

# Introduction to Machine Learning

Types of ML, Scikit-Learn

## AGENDA

1. What is Machine Learning?
2. Types of Machine Learning
3. Machine Learning Process
4. Hands-on with Scikit-Learn

**What's the most interesting place you've traveled to, or the top spot on your travel bucket list?**

## 1. What is Machine Learning?

Machine Learning (ML) is a field within artificial intelligence (AI) that enables systems to learn from data instead of following explicit programming instructions.

ML allows systems to find patterns and make decisions, making it essential in various fields, such as healthcare, finance, and technology.

## Importance of Machine Learning

Machine Learning (ML) is increasingly becoming a vital component in various aspects of our daily lives, influencing how we interact with technology, make decisions, and even experience entertainment. Here's a closer look at some common applications of ML that cater to both newcomers and those with more advanced knowledge:

# 1. Recommendation Systems

Recommendation systems are designed to predict what products or content users might like based on their previous interactions. Companies like **YouTube** and **Netflix** use ML algorithms to analyze your viewing history, preferences, and behavior patterns to suggest new shows, movies, or videos.

## How It Works

- **Collaborative Filtering:** This approach relies on the preferences of other users who share similar interests. For instance, if User A and User B have similar viewing habits, recommendations for User B may be based on what User A enjoyed.
- **Content-Based Filtering:** This method suggests items similar to those the user has liked in the past. If you frequently watch action movies, the system will recommend more titles from that genre.

# 2. Image Recognition

Image recognition technology, such as **object detection** and facial recognition, is utilized across various fields, from social media platforms to security systems. This technology enables machines to interpret and categorize images just like humans do.

## How It Works

- **Convolutional Neural Networks (CNNs):** These are specialized deep learning models designed to process pixel data in images. They automatically learn hierarchical features, such as edges, shapes, and textures, enabling them to identify objects within images accurately.
- **Transfer Learning:** By using pre-trained models, developers can significantly reduce the time and data required to build robust image recognition systems for specific tasks.

