Lecture 6

Ensemble Methods
Boosting & Bagging
Naïve Bayes

**Ensemble methods** are a class of machine learning techniques used in data mining to improve the predictive performance of models by combining the outputs of multiple base models (learners). These base models can be of the same type or different types. The idea behind ensemble methods is to harness the wisdom of the crowd, aggregating the opinions of multiple models to make more accurate and robust predictions. Some commonly used ensemble methods in data mining include:

**Bagging (Bootstrap Aggregating)**

Bagging is a simple ensemble method that involves creating multiple bootstrapped (randomly sampled with replacement) datasets from the original data and training a separate base model on each dataset. The final prediction is typically determined by taking a majority vote for classification problems or averaging the predictions for regression problems.

*Random Forest*, a popular algorithm, is an extension of bagging that applies this approach to decision trees, resulting in a more robust and accurate model.

**Boosting**:

Boosting is a family of ensemble methods that aim to improve the performance of weak learners (models slightly better than random guessing) by iteratively training and combining them.

*AdaBoost (Adaptive Boosting)* is a well-known boosting algorithm that assigns weights to misclassified instances and focuses on difficult examples in each iteration.

***Gradient Boosting*** techniques, such as *Gradient Boosting Machines (GBM)* and *XGBoost*, build models sequentially, minimizing the errors of the previous models.

***Stacking (Stacked Generalization):*** Stacking combines the predictions of multiple base models by training a meta-model on their outputs. The base models make predictions on the dataset, and their predictions are used as input features for the meta-model, which generates the final prediction. Stacking can be used to capture complementary information from different base models.

***Voting and Averaging***: Simple ensemble methods include majority voting and averaging of base model predictions. In majority voting, the most frequently predicted class by the base models is selected for classification problems. In averaging, the predictions of the base models are averaged to produce the final prediction for regression problems.

***Random Subspace Method (Feature Bagging):*** Similar to bagging, this technique applies randomization to feature selection. It trains multiple base models, each on a random subset of features. Feature bagging can help reduce overfitting and increase the diversity of the base models.

**Random Patches**: Random patches combine both bootstrapping and feature bagging. In this approach, random subsets of data and random subsets of features are used to train individual base models. This technique is often used with ensemble methods like Random Forest.

**Blending**: Blending is similar to stacking but involves partitioning the training data into different subsets. Each subset is used to train a distinct base model. The base models' predictions are then combined, typically with a simple model like linear regression.

Ensemble methods are known for their ability to enhance the predictive accuracy and generalization performance of machine learning models. They are widely used in various domains, including classification, regression, and feature selection. The choice of which ensemble method to use often depends on the characteristics of the problem, the base models available, and the desired level of model diversity and performance.

Bagging and boosting are two popular ensemble methods used in machine learning to improve the performance of models by combining multiple base models. While they share the goal of increasing predictive accuracy, they differ in several key aspects, despite sometimes being confused.

**Bagging (Bootstrap Aggregating):**
*Base Model Diversity*: Bagging focuses on reducing the variance of the base models. It aims to create multiple base models with similar predictive abilities. Each base model is trained independently on a bootstrapped (randomly sampled with replacement) subset of the training data. The predictions of individual base models are typically combined through majority voting (classification) or averaging (regression).

*Sequential vs. Parallel*: Bagging base models are trained independently and in parallel. There is no sequential adjustment of weights or focus on misclassified instances.

*Robustness to Overfitting*: Bagging helps reduce overfitting by combining the predictions of multiple base models. Since it averages or votes on the outputs, it tends to provide a more stable and less noisy prediction.

*Examples*: Random Forest is a well-known ensemble method based on bagging. It applies bagging to decision trees, creating multiple decision trees with the goal of reducing variance and improving generalization.

**Boosting**:
*Base Model Diversity*: Boosting focuses on improving the accuracy of weak learners (models that perform slightly better than random guessing). It aims to combine base models in a sequential manner, with each new model addressing the mistakes of the previous ones. The weights of data points are adjusted in each iteration to emphasize incorrectly predicted instances, making boosting particularly useful for difficult examples.

*Sequential Process*: Boosting is an iterative and sequential process. Base models are trained sequentially, and the subsequent models give more weight to instances that were misclassified by previous models. This leads to a strong focus on challenging data points.

*Handling Overfitting*: Boosting can be more prone to overfitting compared to bagging, as it aggressively adjusts the model to minimize errors on the training data. Techniques like early stopping and model complexity control are often used to mitigate overfitting.

*Examples*: AdaBoost (Adaptive Boosting), Gradient Boosting, and XGBoost are well-known boosting algorithms. These algorithms assign weights to instances based on their accuracy and focus on improving the classification or regression performance over iterations.

In summary, bagging aims to reduce the variance of the base models and is characterized by the independence of training instances, while boosting focuses on increasing the accuracy of the ensemble by sequentially correcting the errors made by previous base models. The choice between bagging and boosting depends on the nature of the problem, the quality of base models, and the trade-off between variance reduction and error correction.

**AdaBoost (Adaptive Boosting), Gradient Boosting**, and **XGBoost** are all popular ensemble algorithms used in machine learning for improving predictive accuracy. They are based on boosting, a technique that sequentially combines the predictions of weak base models to create a strong learner. While they share similarities in boosting, they also have key differences in terms of algorithm design and performance optimization:

*AdaBoost (Adaptive Boosting):*
*Base Model*: AdaBoost typically uses decision trees with a depth of one (stumps) as base models. These stumps are often referred to as "weak learners."

