

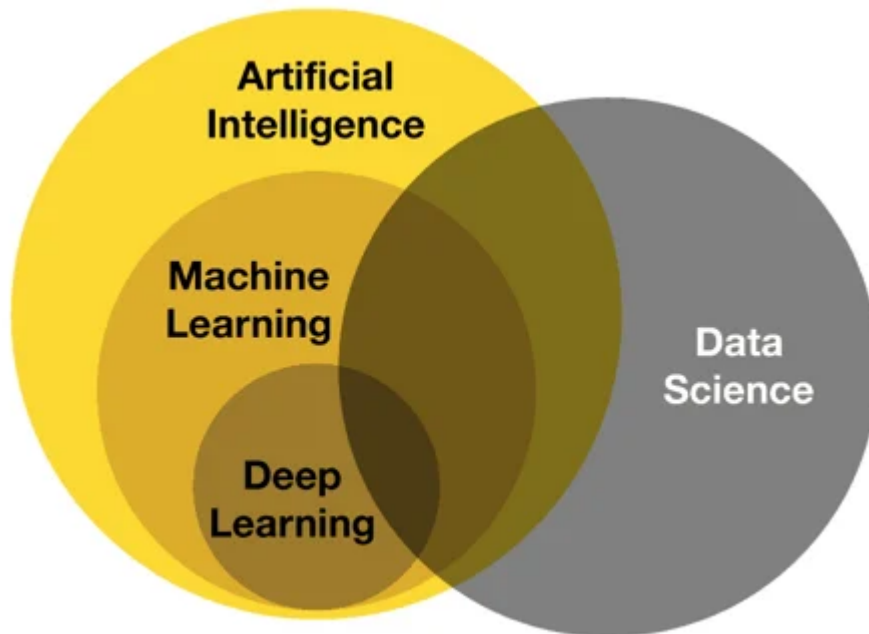
Supervised Learning

AGENDA

1. Recap - Overview Machine Learning
2. Overview of Supervised Learning
3. Various algorithms used in supervised learning
4. Linear Regression in Depth
5. Hands-on Activity - Linear Regression

What's an interesting item in your bag or on your desk right now?

1. Recap - Overview Machine Learning



- **AI (Artificial Intelligence):** The simulation of human intelligence processes by machines, particularly computer systems, enabling them to perform tasks that typically require human intelligence, such as reasoning, learning, and problem-solving.
- **ML (Machine Learning):** A subset of AI that focuses on the development of algorithms and statistical models that allow computers to **improve their performance** on a specific task through experience without being explicitly programmed.
- **DL (Deep Learning):** A specialized area of machine learning that utilizes neural networks with many layers (deep networks) to analyze and learn from large amounts of data, particularly effective in tasks like image and speech recognition.
- **DS (Data Science):** An interdisciplinary field that combines statistical analysis, data mining, and machine learning to extract insights and knowledge from structured and unstructured data.

Learning vs Not Learning

Machine Learning (ML) focuses on algorithms that allow computers to learn from and make predictions.

ML algorithms improve their performance as they are exposed to more data over time, to identify patterns and make informed predictions without explicit programming.

Examples of ML:

- **Spam Detection:** An email filtering system that uses ML to classify emails as spam or not spam based on patterns in the text, sender information, and user behavior.

- **Recommendation Systems:** Services like Netflix or Amazon that suggest movies or products based on user preferences and past behavior using collaborative filtering techniques.
- **Image Recognition:** Applications that can identify objects, people, or scenes in images, such as facial recognition technology used in security systems.
- **Predictive Analytics:** A method that uses statistical algorithms and machine learning techniques to identify the likelihood of future outcomes based on historical data. For instance, businesses use predictive analytics to forecast sales trends by analyzing past sales data, helping them make informed decisions on inventory and marketing strategies.

Traditional Programming (What is Not Machine Learning?)

While many coding tasks involve logic and data processing, not all programming tasks fall under the category of machine learning. Traditional coding involves explicitly defining rules and procedures to perform a specific task without the capability to learn from data.

Examples of Non-ML Coding:

- **Basic Programming Tasks:** Writing a simple program to calculate the sum of two numbers or to sort a list of items. These tasks follow fixed instructions without any learning capability.

