

DATA 5300: Analysis of United Airlines Gain Per Flight

Team Member: Vincent Chan, Muykhim Ing, Jimmy Nam, Alex Song

Introduction

Flight performance is a key factor affecting operational efficiency and customer satisfaction in the airline industry. For United Airlines, understanding how flights recover time in the air known as gain per flight, can provide insights into scheduling efficiency and help improve the passenger experience. Gain is calculated as the difference between a flight's departure delay and arrival delay with positive values indicating that the flight made up time during the journey.

This project uses data from *nycflights13* dataset, which contains detailed records of all flights departing from New York City in 2013. Our analysis focuses specifically on United Airlines flights and examines how gain relates to factors such as departure delays, flight destinations and flight duration. By combining summary statistics, visualisations, confidence intervals and hypothesis tests, this report provides a clear, data driven overview of flight efficiency. The results are presented in a way that is accessible to a non technical audience, offering insights that could support better operational planning and an improved passenger experience.

Methodology

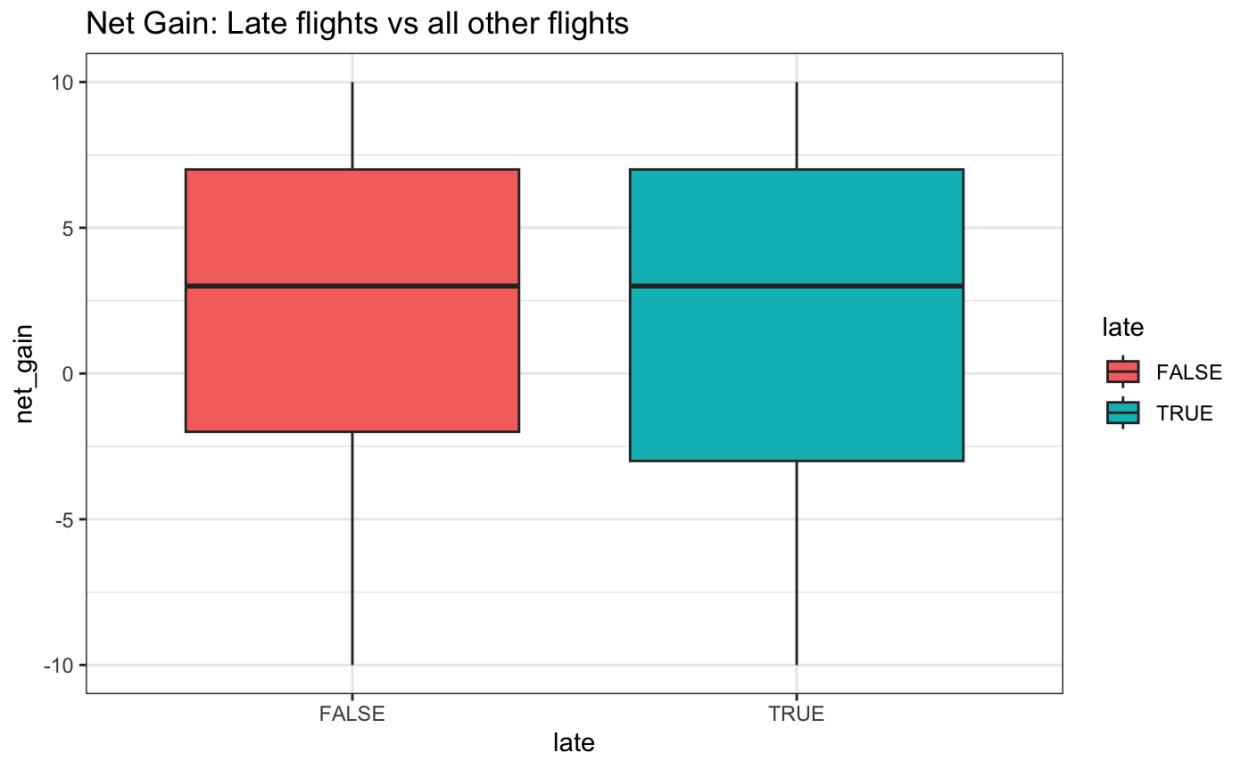
The analysis used the *nycflights13* dataset, which contains detailed information on flights departing from New York City airports in 2013. To focus on United Airlines only, we filtered only flights operated by UA from the *flights* dataset. Because the analysis depends on comparing departure and arrival delays, flights with missing values for either `dep_delay` or `arr_delay` were removed.