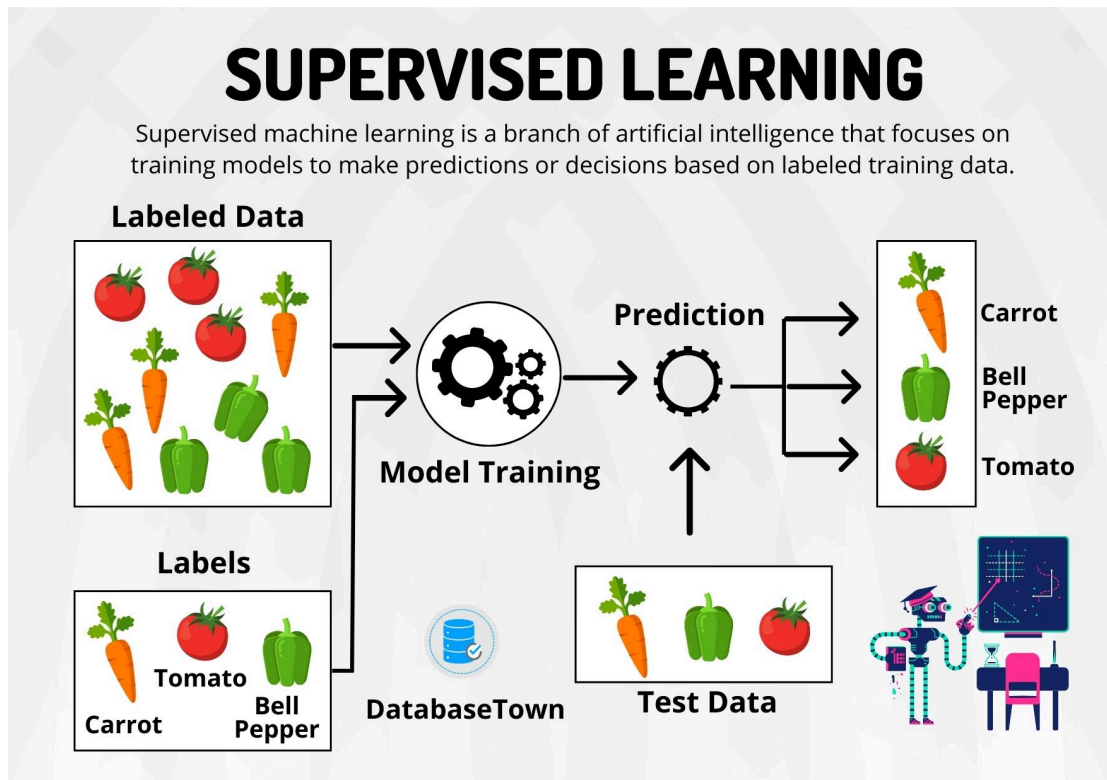
# 3. Predictive Analysis

Predictive analysis involves using historical data to make forecasts about future events, which is valuable in numerous fields like finance, healthcare, and marketing. For example, businesses use ML to predict **stock prices** or forecast **weather patterns**.

## How It Works

- **Time Series Analysis:** This involves statistical techniques to analyze time-ordered data points, enabling businesses to forecast trends and seasonal variations.
- **Regression Algorithms:** Techniques like linear regression, decision trees, and more complex ensemble methods (e.g., Random Forest, Gradient Boosting) can predict numeric outcomes based on input features.

## 2. Types of Machine Learning



Supervised learning is one of the most widely used techniques in machine learning, where the **model learns from labeled data**.

This method is pivotal for tasks that require predictions based on input features, and it plays a crucial role in various real-world applications.

### What is Supervised Learning?

In supervised learning, the model is trained on a dataset that contains both the input features and their corresponding known outputs. This means each example in the training data is paired with a label, which serves as the ground truth for the model to learn from.

#### Example: Predicting Housing Prices

Consider a scenario where we want to predict housing prices based on various features such as:

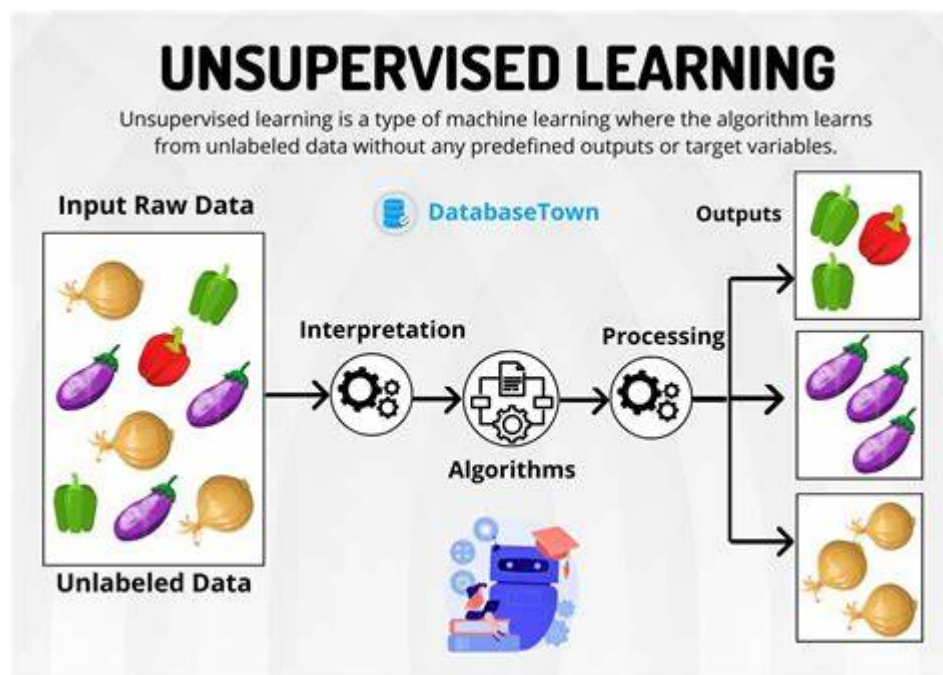
- Square footage
- Number of bedrooms
- Location
- Age of the property

In this case, the input features are the attributes of the houses, while the known output is the price at which each house was sold. The model uses this labeled data to learn the relationship between the input features and the housing prices.