Example 1: A Python script that takes user input and prints the sum:

```
num1 = int(input("Enter first number: "))
num2 = int(input("Enter second number: "))
print("Sum:", num1 + num2)
```

Example 2: A script takes a student's score as input and classifies the grade.

```
score = int(input("Enter your score: ")) # Take input from the user

# Conditional logic to classify the grade based on the score
if score >= 90:
    grade = "A"
elif score >= 80:
    grade = "B"
elif score >= 70:
    grade = "C"
elif score >= 60:
    grade = "D"
else:
    grade = "F"

# Print the classified grade
print("Your grade is:", grade)
```

Algorithm vs Model

- **Algorithm:** A set of rules or instructions that a computer follows to solve a problem or complete a task. In machine learning, algorithms process data and learn from it to make predictions or decisions.
 - **Example:** The **linear regression algorithm** is used to find the best-fit line for a set of data points, which helps in predicting outcomes based on input features.
- **Model:** The output created by an algorithm after it has been trained on a specific dataset. A model represents the learned patterns or relationships in the data and can make predictions on new, unseen data.
 - **Example:** After training the linear regression algorithm on a dataset of house prices and their features (like size and location), the resulting **linear regression model** can predict the price of a new house based on its features.

2. Overview of Supervised Learning

Supervised Learning is a type of machine learning that uses labeled data to train algorithms that classify data or predict outcomes accurately.

Understanding the key components

Labeled Data: Each data point is associated with a **known output** (label), which the model uses to learn patterns.

Target Variable: The variable that the model is trying to predict.

Goal: The objective in supervised learning is to **predict outcomes** based on input features.

Example 1

Can we predict the price of a car prices if we know the car age , mileage , brand , and horsepower ?

Labelled Data:

input features (car age , brand , mileage , horsepower) and a labeled outcome (the car's price)

- Import the dataset from here (without attribute names):
<https://archive.ics.uci.edu/dataset/10/automobile>
- Download dataset from here (with attribute names):
<https://www.kaggle.com/datasets/toramky/automobile-dataset>

- Full description of these attributes on the UCI Automobile Dataset webpage, under the "Attribute Information" section or the Kaggle page.

```
In [1]: import pandas as pd

url = "https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-8

columns = ["symboling", "normalized_losses", "make", "fuel_type", "aspiration",
           "body_style", "drive_wheels", "engine_location", "wheel_base", "length",
           "height", "curb_weight", "engine_type", "num_cylinders", "engine_size",
           "bore", "stroke", "compression_ratio", "horsepower", "peak_rpm", "city_mpg",
           "highway_mpg", "overall_mpg"]

data = pd.read_csv(url, names=columns)

data.info
```

```

Out[1]: <bound method DataFrame.info of
      uel_type aspiration num_doors \
0          3          ?   alfa-romero    gas    std    two
1          3          ?   alfa-romero    gas    std    two
2          1          ?   alfa-romero    gas    std    two
3          2        164      audi      gas    std    four
4          2        164      audi      gas    std    four
..      ...      ...      ...      ...      ...      ...
200       -1         95     volvo      gas    std    four
201       -1         95     volvo      gas    turbo   four
202       -1         95     volvo      gas    std    four
203       -1         95     volvo  diesel    turbo   four
204       -1         95     volvo      gas    turbo   four

      body_style drive_wheels engine_location wheel_base ... engine_size \
0  convertible      rwd      front      88.6 ...      130
1  convertible      rwd      front      88.6 ...      130
2   hatchback      rwd      front      94.5 ...      152
3      sedan      fwd      front      99.8 ...      109
4      sedan      4wd      front      99.4 ...      136
..      ...      ...      ...      ... ...      ...
200      sedan      rwd      front     109.1 ...      141
201      sedan      rwd      front     109.1 ...      141
202      sedan      rwd      front     109.1 ...      173
203      sedan      rwd      front     109.1 ...      145
204      sedan      rwd      front     109.1 ...      141

      fuel_system bore stroke compression_ratio horsepower peak_rpm \
0      mpfi  3.47   2.68           9.0         111      5000
1      mpfi  3.47   2.68           9.0         111      5000
2      mpfi  2.68   3.47           9.0         154      5000
3      mpfi  3.19   3.40          10.0         102      5500
4      mpfi  3.19   3.40           8.0         115      5500
..      ...   ...   ...           ...         ...      ...
200      mpfi  3.78   3.15           9.5         114      5400
201      mpfi  3.78   3.15           8.7         160      5300
202      mpfi  3.58   2.87           8.8         134      5500
203      idi  3.01   3.40          23.0         106      4800
204      mpfi  3.78   3.15           9.5         114      5400

      city_mpg highway_mpg price
0          21          27  13495
1          21          27  16500
2          19          26  16500
3          24          30  13950
4          18          22  17450
..      ...      ...      ...
200         23          28  16845
201         19          25  19045
202         18          23  21485
203         26          27  22470
204         19          25  22625

[205 rows x 26 columns]>

```

Regression vs Classification

Regression is a type of supervised learning used to predict a continuous output variable based on one or more input features.

The output is a continuous value, meaning it can take any numerical value within a range.

Examples:

- Predicting house prices based on features like size, location, and number of bedrooms.
- Estimating a person's weight based on height and age.

Common Algorithms: Linear Regression, Polynomial Regression, Decision Trees (for regression), Support Vector Regression (SVR)

Classification is a type of supervised learning used to categorize input data into predefined classes or categories.

