

Mini Project 2

PSTAT100: Data Science Concepts and Analysis

Instructor: Ali Abuzaid

STUDENT NAME

- Phuc Lu (pdlu)
- Valerie De La Fuente (valeriedelafuente)

Instructions

- This mini project is designed to give you practical experience with real-world data using R and Shiny. You'll create an interactive web application that allows users to explore and visualize a dataset.
- Work in groups of **2 students** from the same discussion section.
- Individual submissions will not be accepted.
- Please use the provided MP 2.qmd file to type your Documentation and Presentation and submit it as a PDF file. You can utilize **RStudio** for this purpose. For guidance, refer to the [Tutorial: Hello, Quarto](#)).
- Please submit a .zip file that includes:

Your **app.R** file (fully working Shiny app).

A short project report (PDF).

Reminder: If your app fails to open or the .zip is incorrect, you will receive a score of **ZERO**. Test everything before submission.

Due Date

Due Date: Sunday, June 1, 2025, 11:59 PM

0.1 2 Tasks:

0.1.1 2.1.1: Data Loading:

```
1 library(tidyverse)
2 data(infert)
```

Variable	Description	Value
education	years of education	0–5 years, 6–11 years, 12+ years
age	age in years	22-24
parity	number of prior live births	1-6
induced	number of prior induced abortions	0, 1, 2 or more
case	case status	case, control
spontaneous	number of prior spontaneous abortions	0, 1, 2 or more
stratum	id for each matched set or group	1–83
pooled.stratum	represent a group or broader classification of strata (may be used for conditional logistic regression)	1 - 63

0.1.2 2.1.2: Data Preparation:

```
1 infert %>% head()
```

	education	age	parity	induced	case	spontaneous	stratum	pooled.stratum
1	0-5yrs	26	6	1	1	2	1	3
2	0-5yrs	42	1	1	1	0	2	1
3	0-5yrs	39	6	2	1	0	3	4
4	0-5yrs	34	4	2	1	0	4	2
5	6-11yrs	35	3	1	1	1	5	32
6	6-11yrs	36	4	2	1	1	6	36

```
1 library(stringr)
2 # rename variables and remove "yrs" from education values. factorize necessary variables
3 infert_ds <- infert %>%
```

```

4  rename(
5    Education = education,
6    Age = age,
7    Parity = parity,
8    InducedAbortions = induced,
9    SpontaneousAbortions = spontaneous,
10   CaseStatus = case,
11   StratumID = stratum,
12   PooledStratumID = pooled.stratum
13 ) %>%
14 mutate(
15   Education = factor(str_remove(Education, "yrs"), levels = c("0-5", "6-11", "12+")),
16   CaseStatus = factor(CaseStatus, levels = c(0, 1), labels = c("control", "case")),
17   Parity = factor(Parity),
18   InducedAbortions = factor(InducedAbortions),
19   SpontaneousAbortions = factor(SpontaneousAbortions),
20   StratumID = factor(StratumID),
21   PooledStratumID = factor(PooledStratumID)
22 )
23
24 factor_vars <- names(infert_ds)[sapply(infert_ds, is.factor)]
25 numeric_vars <- names(infert_ds)[sapply(infert_ds, is.numeric)]

```

```

1  summary(infert_ds)

```

Education	Age	Parity	InducedAbortions	CaseStatus
0-5 : 12	Min. :21.00	1:99	0:143	control:165
6-11:120	1st Qu.:28.00	2:81	1: 68	case : 83
12+ : 0	Median :31.00	3:36	2: 37	
NA's:116	Mean :31.50	4:18		
	3rd Qu.:35.25	5: 6		
	Max. :44.00	6: 8		

