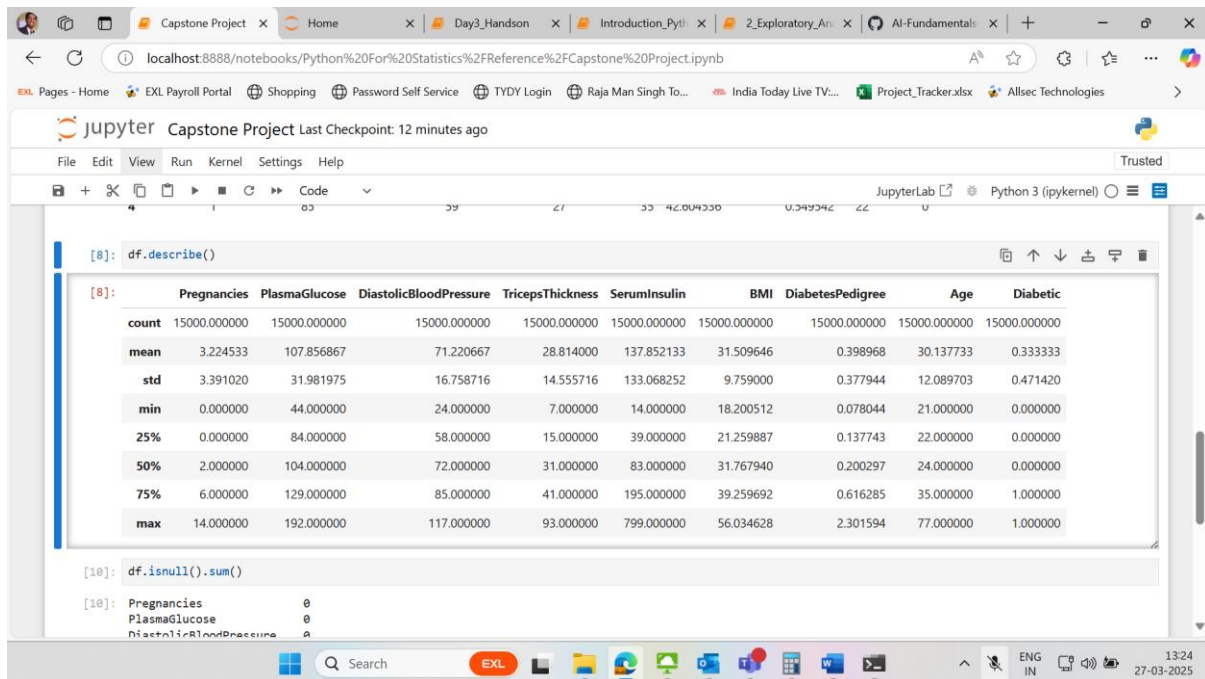# 1. Data Understanding & Cleaning

- **Load the dataset and examine its structure.**

- **Identify missing or inconsistent data and suggest appropriate handling techniques.**

- **Analyze the distribution of each variable using summary statistics.**

- **Discuss potential data transformations (e.g., normalization, scaling).**

**Solution:**

- ID REMOVED since it is not required

- Data is numerical

- Missing value handling:

    - No Null values found

- Since Measures of central tendencies are different, we **normalize** the data

## 2. Descriptive Statistics

Compute central tendency measures (mean, median, mode) for all numerical variables.

Calculate dispersion metrics (variance, standard deviation, interquartile range).

Create frequency distributions for categorical variables.

