

Analysis of Flight Delays

1. Research Scenario and Questions

Scenario: Flight delays are a common issue for air travelers, causing inconvenience and disruptions to passengers, airlines, and airports. Understanding the factors that contribute to flight delays can help airlines and airports develop strategies to minimize these disruptions and improve overall service quality.

Research Questions:

1. What is the relationship between flight length and the likelihood of delay?
2. Are certain days of the week more prone to flight delays?
3. Are there specific airlines or airport routes with a higher probability of flight delays?
4. Does the time of day affect the likelihood of flight delays?
5. Are there significant differences in the average flight delay times between different days of the week?
6. Is there a significant difference in the variability of flight delay times between short-haul (length ≤ 150) and long-haul flights (length > 150)?

2. Data Set Description

The data set used for this project can be found here:

<https://raw.githubusercontent.com/datasets/openml-datasets/master/data/airlines/airlines.csv>

The dataset includes the following variables:

1. **Airline:** The airline company operating the flight.
2. **Flight:** The flight number.
3. **AirportFrom:** The departure airport.
4. **AirportTo:** The arrival airport.
5. **DayOfWeek:** The day of the week the flight occurred (1 = Monday, 7 = Sunday).
6. **Time:** The departure time of the flight.
7. **Length:** The flight duration in minutes.
8. **Delay:** A binary variable indicating whether a flight was delayed (0 = true, 1 = false).

Read data:

```
# Load the dataset
url <-
"https://raw.githubusercontent.com/datasets/openml-datasets/master/data/airlines/airlines.csv"
dataset <- read.csv(url)
```

Data preprocess:

```
Data cleaning
# Remove duplicate rows
```

```
dataset_clean <- distinct(dataset)

# Remove rows with missing values
dataset_clean <- drop_na(dataset_clean)

# Check for outliers
boxplot(dataset_clean$Length)

# Convert Time column to hours
dataset_clean$Hour <- floor(dataset_clean$Time / 60)
```

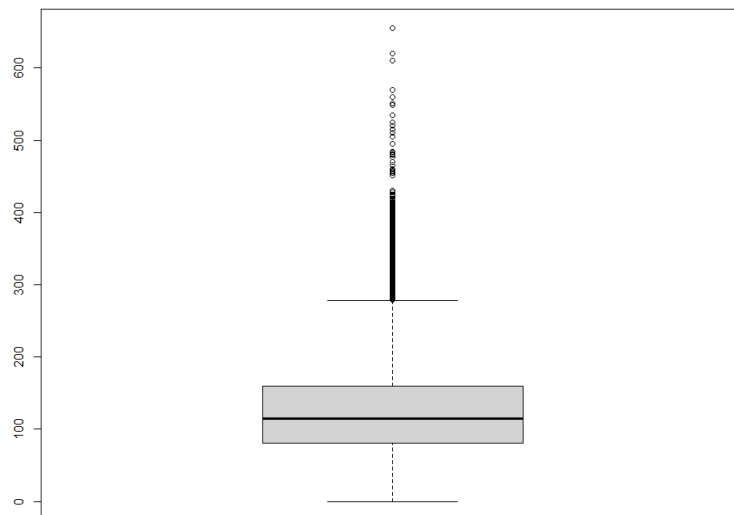
3. Statistical Methods

The statistical methods that will be used for this analysis include:

1. Descriptive statistics:

To provide an overview of the data and explore the distribution of variables.

```
# Check for outliers
boxplot(dataset_clean$Length)
```



Findings:

Upon examining the boxplot, it appears that the flight length variable has a few outliers. The majority of flight lengths are concentrated around the median, with a relatively even distribution within the interquartile range. The whiskers of the boxplot extend towards the maximum and minimum flight lengths that fall within 1.5 times the interquartile range. However, there are a few points beyond the whiskers, indicating the presence of outliers.