The output is a discrete label, meaning it represents distinct classes or categories.

Examples:

- Classifying emails as "spam" or "not spam."
- Identifying whether an image contains a cat or a dog.

Common Algorithms: Logistic Regression, Decision Trees (for classification), Random Forest, Support Vector Machines (SVM).

Key Differences:

Nature of Output:

- Regression predicts *continuous* values (e.g., price, temperature).
- Classification predicts *categorical* labels (e.g., yes/no, spam/not spam).

Evaluation Metrics:

- Regression typically uses metrics like Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared.
- Classification uses metrics like accuracy, precision, recall, F1-score, and confusion matrix.

3. Various Algorithms used in Supervised Learning

Various Regression Algorithms

- **Linear Regression:**
 - Used to predict a continuous target variable based on the linear relationship with one or more predictor variables.

- **Polynomial Regression:**
 - Extends linear regression by fitting a polynomial curve to capture non-linear relationships.
- **Support Vector Regression (SVR):**
 - An extension of SVM for regression problems, it tries to fit as many data points as possible within a specified margin.
- **Decision Trees (Regression Trees):**
 - A tree-based model that splits data into subsets based on feature values, used for predicting continuous outcomes.
- **Random Forest Regression:**
 - An ensemble method that builds multiple decision trees and averages their predictions to improve accuracy and robustness.

Other used algorithms include:

- **Ridge Regression:**
- A linear regression technique that includes L2 regularization to prevent overfitting by penalizing large coefficients.
- **Lasso Regression:**
 - Similar to ridge regression but uses L1 regularization, which can reduce some coefficients to zero, effectively performing variable selection.
- **Elastic Net:**
 - Combines L1 and L2 regularization, balancing between ridge and lasso regression.
- **Gradient Boosting Machines (GBM):**
 - Builds models in a stage-wise fashion by optimizing loss functions and is particularly effective for complex datasets.

Classification Algorithms

- **Logistic Regression:**
 - Despite its name, it's a classification algorithm used to predict binary outcomes based on one or more predictor variables.
- **Support Vector Machines (SVM):**
 - Finds the hyperplane that best separates data points of different classes, effective in high-dimensional spaces.
- **Decision Trees:**

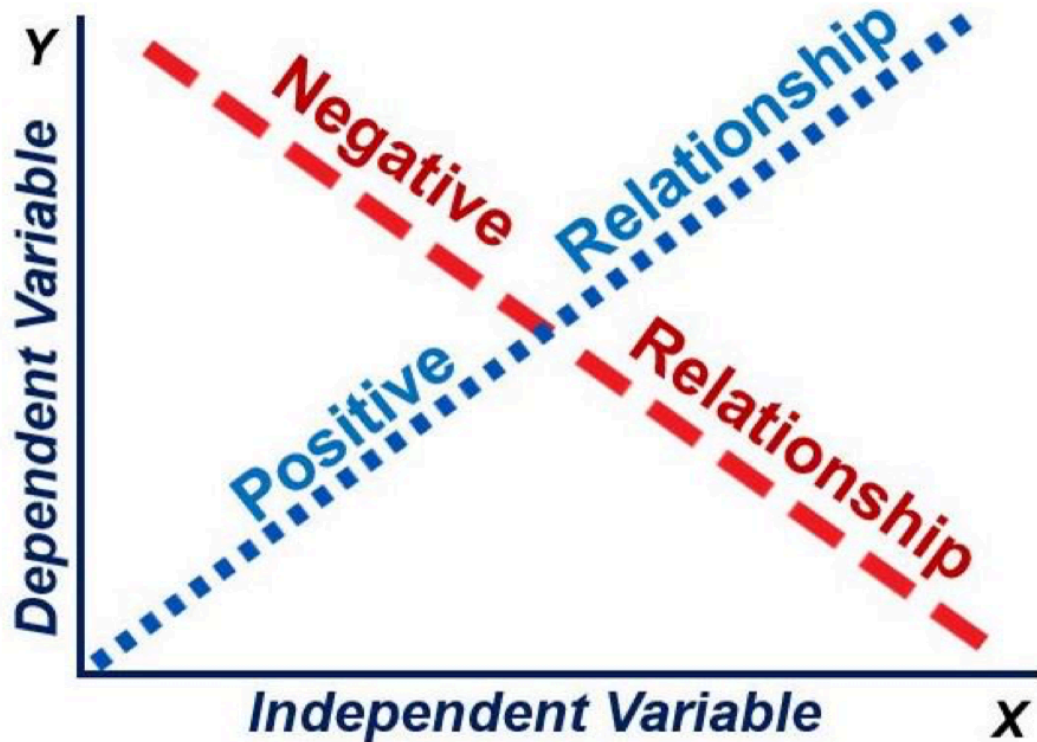
- A tree-like model that makes decisions based on feature values to classify data into categories.
- **Random Forest:**
 - An ensemble method that combines multiple decision trees to improve classification accuracy and reduce overfitting.
- **k-Nearest Neighbors (k-NN):**
 - Classifies a data point based on the majority class of its k-nearest neighbors in the feature space.
- **Naive Bayes:**
 - A probabilistic classifier that applies Bayes' theorem, assuming independence between features, often used for text classification.
- **Gradient Boosting Classifiers:**
 - Builds models iteratively, optimizing the prediction accuracy by focusing on errors from previous iterations.
- **Artificial Neural Networks (ANN):**
 - A flexible model inspired by biological neural networks, used for both classification and regression tasks.

