

## Text Classification

Text Classification is the process of assigning predefined categories or labels to textual data. It is a common task in Natural Language Processing (NLP) with applications such as spam detection, sentiment analysis, topic labeling, and language detection.

task ~~like~~ predefined classes or categories. The goal of text classification is to automatically analyze and assign a label to new text samples.

**Step 1: Data Collection:** The first step is to gather a large dataset of labeled text documents. These documents should be representative of the categories you want to classify. For instance, if you want to perform sentiment analysis, you would need a collection of texts labeled as positive or negative sentiments.

### 1. Data Preprocessing

Text data is unstructured, so preprocessing is necessary to convert it into a usable format.

Common steps include:

- **Tokenization:** Splitting text into words or phrases.
- **Lowercasing:** Converting all text to lowercase to maintain uniformity.
- **Removing Stop Words:** Filtering out common words like "is," "the," "and," which add little value.
- **Stemming/Lemmatization:** Reducing words to their base or root forms (e.g., "running" → "run").
- **Removing Punctuation and Special Characters:** Cleaning up unnecessary symbols.

### 2. Feature Extraction/Representation

Convert text into numerical features using techniques like:

- **Bag of Words (BoW):**
  - Represents text as a frequency distribution of words in a document.
- **TF-IDF (Term Frequency-Inverse Document Frequency):**
  - Weights word importance based on frequency in a document relative to its frequency across all documents.
- **Word Embeddings:**
  - Learn dense vector representations of words (e.g., Word2Vec, GloVe, FastText).
- **Sentence Embeddings:**
  - Represent entire sentences using models like BERT or SentenceTransformers.

### 3. Model Training

Train a classification model using the numerical features. Common algorithms include:

- Naive Bayes: Suitable for text data with BoW or TF-IDF features.
- Support Vector Machine (SVM): Effective for high-dimensional data.
- Logistic Regression: A simple and interpretable baseline.
- Deep Learning Models:
  - Recurrent Neural Networks (RNNs) or Long Short-Term Memory (LSTM) networks for sequential data.
  - Transformers (e.g., BERT, GPT) for context-aware and large-scale text classification.

### 4. Evaluation

Assess the performance of the classifier using metrics like:

- Accuracy: Percentage of correctly classified samples.
- Precision, Recall, F1-score: Evaluate performance for imbalanced datasets.
- Confusion Matrix: Visualize true positives, false positives, true negatives, and false negatives.

### 5. Prediction

Use the trained model to classify new, unseen text into the relevant categories.

## Naive Bayes Model

The Naive Bayes model is a probabilistic machine learning algorithm based on Bayes' Theorem. It is widely used for classification tasks and assumes that the features are conditionally independent given the class label (hence, "naive").

### Bayes' Theorem

Bayes' Theorem calculates the probability of a class  $C$  given a feature vector  $X = (x_1, x_2, \dots, x_n)$ :

$$P(C|X) = \frac{P(X|C) \cdot P(C)}{P(X)}$$

Where:

- $P(C|X)$ : Posterior probability of class  $C$  given the features  $X$ .
- $P(X|C)$ : Likelihood of observing  $X$  given class  $C$ .
- $P(C)$ : Prior probability of class  $C$ .
- $P(X)$ : Evidence or total probability of  $X$  (can often be ignored for classification).

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Annotations for the Bayes' Theorem formula:

- LIKELIHOOD: the probability of "B" being TRUE given that "A" is TRUE
- PRIOR: the probability of "A" being TRUE
- POSTERIOR: the probability of "A" being TRUE given that "B" is TRUE
- marginal: The probability of "B" being TRUE

### Steps in Naive Bayes Classification

#### 1. Calculate Priors:

- Compute  $P(C)$ : The probability of each class in the dataset.  $P(A)$

#### 2. Compute Likelihoods:

- For each feature  $x_i$ , calculate  $P(x_i|C)$  from the training data.

#### 3. Apply Bayes' Theorem:

- Use the formula to calculate  $P(C|X)$  for each class.  $P(A/B)$

#### 4. Make Predictions:

- Assign the class with the highest posterior probability:  $A = \arg \max P(A/B)$

$$C = \arg \max_C P(C|X)$$

## Types of Naive Bayes Classifiers

### 1. Gaussian Naive Bayes:

- Assumes continuous data follows a Gaussian (normal) distribution.
- Likelihood:

$$P(x|C) = \frac{1}{\sqrt{2\pi\sigma_C^2}} \exp\left(-\frac{(x - \mu_C)^2}{2\sigma_C^2}\right)$$

where  $\mu_C$  and  $\sigma_C^2$  are the mean and variance for feature  $x$  in class  $C$ .

### 2. Multinomial Naive Bayes:

- Used for discrete data, e.g., word frequencies in text classification.

### 3. Bernoulli Naive Bayes:

- Used for binary features, e.g., presence or absence of a term in text.

## Advantages of Naive Bayes

- Simple and easy to implement.
- Efficient, even for large datasets.
- Works well with high-dimensional data.
- Robust to irrelevant features.

## Limitations of Naive Bayes

- Assumes conditional independence of features, which might not hold in real-world data.
  - Struggles with continuous data if not modeled properly (e.g., Gaussian assumption may not fit).
  - Sensitive to zero probabilities (can be mitigated with smoothing techniques, such as Laplace smoothing).
- Spam Email Detection:** Classify emails as spam or not spam based on word frequencies.
  - Sentiment Analysis:** Determine whether a text expresses positive or negative sentiment.
  - Topic Categorization:** Assign documents to predefined topics like politics, sports, etc.
  - Language Detection:** Identify the language of a given text.
  - Document Filtering:** Filter or categorize content based on relevance or topic.
  - News Classification:** Classify news articles into categories like technology, politics, etc.
  - Tweet Classification:** Classify tweets into categories such as positive, negative, or neutral.
  - Customer Feedback Analysis:** Classify customer reviews into positive or negative categories.
  - Text Classification for Search Engines:** Categorize web pages based on their content for search engine optimization.
  - Medical Text Classification:** Classify medical records or research papers into specific medical topics or diagnoses.

## Example: Spam Detection

For an email with words  $X = [\text{cheap}, \text{buy}, \text{free}]$ :

- Compute  $P(\text{Spam}|X)$  and  $P(\text{Not Spam}|X)$ .
- Assign the class with the higher posterior probability, classifying the email as **Spam** or **Not Spam**.



applications of  
naive bayes

### 3.6 Classification Algorithms : K Nearest Neighbor

- K-Nearest Neighbour is one of the only Machine Learning algorithms based totally on supervised learning approach.