## How Does the Model Learn?

1. **Training Phase:** During this phase, the model analyzes the training data. It identifies patterns and relationships between the input features and the known outputs. For instance, it might learn that larger houses generally have higher prices, or that houses in certain neighborhoods sell for more.
2. **Evaluation Phase:** After training, the model is evaluated using a separate dataset known as the validation or test set. This set contains input features with known outputs that the model has not seen before. The model's performance is measured by comparing its predictions against the actual known outputs. Metrics such as Mean Squared Error (MSE) for regression tasks or accuracy for classification tasks are commonly used.
3. **Prediction Phase:** Once the model has been trained and validated, it can be used to make predictions on new, unseen data. For example, if a new house is listed for sale with specific features, the model can provide an estimated price based on what it has learned from the training data.

## Unsupervised Learning



Unsupervised learning is another fundamental approach in machine learning, focusing on discovering patterns and structures in unlabeled data. Unlike supervised learning, where each data point has a corresponding label, unsupervised learning seeks to understand the underlying relationships within the data without predefined outcomes.

## What is Unsupervised Learning?

In unsupervised learning, the model is presented with data that lacks labels or known outputs. The goal is to identify patterns, groupings, or structures that may not be immediately apparent. This method is particularly useful in exploratory data analysis, helping to uncover insights that can guide further investigation.

### **Example: Clustering Customers by Purchasing Behavior**

Consider a retail scenario where a company wants to analyze customer purchasing behavior. The dataset might include features such as:

- Age
- Gender
- Purchase history
- Spending habits

Since this data is unlabeled (the company doesn't know how many types of customers exist or what their purchasing categories are), the model can apply clustering algorithms to identify groups of similar customers based on their behaviors.

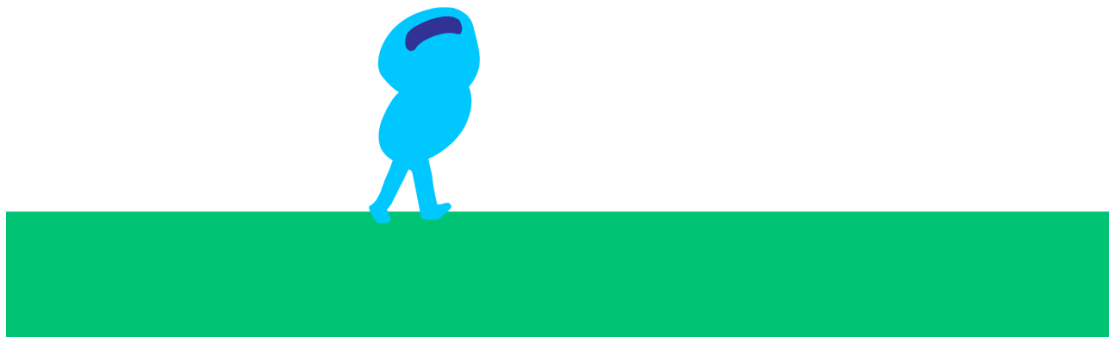
### **How Does the Model Learn?**

1. **Data Analysis:** The model processes the data, examining the features to find similarities and differences among data points. It looks for natural groupings or patterns without any guidance from labels.
2. **Clustering:** A common technique in unsupervised learning is clustering, where the model groups similar data points together. For example, it might segment customers into distinct groups, such as high spenders, occasional buyers, and bargain hunters, based on their purchasing behavior.
3. **Dimensionality Reduction:** Another technique used in unsupervised learning is dimensionality reduction. This method simplifies complex datasets by reducing the number of features while preserving important information. Techniques like Principal Component Analysis (PCA) can help visualize high-dimensional data in two or three dimensions, making it easier to identify patterns.

## **Reinforcement Learning**

Reinforcement learning (RL) is a unique approach within the field of machine learning that focuses on training models to make decisions by interacting with their environment.

Unlike supervised and unsupervised learning, reinforcement learning emphasizes learning through trial and error, using feedback in the form of rewards and penalties to achieve specific objectives.



