# Assignment 4 – Correlation Analysis

**Name:** Aman Nadeem
**Roll No:** 2225165002
**Course:** Applied Data Science with AI
**Week #: 4**
**Project Title:** Customer Churn Prediction

---

# 1. Reading Summary

**Reading Material:**

- Khan Academy – Statistics & Probability

- Introductory Stats for Data Science Notes

**Key Learnings:**

- Mean, median, and mode describe the center of data.

- Variance shows how spread out the data is.

- Correlation explains how two variables are related, either positively or negatively.

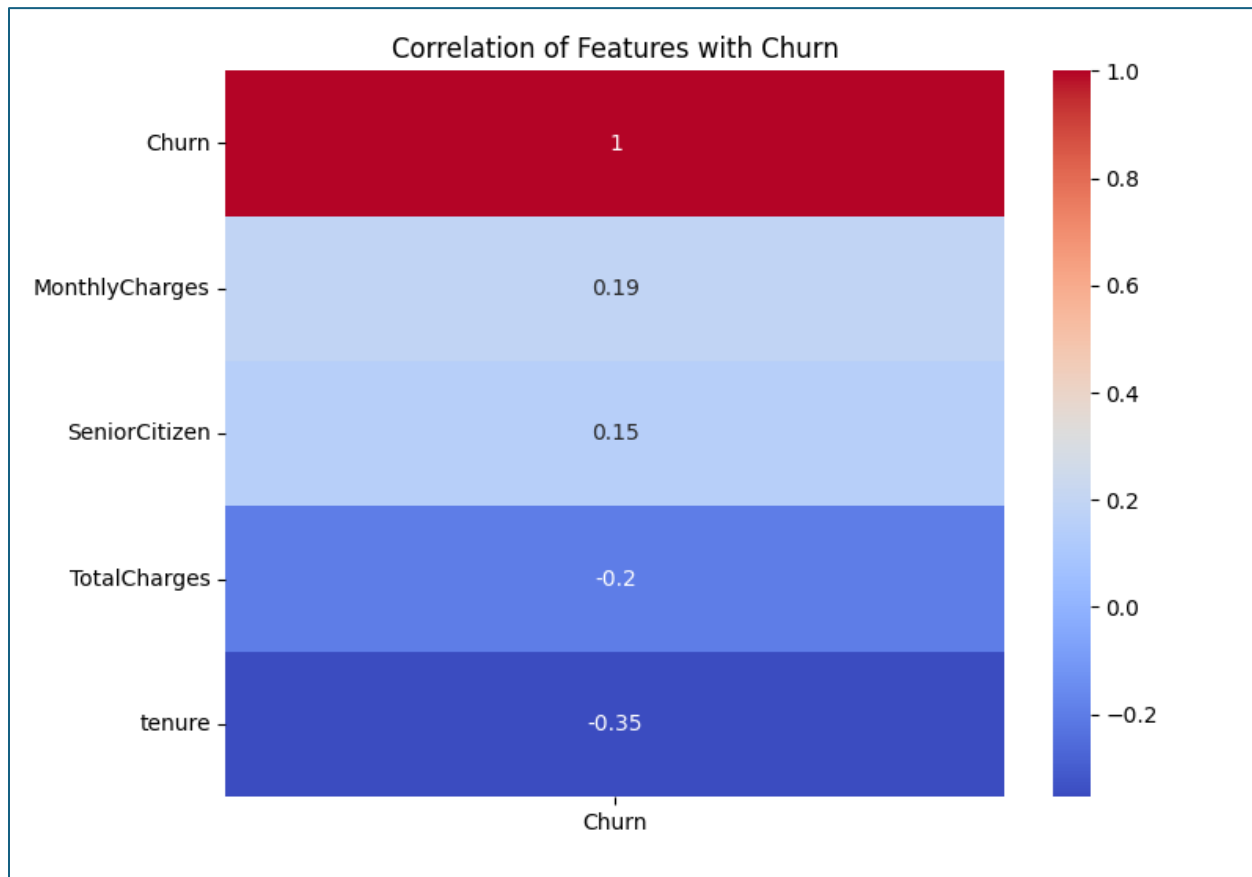- Understanding these measures helps to find which features matter most for predictions.

## Reflection:

This week's readings helped me understand how statistical measures guide us in identifying key factors that affect customer churn.

# 2. Classroom Task Documentation

## Task Performed:

- Practiced calculating mean, median, mode, variance, and correlation in datasets.