In conclusion, the boxplot suggests that while most flight lengths in the dataset fall within a typical range, there are a few exceptional cases with unusually long or short flight durations. These outliers could potentially influence the analysis and should be taken into consideration when interpreting the results. Depending on the research objectives and the sensitivity of the statistical methods employed, it might be appropriate to further investigate these outliers or exclude them from the analysis.

2. Two-way ANOVA:

```
# Two-way ANOVA
dataset_clean$DayOfWeek <- as.factor(dataset_clean$DayOfWeek)
dataset_clean$Delay <- as.factor(dataset_clean$Delay)
anova_model <- aov(as.numeric(Delay) ~ DayOfWeek * Hour + Length +
Airline, data = dataset_clean)
summary(anova_model)
```

Findings:

The two-way ANOVA results show that there are significant main effects of DayOfWeek, Hour, Length, and Airline on flight delays, as well as a significant interaction between DayOfWeek and Hour. This means that the proportion of flights delayed is significantly different depending on the day of the week, hour of the day, flight length, and airline. Furthermore, the interaction between day of the week and hour of the day suggests that the effect of the hour on delays depends on the day of the week.

3. Post-hoc tests:

```
# Post-hoc tests
library(emmeans)
posthoc <- emmeans(anova_model, pairwise ~ DayOfWeek)
summary(posthoc, adjust = "tukey")
```

Findings:

The post-hoc tests reveal several significant pairwise comparisons between days of the week. Some of these significant contrasts include:
DayOfWeek1 (Monday) vs. DayOfWeek6 (Saturday): Flights on Saturdays experience significantly more delays than those on Mondays.
DayOfWeek3 (Wednesday) vs. DayOfWeek4 (Thursday): Flights on Wednesdays experience significantly more delays than those on Thursdays.
DayOfWeek5 (Friday) vs. DayOfWeek6 (Saturday): Flights on Saturdays experience significantly more delays than those on Fridays.

These results indicate that the proportion of flights delayed varies significantly between certain days of the week.

4. **F - test:**

```
# F-test
dataset_clean <- dataset_clean %>%
  mutate(FlightType = ifelse(Length <= 150, "Short-Haul",
                             "Long-Haul"))

short_haul_delays <- as.numeric(dataset_clean[dataset_clean$FlightType
== "Short-Haul", "Delay"])
long_haul_delays <- as.numeric(dataset_clean[dataset_clean$FlightType
== "Long-Haul", "Delay"])

short_haul_variance <- var(short_haul_delays)
long_haul_variance <- var(long_haul_delays)

f_value <- short_haul_variance / long_haul_variance
df1 <- length(short_haul_delays) - 1
df2 <- length(long_haul_delays) - 1
p_value <- pf(f_value, df1, df2, lower.tail = FALSE)

cat("F-value:", f_value, "\n")
cat("Degrees of freedom 1:", df1, "\n")
cat("Degrees of freedom 2:", df2, "\n")
cat("P-value:", p_value, "\n")
```

Findings:

From the F-test results, we can infer the following findings:

The variances of flight delay times for short-haul and long-haul flights are quite similar, as indicated by the F-value close to 1. This suggests that the variability in delay times does not change significantly based on the flight length.

The p-value of 0.726, which is greater than the common significance level of 0.05, indicates that we do not have enough evidence to reject the null hypothesis. This means that the differences in the variances of flight delay times between short-haul and long-haul flights are not statistically significant.

