

# Supervised learning

Supervised learning is a fundamental concept in machine learning, a subfield of artificial intelligence. It is a type of learning where an algorithm is trained on a labeled dataset to make predictions or decisions without explicit programming. In supervised learning, the algorithm learns from examples provided in the form of input-output pairs, where the input represents the data, and the output represents the desired outcome or label.

Here are some key points to understand about supervised learning:

1. **Labeled Data:** In supervised learning, you need a dataset with labeled examples. Labeled data consists of input features and their corresponding correct output or target values. For example, in a spam email classifier, the input features might be the email content, and the labels would indicate whether each email is spam or not.
2. **Training Phase:** During the training phase, the algorithm learns to make predictions by finding patterns and relationships in the labeled data. It adjusts its internal parameters to minimize the difference between its predictions and the actual labels in the training dataset.
3. **Types of Supervised Learning:**
  - **Classification:** In classification tasks, the algorithm learns to categorize input data into discrete classes or categories. Examples include spam detection, image classification, and sentiment analysis.
  - **Regression:** In regression tasks, the algorithm predicts a continuous numerical value. Examples include predicting house prices based on features like square footage and number of bedrooms or forecasting stock prices based on historical data.
4. **Evaluation:** Once the model is trained, it is evaluated on a separate dataset (the test set) to assess its performance. Common evaluation metrics include accuracy, precision, recall, F1-score for classification tasks, and mean squared error (MSE) for regression tasks.
5. **Generalization:** The goal of supervised learning is to build a model that can generalize its learned patterns to new, unseen data. A well-generalized model performs well on both the training data and unseen data, avoiding overfitting (model fitting noise in the training data) and underfitting (model being too simplistic).
6. **Applications:** Supervised learning is widely used across various domains, including healthcare for disease diagnosis, finance for fraud detection, natural language processing for translation and chatbots, autonomous vehicles for object recognition, and many more.
7. **Challenges:** Challenges in supervised learning include the need for high-quality labeled data, potential bias in the training data, selecting appropriate algorithms, and dealing with overfitting or underfitting.
8. **Continual Learning:** In some cases, supervised learning models can be updated over time with new data to adapt to changing conditions, a process known as continual learning.

In summary, supervised learning is a foundational machine learning approach that enables algorithms to learn from labeled data, make predictions or decisions, and generalize their knowledge to new, unseen data. It plays a crucial role in a wide range of real-world applications, making it an essential concept in the field of artificial intelligence and data science.

X=[1,2,3,4]- Features(Columns) y=[0,1,1,0]-Target variable - in form of (True or False)

## Unsupervised Learning

Unsupervised learning is a type of machine learning where the algorithm is trained on a dataset without labeled output or target variables. In unsupervised learning, the model's objective is to find patterns, structures, or relationships within the data without any explicit guidance in the form of correct answers. This makes it different from supervised learning, where the algorithm is trained on labeled data with known outcomes.

There are two primary types of unsupervised learning techniques:

1. **Clustering:** Clustering algorithms group data points into clusters or categories based on their similarities or proximity to each other. The goal is to identify natural groupings within the data without any prior knowledge of what those groups might be. Common clustering algorithms include K-Means clustering, hierarchical clustering, and DBSCAN (Density-Based Spatial Clustering of Applications with Noise).
2. **Dimensionality Reduction:** Dimensionality reduction techniques aim to reduce the number of features or variables in a dataset while preserving the essential information. This is useful for simplifying complex data, removing noise, and visualizing data in lower dimensions. Principal Component Analysis (PCA) and t-distributed Stochastic Neighbor Embedding (t-SNE) are popular dimensionality reduction methods.

Unsupervised learning has various applications, including:

- **Clustering:** Grouping customers for targeted marketing, identifying anomalies in cybersecurity, segmenting image data, etc.
- **Dimensionality Reduction:** Reducing the dimensionality of high-dimensional data for visualization, feature engineering, or speeding up subsequent supervised learning tasks.
- **Anomaly Detection:** Identifying unusual patterns or outliers in data, which can be applied in fraud detection, fault detection, and more.
- **Recommendation Systems:** Creating personalized recommendations for users by clustering or finding patterns in their behavior.
- **Natural Language Processing (NLP):** Word embedding techniques like Word2Vec and Doc2Vec are used for unsupervised learning to represent words or documents in a continuous vector space.
- **Image and Audio Processing:** Feature extraction and clustering in image and audio data analysis.

Unsupervised learning is a fundamental part of machine learning and data analysis, and it plays a crucial role in extracting valuable insights from unstructured or unlabeled datasets. It's often used as a preliminary step before applying supervised learning techniques or for exploratory data analysis.

## Reinforcement learning

Reinforcement learning (RL) is a subfield of artificial intelligence (AI) and machine learning that deals with how agents or decision-makers can learn to make sequences of decisions in an environment to maximize a cumulative reward. It is inspired by behavioral psychology, where an agent interacts with an environment and learns to take actions that lead to the best possible outcomes.

Here are some key concepts and components of reinforcement learning:

1. **Agent:** The entity that learns and makes decisions. This can be a robot, a game-playing AI, or any other autonomous system.
2. **Environment:** The external system or world in which the agent operates and interacts. The environment provides feedback to the agent in the form of rewards or penalties based on the agent's actions.
3. **State (s):** A representation of the current situation or configuration of the environment. The state contains all the relevant information the agent needs to make decisions.
4. **Action (a):** The choices or decisions that the agent can make to interact with the environment. Actions can be discrete or continuous, depending on the problem.
5. **Policy ( $\pi$ ):** A strategy or mapping that the agent uses to determine which action to take in a given state. The policy defines the agent's behavior.
6. **Reward (r):** A numerical value that the environment provides to the agent after each action. The reward serves as feedback to the agent, indicating how good or bad the action was in the current state. The goal of RL is to maximize the cumulative reward over time.
7. **Value Function (V):** A function that estimates the expected cumulative reward an agent can achieve starting from a particular state and following a specific policy. It helps the agent evaluate the desirability of different states.
8. **Q-Value Function (Q):** A function that estimates the expected cumulative reward an agent can achieve starting from a particular state, taking a specific action, and then following a particular policy. Q-values are used to determine the best action to take in a given state (Q-learning).
9. **Exploration vs. Exploitation:** RL agents face a trade-off between exploring new actions to discover potentially better strategies and exploiting known actions to maximize immediate rewards. Balancing exploration and exploitation is a fundamental challenge in RL.
10. **Markov Decision Process (MDP):** A formal mathematical framework used to model RL problems. It includes the concepts of states, actions, rewards, and transition probabilities, and it assumes the Markov property, which means that the future depends only on the current state and action.

Reinforcement learning algorithms can be broadly categorized into model-free and model-based methods. Model-free methods, like Q-learning and deep reinforcement learning (DRL), directly learn the optimal policy from interacting with the environment. Model-based methods involve building a model of the environment and using it to plan and make decisions.

RL has applications in various domains, including robotics, autonomous vehicles, game playing (e.g., AlphaGo), recommendation systems, and healthcare. It has gained significant attention in recent years due to breakthroughs in deep reinforcement learning and its potential to tackle complex problems. However, RL can be challenging, as it often requires a large number of

interactions with the environment and may suffer from issues like exploration difficulties and

# Common application of ML

Machine Learning (ML) has a wide range of common applications across various industries and domains. Here are some of the most prevalent applications of machine learning:

## 1. Image and Video Recognition:

- **Object Detection:** ML algorithms can identify and locate objects within images or videos, which is used in autonomous vehicles, surveillance, and facial recognition.
- **Image Classification:** ML models can classify images into categories, such as in medical imaging (identifying diseases), e-commerce (product recommendation), and content moderation.

## 2. Natural Language Processing (NLP):

- **Text Classification:** ML models can categorize text data into topics, sentiments, or spam detection, commonly used in email filtering, social media analysis, and news categorization.
- **Machine Translation:** Services like Google Translate use ML to translate text between languages.
- **Chatbots and Virtual Assistants:** ML powers chatbots and virtual assistants that can understand and respond to natural language input.

## 3. Recommendation Systems:

- ML algorithms are used in recommendation engines to suggest products, movies, music, and content to users based on their past behavior and preferences, as seen in Netflix and Amazon recommendations.

## 4. Predictive Analytics:

- Businesses use ML for forecasting demand, stock prices, and customer behavior, enabling better decision-making and resource allocation.

## 5. Healthcare:

- ML aids in medical diagnosis, predicting disease outbreaks, and drug discovery. It's also used for personalized treatment recommendations based on patient data.

## 6. Financial Services:

- **Fraud Detection:** ML models can identify fraudulent transactions by analyzing patterns in financial data.
- **Algorithmic Trading:** ML is used to develop trading algorithms that make predictions based on market data.

## 7. Autonomous Vehicles:

- ML plays a crucial role in self-driving cars by processing sensor data (e.g., lidar, cameras) to make real-time driving decisions.

## 8. Gaming:

- ML is used to create intelligent NPCs (non-player characters), improve game physics, and enhance user experiences in video games.

## 9. Manufacturing and Quality Control:

- ML can monitor production lines for defects, optimize processes, and predict maintenance needs to reduce downtime.

## 10. Supply Chain Management:

- ML models can optimize inventory management, route planning, and demand forecasting in logistics and supply chain operations.

## 11. Energy Management:

- ML is used for predicting energy consumption, optimizing grid operations, and improving energy efficiency in smart grids.

## 12. Social Media Analysis:

- ML helps in sentiment analysis, content recommendation, and trend prediction on social media platforms.

## 13. Agriculture:

- ML is used for crop yield prediction, disease detection in plants, and precision farming.

## 14. Security:

- ML is employed in cybersecurity for anomaly detection, identifying malicious activities, and enhancing network security.

## 15. Human Resources:

- ML can assist in candidate screening, employee retention prediction, and workforce planning.

These are just a few examples, and the applications of machine learning continue to expand as the technology advances and businesses find new ways to leverage it for solving complex problems and improving efficiency.

# Regression

Regression is a statistical method used in data analysis and machine learning to model the relationship between a dependent variable and one or more independent variables. It is primarily used for predictive modeling and understanding the relationship between variables.

The key idea behind regression is to find a mathematical equation that best describes the relationship between the variables. This equation is typically represented as:

$$[ Y = aX + b ]$$

Where:

- ( Y ) is the dependent variable (the variable you want to predict or explain).
- ( X ) is the independent variable (the variable that you use to make predictions or explain the variation in the dependent variable).
- ( a ) is the coefficient (slope) that represents how much the dependent variable changes for a unit change in the independent variable.

- ( b ) is the intercept, which represents the value of the dependent variable when the independent variable is zero.

The goal in regression analysis is to estimate the values of ( a ) and ( b ) that best fit the data. There are different types of regression models, including:

1. **Linear Regression:** This is the simplest form of regression, where the relationship between the variables is assumed to be linear. Linear regression seeks to find the best-fitting line through the data points.
2. **Multiple Regression:** In this model, there are multiple independent variables, and the goal is to find the best-fitting linear equation that explains the relationship between the dependent variable and all the independent variables.
3. **Polynomial Regression:** When the relationship between the variables is not linear, polynomial regression can be used. It involves fitting a polynomial equation to the data.
4. **Logistic Regression:** This is used when the dependent variable is binary or categorical. It models the probability of an event occurring based on one or more independent variables.
5. **Ridge and Lasso Regression:** These are techniques used to address overfitting in linear regression by adding regularization terms to the model.
6. **Time Series Regression:** This is used when the data involves time-dependent variables. Time series regression models the relationship between variables over time.

Regression analysis is widely used in various fields, including economics, finance, social sciences, and machine learning. It helps in making predictions, understanding the impact of independent variables on the dependent variable, and drawing insights from data. The choice of regression model depends on the nature of the data and the specific research or analysis.

X- independent variable- In linear regression - one independent variable and one dependent variable  $y = aX + b$  # a- is the coefficient and b is the intercept value if we talk about more than one independent variable - that comes multiple regression  $X_1, X_2, X_3, X_4$ - independent variable which is more than one variable  $y = a_1X_1 + a_2X_2 + a_3X_3 + a_4X_4 + b$  y- dependent variable

Certainly, let's delve deeper into regression analysis and its various aspects:

### 1. Linear Regression:

- **Simple Linear Regression:** In simple linear regression, there is only one independent variable ((X)). The goal is to find a straight line (a linear equation) that best fits the data points. The equation looks like ( $Y = aX + b$ ).
- **Multiple Linear Regression:** In multiple linear regression, there are two or more independent variables ((X<sub>1</sub>, X<sub>2</sub>, \ldots, X<sub>n</sub>)). The goal is to find the best-fitting linear equation that explains how each of these independent variables collectively influences the dependent variable ((Y)). The equation becomes:

$$[Y = a_1X_1 + a_2X_2 + \ldots + a_nX_n + b]$$

### 2. Polynomial Regression:

- When the relationship between the dependent and independent variables is not linear, polynomial regression can be used. In this case, the equation is a polynomial of degree (n), where (n) is a positive integer.

$$[Y = a_0 + a_1X + a_2X^2 + \ldots + a_nX^n]$$

- For example, a quadratic regression ((n = 2)) would have a parabolic relationship between (X) and (Y).