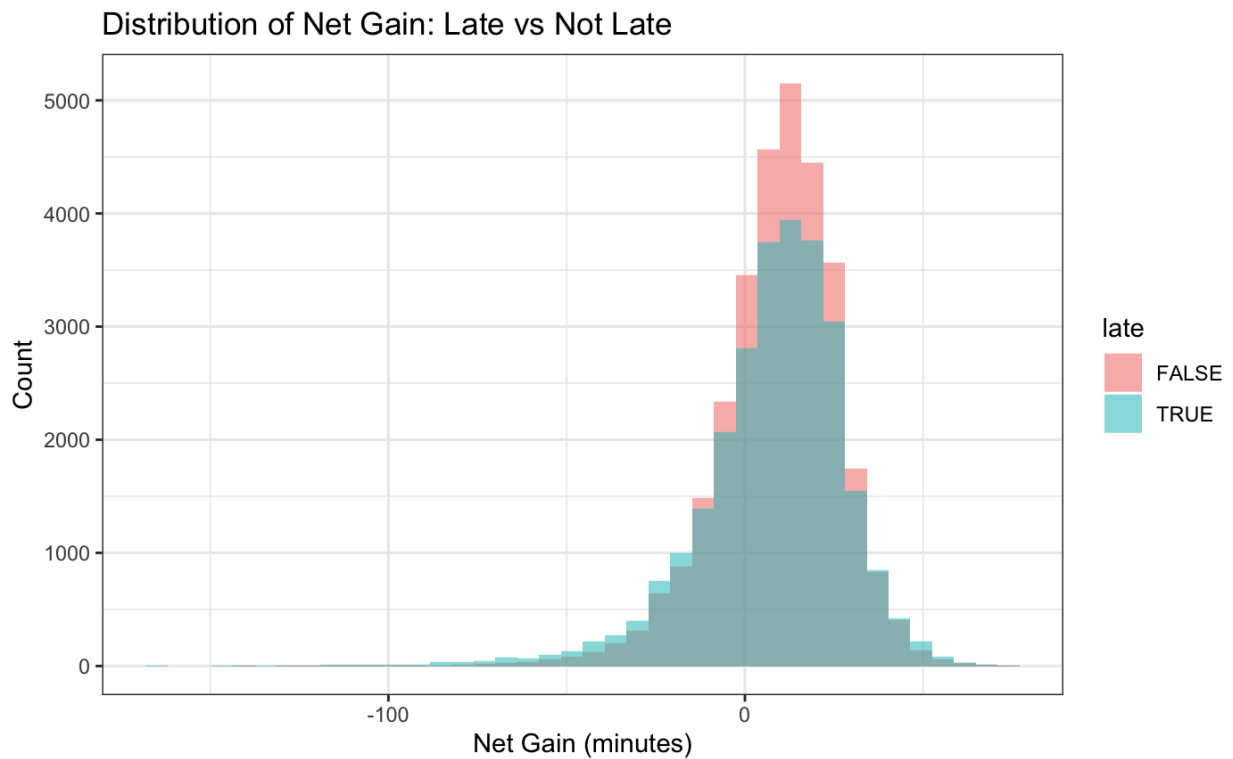
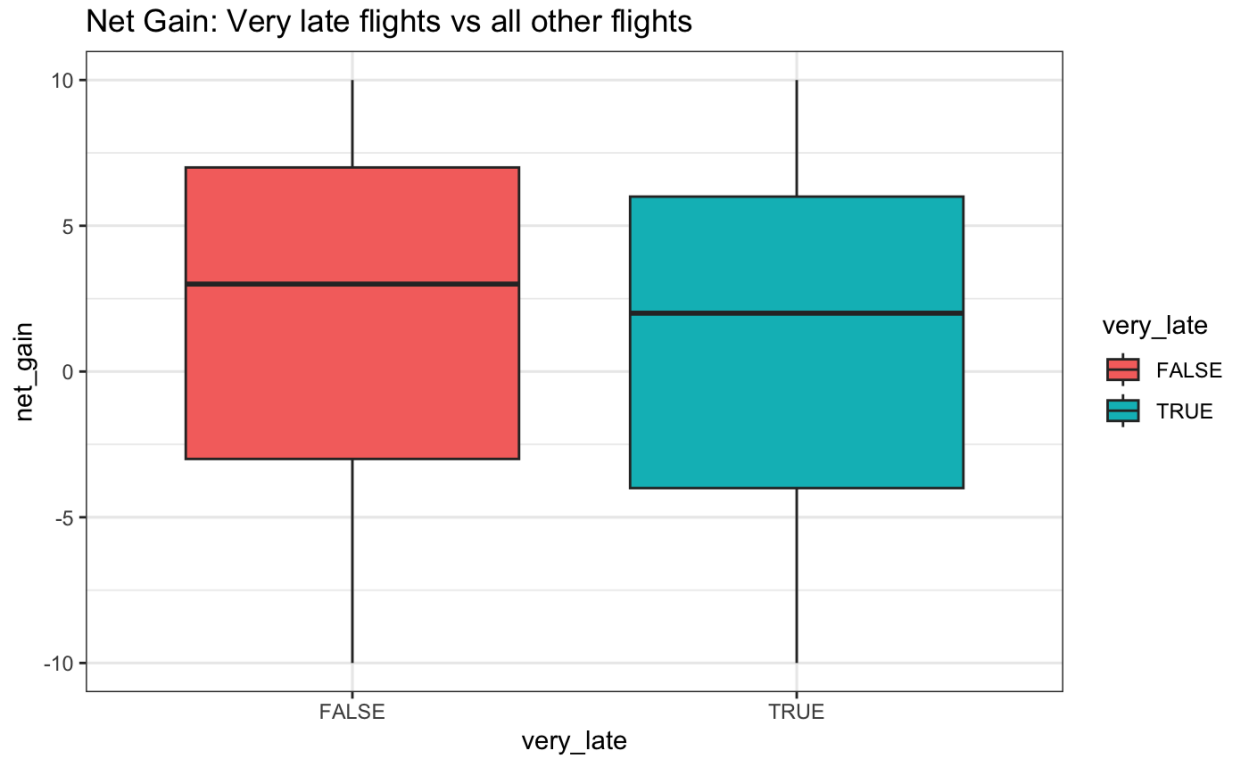
A new variable, net gain, was created to measure how much time each flight made up in the air. This was calculated as the departure delay minus the arrival delay, so positive values indicate time recovered during the flight. To support the comparison questions in the project, several indicator variables were added: `not_late` for flights that departed on time or early, `late` for flights with a departure delay greater than zero, and `very_late` for flights departing more than 30 minutes late.

