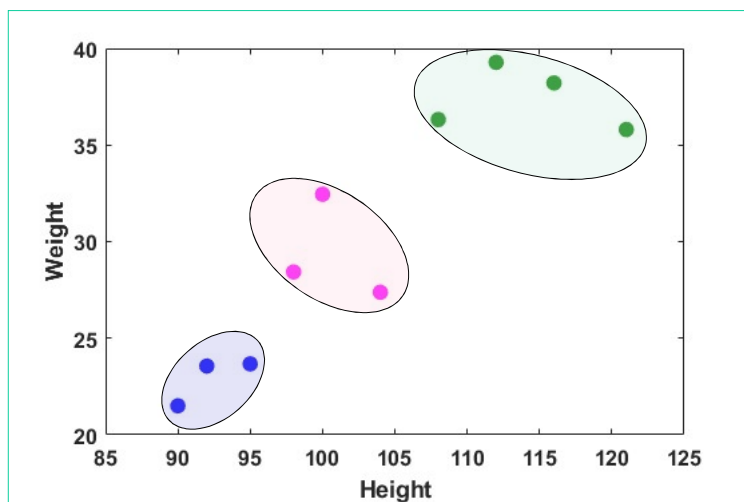


Level: 6

Number of clusters: 4

45

Agglomerative Hierarchical Clustering



Level: 7

Number of clusters: 3

46

Divisive Hierarchical Clustering

- Top-down approach
- Starts with single cluster having all the examples
- It subdivides the cluster into smaller and smaller clusters in the successive step
- At each successive step, a compactness measure is used to choose which cluster to split
 - Compactness measure: Average value of distance between the data points of a cluster
 - Compactness measure (CM_i) of a cluster C_i :

$$CM_i = \frac{1}{N_i^2} \sum_{\mathbf{x} \in C_i} \sum_{\mathbf{x}' \in C_i} \|\mathbf{x} - \mathbf{x}'\|$$

- Where N_i is the number of examples in cluster C_i
- Choose the cluster with larger value of compactness measure to split

47

Divisive Hierarchical Clustering

- To split a cluster, find a pair of examples having maximum Euclidian distance and split around these two examples (keeping them as centroids)
- At each level, number of clusters is increases by one
- The process continues until each example forms a cluster (atomic or singleton cluster) or until it satisfies certain termination condition
 - Termination condition could be
 - Number of clusters
 - Compactness measure of each cluster is within a certain threshold
- **Once two examples are placed in two different clusters at a level, they remain in different clusters at all subsequent levels**
- Example: DIvisive ANALysis (DIANA)

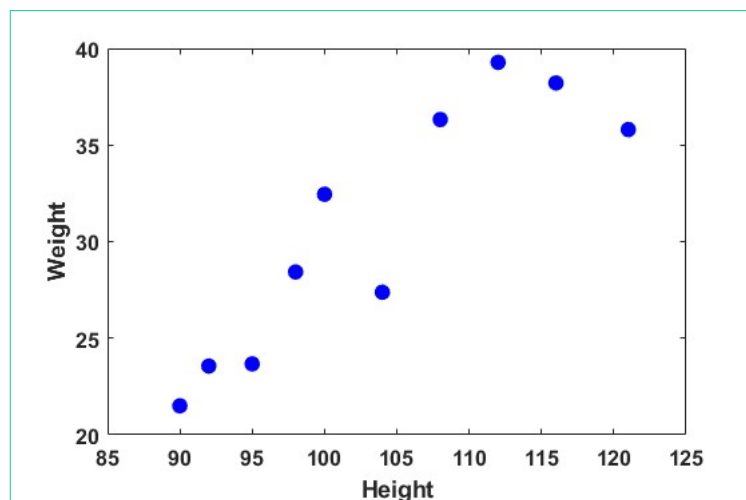
48

Divisive Hierarchical Clustering

- **Given:** Training data, $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N, \mathbf{x}_n \in \mathbb{R}^d$
- **Target:** Partition the data
- **Step1:** Single cluster having all the examples
- **Step2:** Find a pair of examples with in a cluster having maximum Euclidian distance
 - These examples act as centroid
- **Step3:** Split into two clusters by assigning each data point to one of these two examples
- **Step4:** Compute **compactness measure** for each cluster
- **Step5:** Choose the cluster with **larger value** of compactness measure to split
- **Step6:** Repeat **Step2** to **Step5** until **each example forms a cluster (atomic or singleton cluster)** or until it satisfies certain termination condition