*Algorithm Overview*: AdaBoost assigns weights to data points and focuses on instances that are misclassified by previous base models. It iteratively trains a sequence of base models, where each new model tries to correct the errors of the previous ones.

*Weighted Voting*: During each iteration, the weight of incorrectly classified instances is increased, and the weight of correctly classified instances is decreased. The final prediction is made by weighted voting, giving more weight to the most accurate base models.

*Robustness to Outliers*: AdaBoost is sensitive to outliers or noisy data points because it gives them higher importance as it attempts to correct their misclassifications.

*Performance Limitation*: AdaBoost can suffer from overfitting if the base model complexity increases, or if the data has a lot of noise.

*Gradient Boosting*:
*Base Model*: Gradient Boosting typically uses decision trees, but the trees can be deeper and more complex compared to AdaBoost. It's common to use regression trees for regression problems and classification trees for classification problems.

*Algorithm Overview*: Gradient Boosting is based on the gradient descent optimization technique. It minimizes a loss function by iteratively adding base models that move in the direction of the negative gradient of the loss function.

*Sequential Learning*: The base models are trained sequentially, and each new model tries to minimize the error made by the previous ones. Gradient Boosting learns a weighted sum of base models.

*Flexibility*: Gradient Boosting is more flexible in terms of base model choice and can accommodate a variety of loss functions and optimization techniques.

*Handling Outliers*: Gradient Boosting is relatively robust to outliers and can handle them gracefully.

**XGBoost (Extreme Gradient Boosting):**
*Base Model*: XGBoost uses a specialized type of decision tree called "CART" (Classification and Regression Tree) as base models. It employs a parallel and distributed computing framework.

*Algorithm Overview*: XGBoost is designed to optimize performance and computational efficiency. It uses a regularized objective function and includes a unique feature that penalizes complexity.

*Performance Optimizations*: XGBoost incorporates multiple performance optimizations, such as parallelization, out-of-core computing, and hardware-aware computation. These optimizations make it extremely fast and memory efficient.

*Regularization*: XGBoost provides L1 and L2 regularization techniques to prevent overfitting.

*Advanced Features*: XGBoost includes advanced features like handling missing values, monotonic constraints, and built-in support for multi-class classification.

**Comparison**:
All three algorithms are boosting methods, but XGBoost is known for its computational efficiency and scalability.
- AdaBoost focuses on simple base models and can be sensitive to noisy data.
- Gradient Boosting is more flexible and can handle more complex base models.
- XGBoost combines the advantages of Gradient Boosting with significant performance enhancements and regularization techniques.

The choice between AdaBoost, Gradient Boosting, and XGBoost depends on the specific problem, the quality of the data, the need for performance optimization, and the trade-offs between computational resources and predictive accuracy. XGBoost is often a popular choice for many machine learning tasks due to its balance of performance and efficiency.

To implement boosting algorithms in R, you can use various packages that provide implementations of popular boosting algorithms like AdaBoost, Gradient Boosting, and XGBoost. Here's a general overview of how to implement these boosting algorithms in R:

**AdaBoost (Adaptive Boosting):** You can implement AdaBoost in R using the "adabag" package, which provides tools for applying AdaBoost to classification problems. Here's a basic example:
```
# Install and load the adabag package
install.packages("adabag")
library(adabag)

# Load your dataset (e.g., data and labels)
```

```
data <- read.csv("your_data.csv")

# Create a formula for your classification task (e.g., Class ~
Feature1 + Feature2)
formula <- Class ~ Feature1 + Feature2

# Train an AdaBoost model
adaboost_model <- boosting(formula, data = data, mfinal = 100)  #
Adjust 'mfinal' as needed

# Make predictions
predictions <- predict(adaboost_model, newdata = data)

# Evaluate the model (e.g., calculate accuracy, confusion matrix,
etc.)
```

**Gradient Boosting**: For Gradient Boosting, you can use the "gbm" (Generalized Boosted Regression Models) package, which is primarily designed for regression but can also be used for classification problems. Here's an example:

```
# Install and load the gbm package
install.packages("gbm")
library(gbm)

# Load your dataset (e.g., data and labels)
data <- read.csv("your_data.csv")

# Create a formula for your classification task (e.g., Class ~
Feature1 + Feature2)
formula <- Class ~ Feature1 + Feature2

# Train a Gradient Boosting model
gbm_model <- gbm(formula, data = data, distribution = "bernoulli",
n.trees = 100)
# Adjust 'n.trees' as needed

# Make predictions
predictions <- predict(gbm_model, newdata = data, n.trees = 100)
# Adjust 'n.trees' as needed

# Evaluate the model (e.g., calculate accuracy, confusion matrix,
etc.)
```

**XGBoost (Extreme Gradient Boosting):** To use XGBoost in R, you should install and use the "xgboost" package. XGBoost is a popular boosting algorithm known for its speed and efficiency. Here's an example:

```
# Install and load the xgboost package
install.packages("xgboost")
library(xgboost)

# Load your dataset (e.g., data matrix and labels)
data <- read.csv("your_data.csv")
```

```r
# Prepare the data matrix and labels
X <- as.matrix(data[, c("Feature1", "Feature2")])
y <- data$Class

# Train an XGBoost model
xgb_model <- xgboost(data = X, label = y, nrounds = 100)
# Adjust 'nrounds' as needed

# Make predictions
predictions <- predict(xgb_model, newdata = X)

# Evaluate the model (e.g., calculate accuracy, confusion matrix,
etc.)
```

Please note that the above examples provide a basic outline of how to implement boosting algorithms in R. You should replace "your_data.csv" with the actual path to your dataset and adapt the code to your specific dataset, features, and classification problem. Additionally, you can explore further options and tuning parameters provided by these packages to optimize your boosting models. See the package documentation for details.