Types of Regression

Regression analysis is a **statistical modeling technique**, a way of mathematically sorting out a series of variables. We use it to determine which variables have an impact and how they relate to one another.

It tries to determine how strongly related one **dependent variable** is to a series of other changing variables. We usually refer to them as **independent variables**.

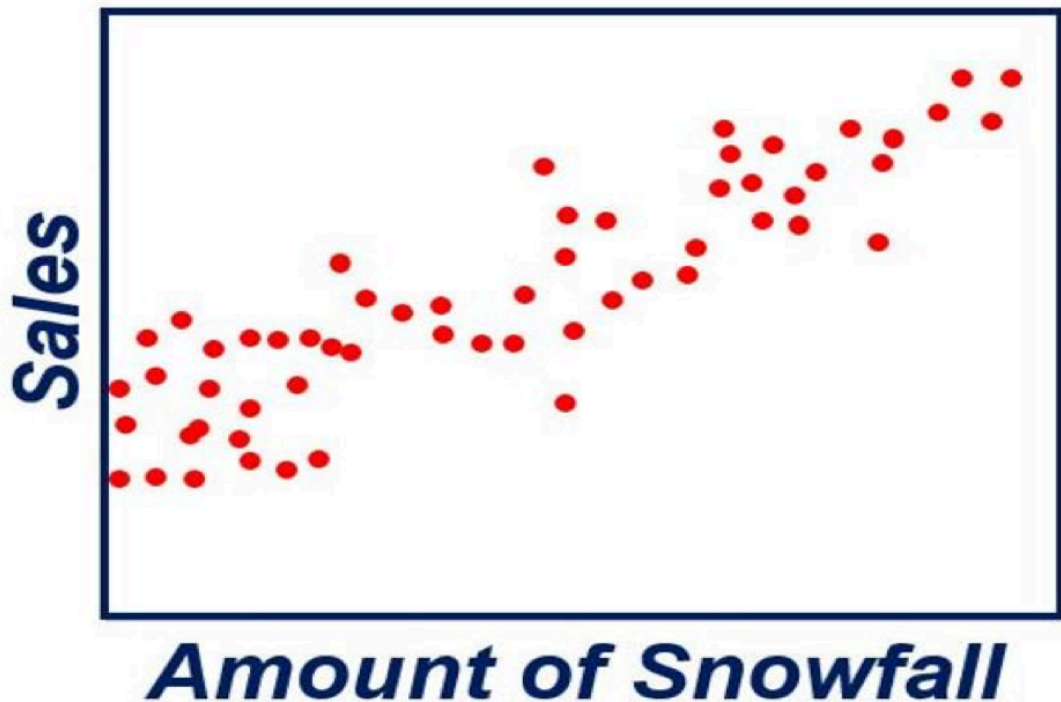
Regression Analysis



Example:

This regression analysis chart relates to the situation i.e., where you are a sales manager. Colleagues' comments that **snowfall** has an impact on **sales** figures appears to be accurate. Each red dot represents one month's worth of data, i.e., sales totals and how much it snowed that same month.

Regression Analysis



In machine learning, **Regression** is a type of supervised learning task where the model's job is to predict a **continuous outcome variable** (like car price or house price).

Based on one or more **input variables (independent)**, the model learns to predict an outcome. The **outcome** can be any **continuous number** (not categories).

Linear Regression: Predicts the target as a linear function of the input.

Linear regression models the relationship between input features (predictors) and the target variable as a **straight line**.

- **Example:** In predicting car prices, we might find that as **mileage** increases, the **price** decreases. A linear model would draw a line showing this trend.

Polynomial Regression: Fits a polynomial line to capture non-linear relationships between input and target variables.

Polynomial regression fits a **curved line** (rather than a straight one) to the data to capture **non-linear** relationships.

- **Example:** In car prices, a **non-linear pattern** might show that cars with a specific age or mileage range have higher prices, creating a curve in the plot of price versus

age/mileage.