After these transformations, the cleaned dataset contained only UA flights with complete delay information and the additional variables needed for comparing groups. All exploratory plots, summary statistics, and confidence interval analyses were based on this cleaned dataset.

Result

1. Does the average gain differ for flights that departed late versus those that did not?
What about for flights that departed more than 30 minutes late?





The boxplot comparing late and not-late flights shows that the two distributions are nearly identical. The medians line up, and both groups include flights that gained time and flights that lost time. The only slight visual difference is that late flights have a bit more variability on the

“lost time” side, but the overall shape and spread remain extremely similar. The second boxplot, which compares flights that departed more than 30 minutes late to all other UA flights, shows a modest downward shift in the distribution for very late flights. These flights tended to have slightly lower net gains overall, though the two groups still overlap substantially.

The overlaid histogram reinforces the lack of a meaningful difference. The late and not-late curves sit directly on top of each other, with the same general shape, center, and tail behavior. The only visible distinction is the height of the bars, which simply reflects the different number of flights in each category. If late flights consistently made up more time in the air, the histogram would show a clear rightward shift for the late-flight distribution, but that shift does not appear.

The numerical results support what the graphs showed. For late versus not-late flights, the observed difference in mean net gain was about -1.73 minutes, indicating that late flights made up slightly less time on average. The 95 percent bootstrap confidence interval, based on 10,000 resamples, ranged from about -0.34 to 0.34 minutes. Since the interval is centered on zero and includes both positive and negative values, the difference is too small to be meaningful. The comparison involving very late flights led to the same outcome. Flights that departed more than 30 minutes late gained about -2.41 fewer minutes than the rest, and the corresponding 95 percent bootstrap confidence interval (roughly -0.51 to 0.53 minutes) also included zero. In both cases, the data provide no evidence that later departures result in more time gained during the flight.

2. What are the five most common destination airports for United Airlines flights from New York City? Describe the distribution and the average gain for each of these five airports.