**Naïve Bayes** is a probabilistic machine learning algorithm used for classification and is based on Bayes' theorem. It's particularly well-suited for text classification and spam filtering but can be applied to a wide range of classification tasks. The "naïve" in its name arises from the assumption that the features used for classification are conditionally independent, which simplifies the model but may not hold in all real-world scenarios. Here's how the Naïve Bayes algorithm works:

1. *Data Representation*: Naïve Bayes starts with a labeled dataset, where each data point is represented as a feature vector and assigned to one of several classes.

2. *Feature Independence (Naïve Assumption)*: The central assumption in Naïve Bayes is that the features used for classification are conditionally independent. In other words, the presence or absence of one feature does not depend on the presence or absence of any other feature. This is a "naïve" simplifying assumption that may not hold in all situations but works surprisingly well in practice.

3. *Bayesian Probabilistic Framework*: The algorithm applies Bayes' theorem to calculate the conditional probability of each class given the observed features. The formula for this is:

Likelihood

Class Prior Probability

$$P(c\,|\,x) = \frac{P(x\,|\,c)\,P(c)}{P(x)}$$

Posterior Probability

Predictor Prior Probability

$$P(c\,|\,X) = P(x_1\,|\,c) \times P(x_2\,|\,c) \times \cdots \times P(x_n\,|\,c) \times P(c)$$

- $P(c|x)$ is the posterior probability of class (target) given predictor (attribute).

- $P(c)$ is the prior probability of class.
- $P(x|c)$ is the likelihood which is the probability of predictor given class.
- $P(x)$ is the prior probability of predictor.

4. *Training*: During the training phase, the algorithm calculates the following probabilities for each class:
$P(Class)$ is the prior probability of each class, which can be estimated by the frequency of each class in the training data. $P(Features \mid Class)$ is the conditional probability of the features given each class. This is typically estimated by counting the occurrences of each feature in each class.

5. *Prediction*: In the prediction phase, the algorithm calculates the conditional probabilities of each class given the observed features for a new data point. It does this for each class.

6. *Classification*: To classify a data point, the algorithm selects the class with the highest conditional probability given the features. This is often referred to as the "maximum a posteriori" or MAP estimation.

**Key Considerations**:
Naïve Bayes is simple and computationally efficient, making it well-suited for large datasets. It works particularly well for text classification tasks, such as spam email detection and sentiment analysis. While the independence assumption simplifies the model, it may not hold in all cases, which can lead to suboptimal results for certain types of data. Naïve Bayes is sensitive to feature selection, so careful choice of relevant features is crucial for good performance.

There are three common variants of the Naïve Bayes algorithm, depending on the distribution of the features:
- Gaussian Naïve Bayes: Assumes that the features follow a Gaussian (normal) distribution.
- Multinomial Naïve Bayes: Designed for discrete data, often used in text classification.
- Bernoulli Naïve Bayes: Suitable for binary data, such as presence or absence of features in a document.

Naïve Bayes is a versatile and widely used algorithm in various applications, including spam filtering, document categorization, and sentiment analysis.

Let's walk through a small example of using the Naïve Bayes algorithm for a basic text classification task. In this example, we'll classify short text messages as either "spam" or "not spam" (ham). We'll use a simple toy dataset of text messages. This type model can be applied to other types of data in addition to text-based data.

**Step 1**: *Import Necessary Libraries* : In R, you can use the "tm" (text mining) package for text preprocessing and the "e1071" package for Naïve Bayes classification. First, install and load the required packages:
```
install.packages("tm")
install.packages("e1071")
library(tm)
library(e1071)
```

**Step 2**: *Prepare the Data:* We'll create a small dataset of text messages and their corresponding labels (spam or ham).

```
# Sample data
text_messages <- c( "Cheap watches for sale!", "Hi there, how are
you?", "Win a free vacation today!", "Meeting at 2 PM.", "URGENT: Call
me ASAP!" )
labels <- c("spam", "ham", "spam", "ham", "spam")
# Create a data frame
data <- data.frame(text = text_messages, label = labels)
```

**Step 3**: *Preprocess the Text Data:* Text data needs to be preprocessed, which includes converting text to lowercase, removing punctuation, and creating a document-term matrix. The "tm" package is useful for these tasks.

```
# Create a Corpus (collection of text documents)
corpus <- Corpus(VectorSource(data$text))
# Preprocess the text
corpus <- tm_map(corpus, content_transformer(tolower))
corpus <- tm_map(corpus, removePunctuation)
corpus <- tm_map(corpus, removeNumbers)
corpus <- tm_map(corpus, stripWhitespace)
# Create a Document-Term Matrix (DTM)
dtm <- DocumentTermMatrix(corpus)
```

**Step 4**: *Train the Naïve Bayes Model:* We'll train a Naïve Bayes model using the preprocessed text data and labels. (dtm is not a data.frame, so this code exactly like this throws an error.)

```
# Train a Naive Bayes model
nb_model <- naiveBayes(x = dtm, y = data$label)
```

**Step 5**: *Make Predictions:* Use the trained Naïve Bayes model to make predictions for new text messages. In this example, we'll predict the labels for the same messages used in training.

```
# Prepare new text data for prediction
new_text <- c("Congratulations! You've won a prize!", "Let's have
lunch tomorrow.")
# Create a new Corpus
new_corpus <- Corpus(VectorSource(new_text))
# Preprocess the new text
new_corpus <- tm_map(new_corpus, content_transformer(tolower))
new_corpus <- tm_map(new_corpus, removePunctuation)
new_corpus <- tm_map(new_corpus, removeNumbers)
new_corpus <- tm_map(new_corpus, stripWhitespace)
# Create a Document-Term Matrix for new text
new_dtm <- DocumentTermMatrix(new_corpus)
# Make predictions
predictions <- predict(nb_model, newdata = new_dtm)
```