- **Mathematical Representation for Polynomial Regression (Advanced/Optional):**

- A polynomial regression equation with a degree of (n) can be represented as:

$$y = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

- Here:
 - (x) represents the input feature,
 - (y) is the target,
 - a_n, a_{n-1}, \dots, a_0 are coefficients that the model learns to fit the curve.
-

4. Linear Regression in Depth

Linear regression is one of the simplest yet most commonly used regression techniques.

It finds the **best-fit line** that predicts the **target variable** as a function of **input features**.

- **Mathematical Representation (Optional/Advanced):**

- The relationship between the input variable (x) and the output (or target) (y) is expressed as:

- **$y = mx + b$**

- Here:
 - (m) is the **slope** of the line (how steep it is),
 - (b) is the **intercept** (where the line crosses the y-axis).
-

Example 2

Can we predict the **price** of a house, if we know the **size** ?

Resources for a case study:

- <https://www.kaggle.com/code/emrearslan123/house-price-prediction>
- <https://www.kaggle.com/code/shreayan98c/boston-house-price-prediction>
- <https://www.kaggle.com/datasets/harishkumardatalab/housing-price-prediction>

Learners will find the project as well as dataset

Let's understand linear regression using a self made example

Step 1: Create the Dataset

Here is our fictional dataset of house sizes and prices:

House Size (sq ft)	Price (\$1000s)
850	120
900	130
1200	150
1500	180
1700	200
1850	210
2000	230
2100	240

Step 2: Code to Implement Linear Regression

1. Import the Required Libraries:

```
import numpy as np
from sklearn.linear_model import LinearRegression
```

2. Create and Visualize the Data:

We'll manually define `X` as our house sizes and `y` as our house prices.

```
# House sizes (sq ft) and prices (in thousands of dollars)
X = np.array([[850], [900], [1200], [1500], [1700], [1850], [2000],
              [2100]])
y = np.array([120, 130, 150, 180, 200, 210, 230, 240])
```

3. Initialize and Train the Model:

```
model = LinearRegression()
model.fit(X, y)
```

The `LinearRegression` class from `sklearn.linear_model` creates a model that finds the best-fit line through data points in a regression task.

4. Make Predictions:

Let's use the model to predict prices for houses in the dataset and visualize the results.

```
predictions = model.predict(X)
print("Predicted prices:", predictions)
```

```
# Create a new input (random size of a house, e.g., 1600 sq ft)
new_input = np.array([[1600]])
```

```
#Make a prediction
predicted_price = model.predict(new_input)
print(f"Predicted price for a house of size {new_input[0][0]} sq ft:
${predicted_price[0] * 1000:.2f}")
```

5. Output Model Parameters:

Display the slope and intercept of the line:

```
print("Slope (m):", model.coef_[0])
print("Intercept (b):", model.intercept_)
```

Fitting a Line

Fitting a line involves **drawing a straight line** through the data points that best represents the overall trend.

The line should be **as close as possible** to each data point, **minimizing the error** between the line's predictions and actual values.

This is achieved through a technique called **least squares**, where the line is adjusted to **minimize the sum of squared differences** between predicted and actual values.

Slope and Intercept

Slope (m): This represents the **steepness of the line** and indicates how much the target variable (y) changes for a unit increase in the input variable (x).

model.coef_

- **Example:** In predicting house prices, if the slope is positive, it might mean that as the area of a house increases, so does its price.

Intercept (b): The intercept is the point where the line crosses the y-axis, representing **the predicted value** of (y) when (x = 0).

model.intercept_

- **Example:** For a house price model, the intercept could represent the base price of a very small or zero-area house.

Best-Fit Line

The best-fit line is the one that minimizes the error (or distance) between the predicted values on the line and the actual values in the dataset.

Linear Regression, `.fit()` method finds the best-fit line, as it minimizes the mean squared error by default

Error Minimization: Linear regression uses **least squares error** to determine the best-fit line, finding the line that results in the smallest sum of squared differences between actual and predicted values.

Result: Once found, this line can be used to predict new values based on input features.

R-squared(R2) and Mean Squared Error (MSE)

Mean Squared Error is the average of the squared differences between the predicted values \hat{y} and actual values y

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

The equation is annotated with boxes and labels: a large box around the fraction and sum is labeled "Mean", a box around the difference $(Y_i - \hat{Y}_i)$ is labeled "Error", and a box around the squared term 2 is labeled "Squared".