The K-Nearest Neighbor (KNN) algorithm is a simple, non-parametric, and instance-based machine learning technique used for both classification and regression tasks. It classifies a data point based on the majority class of its nearest neighbors.

KNN doesn't build a model; it memorizes the training data and makes predictions based on it. KNN uses labeled training data to make predictions. K-Nearest Neighbors (KNN) is considered a lazy learning algorithm because it does not learn or build a model during the training phase. Instead, it simply stores the entire training dataset and does all the work at the time of making predictions.

#### Key Concepts

##### 1. Instance-Based Learning:

KNN memorizes the training dataset and performs classification or regression by comparing a new data point to its neighbors.

##### 2. Similarity Measurement:

The distance between data points is calculated using metrics like:

- Euclidean Distance:

① Euclidean Distance



② Manhattan Distance



- Manhattan Distance:

$$(x_1, y_1) \rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$d(p, q) = \sum_{i=1}^n |p_i - q_i|$$

- Cosine Similarity (for high-dimensional data).

##### 3. K Value:

- $k$ : Number of nearest neighbors to consider.

- Choosing the right  $k$  is critical:

- Small  $k$ : Sensitive to noise, may overfit.

- Large  $k$ : Smoothes decision boundaries, may underfit.

##### 4. Majority Voting (Classification):

- For a given data point, assign the class that is most common among the  $k$  nearest neighbors.

##### 5. Averaging (Regression):

- Predict the average value of the  $k$  nearest neighbors for a new data point.

## Working of KNN (Classification Example)

### 1. Input:

- Training dataset with features and labels.
- New data point to classify.

### 2. Steps:

- Compute the distance between the new data point and all points in the training set.
- Select the  $k$  nearest neighbors based on the smallest distances.
- Determine the majority class label among these neighbors.
- Assign the majority label to the new data point.

### 3. Output:

- Predicted class label for the new data point.

Apply K-Nearest Neighbor Algorithm (KNN) on following data. Predict the student result for values Physics = 6 marks, Chemistry = 8 marks. Consider number of neighbours  $K=3$  and Euclidean Distance as distance measure.

[12]

Physics (marks)	Chemistry (marks)	Results
4	3	Fail
6	7	Pass
7	8	Pass
5	5	Fail
8	8	Pass

### Problem Summary

We need to predict whether the student will pass or fail based on:

- Physics = 6 marks
- Chemistry = 8 marks

Given:

- Number of neighbors  $k=3$
- Distance metric: Euclidean Distance

### Step 1: Calculate Euclidean Distance

The Euclidean distance between two points  $(x_1, y_1)$  and  $(x_2, y_2)$  is given by:

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Distances from the new point  $(6, 8)$  to the given data points:

1. Point 1  $(4, 3)$ :

$$d = \sqrt{(6 - 4)^2 + (8 - 3)^2} = \sqrt{2^2 + 5^2} = \sqrt{4 + 25} = \sqrt{29} \approx 5.39$$

2. Point 2  $(6, 7)$ :

$$d = \sqrt{(6 - 6)^2 + (8 - 7)^2} = \sqrt{0^2 + 1^2} = \sqrt{1} = 1.0$$

3. Point 3  $(7, 8)$ :

$$d = \sqrt{(6 - 7)^2 + (8 - 8)^2} = \sqrt{1^2 + 0^2} = \sqrt{1} = 1.0$$

4. Point 4  $(5, 5)$ :

$$d = \sqrt{(6 - 5)^2 + (8 - 5)^2} = \sqrt{1^2 + 3^2} = \sqrt{1 + 9} = \sqrt{10} \approx 3.16$$

5. Point 5  $(8, 8)$ :

$$d = \sqrt{(6 - 8)^2 + (8 - 8)^2} = \sqrt{2^2 + 0^2} = \sqrt{4} = 2.0$$

## Advantages of KNN

### 1. Simplicity:

Easy to understand and implement.

### 2. No Training Phase:

KNN requires no explicit training process, making it computationally efficient during training.

### 3. Versatile:

Works for both classification and regression tasks.

### 4. Adaptable:

Can handle multi-class problems.

## Disadvantages of KNN

### 1. High Computation:

The algorithm requires calculating distances for all training samples, which can be computationally expensive for large datasets.

### 2. Sensitive to Noise:

Noisy data or irrelevant features can affect the performance.

### 3. Curse of Dimensionality:

Performance degrades as the number of dimensions increases due to sparsity in data.

### 4. Memory Usage:

Requires the entire dataset to be stored for predictions.

### Step 2: Sort Distances

Sort the distances to find the  $k=3$  nearest neighbors.

Point	Physics	Chemistry	Result	Distance
Point 2	6	7	Pass	1.0
Point 3	7	8	Pass	1.0
Point 5	8	8	Pass	2.0
Point 4	5	5	Fail	3.16
Point 1	4	3	Fail	5.39

### Step 3: Majority Voting

The three nearest neighbors are:

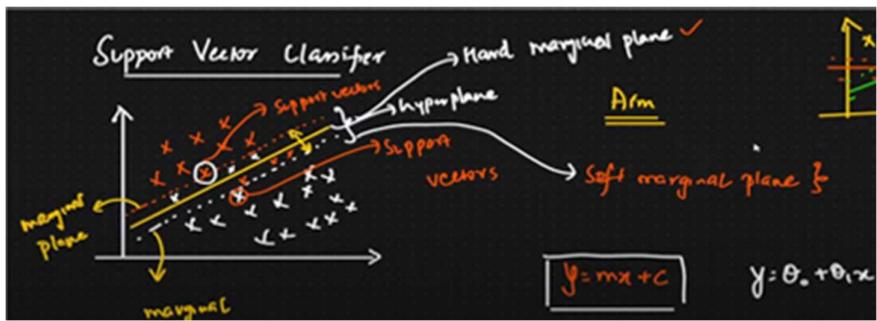
- Point 2 (Pass)
- Point 3 (Pass)
- Point 5 (Pass)

The majority class is Pass.

### Final Prediction

The student will Pass.

CEGPA 3.0  
16 static-238



## Support Vector Classifier (SVC)

Support Vector Classifier (SVC) is a supervised machine learning algorithm used for classification tasks. It is a type of Support Vector Machine (SVM), which aims to find an optimal hyperplane that best separates different classes in the feature space.