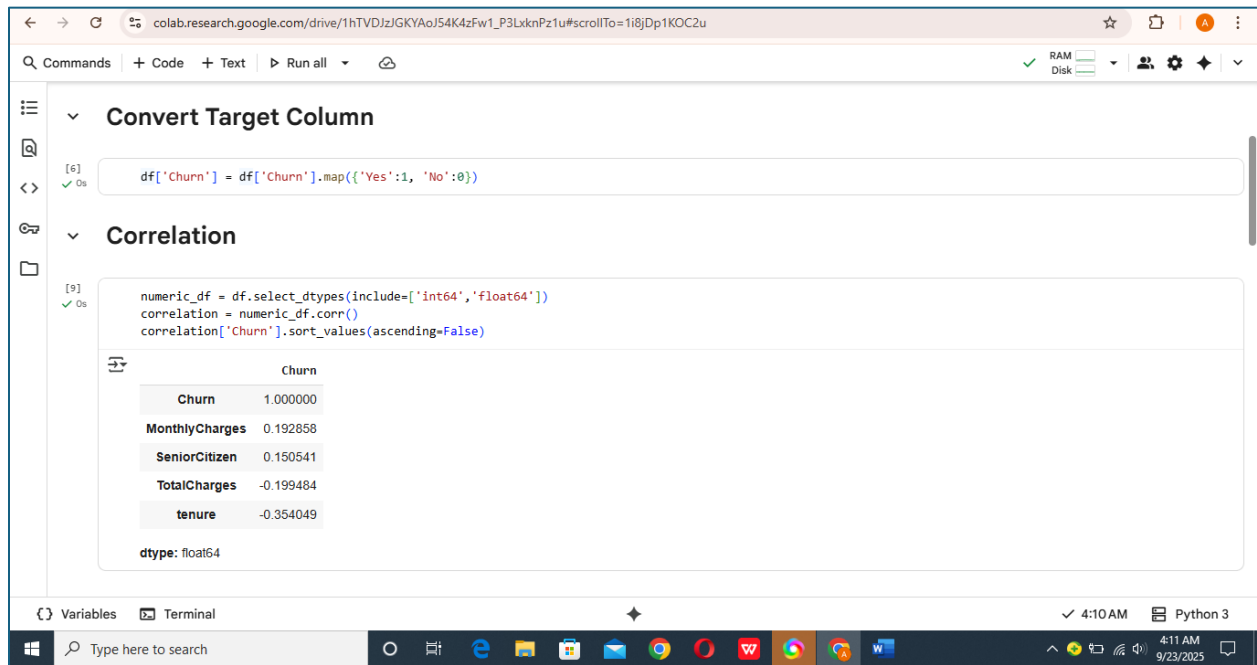- Learned how to build correlation heatmaps using Seaborn.



Correlation of Features with Churn

# 3. Weekly Assignment Submission

**Assignment Title:** Correlation Analysis

**Steps Taken:**

1. Loaded cleaned Churn dataset.

**2.** Converted Churn column to numeric (Yes = 1, No = 0).

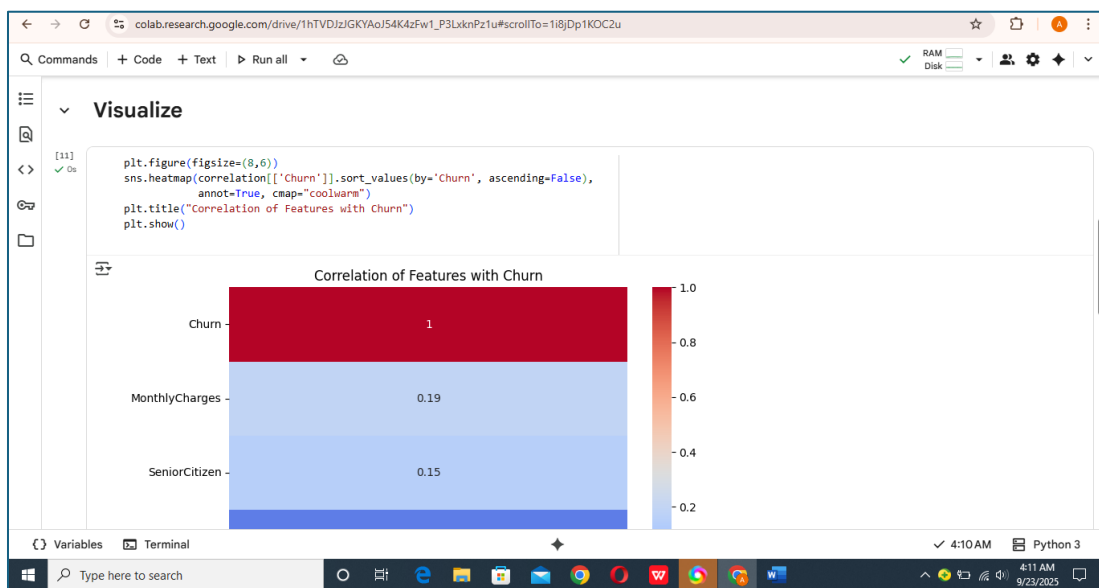**3.** Selected numeric features and computed correlation matrix.



**4.** Created heatmap to visualize correlations.

**5.** Calculated mean, median, mode, and variance for numeric columns.



## Statistical Summary:

- Tenure → Mean ≈ 32, Median ≈ 29, Mode = 1, Variance ≈ 600

- MonthlyCharges → Mean ≈ 64, Median ≈ 70, Mode = 20, Variance ≈ 900

## Output & Key Findings:

- **Tenure (-0.35):** Strong negative correlation with churn. Short-tenure customers are more likely to leave.

- **MonthlyCharges (+0.19):** Higher monthly bills increase the chance of churn.

- **TotalCharges (-0.20):** Customers who paid more overall are less likely to churn.

- **SeniorCitizen (+0.15):** Slightly higher churn tendency in senior citizens.

## Challenges Faced:

At first, I tried correlation on all columns, but it failed due to categorical data. Fixed it by selecting only numeric features.

## GitHub Link:

https://github.com/amannadeem126/Customer-Churn-Prediction

# 4. Project Progress Milestone

- Identified three most important features for churn: **Tenure, MonthlyCharges, TotalCharges**.

- Next week's goal: Build a baseline regression model.

# 5. Self-Evaluation

I completed all tasks on time.