## Data Mining and Machine Learning

Lecture 5: Unsupervised Learning

Dr Sarat Moka

UNSW, Sydney

# Key Topics

- Curse of dimensionality
- Projection Based Dimensionality Reduction
- Manifold Learning
- Principal component analysis (PCA)
- K-Means Clustering
- Hierarchical clustering
- DBSCAN

#### **Books:**

- (A) Data Science and Machine Learning: Mathematical and Statistical Methods, by Kroese, Botev, Taimre, and Vaisman. Click here to download a pdf copy.
- (B) Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow (2nd Ed.) by Aurélien Géron (2019).

Recall that, unlike supervised learning, in unsupervised learning, there is no *response* (or output) variable y. The goal is to extract useful information and pattern from the given feature vectors  $\mathbf{x}$ . Mostly, the objective is to learn about the underlying probability distribution.

In this lecture, we will learn about the curse of dimensionality and how to overcome this problem. We will then learn techniques to reduce the dimensions of datasets. We will also learn about unsupervised learning that will help you to categorise data by identifying the common features.

# 5.1 Curse of dimensionality

• The curse of dimensionality is a term coined by Richard Bellman in 1957 to describe the exponential increase in volume associated with adding extra dimensions to a mathematical space.

#### Illustration

Consider a *d*-dimensional hypercube with side length 2r. The volume V of this hypercube is given by:

$$V = (2r)^d.$$

As d increases, the volume grows exponentially, making it harder to cover the space uniformly with a finite number of data points. For example, if r = 1, the volume of the hypercube is  $2^d$ , which becomes impractically large as d increases.

- In the context of machine learning, the phrase 'curse of dimensionality' typically refers to the various challenges that arise when analyzing and modeling data with a high number of features.
- As the dimensionality of the feature space increases, several problems become more pronounced, impacting the performance and feasibility of learning algorithms.

### Challenges with High-Dimensional Feature Spaces

1. **Data Sparsity**: As the number of features increases, the data points become more sparse in the feature space. For a fixed number of data points, the density of data points decreases exponentially with the increase in dimensionality, making it harder to capture the underlying patterns.

- 2. **Overfitting**: High-dimensional spaces provide more flexibility for models to fit the training data. However, this often leads to overfitting, where the model captures noise in the training data rather than the true underlying distribution. This results in poor generalization to new, unseen data.
- 3. Increased Computational Complexity: The computational cost of many machine learning algorithms increases with the number of features. For example, the time complexity of algorithms such as k-nearest neighbors (KNN) and support vector machines (SVMs) grows with the dimensionality, making them impractical for very high-dimensional data.

### **Hughes Phenomenon**

- An important result related to the curse of dimensionality is the Hughes phenomenon, named after G.F. Hughes who described it in 1968.
- The Hughes phenomenon illustrates that, for a fixed number of training samples, the classification accuracy of a machine learning model initially improves as the number of features increases, but after a certain point, it starts to decline.
- This is due to the increased variance in the estimates of the model parameters as the number of features grows, leading to overfitting.



Figure 5.1: The Hughes phenomenon: classification accuracy vs. number of features.

• Mathematically, suppose we have a classification problem with N training samples and d features. If we denote the classification accuracy by A(d, N), then Hughes' result can be summarized as

$$A(d, N) \approx \begin{cases} \text{Increases with } d & \text{for small } d \\ \text{Decreases with } d & \text{for large } d. \end{cases}$$

This behavior is illustrated in Figure 5.1.

### Distance Measures in High Dimensions

- In high-dimensional spaces, the concept of distance becomes less intuitive. For instance, consider two points  $\mathbf{x}$  and  $\mathbf{x}'$  in  $\mathbb{R}^d$ .
- The Euclidean distance between these points is

$$\|\mathbf{x} - \mathbf{x}'\|_2 = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}.$$

As d increases, the distance between points tends to become more uniform.

#### Example

To illustrate this, let's assume we have two points  $\mathbf{x}$  and  $\mathbf{x}'$  in  $\mathbb{R}^d$ , where each coordinate is drawn independently from a uniform distribution over the interval [0, 1]. The Euclidean distance between these points is given by

$$\|\mathbf{x} - \mathbf{x}'\|_2 = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}.$$

Due to the law of large numbers, as the dimensionality d increases, the sum of squared differences  $\sum_{i=1}^{d} (x_i - x'_i)^2$  will tend to be concentrated around its expected value. For uniformly distributed points, the expected value of  $(x_i - x'_i)^2$  is  $\frac{1}{6}$ . Thus, for large d,  $\sum_{i=1}^{d} (x_i - x'_i)^2 \approx \frac{d}{6}$ .