49

Divisive Hierarchical Clustering

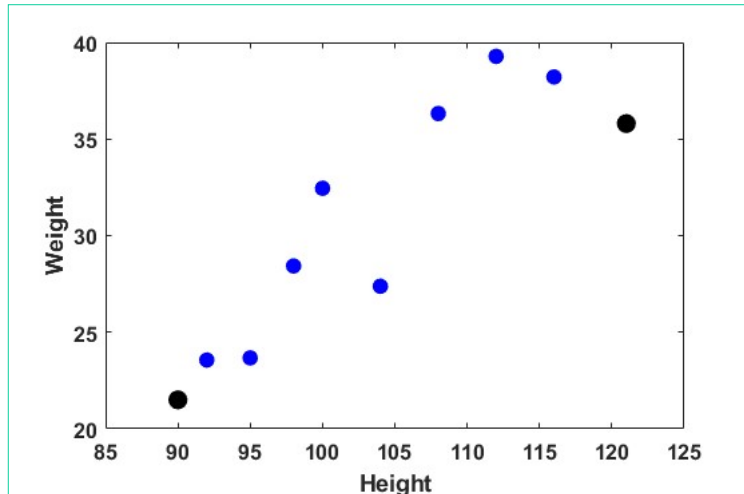


Level: 0

Number of clusters: 1

50

Divisive Hierarchical Clustering

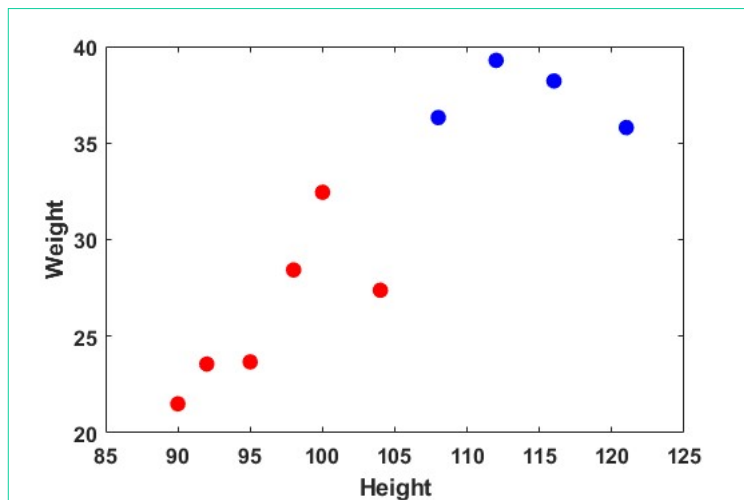


Level: 0

Number of clusters: 1

51

Divisive Hierarchical Clustering

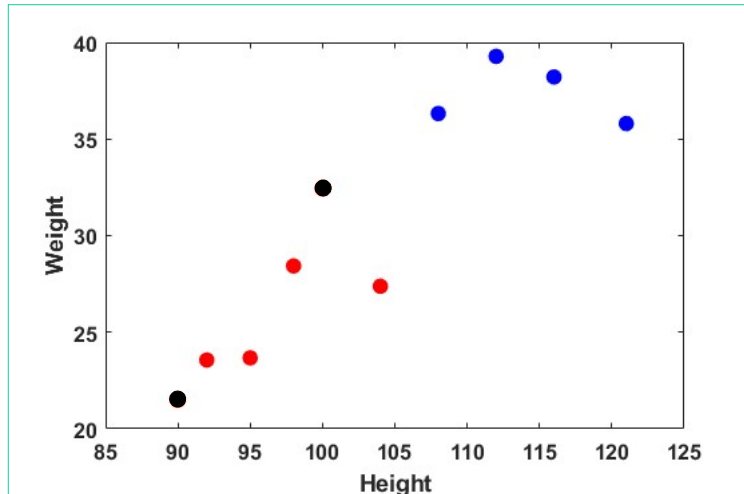


Level: 1

Number of clusters: 2

52

Divisive Hierarchical Clustering

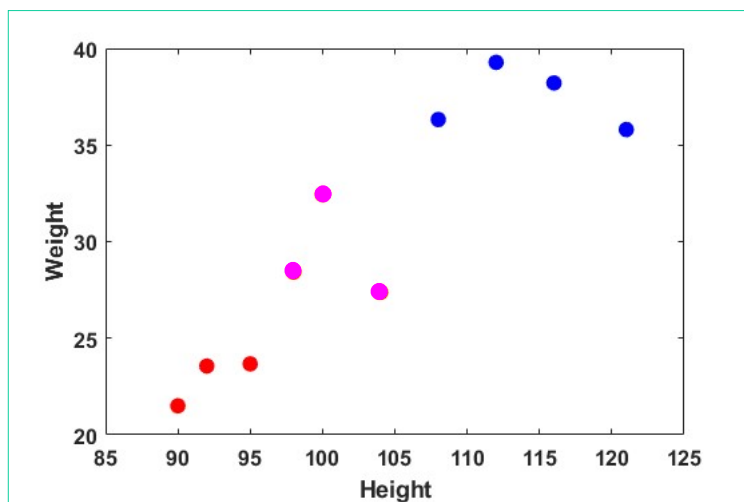


Level: 1

Number of clusters: 2

53

Divisive Hierarchical Clustering



Level: 2

Number of clusters: 3

54

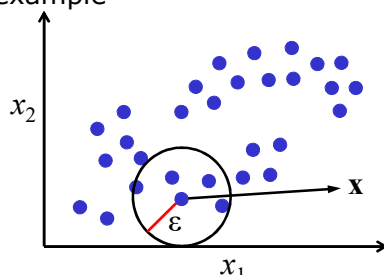
Density-Based Clustering

Density-Based Clustering

- These methods cluster collection of examples based on the notion of **density**
- These methods regard **clusters as dense regions of examples** in the data space that are separated by regions of low density (i.e. noise)
- They discover clusters with **arbitrary shape**
- **They automatically identifies the number of clusters**
- **General idea:** To continue growing the given cluster as long as density (number of examples) in the neighbourhood exceeds some threshold
- Example: **Density-based Special Clustering of Applications with Noise (DBSCAN)**
 - It grows the clusters according to a density-based connectivity analysis

Density-based Special Clustering of Applications with Noise (DBSCAN)

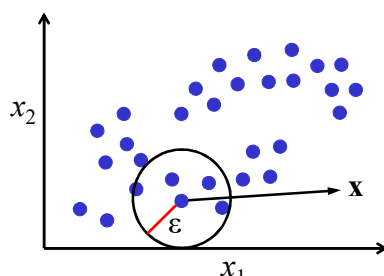