**Step 6**: *Evaluate the Model:* For a simple example like this, you can compare the predicted labels to the actual labels to assess the model's performance.

```
# Combine the new text and predictions
results <- data.frame(text = new_text, predicted_label = predictions,
stringsAsFactors = FALSE)
# Print the results
print(results)
```

This example demonstrates a basic text classification task using Naïve Bayes. In practice, you would typically work with larger and more diverse datasets and perform more thorough preprocessing. Additionally, you would evaluate the model's performance using metrics such as accuracy, precision, recall, and F1 score for a more comprehensive assessment.

**Multinomial Naïve Bayes, Gaussian Naïve Bayes, Bernoulli Naïve Bayes**, and the standard (or regular) Naïve Bayes are variants of the Naïve Bayes algorithm designed for different types of data and feature distributions. They differ in how they handle and model the data. Here's a summary of their key differences, advantages, and disadvantages:

1. **Multinomial Naïve Bayes**:
*Data Type*: Designed for discrete data, especially text data (e.g., word counts, term frequencies).
*Feature Distribution*: Assumes that features follow a multinomial distribution, making it suitable for count-based data.

*Advantages*:
- Works well for text classification tasks, such as spam detection and document categorization.
- Can handle data with multiple discrete categories or classes.

*Disadvantages*:
- Ignores word order and dependencies between terms in text data.
- May not perform well with continuous or real-valued features.

*How it Works*:
*Data Representation*: In text classification, each document is represented as a feature vector, with each feature representing the count or frequency of a term (word) in the document.

*Model Assumption*: Multinomial Naïve Bayes assumes that the features follow a multinomial distribution.

*Probability Estimation*: For each class, the algorithm estimates the conditional probability distribution of feature counts/frequencies. It calculates the likelihood of observing each feature in the documents of that class.

*Prior Probability*: It estimates the prior probability of each class based on the class distribution in the training data.

*Classification*: To classify a new document, the algorithm calculates the posterior probability of each class given the observed feature counts. It selects the class with the highest posterior probability.

2. **Gaussian Naïve Bayes**:
*Data Type*: Designed for continuous data with a Gaussian (normal) distribution.
*Feature Distribution*: Assumes that features follow a Gaussian distribution, making it suitable for real-valued data.

*Advantages*:
- Effective for continuous and normally distributed features.

- Works well for real-valued data, such as measurements.

**Disadvantages**:
- Sensitive to outliers in the data.
- May not perform well if features do not follow a Gaussian distribution.

**How it Works**:
*Data Representation*: Each data point is represented as a feature vector of real-valued features.

*Model Assumption*: Gaussian Naïve Bayes assumes that the features follow a Gaussian distribution within each class.

*Probability Estimation*: For each class, the algorithm estimates the parameters of the Gaussian distribution, including the mean and variance, for each feature.

*Prior Probability*: It estimates the prior probability of each class based on the class distribution in the training data.

*Classification*: To classify a new data point, the algorithm calculates the posterior probability for each class based on the Gaussian distribution parameters. It selects the class with the highest posterior probability.

3. **Bernoulli Naïve Bayes**:
*Data Type*: Designed for binary data (e.g., presence or absence of features).
*Feature Distribution*: Assumes that features follow a Bernoulli distribution.

**Advantages**:
- Useful for binary and presence/absence data, such as text data after binarization.
- Suitable for tasks like sentiment analysis or spam detection.

**Disadvantages**:
- Ignores feature frequency information, only considering the presence/absence.
- May not perform well with continuous or multi-category data.

**How it Works**:
*Data Representation*: Each data point is represented as a binary feature vector, where each feature is either 0 (absent) or 1 (present).

*Model Assumption*: Bernoulli Naïve Bayes assumes that the features follow a Bernoulli distribution, which models binary data.

*Probability Estimation*: For each class, the algorithm estimates the probability of each feature being present (1) or absent (0) in documents of that class.

*Prior Probability*: It estimates the prior probability of each class based on the class distribution in the training data.

*Classification*: To classify a new data point, the algorithm calculates the posterior probability of each class based on the Bernoulli distribution parameters. It selects the class with the highest posterior probability.

4. **Regular (Standard) Naïve Bayes**:
*Data Type*: More general and versatile; can be used for both discrete and continuous data.
*Feature Distribution*: Does not assume a specific feature distribution; it can handle different types of data.

***Advantages***:
- Versatile and can work with a mix of feature types (e.g., both text and numeric features).
- Doesn't make strict distributional assumptions, making it applicable to a wide range of problems.

***Disadvantages***:
- May not perform as well as specialized Naïve Bayes variants on specific types of data.
- Assumes feature independence, which may not hold in all cases.

In summary, the choice of which Naïve Bayes variant to use depends on the nature of the data and the problem you're trying to solve. It's essential to consider the data's characteristics and distribution when selecting the appropriate Naïve Bayes algorithm. Specialized variants like Multinomial, Gaussian, and Bernoulli Naïve Bayes are tailored for specific types of data and can offer better performance in those specific contexts. Regular Naïve Bayes is a more general and flexible option when the data may have mixed types or when you want to avoid making strong distributional assumptions.