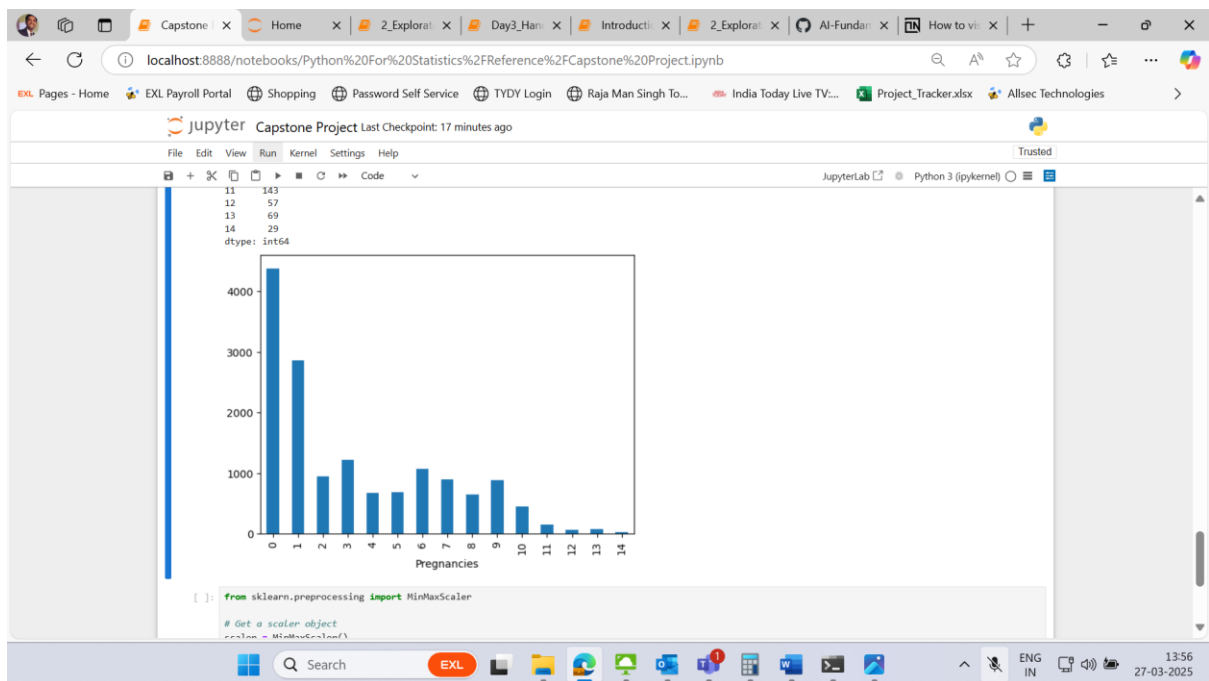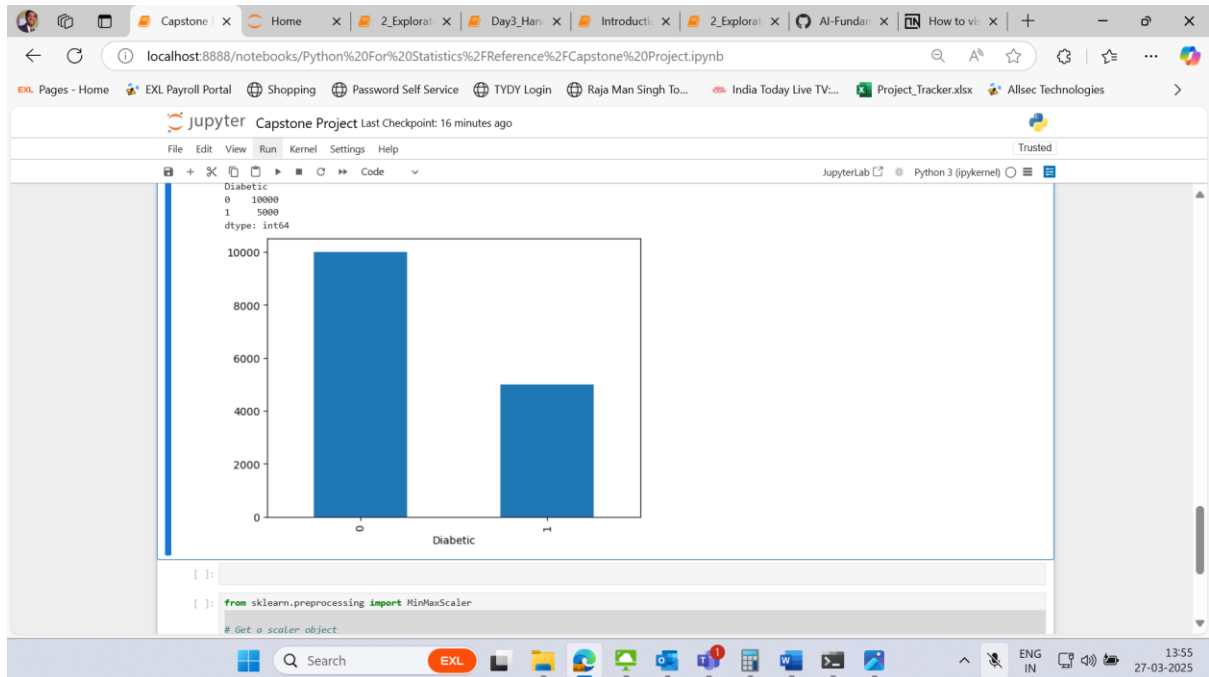Identify outliers using boxplots and discuss their potential impact.