- DBSCAN is a density-based clustering included with noise
- It grows the regions with sufficiently high density (neighbors) into clusters with arbitrary shape
- It defines a cluster as a maximal set of **density-connected points**
- DBSCAN has 5 important components:
 1. **Epsilon (ϵ)**: It is a value of radius of boundary from every example



57

Density-based Special Clustering of Applications with Noise (DBSCAN)

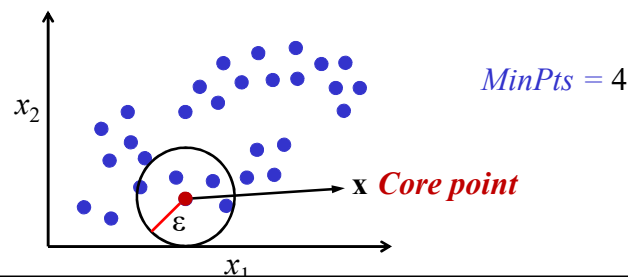
- DBSCAN has 5 important components:
 1. **Epsilon (ϵ)**: It is a value of radius of boundary from every example
 2. **MinPts**: **Minimum number of examples** present inside the boundary with radius of ϵ from an example x
 - These examples with in a boundary are neighbors to x and called as **ϵ -neighborhood of an example, x**



58

Density-based Special Clustering of Applications with Noise (DBSCAN)

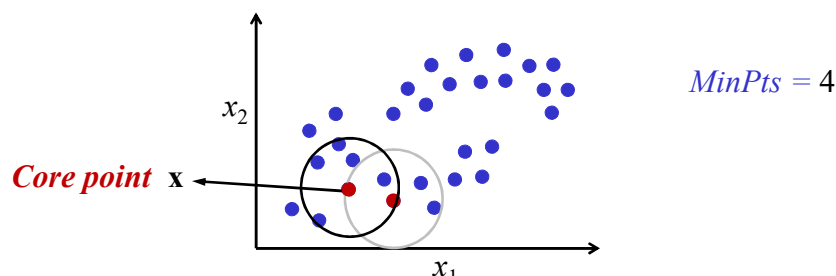
- DBSCAN has 5 important components:
 - Epsilon (ϵ)**: It is a value of radius of boundary from every example
 - MinPts**: Minimum number of examples present inside the boundary with radius of ϵ from an example x
 - Core point**: If there are atleast **MinPts** number of examples are with in ϵ -radius from x , then x is called as **core point**



59

Density-based Special Clustering of Applications with Noise (DBSCAN)

