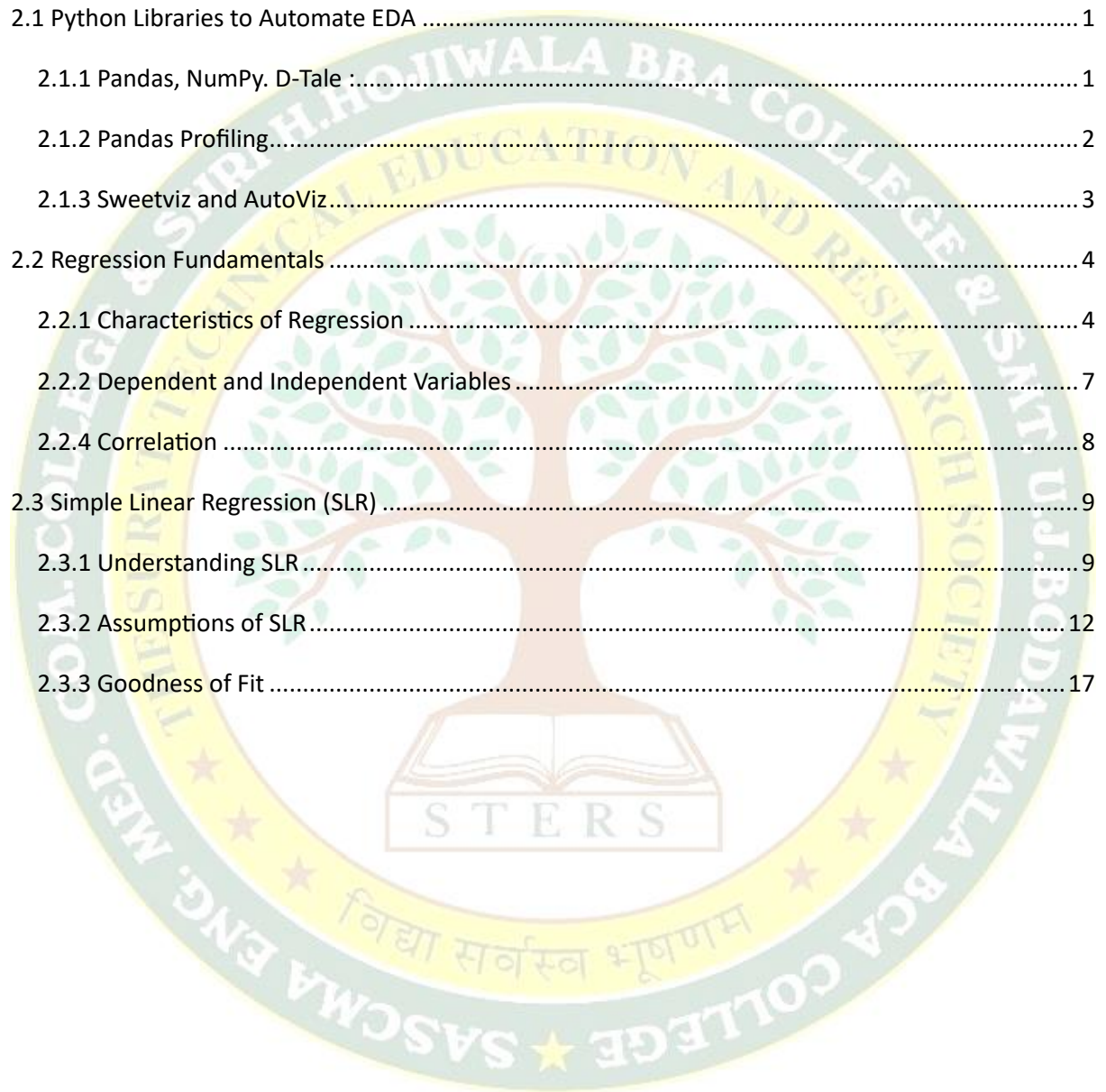
The logo is a circular emblem. The outermost ring is green and contains the text 'SHRI H. HOJIWALA BBA COLLEGE & RESEARCH SOCIETY' in white. Inside this is a yellow ring with the text 'THE SURAT TECHNICAL EDUCATION AND RESEARCH SOCIETY' in green. The center of the logo features a brown tree with green leaves, which is growing out of an open book. Below the book, the word 'STERS' is written in a stylized font. At the bottom of the logo, there is a green banner with the text 'SASCMA ENG. MED. U.J. BODAWALA BCA COLLEGE' in white, and a small yellow star is positioned between 'SASCMA' and 'U.J. BODAWALA'.