In all three variants, the "Naïve" part of Naïve Bayes refers to the independence assumption: the algorithms assume that the features are conditionally independent given the class label. While this assumption simplifies the models, it may not always hold in practice. However, Naïve Bayes algorithms are known for their simplicity and often perform well in various classification tasks, especially when the characteristics of the data align with the assumptions of the chosen variant.

When the naïve assumption of independence fails, the models can perform less well than other classification models. Because it's simplicity, however, it may be a good first choice of model, and then compare the results with other models in an attempt to improve the results (performance).

Resources:
1. https://www.r-bloggers.com/2021/04/naive-bayes-classification-in-r/
2. https://www.geeksforgeeks.org/naive-bayes-classifier-in-r-programming/
3. https://uc-r.github.io/naive_bayes
4. https://www.learnbymarketing.com/tutorials/naive-bayes-in-r/
5. https://www.edureka.co/blog/naive-bayes-in-r/
6. https://www.analyticsvidhya.com/blog/2017/09/naive-bayes-explained/
7. https://www.analyticsvidhya.com/blog/2017/02/introduction-to-ensembling-along-with-implementation-in-r/
8. https://www.pluralsight.com/guides/ensemble-modeling-with-r
9. https://machinelearningmastery.com/machine-learning-ensembles-with-r/
10. https://daviddalpiaz.github.io/r4sl/ensemble-methods.html
11. https://www.listendata.com/2015/08/ensemble-learning-stacking-blending.html
12. https://www.datacamp.com/tutorial/ensemble-r-machine-learning

13. https://www.tmwr.org/ensembles
14. https://rpubs.com/guptadeepak/ensemble
15. https://www.r-bloggers.com/2021/02/machine-learning-with-r-a-complete-guide-to-gradient-boosting-and-xgboost/
16. https://www.geeksforgeeks.org/boosting-in-r/
17. https://bradleyboehmke.github.io/HOML/gbm.html
18. https://datascienceplus.com/gradient-boosting-in-r/

Lecture 7

K-Means
Linear Discriminant Analysis (LDA)
Nearest Centroid Classifier

**The Nearest Centroid Classifier (NCC)** also known as the **Nearest Mean Classifier** or **Centroid-Based Classifier**, is a simple and interpretable classification algorithm used to assign data points to the class whose centroid (mean) is nearest to the data point in feature space. It is particularly suitable for multi-class classification problems and can work well when class distributions are approximately Gaussian or have similar shapes.

*Training Phase*: For each class in the training dataset, calculate the mean (centroid) of the feature vectors of the data points belonging to that class. These centroids represent the "average" data point for each class.

*Classification Phase*: Given a new data point (the one you want to classify), calculate its distance (e.g., Euclidean distance) to each class centroid. Assign the data point to the class whose centroid is closest, i.e., the class with the smallest distance.

Mathematically, during the classification phase, the algorithm calculates the distance between a data point $x$ and each class centroid $c_i$ and assigns the data point to the class with the minimum distance:

$$Class(x) = \underset{i}{\operatorname{argmin}} \|x - c_i\|$$

Where:
$x$ is the data point you want to classify.
$c_i$ is the centroid of class $i$.
$\|x - c_i\|$ represents the distance between $x$ and $c_i$. This distance can be calculated using various distance metrics, such as the Euclidean distance.

***Advantages of the Nearest Centroid Classifier***:
- *Simplicity*: NCC is straightforward and easy to understand, making it a good choice for quick prototyping.
- *Interpretable*: The classifier's decision is based on the proximity to class centroids, providing transparency and interpretability.

***Limitations of the Nearest Centroid Classifier***:

- *Sensitivity to Outliers*: NCC can be sensitive to outliers because it uses the mean (centroid) of each class, and a single outlier can significantly affect the centroid.
- *Assumes Gaussian Distribution*: NCC works best when class distributions are approximately Gaussian or have similar shapes. It may not perform well with highly skewed data.
- *May not capture complex decision boundaries*: NCC assumes that class centroids are representative of the entire class, which may not be the case in situations where classes are complex or overlap.

The Nearest Centroid Classifier is often used in applications where interpretability and simplicity are essential, such as image compression, text classification, and feature selection. However, for more complex classification tasks with non-linear decision boundaries, other algorithms like support vector machines (SVMs), decision trees, or deep learning models may be more suitable.

**Linear Discriminant Analysis (LDA)** is a statistical technique used for dimensionality reduction and classification. It's a supervised learning algorithm that aims to find the linear combinations of features that best separate two or more classes in a dataset. LDA is commonly used in the context of pattern recognition, machine learning, and statistical analysis.

### *Key Concepts of Linear Discriminant Analysis*:
*Objective*: The primary objective of LDA is to maximize the separation between multiple classes while minimizing the variation within each class.

*Assumption*: LDA assumes that the features are normally distributed and that the classes have the same covariance matrix.

*Linear Combinations*: LDA finds linear combinations of the original features that create new axes, known as discriminant functions. These discriminant functions are chosen to maximize the distance between the means of different classes while minimizing the spread (variance) within each class.

*Decision Rule*: The decision rule in LDA involves comparing the posterior probabilities of a data point belonging to different classes. The class with the highest posterior probability is assigned to the data point.

*Covariance Matrix*: LDA calculates a pooled covariance matrix that represents the sum of the covariance matrices of individual classes, weighted by their sample sizes.

### *Steps in Linear Discriminant Analysis*:
*Compute the Within-Class Scatter Matrix ($S_W$)*:
Calculate the scatter matrix for each class and sum them to get the within-class scatter matrix.

$$S_W = \sum_{i=1}^{c} (X_i - M_i)(X_i - M_i)^T$$

where $c$ is the number of classes, $X_i$ is the matrix of data points for class $i$, and $M_i$ is the mean vector for class $i$.

