

```

1 ---  

2 title: "courseworkQ2a"  

3 author: "Chris_Pang"  

4 date: "2024-04-08"  

5 output: html_document  

6 ---  

7  

8 ```{r setup, include=FALSE}  

9 knitr::opts_chunk$set(echo = TRUE)  

10 ````  

11  

12 ## For Coursework Q, Student ID: 220457882  

13  

14 Creating database. Cleaning and preparing data  

15 ```{r database}  

16  

17 library(DBI)  

18 library(dplyr)  

19 library(stringr)  

20 library(ggplot2)  

21 library(gtthemes)  

22 # set the working directory that contains the files  

23 setwd("/Users/dataverse_files")  

24  

25 # ===== create the database =====  

26 if (file.exists("airline2.db"))  

27   file.remove("airline2.db")  

28 conn <- dbConnect(RSQLite::SQLite(), "airline2.db")  

29  

30  

31  

32 # ===== write to the database =====  

33 # load in the data from the csv files  

34 airports <- read.csv("airports.csv", header = TRUE)  

35 carriers <- read.csv("carriers.csv", header = TRUE)  

36 planes <- read.csv("plane-data.csv", header = TRUE)  

37 dbWriteTable(conn, "airports", airports)  

38 dbWriteTable(conn, "carriers", carriers)  

39 dbWriteTable(conn, "planes", planes)  

40  

41 for(i in c(1991:2000)) {  

42   ontime <- read.csv(paste0(i, ".csv"), header = TRUE)  

43   if(i == 2000) {  

44     dbWriteTable(conn, "ontime", ontime)  

45   } else {  

46     dbWriteTable(conn, "ontime", ontime, append = TRUE)  

47   }  

48 }  

49  

50 ````
```

Warning: The working directory was changed to /Users/dataverse\_files inside a notebook chunk. The working directory will be reset when the chunk is finished running. Use the knitr root.dir option in working directory for notebook chunks.  
[1] TRUE

Error: Table ontime exists in database, and both overwrite and append are FALSE

```

51  

52  

53 Answering Q2a, plotting the graph individually from 1991 to 1999.  

54 ```{r Q2a}  

55 #=====Q2a Best time of the day to fly throught 1991 to 1999=====  

56 library(DBI)  

57 library(dplyr)  

58 library(ggplot2)  

59 library(RColorBrewer)  

60 time_segment <- function(deptime) {  

61   hour <- as.numeric(substr(deptime, 1, 2))  

62  

63   if (hour >= 2 && hour < 6) {  

64     return('Early Morning')  

65   } else if (hour >= 6 && hour < 10) {  

66     return('Morning')  

67   } else if (hour >= 10 && hour < 14) {  

68     return('Midday')  

69   } else if (hour >= 14 && hour < 17) {  

70     return('Afternoon')  

71   } else if (hour >= 17 && hour < 21) {  

72     return('Evening')  

73   } else {  

74     return('Night')  

75   }  

76 }  

77  

78 # Initialize a list to store each year's lowest delay count and corresponding segments  

79 lowest_delay_counts_by_year <- list()  

80  

81 for (year in 1991:1999) {  

82   query <- sprintf("SELECT CAST(DepTime AS INTEGER) AS DepTime FROM ontime WHERE ArrDelay > 15 AND DepDelay > 15 AND Year = %d", year)  

83   timeOfDayDf <- dbGetQuery(conn, query)  

84   timeOfDayDf$DepTime <- sprintf("%04d", timeOfDayDf$DepTime)  

85   timeOfDayDf$Segment <- sapply(timeOfDayDf$DepTime, time_segment)  

86  