Most Common United Airlines Destinations

The five most frequent destinations, based on flight counts, are:

1. Chicago O'Hare (ORD) - 17283 flights
2. Atlanta (ATL) - 17215 flights
3. Los Angeles (LAX) - 16174 flights
4. Boston (BOS) - 15508 flights
5. Orlando (MCO) - 14082 flights

These routes represent the highest traffic flights for United Airlines departing NYC.

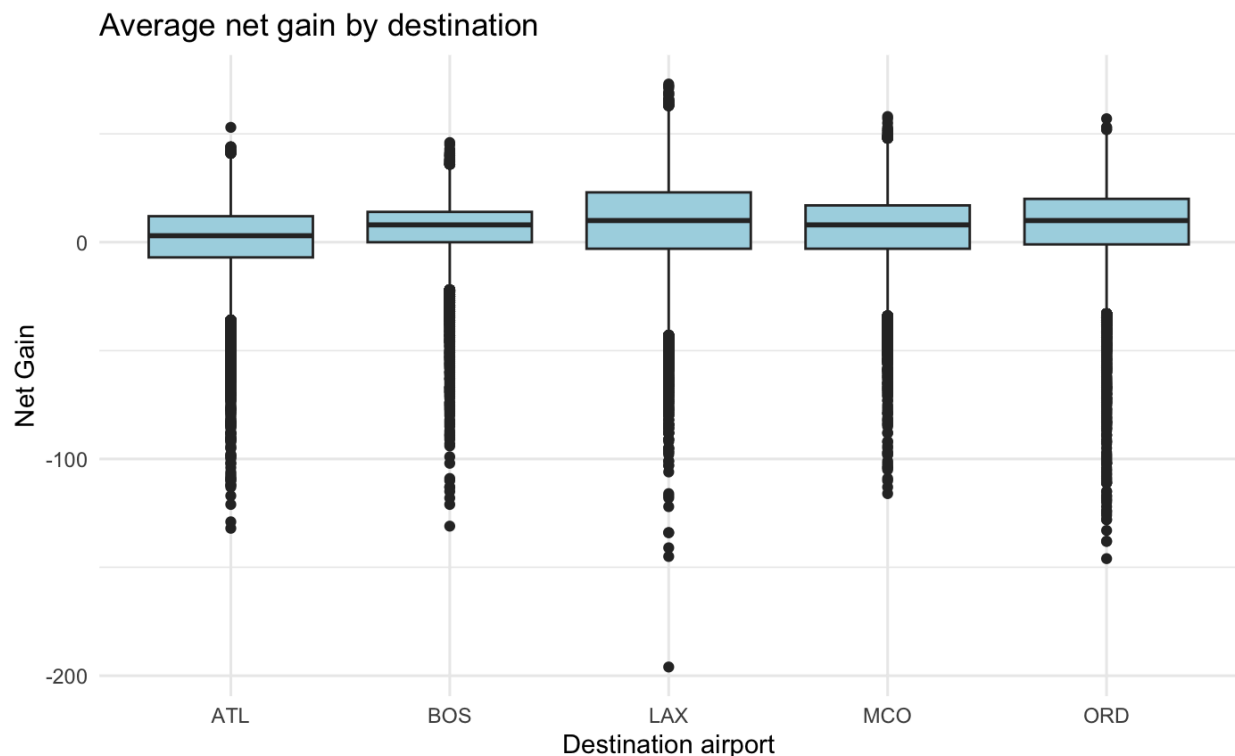
Average Net Gain by Destination

Net gain is calculated as departure delay minus arrival delay, with positive values indicating time recovered during flight. The summary statistics are shown below:

Destination	Flights	Avg Gain	Median Gain	SD
ORD (Chicago)	17283	7.6	10	19.1
ATL (Atlanta)	17215	1.1	3	16.6
LAX (Las Angeles)	16174	8.8	10	21.6
BOS (Boston)	15508	5.8	8	13.9
MCO (Orlando)	14082	5.8	8	16.8

LAX and ORD have the highest average gains (around 8-9 minutes), suggesting that longer flights allow more opportunity to recover lost time. BOS and MCO show moderate gains (around 6 minutes), while ATL has the smallest gain, reflecting limited recovery on this shorter or busier route.