*Compute the Between-Class Scatter Matrix ($S_B$)*:
Calculate the between-class scatter matrix.

$$S_B = \sum_{i=1}^{c} N_i(M_i - M)(M_i - M)^T$$

where $N_i$ is the number of data points in class $i$, $M_i$ is the mean vector for class $i$, $M$ is the overall mean vector.

*Compute the Eigenvalues and Eigenvectors*:
Solve the generalized eigenvalue problem $S_W^{-1}S_B$ to obtain the eigenvalues and eigenvectors.

*Select Discriminant Functions*: The discriminant functions are chosen based on the eigenvalues and eigenvectors. The number of discriminant functions is at most $c^{-1}$, where $c$ is the number of classes.

*Project Data onto Discriminant Functions*: Project the original data onto the selected discriminant functions to obtain a lower-dimensional representation.

*Classification*: Use the projected data for classification, often by applying a threshold or decision rule based on class centroids.

**Pros and Cons of Linear Discriminant Analysis:**
***Pros***:
- *Dimensionality Reduction*: LDA provides a way to reduce the dimensionality of the dataset while preserving class discriminatory information.
- *Feature Extraction*: It helps identify the most informative features for classification.
- *Regularization*: LDA can be regularized to handle cases where the covariance matrices are singular or poorly conditioned.

**Cons**:
- *Assumption of Normality*: LDA assumes that the features are normally distributed within each class.
- *Sensitive to Outliers*: LDA can be sensitive to outliers.
- *Binary Classification*: LDA is inherently designed for binary classification, but extensions exist for multiple classes.
- *Assumption of Equal Covariances*: LDA assumes that the classes have the same covariance matrix, which may not hold in all situations.

In summary, Linear Discriminant Analysis is a powerful technique for dimensionality reduction and classification, particularly when the goal is to separate multiple classes in a dataset. It is widely used in various fields, including pattern recognition and machine learning.

**Data clustering techniques** are essential in data mining for discovering meaningful patterns and structures in datasets. Clustering aims to group similar data points together while separating dissimilar ones. Here are some common data clustering techniques used in data mining:

***K-Means Clustering***:
*Method*: K-Means partitions data into $k$ clusters based on distance from cluster centroids.
*Advantages*: It is computationally efficient and works well with large datasets.
*Considerations*: The number of clusters ($k$) must be specified in advance, and it is sensitive to initial centroid selection.

***Hierarchical Clustering***:
*Method:* Hierarchical clustering creates a tree-like structure of clusters, with data points or groups merging step by step.
*Advantages*: It provides a hierarchy of clusters, allowing different granularity levels for analysis.
*Considerations*: Hierarchical clustering can be computationally intensive, and the choice of linkage method (single, complete, average, etc.) impacts results.

***DBSCAN (Density-Based Spatial Clustering of Applications with Noise):***
*Method*: DBSCAN groups data points into clusters based on their density and connectivity.
*Advantages*: It can discover clusters of arbitrary shapes and is robust to noise.
*Considerations*: DBSCAN requires setting parameters like the minimum number of points in a neighborhood.

***Agglomerative Clustering***:
*Method*: Agglomerative clustering starts with each data point as a single cluster and recursively merges the closest clusters.
*Advantages*: It is straightforward to implement and allows for a flexible number of clusters.
*Considerations*: Agglomerative clustering can be computationally demanding for large datasets.

**Spectral Clustering**:
*Method*: Spectral clustering transforms data into a lower-dimensional space and performs clustering in that space.
*Advantages*: It is effective for data with complex structures and works well when clusters have non-convex shapes.
*Considerations*: Spectral clustering involves eigenvalue decomposition and can be computationally expensive.

***Fuzzy C-Means Clustering***:
*Method*: Fuzzy C-Means assigns data points to clusters with membership degrees, allowing points to belong to multiple clusters to varying degrees.
*Advantages*: It accommodates data points that have ambiguous cluster assignments.
*Considerations*: It requires the tuning of a fuzziness parameter.

***Mean Shift Clustering***:
*Method*: Mean Shift identifies clusters by seeking modes in the density of data points.
*Advantages*: It is robust to initializations and can discover clusters of varying shapes and sizes.
*Considerations*: Parameter selection, especially bandwidth, can impact results.

***BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies):***
*Method*: BIRCH is a hierarchical clustering method optimized for large datasets and online clustering.
*Advantages*: It is memory-efficient and can handle streaming data.
*Considerations*: BIRCH is sensitive to the choice of parameters.

***Self-Organizing Maps (SOM):***
*Method*: SOM is a type of artificial neural network that maps high-dimensional data to a lower-dimensional grid while preserving topological relationships.
*Advantages*: It can visualize complex data structures and clusters.
*Considerations*: SOM requires the tuning of grid size and learning rate.

***OPTICS (Ordering Points To Identify the Clustering Structure):***
*Method*: OPTICS is an extension of DBSCAN that produces a reachability plot to reveal the cluster structure.
*Advantages*: It provides a more detailed view of clusters and density.
*Considerations*: OPTICS can be computationally intensive, and parameter tuning is required.

The choice of clustering technique depends on the characteristics of the data, the desired granularity of clusters, and the goals of the data mining task. It often involves experimenting with different methods and evaluating the quality of clustering results based on domain-specific criteria.
**K-Means** clustering is one of the most widely used and straightforward clustering techniques in data mining and machine learning. It's used to partition a dataset into groups or clusters based on the similarity of data points.

***Basic Procedure***:
*Initialization*: Choose the number of clusters ($k$) you want to create and initialize $k$ cluster centroids randomly or using some predefined method.