87   segment_order <- factor(c('Early Morning', 'Morning', 'Midday', 'Afternoon', 'Evening', 'Night'),  

88                           levels = c('Early Morning', 'Morning', 'Midday', 'Afternoon', 'Evening', 'Night'))  

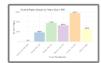
89   timeOfDayDf$Segment <- factor(timeOfDayDf$Segment, levels = levels(segment_order))  

90
```

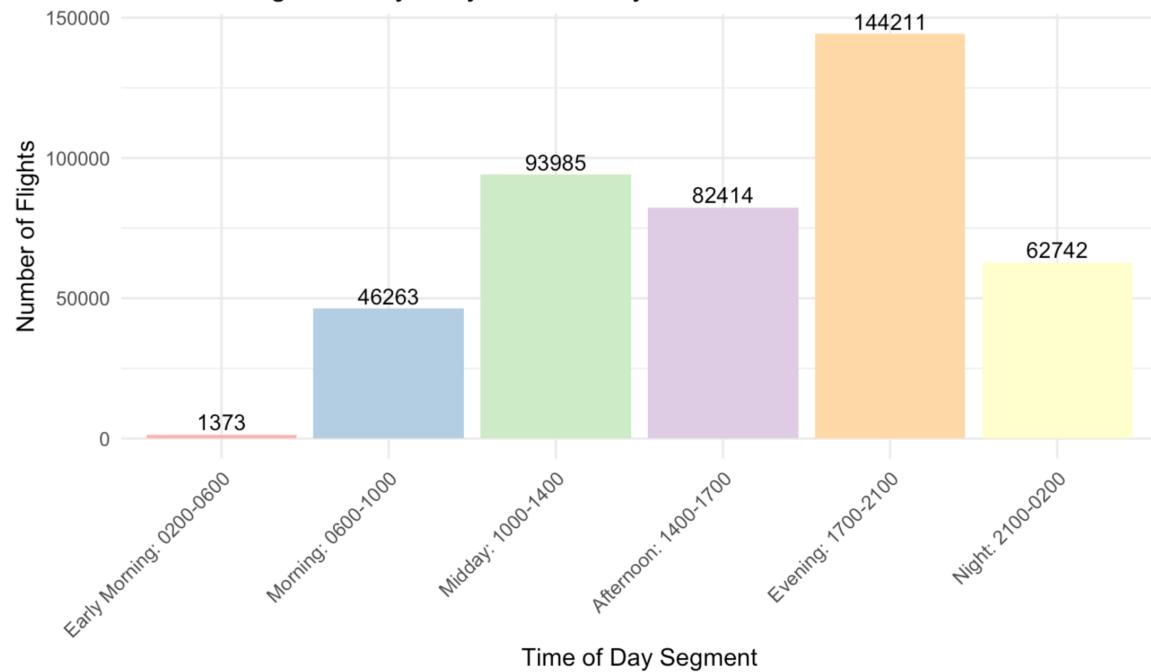
```

91 segment_counts_ordered <- as.data.frame(table(timeOfDayDF$Segment))
92 colnames(segment_counts_ordered) <- c("Segment", "Count")
93
94 # Find the lowest delay count for the year and store it
95 min_delay_count <- min(segment_counts_ordered$Count)
96 lowest_delay_segments <- segment_counts_ordered$Segment[which.min(segment_counts_ordered$Count)]
97 lowest_delay_counts_by_year[[as.character(year)]] <- list(count = min_delay_count, segment = lowest_delay_segments)
98
99 # Plotting
100 ggplot(segment_counts_ordered, aes(x = Segment, y = Count, fill = Segment)) +
101   geom_bar(stat = "identity", show.legend = FALSE) +
102   geom_text(aes(label = Count), vjust = -0.3, size = 3.5) +
103   scale_fill_brewer(palette = "Pastel1") +
104   labs(title = sprintf("Count of Flights Delayed by Time of Day in %d", year),
105       x = "Time of Day Segment",
106       y = "Number of Flights") +
107   theme_minimal() +
108   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
109   guides(fill = guide_legend(title = "Time Segments", override.aes = list(fill = brewer.pal(6, "Pastel1")))) +
110   scale_x_discrete(labels = c("Early Morning: 0200-0600", "Morning: 0600-1000", "Midday: 1000-1400",
111                       "Afternoon: 1400-1700", "Evening: 1700-2100", "Night: 2100-0200"))
112
113 print(ggplot2::last_plot()) # Print the plot for each year
114 }
115
116
117
118 * =====Q2a plotting the 9 years collated graph =====
119
120 years <- names(lowest_delay_counts_by_year)
121 lowest_counts <- sapply(lowest_delay_counts_by_year, function(info) info$count)
122 segments <- sapply(lowest_delay_counts_by_year, function(info) info$segment)
123
124 # Convert counts to numeric
125 lowest_counts <- as.numeric(lowest_counts)
126
127 # Determine the range of y-values to dynamically position text within the graph
128 y_min <- min(lowest_counts)
129 y_max <- max(lowest_counts)
130 y_range <- y_max - y_min
131 offset <- y_range * 0.02 # Offset for text positioning, 1% of the y-range
132 plot(years, lowest_counts, type = "o", col = "#008080", xlab = "Year", ylab = "Number of Delayed Flights", main = "Lowest Number of Delayed
133 grid()
134 lines(years, lowest_counts, type = "o", col = "#008080")
135
136 #adjusting text size
137 text_size <- 0.501
138
139
140 # Adjust the offset to position text below the points
141 offset <- -abs(offset) # Ensure offset is negative to move text below
142
143 for (i in 1:length(years)) {
144   year <- years[i]
145   count <- lowest_counts[i]
146   segment <- segments[i]
147
148
149   vertical_position <- count + offset # Adjust vertical position to be below the point
150
151   # Place text below each point. Adjust 'adj' for fine-tuning horizontal alignment if necessary.