## **Unit-2**

# **Automated EDA (Exploratory Data Analysis)**

## Contents

2.1 Python Libraries to Automate EDA .....	1
2.1.1 Pandas, NumPy, D-Tale :.....	1
2.1.2 Pandas Profiling.....	2
2.1.3 Sweetviz and AutoViz.....	3
2.2 Regression Fundamentals .....	4
2.2.1 Characteristics of Regression .....	4
2.2.2 Dependent and Independent Variables .....	7
2.2.4 Correlation .....	8
2.3 Simple Linear Regression (SLR) .....	9
2.3.1 Understanding SLR.....	9
2.3.2 Assumptions of SLR.....	12
2.3.3 Goodness of Fit .....	17



## 2.1 Python Libraries to Automate EDA

### 2.1.1 Pandas, NumPy, D-Tale :

Exploratory Data Analysis (EDA) is the process of understanding datasets by summarizing their main characteristics using statistical graphics, plots, and information tables. In Python, three foundational libraries used to support and automate EDA tasks are Pandas, NumPy, and D-Tale. Each of these libraries plays a different role but contributes significantly to making data preparation and understanding more efficient.

#### Pandas :

**Pandas** is perhaps the most widely used data manipulation library in Python. It provides powerful data structures such as Series (1D) and DataFrame (2D) which are specifically designed to handle tabular data. With Pandas, you can load data from various file formats like CSV, Excel, JSON, and SQL databases using functions like `read_csv()`, `read_excel()`, and `read_sql()`. Once loaded, the data can be filtered, grouped, sorted, cleaned, and reshaped with ease. Functions like `head()`, `tail()`, `describe()`, `info()`, `value_counts()`, and `groupby()` are particularly useful for EDA. For example, `describe()` gives statistical summaries of numerical columns, and `value_counts()` shows the frequency of unique values in a categorical column.

#### **Key EDA Features in Pandas:**

- **Data Loading:** `read_csv()`, `read_excel()`, `read_json()`
- **Data Cleaning:** handle missing values (`dropna()`, `fillna()`), duplicates
- **Summarization:** `info()`, `describe()`, `value_counts()`
- **Filtering & Slicing:** Select rows/columns using `loc`, `iloc`
- **Grouping and Aggregation:** `groupby()`, `agg()`, `pivot_table()`
- **Visualization Integration:** Works well with Matplotlib and Seaborn

#### NumPy :

NumPy (Numerical Python) complements Pandas by providing efficient operations on large arrays and matrices of numeric data. It supports a wide range of mathematical and statistical operations such as mean, median, standard deviation, dot products, and linear algebra operations. Although Pandas is built on top of NumPy, the latter is especially useful when dealing with low-level numerical data manipulation. For instance, NumPy arrays are faster and more memory-efficient than Python lists for numerical computations.

#### **Role in EDA:**

- **Efficient Computation:** Fast array operations (mean, median, std, etc.)
- **Numerical Summaries:** Used under the hood by Pandas and other libraries
- **Integration:** Works smoothly with Pandas, SciPy, and visualization libraries



**D-Tale:**

D-tale brings interactivity to your data analysis workflow. It combines Pandas with a lightweight Flask server and Dash frontend to provide a web based interface for viewing and editing Pandas DataFrames. You can sort columns, view statistical summaries, filter rows, and even generate charts without writing additional code. The typical workflow involves loading a DataFrame using Pandas and passing it to `dtale.show(df)`, which launches a web interface locally. This makes it easier for analysts and non-technical users to explore the data interactively without diving deep into Python code. D-Tale is a library that lets you view and analyze Pandas DataFrames with an interactive web interface, like Excel + pandas + Plotly.

**Features:**

- Launch a web UI for a DataFrame
- View summary statistics
- Sort, filter, edit data interactively
- Visualize distributions (bar charts, histograms, etc.)
- Detect nulls, outliers, correlations
- Auto-generates Python code for operations

**In summary**, Pandas is the core for data manipulation, NumPy enhances numerical performance, and D-Tale provides an interactive interface for data inspection. Together, they lay the groundwork for effective manual or semi-automated EDA.

```
import pandas as pd
import dtale
df = pd.read_csv("data.csv")
dtale.show(df)
```

### 2.1.2 Pandas Profiling

Pandas Profiling is a robust and widely used tool for automated EDA in Python. As of version 3.0 and above, it is now maintained under the name `ydata-profiling`. This tool allows users to generate a detailed HTML report of a DataFrame with just one line of code. It provides an overview of data types, missing values, correlations, summary statistics, and even warnings for potential data quality issues.

To install the latest version, the recommended command is **`pip install ydata-profiling`**. Once installed, generating a report is straightforward. After importing `ProfileReport` from `ydata_profiling`, you can create a profile by passing your DataFrame to it. The `profile.to_file('eda_report.html')` method generates a fully interactive report saved as an HTML file.

The generated report contains several sections: an overview of the dataset, warnings about problematic features (like high cardinality or constant columns), details of each variable including statistics and histograms, a correlation matrix, and missing value visualizations. The correlation section is particularly useful as it shows both Pearson and Spearman correlation scores between numeric variables, helping identify multicollinearity.

Another strength of Pandas Profiling is its ability to handle both categorical and numerical features seamlessly. It automatically identifies the data types and provides appropriate analysis for each. Additionally, you can customize the report by disabling specific sections (e.g., correlations, duplicates) or enabling advanced options for large datasets (e.g., `minimal=True`).