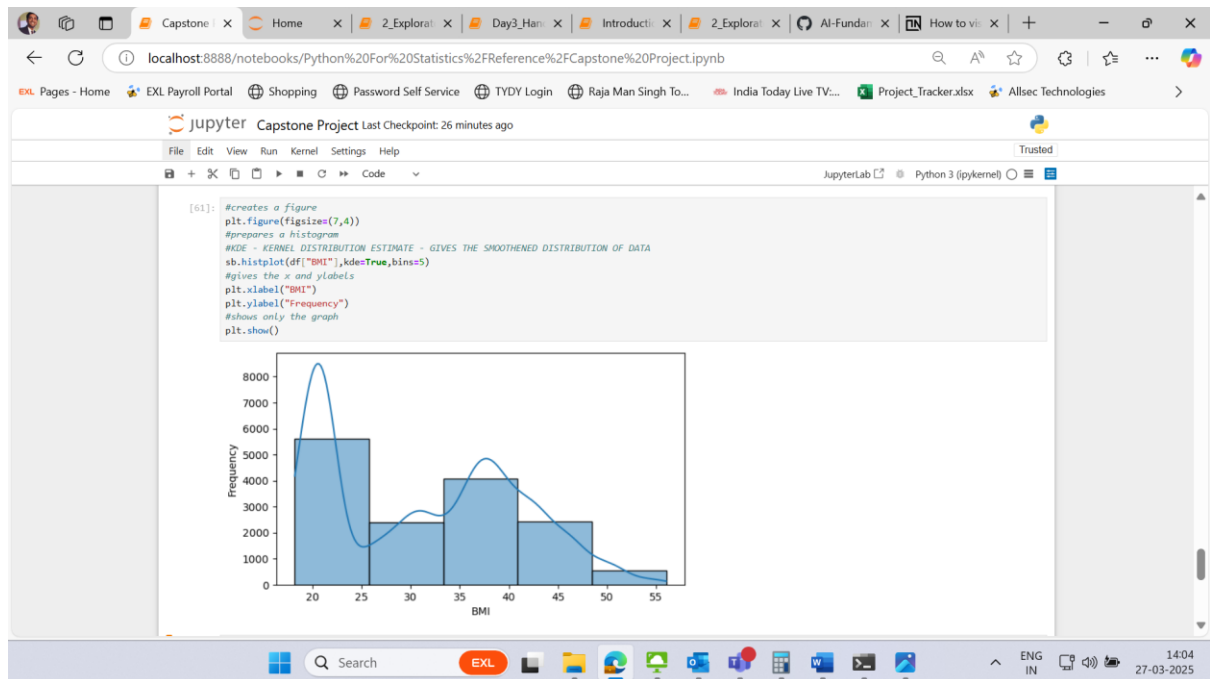Creating pregnencies and diabatic frequency distribution:

Outliers removed using box plot.
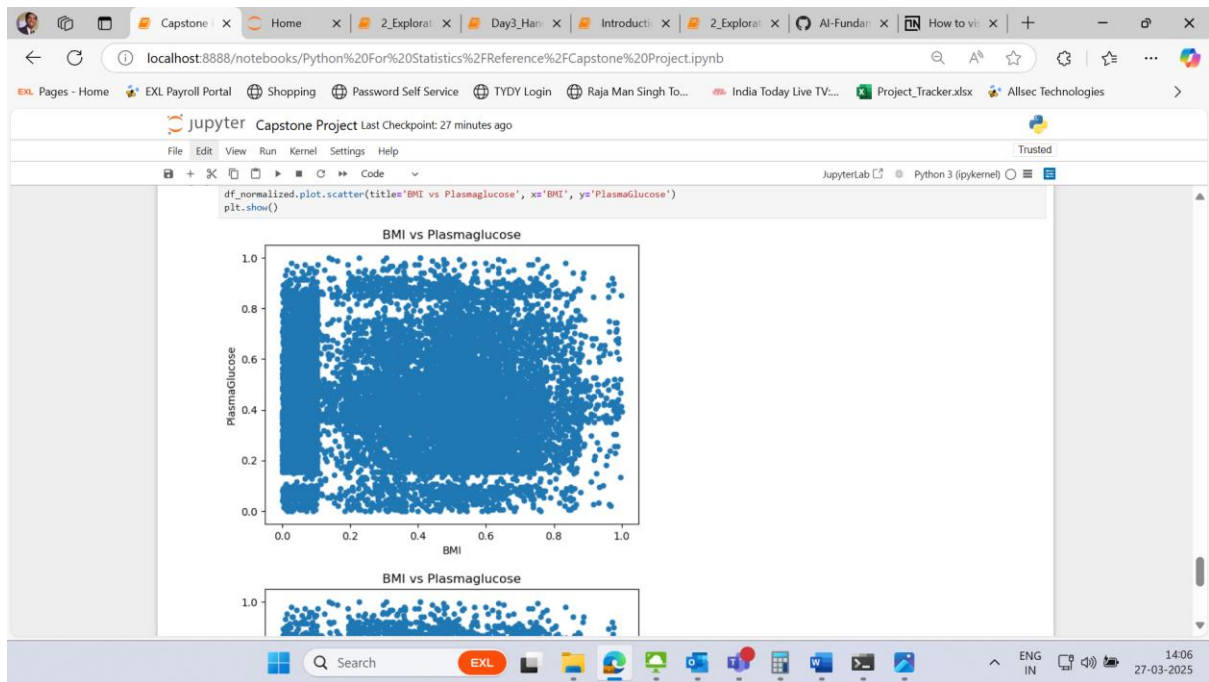
## 3. Exploratory Data Analysis (EDA)

- Generate visualizations (histograms, scatter plots, heatmaps) to explore data distribution and correlations.

- Examine relationships between features such as **BMI vs. PlasmaGlucose**, **Age vs. Diabetes Pedigree**, etc.

- Analyze the correlation matrix to determine the strength of relationships between variables.

