# Cluster Sampling – Population Mean Estimation

**Subject:**

Sampling Techniques

**Submitted By:**

Ann Maria Anil

2341419

6BScDS

**Submitted To:**

Prof.Dibu

**Date of Submission:**

16th December, 2025

**1. Introduction**

Cluster sampling is a probability-based survey technique where the population is divided into mutually exclusive groups called clusters. Instead of sampling individual units from the entire population, a subset of clusters is randomly selected, and all units within these clusters are included in the survey. This approach is especially useful when the population is large, geographically dispersed, or when surveying every unit individually is costly or time-consuming.

This report demonstrates a Shiny-based application for designing a cluster sampling survey, estimating population mean, and visualizing the selected clusters.

**2. Objective**

The main objectives of this cluster sampling exercise are:

- To determine the **number of clusters** to include in the survey based on user-defined parameters such as allowable error and confidence level.

- To calculate **expected cost and time** for surveying the selected clusters.

- To provide a **visual representation** of the survey design for better understanding and planning.

- To create an **interactive tool** that allows users to experiment with different survey parameters and immediately see the effect on cluster selection.

**3. Methodology**

1. **Population and Clusters Setup:**

   o The total population is divided into a chosen number of clusters.

   o Each cluster contains a number of units, and each unit has an associated observed value (simulated in the app for demonstration purposes).

   o Cost and time per cluster are defined to allow estimation of survey resources.

2. **Cluster Summary:**

   o For each cluster, the app computes:

     ▪ Cluster size (number of units)

     ▪ Mean of observed values

     ▪ Variance of observed values

     ▪ Cost and time

3. **Cluster Selection:**

   o Based on the allowable error (sampling bias) and confidence level, the app determines how many clusters should be selected.

   o Clusters are selected randomly to ensure unbiased sampling.

4. **Cost & Time Calculation:**

   o Total expected cost and survey time are calculated by summing the cost and time across the selected clusters.

5. **Visualization:**

   o The design of experiment is plotted using a scatter plot faceted by cluster.

   o Units within selected clusters are highlighted in a distinctive color, while non-selected clusters are shown in a neutral color.

   o This visual representation helps quickly understand which clusters are included and the overall survey coverage.