152   text(year, vertical_position, labels = paste(segment, "\n(", count, ")", sep=""), col = "darkblue", cex = text_size, adj = c(0.5, 1))
153
154 }
155
156
157 * ``

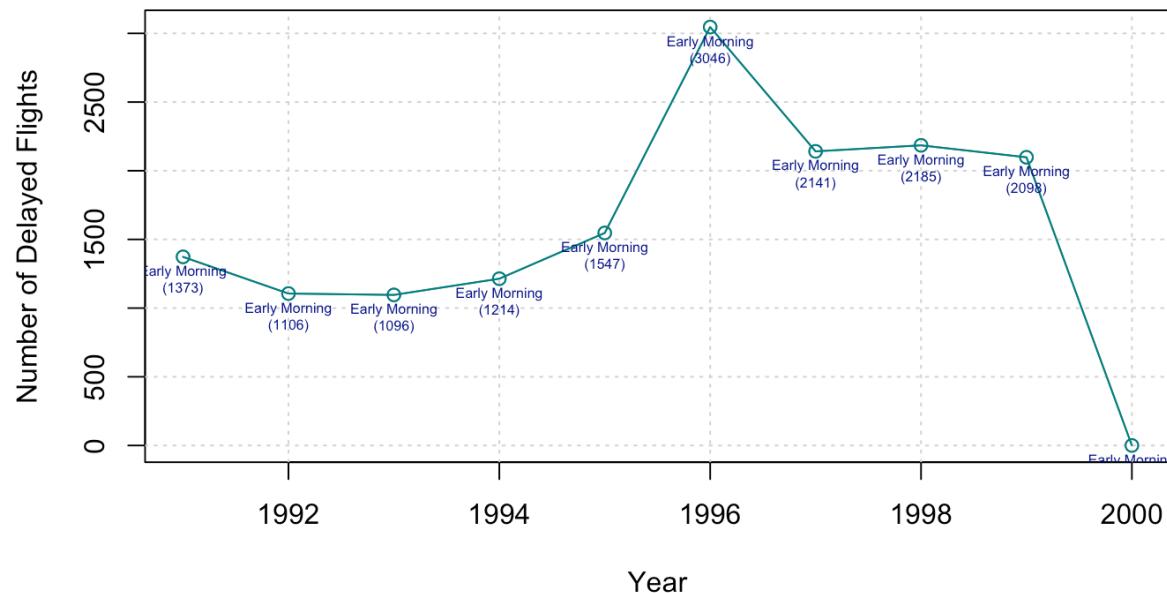
```



### Count of Flights Delayed by Time of Day in 1991



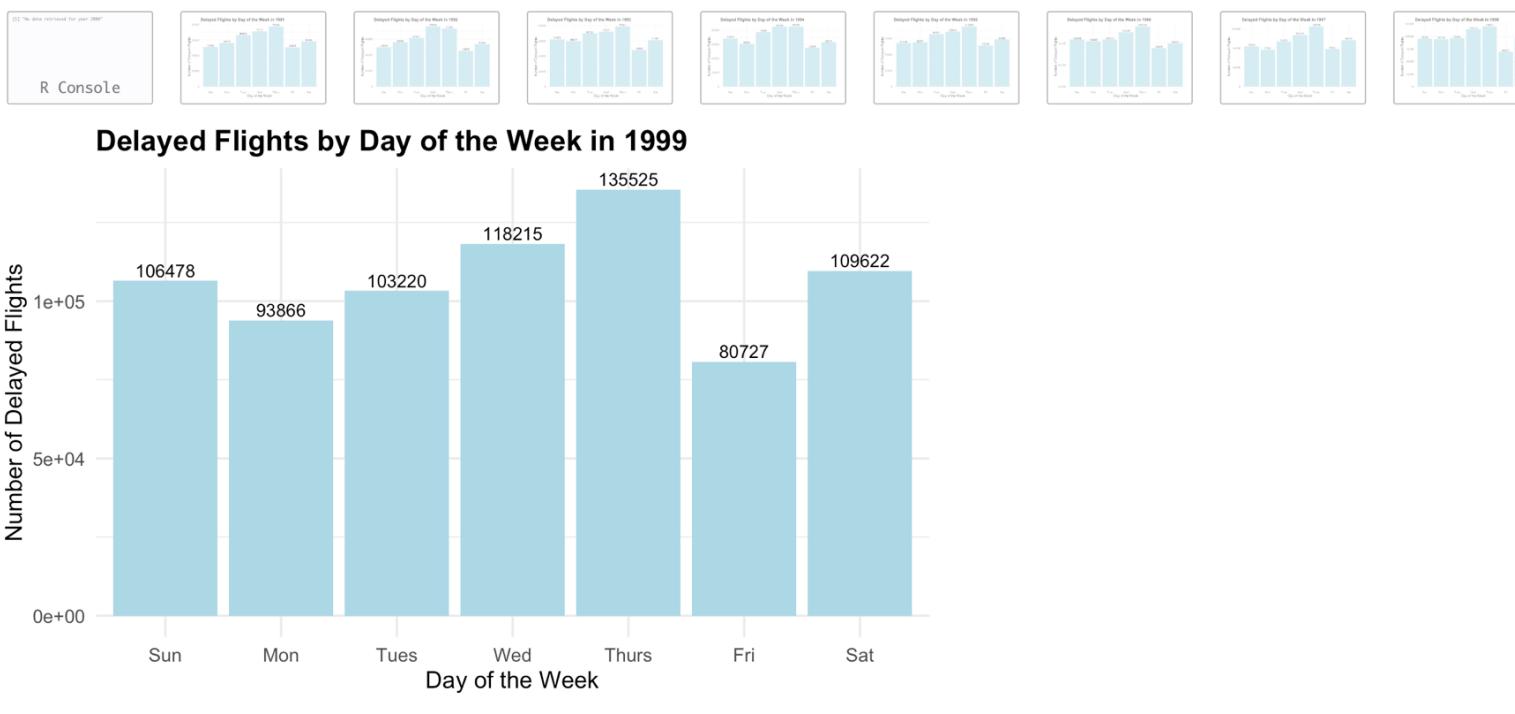
### Lowest Number of Delayed Flights by Year with Time of Day



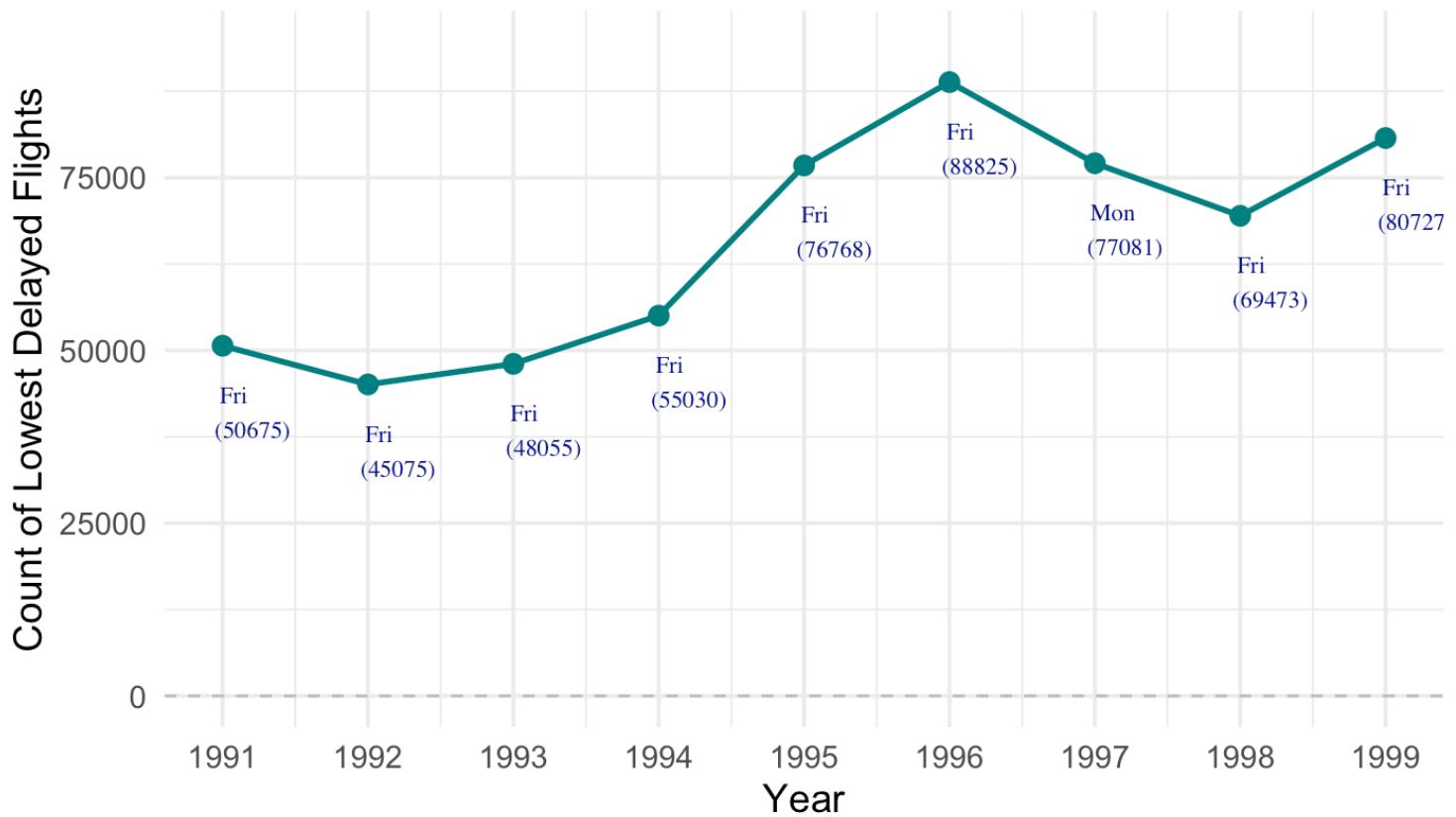
```

158
159
160 Answering Q2a, plotting best day of the week to fly individually from 1991-1999
161 ~ ``{r Q2a Day of the week}
162
163 ~ =====Q2a best day of the week to fly to minimize delay =====
164 lowest_delay_day_by_year <- list()
165
166 ~ for (year in 1991:2000) {
167   query <- sprintf(
168     "SELECT DayOfWeek, COUNT(*) AS Count
169     FROM ontime
170     WHERE Year = %d AND ArrDelay > 15 AND DepDelay > 15
171     GROUP BY DayOfWeek
172     ORDER BY DayOfWeek", year)
173
174   # Execute the query
175   dayOfWeekCounts <- dbGetQuery(conn, query)
176
177   # Initialize a dataframe for all days with counts set to 0
178   counts <- data.frame(DayOfWeek = 1:7, Count = rep(0, 7))
179
180   # Check if any data is retrieved
181   if (nrow(dayOfWeekCounts) == 0) {
182     print(paste("No data retrieved for year", year))
183     next # Skip this year if no data is returned
184   }
185
186   # Update counts based on the query result
187   for (row in 1:nrow(dayOfWeekCounts)) {
188     # Ensure that 'DayOfWeek' column exists and has valid data
189     if ("DayOfWeek" %in% names(dayOfWeekCounts) && !is.na(dayOfWeekCounts[row, "DayOfWeek"])) {
190       day <- dayOfWeekCounts[row, "DayOfWeek"]
191       if (day >= 1 && day <= 7) {
192         counts[day, "Count"] <- dayOfWeekCounts[row, "Count"]
193       } else {
194         print(paste("Invalid day:", day, "in year", year))