### 3. Logistic Regression:

- Logistic regression is used when the dependent variable is categorical, typically binary (e.g., 0 or 1, yes or no). It models the probability of an event occurring.

$$[P(Y=1) = \frac{1}{1 + e^{-(aX + b)}}]$$

- Here, (P(Y=1)) is the probability of the event happening, and (e) is the base of the natural logarithm.

### 4. Regularization in Regression:

- Regularization techniques like Ridge and Lasso are used to prevent overfitting in linear regression models.
- Ridge Regression adds a penalty term to the linear regression equation to shrink the coefficients ((a)) towards zero.
- Lasso Regression not only shrinks coefficients but also has a feature selection property where it can make some coefficients exactly zero, effectively removing those features from the model.

### 5. Time Series Regression:

- Time series regression is used when dealing with time-dependent data, where observations are recorded at regular time intervals.
- This form of regression accounts for the temporal dependencies in the data, making it suitable for forecasting and analyzing time series data.

### 6. Model Evaluation:

- To assess the quality of a regression model, various metrics are used, such as Mean Squared Error (MSE), R-squared ((R<sup>2</sup>)) for linear regression, and Log-Likelihood for logistic regression.
- Cross-validation is often employed to estimate how well the model will perform on unseen data.

### 7. Assumptions in Regression:

- Linear regression assumes that there is a linear relationship between the independent and dependent variables.
- It also assumes that the residuals (the differences between predicted and actual values) are normally distributed, have constant variance (homoscedasticity), and are independent of each other.

### 8. Use Cases:

- Regression analysis is used in a wide range of applications, including predicting stock prices, estimating housing prices, analyzing the impact of marketing campaigns on sales, and understanding the relationship between factors like age and disease risk.

### 9. Machine Learning:

- Regression is a fundamental technique in machine learning. Many machine learning algorithms, such as support vector machines and neural networks, have regression as a component for tasks like regression analysis and predicting numerical values.

In practice, choosing the appropriate regression model depends on the specific problem you are trying to solve and the nature of your data. Understanding the relationships between

## **Simple Linear Regression with scikit-learn:**



```
In [22]: # Import necessary libraries library name is scikit Learn - import sklearn
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression

# Generate some sample data
np.random.seed(0)
X = 2 * np.random.rand(10, 1)
print(X)#One independent variable
Y = 4 + 3 * X + np.random.randn(10, 1)
print(Y)# dependent variable

# Create a LinearRegression model
model = LinearRegression()

# Fit the model to the data
model.fit(X, Y)

# Get the coefficients and intercept
a = model.coef_ # coefficient of linear equation
b = model.intercept_# intercept value of linear equation

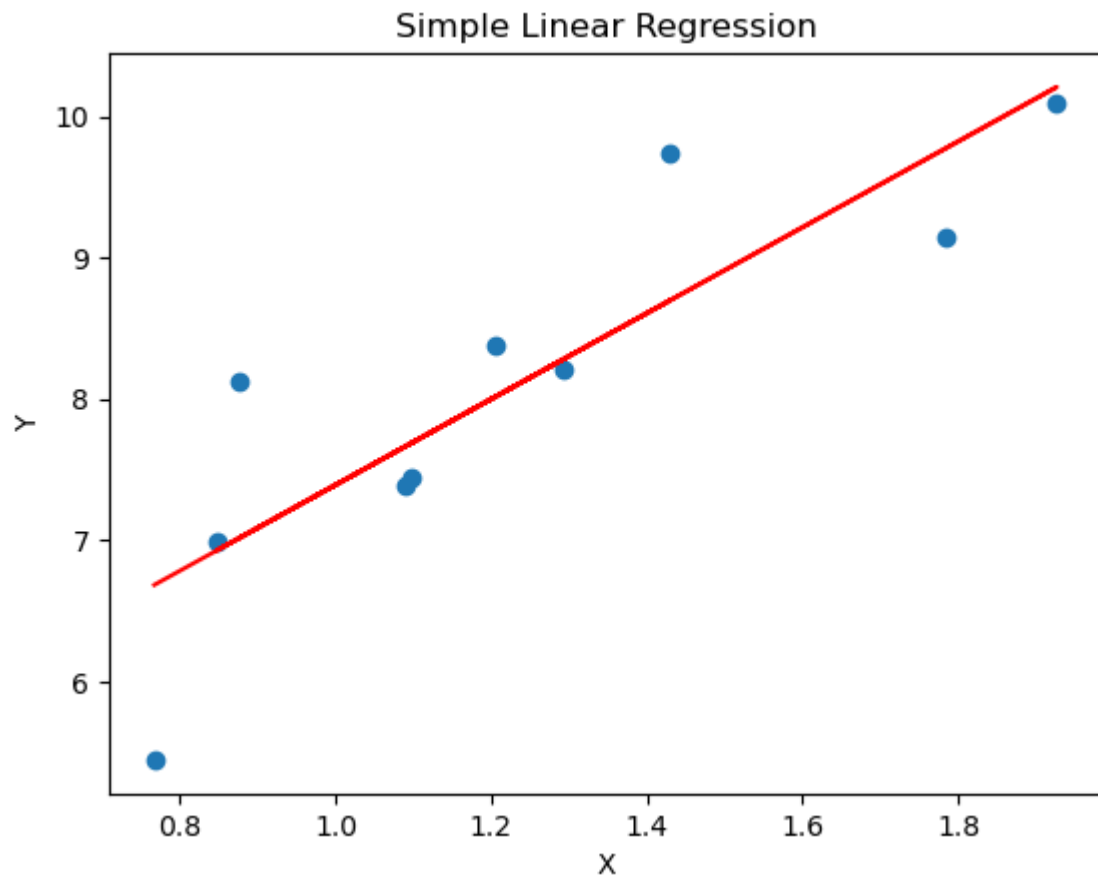
# Make predictions
X_new = np.array([[1.5]]) # New data point for prediction
Y_pred = model.predict(X_new)
# plotting library- matplotlib.pyplot as plt
# Plot the data and regression line
plt.scatter(X, Y, alpha=1)# alpha is used for transparency of your data point
plt.plot(X, a * X + b, color='red')
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Simple Linear Regression')
plt.show()

print(f"Coefficient (a): {a[0][0]}")
print(f"Intercept (b): {b[0]}")
print(f"Predicted Y for X=1.5: {Y_pred[0][0]}")
```

```

[[1.09762701]
 [1.43037873]
 [1.20552675]
 [1.08976637]
 [0.8473096 ]
 [1.29178823]
 [0.87517442]
 [1.783546  ]
 [1.92732552]
 [0.76688304]]
[[ 7.43692459]
 [ 9.74540971]
 [ 8.37761798]
 [ 7.39097411]
 [ 6.98579203]
 [ 8.20903901]
 [ 8.11960234]
 [ 9.14547974]
 [10.09504426]
 [ 5.44655337]]

```



Coefficient (a): 3.037773738915913  
 Intercept (b): 4.354126425524145  
 Predicted Y for X=1.5: 8.910787033898014

## Simple Linear Regression with StatsModels:



```
In [3]: # Import necessary libraries
import numpy as np
import statsmodels.api as sm

# Generate some sample data
np.random.seed(0)
X = 2 * np.random.rand(100, 1)
Y = 4 + 3 * X + np.random.randn(100, 1)

# Add a constant term to the independent variable for the intercept
X_with_const = sm.add_constant(X)

# Create and fit a StatsModels OLS (Ordinary Least Squares) model
model = sm.OLS(Y, X_with_const).fit()

# Get the model summary
print(model.summary())

# Make predictions
X_new = sm.add_constant(np.array([[1.5, 2.5]])) # New data point for prediction
Y_pred = model.predict(X_new)

print(f"Coefficient (a): {model.params[1]}")
print(f"Intercept (b): {model.params[0]}")
print(f"Predicted Y for X=1.5: {Y_pred[0]}")
```

```

=====
                        OLS Regression Results
=====
=
Dep. Variable:          y      R-squared:          0.74
7
Model:                  OLS    Adj. R-squared:      0.74
4
Method:                 Least Squares    F-statistic:      289.
3
Date:                   Sat, 09 Sep 2023    Prob (F-statistic): 5.29e-3
1
Time:                   17:50:55    Log-Likelihood:    -141.5
1
No. Observations:      100    AIC:              287.
0
Df Residuals:          98    BIC:              292.
2
Df Model:               1
Covariance Type:        nonrobust
=====
=
                        coef      std err          t      P>|t|      [0.025      0.97
5]
-----
-
const          4.2222      0.193      21.843      0.000      3.839      4.60
6
x1             2.9685      0.175      17.009      0.000      2.622      3.31
5
=====
=
Omnibus:           11.746    Durbin-Watson:      2.08
3
Prob(Omnibus):     0.003    Jarque-Bera (JB):    4.09
7
Skew:              0.138    Prob(JB):            0.12
9
Kurtosis:          2.047    Cond. No.            3.5
8
=====
=

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correc
tly specified.
Coefficient (a): 2.968467510701018
Intercept (b): 4.222151077447228
Predicted Y for X=1.5: 13.754395392923387

```

**Certainly, to perform regression analysis on a dataset in Python, you'll typically follow these steps:**

1. **Load the Dataset:** First, you need to load your dataset into Python. You can use libraries like Pandas to read data from various file formats (e.g., CSV, Excel) or from databases.
2. **Preprocess the Data:** Data preprocessing is an essential step. You might need to handle missing values, encode categorical variables, and split the data into features (independent variables) and the target variable (dependent variable).
3. **Choose a Regression Model:** Depending on your data and problem type (e.g., linear regression, logistic regression), select an appropriate regression model. Import the necessary libraries.
4. **Split the Data:** Split the dataset into training and testing sets. This helps you evaluate the model's performance on unseen data.
5. **Train the Model:** Fit the regression model to the training data using the `.fit()` method. This step estimates the model parameters.
6. **Evaluate the Model:** Use evaluation metrics to assess how well your model performs. Common metrics for regression include Mean Squared Error (MSE), R-squared ( $R^2$ ), and Root Mean Squared Error (RMSE).
7. **Make Predictions:** Use the trained model to make predictions on the test data or new data points.
8. **Visualize the Results:** Visualize the model's predictions and its relationship with the data using plots and charts.

Here's a general example of how to perform linear regression on a dataset using Python with the scikit-learn library:

```

# Import necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

# Load the dataset (replace 'your_dataset.csv' with the actual dataset file)
data = pd.read_csv('your_dataset.csv')

# Assuming you have a column 'X' as the independent variable and 'Y' as the dependent variable
X = data[['X']]
Y = data['Y']

# Split the data into training and testing sets
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)

```

## Classification

Classification is a fundamental concept in machine learning and statistics. It involves the process of categorizing or labeling data points into predefined classes or categories based on their characteristics or features. The goal of classification is to build a model that can accurately predict the class or category of new, unseen data points.

Here are some key points related to classification:

1. **Supervised Learning:** Classification is a type of supervised learning, where the algorithm is trained on a labeled dataset. This means that each data point in the training dataset is associated with a known class label.
2. **Classes:** In a classification problem, there are typically two or more classes or categories into which the data points are grouped. For example, in a binary classification problem, there are two classes (e.g., spam or not spam, positive or negative sentiment). In multiclass classification, there are more than two classes.
3. **Features:** Features are the characteristics or attributes of the data points that are used to make predictions. These features can be numerical, categorical, or a combination of both.
4. **Classification Algorithms:** There are various classification algorithms, including but not limited to:
  - Logistic Regression
  - Decision Trees
  - Random Forest
  - Support Vector Machines (SVM)
  - k-Nearest Neighbors (k-NN)
  - Naive Bayes

- Neural Networks (Deep Learning)
5. **Training and Testing:** The classification model is trained on a labeled training dataset, and its performance is evaluated on a separate testing dataset to assess its ability to generalize to unseen data.
  6. **Metrics:** Common evaluation metrics for classification models include accuracy, precision, recall, F1 score, and area under the receiver operating characteristic curve (ROC AUC), among others.
  7. **Overfitting and Underfitting:** Like other machine learning models, classification models can suffer from overfitting (model is too complex and fits noise in the data) or underfitting (model is too simple to capture the underlying patterns).
  8. **Imbalanced Data:** In some classification problems, the distribution of classes in the dataset may be imbalanced, with one class significantly outnumbering the others. Handling imbalanced data is a common challenge in classification.
  9. **Feature Engineering:** Feature engineering involves selecting and transforming the most relevant features for the classification problem. It can significantly impact the performance of a classification model.
  10. **Applications:** Classification is widely used in various applications, including:
    - Spam email detection
    - Sentiment analysis
    - Image classification
    - Medical diagnosis
    - Credit risk assessment
    - Fraud detection
    - Object recognition
    - Natural language processing tasks