## 4. Features of the Shiny App

```
1  library(shiny)
2  library(ggplot2)
3  library(dplyr)
4  library(DT)
5
6  ui <- fluidPage(
7
8    titlePanel("Cluster Sampling - Population Mean Estimation"),
9
10   sidebarLayout(
11
12     sidebarPanel(
13       numericInput("N", "Population Size (N):", value = 300, min = 50),
14       numericInput("C", "Total Number of Clusters:", value = 10, min = 2),
15       numericInput("B", "Allowable Error / Sampling Bias (B):", value = 4, min = 1),
16
17       selectInput("Z", "Confidence Level:",
18                   choices = c("90%" = 1.645,
19                               "95%" = 1.96,
20                               "99%" = 2.576),
21                   selected = 1.96),
22
23       numericInput("cost", "Cost per Cluster:", value = 100, min = 1),
24       numericInput("time", "Time per Cluster (hours):", value = 2, min = 0.1),
25
26       actionButton("calc", "Generate Cluster Design")
27     ),
28
```

```r
29        mainPanel(
30          tabsetPanel(
31
32            tabPanel("Cluster Information",
33                    DTOutput("cluster_table")
34            ),
35
36            tabPanel("Sampling Summary",
37                    tableOutput("summary")
38            ),
39
40            tabPanel("Design of Experiment",
41                    plotOutput("cluster_plot", height = "450px")
42            )
43          )
44        )
45      )
46  )
47
48  server <- function(input, output) {
49
50    observeEvent(input$calc, {
51
52      set.seed(123)
53
54      # POPULATION
55      units_per_cluster <- ceiling(input$N / input$C)
56
57      population <- data.frame(
58        Unit = 1:input$N,
```

```r
59        y = rnorm(input$N, mean = 60, sd = 15),
60        Cluster = rep(paste0("C", 1:input$C), each = units_per_cluster)[1:input$N],
61        Cost = rep(input$cost, input$N),
62        Time = rep(input$time, input$N)
63      )
64
65      # Cluster-level summaries
66      cluster_info <- population %>%
67        group_by(Cluster) %>%
68        summarise(
69          Cluster_Size = n(),
70          Mean = mean(y),
71          Variance = var(y),
72          Cost = first(Cost),
73          Time = first(Time),
74          .groups = "drop"
75        )
76
77      M <- nrow(cluster_info)          # total clusters
78      S2 <- var(cluster_info$Mean)     # between-cluster variance
79      Z <- as.numeric(input$Z)
80      B <- input$B
81
82      # REQUIRED NUMBER OF CLUSTERS
83      m <- ceiling((Z^2 * S2) / (B^2))
84      m <- min(m, M)
85
86      # Select clusters (SRS of clusters)
87      selected_clusters <- sample(cluster_info$Cluster, m)
88
```

```r
 89        population$Selected <- ifelse(
 90          population$Cluster %in% selected_clusters,
 91          "Selected", "Not Selected"
 92        )
 93
 94        # COST & TIME CALCULATION
 95        cost_time <- cluster_info %>%
 96          filter(Cluster %in% selected_clusters) %>%
 97          summarise(
 98            Total_Clusters_Selected = m,
 99            Expected_Cost = sum(Cost),
100            Expected_Time = sum(Time)
101          )
102
103        # OUTPUTS
104 -      output$cluster_table <- renderDT({