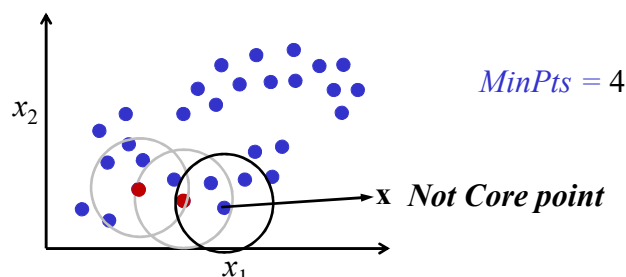
- DBSCAN has 5 important components:
 - Epsilon (ϵ)**: It is a value of radius of boundary from every example
 - MinPts**: Minimum number of examples present inside the boundary with radius of ϵ from an example x
 - Core point**: If there are atleast **MinPts** number of examples are with in ϵ -radius from x , then x is called as **core point**



60

Density-based Special Clustering of Applications with Noise (DBSCAN)

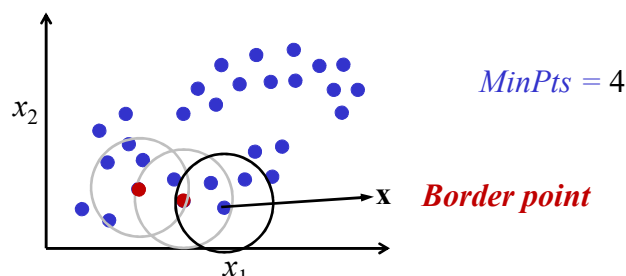
- DBSCAN has 5 important components:
 - Epsilon (ϵ)**: It is a value of radius of boundary from every example
 - MinPts**: Minimum number of examples present inside the boundary with radius of ϵ from an example x
 - Core point**: If there are atleast **MinPts** number of examples are with in ϵ -radius from x , then x is called as **core point**



61

Density-based Special Clustering of Applications with Noise (DBSCAN)

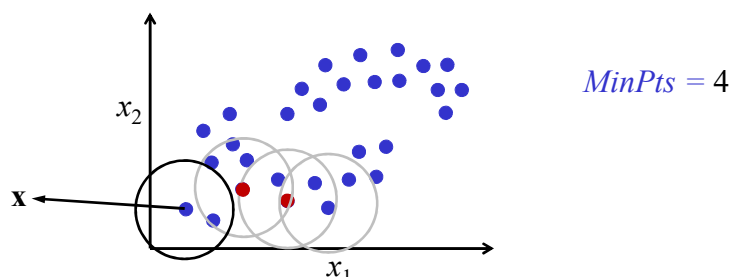
- DBSCAN has 5 important components:
 - Core point**: If there are atleast **MinPts** number of examples are with in ϵ -radius from x , then x is called as **core point**
 - Border point**:
 - The number of examples within ϵ -radius from x is less than **MinPts**
 - Atleast one of the example in neighborhood is **core point**, then x is called as **border point**



62

Density-based Special Clustering of Applications with Noise (DBSCAN)

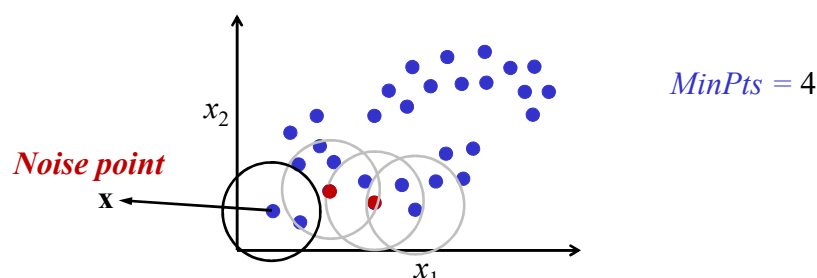
- DBSCAN has 5 important components:
 3. **Core point**: If there are atleast *MinPts* number of examples are with in ϵ -radius from \mathbf{x} , then \mathbf{x} is called as **core point**
 4. **Border point**:
 - The number of examples within ϵ -radius from \mathbf{x} is **less than *MinPts***
 - Atleast one of the example in neighborhood is **core point**, then \mathbf{x} is called as **border point**



63

Density-based Special Clustering of Applications with Noise (DBSCAN)

- DBSCAN has 5 important components:
 5. **Noise point**:
 - The number of examples within ϵ -radius from \mathbf{x} is **less than *MinPts***
 - No example in neighborhood is **core point**
 - The noise point is similar to outlier