Distribution of Net Gain (Boxplot Analysis)



The boxplot highlights the distribution of net gain for each destination:

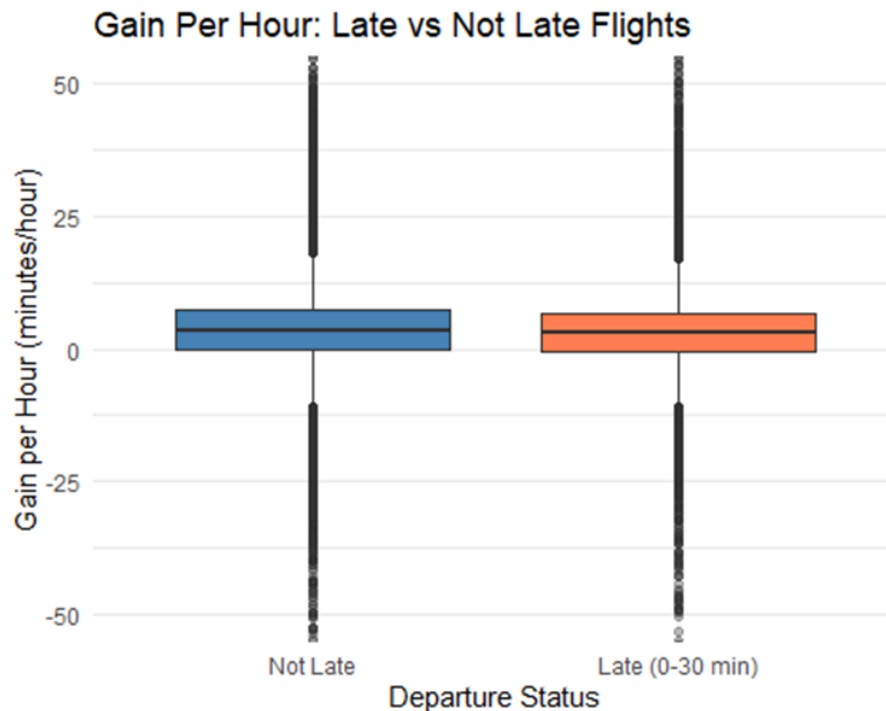
- ATL: has the lowest median and the tightest distribution near zero, visually confirming minimal ability to recover lost time.

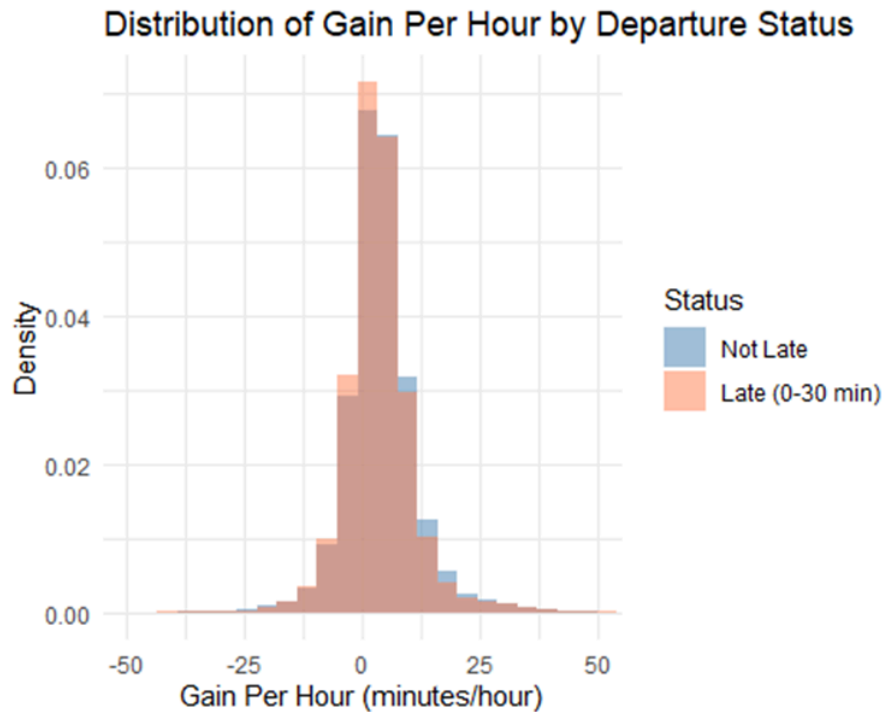
- BOS and MCO: show Moderate, consistent gains with smaller spreads, suggesting more predictable outcomes.
- LAX and ORD: highest median and widest ranges, indicating the greatest potential for time recovery, but also the most variability in flight performance.