- MSE measures the average squared error between predicted and actual values.
- Lower MSE values indicate that predictions are close to the actual values.
- Since MSE squares the errors, it heavily penalizes larger errors, making it sensitive to outliers.
- The `.fit()` method minimizes MSE

R-squared, or the coefficient of determination, represents the proportion of variance in the target variable that can be explained by the independent variables in the model. It is calculated as:

$$R^2 = 1 - \frac{RSS}{TSS}$$

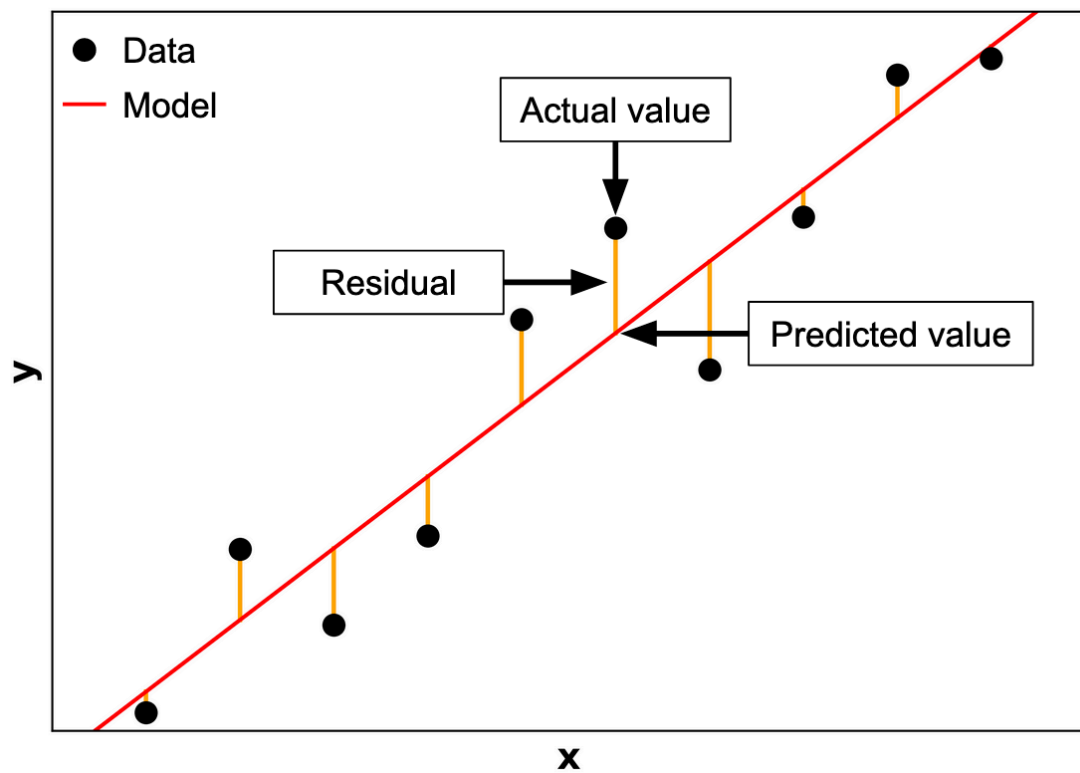
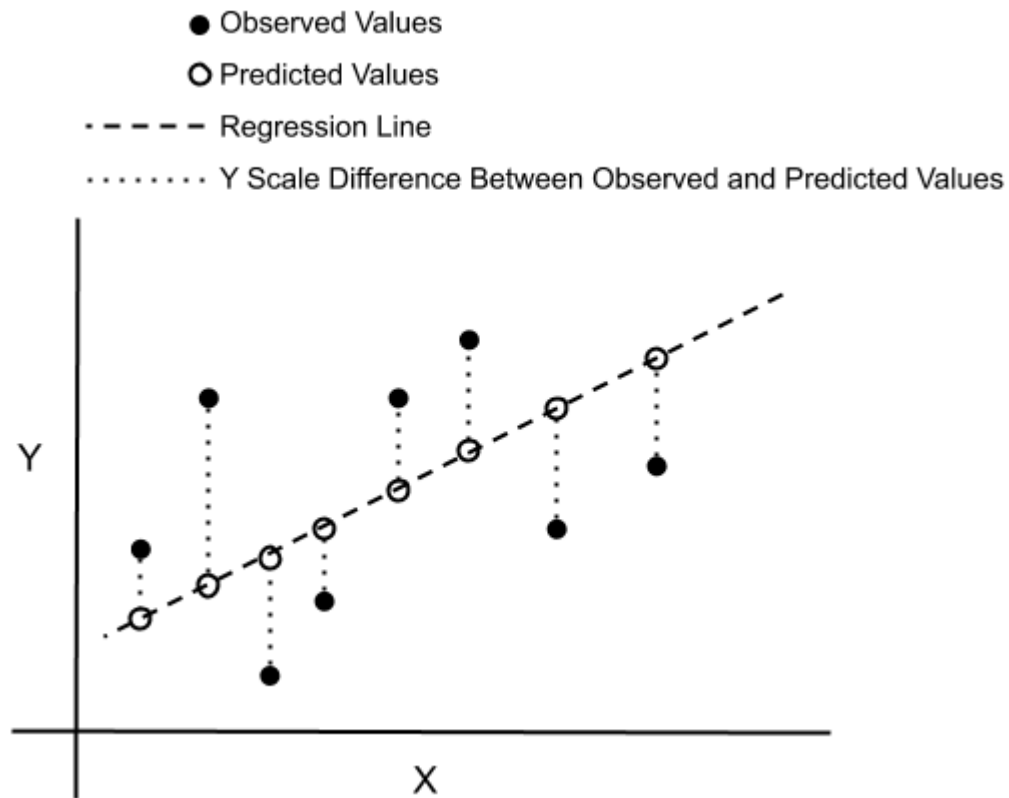
R^2 = coefficient of determination