SpontaneousAbortions	StratumID	PooledStratumID
0:141	1 : 3	41 : 12
1: 71	2 : 3	45 : 9
2: 36	3 : 3	49 : 9
	4 : 3	51 : 9
	5 : 3	12 : 6
	6 : 3	18 : 6
	(Other):230	(Other):197

```

1 str(infert_ds)

'data.frame':  248 obs. of  8 variables:
 $ Education      : Factor w/ 3 levels "0-5","6-11","12+": 1 1 1 1 2 2 2 2 2 2 ...
 $ Age            : num  26 42 39 34 35 36 23 32 21 28 ...
 $ Parity         : Factor w/ 6 levels "1","2","3","4",...: 6 1 6 4 3 4 1 2 1 2 ...
 $ InducedAbortions : Factor w/ 3 levels "0","1","2": 2 2 3 3 2 3 1 1 1 1 ...
 $ CaseStatus     : Factor w/ 2 levels "control","case": 2 2 2 2 2 2 2 2 2 2 ...
 $ SpontaneousAbortions: Factor w/ 3 levels "0","1","2": 3 1 1 1 2 2 1 1 2 1 ...
 $ StratumID      : Factor w/ 83 levels "1","2","3","4",...: 1 2 3 4 5 6 7 8 9 10 ...
 $ PooledStratumID : Factor w/ 63 levels "1","2","3","4",...: 3 1 4 2 32 36 6 22 5 19 ...

```

0.1.3 2.2 Shiny App Development

User Interface (UI):

```

1 # Define UI
2 ui <- fluidPage(
3   titlePanel("Infertility Dataset Explorer"),
4
5   sidebarLayout(
6     sidebarPanel(
7       selectInput("variable", "Select Variable:",
8         choices = names(infert_ds)),
9
10      uiOutput("plotTypeUI"),
11
12      conditionalPanel(
13        condition = "input.plotType == 'Histogram'",
14        sliderInput("bins", "Number of Bins:", min = 5, max = 50, value = 10),
15        checkboxInput("showDensity", "Overlay Density Curve", value = FALSE)
16      ),
17
18      conditionalPanel(
19        condition = "input.plotType == 'Boxplot'",
20        selectInput("boxplotGroup", "View Boxplot Across:",
21          choices = factor_vars),
22        checkboxInput("showJitter", "Add Jittered Data Points", value = FALSE)
23      ),
24
25      checkboxInput("includeNA", "Include NA in Grouped Summary", value = FALSE),

```

```

26
27     h4("Numeric Filters"),
28     sliderInput("ageFilter", "Age Range:",
29                 min(infert_ds$Age), max(infert_ds$Age),
30                 value = c(min(infert_ds$Age), max(infert_ds$Age))),
31
32     selectInput("groupBy", "Group Summary By:",
33                 choices = factor_vars)
34 ),
35
36 mainPanel(
37     plotOutput("mainPlot", hover = "plot_hover"),
38     verbatimTextOutput("hoverInfo"),
39     h4("Summary"),
40     verbatimTextOutput("summaryOutput"),
41     h4("Grouped Summary"),
42     tableOutput("groupedSummary")
43 )
44 )
45 )

```

Server Logic:

```

1  # Define server logic
2  server <- function(input, output, session) {
3
4      # Dynamically adjust plot options
5      output$plotTypeUI <- renderUI({
6          req(input$variable)
7          var <- input$variable
8
9          if (is.numeric(infert_ds[[var]])) {
10             selectInput("plotType", "Select Plot Type:", choices = c("Histogram", "Boxplot"))
11         } else {
12             selectInput("plotType", "Select Plot Type:", choices = c("Bar Chart"))
13         }
14     })
15
16     # Reactive filtered dataset
17     filtered_data <- reactive({
18         infert_ds %>%
19         filter(

```

```

20     Age >= input$ageFilter[1], Age <= input$ageFilter[2]
21   )
22 })
23
24 # Generate plots
25 output$mainPlot <- renderPlot({
26   req(input$variable, input$plotType)
27   df <- filtered_data()
28   var <- input$variable
29
30   if (input$plotType == "Histogram") {
31     p <- ggplot(df, aes_string(x = var)) +
32       geom_histogram(bins = input$bins, fill = "lightpink", color = "black") +
33       labs(title = paste("Histogram of", var), x = var, y = "Count") +
34       theme_minimal()
35
36     if (input$showDensity) {
37       p <- p + geom_density(aes(y = ..count..), color = "blue", size = 1)
38     }
39
40     p
41
42   } else if (input$plotType == "Boxplot") {
43     group_var <- input$boxplotGroup
44
45     p <- ggplot(df, aes_string(x = group_var, y = var)) +
46       geom_boxplot(fill = "royalblue") +
47       labs(title = paste("Boxplot of", var, "across", group_var), x = group_var, y = var) +
48       theme_minimal()
49
50     if (input$showJitter) {
51       p <- p + geom_jitter(width = 0.2, alpha = 0.4)
52     }
53
54     p
55
56   } else if (input$plotType == "Bar Chart") {
57     ggplot(df, aes_string(x = var)) +
58       geom_bar(fill = "lavender", color = "black") +
59       labs(title = paste("Bar Chart of", var), x = var, y = "Count") +
60       theme_minimal()
61   }