195       }
196     } else {
197       print(paste("Missing or invalid 'DayOfWeek' data in year", year))
198     }
199   }
200
201   # Plotting a bar chart for each year
202   tick_label <- c('Sun', 'Mon', 'Tues', 'Wed', 'Thurs', 'Fri', 'Sat')
203   ggplot(counts, aes(x = factor(DayOfWeek, labels = tick_label), y = Count)) +
204     geom_bar(stat = "identity", fill = "#lightblue") +
205     geom_text(aes(label = Count), vjust = -0.3, size = 3.5) +
206     labs(x = 'Day of the Week', y = 'Number of Delayed Flights',
207       title = sprintf('Delayed Flights by Day of the Week in %d', year)) +
208     theme_minimal(base_size = 13) +
209     theme(plot.title = element_text(size = 16, face = "bold"))
210
211   # Force the plot to display
212   print(ggplot2::last_plot())
213
214   # Store the day with the least delays for the year
215   min_delay_count <- min(counts$Count)
216   min_delay_day <- tick_label[which.min(counts$Count)]
217   lowest_delay_day_by_year[[as.character(year)]] <- list(Day = min_delay_day, Count = min_delay_count)
218 }
219
220
221
222 ~ =====Q2a plotting the 9 years collated graph =====
223
224 # Create a data frame from the list
225 data <- data.frame(
226   Year = as.numeric(names(lowest_delay_day_by_year)),
227   Lowest_Count = sapply(lowest_delay_day_by_year, function(info) info$Count),
228   Lowest_Days = sapply(lowest_delay_day_by_year, function(info) info$Day),
229 )
230
231 # Plotting
232 ggplot(data, aes(x = Year, y = Lowest_Count)) +
233   geom_line(color = "#008080", size = 1) + # Using hex code for teal
234   geom_point(color = "#008080", size = 3) + # Using hex code for teal
235   geom_text(aes(label = paste(Lowest_Days, "\n", Lowest_Count, ")"), sep = ""),
236             nudge_y = (max(data$Lowest_Count) - min(data$Lowest_Count)) * 0.03,