So, the Euclidean distance will be approximately  $\|\mathbf{x} - \mathbf{x}'\|_2 \approx \sqrt{\frac{d}{6}}$ .

This means that the distances between different pairs of points will tend to be around  $\sqrt{\frac{d}{6}}$ . The variance of these distances decreases relative to the mean distance as d increases, making the distances between all pairs of points more similar, or "uniform".

- This phenomenon is evident when comparing the distances between the nearest and farthest neighbors of a point in high dimensions.
- The ratio of the distances between the nearest and farthest neighbors approaches 1 as d increases, leading to a loss of contrast in distances.
- In the context of nearest neighbor searches, this uniformity of distances implies that the ratio of the distance to the nearest neighbor to the distance to the farthest neighbor approaches 1.
- Mathematically, if  $dist_{min}$  and  $dist_{max}$  are the minimum and maximum distances from a given point to all other points in the dataset, then

$$\lim_{d \to \infty} \frac{\text{dist}_{\min}}{\text{dist}_{\max}} \to 1$$

• This makes it difficult to distinguish between the closest and farthest neighbors, which can significantly degrade the performance of algorithms that rely on distance measures, such as *k*-nearest neighbors (KNN) and clustering algorithms like *k*-means.

### **Implications and Mitigation Strategies**

The curse of dimensionality implies that adding more features to a dataset can lead to diminishing returns and potentially worse performance if the model overfits the data. To mitigate these effects, several strategies can be employed:

- 1. Feature Selection: Select a subset of relevant features that contribute most to the predictive power of the model. Techniques include mutual information, recursive feature elimination, and regularization methods like Lasso.
- 2. Dimensionality Reduction: Transform the high-dimensional feature space into a lower-dimensional space while preserving the important information. Common techniques include Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA), and t-Distributed Stochastic Neighbor Embedding (t-SNE).
- 3. **Regularization**: Add a penalty term to the learning algorithm to constrain the model complexity, thereby preventing overfitting. Examples include Ridge regression, Lasso regression, and Elastic Net.
- 4. **Increase Training Data**: Collect more training samples to provide sufficient coverage of the high-dimensional space. This can be challenging and expensive but is often necessary for high-dimensional data.

## 5.2 Projection Based Dimensionality Reduction

- Often training instances are not distributed equally across all dimensions, and it's possible to transform them from a high-dimensional space to a low-dimensional space such that the low-dimensional space representation keeps most of the properties available in the original dimensions.
- For example, converting the following 3D data set to a new 2D subspace after projection retains most of the properties in the data set. See Figure 5.2.



Figure 5.2: The new 2D data set after projection. Adapted from Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow by A. Géron, 2019, Sebastopol; CA: O'Reilly Media.

• However, the above approach may not prove useful for some data sets. For example, for the Swiss roll toy data set shown in Figure 5.3, its projection onto a plane doesn't retain properties from the original data set properties.

## 5.3 Manifold Learning and Dimensionality Reduction

- Manifold learning is a technique in machine learning used for dimensionality reduction. The goal of manifold learning is to discover the low-dimensional structure (manifold) embedded within high-dimensional data.
- Unlike traditional linear dimensionality reduction methods like Principal Component Analysis (PCA) (which is covered later), manifold learning techniques are capable of capturing the nonlinear structure of the data.
- Manifold learning techniques are powerful tools for uncovering the underlying structure in high-dimensional datasets, making them useful in various applications such as visualization, clustering, and classification.

### **Concept of Manifold**



Figure 5.3: Swiss roll data set. Adapted from Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow by A. Géron, 2019, Sebastopol; CA: O'Reilly Media.

- A manifold is a topological space that locally resembles Euclidean space.
- For example, the surface of a sphere is a 2-dimensional manifold embedded in 3-dimensional space.
- In the context of high-dimensional data, we assume that the data points lie on or near a low-dimensional manifold within the high-dimensional space; see Figure 5.3.