Although powerful, Pandas Profiling has limitations. It can be slow or memory-intensive on very large datasets (e.g., more than 1 million rows or hundreds of columns). It is best suited for small to medium datasets, or as a first-pass EDA tool before performing deeper custom analysis.

**Code:**

```
import pandas as pd
from ydata_profiling import ProfileReport # for version >= 3.0
df = pd.read_csv("train.csv")
profile = ProfileReport(df, title="EDA Report", explorative=True)
profile.to_file("eda_report.html")
```

### 2.1.3 Sweetviz and AutoViz

Sweetviz and AutoViz are two modern Python libraries designed to perform EDA automatically while focusing on visual storytelling. Both tools create high-quality, interactive, and informative visual reports with minimal code.

Sweetviz is known for its aesthetically pleasing and feature-rich reports. It works particularly well for both single dataset analysis and comparing two datasets (e.g., train vs test). It can also compare values based on a target column, which is useful in classification tasks. Installation is simple via `pip install sweetviz`. After importing the library, you can analyze a DataFrame using `sv.analyze(df)` and generate the report using `report.show_html('sweetviz_report.html')`. The resulting report includes distributions, target associations, correlations, and warnings, all in a clean and interactive layout.

Sweetviz highlights interesting features automatically. For example, it identifies columns with missing values, low variance, or high cardinality, and summarizes them using visual cues. It supports both regression and classification problems by showing how each feature relates to the target.

```
import sweetviz as sv
```

```
import pandas as pd
df = pd.read_csv("train.csv")
report = sv.analyze(df)
report.show_html("sweetviz_report.html")
```

AutoViz, developed by Ram Seshadri, offers another automated visualization tool that works well in Google Colab and Jupyter environments without strict dependency issues. You can install it via `pip install autoviz`. It supports input in the form of CSV files or Pandas DataFrames and automatically detects variable types, selects appropriate charts, and builds dashboards. AutoViz produces histograms, box plots, scatter plots, violin plots, heatmaps, pair plots, and more. It's especially useful when working on classification and regression problems, as it builds charts based on the target variable (if provided). The plots are rendered directly in the notebook, making it more hands-on than Sweetviz, which relies on an external HTML report. While it does not provide deep statistical summaries like Pandas Profiling, its visualization capabilities make it a go-to tool for fast, rich data exploration.

In summary, Sweetviz is better for interactive and comparative HTML reports with strong visual storytelling, while AutoViz excels in quick, notebook-based visualization with minimal code and fewer compatibility problems.

```
from autoviz.AutoViz_Class import AutoViz_Class
AV = AutoViz_Class()
df = AV.AutoViz("train.csv")
```

## 2.2 Regression Fundamentals

### Introduction to Regression

Regression is a fundamental statistical and machine learning technique used to model the relationship between a dependent variable (target) and one or more independent variables (predictors or features). It allows us to make predictions, understand patterns, and quantify the strength and direction of relationships between variables.

In real life, regression is like predicting someone's electricity bill based on how many hours they use an air conditioner. As one variable changes, we try to estimate how another will change with it.

#### 2.2.1 Characteristics of Regression

##### 1. Relationship Analysis

###### What it means:

Regression helps us see how two things are connected. It shows if changing one thing affects the other.



**Example:**

Suppose you're trying to understand how the number of hours you study affects your exam marks. Regression will help you see if studying more really leads to higher marks.

**Analogy:**

Think of a remote control and a TV. When you press the volume button (independent variable), the volume changes (dependent variable). Regression tries to understand and measure this type of cause-effect relationship.

**2. Prediction****What it means:**

Once regression understands the relationship, it can **predict future values**.

**Example:**

If you studied 4 hours and scored 60 marks, and when you studied 8 hours you scored 90 marks, then regression can help **predict your marks if you study for 6 hours**.

**Analogy:**

Think of Google Maps. When you enter your starting point and speed, it predicts how long your trip will take. Regression does something similar: **using past patterns to make future predictions**.

**3. Direction and Strength****What it means:**

Regression tells you **how two things move together**, and **how strongly they're connected**.

**There are two directions:**

- **Positive Relationship:**

If one increases, the other increases.

Example: The more hours you sleep, the more energy you have.

- **Negative Relationship:**

If one increases, the other decreases.

Example: The more time you spend on your phone at night, the less sleep you get.

**Strength of Relationship:**

- If studying time and marks are **very closely related**, that's a **strong relationship**.
- If there's no pattern, the relationship is **weak or none**.

**Analogy:**

Think of holding two magnets:

- If they strongly pull toward each other → strong relationship.
- If they barely move → weak relationship.
- If they push away → negative relationship.

#### 4. Line of Best Fit

What it means:

This is a **line drawn through your data points** that best shows the trend or pattern.

- It is not just a random line — it's **mathematically calculated** to reduce the difference between actual values and predicted values.
- It gives you a **formula** that you can use to make predictions.