3. Another common measure of interest, in addition to total gain, is the gain relative to the duration of the flight. Calculate the gain per hour by dividing the total gain by the duration in hours of each flight. Does the average gain per hour differ for flights that departed late versus those that did not? What about for flights that departed more than 30 minutes late?

The following analysis uses data from United Airlines flights along with weather information to examine how departure delays relate to a flight's "Gain Per Hour".

I. Gain Per Hour Analysis: Late (0–30 minute delay) vs. Not Late Flights

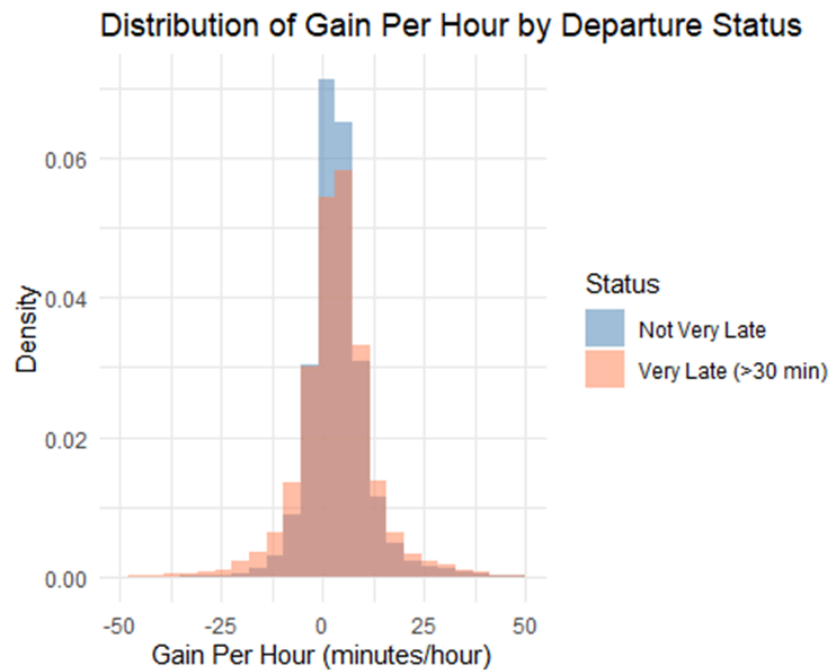
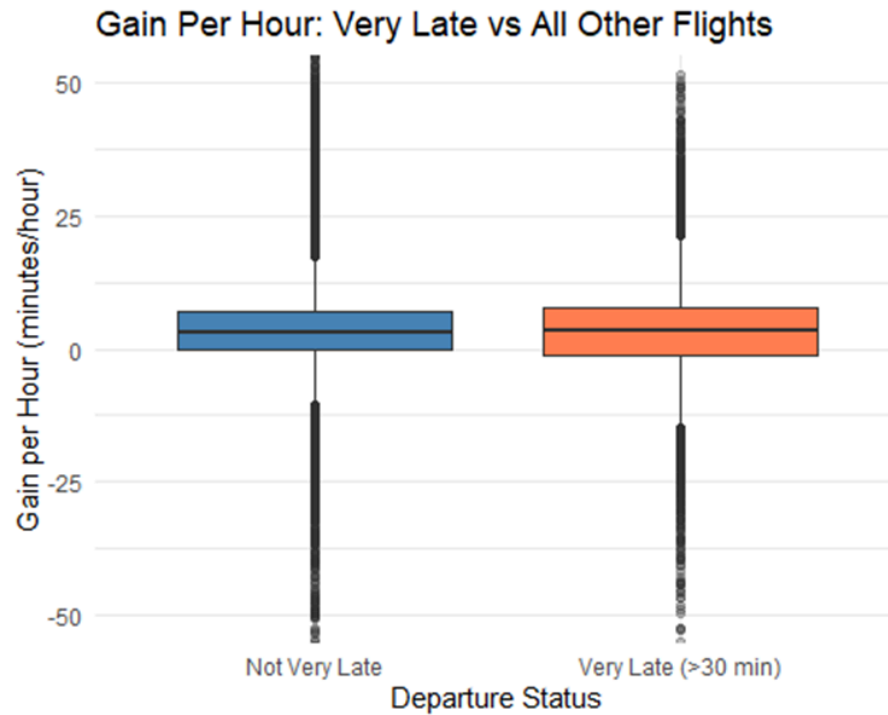




The first comparison grouped flights into “Late” (departing 0 to 30 minutes after schedule) and “Not Late.” Not Late flights formed the larger group (N = 38,035) and showed a mean gain per hour of 3.80 minutes per hour (Median: 3.33, SD: 8.46). Late flights (N = 19,454) had a slightly lower mean gain per hour of 3.23 minutes per hour (Median: 2.91, SD: 8.53).