### **Common Manifold Learning Techniques**

There are several techniques for manifold learning, including:

- **Isomap**: This method seeks to preserve the geodesic distances between all pairs of data points. It uses the shortest path algorithm to approximate the geodesic distance on the manifold.
- Locally Linear Embedding (LLE): LLE preserves local relationships by attempting to keep each point's local neighborhood similar to that in the high-dimensional space.
- t-Distributed Stochastic Neighbor Embedding (t-SNE): t-SNE converts similarities between data points to joint probabilities and minimizes the Kullback-Leibler divergence between the joint probabilities of the low-dimensional embedding and the high-dimensional data.

• Laplacian Eigenmaps: This method uses the graph Laplacian to preserve locality by constructing a weighted graph of the data points and then finding the low-dimensional representation that best preserves the graph structure.

**Mathematical Formulation:** Given a high-dimensional dataset  $\mathbf{X} = \{x_1, x_2, \ldots, x_N\}$  with  $x_i \in \mathbb{R}^D$ , the aim is to find a low-dimensional representation  $\mathbf{Y} = \{y_1, y_2, \ldots, y_N\}$  with  $y_i \in \mathbb{R}^d$  where  $d \ll D$ .

• Isomap:

Minimize 
$$\sum_{i < j} (\delta_{ij}^{(G)} - \|y_i - y_j\|)^2$$

where  $\delta_{ij}^{(G)}$  is the geodesic distance between points  $x_i$  and  $x_j$  in the high-dimensional space.

• LLE:

Minimize 
$$\sum_{i} \|y_i - \sum_{j \in \mathcal{N}(i)} w_{ij} y_j\|^2$$

subject to  $\sum_{i \in \mathcal{N}(i)} w_{ij} = 1$ , where  $\mathcal{N}(i)$  denotes the set of neighbors of  $x_i$ .

• t-SNE:

Minimize 
$$KL(P||Q) = \sum_{i} \sum_{j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

where  $p_{ij}$  and  $q_{ij}$  are the probabilities that points  $x_i$  and  $x_j$  are neighbors in highdimensional and low-dimensional spaces, respectively.

• Laplacian Eigenmaps:

Minimize 
$$\sum_{i,j} \|y_i - y_j\|^2 w_{ij}$$

where  $w_{ij}$  is the weight of the edge connecting  $x_i$  and  $x_j$  in the graph representation.

#### Remark

Manifold learning often operates under the implicit assumption that the task at hand (e.g., classification or regression) will be simpler when expressed in the lowerdimensional space of the manifold. For instance, in the top row of Figure 5.4, the Swiss roll is divided into two classes: in the 3D space (left), the decision boundary is quite complex, but in the 2D unrolled manifold space (right), the decision boundary is a straight line.

However, this assumption does not always hold true. For example, in the bottom row of Figure 5.4, the decision boundary is at  $x_1 = 5$ . This boundary is very simple in the original 3D space (a vertical plane), but it becomes more complex in the unrolled manifold (a collection of four independent line segments).



Figure 5.4: The decision boundary may not always be simpler with lower dimensions. Adapted from Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow by A. Géron, 2019, Sebastopol; CA: O'Reilly Media

## 5.4 Principal Component Analysis (PCA)

- The main goal of PCA is to reduce the dimensionality of a given data set consists of many features. This is precisely the reason why the PCA is called *feature reduction* (or feature extraction) mechanism.
- In PCA we start with a *d*-dimensional data set  $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathbb{R}^d$ .

- PCA makes no assumptions on how the data is generated.
- Our goal is to represent this d-dimensional input data using n feature vectors of dimension k < d. This is achieved as follows.
- First write the input data as a matrix:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{1}^{\top} \\ \mathbf{x}_{2}^{\top} \\ \vdots \\ \mathbf{x}_{n}^{\top} \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,d} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,d} \end{bmatrix}.$$

• Assume that the underlying distribution of the data has mean **0**.

#### Remark

In practice, this assumption is satisfied by centering the data before applying PCA, that is, subtract the column mean for every column. In other words, the input to the PCA is of the form:

$$x'_{i,j} = x_{i,j} - \frac{1}{n} \sum_{j=1}^{n} x_{i,j}$$

for all i = 1, 2, ..., d.

• Further, let  $\Sigma$  be the covariance matrix of the underlying distribution of the data. By definition,

$$\mathbf{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^ op]$$

and it can be estimated using

$$\widehat{\mathbf{\Sigma}} = rac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^{ op}$$

• Since  $\widehat{\Sigma}$  is a symmetric matrix, suppose

$$\widehat{\boldsymbol{\Sigma}} = \mathbf{U}\mathbf{D}\mathbf{U}^{\top},$$

is the singular value decomposition of  $\Sigma$ .