105          datatable(cluster_info, options = list(pageLength = 10))
106 ^      })
107
108 -      output$summary <- renderTable({
109          data.frame(
110            Parameter = c("Total Clusters",
111                          "Clusters Selected",
112                          "Allowable Error (B)",
113                          "Confidence Level (Z)",
114                          "Expected Cost",
115                          "Expected Time"),
116            Value = c(M,
117                      m,
118                      B,
119                      Z
```

```r
125 -      output$cluster_plot <- renderPlot({
126          ggplot(population,
127                 aes(x = Unit, y = y, color = Selected)) +
128            geom_point(size = 2) +
129            facet_wrap(~Cluster, scales = "free_x") +
130            scale_color_manual(
131              values = c("Selected" = "#1F7884",
132                         "Not Selected" = "#D9D9D9")
133            ) +
134            labs(
135              title = "Cluster Sampling Design",
136              subtitle = paste("Selected Clusters:", paste(selected_clusters, collapse = ", ")),
137              x = "Population Units",
138              y = "Observation Value",
139              color = "Cluster Status"
140            ) +
141            theme_minimal(base_size = 14) +
142            theme(
143              plot.title = element_text(face = "bold", size = 16),
144              plot.subtitle = element_text(size = 12),
145              legend.position = "bottom",
146              legend.title = element_text(face = "bold")
147            )
148 -      })
149
150 -  })
151
152 - }
153
154  shinyApp(ui, server)
```
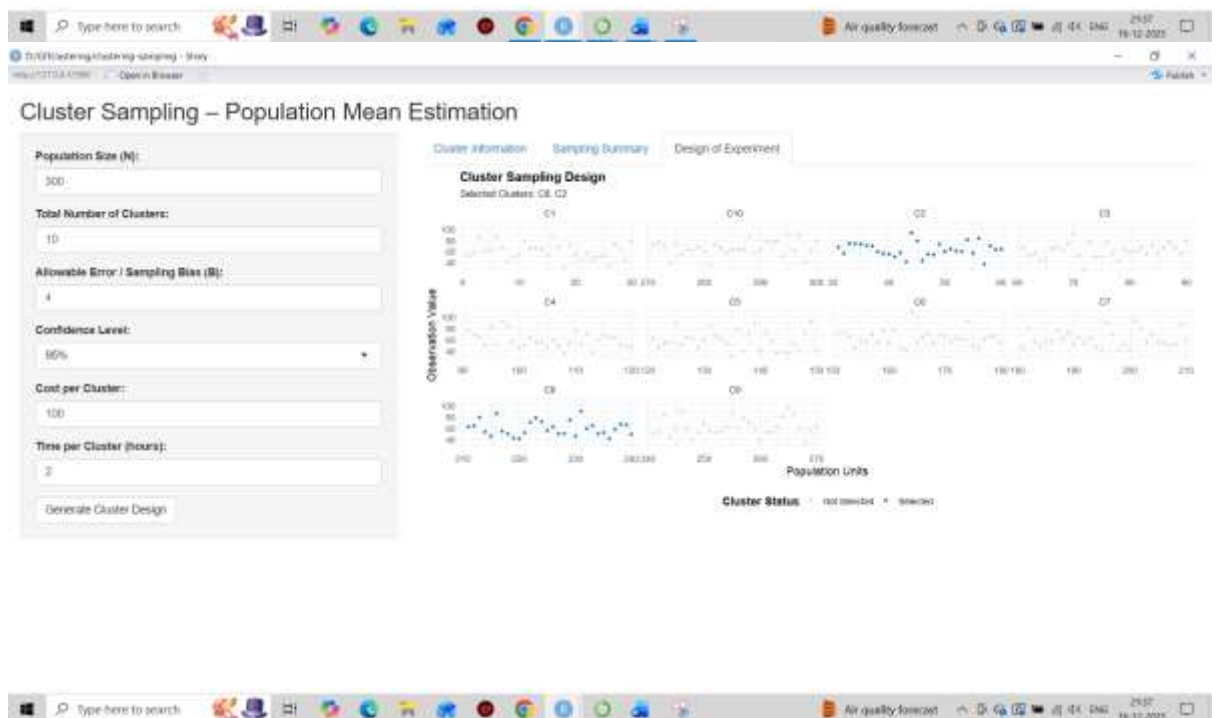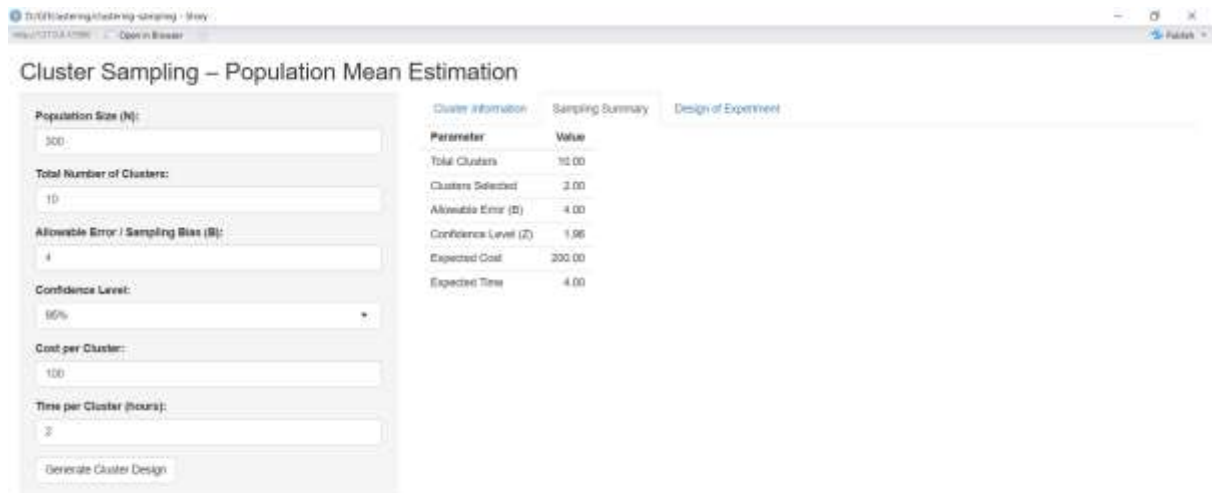
- **Interactive Inputs:** Allows users to define population size, number of clusters, allowable error, confidence level, cost, and time per cluster.

- **Cluster Information Tab:** Displays detailed summary for each cluster in an interactive table.

- **Sampling Summary Tab:** Provides overall information including total clusters, clusters selected, allowable error, confidence level, expected cost, and time.

- **Design of Experiment Plot:** Shows a clear visualization of selected vs. non-selected clusters, facilitating planning and communication.

- **Reactivity:** Changing any input updates all outputs immediately, allowing users to experiment and observe effects dynamically.

**5.Output screenshots:**

Cluster Sampling – Population Mean Estimation



Cluster Sampling – Population Mean Estimation

## 6. Conclusion

- The Shiny app provides an **efficient and interactive tool** for designing cluster-based surveys.

- It allows users to **estimate required clusters**, **calculate expected resources**, and **visualize survey design** quickly.

- The approach ensures **efficient survey planning**, minimizing cost and time while maintaining statistical reliability.

- The visualization component makes it easier for survey designers and stakeholders to interpret and communicate the survey plan.

**7.Professional github link**

Including the git repository link : https://github.com/ann-maria-anil/clustering-sampling-r-shiny