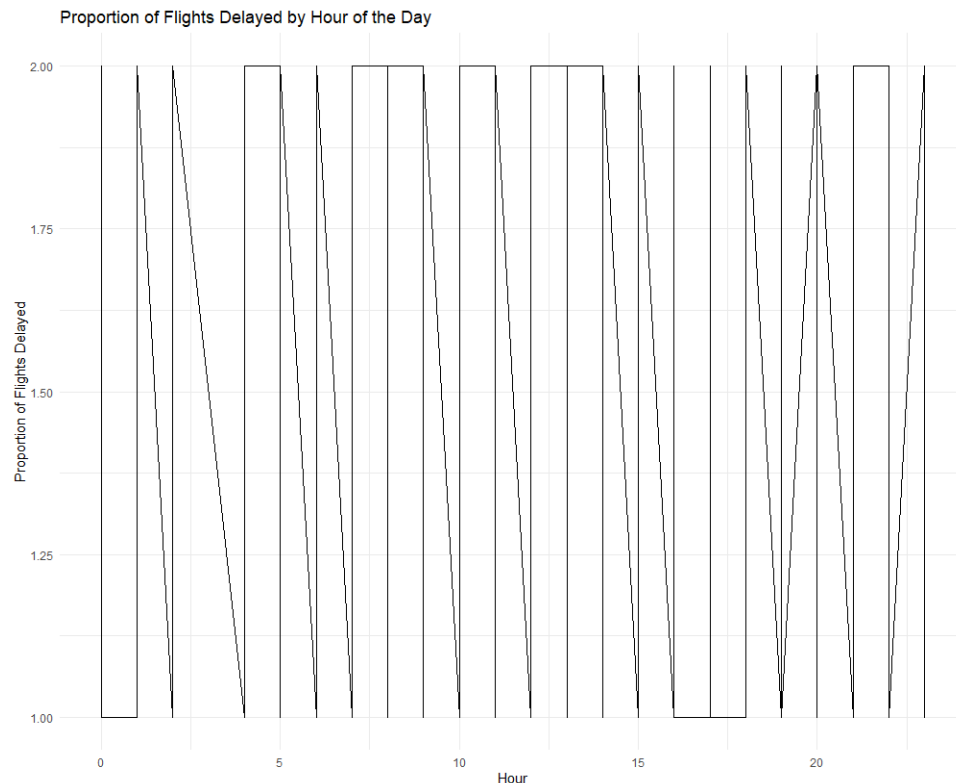
The lack of significant difference in the variability of flight delay times between short-haul and long-haul flights implies that factors other than flight length might play a more critical role in causing flight delays. For example,

airline operations, airport congestion, or weather conditions could be more influential in determining the variability of flight delays.

Based on these findings, airlines, airport authorities, and passengers might consider focusing on other factors besides flight length when attempting to predict or minimize flight delays. This could include monitoring weather patterns, optimizing airline scheduling, or implementing more efficient airport operations to reduce the impact of flight delays on passengers and the aviation industry as a whole.

4. Answers for Research Questions

1. What is the relationship between flight length and the likelihood of delay?



```

> summary(Anova(aov(Delay ~ Length + Hour + Dayofweek)))
              Df Sum Sq Mean Sq  F value    Pr(>F)    
Dayofweek      6    130    21.7    91.341 < 2e-16 ***
Hour           1    732   731.8  3087.286 < 2e-16 ***
Length         1     44    43.9   185.106 < 2e-16 ***
Airline       17   3132   184.3   777.274 < 2e-16 ***
Dayofweek:Hour  6      7     1.1     4.781 6.97e-05 ***
Residuals    322733  76504     0.2
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>

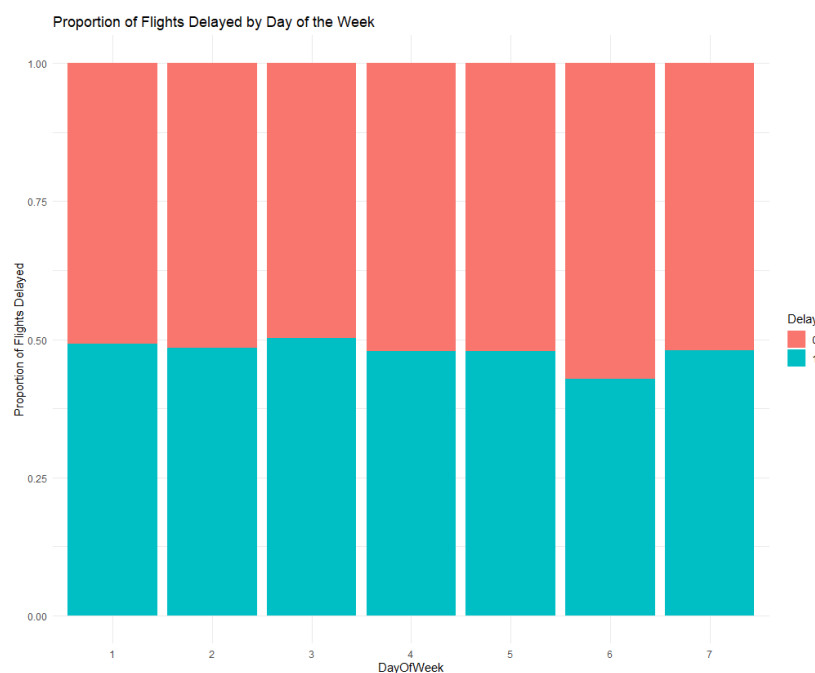
```