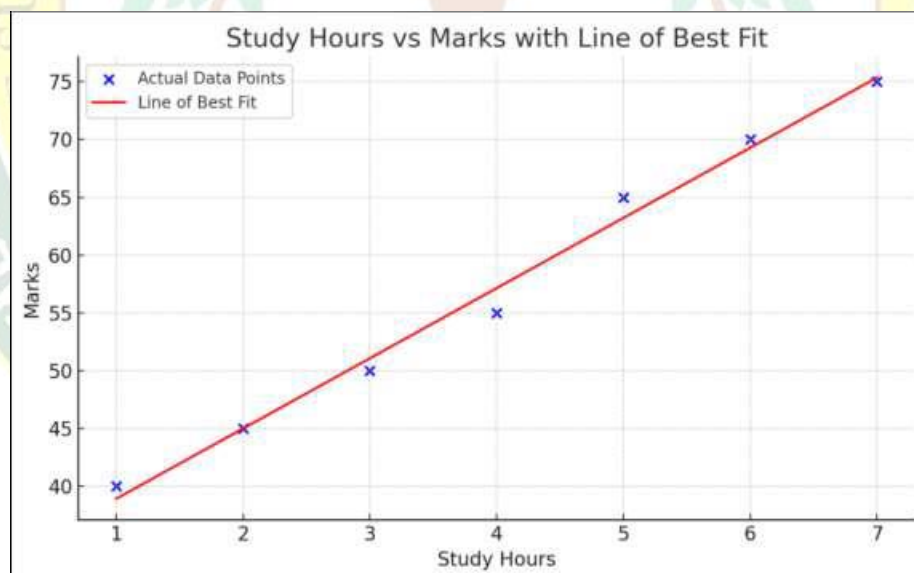
**Example:**

Imagine you have a bunch of **dots on a graph**.

Each dot shows something like:

- How many hours you studied
- What marks you got

Now, you want to draw one straight line that best shows the general pattern of all those dots.





**Recap in Table Form (characterstics)**

Feature	In Simple Words	Real-Life Analogy
Relationship Analysis	Checks if two things are connected	Does more studying lead to better marks?
Prediction	Helps guess the future based on past data	Like predicting rain based on cloud patterns
Direction & Strength	Tells if change is together/opposite, and how strong	Magnets pulling/pushing each other
Line of Best Fit	A line that best shows the trend	A straight path through footprints on sand

**2.2.2 Dependent and Independent Variables****Independent Variable (also called Input, Predictor, or Feature):**

The variable(s) used to make predictions. This is what you already know.

Example: Hours of study, temperature, years of experience.

**Dependent Variable (also called Output or Target):**

The variable you want to predict or explain.

Example: Exam score, electricity bill, salary.

**Analogy:**

Think of making tea:

- The amount of sugar and milk you use are independent variables.
- The sweetness of the tea is the dependent variable.

Changing the input changes the output!

**2.2.3 Covariance****Covariance measures how two variables change together:**

- Positive Covariance: Both variables increase or decrease together.
- Negative Covariance: One variable increase while the other decreases.

**Analogy:**

Imagine two friends who jog together:

- If both increase their speed at the same time, the covariance is positive.
- If one speeds up and the other slows down, the covariance is negative.

**Limitations:**

Covariance only tells us the direction of the relationship, not the strength.

Also, the magnitude of covariance is not standardized, so it's hard to compare between datasets.

**2.2.4 Correlation**

Correlation is a normalized version of covariance that shows both the direction and strength of a linear relationship. - Values range from -1 to +1:

+1: Perfect positive correlation

0: No correlation

-1: Perfect negative correlation

**Analogy:**

Think of correlation as a friendship scale:

- +1: Best friends (do everything together)
- 0: Strangers (no relation)
- -1: Enemies (when one wins, the other loses)

Correlation makes it easier to interpret and compare relationships because it is unitless and scaled.

**Summary Table**

Concept	Meaning	Analogy
Regression	Predicting a value using known variables	Predicting salary using years of experience
Dependent Variable	The value you want to predict	Exam score
Independent Variable	The value you use to make predictions	Hours studied

Covariance	Direction of joint variation between two variables	Two friends jogging at the same or opposite speeds
Correlation	Direction and strength of linear relationship between variables	Closeness of friendship (-1 to +1)

## 2.3 Simple Linear Regression (SLR)

Simple Linear Regression (SLR) is a technique used in statistics and machine learning to model the relationship between two variables: one independent variable ( $x$ ) and one dependent variable ( $y$ ). It helps us understand how changes in the independent variable influence the dependent variable and allows us to make predictions.

### 2.3.1 Understanding SLR

In Simple Linear Regression, the goal is to fit a straight line through a set of data points such that the line best represents the trend in the data. The line can then be used to predict unknown values.

We use SLR when we want to answer questions like:

- "Does the number of hours study affect exam marks?"
- "Does the outside temperature affect the number of ice creams sold?"

In such cases, we assume a linear relationship between the input ( $x$ ) and output ( $y$ ), and SLR helps us find the best-fitting straight line.