### Key Concepts in SVC

- Hyperplane:** A hyperplane is a decision boundary that separates data points of different classes. In a 2D space, it's a line, and in higher dimensions, it's a plane or an n-dimensional hyperplane. The goal of SVC is to find the hyperplane that best separates the classes.
- Margin:** The margin is the distance between the hyperplane and the nearest data points from both classes. SVC aims to maximize this margin, as a larger margin generally leads to better generalization and reduced overfitting.
- Support Vectors:** Support vectors are the data points that are closest to the hyperplane and are critical in defining the position and orientation of the hyperplane. These points are crucial for the classification decision and are used to maximize the margin.

#### 4. Kernel Trick:

For non-linearly separable data, SVM uses kernel functions to transform the input data into a higher-dimensional space where a linear hyperplane can separate the classes. Common kernels include:

- **Linear Kernel:** Suitable for linearly separable data.
- **Polynomial Kernel:** Captures polynomial relationships.
- **Radial Basis Function (RBF) or Gaussian Kernel:** Suitable for non-linear data.
- **Sigmoid Kernel:** Similar to a neural network activation function.

#### 4. Linearly Separable Data:

- If the data can be separated into two classes using a straight line (2D), plane (3D), or hyperplane (higher dimensions), SVM constructs a linear decision boundary (hyperplane) between the classes.

## Mathematical Formulation of SVC

Given a set of training data  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , where:

- $x_i$  is the feature vector (input data),
- $y_i$  is the class label (+1 or -1 for binary classification).

The SVC aims to find a hyperplane defined by:

$$w^T x + b = 0$$

where:

- $w$  is the weight vector normal to the hyperplane,
- $b$  is the bias term.

The objective is to **maximize the margin**, which is defined as  $\frac{2}{\|w\|}$ , while ensuring that all data points are correctly classified. The constraints for the classification are:

$$y_i(w^T x_i + b) \geq 1, \quad \forall i$$

This constraint ensures that each point is correctly classified and lies on the correct side of the margin.

The optimization problem can be formulated as:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

subject to:

$$y_i(w^T x_i + b) \geq 1, \quad \forall i$$

This is a convex optimization problem that can be solved using methods like **Lagrange Multipliers** and **Quadratic Programming**.



## Advantages of SVC

1. **Effective in High Dimensional Spaces:** SVC performs well when the number of dimensions (features) is high, which makes it suitable for complex problems such as text classification and image recognition.
2. **Robust to Overfitting:** Due to the margin maximization, SVC tends to generalize well, especially when the data is not noisy.
3. **Versatile:** By using different kernels, SVC can handle both linearly separable and non-linearly separable data.

## Limitations of SVC

1. **Choice of Kernel:** The performance of SVC depends significantly on the choice of kernel and the parameters (e.g.,  $C, \gamma$  for the RBF kernel). Tuning these parameters can be computationally expensive.
2. **Computational Complexity:** SVC can be slow and memory-intensive, especially for large datasets, since it involves quadratic programming.
3. **Sensitive to Noisy Data:** SVC can be sensitive to noisy data and outliers, which can affect the margin and, consequently, the decision boundary.

## What is the hypothesis used in the vector space model for classification?

In the **Vector Space Model (VSM)** for classification, the **hypothesis** is typically represented as a **decision function** that maps an input vector (representing a data point) to a predicted class label. The hypothesis is designed to classify data points based on their position in the feature space (a vector space), where each dimension corresponds to a feature of the data.

### Hypothesis in the Vector Space Model

The general form of the hypothesis in the vector space model is:

$$h(x) = \text{sign}(w \cdot x + b)$$

Where:

- $h(x)$  is the predicted class label for the input data point  $x$ ,
- $w$  is the weight vector (which represents the direction of the hyperplane),
- $x$  is the feature vector of the data point,
- $b$  is the bias term, and
- $\text{sign}()$  is the activation function that outputs the class label based on the sign of the expression  $w \cdot x + b$ .

### Explanation of the Hypothesis Components

#### 1. Weight Vector ( $w$ ):

- The weight vector  $w$  determines the orientation of the hyperplane in the feature space. It is learned during the training process, and it is a crucial component in determining how the hyperplane separates different classes.

#### 2. Feature Vector ( $x$ ):

- The feature vector  $x$  represents the data point in the vector space, typically with each feature corresponding to a different dimension in the space.

#### 3. Bias ( $b$ ):

- The bias term  $b$  shifts the hyperplane, allowing it to better fit the data. It determines how far the hyperplane is from the origin.

#### 4. Activation Function ( $\text{sign}()$ ):

- The activation function  $\text{sign}(w \cdot x + b)$  produces the predicted class label. If the result of  $w \cdot x + b$  is positive, the class label will be one class (e.g., 1), and if it is negative, the label will be the other class (e.g., -1).

### Application in Classification Models

In the context of classification, the **hypothesis function** is used to separate classes based on their feature vectors. For example, in **Support Vector Machines (SVM)**, the hypothesis is based on the decision boundary (a hyperplane) that maximizes the margin between the classes. The data points closest to the hyperplane are called support vectors, and they play a critical role in the classification process.

#### Example:

For a binary classification problem, the hypothesis can be written as:

$$h(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ -1 & \text{if } w \cdot x + b < 0 \end{cases}$$

Here:

- If the dot product of the weight vector  $w$  and the feature vector  $x$ , plus the bias  $b$ , is greater than zero, the data point is classified as Class 1.

## Vector Space Classification Using Hyperplanes

Vector Space Classification using hyperplanes is a method of classifying data points in a vector space by separating them using hyperplanes. A **hyperplane** is a decision boundary used to separate different classes in a feature space. This approach is used in algorithms like **Support Vector Machines (SVM)**.

### Concept of a Hyperplane

A **hyperplane** is a flat affine subspace of one dimension less than the space it is embedded in:

- In 2D, a hyperplane is a line that divides the space into two halves.
- In 3D, it is a plane.
- In higher dimensions, it's called a hyperplane.

Given a dataset in a vector space, the goal is to find the best hyperplane that separates the data points belonging to different classes.

### Mathematical Representation of a Hyperplane

In  $n$ -dimensional space, a hyperplane can be mathematically represented by the following equation:

$$w \cdot x + b = 0$$

Where:

- $w$  is the weight vector (normal to the hyperplane),
- $x$  is a vector representing a data point,
- $b$  is the bias term (offset of the hyperplane from the origin),
- $\cdot$  represents the dot product.