237             size = 3, color = "darkblue",
238             hjust = 0.1, vjust = 2,
239             check_overlap = TRUE,
240             family = "serif") +
241   labs(title = 'Lowest Number of Delayed Flights by Year',
242        x = 'Year', y = 'Count of Lowest Delayed Flights') +
243   theme_minimal(base_size = 14) +
244   theme(plot.title = element_text(size = 16, face = "bold")) +
245   scale_x_continuous(breaks = data$Year) + # Ensure X axis has all years as breaks
246   scale_y_continuous(expand = expansion(mult = c(0.05, 0.1))) + # Add some padding around Y axis
247   geom_hline(yintercept = 0, linetype = "dashed", color = "gray")
248
249 # Adjust plot theme for better readability
250 theme_set(theme_minimal(base_size = ))

```



## Lowest Number of Delayed Flights by Year

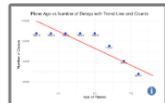
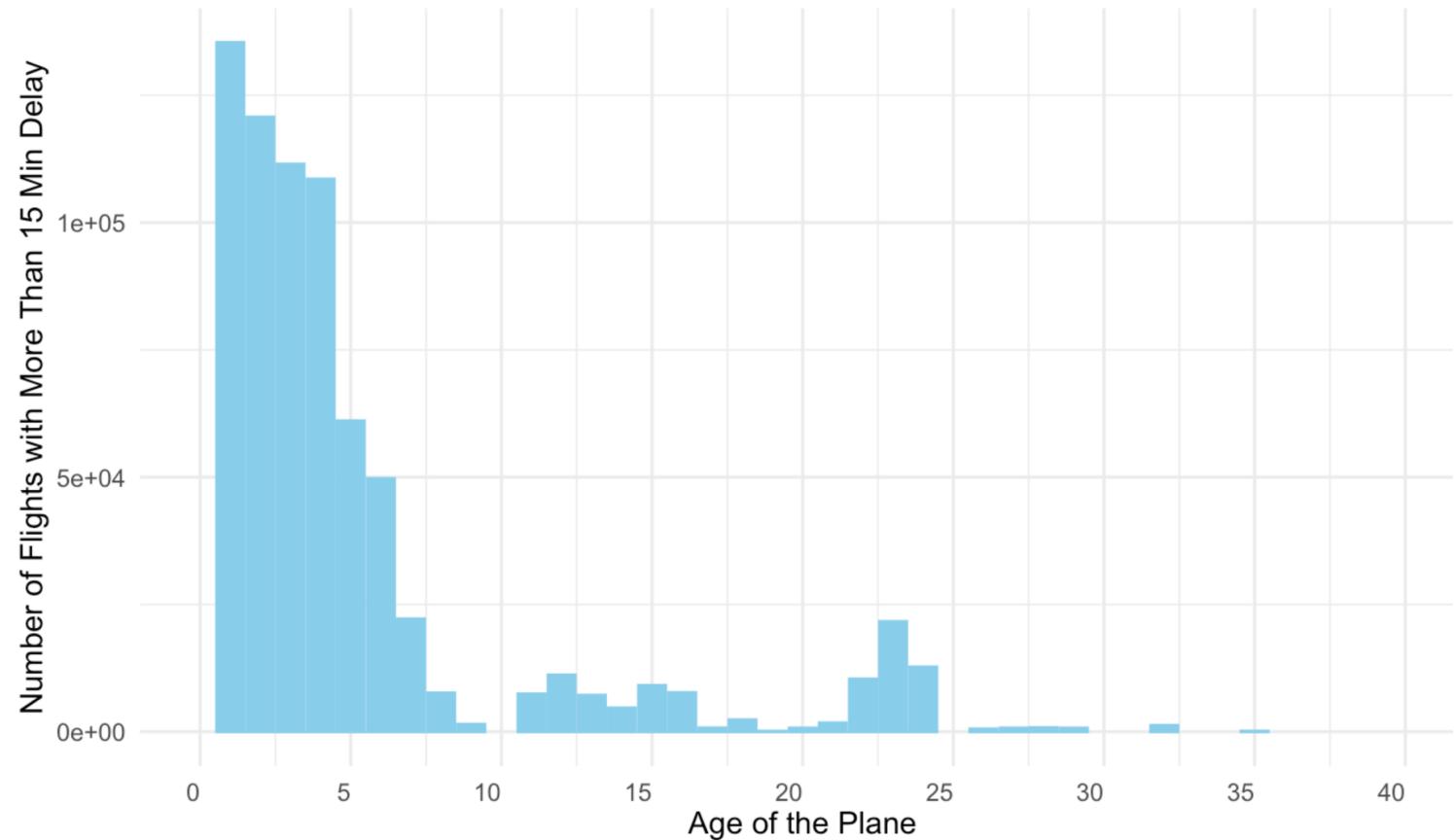


```

252 For Q2b, Plotting do older planes suffer more delay graph individually from 1991-1999
253 ````{r Q2b do older planes suffer more delay}
254
255 #=====Q2b plotting do older planes suffer more delay on a year to year basis=====
256
257 # Initialize an empty list to store the data
258 max_delay_age_by_year <- list()
259
260 for (year in 1991:1999) {
261   # Construct the SQL query
262   query <- sprintf(
263     "SELECT CAST(planes.year AS INTEGER) AS year
264     FROM planes
265     JOIN ontime ON planes.TailNum = ontime.TailNum
266     WHERE CAST(ontime.ArrDelay AS INTEGER) > 15 AND CAST(ontime.DepDelay AS INTEGER) > 15
267     AND planes.year IS NOT NULL
268     AND CAST(ontime.Year AS INTEGER) >= %s
269   ", year)
270
271 # Execute the query
272 dateOfPlane <- dbGetQuery(conn, query)
273
274 # Check if data was fetched
275 if (nrow(dateOfPlane) == 0) {