64

Clustering using DBSCAN

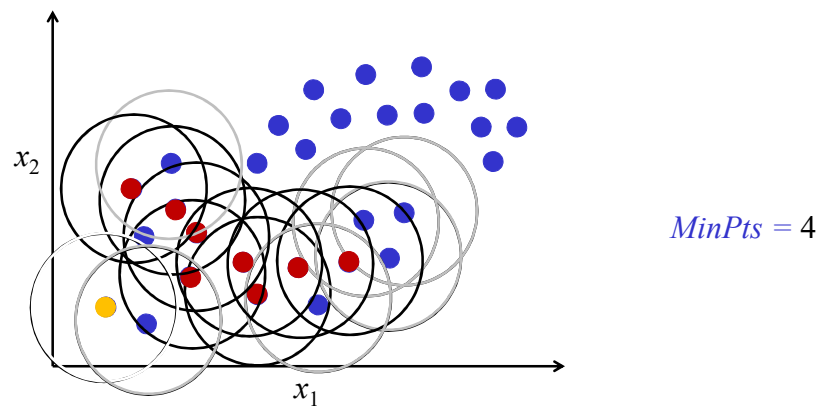
- **Given:** Training data, $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N, \mathbf{x}_n \in \mathbb{R}^d$
- First step is to identify **core points**, **border points** and **noise points**
 - Only **core points** and **border points** are considered inside the cluster
 - **Noise points** are not taken into the cluster
 - Thus DBSCAN is robust to outliers
- Next step is to find the **connected components** of core points
- **Connected component of core points:** Connecting the core points that are reachable from any point
 - All the connected (reachable) core points form a cluster

Clustering using DBSCAN

- The **connected component of core points** is obtained by understanding following *two* definitions.
- **Directly density-reachable:** A core point \mathbf{x}_i is directly density-reachable to a core point \mathbf{x}_j , if the core point \mathbf{x}_j is within ϵ -distance from core point \mathbf{x}_i
- **Density-reachable:** A core point \mathbf{x}_i is indirectly reachable to another core point \mathbf{x}_j through other core points, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_K$ such that
 - \mathbf{x}_i is directly density-reachable to \mathbf{x}_1
 - \mathbf{x}_1 is directly density-reachable to \mathbf{x}_2
 - \mathbf{x}_k is directly density-reachable to \mathbf{x}_{k+1}
 - \mathbf{x}_K is directly density-reachable to \mathbf{x}_j

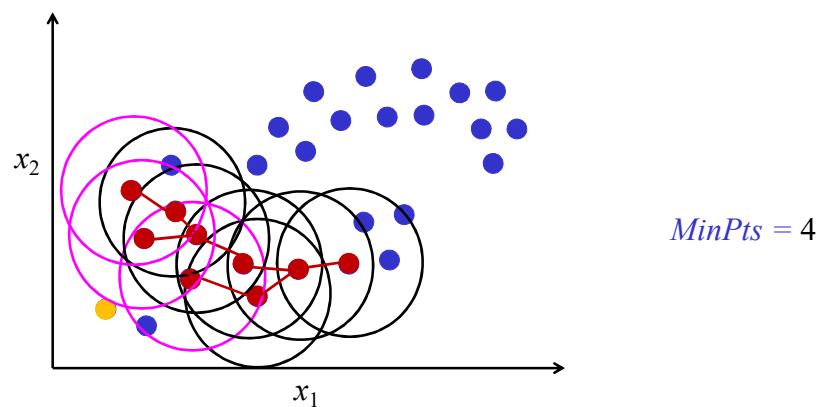
Clustering using DBSCAN

- Directly density-reachable:
- Density-reachable:



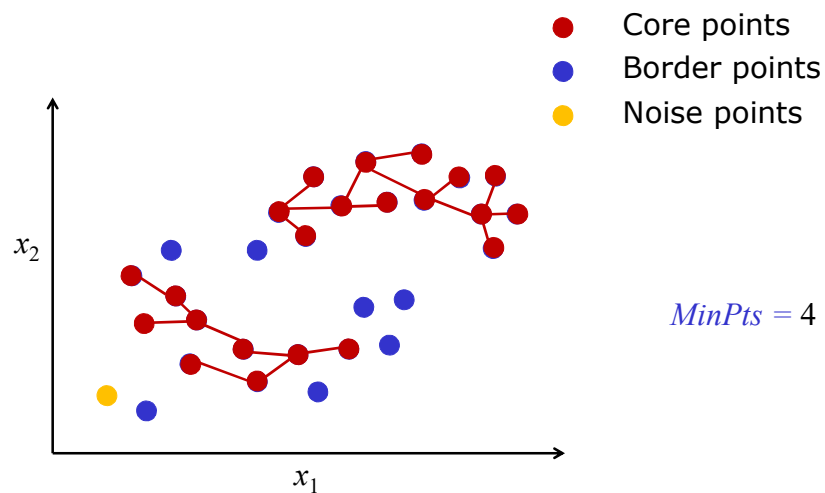
Clustering using DBSCAN

- Directly density-reachable:
- Density-reachable:



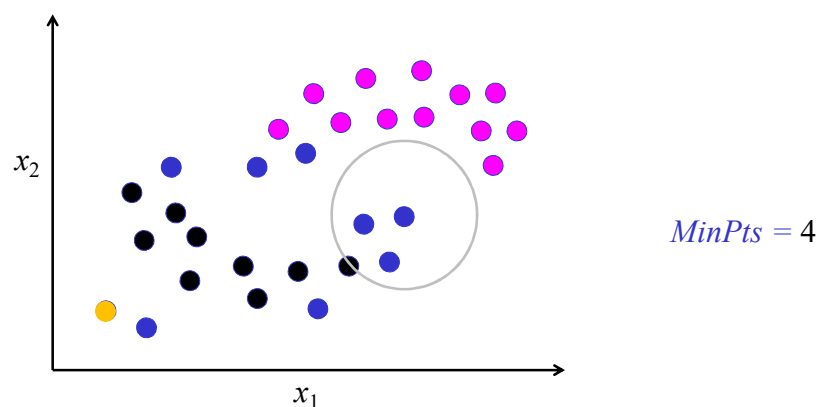
Clustering using DBSCAN

- All the core points with connected component forms a cluster



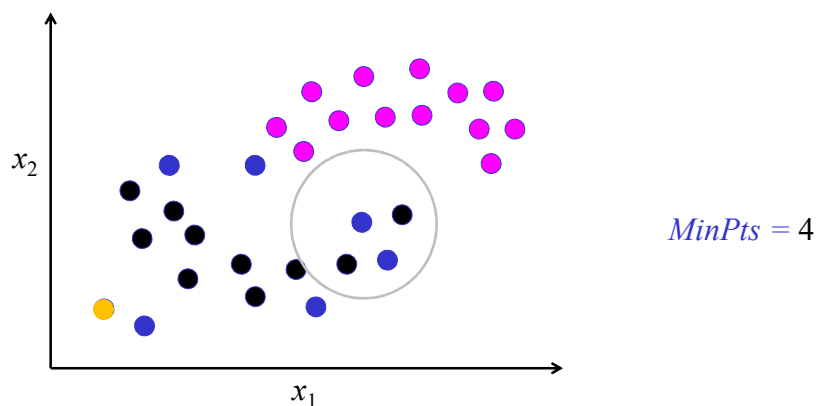
Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



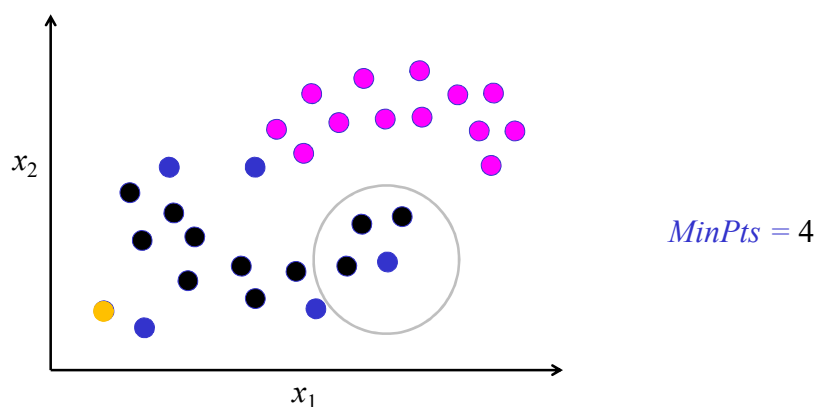
Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



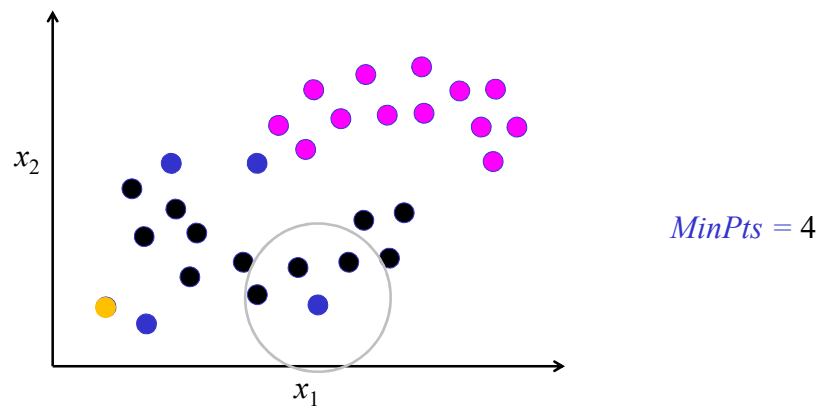
Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