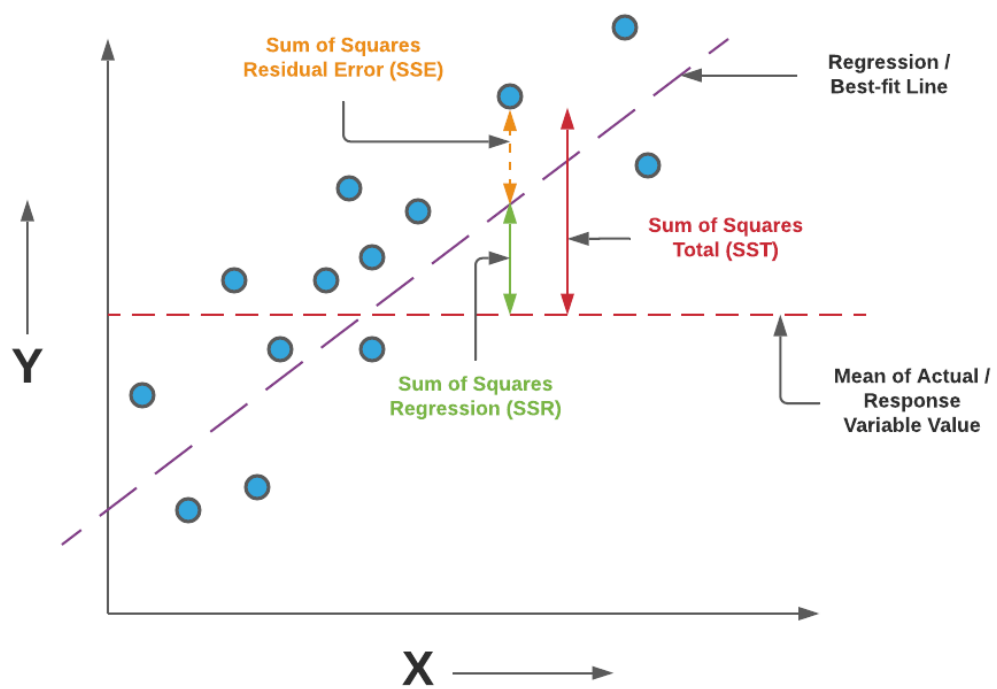
RSS = sum of squares of residuals

TSS = total sum of squares

Understand Residuals:

Residuals are the differences between observed and predicted values in regression analysis.





5. Hands-On Activity

Refer to the notebook: Simple Linear Regression

Live Exercise

Now it's your turn!

Task 1: Predicting Monthly Electricity Bill

You work as a data analyst for an energy consulting firm. The firm wants to help households predict their monthly electricity bill based on the number of hours they use electricity each day on average. This will allow homeowners to manage their energy use more effectively.

Using historical data collected from various households, your task is to build a simple linear regression model that predicts the monthly electricity bill based on daily electricity usage in hours. Here is a small dataset to help you get started:

Average Daily Usage (hours)	Monthly Bill (\$)
3	50
4	60
6	80
7	95

Average Daily Usage (hours)	Monthly Bill (\$)
8	110
10	130
12	150

Using this data, build a linear regression model and answer the following questions:

1. **What is the expected monthly bill for a household that uses electricity for 5 hours a day?**
2. **How much does the monthly bill increase, on average, for each additional hour of daily usage?**
3. **Create a visualisation that shows the data, the linear regressor, and the result value.**

By completing this exercise, you'll better understand how linear regression can be used to predict outcomes based on input features.

END

THANK YOU!