- Principal Components: The k columns of U corresponding to k largest diagonal elements of D are called k principal components.
- Let  $\mathbf{U}_k$  is the matrix whose columns are the k principal components.

• Then, the map

$$\mathbf{z} := \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}$$

maps  $\mathbf{x}$  to a vector  $\mathbf{z} \in \mathbb{R}^d$  lying on the subspace spanned by the k principal components. This allows us to represent the data with respect to k principal components. For that we take another mapping:

$$\mathbf{z}' = \mathbf{U}_k^\top \mathbf{z} = \mathbf{U}_k^\top \left( \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x} \right) = \mathbf{U}_k^\top \mathbf{x}.$$

- That is  $\mathbf{z}'_1, \ldots, \mathbf{z}'_n \in \mathbb{R}^k$  are the k dimensional representation of the given d-dimensional data  $\mathbf{x}_1, \ldots, \mathbf{x}_n$  with respect the standard basis.
- This procedure is called PCA!
- Let **D** be a diagonal matrix with diagonal being  $(d_{1,1}, d_{2,2}, \ldots, d_{n,n})$ . Then,  $d_{l,l}$  can be interpreted as variance of the data in the *l*-th principal component and hence

$$\nu = \sum_{l=1}^{d} d_{l,l}$$

can be used as a measure of the variance in the data.

#### Question

Why to select k principal components? Why not any set of k columns of **U**?

A: It minimizes the total squared distance error between the projected points and the original points, i.e., keeps as much variance of the given data as possible. In other words, for any  $k \leq d$ , the eigenvectors  $\mathbf{u}_1, \ldots, \mathbf{u}_k$  corresponds to the k principal components are the solution of the following optimization:

$$\min_{\mathbf{v}_1,\dots,\mathbf{v}_k} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{V}_k \mathbf{V}_k^\top \mathbf{x}_i\|^2,$$
(5.1)

where the minimization is taken over all k unit length orthogonal vectors  $\{\mathbf{v}_1, \ldots, \mathbf{v}_k\}$ in  $\mathbb{R}^d$  and  $\mathbf{V}_k = [\mathbf{v}_1, \ldots, \mathbf{v}_k]$ .

## 5.5 *K*-means Clustering

• In *clustering*, the goal is to divide the given unlabelled feature vectors  $D = {\mathbf{x}_1, ..., \mathbf{x}_n}$  into a set of K groups (or clusters), so that the samples belong to the same group are more *similar* to each other than the samples belong to different groups.

Q: What is difference between classification and clustering?

- *K*-means clustering is a simple heuristic method which ignores the distributional properties of the data.
- In K-means clustering, the entire feature vector space is divided into K regions using a distance function  $\text{Dist}(\mathbf{x}, \mathbf{x}')$ . Some well-known choice of the distance functions are

$$\operatorname{Dist}(\mathbf{x}, \mathbf{x}') = \begin{cases} \|\mathbf{x} - \mathbf{x}'\| = \sqrt{\sum_{i=1}^{d} (x_i - x'_i)^2}, & (\operatorname{Euclidean}) \\\\ \sum_{i=1}^{d} |x_i - x'_i|, & (\operatorname{Manhattan}) \\\\ \max_{i=1,\dots,d} |x_i - x'_i|, & (\operatorname{Maximum}) \\\\ \sum_{i=1}^{d} \mathbb{I}(x_i \neq x'_i), & (\operatorname{Hamming distance for binary features}) \end{cases}$$

• Once we fix the distance function, we can choose K points  $\mathbf{c}_1, \ldots, \mathbf{c}_K$  as cluster centers. Further, we partition the whole feature space into regions  $\mathcal{R}_1, \ldots, \mathcal{R}_K$ , where

$$\mathcal{R}_k = \{ \mathbf{x} : \text{Dist}(\mathbf{x}, \mathbf{c}_k) \le \text{Dist}(\mathbf{x}, \mathbf{c}_j), \quad \forall j \neq k \}.$$



Figure 5.5: Voronoi diagram

Q: How to find the centers 
$$\mathbf{c}_1, \ldots, \mathbf{c}_K$$
?