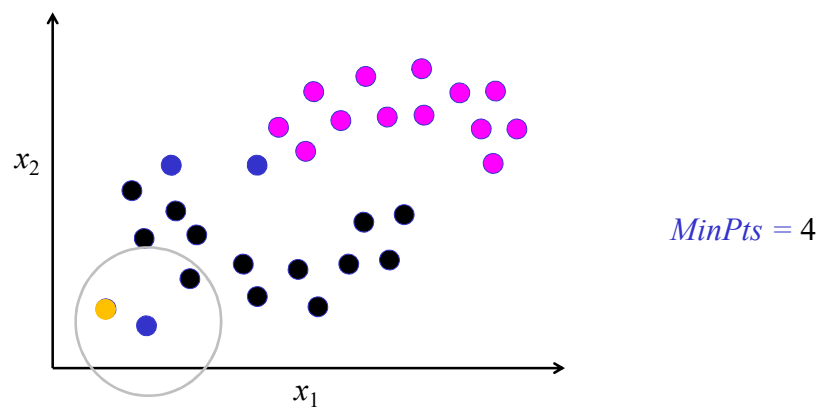
Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



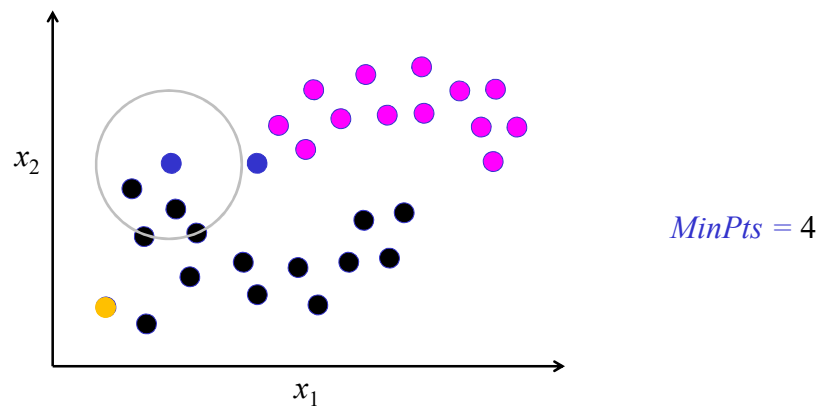
Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



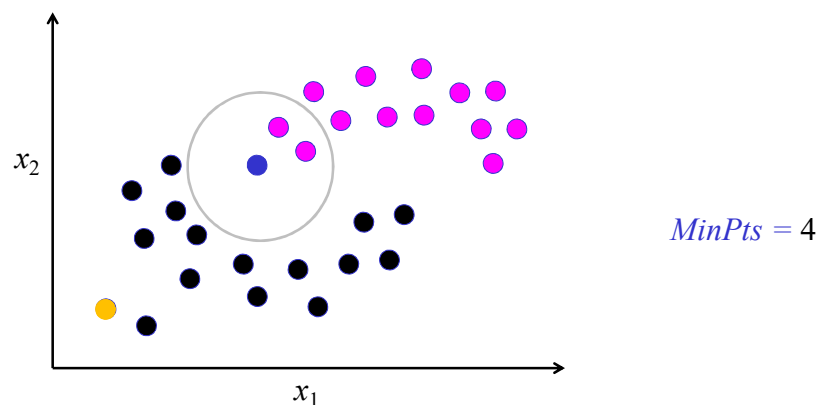
Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



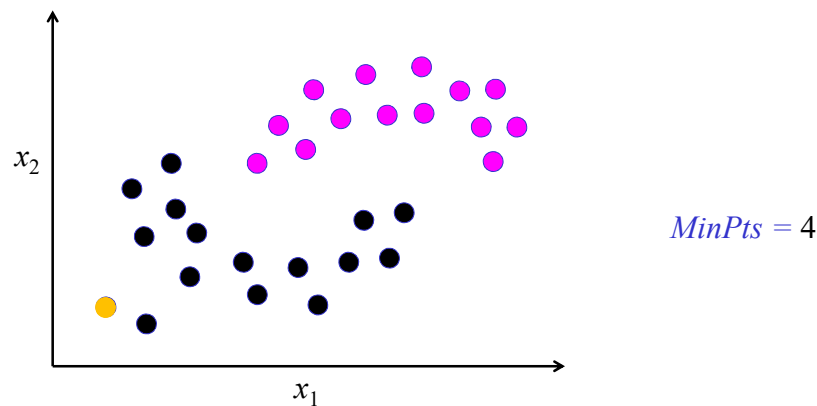
Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