#### 2.3.1.1 Fitting a Straight Line ( $y = mx + b$ )

The straight-line equation used in Simple Linear Regression is:  $y = mx + b$

Where:

- $y$  is the dependent variable (the output or predicted value)
- $x$  is the independent variable (the input)
- $m$  is the slope of the line (how much  $y$  changes when  $x$  increases by 1)
- $b$  is the intercept (the value of  $y$  when  $x$  is 0)

#### Example:

Let's say we want to predict marks based on study hours. After calculating, we find:

$$y = 10x + 30$$

This means:

- For every extra hour studied, marks increase by 10 points (slope = 10)
- If no studying is done, the student still scores 30 marks (intercept = 30)

#### How to find $m$ and $b$ :

Given  $n$  data points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ :

Here,



1. **Slope (m):**  $m = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$
2. **Intercept (b):**  $b = \frac{\sum y - m(\sum x)}{n}$

Here,

- $\sum x$ : sum of all x-values
- $\sum y$ : sum of all y-values
- $\sum xy$ : sum of products of x and y
- $\sum x^2$ : sum of squares of x-values
- $n$ : number of data points

Given this data:

x (Hours Studied)	y (Marks Scored)
1	40
2	50
3	60
4	70

So the final regression line is:

$$y = 10x + 30$$

Step 1: Find necessary values

- $n = 4$
- $\sum x = 1 + 2 + 3 + 4 = 10$
- $\sum y = 40 + 50 + 60 + 70 = 220$
- $\sum xy = 1 \times 40 + 2 \times 50 + 3 \times 60 + 4 \times 70 = 600$
- $\sum x^2 = 1^2 + 2^2 + 3^2 + 4^2 = 30$

Step 2: Calculate the slope (m)  $m = \frac{4(600) - (10)(220)}{4(30) - 10^2} = \frac{2400 - 2200}{120 - 100} = \frac{200}{20} = 10$

Step 3: Calculate the intercept (b)  $b = \frac{220 - 10(10)}{4} = \frac{120}{4} = 30$

So the final regression line is:  $y = 10x + 30$

### 2.3.1.2 Identifying Dependent and Independent Variables

Understanding which variable is dependent and which is independent is key:

**Independent Variable (x):**

- o The input
- o The cause or factor we can control
- o Examples: Hours studied, advertising budget, temperature

**Dependent Variable (y):**

- o The output
- o The effect or result that depends on x
- o Examples: Exam marks, sales, electricity usage

**Example:**

If we are analyzing how temperature affects ice cream sales:

- o Temperature is independent (x)
- o Ice cream sales are dependent (y)

Code:

```
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
# Step 1: Load the CSV file
df = pd.read_csv("study_hours_vs_marks.csv")
# Step 2: Separate features and target
X = df[['Study_Hours']] # Independent variable (2D)
y = df['Marks'] # Dependent variable (1D)
model = LinearRegression()
model.fit(X, y)
y_pred = model.predict(X)
```

### 2.3.2 Assumptions of SLR

An assumption is like a rule that should be true so that linear regression can give us correct and trustworthy results. When we use Simple Linear Regression (SLR), we are trying to draw a straight line that best fits the data between an independent variable (X) and a dependent variable (Y). But for the results to be correct and meaningful, certain assumptions must be true. If these assumptions are violated, our predictions may be wrong.

#### What does it mean?

In Simple Linear Regression, we are trying to model the relationship between:

- X (the independent variable — like study hours), and
- Y (the dependent variable — like exam marks).

The assumption of linearity means that Y increases or decreases in a straight-line fashion as X increases or decreases. This means the relationship can be written as:

$$Y = mX + b$$

Where:

- m is the slope (how much Y changes for each unit change in X)
- b is the intercept (Y value when X = 0)

#### In Simple Words,

Imagine you are plotting a graph with how many hours you study on the x-axis and how many marks you get on the y-axis.

- If studying 1 hour gives 50 marks
- Studying 2 hours gives 60 marks
- Studying 3 hours gives 70 marks

Then the data points form a straight line — this is linearity. It means the pattern is consistent and predictable: for every extra hour, you gain 10 more marks.

But if the graph looks like:

- 1 hour → 50 marks
- 2 hours → 90 marks
- 3 hours → 55 marks

Then the pattern is not linear — it's zigzag or curved. That breaks the linearity assumption.



### 2.3.2.1. Linearity

#### What does it mean?

In Simple Linear Regression, we are trying to model the relationship between:

- X (the independent variable — like study hours), and
- Y (the dependent variable — like exam marks).

The assumption of linearity means that Y increases or decreases in a straight-line fashion as X increases or decreases.

This means the relationship can be written as:

$$Y = mX + b$$

Where:

- m is the slope (how much Y changes for each unit change in X)
- b is the intercept (Y value when X = 0)

#### In Simple Words,

Imagine you are plotting a graph with how many hours you study on the x axis and how many marks you get on the y-axis.