- Starting from an initial guess for the centers, it iteratively finds the new centers as the centroids of the points in each Voronoi cell formed by the current centers.
- Let's see the algorithm:

#### Algorithm 1: *K*-Means Algorithm

Input: The training data, number of clusters K, initial centers  $\mathbf{c}_1, \ldots, \mathbf{c}_K$ . Output: Cluster centers and cells while a stopping criteria is not met do  $\begin{bmatrix} \mathcal{R}_1, \ldots, \mathcal{R}_K \leftarrow \varnothing \\ \mathbf{for} \quad i = 1 \text{ to } n \text{ do} \\ & \mathbf{d} \leftarrow (\text{Dist}(\mathbf{x}_i, \mathbf{c}_1), \ldots, \text{Dist}(\mathbf{x}_i, \mathbf{c}_K)) \\ & k \leftarrow \arg \min_j d_j \\ & \mathcal{R}_k \leftarrow \mathcal{R}_k \cup \{\mathbf{x}_i\} \\ \text{ end} \\ & \mathbf{for} \quad k = 1 \text{ to } K \text{ do} \\ & & \left| \begin{array}{c} \mathbf{c}_k \leftarrow \frac{\sum_{\mathbf{x} \in \mathcal{R}_k} \mathbf{x}}{|\mathcal{R}_k|} \\ & \mathbf{end} \\ & \mathbf{end} \\ & \mathbf{end} \\ \end{array} \right|$ end end Return:  $\{\mathcal{R}_k\}$  and  $\{\mathbf{c}_k\}$ ;

• An example of K-means clustering can be seen in the following figure. The thick dark circles are the final centroids.



## 5.6 Hierarchical clustering

Hierarchical clustering is a method of cluster analysis which seeks to build a hierarchy of clusters. It is particularly useful in fields such as evolutionary biology for constructing phylogenetic trees. There are two primary types of hierarchical clustering:

- 1. Agglomerative Clustering (Bottom-Up Approach): This method starts with each data point as its own cluster and iteratively merges the closest pairs of clusters until only one cluster remains or a certain number of clusters is achieved.
- 2. Divisive Clustering (Top-Down Approach): This method starts with all data points in a single cluster and iteratively splits the clusters into smaller clusters.

## Agglomerative Clustering

Steps in Agglomerative Clustering:

- 1. Initialization: Each data point is considered as a single cluster.
- 2. Distance Calculation: Compute the distance between each pair of clusters.
- 3. Merge Clusters: Merge the pair of clusters with the smallest distance.
- 4. Update Distances: Recompute the distances between the new cluster and the remaining clusters.
- 5. **Repeat**: Continue the process until a single cluster remains or the desired number of clusters is reached.

#### Example

Consider a dataset with points labeled from 1 to 8. Initially, each point is its own cluster. As the algorithm proceeds, clusters are merged step-by-step. For example, points 1 and 2 might merge to form a new cluster, labeled 9. Points 3 and 4 merge to form cluster 10. The process continues until a hierarchy is formed, which can be represented using a **dendrogram**.

## Dendrogram

A dendrogram is a tree-like diagram that records the sequences of merges or splits. It provides a visual representation of the hierarchical clustering process. The vertical axis represents the distance or dissimilarity between clusters, while the horizontal axis represents the individual data points or clusters. See Figure 5.6 for an illustration.



Figure 5.6: Dendrogram Example: Left: a cluster hierarchy of 15 clusters. Right: the corresponding dendrogram.

## Linkage Criteria

The way distances between clusters are defined can vary, leading to different types of linkage criteria:

• Single Linkage: The distance between two clusters is defined as the minimum distance between any single data point in the first cluster and any single data point in the second cluster.

$$d_{\min}(I,J) = \min_{i \in I, j \in J} \operatorname{dist}(x_i, x_j).$$

• **Complete Linkage**: The distance between two clusters is defined as the maximum distance between any single data point in the first cluster and any single data point in the second cluster.

$$d_{\max}(I,J) = \max_{i \in I, i \in J} \operatorname{dist}(x_i, x_j).$$

• Average Linkage: The distance between two clusters is defined as the average distance between all pairs of data points from the two clusters.

$$d_{\text{avg}}(I,J) = \frac{1}{|I||J|} \sum_{i \in I} \sum_{j \in J} \operatorname{dist}(x_i, x_j).$$

• Ward's Minimum Variance Method: This method aims to minimize the total within-cluster variance. The distance between two clusters is defined as the increase in the total within-cluster variance after merging the two clusters.

$$d_{\text{Ward}}(I,J) = \frac{|I||J|}{|I|+|J|} \|\mathbf{x}_I - \mathbf{x}_J\|^2$$

where  $\mathbf{x}_I$  and  $\mathbf{x}_J$  are the centroids of clusters I and J respectively.