The vector  $w$  is perpendicular to the hyperplane, and the value of  $b$  determines the distance of the hyperplane from the origin.

### Classification with Hyperplanes

In classification problems, the goal is to find a hyperplane that best separates data points belonging to different classes. **Support Vector Machines (SVMs)** are one of the most common algorithms that use hyperplanes for classification. The key idea is to find the hyperplane that **maximizes the margin** between the two classes.

#### Steps for Vector Space Classification using Hyperplanes:

##### 1. Represent Data as Vectors:

- Represent each data point as a vector in an  $n$ -dimensional space. For example, if each data point has  $n$  features, it will be represented as an  $n$ -dimensional vector.

##### 2. Find the Optimal Hyperplane:

- The goal is to find the hyperplane that separates the classes with the **maximum margin**.  
The margin is the distance between the closest points (support vectors) of each class to the hyperplane.

##### 3. Construct the Hyperplane:

- The equation of the optimal hyperplane is defined as:

$$w \cdot x + b = 0$$

Here,  $w$  is the vector normal to the hyperplane, and the hyperplane is positioned such that the margin between the data points and the hyperplane is maximized.

##### 4. Classify New Data Points:

- Once the hyperplane is determined, classify new data points by calculating the value of  $w \cdot x + b$  for each new point:
  - If  $w \cdot x + b > 0$ , classify the point as **Class 1**.
  - If  $w \cdot x + b < 0$ , classify the point as **Class -1**.

## Kernel Function

In machine learning and statistics, kernel functions are used to transform data into a higher-dimensional space to make it easier to find patterns or relationships, particularly for algorithms like Support Vector Machines (SVMs) and Kernel Principal Component Analysis (Kernel PCA).

A kernel is a function that computes the inner product (dot product) of two vectors in some high-dimensional feature space without explicitly mapping the data to that space, which can save on computational cost. This is often referred to as the "kernel trick." The kernel trick allows algorithms that rely on inner products, such as SVMs, to work in higher-dimensional spaces without needing to calculate the coordinates of data points in that space directly.

The Linear Kernel is one of the simplest and most commonly used kernel functions in machine learning, particularly for algorithms like Support Vector Machines (SVMs) and Kernel Principal Component Analysis (Kernel PCA). It is used when the data is already linearly separable, meaning that a linear decision boundary (hyperplane) can effectively separate the different classes in the feature space.

### Formula:

The Linear Kernel function computes the dot product of two input vectors  $x$  and  $y$ , which are typically the feature vectors of data points.

$$K(x, y) = x^T y$$

Where:

- $x$  and  $y$  are the input vectors (data points).
- $x^T y$  represents the dot product of the two vectors.

### How It Works:

- The kernel measures the similarity between two data points by taking the dot product of their feature vectors.
- If the data points are closer together in the feature space, the dot product will be larger, indicating higher similarity.
- The result is a scalar value that represents how similar the two vectors are. A higher value indicates that the vectors (data points) are more similar.

The Radial Basis Function (RBF) Kernel, also known as the Gaussian Kernel, is one of the most popular and widely used kernel functions in machine learning, particularly in Support Vector Machines (SVMs). It is especially useful for dealing with non-linearly separable data.

### Formula:

The RBF Kernel measures the similarity between two data points  $x$  and  $y$  by calculating the exponential of the negative squared Euclidean distance between them, scaled by a parameter  $\sigma$ .

$$K(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

Where:

- $x$  and  $y$  are two input vectors (data points).
- $\|x - y\|^2$  is the squared Euclidean distance between the vectors  $x$  and  $y$ .
- $\sigma$  is a scale parameter (also called bandwidth or width), which controls the spread of the Gaussian function.

The Polynomial Kernel is another popular kernel function used in machine learning algorithms like Support Vector Machines (SVMs) and Kernel Principal Component Analysis (Kernel PCA). It is useful when the relationship between the features of the data is not linear but can be modeled by a polynomial function.

### Formula:

The Polynomial Kernel function computes the similarity between two data points  $x$  and  $y$  by taking the dot product of the input vectors and raising it to a given power  $d$ , with an optional constant  $c$  added to the dot product before raising it to the power.

$$K(x, y) = (x^T y + c)^d$$

Where:

- $x$  and  $y$  are input vectors (data points).
- $x^T y$  is the dot product of  $x$  and  $y$ .
- $c$  is a constant (usually  $c \geq 0$ ) that adds bias to the kernel function.
- $d$  is the degree of the polynomial, which controls the complexity of the decision boundary.

The parameter  $d$  controls the degree of the polynomial, which determines the complexity of the decision boundary. A higher  $d$  increases the flexibility of the decision boundary, while a lower  $d$  results in a simpler, less complex boundary.

The **Sigmoid Kernel**, also known as the **Hyperbolic Tangent Kernel**, is another type of kernel function used in machine learning algorithms such as **Support Vector Machines (SVMs)** and **Neural Networks**. It is based on the **sigmoid function**, which is commonly used in the activation functions of artificial neural networks.

#### Formula:

The formula for the Sigmoid Kernel is given by:

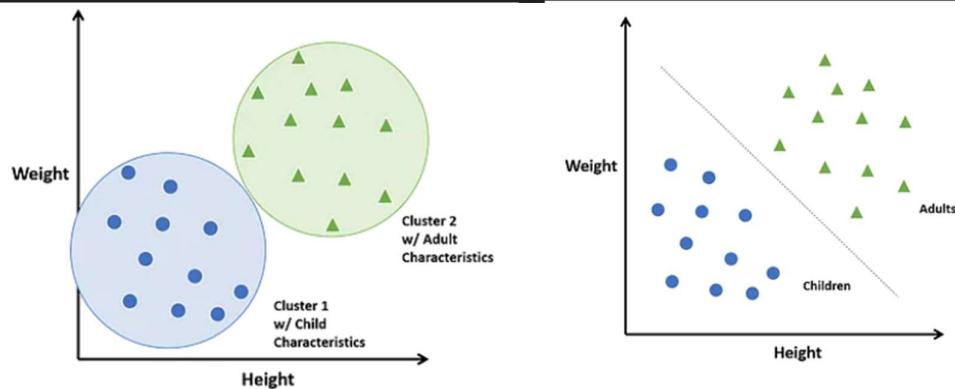
$$K(x, y) = \tanh(\alpha x^T y + c)$$

Where:

- $x$  and  $y$  are input vectors (data points).
- $x^T y$  is the dot product of  $x$  and  $y$ , i.e., the similarity between the data points.
- $\alpha$  is a scaling parameter that controls the steepness of the sigmoid curve.
- $c$  is a constant that shifts the curve along the horizontal axis (also known as the **bias term**).