## What is Reinforcement Learning?

In reinforcement learning, an agent (the model) learns how to navigate an environment by taking actions and receiving feedback based on those actions. The primary goal is to maximize cumulative rewards over time, effectively training the agent to make better decisions based on past experiences.

### Example: Robotics and Gaming

One of the most common applications of reinforcement learning is in robotics, where robots learn to perform tasks through interaction with their surroundings. For instance, a robot might learn to pick up objects or navigate a maze. Similarly, in gaming, reinforcement learning has been used to develop intelligent agents that can play video games, making decisions to achieve high scores or win against opponents.



### How Does the Model Learn?

1. **Agent and Environment:** The reinforcement learning framework consists of an agent and an environment. The agent takes actions within the environment, and the environment responds by providing feedback in the form of rewards (positive feedback) or penalties (negative feedback).
2. **Exploration vs. Exploitation:** A key aspect of reinforcement learning is the balance between exploration and exploitation. The agent must explore new actions to discover potentially better strategies (exploration) while also leveraging known actions that yield good rewards (exploitation). Striking the right balance is crucial for effective learning.
3. **Reward Signal:** The agent receives a reward signal from the environment after each action. Positive rewards reinforce behaviors that lead to desirable outcomes, while penalties discourage actions that result in negative outcomes. Over time, the agent learns to associate specific actions with their corresponding rewards, optimizing its decision-making strategy.
4. **Policy and Value Function:** In reinforcement learning, a policy defines the agent's behavior, dictating which action to take in a given state. The value function estimates the expected reward for being in a particular state, helping the agent evaluate the potential long-term benefits of its actions.

## 3. Machine Learning Process

The machine learning process is a systematic approach that outlines the steps involved in developing, training, and deploying machine learning models.

Understanding this process is crucial for effectively leveraging machine learning to solve real-world problems. Below are the key stages of the machine learning process:

### 1. Problem Definition

The first step is to clearly define the problem you are trying to solve. This involves understanding the business objectives and determining how machine learning can provide a solution. It's essential to articulate the expected outcomes and success criteria for the project.

#### Example:

- **Objective:** Predicting customer churn for a subscription-based service.
- **Success Criteria:** Reducing churn rate by 15% over the next quarter.
- **4Ws:** Who, What, When, Why.

### 2. Data Collection

Once the problem is defined, the next step is to gather relevant data. Data can be collected from various sources, such as databases, APIs, or web scraping. The quality and quantity of data are critical as they directly impact model performance.

### Considerations:

- Ensure that the data is representative of the problem domain.
- Address any privacy or ethical concerns related to data usage.

## 3. Data Preprocessing

After collecting the data, it often requires preprocessing to clean and prepare it for analysis. This step involves several tasks:

- **Handling Missing Values:** Decide how to deal with any missing data points, either by removing them or imputing values.
- **Feature Engineering:** Create new features or modify existing ones to improve model performance.
- **Normalization/Standardization:** Scale features to ensure that they contribute equally to the model training.

## 4. Model Selection

In this stage, you choose the appropriate machine learning algorithms based on the problem type (e.g., classification, regression, clustering) and the data characteristics.

Common algorithms include:

- **Supervised Learning:** Decision trees, logistic regression, support vector machines (SVM).
- **Unsupervised Learning:** K-means clustering, hierarchical clustering, principal component analysis (PCA).

## 5. Model Training

Once you've selected a model, the next step is to train it using the prepared dataset. During training, the model learns patterns from the data by adjusting its parameters to minimize errors in predictions.

### Tips:

- Use techniques such as cross-validation to evaluate model performance during training.
- Monitor for overfitting, where the model performs well on training data but poorly on unseen data.



## 6. Model Evaluation

After training, the model's performance is evaluated using a separate test dataset. Metrics such as accuracy, precision, recall, F1 score, or mean squared error (MSE) are calculated to assess how well the model is performing.

### Example:

- For a classification problem, use confusion matrices to understand true positives, false positives, true negatives, and false negatives.

## 7. Model Deployment

Once the model is validated and meets the success criteria, it is deployed into a production environment where it can be used for making predictions on new data.