#### Algorithm for Agglomerative Clustering

Here is a pseudocode representation of the agglomerative clustering algorithm:

#### Algorithm 4.7.1: Greedy Agglomerative Clustering

- 1. Initialize the set of cluster identifiers:  $I = \{1, ..., n\}$ .
- 2. Initialize the corresponding label sets:  $L_i = \{i\}, \forall i \in I$ .
- 3. Initialize a distance matrix  $D = [d_{ij}]$  with  $d_{ij} = d(\{i\}, \{j\})$ .
- 4. For k = n + 1 to 2n K do:
  - (a) Find i and j > i in I such that  $d_{ij}$  is minimal.
  - (b) Create a new label set  $L_k := L_i \cup L_j$ .
  - (c) Add the new identifier k to I and remove the old identifiers i and j from I.
  - (d) Update the distance matrix D with respect to the identifiers i, j, and k.
- 5. Return  $L_i, \forall i = 1, \ldots, 2n K$ .

#### Lance-Williams Update

- The Lance-Williams update formula is a recursive method for calculating the distance between clusters in hierarchical clustering. It allows for efficient computation of cluster distances when clusters are merged during the agglomerative clustering process.
- Suppose at some stage in the clustering algorithm, clusters I and J, with identifiers i and j respectively, are merged into a new cluster  $K = I \cup J$  with identifier k. Let M, with identifier m, be another existing cluster. The Lance-Williams update rule provides a way to update the distance  $d_{km}$  between the new cluster K and the cluster M using the distances  $d_{im}$ ,  $d_{jm}$ , and  $d_{ij}$ .
- The general form of the Lance-Williams update formula is:

$$d_{km} = \alpha d_{im} + \beta d_{jm} + \gamma d_{ij} + \delta |d_{im} - d_{jm}|,$$

where  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  are constants that depend on the chosen linkage criterion. These constants are derived based on the specific characteristics of the clusters involved, such as their sizes.

• Table 5.1 lists the constants  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  for some common linkage criteria.

Table 5.1: Constants for the Lance-Williams update rule for various linkage functions, where  $n_i$ ,  $n_j$ , and  $n_m$  denote the number of elements in the corresponding clusters.

Linkage	$\alpha$	$\beta$	$\gamma$	$\delta$
Single	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\frac{1}{2}$
Complete	$\frac{\overline{1}}{2}$	$\frac{\overline{1}}{2}$	0	$\frac{1}{2}$
Group Average	$\frac{ ilde{n_i}}{ ilde{n_i+n_j}}$	$\frac{\tilde{n_j}}{n_i+n_j}$	0	Õ
Ward	$\frac{n_i + n_m}{n_i + n_i + n_m}$	$\frac{n_j + n_m}{n_i + n_i + n_m}$	$-\frac{n_m}{n_i+n_i+n_m}$	0

• For example, in the single linkage criterion, the distance between two clusters is defined as the minimum distance between any pair of points from the two clusters. Using the Lance-Williams update, this is expressed as:

$$d_{km} = \min\{d_{im}, d_{jm}\} = \frac{d_{im}}{2} + \frac{d_{jm}}{2} - \frac{|d_{im} - d_{jm}|}{2}$$

• In the complete linkage criterion, the distance between two clusters is defined as the maximum distance between any pair of points from the two clusters. The update formula for complete linkage is:

$$d_{km} = \max\{d_{im}, d_{jm}\} = \frac{d_{im}}{2} + \frac{d_{jm}}{2} + \frac{|d_{im} - d_{jm}|}{2}$$

• In practical applications, the Lance-Williams update formula significantly reduces the computational complexity of hierarchical clustering algorithms by enabling efficient recalculation of distances as clusters are merged. This makes it feasible to apply hierarchical clustering to larger datasets.