In summary, classification is a machine learning technique that plays a crucial role in solving a wide range of real-world problems by automatically assigning data points to predefined categories or classes based on their features. It is a key component of many practical applications and is supported by a variety of algorithms and evaluation methods.

```

In [17]: # Import necessary libraries
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
import pandas as pd

# Load the Iris dataset
iris = datasets.load_iris()
X = iris.data # Features- petal length, petal width, sepal length, sepal width
print(X[0])
y = iris.target # Target Labels (species)
print(y[0])# 0- species name -1 virginica , 1- vesicular , 3- another species
# Split the dataset into training and testing sets (70% training, 30% testing)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)
# X_train, y_train- Training data-70%
# X_test, y_test-Testing data-30 %
# Preprocess the data (standardize features)
scaler = StandardScaler()# Normalization- is technique where we normalize the
X_train = scaler.fit_transform(X_train) # to transform the value in data set
X_test = scaler.transform(X_test)

# Create a classification model (Logistic Regression in this case)
model = LogisticRegression()# is algorithm comes under in classification - i

# Train the model on the training data
model.fit(X_train, y_train)# this is used for fitting training data

# Make predictions on the test data
y_pred = model.predict(X_test) # for prediction - simply we pass X_test
print("Predict:",y_pred)
print("testing",y_test)

# Evaluate the model's performance
accuracy = accuracy_score(y_test, y_pred) # it is used to calculate accuracy of
print("Accuracy:", accuracy)

# Print a classification report with additional metrics
print(classification_report(y_test, y_pred))# ALL the precision attribute

# Print the names of the target classes
target_names = iris.target_names
print("Target Class Names:", target_names)
y_pred1=model.predict(np.array([[1.5,2.5,3.5,3]]))# features -random
print("predict1",y_pred1)

```



```

[5.1 3.5 1.4 0.2]
0
Predict: [2 2 1 1 1 2 0 2 2 1 2 2 1 0 1 1 1 0 0 2 2 2 1 2 1 2 2 1 2 1 0 0 1 1
1 1 1
1 1 0 1 2 0 1 1]
testing [2 2 1 1 1 2 0 2 2 1 1 2 1 0 1 1 1 0 0 2 2 2 1 2 1 1 2 1 2 1 0 0 1 1
1 1 1
1 1 0 1 2 0 1 1]
Accuracy: 0.9555555555555556

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	8
1	1.00	0.92	0.96	24
2	0.87	1.00	0.93	13
accuracy			0.96	45
macro avg	0.96	0.97	0.96	45
weighted avg	0.96	0.96	0.96	45

```

Target Class Names: ['setosa' 'versicolor' 'virginica']
predict1 [2]

```

## Common business application

Common business applications are software programs or tools that are widely used by businesses to perform various tasks and operations efficiently. These applications help streamline processes, improve productivity, and make data-driven decisions. Here are some common types of business applications:

1. **Office Suites:** Office suites like Microsoft Office (Word, Excel, PowerPoint) and Google Workspace (Docs, Sheets, Slides) are essential for creating documents, spreadsheets, presentations, and email communication.
2. **Accounting Software:** Tools like QuickBooks, Xero, and FreshBooks help businesses manage their finances, track expenses, create invoices, and generate financial reports.
3. **Customer Relationship Management (CRM):** CRM software such as Salesforce, HubSpot, and Zoho CRM helps businesses manage customer interactions, track sales leads, and improve customer satisfaction.
4. **Enterprise Resource Planning (ERP):** ERP systems like SAP, Oracle, and Microsoft Dynamics offer comprehensive solutions for managing various aspects of a business, including finance, HR, inventory, and supply chain.
5. **Project Management:** Tools like Trello, Asana, and Basecamp assist in planning, tracking, and collaborating on projects, ensuring that they are completed on time and within budget.
6. **Human Resources Management:** HR software such as BambooHR, ADP, and Workday automates HR processes, including employee onboarding, payroll, benefits administration, and performance management.
7. **Inventory Management:** Inventory software like Fishbowl and TradeGecko helps businesses track and manage their inventory levels, orders, and stock replenishment.
8. **Communication and Collaboration:** Tools like Slack, Microsoft Teams, and Zoom facilitate internal and external communication, file sharing, and virtual meetings.

9. **Business Intelligence (BI):** BI software such as Tableau, Power BI, and QlikView allows businesses to analyze and visualize data to make informed decisions.
10. **Marketing Automation:** Marketing automation platforms like HubSpot, Marketo, and Mailchimp help businesses automate marketing campaigns, lead nurturing, and customer engagement.
11. **E-commerce Platforms:** Platforms like Shopify, WooCommerce, and Magento enable businesses to create and manage online stores for selling products and services.
12. **Customer Support and Help Desk:** Customer support tools like Zendesk, Freshdesk, and Intercom assist businesses in providing efficient customer support and managing inquiries.
13. **Content Management Systems (CMS):** CMS like WordPress, Drupal, and Joomla are used to create and manage websites and digital content.
14. **Security and Antivirus Software:** Security software such as McAfee, Norton, and Bitdefender helps protect businesses from cyber threats and ensure data security.
15. **Analytics and Reporting:** Analytics tools like Google Analytics and Adobe Analytics provide insights into website and app performance, user behavior, and marketing effectiveness.
16. **Document Management:** Document management systems like SharePoint and DocuWare help businesses organize, store, and retrieve documents securely.
17. **Time and Expense Tracking:** Time tracking and expense management tools like Toggl and Expensify assist in tracking employee hours and expenses for accurate billing and payroll.
18. **Collaborative Document Editing:** Tools like Google Docs and Microsoft SharePoint enable real-time collaboration on documents, spreadsheets, and presentations.
19. **Video Conferencing and Webinars:** Platforms like Zoom, Webex, and GoToMeeting facilitate video conferencing, webinars, and online presentations.
20. **Social Media Management:** Social media management tools like Hootsuite and Buffer help businesses schedule posts, monitor social media activity, and engage with their audience.

These are just a few examples of common business applications. The choice of specific applications depends on the needs and goals of each individual business. Many businesses also integrate multiple applications to create a customized technology stack that meets their

## Distance vectors

In machine learning, "distance vectors" is not a standard term. However, it's possible that you might be referring to vectors that represent the distances between data points in a dataset. These vectors can be used for various purposes, such as clustering, anomaly detection, or dimensionality reduction. Let me explain this concept in more detail.

### Distance Vectors in Clustering:

One common use of distance vectors in machine learning is in clustering algorithms. In clustering, you aim to group similar data points together based on some distance metric. Here's how distance vectors are used in clustering:

1. **Data Representation:** Each data point in your dataset is represented as a feature vector. These feature vectors can be high-dimensional and include various attributes or features describing the data points.

2. **Distance Calculation:** A distance metric (e.g., Euclidean distance, cosine similarity, etc.) is used to compute the distance between each pair of data points in the dataset. These distances are often stored in a distance matrix.
3. **Clustering:** Clustering algorithms, such as K-Means, Hierarchical Clustering, or DBSCAN, use the distance matrix to group similar data points into clusters. The distance vectors help determine which data points are close to each other and should belong to the same cluster.

### Example:

Let's consider an example where you have a dataset of customer data, and you want to group customers based on their purchasing behavior. You represent each customer as a feature vector with attributes like "total spending," "average purchase frequency," and "geographic location."

You calculate the pairwise distances between customers using a distance metric (e.g., Euclidean distance) and create a distance matrix. This matrix contains distance vectors that represent the dissimilarity between each pair of customers.

```
import numpy as np
from sklearn.metrics import pairwise_distances

# Sample customer data
customer_data = np.array([
    [100, 5, 1], # Customer 1
    [200, 4, 2], # Customer 2
    [50, 6, 3], # Customer 3
])

# Calculate pairwise Euclidean distances
distance_matrix = pairwise_distances(customer_data, metric='euclidean')

print(distance_matrix)
```

In this example, the distance matrix contains distance vectors that indicate how different each customer is from every other customer. Clustering algorithms can then use this matrix to group customers with similar purchasing behavior.

So, while "distance vectors" may not be a standard term, the concept of using vectors to represent distances between data points is a crucial aspect of various machine learning tasks, especially in clustering and similarity-based analysis.

## distances used in ML

In machine learning, various distance metrics are used to quantify the similarity or dissimilarity between data points, which can be crucial for tasks such as clustering, classification, recommendation systems, and dimensionality reduction. Here are some commonly used distance metrics in machine learning:

### 1. Euclidean Distance:

- Euclidean distance is one of the most common distance metrics.
- It measures the straight-line distance between two points in a Euclidean space.
- It is defined as the square root of the sum of squared differences between corresponding elements of two vectors.
- Formula: 
$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

### 2. Manhattan Distance (L1 Norm):

- Manhattan distance is also known as the L1 norm or Taxicab distance.
- It measures the distance between two points as the sum of the absolute differences of their coordinates.
- Formula: 
$$\text{Manhattan Distance} = \sum_{i=1}^n |x_i - y_i|$$

### 3. Cosine Similarity:

- Cosine similarity measures the cosine of the angle between two vectors.
- It is often used for text and document analysis but is applicable to any data represented as vectors.
- Values range from -1 (perfectly dissimilar) to 1 (perfectly similar).
- Formula: 
$$\text{Cosine Similarity} = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \cdot \|\mathbf{B}\|}$$

### 4. Jaccard Similarity:

- Jaccard similarity is used to measure the similarity between two sets.
- It is defined as the size of the intersection of the sets divided by the size of their union.
- Commonly used in text analysis and recommendation systems.
- Formula: 
$$\text{Jaccard Similarity} = \frac{|A \cap B|}{|A \cup B|}$$

### 5. Minkowski Distance:

- Minkowski distance is a generalization of both Euclidean and Manhattan distances.
- It allows you to control the degree of emphasis placed on each dimension by adjusting the exponent (p).
- Euclidean distance corresponds to (p = 2), and Manhattan distance corresponds to (p = 1).

### 6. Hamming Distance:

- Hamming distance measures the difference between two binary strings of equal length.
- It counts the number of positions at which the corresponding elements are different.

### 7. Mahalanobis Distance:

- Mahalanobis distance considers the correlations between dimensions in the data.
- It is particularly useful when dealing with data with different units or scales.

### 8. Kullback-Leibler Divergence (KL Divergence):

- KL divergence measures the difference between two probability distributions.
- It is often used in probabilistic models, such as in information retrieval and topic modeling.

The choice of distance metric depends on the nature of the data and the specific problem you are trying to solve. Different distance metrics are more appropriate for different types of data

## Classification using Nearest Neighbors

Classification using Nearest Neighbors is a simple but effective machine learning algorithm that falls under the category of instance-based or lazy learning. It is often referred to as k-Nearest Neighbors (k-NN) because it classifies data points based on the majority class among their k-nearest neighbors in the training dataset. Here's how the k-NN classification algorithm works:

1. **Choose the Number of Neighbors (k):** The first step is to select the number of neighbors (k) to consider when making a prediction. This is typically a hyperparameter that you can tune to achieve the best performance for your specific problem.
2. **Calculate Distances:** For a given data point in the test dataset that you want to classify, calculate the distances between that point and all data points in the training dataset. Common distance metrics include Euclidean distance, Manhattan distance, or cosine similarity.
3. **Find Nearest Neighbors:** Sort the calculated distances in ascending order and select the k data points (neighbors) with the smallest distances to the test data point.
4. **Majority Vote:** Determine the class labels of the k nearest neighbors. In a binary classification problem, for example, count the number of neighbors belonging to each class. The class that appears most frequently among the k neighbors is the predicted class for the test data point. In multi-class classification, you can use methods like weighted voting to handle ties.
5. **Make a Prediction:** Assign the predicted class label to the test data point based on the majority class among its k-nearest neighbors.
6. **Repeat for all Test Data:** Repeat the above steps for all data points in the test dataset to obtain predictions for the entire dataset.

Here's a Python code example using the scikit-learn library to perform k-NN classification:

```
from sklearn.neighbors import KNeighborsClassifier

# Create a k-NN classifier with k=3 (you can adjust k)
knn = KNeighborsClassifier(n_neighbors=3)

# Fit the classifier to your training data
knn.fit(X_train, y_train)

# Make predictions on test data
y_pred = knn.predict(X_test)

# Evaluate the model's performance (e.g., accuracy, precision, recall, etc.)
```

Remember that the choice of the distance metric, the value of k, and other hyperparameters can significantly impact the performance of the k-NN algorithm. It's essential to experiment with different values and possibly perform hyperparameter tuning to optimize the model for your

specific dataset. Additionally, preprocessing the data (e.g., scaling features) can also influence the results.

Certainly, let's delve deeper into the k-Nearest Neighbors (k-NN) classification algorithm:

**1. Choose the Number of Neighbors (k):** The choice of the number of neighbors, represented by 'k', is a crucial decision in the k-NN algorithm. A smaller value of k (e.g., 1 or 3) makes the algorithm more sensitive to noise in the data, potentially leading to overfitting. A larger value of k (e.g., 10 or 20) makes the algorithm more robust to noise but may smooth out decision boundaries, potentially leading to underfitting. You should experiment with different values of k to find the one that works best for your dataset.

**2. Calculate Distances:** To determine the k-nearest neighbors for a test data point, the algorithm computes a distance metric between that point and all data points in the training set. Common distance metrics include:

- **Euclidean Distance:** This is the straight-line distance between two points in Euclidean space.
- **Manhattan Distance:** Also known as city block distance or L1 norm, it is the sum of the absolute differences between the coordinates of two points.
- **Cosine Similarity:** Used when dealing with text data or high-dimensional data, it measures the cosine of the angle between two vectors.

The choice of distance metric depends on the nature of your data and the problem you're solving.

**3. Find Nearest Neighbors:** After calculating distances, the algorithm selects the k data points with the smallest distances to the test data point. These are considered the "nearest neighbors" of the test point.

**4. Majority Vote:** Once you have identified the k-nearest neighbors, you need to determine their class labels. In a binary classification scenario (two classes), you can simply count the number of neighbors belonging to each class. The class with the majority of neighbors is assigned as the predicted class for the test data point. For example, if you have k=5 and 3 neighbors belong to class A while 2 neighbors belong to class B, the test point is classified as class A. In multi-class classification, you can use a similar approach but may employ techniques like weighted voting if ties occur.

**5. Make a Prediction:** Based on the majority class among the k-nearest neighbors, the algorithm assigns the predicted class label to the test data point.

**6. Repeat for all Test Data:** The process is repeated for all data points in the test dataset, generating predictions for each one.

**Evaluation:** After making predictions for the entire test dataset, you should evaluate the performance of your k-NN classifier using appropriate metrics such as accuracy, precision, recall, F1-score, or confusion matrices. This evaluation helps you assess how well your model is performing on unseen data.

It's important to note that k-NN is a non-parametric algorithm, meaning it doesn't make assumptions about the underlying data distribution. It can be used for both classification and regression tasks. However, it can be computationally expensive, especially with large datasets,

because it requires calculating distances to all training data points for each test data point. To mitigate this, various optimization techniques and data structures like KD-trees or Ball Trees are used in practice.

Certainly, let's delve deeper into the k-Nearest Neighbors (k-NN) classification algorithm:

**1. Choose the Number of Neighbors (k):** The choice of the number of neighbors, represented by 'k', is a crucial decision in the k-NN algorithm. A smaller value of k (e.g., 1 or 3) makes the algorithm more sensitive to noise in the data, potentially leading to overfitting. A larger value of k (e.g., 10 or 20) makes the algorithm more robust to noise but may smooth out decision boundaries, potentially leading to underfitting. You should experiment with different values of k to find the one that works best for your dataset.

**2. Calculate Distances:** To determine the k-nearest neighbors for a test data point, the algorithm computes a distance metric between that point and all data points in the training set. Common distance metrics include:

- **Euclidean Distance:** This is the straight-line distance between two points in Euclidean space.
- **Manhattan Distance:** Also known as city block distance or L1 norm, it is the sum of the absolute differences between the coordinates of two points.
- **Cosine Similarity:** Used when dealing with text data or high-dimensional data, it measures the cosine of the angle between two vectors.

The choice of distance metric depends on the nature of your data and the problem you're solving.

**3. Find Nearest Neighbors:** After calculating distances, the algorithm selects the k data points with the smallest distances to the test data point. These are considered the "nearest neighbors" of the test point.

**4. Majority Vote:** Once you have identified the k-nearest neighbors, you need to determine their class labels. In a binary classification scenario (two classes), you can simply count the number of neighbors belonging to each class. The class with the majority of neighbors is assigned as the predicted class for the test data point. For example, if you have k=5 and 3 neighbors belong to class A while 2 neighbors belong to class B, the test point is classified as class A. In multi-class classification, you can use a similar approach but may employ techniques like weighted voting if ties occur.

**5. Make a Prediction:** Based on the majority class among the k-nearest neighbors, the algorithm assigns the predicted class label to the test data point.

**6. Repeat for all Test Data:** The process is repeated for all data points in the test dataset, generating predictions for each one.

**Evaluation:** After making predictions for the entire test dataset, you should evaluate the performance of your k-NN classifier using appropriate metrics such as accuracy, precision, recall, F1-score, or confusion matrices. This evaluation helps you assess how well your model is performing on unseen data.

It's important to note that k-NN is a non-parametric algorithm, meaning it doesn't make assumptions about the underlying data distribution. It can be used for both classification and regression tasks. However, it can be computationally expensive, especially with large datasets, because it requires calculating distances to all training data points for each test data point. To

```
In [11]: # Import necessary libraries
from sklearn.neighbors import KNeighborsClassifier
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load a sample dataset (Iris dataset in this case)
data = load_iris()
X = data.data # Features (sepal length, sepal width, petal length, petal width)
y = data.target # Target variable (species)

# Split the data into a training set and a test set (80% training, 20% testing)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create a k-NN classifier with k=3 (you can adjust k)
knn = KNeighborsClassifier(n_neighbors=3)

# Fit the classifier to your training data
knn.fit(X_train, y_train)

# Make predictions on the test data
y_pred = knn.predict(X_test)

# Evaluate the model's performance using accuracy as an example metric
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.2f}")

# You can also use other evaluation metrics as needed (e.g., precision, recall)
```

Accuracy: 1.00

In this example:

We import the necessary libraries, including scikit-learn for the k-NN classifier, a sample dataset (the Iris dataset), and functions for data splitting and evaluation.

We load the Iris dataset, which contains features (sepal length, sepal width, petal length, and petal width) and target labels (species of Iris flowers).

We split the dataset into training and testing sets using the `train_test_split` function.

We create a k-NN classifier with `n_neighbors=3`. You can adjust the value of `n_neighbors` to experiment with different values of `k`.

We fit the k-NN classifier to the training data using the `fit` method.

We make predictions on the test data using the `predict` method.



We evaluate the model's performance using accuracy as an example metric. You can use other evaluation metrics like precision, recall, or F1-score depending on your specific classification problem.

This example demonstrates the basic steps of using k-NN for classification in Python. You can adapt this code to your own dataset and classification problem by replacing the Iris dataset with your data and adjusting the classifier and evaluation metrics accordingly.

## Understanding train test dataset

Understanding the concepts of train and test datasets is fundamental in machine learning, as they are essential components of the model development and evaluation process. Let's explore these concepts:

### 1. Training Dataset:

- The training dataset is a subset of your data that is used to train your machine learning model. It contains examples of input data (features) and their corresponding known target values (labels or outcomes).
- The purpose of the training dataset is to allow the model to learn the patterns and relationships within the data. During training, the model adjusts its internal parameters to minimize the difference between its predictions and the actual target values in the training data.
- Typically, a larger training dataset leads to better model performance, but it also requires more computational resources and time.

### 2. Test Dataset:

- The test dataset is a separate subset of your data that is not used during the training process. It contains examples of input data, just like the training dataset, but the corresponding target values are kept hidden from the model.
- The purpose of the test dataset is to evaluate how well your trained machine learning model generalizes to new, unseen data. It helps you estimate the model's performance on real-world, future data.
- It's important that the test dataset is not used for training in any way, as using it during training can lead to overfitting, where the model performs well on the test data but poorly on new, unseen data.

### 3. Validation Dataset (Optional):

- In addition to the training and test datasets, some machine learning workflows use a validation dataset. This dataset is used during the model development and tuning process.
- The validation dataset helps you fine-tune model hyperparameters and make decisions about model architecture. It's often used to prevent overfitting by providing an independent dataset for assessing the model's performance during training.
- After hyperparameter tuning is complete, you typically retrain your model using the full training dataset (including the validation data) before evaluating it on the test dataset.

Here's a typical workflow involving these datasets:

1. **Data Splitting:** You start with a single dataset that is split into three subsets: training, validation (if used), and test datasets. The most common split ratio is something like 70-80% for training, 10-15% for validation, and 10-15% for testing, but this can vary based on the specific project and dataset size.
2. **Model Training:** You use the training dataset to train your machine learning model. The model learns from the features and target values in this dataset.
3. **Model Validation (Optional):** If a validation dataset is used, you evaluate your model's performance on it during training to make hyperparameter and architecture adjustments.
4. **Model Testing:** Once you are satisfied with your model's performance on the validation dataset (if used), you evaluate it on the test dataset. This provides an unbiased assessment of how well your model generalizes to new, unseen data.

The goal of this process is to build a model that not only performs well on the training data (which it has seen before) but also on new data it hasn't encountered during training. This ensures that your machine learning model can make accurate predictions in real-world scenarios.

## Training and predicting result

Training and predicting are fundamental steps in machine learning and predictive modeling. Here's a high-level overview of these processes:

### 1. Training:

- **Data Collection:** Start by collecting a dataset that contains examples of the problem you want to solve. This dataset should include input features (variables) and corresponding target labels (the values you want to predict).
- **Data Preprocessing:** Clean and preprocess the data to handle missing values, outliers, and format issues. This may involve techniques like data imputation, scaling, encoding categorical variables, and feature engineering.
- **Splitting the Data:** Divide your dataset into two or more subsets. Common splits are training, validation, and test sets. The training set is used to train your model, the validation set helps you tune hyperparameters, and the test set is used to evaluate the final model's performance.
- **Choosing a Model:** Select an appropriate machine learning algorithm or model for your problem. The choice depends on the nature of your data and the problem type (classification, regression, clustering, etc.).
- **Training the Model:** Feed the training data into your chosen model, which will learn patterns and relationships in the data. The model iteratively adjusts its parameters to minimize a loss or error function.

### 2. Predicting:

- **Data Preprocessing:** Similar to the training phase, preprocess any new data you want to make predictions on to ensure it has the same format as your training data.
- **Model Inference:** Use the trained model to make predictions on new, unseen data. For example, if you're doing classification, the model will output class probabilities or labels. In regression, it will produce numeric predictions.

- **Post-processing:** Depending on your application, you may need to post-process the model's predictions. This can involve rounding, thresholding, or converting predictions into actionable recommendations.

### 3. Evaluating Results:

- Measure the performance of your model using appropriate evaluation metrics. Common metrics include accuracy, precision, recall, F1-score, mean squared error (MSE), and many others. The choice of metric depends on the problem you're solving.
- Use the test set (or validation set if you have one) to assess how well your model generalizes to new data.
- Iterate: If the model's performance is not satisfactory, you may need to revisit the training phase, adjust hyperparameters, collect more data, or try different algorithms.

### 4. Deployment:

- If your model performs well, you can deploy it in a real-world application. This may involve integrating it into a software system, creating a user interface, or setting up automated pipelines for continuous prediction.
- Monitor the model's performance in the deployed environment and retrain it periodically with new data to ensure it remains accurate over time.

Remember that machine learning is an iterative process, and fine-tuning models and data preprocessing are often required to achieve the best results. Additionally, ethical considerations, bias mitigation, and data privacy should also be taken into account throughout

```
In [12]: from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = datasets.load_iris()
X = iris.data # Features
y = iris.target # Target Labels

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Choose and train a model (RandomForestClassifier in this example)
model = RandomForestClassifier()
model.fit(X_train, y_train)

# Make predictions on the test data
y_pred = model.predict(X_test)

# Evaluate the model's performance
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
```

Accuracy: 1.0

If you want to use built-in datasets that come with popular Python libraries like scikit-learn or seaborn, you can do so easily. I'll show you how to use a built-in dataset from scikit-learn as an example. Scikit-learn provides several datasets for practice and learning purposes. In this example, we'll use the Iris dataset, which is a classic dataset for classification.

Here's how you can load and work with a built-in dataset like Iris:

```
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = datasets.load_iris()
X = iris.data # Features
y = iris.target # Target Labels

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
0.2, random_state=42)

# Choose and train a model (RandomForestClassifier in this example)
model = RandomForestClassifier()
model.fit(X_train, y_train)

# Make predictions on the test data
y_pred = model.predict(X_test)

# Evaluate the model's performance
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
```

In this code:

1. We import the necessary libraries, including scikit-learn and its datasets module.
2. We load the Iris dataset using `datasets.load_iris()`.
3. We split the dataset into training and testing sets using `train_test_split`.
4. We choose a machine learning model (RandomForestClassifier in this case) and train it using the training data.
5. We use the trained model to make predictions on the test data.
6. Finally, we evaluate the model's accuracy on the test data.

You can replace the Iris dataset with other built-in datasets provided by scikit-learn or use datasets from other libraries like seaborn or Matplotlib for different types of analyses or tasks. The process is similar; you load the dataset, split it into training and testing sets, train a model, make predictions, and evaluate its performance.

# Understanding matrices in ml

Matrices are fundamental mathematical structures used extensively in machine learning and various other fields of mathematics and science. In machine learning, matrices play a critical role in representing and manipulating data, as well as in defining the parameters and operations of many algorithms. Here's an overview of how matrices are used in machine learning:

## 1. Matrix Basics:

- A matrix is a two-dimensional array of numbers, symbols, or expressions arranged in rows and columns.
- The size of a matrix is represented as "m x n," where "m" is the number of rows, and "n" is the number of columns.
- Elements in a matrix are often denoted as  $A[i][j]$ , where "i" is the row index, and "j" is the column index.

## 2. Data Representation:

- In machine learning, datasets are often represented as matrices. Each row of the matrix represents a data point, and each column represents a feature or attribute of that data point.
- For example, in a dataset of house prices, each row could represent a house, and each column could represent features like the number of bedrooms, square footage, and price.

## 3. Matrix Operations:

- Matrices support various operations that are crucial in machine learning, such as addition, subtraction, multiplication, and transposition.
- Matrix addition and subtraction are performed element-wise, meaning corresponding elements in two matrices are added or subtracted from each other.

## 4. Matrix Multiplication:

- Matrix multiplication is a fundamental operation in machine learning, used in various contexts, including linear transformations, neural networks, and more.
- In matrix multiplication, the dot product of rows and columns is calculated to produce a new matrix. It's important to ensure that the dimensions of the matrices are compatible for multiplication (the number of columns in the first matrix must match the number of rows in the second matrix).

## 5. Matrix Transposition:

- Transposing a matrix means swapping its rows and columns. It's denoted as  $A^T$ .
- Transposition is useful for operations like calculating the inner product of two matrices or reshaping data.

## 6. Matrix Inversion and Determinant:

- In some machine learning algorithms, matrix inversion and determinant calculations are required. However, not all matrices are invertible.
- Determinants provide information about the scaling factor of the transformation represented by a matrix.

## 7. Eigenvalues and Eigenvectors:

- Eigenvalues and eigenvectors of a matrix are important in various machine learning algorithms, such as Principal Component Analysis (PCA) and spectral clustering.
- They represent specific characteristics of the matrix and can be used for dimensionality reduction and feature extraction.

## 8. Solving Linear Equations:

- Matrices are used to solve systems of linear equations, which often arise in machine learning problems, such as linear regression.

## 9. Matrix Decompositions:

- Techniques like Singular Value Decomposition (SVD) and QR decomposition are used to break down matrices into simpler components, which can be useful for dimensionality reduction and data compression.

## 10. Convolutional Neural Networks (CNNs):

- In deep learning, matrices are central to convolutional layers in CNNs. These layers perform operations like convolution and pooling on multi-dimensional matrices (image data).

Understanding matrices and their operations is crucial for working with data and algorithms in machine learning. It forms the foundation for linear algebra, which is a key mathematical discipline in the field. If you're getting started with machine learning, it's essential to have a solid grasp of matrices and their applications.

In Python, you can work with matrices using various libraries, with NumPy being one of the most popular and widely used for numerical and matrix operations. Here's how you can work with matrices in Python using NumPy:

## Installation:

If you haven't already installed NumPy, you can do so using pip:

```
pip install numpy
```

## Importing NumPy:

```
import numpy as np
```

## Creating Matrices:

You can create matrices in Python using NumPy in several ways:

### 1. From Lists or Arrays:

```
# Creating a 2x3 matrix from a list of lists  
matrix = np.array([[1, 2, 3], [4, 5, 6]])
```

### 2. Zeros and Ones Matrices:

```
# Creating a 3x3 matrix of zeros
```

```
zeros_matrix = np.zeros((3, 3))
```

```
# Creating a 2x2 matrix of ones
```

```
ones_matrix = np.ones((2, 2))
```

### 3. Identity Matrix:

```
# Creating a 4x4 identity matrix
```

```
identity_matrix = np.identity(4)
```

## Basic Matrix Operations:

Once you have created matrices, you can perform various operations on them:

### 1. Matrix Addition and Subtraction:

```
A = np.array([[1, 2], [3, 4]])
```

```
B = np.array([[5, 6], [7, 8]])
```

```
C = A + B # Matrix addition
```

```
D = A - B # Matrix subtraction
```

### 2. Matrix Multiplication:

```
A = np.array([[1, 2], [3, 4]])
```

```
B = np.array([[5, 6], [7, 8]])
```

```
E = np.dot(A, B) # Matrix multiplication
```

### 3. Matrix Transposition:

```
A = np.array([[1, 2], [3, 4]])
```

```
A_transpose = np.transpose(A) # Transpose A
```

### 4. Matrix Inversion:

NumPy provides functions for matrix inversion, such as `np.linalg.inv()`, but be cautious as not all matrices are invertible.

## Additional Matrix Operations:

### 1. Eigenvalues and Eigenvectors:

NumPy can compute eigenvalues and eigenvectors of a matrix using `np.linalg.eig()`.

### 2. Solving Linear Equations:

NumPy can solve systems of linear equations using `np.linalg.solve()`.

### 3. Matrix Decompositions:

NumPy provides functions for various matrix decompositions, such as SVD (`np.linalg.svd()`), QR decomposition (`np.linalg.qr()`), and Cholesky decomposition (`np.linalg.cholesky()`).

## Example:

Here's a complete example of creating, multiplying, and transposing matrices in Python using NumPy:

```
import numpy as np

# Create two matrices
A = np.array([[1, 2, 3], [4, 5, 6]])
B = np.array([[7, 8], [9, 10], [11, 12]])

# Multiply matrices A and B
C = np.dot(A, B)

# Transpose matrix C
C_transpose = np.transpose(C)

print("Matrix C:")
print(C)
print("\nTransposed Matrix C:")
print(C_transpose)
```

This is just a basic introduction to working with matrices in Python using NumPy. NumPy provides a wide range of functions and capabilities for advanced matrix operations, making it a

## Classification matrices in ml

Classification matrices are essential tools in machine learning for evaluating the performance of classification models. They provide insights into how well a model is performing in terms of making correct and incorrect predictions. The most commonly used classification matrices include:

### 1. Confusion Matrix:

- A confusion matrix is the foundation of classification matrices.
- It is a square matrix that compares the actual (ground truth) class labels with the predicted class labels.
- It consists of four components:
  - True Positives (TP): The number of correctly predicted positive instances.
  - True Negatives (TN): The number of correctly predicted negative instances.
  - False Positives (FP): The number of negative instances predicted as positive.
  - False Negatives (FN): The number of positive instances predicted as negative.

	Predicted Positive	Predicted Negative
Actual Positive	True Positives (TP)	False Negatives (FN)
Actual Negative	False Positives (FP)	True Negatives (TN)

### 2. Accuracy:



- Accuracy is a straightforward metric that calculates the ratio of correctly predicted instances to the total number of instances.
- $\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$
- It is not suitable for imbalanced datasets.

### 3. Precision:

- Precision measures the accuracy of positive predictions. It answers the question: "Of all the instances predicted as positive, how many were actually positive?"
- $\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$

### 4. Recall (Sensitivity or True Positive Rate):

- Recall measures the ability of the model to correctly identify positive instances. It answers the question: "Of all the actual positive instances, how many were correctly predicted?"
- $\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$

### 5. F1-Score:

- The F1-score is the harmonic mean of precision and recall, providing a balance between the two.
- $\text{F1-Score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$

### 6. Specificity (True Negative Rate):

- Specificity measures the ability of the model to correctly identify negative instances.
- $\text{Specificity} = \text{TN} / (\text{TN} + \text{FP})$

### 7. False Positive Rate (FPR):

- FPR is the complement of specificity and measures the proportion of negative instances incorrectly classified as positive.
- $\text{FPR} = \text{FP} / (\text{TN} + \text{FP})$

### 8. Receiver Operating Characteristic (ROC) Curve:

- ROC curve is a graphical representation of a classifier's performance across various thresholds.
- It plots the True Positive Rate (Recall) against the False Positive Rate at different threshold values.
- The area under the ROC curve (AUC-ROC) is a summary metric, where a higher value indicates better performance.

### 9. Precision-Recall Curve:

- The precision-recall curve plots precision against recall at different threshold values.
- It is especially useful for imbalanced datasets.

### 10. Matthews Correlation Coefficient (MCC):

- MCC takes into account all four values from the confusion matrix and produces a score that ranges from -1 (completely wrong predictions) to +1 (perfect predictions).
- It is robust for imbalanced datasets.
- $\text{MCC} = (\text{TP} * \text{TN} - \text{FP} * \text{FN}) / \sqrt{(\text{TP} + \text{FP}) * (\text{TP} + \text{FN}) * (\text{TN} + \text{FP}) * (\text{TN} + \text{FN})}$

These classification matrices and metrics help assess the performance of classification models, allowing you to choose the most appropriate evaluation criteria based on the problem's context and the importance of precision and recall in your application.

In [13]: `from sklearn.metrics import confusion_matrix, accuracy_score, precision_score,`

```
# Example ground truth and predicted labels
y_true = [1, 0, 1, 1, 0, 0, 1, 0] # y_test
y_pred = [1, 0, 0, 1, 1, 0, 1, 1] # prediction value

# Confusion Matrix
confusion = confusion_matrix(y_true, y_pred)
print("Confusion Matrix:\n", confusion)

# Accuracy
accuracy = accuracy_score(y_true, y_pred)
print("Accuracy:", accuracy)

# Precision
precision = precision_score(y_true, y_pred)
print("Precision:", precision)

# Recall
recall = recall_score(y_true, y_pred)
print("Recall:", recall)

# F1-Score
f1 = f1_score(y_true, y_pred)
print("F1-Score:", f1)

# ROC Curve and AUC
fpr, tpr, thresholds = roc_curve(y_true, y_pred)
roc_auc = auc(fpr, tpr)
print("AUC-ROC:", roc_auc)

# Matthews Correlation Coefficient
mcc = matthews_corrcoef(y_true, y_pred)
print("Matthews Correlation Coefficient:", mcc)
```

Confusion Matrix:

```
[[2 2]
```

```
[1 3]]
```

Accuracy: 0.625

Precision: 0.6

Recall: 0.75

F1-Score: 0.6666666666666665

AUC-ROC: 0.625

Matthews Correlation Coefficient: 0.2581988897471611

```
In [14]: import numpy as np

def nearest_neighbors(query_point, data_points, k=1):
    """
    Find the k nearest neighbors of a query point among a set of data points.

    Parameters:
    - query_point: The query point for which nearest neighbors are to be found
    - data_points: An array of data points to search for neighbors in.
    - k: The number of nearest neighbors to find (default is 1).

    Returns:
    - indices: Indices of the k nearest neighbors in the data_points array.
    - distances: Distances from the query point to the k nearest neighbors.
    """
    # Calculate distances between the query point and all data points
    distances = np.linalg.norm(data_points - query_point, axis=1)

    # Find the indices of the k smallest distances
    indices = np.argsort(distances)[:k]

    return indices, distances[indices]

# Example usage:
if __name__ == "__main__":
    # Generate some random data points
    np.random.seed(0)
    data_points = np.random.rand(10, 2)

    # Query point for which we want to find the nearest neighbors
    query_point = np.array([0.5, 0.5])

    # Find the 3 nearest neighbors
    k = 3
    nearest_indices, nearest_distances = nearest_neighbors(query_point, data_p

    print(f"The {k} nearest neighbors to the query point {query_point} are:")
    for i, idx in enumerate(nearest_indices):
        print(f"Neighbor {i+1}: Data point {data_points[idx]}, Distance: {near
```

```
The 3 nearest neighbors to the query point [0.5 0.5] are:
Neighbor 1: Data point [0.60276338 0.54488318], Distance: 0.11213746732281339
Neighbor 2: Data point [0.4236548 0.64589411], Distance: 0.16466232687378476
Neighbor 3: Data point [0.5488135 0.71518937], Distance: 0.22065634268116185
```

## Regression in nearest neighbors

Regression in nearest neighbors is a machine learning technique used for predicting continuous numerical values, such as prices, temperatures, or scores, based on the similarity between data points. It is an extension of the k-nearest neighbors (KNN) algorithm, which is primarily used for classification tasks. Here's a brief overview of how regression works in nearest neighbors:

1. **K-Nearest Neighbors (KNN):** KNN is a supervised learning algorithm used for both classification and regression tasks. In the context of regression, it works as follows:
  - Given a new data point that you want to predict a numerical value for, KNN identifies the  $k$  nearest data points in the training dataset based on a distance metric (usually Euclidean distance).
  - It then calculates the average (or weighted average) of the target values of these  $k$  neighbors. This average is used as the predicted value for the new data point.
2. **Distance Metric:** The choice of distance metric is crucial in nearest neighbors regression. Common distance metrics include Euclidean distance, Manhattan distance, and Minkowski distance. The metric determines how the similarity or distance between data points is calculated.
3. **Hyperparameter Selection:** You need to specify the value of the hyperparameter ' $k$ ' in KNN. The choice of ' $k$ ' can significantly impact the model's performance. Smaller values of ' $k$ ' can lead to noisy predictions, while larger values can result in overly smoothed predictions.
4. **Prediction:** Once the nearest neighbors are identified, and their target values are averaged, this average becomes the predicted numerical value for the input data point. In the case of weighted KNN, closer neighbors might contribute more to the prediction than those farther away.
5. **Model Evaluation:** To assess the performance of a nearest neighbors regression model, you typically use evaluation metrics like Mean Squared Error (MSE), Mean Absolute Error (MAE), or R-squared ( $R^2$ ). These metrics help you quantify the accuracy of the model's predictions.
6. **Scalability:** One limitation of KNN regression is that it can be computationally expensive, especially with large datasets, as it requires calculating distances to all data points. Various techniques, such as KD-trees and ball trees, can be used to speed up the search for nearest neighbors.

In summary, nearest neighbors regression is a simple yet effective algorithm for predicting numerical values based on the similarity between data points. It is particularly useful when you have limited features and a relatively small dataset. However, it may not perform well in high-dimensional spaces or with noisy data, and the choice of the ' $k$ ' value and distance metric should be carefully considered.

Certainly! You can use built-in datasets from scikit-learn to perform nearest neighbors regression. One of the commonly used datasets for regression tasks in scikit-learn is the Boston Housing dataset. Here's how you can perform nearest neighbors regression using this dataset:

```

# Import necessary libraries
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_boston
from sklearn.neighbors import KNeighborsRegressor
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import train_test_split

# Load the Boston Housing dataset
data = load_boston()
X = data.data
y = data.target

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
0.2, random_state=42)

# Create a K-Nearest Neighbors Regressor
k = 5 # Number of neighbors
knn_regressor = KNeighborsRegressor(n_neighbors=k)

# Fit the model on the training data
knn_regressor.fit(X_train, y_train)

# Make predictions on the test data
y_pred = knn_regressor.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
print(f"Mean Squared Error: {mse:.2f}")