*Assignment Step*: Assign each data point to the nearest cluster centroid based on a distance metric, typically Euclidean distance. Each data point belongs to the cluster with the closest centroid.

*Update Step*: Recalculate the cluster centroids as the mean (average) of all data points assigned to that cluster.

*Iteration*: Repeat the assignment and update steps until the clusters stabilize or until a convergence criterion is met. Common convergence criteria include a maximum number of iterations, minimal centroid movement, or a predefined error threshold.

***Key Points***:
*Choosing the Number of Clusters* ($k$): Determining the appropriate number of clusters is often a critical decision. Common methods for selecting $k$ include the Elbow Method and the Silhouette Score. It's often necessary to experiment with different values of $k$ to find the most suitable one for your data.

*Initialization*: The choice of initial cluster centroids can impact the final clustering results. Random initialization can lead to different results in each run. The K-Means++ initialization method is often preferred as it distributes the initial centroids more effectively.

*Distance Metric*: The choice of distance metric (usually Euclidean distance) can be modified to suit the data and problem, such as using Manhattan distance or other custom distance functions.

*Sensitivity to Initializations*: K-Means can converge to local optima, which means the results can vary with different initializations. Running the algorithm multiple times with different initializations and selecting the best result is a common practice.

*Scalability*: K-Means can be computationally expensive, especially for large datasets, as it requires calculating distances between data points and centroids. Techniques like Mini-Batch K-Means can be used to make it more scalable.

*Cluster Shape*: K-Means assumes that clusters are spherical, equally sized, and isotropic, which means it may not perform well when dealing with non-spherical or unevenly sized clusters.

*Outliers*: K-Means can be sensitive to outliers, as they can significantly affect cluster centroids.

***Applications***: K-Means clustering is used in various domains and applications, including:
- Customer segmentation in marketing.
- Image compression and color quantization.
- Document classification.
- Anomaly detection.
- Genomic data analysis.
- Natural language processing.
- Recommendation systems.
- Identifying species in biology.
- Image and video analysis.

K-Means is a fundamental and widely used clustering technique, but it may not be suitable for all types of data and clustering tasks, especially when dealing with complex cluster shapes or varying cluster sizes. In such cases, more advanced clustering methods like hierarchical clustering or density-based clustering (e.g., DBSCAN) may be more appropriate.

In this example, we'll use the built-in *iris* dataset to perform K-Means clustering. This dataset contains measurements of sepal length, sepal width, petal length, and petal width for 150 iris flowers, with three species: setosa, versicolor, and virginica. We'll aim to cluster the flowers into three groups based on these measurements.
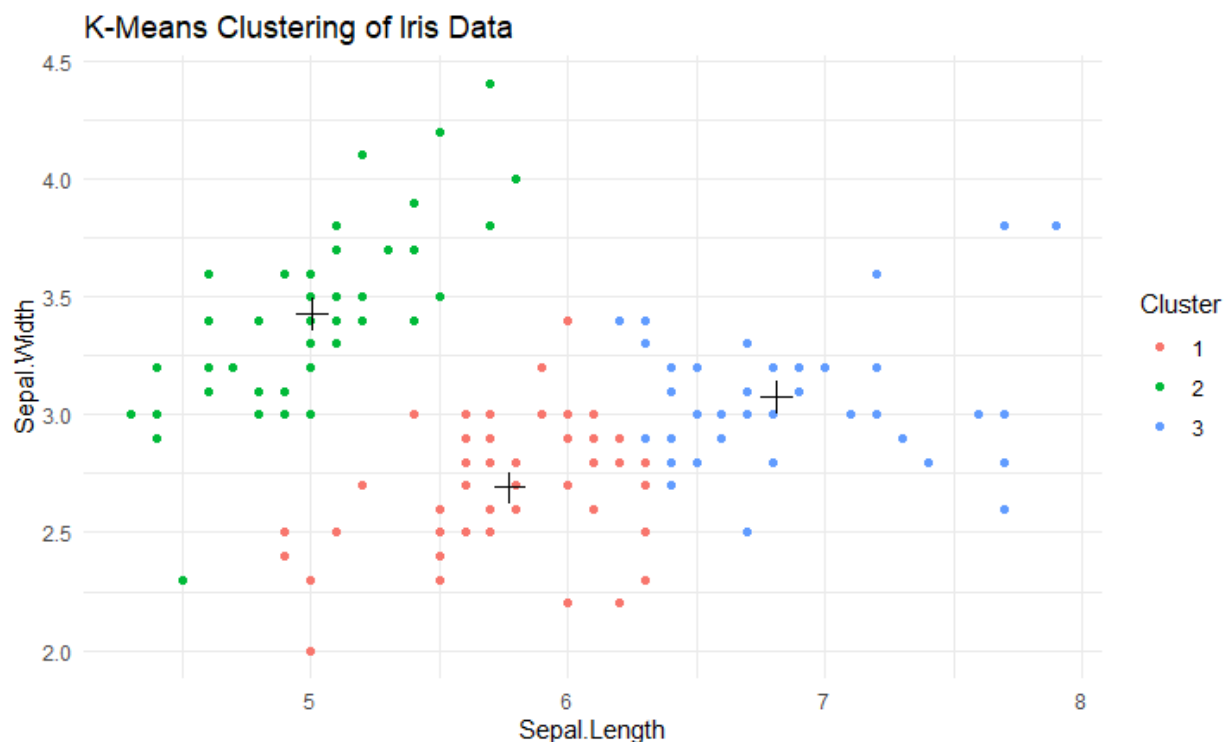
Here's how to perform K-Means clustering in R:
```
# Load the iris dataset
data(iris)
# Select the relevant features (sepal length and sepal width)
iris_features <- iris[, c("Sepal.Length", "Sepal.Width")]
# Set the number of clusters (k)
k <- 3
# Perform K-Means clustering
kmeans_result <- kmeans(iris_features, centers = k)
# Display the cluster assignments
cluster_assignments <- kmeans_result$cluster
print(cluster_assignments)
# Display the cluster centers
cluster_centers <- kmeans_result$centers
print(cluster_centers)
# Visualize the data and cluster centers
library(ggplot2)
# Plot the data points
ggplot(iris, aes(x = Sepal.Length, y = Sepal.Width, color =
factor(cluster_assignments))) + geom_point() + geom_point(data =
as.data.frame(cluster_centers), aes(x = Sepal.Length, y =
Sepal.Width), color = "black", size = 4, shape = 3) + labs(title = "K-
Means Clustering of Iris Data", color = "Cluster") + theme_minimal()
```

