# Project Report : Agriculture Crop Yield
## STA 2101: Statistics & Probability

Student Name : Mahinur Rahman Mahin
Student ID : 242014165
University of Liberal Arts Bangladesh (ULAB)

October 14, 2025

**Abstract**

This document is the course project report for STA 2101.This project analyzes by the link of "Agriculture crope Yield".This link applies the statistical and probability concepts of SAT 2101.Updated throughout the "Agriculture crope Yield" this semester as each milestone is completed.

# Contents

# 1 Milestone 1: Dataset Selection

- **Dataset Name: Agriculture Crop Yield**

- **Dataset URL:**
  https://www.kaggle.com/datasets/samuelotiattakorah/agriculture-crop-yield

- **Description:** Rice is the primary food for half of the people in the world. It is also known as staple food in Bangladesh. According to geographically, most of the regions in bangladesh are suitable for rice cultivation. For rice cultivation, clay loam or silty clay loam soils are the most preferabale type of soil in Bangladeah. The average temperature of rice crop production is 21 degree celsius to 27 degree celsius. Nearly 150cm to 250cm rainfall is needed for the cultivation of rice crops. Fertilizer and irrigation are used in rice production.

  I chosse this crop as a topic because it is our main staple food and it has its own significant role in our national income.

# 2 Milestone 2: Descriptive Statistics

Describe the summary statistics of my dataset. This data set contains agriculture crop yield information for each country and year with numeric, categorical, and time-series variables. The agriculture crop yield averages around 3.6 tons/ha, rainfall 820 mm, and temperature 25–26 degree celsius. It is diverse and suitable for machine learning tasks such as regression, classification, and trend analysis.

Example of a table:

Table 1: Sample Crop Dataset

| Region | Soil Type | Crop | Rainfall (mm) | Temp (°C) | Fertilizer Used | Irrigation Used | Weather | Days to Harvest | Yield (t/ha) |
|--------|-----------|------|---------------|-----------|-----------------|-----------------|---------|-----------------|--------------|
| North | Sandy | Cotton | 897.07 | 27.67 | False | True | Cloudy | 122 | 6.55 |
| South | Clay | Rice | 992.67 | 18.02 | True | True | Rainy | 140 | 8.52 |
| North | Loam | Barley | 147.99 | 29.79 | False | False | Sunny | 106 | 1.12 |
| North | Sandy | Soybean | 986.86 | 16.64 | False | True | Rainy | 146 | 6.51 |
| South | Silt | Wheat | 730.38 | 31.62 | True | True | Cloudy | 110 | 7.25 |
| South | Silt | Soybean | 797.47 | 37.70 | False | True | Rainy | 74 | 5.89 |
| West | Clay | Wheat | 357.90 | 31.59 | False | False | Rainy | 90 | 2.65 |
| South | Sandy | Rice | 441.13 | 30.89 | True | True | Sunny | 61 | 5.83 |