If studying 1 hour gives 50 marks

Studying 2 hours gives 60 marks

Studying 3 hours gives 70 marks

Then the pattern is not linear — it's zigzag or curved. That breaks the linearity assumption.

#### Why does this matter?

If the relationship between X and Y is not linear, then using a straight-line equation (SLR) to predict Y will not work well — it may give wrong or misleading results.

You might then need to use:

- Polynomial Regression (for curves)
- Transformations (like log or square root)
- Non-linear models (for more complex patterns)

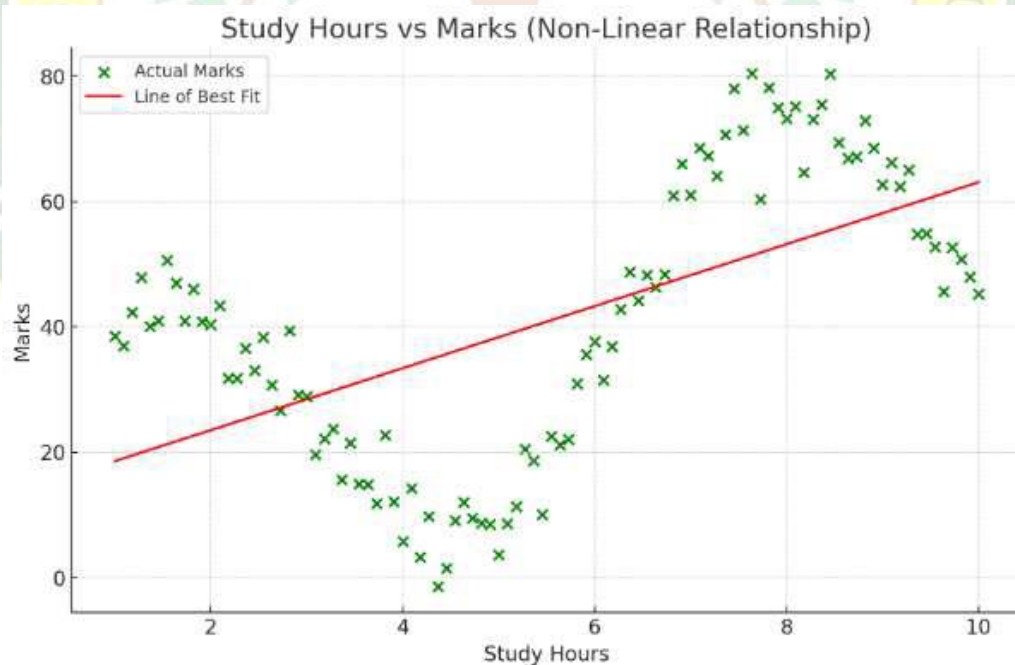
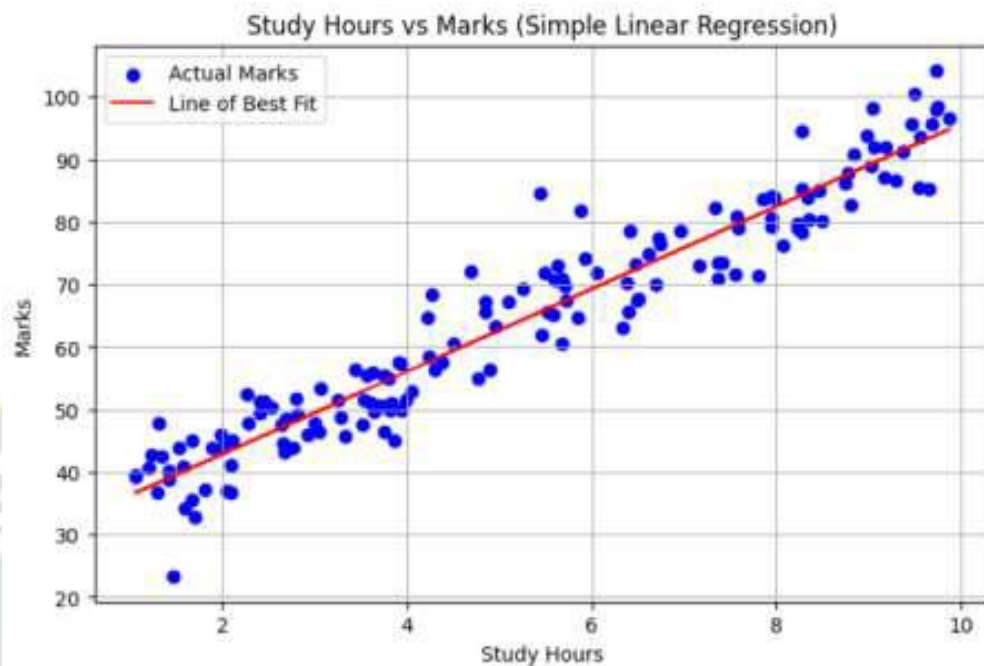
#### How to check linearity?

You can plot a scatter plot of X vs. Y.

