# Course Project Report
## STA 2101: Statistics & Probability

## Project Title : Dhaka AQ

Student Name : Md. Shahriar Alam
Student ID : 242014180
University of Liberal Arts Bangladesh (ULAB)

October 14, 2025

### Abstract

This project analyzes the link between weather and air quality in Dhaka. It uses the "Dhaka Daily Air Quality and Weather" dataset. The study identifies key weather factors that impact the Air Quality Index (AQI). The analysis applies the statistical and probability concepts of STA 2101.

# Contents

# 1 Milestone 1: Dataset Selection

- **Dataset Name :** Dhaka Daily Air Quality and Weather Dataset

- **Dataset URL :** `https://www.kaggle.com/datasets/albab12/dhaka-daily-air-quality-and-weather-dataset`

- **Description :** This project uses the "Dhaka Daily Air Quality & Weather" dataset. The dataset provides daily records for Dhaka, Bangladesh. It contains two main types of information: air quality and weather. The air quality data includes the Air Quality Index (AQI). It also measures several pollutants. These pollutants include PM2.5, PM10, nitrogen dioxide, ozone, carbon monoxide, and sulfur dioxide. The weather data includes daily temperature. It also has information on humidity, barometric pressure, and wind speed.

  This dataset was chosen because it provides comprehensive variables for both air quality and weather. This makes it ideal for studying the relationship between these factors in Dhaka using Statistics & Probability concept.
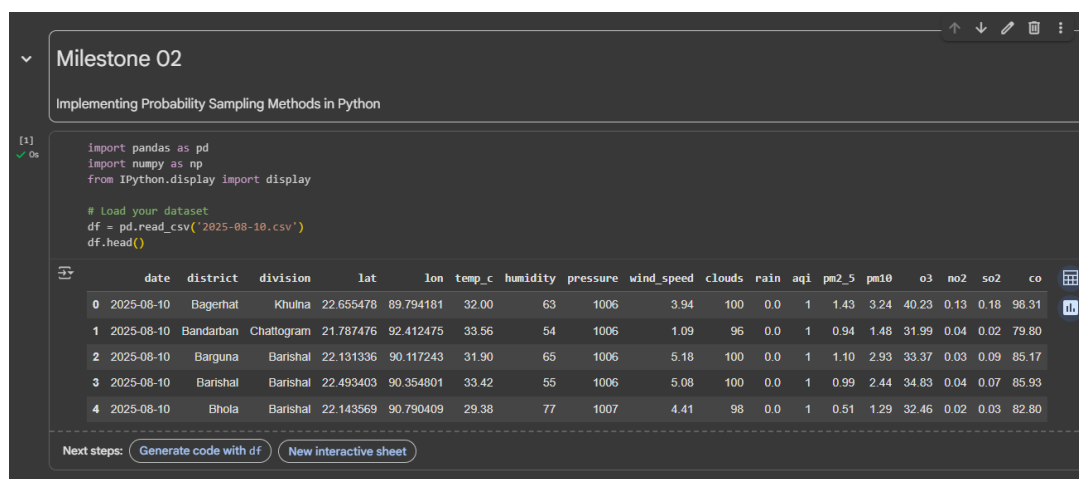
# 2 Milestone 02: Probability Sampling Methods



Figure 1: Overview

# Part A — Setup



Figure 2: Setup

# Part B — Simple Random Sampling



Figure 3: Simple Random Sampling

# Part C — Systematic Sampling



Figure 4: Systematic Sampling

# Part D — Stratified Sampling



Figure 5: Stratified Sampling

# Part E — Cluster Sampling



Figure 6: Cluster Sampling
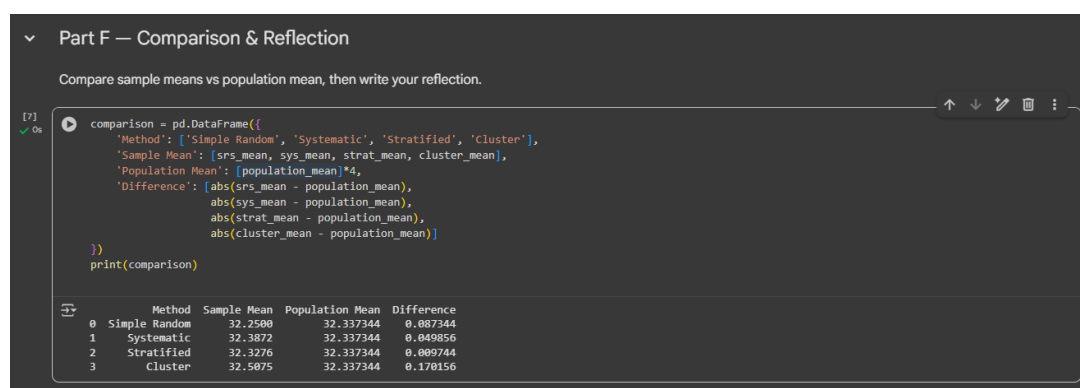
## Part F — Comparison & Reflection



Figure 7: Comparison and Reflection

In this milestone, I used four probability sampling methods. The population mean was 32.337344. All methods gave sample means close to it but not the same. Stratified sampling gave the closest mean 32.3276. The difference was very small. This happened because it kept the same group ratio as the dataset. Simple random sampling gave 32.25 which is a bit lower than the population mean. Systematic sampling gave 32.3872 which is slightly higher. Cluster sampling gave 32.5075 which is the farthest from the population mean.

Simple random sampling was the easiest. I wrote only one line of code and got the result. It did not need any group or pattern. Systematic sampling was also easy. I just needed k and a start point. Stratified sampling was harder. I had to use a division column and take samples from each group by proportion. Cluster sampling was easy to code but tricky to pick clusters.

Each method works for different goals. Simple random sampling is good for small or mixed data. Systematic sampling is good when data has no clear order. Stratified sampling is best for datasets with clear groups. Cluster sampling is useful for large data that is grouped by place or type. From this, I saw stratified sampling gave the most accurate result. Simple random sampling was the easiest to use.

# 3 Milestone 3: Data Visualization

Add graphs and figures using LaTeX. Example:

# 4 Milestone 4: Probability Distributions

Identify probability distributions in your dataset. Perform fitting, plots, and discuss results.

```
example-figure.png
```

Figure 8: Sample dataset visualization (replace with your figure)

# 5 Milestone 5: Hypothesis Testing

State hypotheses, perform tests, and report conclusions.

# 6 Milestone 6: Regression Analysis

Fit regression models, explain coefficients, and evaluate model fit.

# 7 Milestone 7–12: Further Analysis

Continue documenting each milestone here as instructed in class.

# 8 Final Conclusion

Summarize the overall findings of your project. Mention challenges, learning outcomes, and possible future work.

# References

List your references here in proper citation format. If you prefer, you may use BibTeX.