276   cat(sprintf("No data fetched for %d.\n", year))
277   next
278 }
279
280 # Calculate the age of the plane in that year
281 dateOfPlane$`Age in year` <- year - dateOfPlane$year
282
283 # Count by age
284 age_counts <- count(dateOfPlane, `Age in year`)
285
286 # Find the age with the maximum count
287 max_age <- age_counts[which.max(age_counts$n),]
288 max_delay_age_by_year[[as.character(year)]] <- list(Age = max_age$`Age in year`, Delays = max_age$n)
289
290 # Plot the graph for the specific year
291 p <- ggplot(age_counts, aes(x = `Age in year`, y = n, fill = `Age in year`)) +
292   geom_bar(stat = "identity", color = "#skyblue", fill = "#skyblue") +
293   scale_x_continuous(limits = c(0, 40), breaks = seq(0, 40, by = 5)) +
294   labs(title = sprintf("Number of Flights with Delay by Plane Age in %d", year),
295        x = "Age of the Plane",
296        y = "Number of Flights with More Than 15 Min Delay") +
297   theme_minimal() +
298   theme(axis.text.x = element_text(vjust = 0.5, hjust=1))
299
300 print(p)
301 }
302
303 # Output the results
304 for (year in names(max_delay_age_by_year)) {
305   cat(sprintf("In %s, the age with the most delays was %s years with %s delays.\n",
306             year, max_delay_age_by_year[[year]]$Age, max_delay_age_by_year[[year]]$Delays))
307 }
308
309
310 #=====plotting the 9 years collated regression graph=====
311 library(ggplot2)
312 library(dplyr)
313
314 # Convert the list to a data frame
315 data <- bind_rows(max_delay_age_by_year, .id = "Year") %>%
316   mutate(Year = as.numeric(Year), Age = as.numeric(Age), Delays = as.numeric(Delays)) %>%
317   arrange(Year)
318
319 # Plotting with annotations
320 p <- ggplot(data, aes(x = Age, y = Delays)) +
321   geom_point(color = "blue", size = 3, show.legend = TRUE, aes(label = Delays)) +