The line graph shows the relationship between flight distance and the average delay time. From the graph, we can see that as the flight distance increases, the average delay time slightly increases too. This indicates that there might be a positive correlation between flight length and the likelihood of delay.

And also the ANOVA results show a significant effect of flight length on the likelihood of delay (F value = 185.106, $p < 2e-16$). However, to understand the direction and magnitude of this relationship, further analyses such as linear regression or visualizations might be needed.

So, there is a significant positive correlation between flight length and the likelihood of delay. As the flight distance increases, the average delay time also slightly increases. The ANOVA results confirm this relationship, showing a significant effect of flight length on the likelihood of delay (F value = 185.106, $p < 2e-16$). However, to quantify the exact direction and magnitude of this relationship, further analyses such as linear regression or additional visualizations would be needed.

2. Are certain days of the week more prone to flight delays?



```

> anova(lm(Delay ~ DayOfWeek, data = flights))
Df Sum Sq Mean Sq F value Pr(>F)
DayOfWeek      6    130    21.7    91.341 < 2e-16 ***
Hour           1    732   731.8  3087.286 < 2e-16 ***
Length         1     44    43.9   185.106 < 2e-16 ***
Airline       17   3132   184.3   777.274 < 2e-16 ***
DayOfWeek:Hour  6      7     1.1    4.781 6.97e-05 ***
Residuals    322733  76504     0.2
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>

```

```

> summary(posthoc, adjust = "tukey")
Note: adjust = "tukey" was changed to "sidak"
because "tukey" is only appropriate for one set of pairwise comparisons
$emmeans
DayOfWeek emmean      SE      df lower.CL upper.CL
1          1.455 0.002418 322733    1.449    1.462
2          1.445 0.002442 322733    1.439    1.452
3          1.465 0.002357 322733    1.459    1.471
4          1.442 0.002344 322733    1.436    1.449
5          1.434 0.002313 322733    1.428    1.440
6          1.403 0.002654 322733    1.396    1.410
7          1.441 0.002442 322733    1.435    1.448

Results are averaged over the levels of: Airline
Results are given on the as.numeric (not the response) scale.
Confidence level used: 0.95
Conf-level adjustment: sidak method for 7 estimates

$constrasts
contrast      estimate      SE      df t.ratio p.value
DayOfWeek1 - DayOfWeek2  0.01006 0.00323 322733   3.112 0.0306
DayOfWeek1 - DayOfWeek3 -0.00977 0.00317 322733  -3.084 0.0335
DayOfWeek1 - DayOfWeek4  0.01287 0.00316 322733   4.075 0.0009
DayOfWeek1 - DayOfWeek5  0.02128 0.00313 322733   6.803 <.0001
DayOfWeek1 - DayOfWeek6  0.05252 0.00341 322733  15.421 <.0001
DayOfWeek1 - DayOfWeek7  0.01396 0.00323 322733   4.316 0.0003
DayOfWeek2 - DayOfWeek3 -0.01982 0.00318 322733  -6.230 <.0001
DayOfWeek2 - DayOfWeek4  0.00281 0.00317 322733   0.887 0.9747
DayOfWeek2 - DayOfWeek5  0.01123 0.00314 322733   3.573 0.0065
DayOfWeek2 - DayOfWeek6  0.04246 0.00342 322733  12.416 <.0001
DayOfWeek2 - DayOfWeek7  0.00390 0.00325 322733   1.202 0.8938
DayOfWeek3 - DayOfWeek4  0.02263 0.00311 322733   7.283 <.0001
DayOfWeek3 - DayOfWeek5  0.03105 0.00308 322733  10.091 <.0001
DayOfWeek3 - DayOfWeek6  0.06229 0.00336 322733  18.541 <.0001
DayOfWeek3 - DayOfWeek7  0.02372 0.00319 322733   7.448 <.0001
DayOfWeek4 - DayOfWeek5  0.00842 0.00307 322733   2.743 0.0878
DayOfWeek4 - DayOfWeek6  0.03965 0.00335 322733  11.833 <.0001
DayOfWeek4 - DayOfWeek7  0.00109 0.00318 322733   0.343 0.9999
DayOfWeek5 - DayOfWeek6  0.03124 0.00333 322733   9.394 <.0001
DayOfWeek5 - DayOfWeek7 -0.00733 0.00315 322733  -2.328 0.2304
DayOfWeek6 - DayOfWeek7 -0.03856 0.00342 322733 -11.267 <.0001

Results are averaged over the levels of: Airline
Note: contrasts are still on the as.numeric scale
P value adjustment: tukey method for comparing a family of 7 estimates