- Perform segmentation analysis by dividing patients into age groups and evaluating diabetes risk.
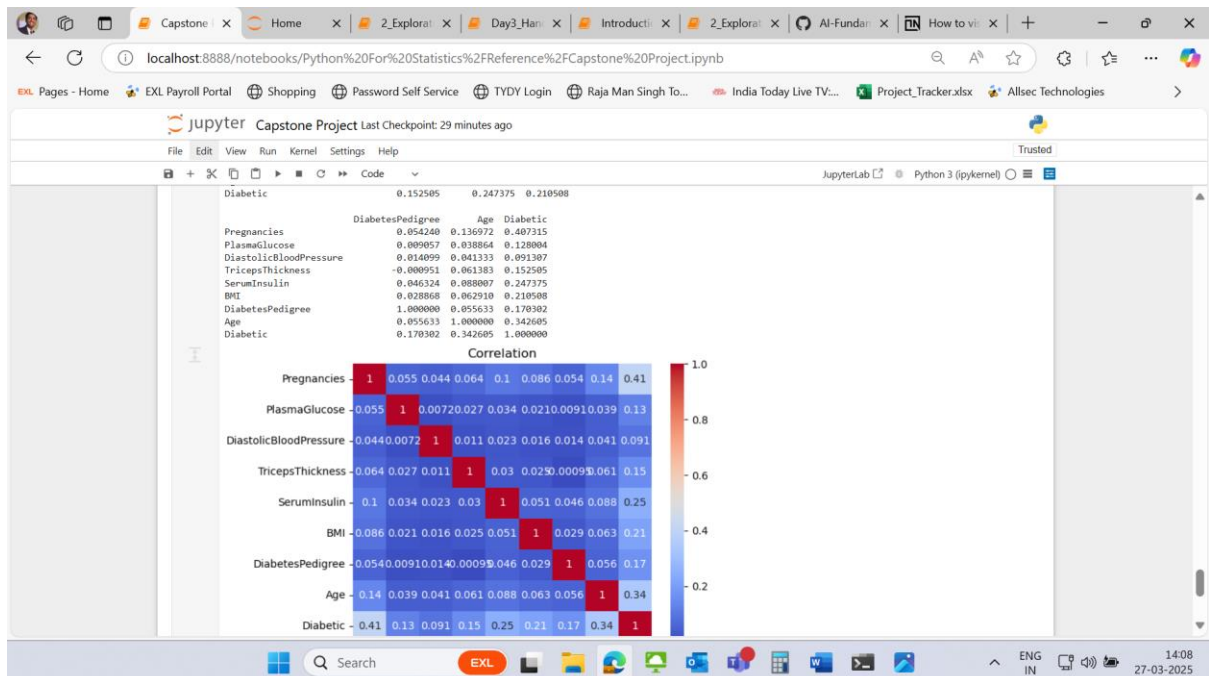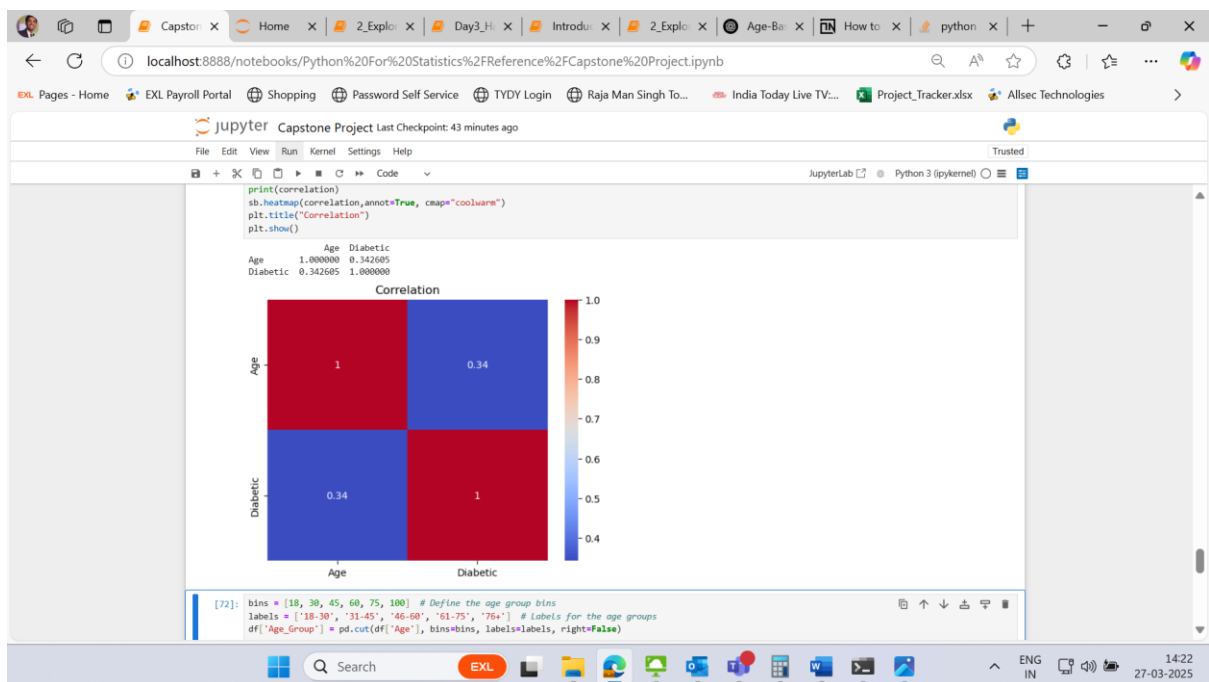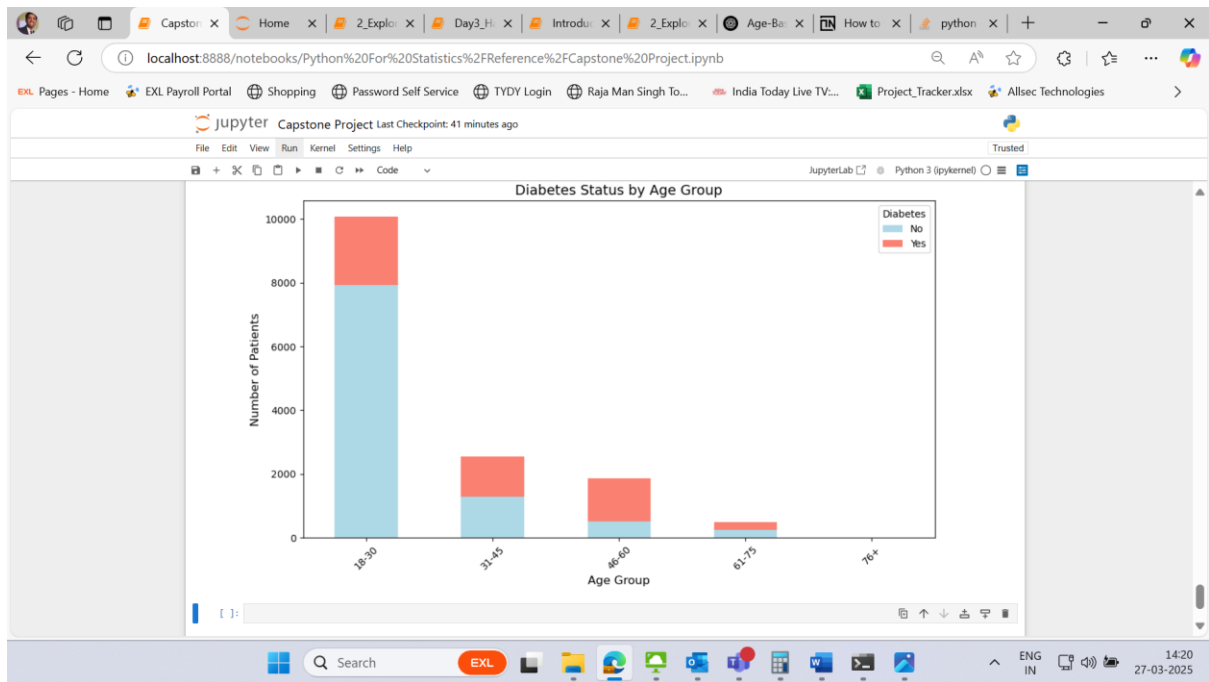
**Historgram:**



**BMI Vs Plasma Glucose – Scatterplot**

No clear relationship between the two.

**Correlation Matrix**:

## 4. Inferential Statistics

- Conduct hypothesis testing:
  - Compare the mean BMI of diabetic vs. non-diabetic patients (t-test).
  - Assess the relationship between **Pregnancies** and **Diabetes** using a chi-square test.

- Perform an ANOVA test to compare **PlasmaGlucose** levels across different age groups.

- Interpret p-values and confidence intervals to draw meaningful conclusions.



**Difference between BMI of Diabetic vs non diabetic patients is statistically significant.**

- Perform an ANOVA test to compare **PlasmaGlucose** levels across different age groups.

**PlasmaGlucose is statistically different across agegroups.**



**P Value is statistically significant. Hence, there is a relationship between pregnancies and diabetes.**

**Conclusions and Recommendations:**

Pregnencies, age and seruminsulin impact Diabetese positively. Hence, close attention should be given on these.