This process may involve integrating the model with existing software applications or setting up APIs for access.

## 8. Monitoring and Maintenance (ML Ops)

Machine Learning Operations (ML Ops) is a critical aspect of deploying machine learning models. It refers to the practices and tools used to manage and monitor machine learning models in production.

Key components include:

- **Model Monitoring:** Continuously track the model's performance and data drift over time to ensure it remains effective.
- **Versioning:** Maintain versions of models, datasets, and code to ensure reproducibility and manage updates.
- **Automated Retraining:** Set up systems to automatically retrain models as new data becomes available or as performance declines.
- **Collaboration:** Foster collaboration between data scientists, engineers, and operations teams to streamline the ML lifecycle.

## 4. Hands-on with Sci-kit Learn

### Install Necessary Libraries

- Open your Jupyter Notebook and install Scikit-Learn by running:

```
pip install scikit-learn
```

# Basic Python Refresher for Data Handling

Load Data: Use Python libraries like Pandas to load and view data.

```
import pandas as pd
```

## Exploring Data

- Load a Simple Dataset
- Load the Iris dataset from Scikit-Learn.
- This dataset is popular for beginners, as it contains labeled data about flower species and related features (sepal length, petal length, etc.).

```
# Load a sample dataset from Scikit-Learn
from sklearn.datasets import load_iris
data = load_iris()
df = pd.DataFrame(data.data, columns=data.feature_names)
df['target'] = data.target
df.head()
```

```
# Explore Data: Check the first few rows, summary statistics, and column types.
```

## Visualize Basic Features

- Use basic visualizations to understand the data.
- Plotting example:

```
import matplotlib.pyplot as plt

# Scatter plot of sepal length and width
plt.scatter(df['sepal length (cm)'], df['sepal width (cm)'],
            c=data.target)
plt.xlabel('Sepal Length (cm)')
plt.ylabel('Sepal Width (cm)')
plt.title('Iris Dataset - Sepal Length vs Sepal Width')
plt.show()
```

## Program Breakdown

### Pip installs

```
pip install scikit-learn
```

- Pip is the package installer for Python, allowing you to add libraries that aren't part of the default Python installation.
- Variations of use:

- `pip install <package>` (Standard Command Line) to directly download a package from terminal/command prompt. Most common method but might not work directly in Jupyter notebooks or interactive environments
- `!pip install <package>` (Shell Magic in Jupyter) to run shell commands. `!` is a "magic command" in Jupyter that allows running commands as if in a shell. This approach is convenient but can sometimes lead to issues with environment mismatches (e.g., installing in one Python environment while Jupyter is using another).
- `%pip install <package>` (IPython Magic) to run in interactive environment. `%pip` is environment-aware within Jupyter and ensures the package is installed in the exact Python environment that Jupyter is using.
- `python -m pip install <package>` (Python Module Invocation) to run `pip` as a module ( `-m` ). When there are multiple Python versions installed, as it ensures `pip` runs in the same environment as the specified Python interpreter.

## Import Commands

```
from sklearn.datasets import load_iris
```

- using `from`, Python knows you are accessing only part of a package, rather than importing the entire package
- `sklearn.datasets` targets the submodule `datasets` from library `sklearn`. It contains small, ready-to-use datasets (like Iris, Digits, and Boston datasets) that are helpful for testing algorithms or learning purposes.
  - `load_digits`: A dataset of handwritten digits, ideal for practicing image classification tasks.
  - `load_boston`: Contains housing prices in Boston; often used for regression tasks (note: has been deprecated in newer versions).
  - `load_wine`: Data on wine samples, with features describing chemical properties, often used in classification.
  - `load_breast_cancer`: A dataset with features from breast cancer cases, used for binary classification.
- Other important modules in Scikit-Learn include `sklearn.model_selection`, `sklearn.preprocessing`, `sklearn.metrics`, `sklearn.ensemble`
- `import` keyword makes the specified code accessible without needing to type out the full module path each time

## Pandas Dataframe

```
df = pd.DataFrame(data.data, columns=data.feature_names)
```

- Actual syntax:

```
pd.DataFrame(data=None, index=None, columns=None, dtype=None, copy=False)
```