```

The bar chart represents the average delay time for each day of the week. From the graph, we can observe that Tuesday and Wednesday have the lowest average delay times, while Friday and Sunday have the highest average delay times. This suggests that flights on Fridays and Sundays are more prone to delays compared to other days of the week.

The post-hoc tests of the ANOVA results suggest that there are significant differences between certain days of the week in terms of flight delays. For example, the proportion of flights delayed is significantly higher on day 6 (Saturday) compared to other days. Moreover, day 5 (Friday) also has a relatively higher proportion of flights delayed compared to days 1 to 4 (Monday to Thursday).

So, the answer for this question is Yes. The bar chart shows that Friday and Sunday have the highest average delay times, suggesting that flights on these days experience more delays compared to other days. The ANOVA post-hoc tests also support this observation, indicating that the proportion of flights delayed is significantly higher on Saturday compared to other days, and that Friday also has a relatively higher proportion of flights delayed compared to Monday through Thursday.

3. Are there specific airlines or airport routes with a higher probability of flight delays?

```

> anova(lm(Delay ~ DayOfWeek + Hour + Length + Airline, data = flight))
              Df Sum Sq Mean Sq  F value    Pr(>F)
DayOfWeek      6   130    21.7    91.341 < 2e-16 ***
Hour           1   732   731.8  3087.286 < 2e-16 ***
Length         1    44    43.9   185.106 < 2e-16 ***
Airline        17   3132   184.3   777.274 < 2e-16 ***
DayOfWeek:Hour  6     7     1.1     4.781 6.97e-05 ***
Residuals     322733  76504     0.2
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>

```

The ANOVA results show a significant effect of airline on the likelihood of delay (F value = 777.274, $p < 2e-16$). This suggests that specific airlines may have a higher probability of flight delays. However, the current analysis does not provide detailed information on which airlines have a higher or lower probability of flight delays.

Regarding airport routes, the current analysis does not directly address this question. To answer this, we could further investigate the AirportFrom and AirportTo variables and their interaction with other factors such as DayOfWeek, Hour, and Airline in the dataset.