What is the difference between clustering and classification? Can clustering be used for classification purposes?

Aspect	Clustering	Classification
Definition	Unsupervised learning method that groups data into clusters based on similarity.	Supervised learning method that assigns predefined labels to input data based on training.
Data Requirement	Does not require labeled data for learning.	Requires labeled data for training.
Output	Produces clusters with similar data points grouped together.	Produces class labels for the given input data.
Goal	To discover underlying patterns or structures in data.	To predict the category or label of new data based on prior knowledge.
Nature of Learning	Unsupervised (no prior knowledge of data labels).	Supervised (requires prior knowledge of data labels).
Algorithms Used	K-Means, DBSCAN, Agglomerative Clustering, Mean Shift.	Logistic Regression, Decision Trees, Random Forest, SVM, Neural Networks.
Evaluation Metrics	Silhouette score, Davies-Bouldin index, or manual validation of clusters.	Accuracy, precision, recall, F1-score, confusion matrix.
Applications	Market segmentation, anomaly detection, customer profiling, image compression.	Spam filtering, sentiment analysis, disease diagnosis, credit risk prediction.
Human Intervention	May require domain expertise to interpret and label clusters post-hoc.	Labels are predefined, so minimal interpretation is needed during model training.
Adaptability	Useful when no prior knowledge of data categories exists; focuses on data exploration.	Requires prior knowledge; focuses on assigning data to predefined categories.



## **Can Clustering Be Used for Classification?**

Yes, clustering can be used as a **preprocessing step** or in a hybrid approach for classification purposes. Here's how:

### **1. Label Generation:**

- If labeled data is unavailable, clustering can generate pseudo-labels by grouping similar data points. These labels can then be used to train a classification model.

### **2. Feature Engineering:**

- Clustering can help identify patterns or groupings in data that can be used as features for classification models.

### **3. Anomaly Detection:**

- Clustering models can identify outliers, which might correspond to specific labels (e.g., anomalies vs. normal behavior).

### **4. Semi-Supervised Learning:**

- Clustering can combine with labeled data to enhance learning when only a small portion of the data is labeled.

However, clustering is inherently unsupervised, so its use for classification relies on additional steps to interpret and convert the clusters into meaningful labels.

What are different types of clustering algorithm? Explain any one of them.

### K-Means Clustering Algorithm

The K-Means algorithm is one of the most popular and widely used clustering techniques. It is a **centroid-based clustering method** that aims to partition the dataset into  $k$  clusters, where  $k$  is a predefined number. Each cluster is represented by its centroid, which is the average of all the points in that cluster.

K-means clustering is a **distance-based unsupervised clustering algorithm** where data points that are close to each other are grouped in a given number of clusters/groups

### Steps of the K-Means Algorithm

#### 1. Initialization:

- Select the number of clusters ( $k$ ).
- Randomly initialize  $k$  cluster centroids or use an advanced initialization technique like  $k$ -Means++.

#### 2. Assignment Step:

- For each data point, calculate the distance to all  $k$  centroids.
- Assign each data point to the cluster whose centroid is closest (typically using Euclidean distance).

#### 3. Update Step:

- Recalculate the centroids of the clusters by taking the mean of all data points assigned to each cluster.

#### 4. Repeat:

- Alternate between the assignment and update steps until:
  - The centroids no longer change significantly, or
  - A maximum number of iterations is reached.

#### 5. Output:

- The algorithm outputs  $k$  clusters and their respective centroids.

### Key Features of K-Means:

- It requires the number of clusters  $k$  as input.
- The objective is to minimize the **intra-cluster sum of squares (WCSS)**:

$$WCSS = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

where  $C_i$  is the  $i$ -th cluster and  $\mu_i$  is its centroid.

## Advantages

1. **Simplicity:**
  - The algorithm is easy to understand and implement.
2. **Scalability:**
  - Computationally efficient for large datasets with relatively low-dimensional features.
3. **Speed:**
  - The iterative approach converges quickly for well-behaved datasets.
4. **Wide Applicability:**
  - Effective for linearly separable and well-defined spherical clusters.

## Disadvantages

1. **Predefined  $k$ :**
  - The number of clusters ( $k$ ) must be specified in advance, which can be challenging when the optimal value is unknown.
2. **Sensitive to Initialization:**
  - Poor initialization of centroids can lead to suboptimal results or convergence to local minima.  $k$ -Means++ can address this issue.
3. **Limited to Spherical Clusters:**
  - Struggles with datasets containing clusters of irregular shapes or varying densities.
4. **Impact of Outliers:**
  - Outliers can distort the position of centroids and degrade clustering quality.
5. **Curse of Dimensionality:**
  - Performance may degrade in high-dimensional spaces.

## Applications

1. **Customer Segmentation:**
  - Group customers based on purchasing behavior or demographics for targeted marketing.
2. **Image Compression:**
  - Reduce the number of colors in an image by clustering pixel values.
3. **Document Clustering:**
  - Categorize textual documents based on content similarity.
4. **Market Segmentation:**
  - Identify distinct consumer groups for strategic business planning.
5. **Anomaly Detection:**
  - Detect data points that do not fit into any cluster as potential anomalies.

## Agglomerative Hierarchical Clustering (AHC) Algorithm

Agglomerative Hierarchical Clustering (AHC) is a **bottom-up approach** to hierarchical clustering where each data point starts as its own cluster, and pairs of clusters are merged iteratively based on similarity until a single cluster or a desired number of clusters is formed.

### Algorithm Steps

1. Initialize Clusters:
  - Treat each data point as an individual cluster.
  - If there are  $n$  data points, start with  $n$  clusters.
2. Compute Similarity/Distance Matrix:
  - Calculate the pairwise distances (e.g., Euclidean, Manhattan) between all clusters.
  - Store these distances in a similarity/distance matrix.
3. Merge Closest Clusters:
  - Identify the two clusters with the smallest distance in the matrix.
  - Merge these two clusters into a single cluster.
4. Update the Distance Matrix:
  - Recalculate the distances between the new cluster and all other clusters based on the chosen linkage criterion:
    - **Single Linkage:** Distance between the closest pair of points in two clusters.
    - **Complete Linkage:** Distance between the farthest pair of points in two clusters.
    - **Average Linkage:** Average distance between all pairs of points in two clusters.
    - **Centroid Linkage:** Distance between the centroids of two clusters.
5. Repeat Steps 3 and 4:
  - Continue merging clusters and updating the distance matrix until all data points are in one cluster or a predefined number of clusters is reached.
6. Construct a Dendrogram:
  - Visualize the hierarchy of clusters using a dendrogram, where each merge is represented as a branch.