1. Data - Core data for dataframe can be types like dictionaries, lists, arrays, series or other dataframes

2. Index - Labels for rows, that can be a list or array, example:

```
index=['Row1', 'Row2']
```

3. Columns - Labels for columns, that can be a list or array, example:

```
columns=['Column1', 'Column2']
```

4. Datatype - `dtype` to set the data type of each column, example: `float` , `int` , `str` .

5. Copy - Could be `True` or `False` . When set to `True` creates a copy of the data. `False` changes the original data.

## Example:

Creating a DataFrame from a Dictionary:

```
import pandas as pd
data = {'Name': ['Alice', 'Bob', 'Catlyn'], 'Age': [24, 30, 45], 'Calls': [24, 30, 45]}
df = pd.DataFrame(data, index=['Person1', 'Person2', 'Person3'], columns=['Name', 'Age', 'Calls'], dtype='str', copy=True)
```

```
In [11]: import pandas as pd
data = {'Name': ['Alice', 'Bob', 'Catlyn'], 'Age': [24, 30, 45], 'Calls': [24, 30, 45]}
df = pd.DataFrame(data, index=['Person1', 'Person2', 'Person3'], columns=['Name', 'Age', 'Calls'], dtype='str', copy=True)
```

```
In [ ]: df
```

## Visualisation Commands:

```
plt.scatter(df['sepal length (cm)'], df['sepal width (cm)'], c=data.target)
```

- Syntax:

```
plt.scatter(x, y, **kwargs)
```

- `x` : The data for the x-axis (e.g., a list or array of values).
- `y` : The data for the y-axis, with the same length as x.
- `**kwargs` : Keyword arguments, can take various values:
  - `c`: Color of each point
  - `s`: size of each point
  - `alpha`: Transparency from 0 to 1
  - `label`: label for the plot

```
In [ ]: # Example usage:
import matplotlib.pyplot as plt

plt.scatter(df['Age'], df['Calls'], c='blue', s=50, alpha=0.7)
```

```
In [ ]: plt.xlabel('Age this is')
```

```
In [ ]: plt.ylabel('This is Calls')
```

```
In [ ]: plt.title("It's my life")
```

```
In [17]: # In non-interactive environments (like a Python script), plt.show() is necessary
plt.show()
```

```
In [ ]: # Example usage: All at once
import matplotlib.pyplot as plt

plt.scatter(df['Age'], df['Calls'], c='blue', s=50, alpha=0.7)
plt.xlabel('Age this is')
plt.ylabel('This is Calls')
plt.title("It's my life")
plt.show()
```

### Display Scikit-Learn Plots in a New Window

Use the %matplotlib magic command to select a GUI backend.

```
%matplotlib tk # or %matplotlib qt
```

```
In [19]: # %matplotlib tk
```

## Live Exercise (Advanced/Optional)

Now it's your turn!

## Task 1: Simple Linear Regression with Scikit-Learn (Optional Activity)

### Objective

In this exercise, students will learn the basics of linear regression, a fundamental machine learning technique, by building a simple linear regression model using Scikit-Learn. They will also visualize the results to understand how the model fits the data.

### Dataset

We will use a synthetic dataset for this exercise, where we predict the price of a house based on its size (in square feet).

### Steps

## 1. Set Up the Environment

Make sure you have the following libraries installed:

- `numpy`
- `pandas`
- `matplotlib`
- `scikit-learn`

You can install them using `pip install numpy pandas matplotlib scikit-learn`.

## 2. Import Libraries

Start by importing the necessary libraries in your Jupyter Notebook.

## 3. Create a Synthetic Dataset

Generate a simple dataset: Create a synthetic dataset with 100 samples, where house sizes are generated randomly within a range of 0 to 2.5K square feet. The house prices are calculated using a formula that includes some noise to simulate real-world conditions.

## 4. Split the Dataset

Divide the dataset into training and testing sets. Use 80% of the data for training and 20% for testing to evaluate the model's performance.

## 5. Build and Train the Model

Create and train the linear regression model using the training data.

## 6. Make Predictions

Use the model to make predictions on the test set.