# Visualize the predictions
plt.scatter(y_test, y_pred)
plt.xlabel("Actual Prices")
plt.ylabel("Predicted Prices")
plt.title("Actual Prices vs. Predicted Prices")
plt.show()

```

In this code:

- We load the Boston Housing dataset using scikit-learn's `load_boston` function.
- We split the data into training and testing sets

```

In [15]: # Import necessary libraries
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_boston
from sklearn.neighbors import KNeighborsRegressor
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import train_test_split

# Load the Boston Housing dataset
data = load_boston()
X = data.data
y = data.target

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create a K-Nearest Neighbors Regressor
k = 5 # Number of neighbors
knn_regressor = KNeighborsRegressor(n_neighbors=k)

# Fit the model on the training data
knn_regressor.fit(X_train, y_train)

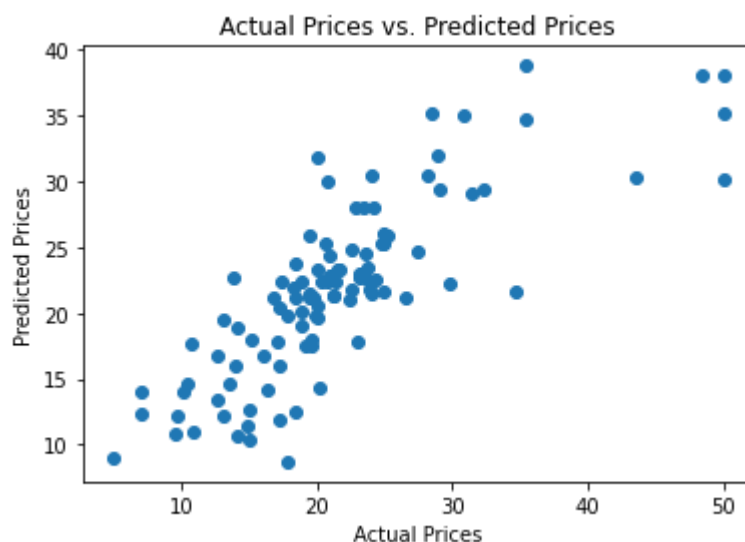
# Make predictions on the test data
y_pred = knn_regressor.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
print(f"Mean Squared Error: {mse:.2f}")

# Visualize the predictions
plt.scatter(y_test, y_pred)
plt.xlabel("Actual Prices")
plt.ylabel("Predicted Prices")
plt.title("Actual Prices vs. Predicted Prices")
plt.show()

```

Mean Squared Error: 25.86



It seems like there might be a typo in your question. I believe you are referring to "scikit-learn," not "scikit learning." Scikit-learn is a popular machine learning library in Python that provides a wide range of tools for working with machine learning algorithms, data preprocessing, model evaluation, and more. It's often abbreviated as "sklearn."

Here's an overview of key concepts and steps involved in using scikit-learn for machine learning tasks:

1. **Importing the Library:** To use scikit-learn, you first need to import it:

```
import sklearn
```

2. **Loading the Dataset:** Scikit-learn provides some built-in datasets for practice. You can also load your own dataset using tools like Pandas or NumPy. Commonly used functions to load datasets are `load_iris`, `load_digits`, and `fetch_openml`, among others.

```
from sklearn.datasets import load_iris
data = load_iris()
X, y = data.data, data.target
```

3. **Data Preprocessing:** Before feeding the data into machine learning models, you often need to preprocess it. This can include tasks like handling missing values, scaling features, encoding categorical variables, and splitting the data into training and testing sets.

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

4. **Choosing a Model:** Scikit-learn offers a wide range of machine learning algorithms for classification, regression, clustering, and more. You'll need to choose an appropriate model for your specific task.

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
```

5. **Training the Model:** You train the model using the training data by calling the `fit` method.

```
model.fit(X_train, y_train)
```

6. **Making Predictions:** After training, you can use the trained model to make predictions on new data.

```
y_pred = model.predict(X_test)
```

7. **Evaluating the Model:** You should assess the model's performance using various metrics such as accuracy, precision, recall, F1-score, and more, depending on the problem type (classification, regression, etc.).

```
from sklearn.metrics import accuracy_score
accuracy = accuracy_score(y_test, y_pred)
```

8. **Hyperparameter Tuning:** You can fine-tune the model's hyperparameters using techniques like grid search or randomized search to improve performance.

```
from sklearn.model_selection import GridSearchCV
param_grid = {'n_estimators': [50, 100, 200]}
grid_search = GridSearchCV(model, param_grid, cv=5)
grid_search.fit(X_train, y_train)
best_model = grid_search.best_estimator_
```

9. **Deploying the Model:** Once satisfied with the model's performance, you can deploy it for making predictions on new, unseen data.

These are the fundamental steps to use scikit-learn for machine learning tasks. Depending on your specific problem, you may need to explore additional features and techniques provided by scikit-learn, such as feature selection, dimensionality reduction, and more advanced model evaluation methods. Additionally, scikit-learn's documentation and online tutorials are excellent

Scikit-learn provides several built-in datasets that you can use for practice and experimentation. These datasets cover a variety of machine learning tasks, including classification, regression, and clustering. Here are some additional datasets available in scikit-learn:

1. **Breast Cancer Wisconsin (Diagnostic) Dataset:**

- Task: Binary classification
- Description: Contains features computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. The goal is to predict whether the tumor is malignant or benign.
- Loading:

```
from sklearn.datasets import load_breast_cancer
data = load_breast_cancer()
```

2. **Diabetes Dataset:**

- Task: Regression
- Description: Contains ten baseline variables (age, sex, BMI, average blood pressure, etc.) and a quantitative measure of disease progression one year after baseline. The goal is to predict disease progression.
- Loading:

```
from sklearn.datasets import load_diabetes
data = load_diabetes()
```

3. **Iris Dataset** (You mentioned this one already):

- Task: Multiclass classification
- Description: Contains measurements of sepal length, sepal width, petal length, and petal width for three species of iris flowers. The goal is to classify iris flowers into species.
- Loading:



```
from sklearn.datasets import load_iris
data = load_iris()
```

#### 4. Digits Dataset:

- Task: Multiclass classification
- Description: Consists of 8x8 pixel images of handwritten digits (0 through 9). The goal is to recognize and classify the digits.
- Loading:

```
from sklearn.datasets import load_digits
data = load_digits()
```

#### 5. Wine Dataset:

- Task: Multiclass classification
- Description: Contains measurements of various chemical constituents found in different types of wines. The goal is to classify the wines into one of three classes.
- Loading:

```
from sklearn.datasets import load_wine
data = load_wine()
```

#### 6. Labeled Faces in the Wild (LFW) Dataset:

- Task: Face recognition (can be used for clustering as well)
- Description: Contains thousands of labeled images of faces collected from the web. The goal is to recognize and identify faces.
- Loading:

```
from sklearn.datasets import fetch_lfw_people
data = fetch_lfw_people(min_faces_per_person=70, resize=0.4)
```

#### 7. Boston Housing Dataset:

- Task: Regression
- Description: Contains housing-related information for various neighborhoods in Boston. The goal is to predict the median value of owner-occupied homes.
- Loading:

```
from sklearn.datasets import load_boston
data = load_boston()
```

#### 8. Olivetti Faces Dataset:

- Task: Face recognition
- Description: Contains a collection of 400 images of 40 distinct subjects, each with 10 different poses. The goal is to recognize and identify individuals' faces.
- Loading:

```
from sklearn.datasets import fetch_olivetti_faces
data = fetch_olivetti_faces(shuffle=True, random_state=42)
```

These datasets provide a diverse set of tasks and challenges for machine learning and are useful for practicing and testing various algorithms and techniques. You can load these

## Application using regression in NN

Regression is a type of supervised machine learning technique used to predict a continuous output variable based on one or more input features. Neural networks (NN) can be applied to regression tasks just as effectively as they can be applied to classification tasks. Here's how you can create an application using regression in a neural network:

### 1. Collect and Preprocess Data:

- Gather a dataset that contains examples with input features and corresponding continuous target values (i.e., the values you want to predict).
- Preprocess the data by handling missing values, scaling features, and splitting the dataset into training and testing sets.

### 2. Design Your Neural Network:

- Decide on the architecture of your neural network. For regression, a simple feedforward neural network with one or more hidden layers can work well.
- The number of input neurons should match the number of features in your dataset, and the output layer should have a single neuron since you're predicting a single continuous value.
- Choose an appropriate activation function for the output layer (often linear or ReLU).

### 3. Compile Your Model:

- Select an appropriate loss function for regression. Common choices include Mean Squared Error (MSE) and Mean Absolute Error (MAE).
- Choose an optimizer (e.g., Adam, SGD) and specify evaluation metrics (e.g., Mean Absolute Error) to monitor during training.

### 4. Train Your Model:

- Feed your training data into the neural network and adjust the model's weights through backpropagation to minimize the chosen loss function.
- Monitor the training process by evaluating the model on the validation set and watching for overfitting.

### 5. Evaluate Your Model:

- Once the model is trained, evaluate its performance on the test dataset using the chosen evaluation metric(s).
- Visualize the predictions versus the actual target values to understand the model's accuracy.

### 6. Make Predictions:

- Use the trained neural network to make predictions on new, unseen data. Input the features into the model, and it will produce a continuous output value as a prediction.

### 7. Fine-Tuning and Hyperparameter Optimization (Optional):

- Experiment with different architectures, activation functions, learning rates, and regularization techniques to improve the model's performance.
- Perform cross-validation and hyperparameter tuning to find the best combination of hyperparameters.

#### 8. Deployment:

- If your model performs well, you can deploy it in your application for real-world predictions.
- Ensure that the input data in your application is preprocessed in the same way as during training.

#### 9. Monitoring and Maintenance:

- Continuously monitor your deployed model's performance and retrain it periodically with new data to keep it up to date.

Here's a simplified Python code snippet using TensorFlow and Keras for regression with a neural network:

```
import tensorflow as tf
from tensorflow import keras

# Load and preprocess your data here

# Define the neural network architecture
model = keras.Sequential([
    keras.layers.Dense(64, activation='relu', input_shape=(num_features,)),
    keras.layers.Dense(32, activation='relu'),
    keras.layers.Dense(1, activation='linear') # Linear activation for regression
])

# Compile the model
model.compile(loss='mean_squared_error', optimizer='adam', metrics=['mean_absolute_error'])

# Train the model
model.fit(X_train, y_train, epochs=100, batch_size=32, validation_data=(X_val, y_val))

# Evaluate the model
test_loss, test_mae = model.evaluate(X_test, y_test)

# Make predictions
predictions = model.predict(X_new_data)
```

This code provides a basic outline for building a regression model using a neural network. You

```

In [16]: import tensorflow as tf
from tensorflow import keras
import numpy as np
import matplotlib.pyplot as plt

# Generate some example data
np.random.seed(0)
X = np.random.rand(100, 1) # Input feature (example: random values between 0
y = 3 * X + 2 + 0.1 * np.random.randn(100, 1) # Target variable with some noi

# Split the data into training and testing sets
X_train, X_test = X[:80], X[80:]
y_train, y_test = y[:80], y[80:]

# Define a simple neural network model for regression
model = keras.Sequential([
    keras.layers.Dense(1, input_shape=(1,), activation='linear')
])

# Compile the model
model.compile(optimizer='adam', loss='mean_squared_error')

# Train the model
history = model.fit(X_train, y_train, epochs=100, verbose=0)

# Evaluate the model on the test data
loss = model.evaluate(X_test, y_test)
print(f"Test Loss: {loss:.4f}")

# Make predictions
predictions = model.predict(X_test)

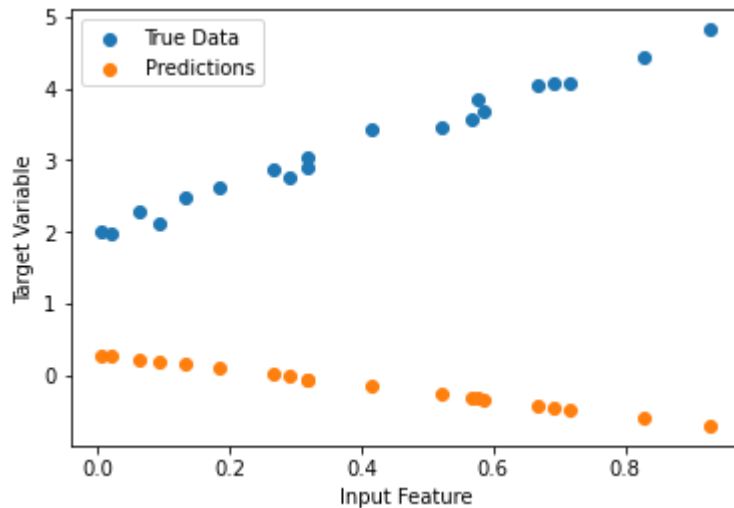
# Plot the results
plt.scatter(X_test, y_test, label='True Data')
plt.scatter(X_test, predictions, label='Predictions')
plt.xlabel('Input Feature')
plt.ylabel('Target Variable')
plt.legend()
plt.show()

```

```

1/1 [=====] - 0s 93ms/step - loss: 12.6244
Test Loss: 12.6244
1/1 [=====] - 0s 84ms/step

```



## Foundations of Clustering

Clustering is a fundamental technique in data analysis and machine learning that involves grouping similar data points together based on some similarity metric. It's widely used in various applications, such as data compression, pattern recognition, customer segmentation, and recommendation systems. Here are some foundational concepts and techniques related to clustering:

1. **Definition of Clustering:** Clustering is the process of partitioning a dataset into groups, called clusters, where data points within the same cluster are more similar to each other than to those in other clusters. The goal is to discover hidden patterns or structures in the data.
2. **Types of Clustering:** There are two main types of clustering:
  - **Hard Clustering:** Each data point belongs exclusively to one cluster.
  - **Soft Clustering (Fuzzy Clustering):** Data points can belong to multiple clusters with varying degrees of membership.
3. **Distance Metrics:** Clustering algorithms rely on distance or similarity metrics to measure the similarity between data points. Common distance metrics include Euclidean distance, Manhattan distance, cosine similarity, and more. The choice of metric depends on the nature of the data and the problem at hand.
4. **Centroid-based Clustering:** In centroid-based clustering, each cluster is represented by a centroid (a central point). K-means is a popular centroid-based clustering algorithm that aims to minimize the sum of squared distances between data points and their cluster centroids.
5. **Hierarchical Clustering:** Hierarchical clustering builds a tree-like structure of clusters, known as a dendrogram, by successively merging or splitting clusters. Agglomerative and divisive are the two main approaches in hierarchical clustering.
6. **Density-based Clustering:** Density-based clustering algorithms, such as DBSCAN (Density-Based Spatial Clustering of Applications with Noise), group data points based on their density within the feature space. It can discover clusters of arbitrary shapes and is robust to noise.