### Example of AHC

Let's consider 4 data points:  $A, B, C, D$ .

1. Initially, each point is a separate cluster:  
 $\{A\}, \{B\}, \{C\}, \{D\}$ .
2. Compute the pairwise distances and merge the closest pair, say  $\{A\}$  and  $\{B\}$ .
3. Update the distance matrix for the new cluster  $\{A, B\}$  and merge the next closest pair.
4. Repeat until all points are merged.

## Advantages

1. No Predefined Clusters:
  - AHC doesn't require specifying the number of clusters beforehand (unlike K-Means).
2. Dendrogram Interpretation:
  - The dendrogram provides a visual representation of the clustering process and relationships between data points.
3. Versatility:
  - Works with different distance metrics and linkage criteria, making it adaptable to various data types.
4. Hierarchical Insight:
  - Reveals a nested structure of clusters, which can be useful for exploratory analysis.

## Disadvantages

1. Scalability:
  - Computationally expensive for large datasets:
    - Time complexity:  $O(n^3)$
    - Space complexity:  $O(n^2)$
2. Sensitivity to Noise:
  - Outliers or noisy data can significantly affect cluster formation.
3. Non-Optimal Merges:
  - Once two clusters are merged, the process cannot be undone, leading to potential suboptimal results.
4. Interpretability for Large Data:
  - Dendograms become difficult to interpret when the dataset is very large.

## Mixture of Gaussians Model

### Gaussian Mixture Model (GMM)

The Gaussian Mixture Model (GMM) is a probabilistic model that assumes data points are generated from a mixture of several Gaussian distributions. Each Gaussian distribution represents a cluster, and the data is modeled as a combination of these distributions, with each cluster having its own mean and variance.

#### Key Concepts:

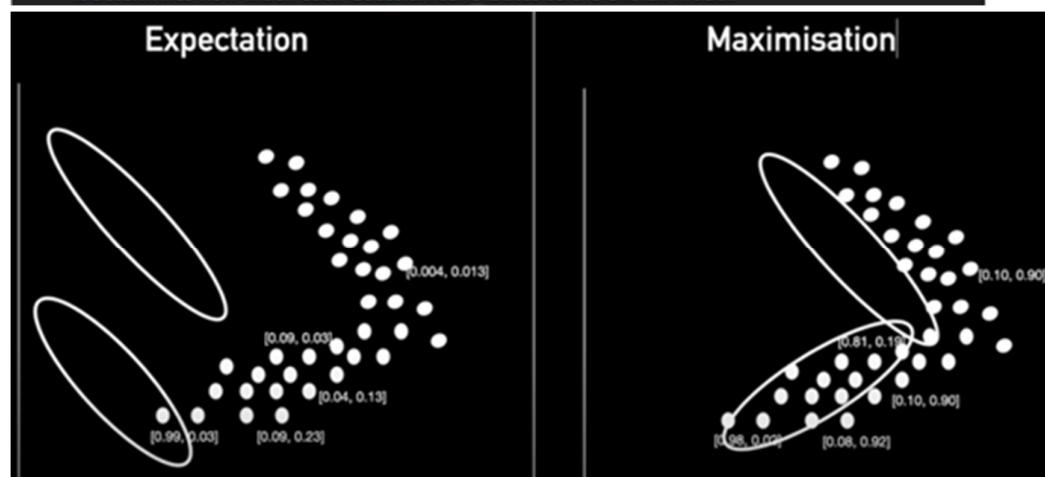
1. **Gaussian Distribution:** A normal distribution, characterized by its mean (average) and variance (spread). In a GMM, each cluster is assumed to follow a Gaussian distribution.
2. **Mixture Model:** A model that assumes the data is generated from a mixture of several different distributions (in this case, Gaussian distributions).

#### How GMM Works:

- **Initialization:** The number of clusters (components)  $K$  is specified, and the initial parameters of the Gaussian distributions (mean, covariance, and weight) are estimated. This can be done using methods like k-means or random initialization.
- **Expectation-Maximization (EM) Algorithm:**
  1. **Expectation Step (E-step):** Given the current parameters, compute the probability that each data point belongs to each cluster (Gaussian distribution). This is done using Bayes' theorem, and the result is the "responsibility" each cluster has for each data point.
  2. **Maximization Step (M-step):** Update the parameters (mean, covariance, and mixture weights) of the Gaussian distributions based on the responsibilities computed in the E-step. The mean is updated as a weighted average of the data points, the covariance matrix is updated based on the spread of the points, and the weights are updated based on the probability that data points belong to each cluster.
  3. Repeat the E-step and M-step until convergence, meaning the parameters of the Gaussian distributions no longer change significantly.

#### Parameters of GMM:

- **Means ( $\mu_k$ ):** The center of each Gaussian distribution (the average value).
- **Covariances ( $\Sigma_k$ ):** The spread or shape of each Gaussian distribution, which controls the size and orientation of the cluster.
- **Mixing Coefficients ( $\pi_k$ ):** The weight or proportion of each Gaussian distribution in the mixture. It determines how much each Gaussian contributes to the overall model.



## Expectation-Maximization (EM) Algorithm in Clustering

The Expectation-Maximization (EM) algorithm is a **model-based clustering** method that assumes data points are generated by a mixture of probability distributions, such as Gaussian distributions. It aims to assign data points to clusters based on probabilistic models.

### Key Steps of EM Algorithm

#### 1. Initialization:

- Assume the number of clusters  $k$ .
- Initialize the parameters of the probability distributions (e.g., mean, variance, and mixing coefficients for Gaussian Mixture Models).

#### 2. Expectation Step (E-step):

- Calculate the probability (or responsibility) that each data point belongs to each cluster.
- Formula for responsibility:

$$r_{ij} = \frac{\pi_j \cdot \mathcal{N}(x_i | \mu_j, \Sigma_j)}{\sum_{l=1}^k \pi_l \cdot \mathcal{N}(x_i | \mu_l, \Sigma_l)}$$

where:

$r_{ij}$ : Responsibility of cluster  $j$  for data point  $i$ .