## 7. Visualize the Results

Plot the training data, the test data, and the regression line to visually assess the model's performance. The training data points will be displayed in blue, the testing data points in green, and the regression line in red.

## Discussion Points

- What do you observe about the regression line?
- How well does the model predict the house prices?
- What could be improved in this model?

END

# THANK YOU!

## Live Exercise Solutions

```
In [22]: # Import necessary libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

```
In [23]: # Step 1: Create a synthetic dataset
# Generate random house sizes (in square feet) between 500 and 2500
np.random.seed(0) # For reproducibility
house_sizes = np.random.randint(500, 2500, size=100)

# Generate house prices based on sizes with some noise
# Price = 50 * size + some random noise
house_prices = 50 * house_sizes + np.random.normal(0, 10000, size=100)
```

```
In [24]: # Step 2: Create a DataFrame
data = pd.DataFrame({
    'Size': house_sizes,
    'Price': house_prices
})
```

```
In [25]: # Step 3: Split the dataset into training and testing sets
X = data[['Size']] # Features (independent variable)
y = data['Price'] # Target (dependent variable)

# 80% of the data for training and 20% for testing
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_
```

```
In [26]: # Step 4: Build and train the linear regression model
model = LinearRegression() # Create an instance of LinearRegression
model.fit(X_train, y_train) # Train the model
```

```
Out[26]: ▼ LinearRegression ⓘ ?
LinearRegression()
```

```
In [27]: # Step 5: Make predictions on the test set
predictions = model.predict(X_test)
```

```
In [28]: # Step 6: Visualize the results
plt.figure(figsize=(10, 6))
```

```

# Plot training data
plt.scatter(X_train, y_train, color='blue', label='Training Data')

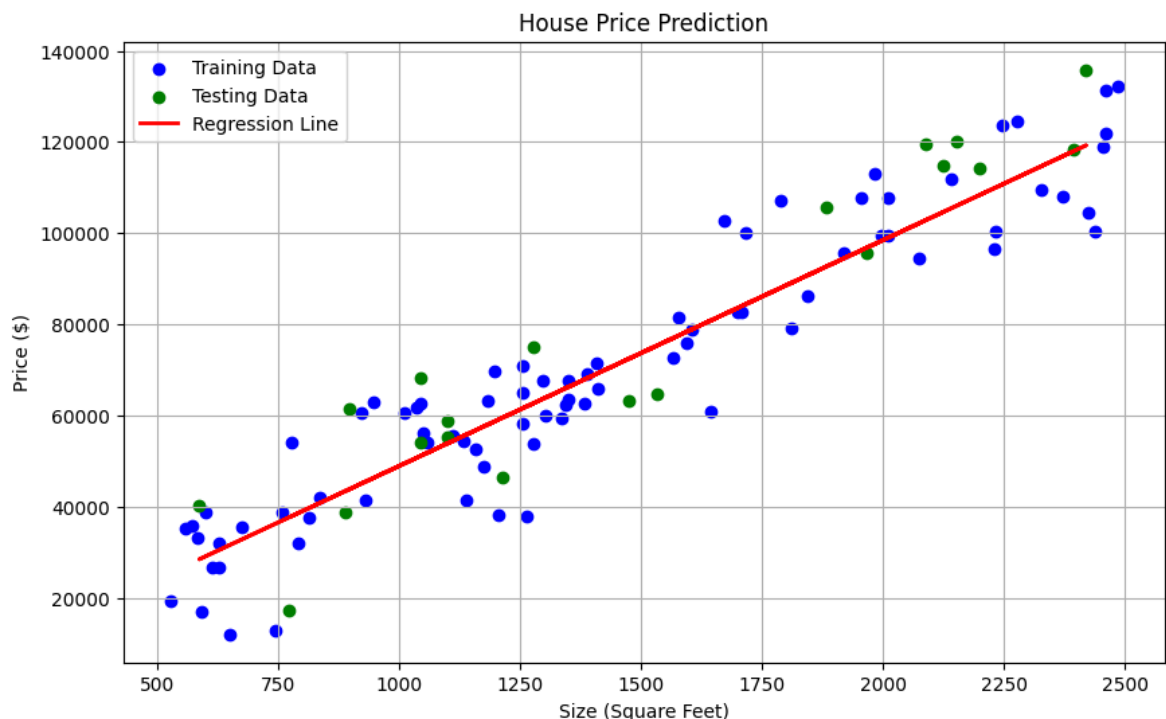
# Plot testing data
plt.scatter(X_test, y_test, color='green', label='Testing Data')