7. **Model-based Clustering:** Model-based clustering algorithms assume that the data is generated from a probabilistic model. Gaussian Mixture Models (GMM) is a popular model-based clustering technique that assumes data points are generated from a mixture of Gaussian distributions.
8. **Evaluation Metrics:** To assess the quality of clustering results, various metrics like Silhouette Score, Davies-Bouldin Index, and internal/external validation measures can be used. These metrics help quantify the separation and cohesion of clusters.
9. **Choosing the Number of Clusters (K):** One of the critical tasks in clustering is determining the optimal number of clusters. Methods like the Elbow Method, Silhouette Analysis, and Gap Statistics can assist in choosing an appropriate value for K.
10. **Preprocessing and Feature Engineering:** Data preprocessing, including handling missing values, scaling features, and reducing dimensionality, can significantly impact clustering results. Proper preprocessing can improve the performance of clustering algorithms.
11. **Cluster Interpretation:** After clustering, interpreting the meaning of the clusters is essential. It involves analyzing the characteristics of data points within each cluster to gain insights and make decisions or recommendations.
12. **Scalability and Efficiency:** The scalability of clustering algorithms is an important consideration for large datasets. Some algorithms are more efficient and suitable for high-dimensional data than others.
13. **Clustering Applications:** Clustering has a wide range of applications, such as customer segmentation in marketing, anomaly detection in cybersecurity, image segmentation in computer vision, and document clustering in natural language processing.

Understanding these foundational concepts and techniques is crucial for effectively applying clustering algorithms to real-world data analysis and machine learning tasks. Different clustering methods should be chosen based on the nature of the data and the objectives of the

In Python, you can implement clustering algorithms using various libraries and packages, such as Scikit-Learn, SciPy, and NumPy. Here, I'll provide a basic example of how to perform K-means clustering using Scikit-Learn, one of the most commonly used machine learning libraries for clustering tasks.

Assuming you have your data in a NumPy array or a Pandas DataFrame, follow these steps to perform K-means clustering:

1. **Import the necessary libraries:**

```
import numpy as np
import pandas as pd
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
```

2. **Load your data:**

You should load your dataset into a NumPy array or a Pandas DataFrame. For example:

```
# Assuming you have a Pandas DataFrame 'data'
X = data.values
```

3. **Choose the number of clusters (K):**

Decide on the number of clusters you want to create. You can use methods like the Elbow Method to help you choose an appropriate value for K.

#### 4. Instantiate the K-means model:

```
kmeans = KMeans(n_clusters=K)
```

#### 5. Fit the model to your data:

```
kmeans.fit(X)
```

#### 6. Get cluster labels and centroids:

After fitting the model, you can obtain the cluster labels for each data point and the cluster centroids.

```
cluster_labels = kmeans.labels_  
cluster_centers = kmeans.cluster_centers_
```

#### 7. Visualize the clusters (optional):

If your data is 2D, you can visualize the clusters by plotting the data points and cluster centroids.

```
plt.scatter(X[:, 0], X[:, 1], c=cluster_labels, cmap='viridis')  
plt.scatter(cluster_centers[:, 0], cluster_centers[:, 1], marker='o',  
            s=200, c='red', label='Centroids')  
plt.legend()  
plt.show()
```

#### 8. Interpret and analyze the clusters:

Analyze the clusters to gain insights from your data. You can examine the characteristics of data points within each cluster to understand their meaning and make decisions or recommendations.

Here's a complete example with some sample data:

```

import numpy as np
import pandas as pd
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

# Sample data
data = np.array([[1, 2], [5, 8], [1.5, 1.8], [8, 8], [1, 0.6], [9, 1
1]])

# Choose the number of clusters
K = 2

# Instantiate the K-means model
kmeans = KMeans(n_clusters=K)

# Fit the model to the data
kmeans.fit(data)

# Get cluster labels and centroids
cluster_labels = kmeans.labels_
cluster_centers = kmeans.cluster_centers_

```

## K means Clustering

K-means clustering is a popular unsupervised machine learning algorithm used for partitioning a dataset into groups or clusters based on the similarity of data points. It is commonly employed for tasks such as customer segmentation, image compression, anomaly detection, and more. Here's a basic overview of how K-means clustering works:

1. **Initialization:** K-means begins by randomly selecting K initial cluster centroids. These centroids represent the center points of the K clusters you want to create.
2. **Assignment:** For each data point in the dataset, the algorithm calculates the distance between that point and each of the K centroids. It assigns the data point to the cluster whose centroid is closest to it. This assignment is based on a distance metric, typically Euclidean distance.
3. **Update:** After all data points have been assigned to clusters, the algorithm calculates new cluster centroids by taking the mean of all the data points in each cluster. These new centroids represent the updated cluster centers.
4. **Repeat:** Steps 2 and 3 are repeated iteratively until a convergence criterion is met. The most common convergence criterion is when the centroids no longer change significantly or when a fixed number of iterations is reached.
5. **Result:** Once the algorithm converges, it returns the final K cluster centroids and the assignment of data points to clusters.

Key points to consider when using K-means clustering:



- **Choosing K:** Selecting the right number of clusters, K, is a crucial step. Various methods, such as the elbow method or silhouette analysis, can help you determine the optimal K value.
- **Initialization:** The choice of initial cluster centroids can affect the algorithm's convergence and the quality of the clusters. Random initialization is common, but more sophisticated methods like K-means++ exist to improve initialization.
- **Scaling:** It's important to scale the features of your dataset before applying K-means, as it's sensitive to the scale of the variables.
- **Sensitive to Initialization:** K-means can converge to local optima depending on the initial centroids. Running the algorithm multiple times with different initializations and selecting the best result is a common practice.
- **Non-Convex Clusters:** K-means assumes that clusters are convex and equally sized, which may not hold in all cases. Other clustering algorithms like DBSCAN or hierarchical clustering can handle more complex cluster shapes.
- **Outliers:** K-means can be sensitive to outliers, as they can significantly impact the centroids and cluster assignments. Consider preprocessing or using robust variants of K-means if outliers are present.
- **Scalability:** K-means can be computationally expensive for large datasets or a large number of clusters. Alternative algorithms like Mini-batch K-means can be more efficient.

In summary, K-means clustering is a widely used method for partitioning data into clusters

```
In [17]: # Import necessary libraries
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs

# Generate synthetic data for demonstration
data, _ = make_blobs(n_samples=300, centers=4, random_state=0, cluster_std=0.6)

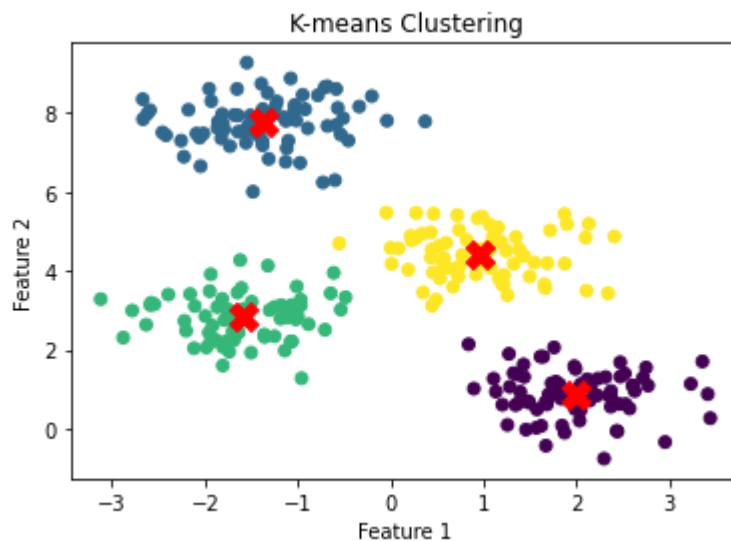
# Create a K-means instance with the desired number of clusters (K)
kmeans = KMeans(n_clusters=4)

# Fit the K-means model to your data
kmeans.fit(data)

# Get cluster assignments for each data point
labels = kmeans.labels_

# Get cluster centroids
centroids = kmeans.cluster_centers_

# Plot the data points and cluster centroids
plt.scatter(data[:, 0], data[:, 1], c=labels)
plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=200, c='red')
plt.title('K-means Clustering')
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.show()
```



To perform K-means clustering in Python, you can use libraries such as scikit-learn, which provide easy-to-use implementations of the K-means algorithm. Here's a step-by-step example of how to perform K-means clustering using scikit-learn:

```

# Import necessary libraries
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs

# Generate synthetic data for demonstration
data, _ = make_blobs(n_samples=300, centers=4, random_state=0, cluster_std=0.60)

# Create a K-means instance with the desired number of clusters (K)
kmeans = KMeans(n_clusters=4)

# Fit the K-means model to your data
kmeans.fit(data)

# Get cluster assignments for each data point
labels = kmeans.labels_

# Get cluster centroids
centroids = kmeans.cluster_centers_

# Plot the data points and cluster centroids
plt.scatter(data[:, 0], data[:, 1], c=labels)
plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=200, c='red')
plt.title('K-means Clustering')
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.show()

```

In this example:

1. We import the necessary libraries, including NumPy, Matplotlib, and scikit-learn.
2. We generate synthetic data using `make_blobs` , but you can replace this with your own dataset.
3. We create a `KMeans` instance with the desired number of clusters (K) using `KMeans(n_clusters=4)` .
4. We fit the K-means model to our data using `kmeans.fit(data)` .
5. We obtain the cluster assignments for each data point using `kmeans.labels_` .
6. We get the cluster centroids using `kmeans.cluster_centers_` .
7. Finally, we plot the data points with their cluster assignments and the cluster centroids.