### **Divisive Approach**

- In contrast to the bottom-up (agglomerative) approach to hierarchical clustering, the divisive approach starts with a single cluster containing all data points. This cluster is then recursively divided into smaller clusters. The goal of each division is to create two clusters that are as "dissimilar" as possible, which can then be further divided in subsequent steps.
- The divisive approach can utilize the same linkage criteria as agglomerative clustering to determine the best way to divide a parent cluster into two child clusters. These linkage criteria include single linkage, complete linkage, average linkage, and Ward's minimum variance linkage. The choice of linkage criterion affects how the distance between clusters is calculated and, consequently, how the clusters are split.

• One challenge with the divisive approach is that its implementation tends to scale poorly with the number of data points, n. This issue is related to the well-known max-cut problem, which is defined as follows: given an  $n \times n$  matrix of positive costs  $c_{ij}$  for  $i, j \in \{1, \ldots, n\}$ , partition the index set  $I = \{1, \ldots, n\}$  into two subsets J and K such that the total cost across the sets,  $\sum_{j \in J} \sum_{k \in K} c_{jk}$ , is maximal. Solving this problem exactly is computationally expensive, but approximate solutions can be found using heuristic methods such as the cross-entropy algorithm.

# 5.7 Density-based Spatial Clustering of Applications with Noise (DBSCAN)

- This is a simple data clustering algorithm proposed by Martin Ester, Hans-Peter Kriegel, Jörg Sander and Xiaowei Xu in 1996 (see Reference (B)). It is one of the most cited articles in scientific literature. The algorithm was awarded the test of time award in 2014 at the leading data mining conference, ACM SIGKDD.
- Let  $D = {\mathbf{x}_1, \ldots, \mathbf{x}_n}$  is the set of given unlabelled *d*-dimensional data points. The goal is again to divide the data into a set of groups (or clusters) so that the points belong to the same group are more similar to each other than the samples belong to different groups.
- This algorithm takes two parameters: radius parameter r that specifies a neighbourhood around a point, and a parameter  $n_{min}$  to determine a point is *core point* or not.

#### • Definitions:

- Core point: A point  $\mathbf{x} \in D$  is called core point if there are  $n_{min}$  in D that are within r distance from  $\mathbf{x}$  (excluding  $\mathbf{x}$ ).
- Directly reachable: If  $\mathbf{x}$  is a core point, a point  $\mathbf{y}$  is said to be directly reachable from  $\mathbf{x}$  if  $\mathbf{y}$  is within a distance r from  $\mathbf{x}$ . Note that we do not say  $\mathbf{x}$  and  $\mathbf{y}$  are directly reachable even if they are within a distance r from each other if none of them is a core point.
- Reachable: A point  $\mathbf{x}'$  is reachable from a core point  $\mathbf{x}$ , if there is a path of points  $\mathbf{x}_{(1)}, \ldots, \mathbf{x}_{(k)}$  for some k such that  $\mathbf{x}_{(1)} = \mathbf{x}$ ,  $\mathbf{x}_{(k)} = \mathbf{x}'$ , and  $\mathbf{x}_{(i+1)}$  is directly reachable from  $\mathbf{x}_{(i)}$  for all  $i = 1, \ldots, k 1$ . Note that all the points  $\mathbf{x}_{(i)}$  for  $i = 1, \ldots, k 1$  are core points.
- Outliers or Noise: All the points D that are not reachable from any other points are called outliers or noise.
- The number of clusters obtaining by algorithm depends on r and  $n_{min}$ .

```
Data: D, Dist, r, and n_{min}
Result: Clustered data
C \leftarrow 0;
for each \mathbf{x} in D do
    if x is unlabelled then
         Find neighbours N_{\mathbf{x}} of \mathbf{x} that are within r distance;
         if |N_{\mathbf{x}}| < n_{min} then
             label(\mathbf{x}) = Noise;
         else
             C \leftarrow C + 1;
             label(\mathbf{x}) = C;
             for every point \mathbf{x}' reachable from \mathbf{x} do
                 label(\mathbf{x}') = C;
             end
         end
    end
end
return Clustered data ;
```

- The above algorithm sequentially goes through all the points in D.
- If a point **x** is not a core point, it labels **x** as a outlier (or noise).
- On the other hand, if **x** is a core point, it increases the cluster label C by 1 and labels **x** and every point **y** that is reachable from **x** as points of cluster C. In this process, the points that are already labelled as noise can be relabelled as cluster C points.

#### Exercise

Use your choice of programming language to implement the above algorithm.