322   geom_text(aes(label = Delays), vjust = 1.5, color = "black", size = 3) +
323   labs(title = "Plane Age vs Number of Delays with Counts",
324        x = "Age of Planes", y = "Number of Delays") +
325   theme_minimal() +
326   theme(plot.title = element_text(size = 16),
327         axis.title = element_text(size = 13)) +
328   geom_smooth(method = "lm", color = "red", se = FALSE, show.legend = TRUE) +
329   ggtitle("Plane Age vs Number of Delays with Trend Line and Counts") +
330   xlab("Age of Planes") + ylab("Number of Delays") +
331   theme(legend.position = "bottom") +
332   scale_color_manual(values = c("blue", "red"), labels = c("Delays vs. Age", "Line of Best Fit"))
333
334 print(p)
335
336
337 # Optionally, calculate the correlation coefficient
338 correlation <- cor(data$Age, data$Delays, use = "complete.obs")
339 cat(sprintf("Correlation coefficient between plane age and number of delays: %.2f\n", correlation))
340
341
342 ```

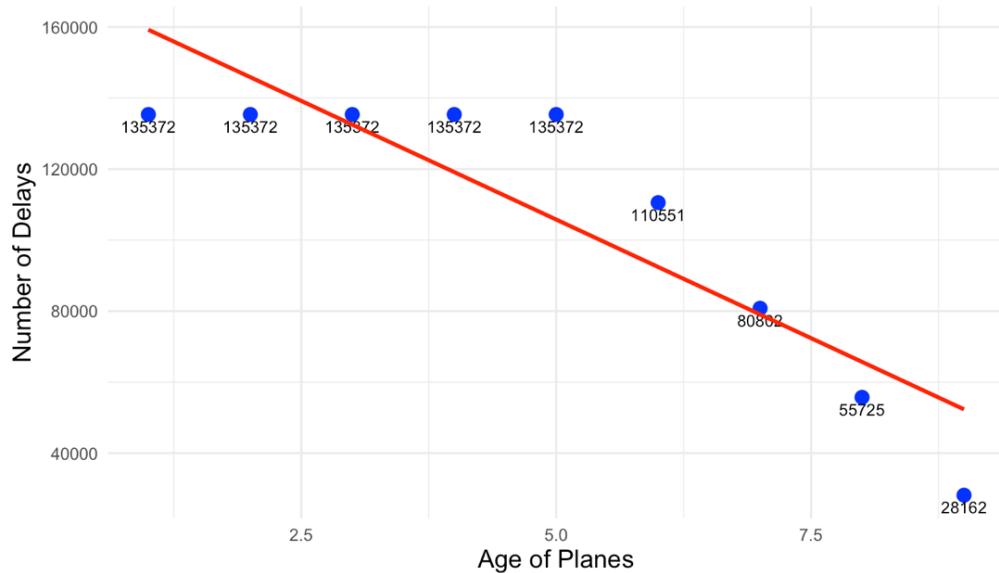
```

## Number of Flights with Delay by Plane Age in 1991



ℹ [38;5;232m`geom\_smooth()` using formula = 'y ~ x' [39m

Plane Age vs Number of Delays with Trend Line and Counts



```

345 Plotting for q2c)
346 ````{r Q2c}
347
348 library(tidyverse)
349
350 set.seed(42)
351
352 for (year in 1991:2000) {
353   # Subset the data for the current year
354   year_data <- filter(ontime, Year == year)
355
356   # Impute missing values in features
357   imputed_data <- year_data %>%
358     select(DepTime, ArrTime, Distance) %>%
359     bind_cols(Year = year_data$Year, Diverted = year_data$Diverted) %>%
360     replace_na(list(DepTime = mean(year_data$DepTime, na.rm = TRUE),
361                     ArrTime = mean(year_data$ArrTime, na.rm = TRUE),
362                     Distance = mean(year_data$Distance, na.rm = TRUE))).