```

```

62 })
63
64 # Show hover info
65 output$hoverInfo <- renderPrint({
66   hover <- input$plot_hover
67
68   if (input$plotType == "Histogram" || input$plotType == "Boxplot") {
69     if (!is.null(hover)) {
70       cat("Hovered at:\n")
71       cat("x:", round(hover$x, 2), "\n")
72       cat("y:", round(hover$y, 2), "\n")
73     } else {
74       cat("Hover over the plot to see coordinates.")
75     }
76   } else {
77     cat("Hover info disabled for bar charts (categorical x-axis).")
78   }
79 })
80
81 # Show summary statistics
82 output$summaryOutput <- renderPrint({
83   req(input$variable)
84   summary(filtered_data()[[input$variable]])
85 })
86
87 # Grouped summary
88 output$groupedSummary <- renderTable({
89   df <- filtered_data()
90   group_var <- input$groupBy
91
92   if (!input$includeNA) {
93     df <- df %>% drop_na(all_of(group_var))
94   }
95
96   df %>%
97     group_by(across(all_of(group_var))) %>%
98     summarise(across(where(is.numeric), mean, na.rm = TRUE), .groups = "drop")
99   }, rownames = TRUE)
100 }
101
102 # Run the app
103 shinyApp(ui = ui, server = server)

```

1 Documentation and Presentation

1.1 App Purpose

The purpose of this app is make it easy for the users to understand the infertility data set without having to do coding on their own. All it takes is point, click, and drag to get beautiful graphs or summary statistics of the data.

1.2 How it Works

The select variable feature allow the users to select one of eight variables from the `infert` data set to analyze. The app will automatically select a graph that works with the variable of interest. The user can also adjust the number of bins there are in the graph. If a boxplot is chosen, the user can select another variable to compare with. The user can also add jitters to the boxplot, as well include missing (NA) values in the boxplot. The user can also group the summary by a certain variable as well. The grouped summary will be displayed at the very bottom of the page. On the graphs, the user can hover their mouse over the graph to view the coordinates (x, y) for more precise information. Below the graph feature are summary statistics. These include the minimum, 1st quartile, median, mean, 3rd quartile, and maximum value.

1.3 Insights

The initial questions that this data set tries to answer is what factors are associated with secondary infertility in women. Two of the main components to this questions are spontaneous versus induced abortions. Spontaneous abortions is abortion due to natural causes, while induced abortion is to deliberately terminate the pregnancy. From the data gathered in this data set, the distribution of spontaneous and induced abortions are similar, in that the majority of women in this data reported to having 0 of either types of abortions. On the other hand, about 1/3 of the total women in the data set reported having 2 abortions in the respective type of abortion.

It's valuable look at the distribution of the age of women in this data set. The surveyed women are from 20 to 45 years old. The distribution of the age variable has center at around 28-30 years old and follow as a bell shape with heavier tail on both sides. Another interesting variable is parity, which is the number of prior live births. This variable has a negative slope where the majority of women reported to having 1 live birth and as the number of live birth increase, the number of women reported having said number of live births decreases.

1.4 Reflections

Since one person has more experienced with Rshiny, rather than having the person learn RShiny entirely, we decided to just let the RShiny expert person code the app with all of its features, while the less experienced person handle other tasks. The second person is very confident in writing, so they're tasked with writing and testing the app. Their job is to tinker with the app to figure out and document how the app works, how easy the app is to use, and to catch any potential app quirks from the user experience. This division of tasks is more time efficient, minimizes friction, and lets both people focus on tasks that best suits their expertise.

As of now, this app cannot compare two variables to learn about their relationships. This feature can be implemented in a future version of this app.