This code does the following:
- Loads the iris dataset and selects the relevant features (sepal length and sepal width).
- Sets the number of clusters ($k$) to 3.
- Performs K-Means clustering using the kmeans function.
- Displays the cluster assignments, which indicate which data point belongs to which cluster.
- Displays the cluster centers, which represent the mean values of each cluster's data points.
- Visualizes the data points with different colors representing the assigned clusters and marks the cluster centers on the plot.

When you run this code, you'll see a plot with data points colored by cluster, and the cluster centers marked as black triangles. K-Means clustering has grouped the data points into three clusters based on their sepal length and sepal width measurements.



You can adjust the value of k to experiment with different numbers of clusters and observe how it affects the clustering results. K-Means is sometimes used as a semi-supervised classification model, with k chosen according to the number of classes that are needed. Identification of which clusters coincide with which of the initial classes will have to be determined manually.

Resources:
1. https://www.datanovia.com/en/lessons/k-means-clustering-in-r-algorith-and-practical-examples/
2. https://uc-r.github.io/kmeans_clustering
3. https://www.datacamp.com/tutorial/k-means-clustering-r
4. https://www.geeksforgeeks.org/k-means-clustering-in-r-programming/

5.  https://www.guru99.com/r-k-means-clustering.html
6.  https://www.statology.org/k-means-clustering-in-r/
7.  https://towardsdatascience.com/k-means-clustering-in-r-feb4a4740aa
8.  https://www.geeksforgeeks.org/linear-discriminant-analysis-in-r-programming/
9.  https://www.r-bloggers.com/2021/05/linear-discriminant-analysis-in-r/
10. https://towardsdatascience.com/linear-discriminant-analysis-lda-101-using-r-6a97217a55a6
11. https://www.statology.org/linear-discriminant-analysis-in-r/
12. http://www.sthda.com/english/articles/36-classification-methods-essentials/146-discriminant-analysis-essentials-in-r/
13. https://uw.pressbooks.pub/appliedmultivariatestatistics/chapter/discriminant-analysis/
14. https://medium.com/edureka/linear-discriminant-analysis-88fa8ad59d0f
15. https://www.geeksforgeeks.org/ml-nearest-centroid-classifier/
16. https://cran.r-hub.io/web/packages/lolR/vignettes/nearestCentroid.html
17. https://idc9.github.io/stor390/notes/classification/classification.html

FYI, there is another LDA that comes up in data science contexts. We won't be learning this particular algorithm, but I've included some optional information about it here to help you distinguish the two.

LDA, which stands for *Latent Dirichlet Allocation*, is a popular probabilistic model used for topic modeling and document classification in natural language processing and text mining. LDA is a generative statistical model that helps discover underlying topics in a collection of documents and assigns topics to individual documents. Here's how LDA works:

*Assumptions*: LDA assumes that documents are mixtures of topics, and topics are mixtures of words. In other words, it assumes that there is a hidden (latent) structure of topics that generates the observed text data.

*Initialization*: LDA begins with a set of documents and a predefined number of topics (a hyperparameter). Each document is represented as a bag of words, where the order of words does not matter.

*Random Assignment*: Initially, LDA randomly assigns words in each document to one of the topics. This is called the "initialization phase."

*Iterations*: LDA iterates through the following steps until it converges to a stable state:
1.  For Each Word: LDA goes through each word in each document and reassigns it to a topic based on a probability distribution.
2.  For Each Topic: LDA updates the topic distribution for each document based on the words currently assigned to the topic.

*Mathematics*: LDA uses Bayesian statistics to compute the probability of a word belonging to a topic and the probability of a document containing a mixture of topics. It uses the Dirichlet distribution to model these probabilities.

*Convergence*: The iterations continue until the model converges or reaches a predetermined number of iterations.

*Output*: Once the model has converged, LDA provides two key outputs:
1. Topic-Word Distribution: This shows the probability of each word in the vocabulary belonging to each topic.
2. Document-Topic Distribution: This shows the probability of each document containing a mixture of topics.

*Interpretation*: After running LDA, you can interpret the results by examining the most probable words for each topic. This allows you to assign human-readable labels to the discovered topics.

***Key Points***:
- LDA is a generative model that uncovers the underlying structure of topics in a collection of documents.
- It is based on the idea that each document is a mixture of topics, and each topic is a mixture of words.
- LDA uses probability distributions to iteratively update the assignment of words to topics and the distribution of topics in documents.
- LDA is commonly used for document clustering, topic modeling, and content recommendation.
- LDA is a valuable tool for understanding and organizing large text datasets. It has applications in information retrieval, recommendation systems, and content analysis, among others.