- If the points form a cloud that fits around a straight line → good linearity

- If the points curve or form strange patterns → linearity is violated You can also look at residual plots (errors vs. X) — they should look like a

random cloud, not a curve. n holds. If they form a curve or inconsistent shape, this assumption is broken.



### 2.3.2.2 Independence of Errors

When we use Simple Linear Regression to make predictions, it often doesn't match the actual value exactly. The gap between the actual value and the predicted value is called an **error** (or **residual**).

Now, the **assumption of independence** says:

Each error should be on its own — it should not be connected or influenced by any other error.

**Let's break it down with a simple example:**

Imagine you're predicting students' marks based on how many hours they study.

You predict 60 marks for Riya, but she actually got 65 → error = -5

- Then you predict 70 for Aarav, but he got 75 → error = -5 again
- Then for Siya → error = -5 again

This kind of repeating pattern means the errors are connected. They are **not independent**. This is bad, because it shows the model is making similar mistakes over and over — maybe missing something important.

**Real-Life Analogy:**

Think of students taking a quiz **alone**. One student makes a silly mistake — this doesn't affect the next student. That's **independent** behaviour.

Now imagine students **copying** each other. If the first student makes a mistake, others do the same. That's **dependent** behaviour.

**Why does this matter?**

If the errors are connected, your model might:

- Be missing something important,
- Make wrong predictions in the future,
- Or give you false confidence in the results.

That's why this assumption is so important in linear regression

### 2.3.2.3 Homoscedasticity

**What Is It?**

**Homoscedasticity** means that the **spread (or variance) of errors** in a regression model should remain **constant** across all values of the independent variable (X).



In other words:

- Whether X is small, medium, or large, the **errors between actual and predicted Y values should not increase or decrease significantly.**

If the **spread of errors increases or decreases** as X increases, we say the model suffers from **heteroscedasticity** (the opposite of homoscedasticity), and this violates the assumption.

### Why Is It Important ?

When errors are **unevenly spread** (i.e., heteroscedastic), the model might:

- Make **unreliable predictions** in certain ranges of data.
- Produce **incorrect confidence intervals** and **hypothesis test results.**

This means you might wrongly conclude that a variable is important (or not) based on bad statistics.

### Example:

Suppose you're predicting **students' marks based on study hours:**

- For students who study 1–3 hours: your prediction is off by  $\pm 5$  marks.
- For students who study 8–10 hours: your prediction is off by  $\pm 20$  marks.

That increasing error range shows **heteroscedasticity.**

But if the prediction is always off by about the range of near marks regardless of study hours, then the model is **homoscedastic.**

### Analogy:

Think of **throwing darts** at a board from different distances:

- If your darts always land with the same spread — whether you're close or far — that's **homoscedasticity.**
- If your spread becomes wider the farther you go — your aim gets worse — that's **heteroscedasticity.**

### 2.3.2.4 Normality of Errors

#### What Is It?

This assumption says the **errors (residuals)** in the regression model should follow a **normal distribution** — that is, a bell-shaped curve.

This doesn't mean your data (X or Y) must be normal, just the errors (differences between predicted and actual Y values).

#### Why Is It Important ?

- Many regression-related calculations (like p-values, t-tests, and confidence intervals) assume **normality of errors.**

- If errors are not normally distributed, these results can be misleading, especially in **small datasets**.

**Example:**

You're predicting house prices:

- Most predictions are very close (off by \$5,000 or less).
- A few predictions are off by \$15,000.
- Very few are off by \$50,000.

This is a **normal error pattern**: most errors are small, few are large. But if errors are all over the place, or heavily skewed in one direction, they **don't follow a normal distribution** — violating this assumption.

### 2.3.3 Goodness of Fit

**What is "Goodness of Fit"?**

When you build a **regression model**, you are drawing a line (or equation) that tries to predict something based on input data. For example:

- Predicting **marks** based on **study hours**
- Predicting **sales** based on **advertising budget**
- Predicting **house price** based on the **number of rooms**

But once you build the model, a very important question arises:

"How well is my model actually performing?"

"Is this model truly understanding the pattern in the data, or is it just guessing?"

This is where Goodness of Fit comes in.

It tells you:

- How **accurate** your model's predictions are
- How much of the **variation in the result (Y)** can be explained by your input variable (X)
- Whether your model is a **good match for the data**

**Why is Goodness of Fit Important?**

Not every model is useful. Sometimes:

- A model might fit the training data perfectly, but perform poorly on new data.
- A model might include too many variables just to improve accuracy slightly.
- A model might accidentally capture patterns that happened only by luck.

**Goodness of Fit helps you:**

- Understand if your model can be trusted
- Avoid being fooled by random or noisy data
- Decide whether the model is good enough for real-world predictions or decisions.

**Summary**

- Goodness of Fit tells us how well our model fits the actual data.
- It's a way to measure how much we can trust the model's predictions.
- It's the foundation for evaluating any regression model.
- It can be calculated using measures like  $R^2$ , Adjusted  $R^2$ , Significance tests, and RSS (which we will explore next).