# Plot regression line
plt.plot(X_test, predictions, color='red', linewidth=2, label='Regression Line')

# Add Labels and title
plt.title('House Price Prediction')
plt.xlabel('Size (Square Feet)')
plt.ylabel('Price ($)')
plt.legend()
plt.grid()

# Show the plot
plt.show()

```



### Explanation of some of the uncovered parts

`np.random.seed(0)` sets the seed for NumPy's random number generator. Setting a seed ensures that the random numbers generated are the same each time the code is run.

Specifying a starting point for the sequence of random numbers that will be generated.

`np.random.randint(500, 2500, size=100)` : This command generates 100 random integers between 500 and 2500.

- 500 is the minimum size of the house in square feet.
- 2500 is the maximum size.
- size=100 indicates that we want to create 100 different house sizes.



`np.random.rand(5)` : Generates 5 random numbers

```
In [37]: np.random.seed(1000)
```

```
In [41]: np.random.randint(500, 2500, size=100)
```

```
Out[41]: array([1514, 1736, 1667, 1643, 2393, 1668, 2011,  861, 1385, 2162, 2362,
                2191, 1614, 1447, 2352, 2205,  625, 1470, 2492,  930, 2011, 1331,
                2322, 2211, 1412, 1398, 1943, 1959,  945, 2481,  539,  661, 1466,
                976, 1479,  721,  881, 1592, 1216, 1254, 1676, 1250, 2039,  559,
                1904,  948, 2372, 1765, 1698, 1175, 2376,  649, 1859, 2027, 1949,
                1592, 1718,  869, 1608, 2310, 1223, 2380, 1831, 1123,  631, 1792,
                635, 1978, 2399, 1857, 1378, 2029, 2477, 1203, 1612, 1694, 1675,
                1106,  746, 1636, 1682, 1624, 1499, 2291,  539, 2345, 1836, 1336,
                1355, 1661,  649, 1799,  577, 1166, 1653, 1235, 2175,  752, 2226,
                2003], dtype=int32)
```

### On train test split

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.2, random_state=0)
```

utilizes the `train_test_split` function from the `sklearn.model_selection` module in Python

`train_test_split()` is a function that splits arrays or matrices into random train and test subsets.

This function takes 4 inputs and return 4 outputs

Input parameters:

- `x` : represents features
- `y` : represents target variables
- `test_size` indicates the amount of data for testing
- `random_state` sets the seed to 0 so the split is same, everytime we run the code

The output variables are:

- `X_train` : will contain the training features
- `X_test` : will contain the testing features
- `y_train` : will contain the training target values
- `y_test` : will contain the testing target values

### Simulate a similar function:

```
In [42]: def idreessplit(a,b,c,d):
          w,x,y,z = 20,40, 50, 60

          return w,x,y,z
```

```
In [43]: a,b,c,d= 'idrees1','idrees2','idrees3','idrees4'
```

```
In [44]: person1,person2,person3,person4=idreessplit(a,b,c,d)
print(f'{person1} {person2} {person3} {person4}')
```

20 40 50 60

```
In [50]: # another example

def idreessplit(a,b,c,d):

    a,b,c,d= 'idrees train 20% -','idrees train 80% -','idrees test 20% -','idrees test 80% -'

    return a,b,c,d

p1,p2,p3,p4 = idreessplit(a,b,c,d)

print(f'{p1} {p2} {p3} {p4}')
```

idrees train 20% - idrees train 80% - idrees test 20% - idrees test 80%-

In [ ]:

## Programming Interview Questions

1. topic:

- question

**Mohammad Idrees Bhat**

Tech Skills Trainer | AI/ML Consultant