➤ [K-Means Clustering with Python](#)

```

In [18]: # Import necessary libraries
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.datasets import load_iris

# Load the Iris dataset
iris = load_iris()
data = iris.data

# Create a K-means instance with the desired number of clusters (K)
kmeans = KMeans(n_clusters=3) # In this case, we assume there are 3 clusters

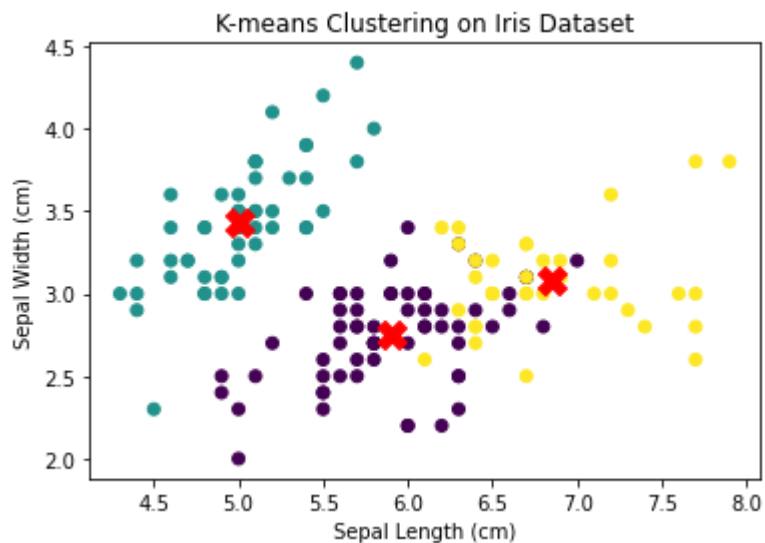
# Fit the K-means model to your data
kmeans.fit(data)

# Get cluster assignments for each data point
labels = kmeans.labels_

# Get cluster centroids
centroids = kmeans.cluster_centers_

# Plot the data points and cluster centroids
plt.scatter(data[:, 0], data[:, 1], c=labels)
plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=200, c='red')
plt.title('K-means Clustering on Iris Dataset')
plt.xlabel('Sepal Length (cm)')
plt.ylabel('Sepal Width (cm)')
plt.show()

```





```

In [19]: import numpy as np

class KMeans:
    def __init__(self, n_clusters, max_iters=100, random_state=None):
        self.n_clusters = n_clusters
        self.max_iters = max_iters
        self.random_state = random_state
        self.centroids = None

    def fit(self, X):
        if self.random_state is not None:
            np.random.seed(self.random_state)

        # Randomly initialize centroids
        initial_indices = np.random.choice(len(X), self.n_clusters, replace=False)
        self.centroids = X[initial_indices]

        for _ in range(self.max_iters):
            # Assign each data point to the nearest centroid
            labels = self._assign_labels(X)

            # Update centroids based on the mean of data points in each cluster
            new_centroids = self._update_centroids(X, labels)

            # Check for convergence
            if np.all(self.centroids == new_centroids):
                break

            self.centroids = new_centroids

        return labels

    def _assign_labels(self, X):
        distances = np.linalg.norm(X[:, np.newaxis] - self.centroids, axis=2)
        return np.argmin(distances, axis=1)

    def _update_centroids(self, X, labels):
        new_centroids = np.array([X[labels == k].mean(axis=0) for k in range(self.n_clusters)])
        return new_centroids

if __name__ == "__main__":
    # Example usage
    # Generate random data for demonstration purposes
    np.random.seed(0)
    X = np.random.randn(300, 2)

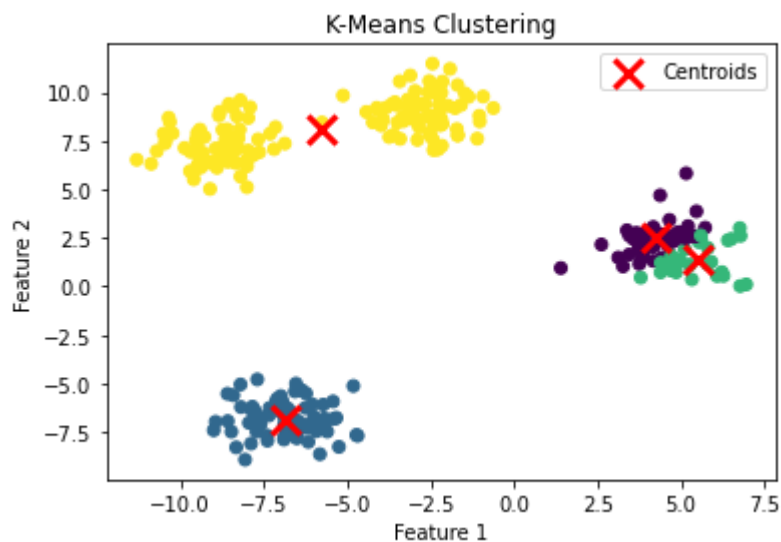
    # Initialize and fit KMeans
    kmeans = KMeans(n_clusters=3, max_iters=100, random_state=42)
    labels = kmeans.fit(X)

    # Print cluster centers
    print("Cluster centers:")
    print(kmeans.centroids)

```

```
Cluster centers:  
[[-0.79665105  0.76370532]  
 [ 0.99439046  0.25558988]  
 [-0.4297535  -0.96921471]]
```

```
In [20]: import numpy as np  
import matplotlib.pyplot as plt  
from sklearn.datasets import make_blobs  
  
# Generate synthetic data for demonstration  
X, _ = make_blobs(n_samples=300, centers=4, random_state=42)  
  
# Initialize and fit KMeans  
kmeans = KMeans(n_clusters=4, max_iters=100, random_state=42)  
labels = kmeans.fit(X)  
  
# Get cluster centers  
centroids = kmeans.centroids  
  
# Plot the data points and cluster centers  
plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis')  
plt.scatter(centroids[:, 0], centroids[:, 1], marker='x', s=200, linewidths=3,  
plt.legend()  
plt.title('K-Means Clustering')  
plt.xlabel('Feature 1')  
plt.ylabel('Feature 2')  
plt.show()
```



K-means clustering is an iterative algorithm that groups a set of data points into K clusters, where each data point belongs to the cluster with the nearest mean. Here's a mathematical explanation of how the algorithm works:

### Step 1: Initialization

- Choose the number of clusters, K, that you want to create.

- Initialize K cluster centroids randomly. These centroids are the initial guesses for the center of each cluster.

### Step 2: Assignment

- For each data point in the dataset, calculate the distance between the data point and each of the K centroids. Common distance metrics include Euclidean distance and Manhattan distance.
- Assign the data point to the cluster whose centroid is the closest (i.e., has the smallest distance).

Mathematically, if you have N data points  $\{x_1, x_2, \dots, x_i, \dots, x_n\}$ , and K centroids  $\{\mu_1, \mu_2, \dots, \mu_j, \dots, \mu_k\}$ , you assign each data point  $x_i$  to the cluster  $j$  that minimizes the distance measure:

$$[j = \arg\min_{\{j\}} \text{distance}(x_i, \mu_j)]$$

### Step 3: Update

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

Mathematically, for each cluster  $j$ , you compute the new centroid  $\mu_j$  as:

$$[\mu_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_i]$$

Where:

- $(\mu_j)$  is the new centroid for cluster  $j$ .
- $(n_j)$  is the number of data points in cluster  $j$ .
- $(x_i)$  are the data points assigned to cluster  $j$ .

### Step 4: Convergence

- Repeat steps 2 and 3 until one of the convergence criteria is met. Common convergence criteria include:
  - A maximum number of iterations is reached.
  - The centroids no longer change significantly (i.e., they converge).

### Step 5: Result

- After convergence, you have K clusters, each represented by its centroid.
- Each data point belongs to one of these clusters.

The algorithm aims to minimize the total within-cluster variance, which is the sum of the squared distances between each data point and its assigned centroid within its cluster. This objective function is also known as the "inertia" or "sum of squared distances."

Mathematically, the objective function is:

$$[J = \sum_{j=1}^K \sum_{i=1}^{n_j} \text{distance}(x_i, \mu_j)^2]$$

K-means is a heuristic algorithm, and there is no guarantee that it will find the globally optimal solution. The final clusters obtained can depend on the initial random centroids and the order in which data points are processed. To mitigate this, the algorithm is often run multiple times with



different initializations, and the best result is chosen based on the lowest objective function

Sure, let's go through a step-by-step example of K-means clustering using a small dataset. We'll perform K-means clustering on a 2D dataset with four data points, and we'll cluster them into two clusters ( $K=2$ ).

### Step 1: Initialization

- Choose  $K$  (the number of clusters) to be 2.
- Initialize the centroids randomly. For simplicity, let's assume our initial centroids are (2, 2) and (5, 5).

### Step 2: Assignment

- Calculate the distance between each data point and both centroids.

Data Point	Distance to Centroid 1	Distance to Centroid 2	Assigned Cluster
A (1, 1)	1.41	5.66	1
B (2, 2)	0.00	4.24	1
C (5, 5)	4.24	0.00	2
D (6, 6)	5.66	1.41	2

Based on the distances, data points A and B are assigned to cluster 1 (closer to centroid 1), while data points C and D are assigned to cluster 2 (closer to centroid 2).

### Step 3: Update

- Recalculate the centroids for each cluster by taking the mean of data points in each cluster.

For Cluster 1:

- New centroid = Mean of (1, 1) and (2, 2) =  $((1+2)/2, (1+2)/2) = (1.5, 1.5)$

For Cluster 2:

- New centroid = Mean of (5, 5) and (6, 6) =  $((5+6)/2, (5+6)/2) = (5.5, 5.5)$

### Step 4: Convergence

- Check if the centroids have changed significantly. If they have, go back to step 2; otherwise, we have convergence.

In this case, the centroids have changed, so we continue to the next iteration.

### Step 2 (Iteration 2): Assignment

- Calculate the distance between each data point and both centroids based on the new centroids.

Data Point	Distance to New Centroid 1	Distance to New Centroid 2	Assigned Cluster
A (1, 1)	0.71	5.66	1

Data Point	Distance to New Centroid 1	Distance to New Centroid 2	Assigned Cluster
B (2, 2)	0.71	4.95	1
C (5, 5)	4.95	0.71	2
D (6, 6)	5.66	1.41	2

### Step 3 (Iteration 2): Update

- Recalculate the centroids for each cluster.

For Cluster 1:

- New centroid = Mean of (1, 1) and (2, 2) =  $((1+2)/2, (1+2)/2) = (1.5, 1.5)$

For Cluster 2:

- New centroid = Mean of (5, 5) and (6, 6) =  $((5+6)/2, (5+6)/2) = (5.5, 5.5)$

### Step 4 (Iteration 2): Convergence

- Check if the centroids have changed significantly. If they have, go back to step 2; otherwise, we have convergence.

In this case, the centroids have not changed between iterations. Therefore, we have convergence.

### Step 5: Result

- The final clusters are as follows:
  - Cluster 1: Data points A and B.
  - Cluster 2: Data points C and D.

The centroids of the two clusters are (1.5, 1.5) for Cluster 1 and (5.5, 5.5) for Cluster 2.

This is a simplified example, but it illustrates the basic steps of the K-means clustering algorithm. In practice, K-means is used for larger datasets with more dimensions, and the

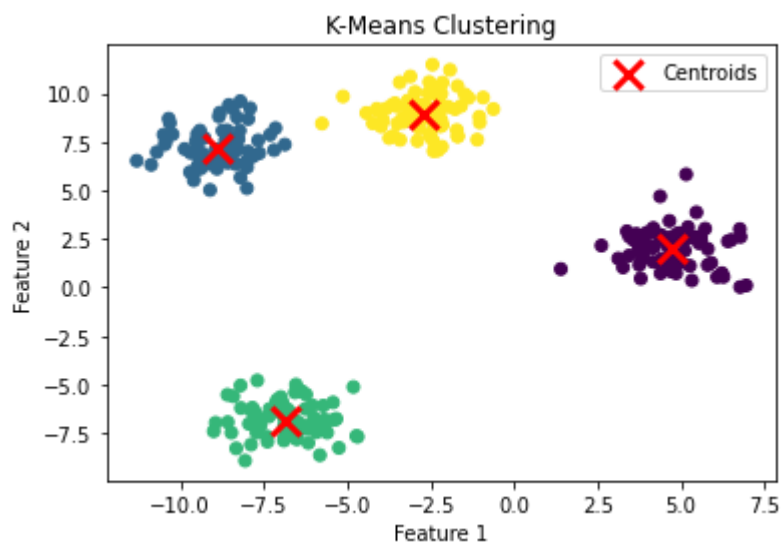
```
In [22]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs

# Generate synthetic data for demonstration
X, _ = make_blobs(n_samples=300, centers=4, random_state=42)

# Initialize and fit KMeans
kmeans = KMeans(n_clusters=4, random_state=42)
kmeans.fit(X)

# Get cluster labels and centroids
labels = kmeans.labels_
centroids = kmeans.cluster_centers_

# Plot the data points and cluster centers
plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis')
plt.scatter(centroids[:, 0], centroids[:, 1], marker='x', s=200, linewidths=3,
plt.legend()
plt.title('K-Means Clustering')
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.show()
```



```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import silhouette_score
```

## Load the Iris dataset

```
iris = load_iris()
X = iris.data
```

## Standardize the features (mean=0, std=1)

```
scaler = StandardScaler() X_scaled = scaler.fit_transform(X)
```

## Determine the optimal number of clusters using the Elbow method

```
from sklearn import config_context
```

```
with config_context(assume_finite=True): kmeans = KMeans(n_clusters=n_clusters,  
random_state=42) kmeans.fit(X_scaled)
```

## Plot the Elbow curve to find the optimal number of clusters

```
plt.figure(figsize=(8, 6)) plt.plot(range(1, 11), inertia, marker='o', linestyle='-', color='b')  
plt.xlabel('Number of Clusters') plt.ylabel('Inertia (Within-Cluster Sum of Squares)')  
plt.title('Elbow Method for Optimal Cluster Number') plt.grid(True) plt.show()
```

## Based on the Elbow method, let's choose K=3

```
n_clusters = 3
```

## Perform K-means clustering with the chosen number of clusters

```
kmeans = KMeans(n_clusters=n_clusters, random_state=42) kmeans.fit(X_scaled) labels =  
kmeans.labels_
```

## Evaluate the quality of the clusters using Silhouette Score

```
silhouette_avg = silhouette_score(X_scaled, labels) print(f'Silhouette Score:  
{silhouette_avg:.2f}')
```

## Visualize the clustering results

```
plt.scatter(X_scaled[:, 0], X_scaled[:, 1], c=labels, cmap='viridis')  
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], marker='x', s=200,  
linewidths=3, color='red', label='Centroids') plt.legend() plt.title('K-Means Clustering (Scaled  
Features)') plt.xlabel('Feature 1 (Standardized)') plt.ylabel('Feature 2 (Standardized)') plt.show()
```

```
In [23]: import pandas as pd
df=pd.read_csv(r"C:\Users\DELL\Downloads\IRIS.csv")
df
```

```
Out[23]:
```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
...	...	...	...	...	...
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

150 rows × 5 columns

```
In [ ]:
```