### 2.3.3.1 $R^2$ and Adjusted $R^2$

#### What is $R^2$ ?

$R^2$ , also known as R-squared or the coefficient of determination, is a number that tells us:

How well does the regression model explain the variation in the output (Y) using the input (X)?

It measures the strength of the relationship between the independent variable and the dependent variable in a regression model.

#### The Range of $R^2$

$R^2$  values range from 0 to 1 (or 0% to 100%).

- $R^2 = 0$  means the model explains none of the variation in Y. The model is no better than simply guessing the average.
- $R^2 = 1$  means the model explains all of the variation in Y perfectly. Every prediction lies exactly on the regression line.
- Values closer to 1 mean the model is very good, while values closer to 0 mean the model is not helpful.

#### Mathematical View (Simplified)

$R^2$  is calculated as:

Formula:

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

Where:

- **RSS** = Residual Sum of Squares (total prediction error)
- **TSS** = Total Sum of Squares (total variation in the actual values)

So, if your prediction errors (RSS) are small,  $R^2$  becomes high — showing a good model.

#### What is Adjusted $R^2$ ?

Adjusted  $R^2$  is an improved version of  $R^2$  that **corrects** for the number of predictors (independent variables) in a regression model.

#### Why?

Because  **$R^2$  always increases** (or stays the same) when you add more variables — even if those new variables are **useless**.

Adjusted  $R^2$  tells you:



“Is this extra variable **really helping**, or is it just **adding noise**?”

### Formula (Conceptually):

$$\text{Adjusted } R^2 = 1 - \left( \frac{(1 - R^2)(n - 1)}{n - k - 1} \right)$$

Where:

- **n** = number of observations (rows)
- **k** = number of independent variables (X's)
- **R<sup>2</sup>** = regular R-squared

This formula **penalizes the model** for adding variables. If a new variable **does not improve the model enough**, Adjusted R<sup>2</sup> **goes down**.

### Real-Life Analogy:

Imagine you're preparing for an exam.

- **R<sup>2</sup>** is like counting how many **books** you've added to your study desk.
- **Adjusted R<sup>2</sup>** checks if each book is actually **helping you score better**.

If you pile on books just to feel productive, R<sup>2</sup> goes up. But if only a few books are truly helpful, Adjusted R<sup>2</sup> shows that.

### When to Use Adjusted R<sup>2</sup>?

Always use Adjusted R<sup>2</sup> when:

- You have **multiple independent variables**
- You are comparing **different models**
- You want to avoid **overfitting**

### R<sup>2</sup> vs Adjusted R<sup>2</sup> – Key Differences

Feature	R <sup>2</sup>	Adjusted R <sup>2</sup>
Increases with more variables	Always (or stays the same)	Only if the new variable helps
Penalizes useless variables	No	Yes
Usefulness	Good for simple models	Best for multiple-variable regression
Helps avoid overfitting?	No	Yes

### Summary

- Adjusted R<sup>2</sup> helps you judge if adding new variables **improves your model or not**.
- It's a **more honest measure** of model quality in multiple regression.
- If Adjusted R<sup>2</sup> increases, the new variable is likely **useful**.
- If Adjusted R<sup>2</sup> decreases, the variable likely just **adds noise**.

### 2.3.3.3 Residual Sum of Squares (RSS)

#### What is RSS?

**Residual Sum of Squares (RSS)** is a measure of how **far off your model's predictions** are from the actual values.

In other words, it adds up all the **squared errors** between what your model predicted and what really happened.

$$RSS = \sum (Y_i - \hat{Y}_i)^2$$

Where

- $Y$  = actual value
- $\hat{Y}$  = predicted value

We **square the errors** to:

- Avoid negatives canceling out positives
- Give **larger errors more impact**

**Example:**

Actual Marks (Y)	Predicted Marks ( $\hat{Y}$ )	Error (Y - $\hat{Y}$ )	Squared Error
60	65	-5	25
70	72	-2	4
80	75	5	25
<b>Total</b>			<b>54 (RSS)</b>

The total error (RSS) is 54.

#### Why is RSS important?

- A **smaller RSS** means the predictions are **closer to the actual values** — the model fits the data well.
- A **larger RSS** means the model is **making bigger mistakes** — the fit is poor.

When we calculate  $R^2$ , we actually use RSS like this:

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

So RSS directly affects **how good your  $R^2$  score** will be.

#### Analogy:

Imagine you're throwing darts at a target.

The bullseye is the actual value, and where your dart lands is the predicted

value.

RSS is like the **total squared distance** of all your throws from the center.

- The closer you are, the **lower the RSS**.
- The more off-target you are, the **higher the RSS**.

**Summary:**

Concept	Meaning	Analogy
Significance of $R^2$	Tells us if $R^2$ happened by chance or reflects a real relationship	Guessing 10 matches in a row — lucky or real?
Residual Sum of Squares (RSS)	Total squared prediction errors — lower is better	Dart game — total distance from bullseye