$\pi_j$ : Mixing coefficient (prior probability) of cluster  $j$ .

$\mathcal{N}(x_i | \mu_j, \Sigma_j)$ : Gaussian probability density function of cluster  $j$  for data point  $x_i$ .

#### 3. Maximization Step (M-step):

- Update the parameters of the distributions ( $\mu_j, \Sigma_j, \pi_j$ ) based on the responsibilities computed in the E-step.
- Update formulas:
  - Mean:

$$\mu_j = \frac{\sum_{i=1}^n r_{ij} \cdot x_i}{\sum_{i=1}^n r_{ij}}$$

#### Covariance:

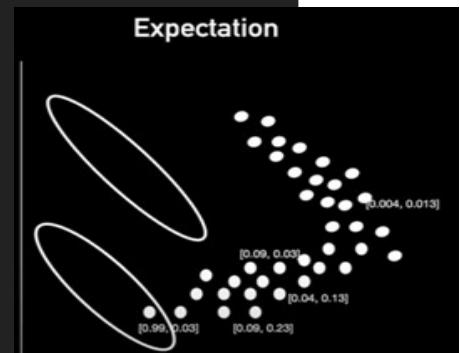
$$\Sigma_j = \frac{\sum_{i=1}^n r_{ij} \cdot (x_i - \mu_j)(x_i - \mu_j)^T}{\sum_{i=1}^n r_{ij}}$$

#### Mixing Coefficient:

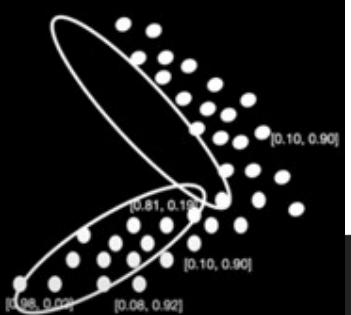
$$\pi_j = \frac{\sum_{i=1}^n r_{ij}}{n}$$

#### 4. Repeat:

- Alternate between E-step and M-step until convergence (when changes in parameters or log-likelihood become negligible).



Expectation



Maximisation

### Advantages of EM in Clustering

- Handles overlapping clusters effectively.
- Provides a probabilistic measure of cluster membership.
- Works with flexible cluster shapes (e.g., ellipsoidal).

### Limitations

- Sensitive to initialization of parameters.
- Can converge to local optima.
- Assumes underlying data distribution (e.g., Gaussian), which may not always be valid.
- Computationally expensive for large datasets.

## Types of Agglomerative Hierarchical Clustering

Agglomerative Hierarchical Clustering (AHC) is a bottom-up approach where each data point starts as its own cluster, and pairs of clusters are merged iteratively based on a certain distance metric. The primary aim is to build a hierarchy of clusters in a tree-like structure known as a **dendrogram**.

There are several methods used to determine how clusters are merged at each step. These methods are based on different ways of calculating the **distance between clusters**. The main types of Agglomerative Hierarchical Clustering are:

### 1. Single Linkage (Nearest Point Linkage)

In Single Linkage, the distance between two clusters is defined as the **shortest distance** between any single point in one cluster and any single point in the other cluster.

- **Formula:**

$$D(C_1, C_2) = \min_{x \in C_1, y \in C_2} d(x, y)$$

where  $C_1$  and  $C_2$  are two clusters, and  $d(x, y)$  is the distance between points  $x$  and  $y$ .

- **Characteristics:**

- Tends to create "long, chain-like" clusters.
- Sensitive to noise and outliers, as it relies on the shortest distance.

### 2. Complete Linkage (Farthest Point Linkage)

In Complete Linkage, the distance between two clusters is defined as the **maximum distance** between any point in one cluster and any point in the other cluster.

- **Formula:**

$$D(C_1, C_2) = \max_{x \in C_1, y \in C_2} d(x, y)$$

- **Characteristics:**

- Tends to produce more compact, spherical clusters.
- Less sensitive to outliers compared to Single Linkage.
- More strict in terms of cluster similarity.

### 3. Average Linkage (Group Average Linkage)

In Average Linkage, the distance between two clusters is defined as the **average distance** between all pairs of points, where one point is taken from each cluster.

- **Formula:**

$$D(C_1, C_2) = \frac{1}{|C_1| \times |C_2|} \sum_{x \in C_1, y \in C_2} d(x, y)$$

where  $|C_1|$  and  $|C_2|$  are the sizes of clusters  $C_1$  and  $C_2$ , respectively.

- **Characteristics:**

- Balances the characteristics of Single and Complete Linkage.
- Creates more balanced clusters in terms of size and shape.

### 5. Centroid Linkage

In Centroid Linkage, the distance between two clusters is calculated based on the **distance between their centroids** (the center of mass of the points in each cluster).

- **Formula:**

$$D(C_1, C_2) = d(\mu_1, \mu_2)$$

where  $\mu_1$  and  $\mu_2$  are the centroids of clusters  $C_1$  and  $C_2$ , and  $d(\mu_1, \mu_2)$  is the distance between them.

- **Characteristics:**

- Typically forms clusters with more balanced shapes.
- Can occasionally lead to counter-intuitive merges, especially in the presence of outliers.

Solve Agglomerative hierarchical clustering for single link with example

A (1, 2), B (2, 3), C (6, 5), D (8, 8)

dataset

#### Step 1: Initialize Clusters

We treat each data point as its own cluster:

- Clusters: {A}, {B}, {C}, {D}

#### Step 2: Calculate Pairwise Distances

We calculate the Euclidean distance between each pair of points:

- Distance between A and B:

$$d(A, B) = \sqrt{(2 - 1)^2 + (3 - 2)^2} = \sqrt{1^2 + 1^2} = \sqrt{2} \approx 1.414$$

- Distance between A and C:

$$d(A, C) = \sqrt{(6 - 1)^2 + (5 - 2)^2} = \sqrt{5^2 + 3^2} = \sqrt{25 + 9} = \sqrt{34} \approx 5.831$$

- Distance between A and D:

$$d(A, D) = \sqrt{(8 - 1)^2 + (8 - 2)^2} = \sqrt{7^2 + 6^2} = \sqrt{49 + 36} = \sqrt{85} \approx 9.220$$

- Distance between B and C:

$$d(B, C) = \sqrt{(6 - 2)^2 + (5 - 3)^2} = \sqrt{4^2 + 2^2} = \sqrt{16 + 4} = \sqrt{20} \approx 4.472$$

- Distance between B and D:

$$d(B, D) = \sqrt{(8 - 2)^2 + (8 - 3)^2} = \sqrt{6^2 + 5^2} = \sqrt{36 + 25} = \sqrt{61} \approx 7.810$$

- Distance between C and D:

$$d(C, D) = \sqrt{(8 - 6)^2 + (8 - 5)^2} = \sqrt{2^2 + 3^2} = \sqrt{4 + 9} = \sqrt{13} \approx 3.606$$

Pair	Distance
A, B	1.414
A, C	5.831
A, D	9.220
B, C	4.472
B, D	7.810
C, D	3.606

### Step 3: Merge Closest Clusters

The closest pair is A and B with a distance of 1.414. We merge them into a new cluster {A, B}.