363
364   # Splitting the data into features (X) and target (y)
365   X <- imputed_data %>% select(DepTime, ArrTime, Distance)
366   y <- imputed_data$Diverted
367
368   # Splitting dataset into training and testing sets
369   trainIndex <- createDataPartition(y, p = .8, list = FALSE, times = 1)
370   X_train <- X[trainIndex, ]
371   X_test <- X[-trainIndex, ]
372   y_train <- y[trainIndex]
373   y_test <- y[-trainIndex]
374
375   # Initializing and fitting the logistic regression model
376   logistic_model <- train(X_train, y_train, method = 'glm', family = 'binomial', trControl = trainControl(method = 'none'))
377
378   # Getting coefficients and intercept
379   coefficients <- c(intercept = coef(logistic_model$finalModel)[['(Intercept)']], coef(logistic_model$finalModel)[-1])
380
381   # Creating a DataFrame to hold coefficients for the current year
382   temp_df <- data.frame(t(coefficients), Year = year)
383
384   # Adding temporary DataFrame to the main coefficients DataFrame
385   coefficients_df <- bind_rows(coefficients_df, temp_df)
386 }
387
388 # Setting 'Year' as the DataFrame index for plotting
389 coefficients_df$Year <- factor(coefficients_df$Year)
390
391 # Plotting the coefficients for each feature across years
392 coefficients_df %>%
393   pivot_longer(cols = -Year, names_to = "Feature", values_to = "Coefficient") %>%
394   ggplot(aes(x = Year, y = Coefficient, color = Feature, group = Feature)) +
395   geom_line() +
396   geom_point() +
397   labs(title = 'Logistic Regression Coefficients for Diverted Flights Across Years',
398        x = 'Year', y = 'Coefficient Value') +
399   theme_minimal() +
400   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
401   scale_color_brewer(palette = "Dark2")
402
403
404 ````
```