Live Exercise Solutions

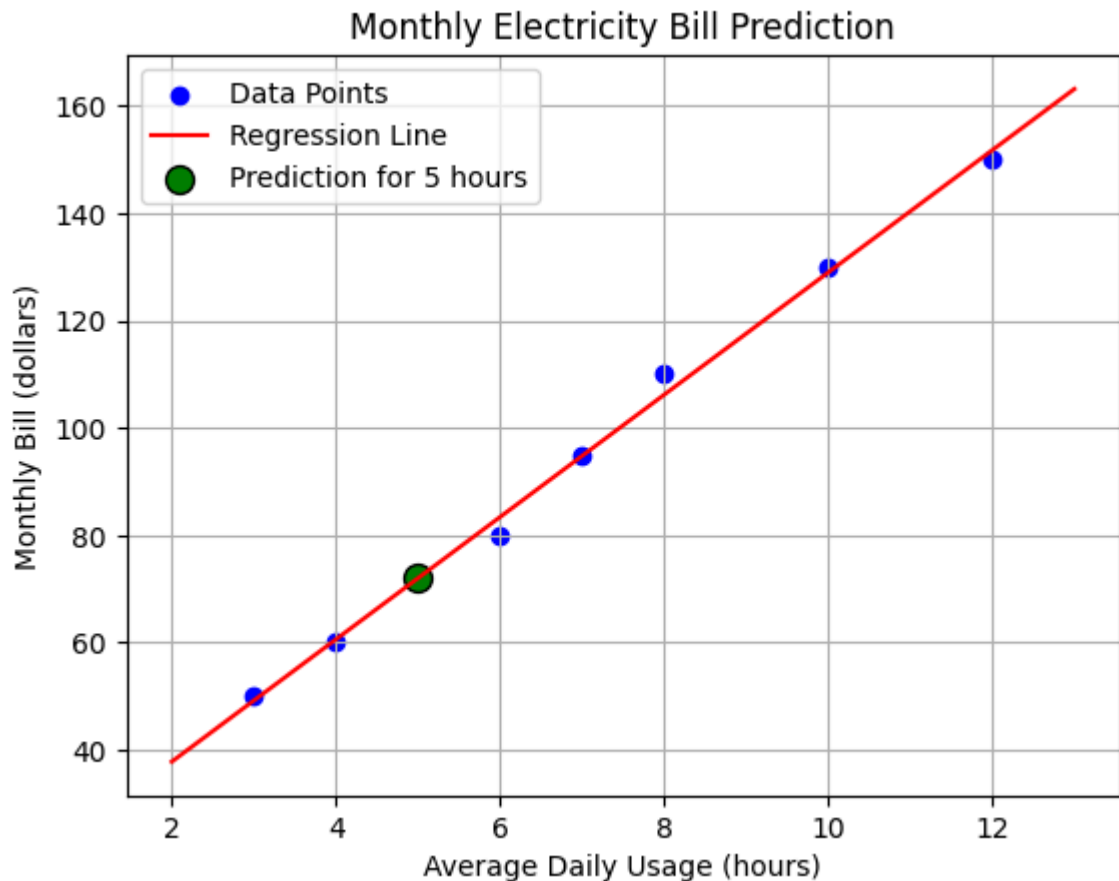
Solution: Predicting Monthly Electricity Bill Using Linear Regression

In this solution, we'll build a simple linear regression model using Python code to predict the monthly electricity bill based on average daily electricity usage. Comments are provided to explain each line and concept.

In []: `# your code`

```
Predicted monthly bill for 5 hours of daily usage: 72.00704225352112
Slope (increase per additional hour): 11.396713615023478
Intercept (base bill): 15.023474178403731
```

```
In [ ]: # your code
```



Interpretation of results

- Using the equation of the line:

bill = (slope * usage) + intercept

- we can interpret that for every additional hour of usage, the bill increases by the slope value.
- This gives us the flexibility to predict bills for any daily usage value.

Final answers based on model output:

- Predicted monthly bill for 5 hours of daily usage: \$72.00704225352112
- Increase in monthly bill for each additional hour of daily usage

Programming Interview Questions

1. What is Machine Learning, and how does it differ from traditional programming?

- Explain the concept of ML and how it uses data to learn patterns rather than being explicitly programmed for specific tasks.

2. Can you describe the different types of Machine Learning?

- Discuss the three main categories: supervised, unsupervised, and reinforcement learning.

3. What are some common applications of Machine Learning?

- Provide examples such as spam detection, recommendation systems, and image recognition.

4. What is the role of data in Machine Learning?

- Explain how data is used to train models and the importance of data quality and quantity.

5. What is overfitting, and how can it be prevented? (Optional/Advanced)

- Define overfitting and discuss techniques like cross-validation, regularization, and pruning.

6. What is supervised learning, and how does it work?

- Describe the process of training a model using labeled data to predict outcomes.

7. What is the difference between regression and classification tasks?

- Explain the distinction between predicting continuous values (regression) and categorical outcomes (classification).

8. What are some common algorithms used in supervised learning?

- List algorithms such as linear regression, logistic regression, decision trees, and support vector machines.

9. How do you evaluate the performance of a supervised learning model?

(Optional/Advanced)

- Discuss metrics such as accuracy, precision, recall, F1-score, and confusion matrix.

10. Can you explain the concept of bias-variance tradeoff in supervised learning?

(Optional/Advanced)

- Describe how bias and variance affect model performance and the importance of finding a balance between the two.

Answers: Introduction to Machine Learning

Mohammad Idrees Bhat

Tech Skills Trainer | AI/ML Consultant