- New clusters: {A, B}, {C}, {D}

### Step 4: Update Distance Matrix

Now, we update the distance matrix by recalculating the distance between the new cluster {A, B} and the remaining clusters {C} and {D}.

For Single Link, the distance between two clusters is the minimum distance between their points.

- Distance between {A, B} and C: minimum of  $d(A, C)$  and  $d(B, C) \rightarrow \min(5.831, 4.472) = 4.472$
- Distance between {A, B} and D: minimum of  $d(A, D)$  and  $d(B, D) \rightarrow \min(9.220, 7.810) = 7.810$

Now, the updated distance matrix is:

Pair	Distance
{A, B}, C	4.472
{A, B}, D	7.810
C, D	3.606

### Step 5: Merge Closest Clusters

The closest pair is C and D with a distance of 3.606. We merge C and D into a new cluster {C, D}.

- New clusters: {A, B}, {C, D}

### Step 6: Update Distance Matrix

Now, we update the distance matrix by recalculating the distance between the two remaining clusters {A, B} and {C, D}.

- Distance between {A, B} and {C, D}: minimum of the distances between all points in {A, B} and {C, D}.
  - Minimum of  $d(A, C)$ ,  $d(A, D)$ ,  $d(B, C)$ , and  $d(B, D)$ :

$$\min(5.831, 9.220, 4.472, 7.810) = 4.472$$

Now, the final distance matrix is:

Pair	Distance
{A, B}, {C, D}	4.472

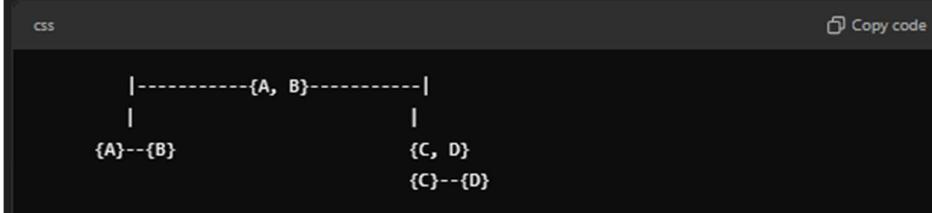
## Step 7: Merge the Final Clusters

The closest pair is {A, B} and {C, D} with a distance of 4.472. We merge them into the final cluster {A, B, C, D}.

- Final cluster: {A, B, C, D}

## Dendrogram for the Example

Here's the dendrogram representation of the clustering process:



- First merge: A and B at a distance of 1.414.
- Second merge: C and D at a distance of 3.606.
- Final merge: {A, B} and {C, D} at a distance of 4.472.

## Spam Filtering

Spam filtering is the process of identifying and blocking unwanted or malicious messages (spam), typically in the context of email communication. The goal of spam filtering is to prevent spam messages from reaching the inbox of users and to ensure that only legitimate, wanted messages are delivered.

Spam filters use a variety of techniques, often based on machine learning and natural language processing, to classify messages as either spam or non-spam (ham). This classification process relies on analyzing the content, metadata, and other characteristics of the messages.

## Types of Spam Filters

1. **Content-based Filters:** These filters analyze the actual content of the email, such as the subject, body, and attachments. They look for certain keywords, phrases, or patterns commonly found in spam messages, such as:
  - Use of excessive capital letters or exclamation marks.
  - Certain phrases like "Buy now!", "Free gift!", or "Congratulations!"
  - Suspicious links or attachments.
2. **Blacklist-based Filters:** These filters block emails from known malicious sources. A **blacklist** is a list of IP addresses, domains, or email addresses known to send spam. If an incoming email matches an entry on the blacklist, it is flagged as spam.
3. **Bayesian Filters:** Bayesian spam filters use probabilistic methods to classify messages. They learn from past examples of spam and non-spam emails. Based on the frequency of words appearing in these messages, they calculate the probability that a new email is spam. For example, words that frequently appear in spam emails are assigned higher probabilities of being spam.
4. **Heuristic Filters:** These filters use predefined rules or heuristics to identify spam. They may rely on things like:
  - The number of links in the email.
  - The presence of certain words in the subject or body.
  - The email's header information, such as the sender's address or the country of origin.
6. **Header-based Filters:** These filters analyze the header information of the email, such as the sender's email address, the reply-to address, and the originating IP address. If the header information looks suspicious or matches known spam patterns, the message may be flagged as spam.
7. **Machine Learning-based Filters:** Advanced spam filters use machine learning (ML) algorithms to classify emails based on features extracted from both the content and metadata. These filters are trained on labeled datasets (spam vs. non-spam) and continuously improve over time. Some common ML algorithms used for spam filtering include:
  - **Naive Bayes Classifier:** A probabilistic algorithm based on Bayes' theorem. It is often used in content-based filtering.
  - **Support Vector Machines (SVM):** A supervised learning algorithm that finds a hyperplane to separate spam from non-spam messages.
  - **Decision Trees:** A tree-like model that is used to classify emails based on certain features.

## Spam Filtering Process

1. **Feature Extraction:** The first step is to extract relevant features from the email. These features could be:
  - Word frequencies (e.g., "free", "win", "prize").
  - Email metadata (e.g., sender, subject).
  - Links and attachments (e.g., suspicious URLs or file types).
  - Text patterns (e.g., capital letters, long subject lines).
2. **Classification:** Once the features are extracted, the spam filtering model classifies the email based on those features. This step involves using one of the filtering techniques mentioned earlier (Bayesian, machine learning, etc.) to determine whether the email is spam or not.
3. **Action:** If the email is classified as spam, it may be:
  - Moved to a spam folder.
  - Deleted.
  - Marked with a warning label (e.g., "This message may be spam"). If the email is classified as non-spam (ham), it is delivered to the inbox.