Clustering using DBSCAN

- All the core points with connected component forms a cluster
- Assign the **border points** to nearby cluster which is at ϵ -radius from that border point



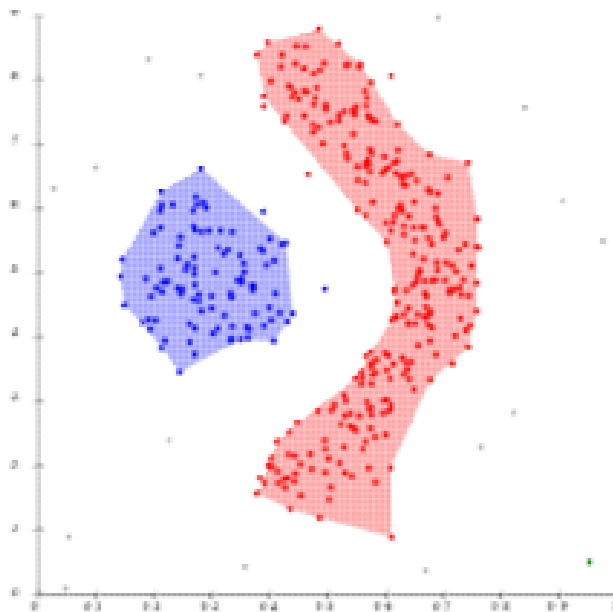
Clustering using DBSCAN

- **Training process:**
- **Given:** Training data, $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N, \mathbf{x}_n \in \mathbb{R}^d$
- Identify the core points, border points and noise points
- Find the connected components of core points
- Each connected component forms a cluster
- Assign each of the border points to a nearby cluster which is at ϵ -radius from that border point
- Noise points are not assigned to any clusters
- Training process, stores the core points as model

Clustering using DBSCAN

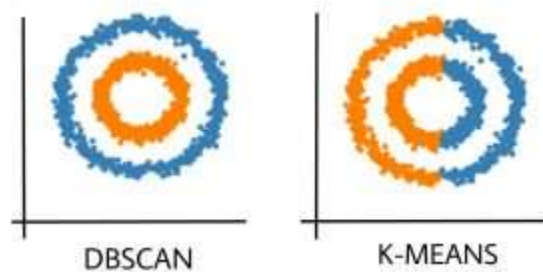
- **Test process:**
- For a test example, identify it as **core point** or **border point** or **noise point**
- If it is a **core point**, assign it to a cluster to which it is **directly density-reachable** or **density-reachable**
- If it is a **border point**, assign it to a **nearby cluster** which is at ϵ -radius from that border point
- If it is a **noise point**, do not assign to any cluster

Clustering using DBSCAN



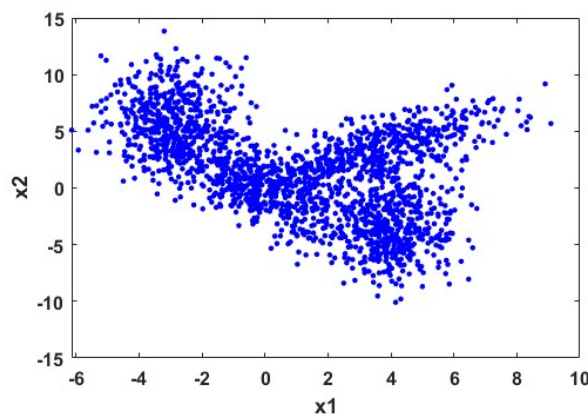
Advantages of DBSCAN

- DBSCAN **does not** require to specify the number of clusters in the data a priori
- DBSCAN can find **arbitrarily shaped clusters**
- DBSCAN has a notion of noise, and is robust to **outliers**
- The parameters ϵ and *MinPts* are experimentally set by the users



Limitation of DBSCAN

- DBSCAN is not suitable when the data is completely dense and there is no low dense area to separate



- The parameters ϵ and *MinPts* should be chosen carefully

Text Books

1. J. Han and M. Kamber, *Data Mining: Concepts and Techniques*, Third Edition, Morgan Kaufmann Publishers, 2011.
2. S. Theodoridis and K. Koutroumbas, *Pattern Recognition*, Academic Press, 2009.