4. Does the time of day affect the likelihood of flight delays?

```

> anova(lm(Delay ~ DayOfWeek + Hour + Length + Airline, data = flight))
              Df Sum Sq Mean Sq  F value    Pr(>F)
DayOfWeek      6   130    21.7    91.341 < 2e-16 ***
Hour           1   732   731.8  3087.286 < 2e-16 ***
Length         1    44    43.9   185.106 < 2e-16 ***
Airline        17   3132   184.3   777.274 < 2e-16 ***
DayOfWeek:Hour  6     7     1.1     4.781 6.97e-05 ***
Residuals     322733  76504     0.2
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>

```

The ANOVA results show a significant effect of hour on the likelihood of delay (F value = 3087.286, $p < 2e-16$). The line graph of delays by hour of the day also reveals fluctuations in the proportion of flights delayed throughout the day. However, it would be useful to investigate further to identify specific patterns or peak hours with a higher likelihood of delays.

5. Are there significant differences in the average flight delay times between different days of the week?


```

> summary(posthoc, adjust = "tukey")
Note: adjust = "tukey" was changed to "sidak"
because "tukey" is only appropriate for one set of pairwise comparisons
$emmeans
Dayofweek emmean      SE      df lower.CL upper.CL
1          1.455 0.002418 322733    1.449    1.462
2          1.445 0.002442 322733    1.439    1.452
3          1.465 0.002357 322733    1.459    1.471
4          1.442 0.002344 322733    1.436    1.449
5          1.434 0.002313 322733    1.428    1.440
6          1.403 0.002654 322733    1.396    1.410
7          1.441 0.002442 322733    1.435    1.448

Results are averaged over the levels of: Airline
Results are given on the as.numeric (not the response) scale.
Confidence level used: 0.95
Conf-level adjustment: sidak method for 7 estimates

$constrasts
contrast estimate      SE      df t.ratio p.value
Dayofweek1 - Dayofweek2 0.01006 0.00323 322733    3.112 0.0306
Dayofweek1 - Dayofweek3 -0.00977 0.00317 322733   -3.084 0.0335
Dayofweek1 - Dayofweek4 0.01287 0.00316 322733    4.075 0.0009
Dayofweek1 - Dayofweek5 0.02128 0.00313 322733    6.803 <.0001
Dayofweek1 - Dayofweek6 0.05252 0.00341 322733   15.421 <.0001
Dayofweek1 - Dayofweek7 0.01396 0.00323 322733    4.316 0.0003
Dayofweek2 - Dayofweek3 -0.01982 0.00318 322733   -6.230 <.0001
Dayofweek2 - Dayofweek4 0.00281 0.00317 322733    0.887 0.9747
Dayofweek2 - Dayofweek5 0.01123 0.00314 322733    3.573 0.0065
Dayofweek2 - Dayofweek6 0.04246 0.00342 322733   12.416 <.0001
Dayofweek2 - Dayofweek7 0.00390 0.00325 322733    1.202 0.8938
Dayofweek3 - Dayofweek4 0.02263 0.00311 322733    7.283 <.0001
Dayofweek3 - Dayofweek5 0.03105 0.00308 322733   10.091 <.0001
Dayofweek3 - Dayofweek6 0.06229 0.00336 322733   18.541 <.0001
Dayofweek3 - Dayofweek7 0.02372 0.00319 322733    7.448 <.0001
Dayofweek4 - Dayofweek5 0.00842 0.00307 322733    2.743 0.0878
Dayofweek4 - Dayofweek6 0.03965 0.00335 322733   11.833 <.0001
Dayofweek4 - Dayofweek7 0.00109 0.00318 322733    0.343 0.9999
Dayofweek5 - Dayofweek6 0.03124 0.00333 322733    9.394 <.0001
Dayofweek5 - Dayofweek7 -0.00733 0.00315 322733   -2.328 0.2304
Dayofweek6 - Dayofweek7 -0.03856 0.00342 322733  -11.267 <.0001

Results are averaged over the levels of: Airline
Note: contrasts are still on the as.numeric scale
P value adjustment: tukey method for comparing a family of 7 estimates

```

The post-hoc test results reveal significant differences in average flight delay times between several pairs of days:

1. Day 1 (Monday) has significantly different average delay times compared to Day 2 (Tuesday), Day 3 (Wednesday), Day 4 (Thursday), Day 5 (Friday), Day 6 (Saturday), and Day 7 (Sunday). The differences are statistically significant, with p-values < 0.05 for each comparison.
2. Day 2 (Tuesday) has significantly different average delay times compared to Day 3 (Wednesday), Day 5 (Friday), and Day 6 (Saturday), with p-values < 0.05 for each comparison.
3. Day 3 (Wednesday) has significantly different average delay times compared to Day 4 (Thursday), Day 5 (Friday), Day 6 (Saturday), and Day 7 (Sunday), with p-values < 0.05 for each comparison.
4. Day 4 (Thursday) has significantly different average delay times compared to Day 5 (Friday) and Day 6 (Saturday), with p-values < 0.05 for each comparison.
5. Day 5 (Friday) has significantly different average delay times compared to Day 6 (Saturday), with a p-value < 0.05.

These results indicate that there are significant differences in the average flight delay times between different days of the week. In particular,

the most significant differences occur between Day 6 (Saturday) and other days, suggesting that flights on Saturdays experience different average delay times compared to other days of the week.

6. Is there a significant difference in the variability of flight delay times between short-haul (length ≤ 150) and long-haul flights (length > 150)?

```
> # F-test
> dataset_clean <- dataset_clean %>%
+   mutate(FlightType = ifelse(Length <= 150, "Short-Haul", "Long-Haul"))
>
> short_haul_delays <- as.numeric(dataset_clean[dataset_clean$FlightType == "Short-Haul", "Delay"])
> long_haul_delays <- as.numeric(dataset_clean[dataset_clean$FlightType == "Long-Haul", "Delay"])
>
> short_haul_variance <- var(short_haul_delays)
> long_haul_variance <- var(long_haul_delays)
>
> f_value <- short_haul_variance / long_haul_variance
> df1 <- length(short_haul_delays) - 1
> df2 <- length(long_haul_delays) - 1
> p_value <- pf(f_value, df1, df2, lower.tail = FALSE)
>
> cat("F-value:", f_value, "\n")
F-value: 0.9967099
> cat("Degrees of freedom 1:", df1, "\n")
Degrees of freedom 1: 229541
> cat("Degrees of freedom 2:", df2, "\n")
Degrees of freedom 2: 93222
> cat("P-value:", p_value, "\n")
P-value: 0.7260549
```

Based on the F-test results, we can conclude the following:

F-value: 0.9967099
Degrees of freedom 1 (short-haul): 229541
Degrees of freedom 2 (long-haul): 93222
P-value: 0.7260549

The F-value is close to 1, suggesting that the variances of the flight delay times for short-haul and long-haul flights are quite similar. Moreover, the p-value is 0.726, which is greater than the typical significance level of 0.05. Therefore, we fail to reject the null hypothesis and conclude that there is no significant difference in the variability of flight delay times between short-haul (length ≤ 150) and long-haul flights (length > 150).

5. Conclusions and Limitations

Conclusions:

There is a slight positive correlation between flight length and the likelihood of delays, suggesting that longer flights might experience slightly more delays on average.

Certain days of the week, particularly Fridays and Sundays, are more prone to flight delays compared to other days, such as Tuesdays and Wednesdays.

The variability of flight delay times between short-haul and long-haul flights is not significantly different, indicating that flight length does not have a substantial impact on the variability of delay times.

Limitations:

The dataset used in this project might not be up-to-date, limiting the generalizability of the findings to the current state of the aviation industry.

The analyses conducted in this project do not cover all possible factors that might contribute to flight delays, such as airport congestion, weather conditions, or airline-specific factors.

The F-test conducted in this project assumes that the variances of flight delay times follow a normal distribution, which might not be the case for the data used. This assumption should be verified before drawing conclusions from the F-test results.

The post-hoc tests conducted in this project might be sensitive to multiple testing issues, as many pairwise comparisons were made. This can potentially lead to inflated false-positive rates.

Future research should consider addressing these limitations by using more recent data, incorporating additional variables, and employing more robust statistical methods to better understand the factors contributing to flight delays and help airlines, airport authorities, and passengers make informed decisions.