# Introduction to Regression Analysis

# Regression Analysis

# Regression Analysis
## Definition

- Regression analysis, a statistical method used for estimating the relationships among variables.
- It's widely used in a variety of fields like economics, social sciences, and engineering.
- The main goal of regression analysis is to understand how the typical value of a dependent variable (the outcome you are trying to predict) changes when any one of the independent variables is varied, while the other independent variables are held fixed.

# Importance of Regression Analysis

## Decision-Making

- **Data-Driven Insights**
- **Risk Assessment**
- **Resource Allocation**

# Importance of Regression Analysis

## Forecasting

- **Predictive Power**
- **Trend Analysis**
- **Time Series Analysis**

# Importance of Regression Analysis

## Understanding Relationships Between Variables

- **Identifying Correlations**
- **Causality Insights**
- **Modeling Complex Relationships**

# Importance of Regression Analysis

## Broader Implications

- **Policy Making**
- **Innovation and Improvement**
- **Education and Training**

Terminology

# Key Terminologies

## Dependent Variable

❖ **Definition:**
The dependent variable, often denoted as $Y$, is the outcome or the variable that you are trying to predict or explain. It's called "dependent" because its values depend on the values of other variables.

❖ **Examples:**
In a study analyzing the impact of study hours on exam scores, the exam score is the dependent variable.
In business, the dependent variable could be annual sales, which might depend on factors like marketing budget, price, etc.

# Key Terminologies

## Independent Variable

❖ **Definition:**

Independent variables, denoted as $X$, are the predictors or explanatory variables that are presumed to influence or predict the dependent variable. These are the variables that you manipulate or observe to see how they affect the dependent variable.

❖ **Examples:**

In the study hours vs. exam scores example, the number of study hours is the independent variable.

# Key Terminologies

## Linear Relationship

❖ **Definition:**

A linear relationship between two variables is one where the change in one variable is associated with a proportional change in the other variable. This can be represented graphically as a straight line in a scatter plot, where one variable is plotted on the x-axis and the other on the y-axis.

❖ **Characteristics:**

**Proportionality**: In a linear relationship, for a unit change in the independent variable, there is a consistent change in the dependent variable. This is represented by the slope of the line in the graph.Z

**Line Equation:** The relationship can typically be described by the linear equation $Y=a+bX$, where is the dependent variable, $X$ is the independent variable, $b$ is the slope, and $a$ is the y-intercept

# Key Terminologies

## Linear Relationship

**Examples**:

- Height and weight can have a linear relationship; as height increases, weight also tends to increase in a proportional manner.
- In economics, there might be a linear relationship between the price of a product and its demand.

# Linear Relationship

**Scenario**

Imagine a small business that sells ice cream. We want to understand how the weather affects ice cream sales.

**Independent Variable: Temperature**

- In this case, the temperature is the independent variable (X).
- It's what we think might influence ice cream sales.
- As temperature changes, we observe how it affects sales.

**Dependent Variable: Ice Cream Sales**

- Ice cream sales are the dependent variable (Y).
- This is what we are trying to predict or explain.
- We expect sales to depend on the temperature.

**Linear Relationship**

- Hypothesis: We hypothesize that as the temperature increases, ice cream sales also increase.

# Simple Linear Regression vs  Multiple Linear Regression

# Simple Linear Regression

Formula: $y=mx+c$

- $y$ (dependent variable),
- $m$ (slope),
- $x$ (independent variable)
- $c$ (y-intercept).

Usage: It's used when the outcome variable is thought to have a linear relationship with a single predictor. For example, predicting a person's weight based on their height.

# Multiple Linear Regression

**Basic Formula:**

$Y = b_0 + b_1 X_1 + b_2 X_2 + ... + b_n X_n$

$b_0$ **is the y-intercept while** $b_1$ $b_2,...,b_n$ **are the coefficients of the respective independent variables.**

**Usage: Used when the outcome is thought to be influenced by more than one factor. For instance, predicting a house's price based on its size, location, and age.**

# Key Differences

**01** **Number of Predictors**

**02** Complexity

**03** Applications

# Assumptions in Regression Analysis

# Linearity

The relationship between the independent and dependent variables is linear. This means that any change in the independent variable is associated with a proportional change in the dependent variable.

# Independence

The observations (data points) in the dataset are independent of each other. This implies that the value of one observation does not influence or predict the value of another observation.

# Homoscedasticity

This refers to the assumption that the variance (or "spread") of the residuals (errors) is constant across all levels of the independent variables. In other words, the size of the error does not change significantly across the range of the independent variable(s).

# Normal Distribution of Errors

The residuals (differences between observed and predicted values) of the model are normally distributed. This assumption is crucial for conducting various statistical tests on the residuals, as many such tests assume normality.

# Common Pitfalls and Misunderstandings

PRB
INFORM
EMPOWER
ADVANCE

# Overfitting

- **Overfitting**: This occurs when a model is too complex and fits not only the underlying pattern in the data but also the noise. It's like a model learning the training data too well, including its anomalies and random fluctuations. This leads to poor generalization to new, unseen data.
- **Characteristics:** Overfitted models perform exceptionally well on training data but poorly on unseen data (testing data) because they've essentially "memorized" the data rather than learning the underlying trends.
- **Consequences**: Leads to inaccurate predictions and poor generalization to new data.
- **Example:** A model that accounts for every single fluctuation in stock market data, including random, non-predictive changes, will likely fail to predict future market movements accurately.
- **Prevention:** Use techniques like cross-validation, simplify the model, or collect more data.

# Underfitting

- **Underfitting**: This happens when a model is too simple to capture the underlying pattern in the data. Such a model fails to learn enough from the training data, resulting in poor performance both on the training and new data.
- **Characteristics:** Underfitted models perform poorly both on training data and unseen data. They fail to capture essential aspects of the data's structure.
- **Consequences**: Leads to inaccurate predictions due to the model's inability to recognize the underlying trend.
- **Example:** Using a linear model to predict sales when the actual relationship is nonlinear and more complex
- **Prevention:** Increase model complexity, feature engineering, or consider different assumptions.

# Correlation vs Causation:

- **Correlation:**

**Definition:** Correlation indicates a relationship or association between two variables. When one variable changes, the other tends to change in a specific way (positively or negatively).

**Example:** There is a correlation between ice cream sales and temperature—sales tend to increase when the temperature rises.

- **Causation:**

**Definition:** Causation implies that changes in one variable cause changes in another. It's a much stronger assertion than correlation.

**Example:** Increasing the temperature does not cause an increase in ice cream sales; they are correlated because both are influenced by a third factor (e.g., seasonality).

## Misunderstanding:

- **The Fallacy:** A common misunderstanding is assuming that because two variables are correlated, one must cause the other. This can lead to incorrect conclusions.
- **Importance**: Understanding the difference is crucial in research and data analysis to avoid false assumptions and inaccurate predictions.
- **Prevention:** To establish causation, more rigorous testing and experimental designs are needed, beyond just finding correlations.