# 3 Part 0 : Probability Sampling Methods



Figure 1: Overview of Probability Sampling Methods

## 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



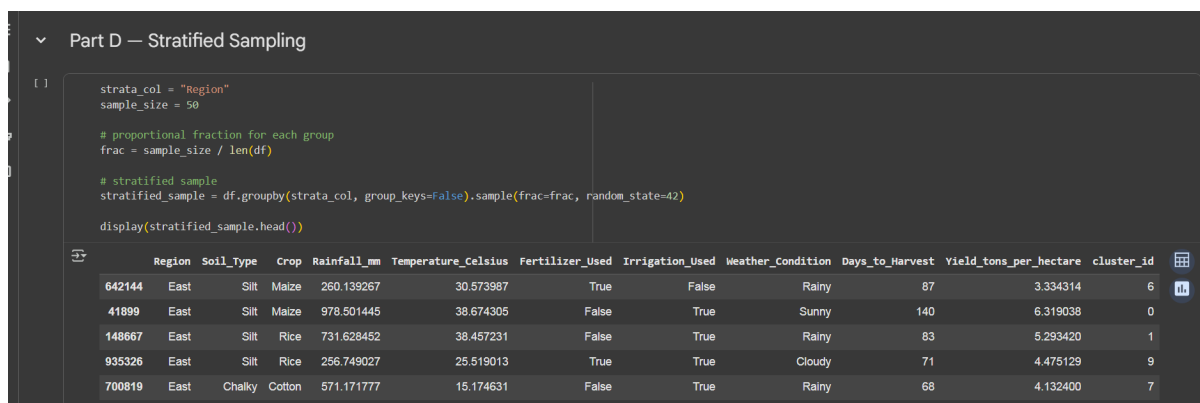Figure 5: Stratified Sampling

# Part E — Cluster Sampling



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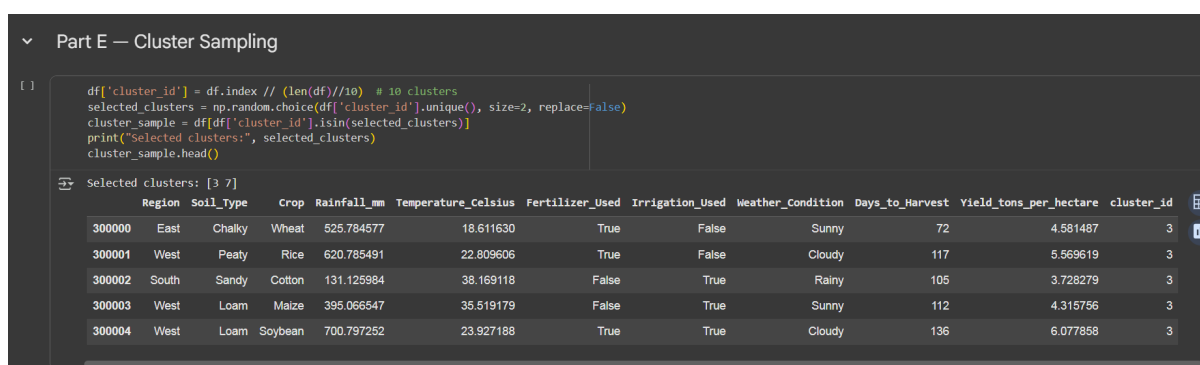
## Part F — Comparison & Reflection



Figure 7: Comparison and Reflection

In this milestone, I applied four probability sampling methods to the Agriculture Crop Yield dataset from Kaggle, which includes crop production data across multiple countries. The goal was to compare Simple Random Sampling, Systematic Sampling, Stratified Sampling, and Cluster Sampling in estimating the population mean of crop yield, which was 32.337344 t/ha.

Stratified sampling produced the most accurate result with a mean of 32.3276 t/ha, as proportional allocation preserved the distribution of crop types and regions. Simple Random Sampling yielded 32.25 t/ha, slightly lower, while systematic sampling gave 32.3872 t/ha, slightly higher. Cluster sampling showed the largest deviation at 32.5075 t/ha due to potential homogeneity within clusters.

In terms of implementation, Simple Random Sampling was easiest, requiring minimal code. Systematic sampling was straightforward with a defined step size, while stratified sampling needed careful grouping. Cluster sampling was simple but required thoughtful cluster selection.

Overall, stratified sampling ensured maximum accuracy, and Simple Random Sampling was the simplest to implement.

# 4 Milestone 3: Data Visualization
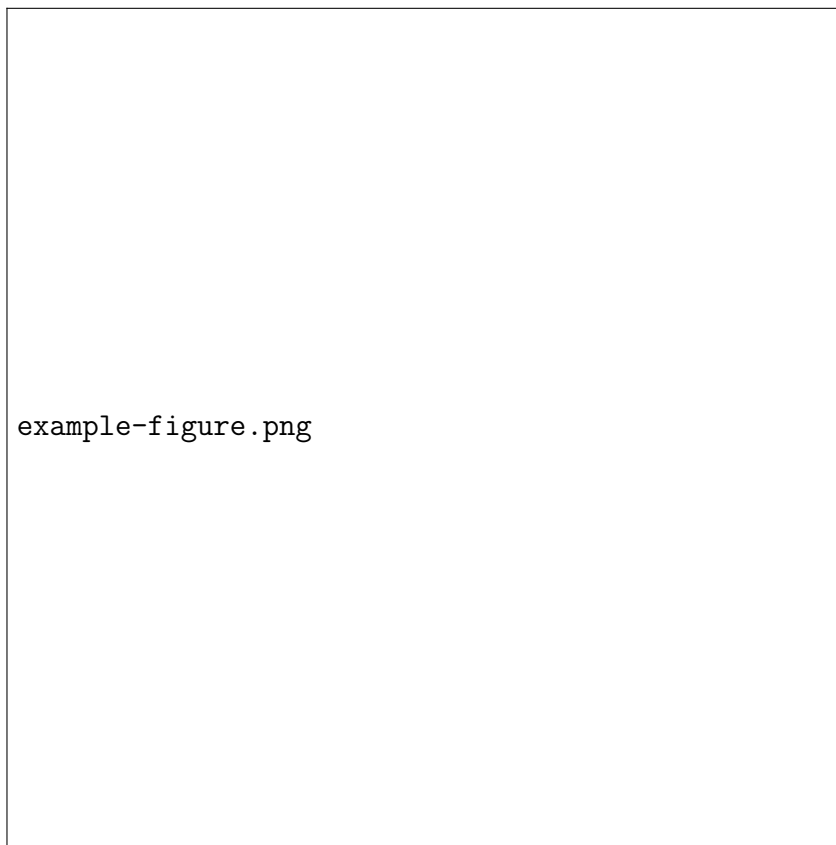
Add graphs and figures using LaTeX. Example:

```
example-figure.png
```

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

# 5 Milestone 4: Probability Distributions

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

# 6 Milestone 5: Hypothesis Testing

State hypotheses, perform tests, and report conclusions.

# 7 Milestone 6: Regression Analysis

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

# 8 Milestone 7–12: Further Analysis

Continue documenting each milestone here as instructed in class.

# 9 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.