The observed difference in mean gain per hour (Late minus Not Late) was -0.57 minutes per hour. The 95 percent Bootstrap Confidence Interval for this difference was [-0.72, -0.43] minutes per hour. Since the entire interval is below zero, the results show that flights with a 0–30 minute departure delay tend to have a statistically significant lower mean rate of time recovery during flight compared to flights that departed on time or early. Boxplots of the two groups confirm that the distribution for Not Late flights is shifted slightly higher, although both groups show a similar wide spread of values.

II. Gain Per Hour Analysis: Very Late (more than 30 minute delay) vs. All Other Flights



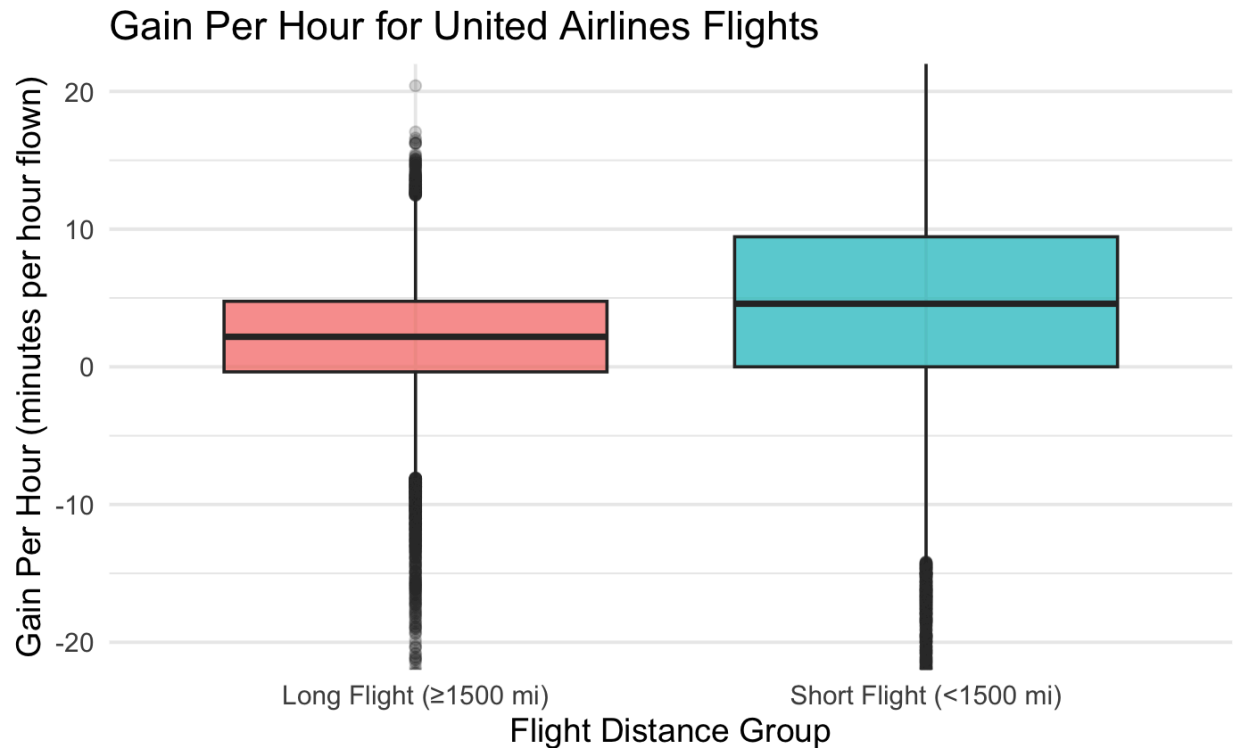
The next comparison focused on flights that were “Very Late” (departing more than 30 minutes behind schedule) versus all other flights, labeled “Not Very Late.” The Not Very Late group (N = 49,987) had a mean gain per hour of 3.69 minutes per hour (Median: 3.18, SD: 8.06). The Very

Late group (N = 7,502) had a lower mean gain per hour of 3.04 minutes per hour (Median: 3.33, SD: 10.9), along with a notably higher standard deviation, indicating more variability.

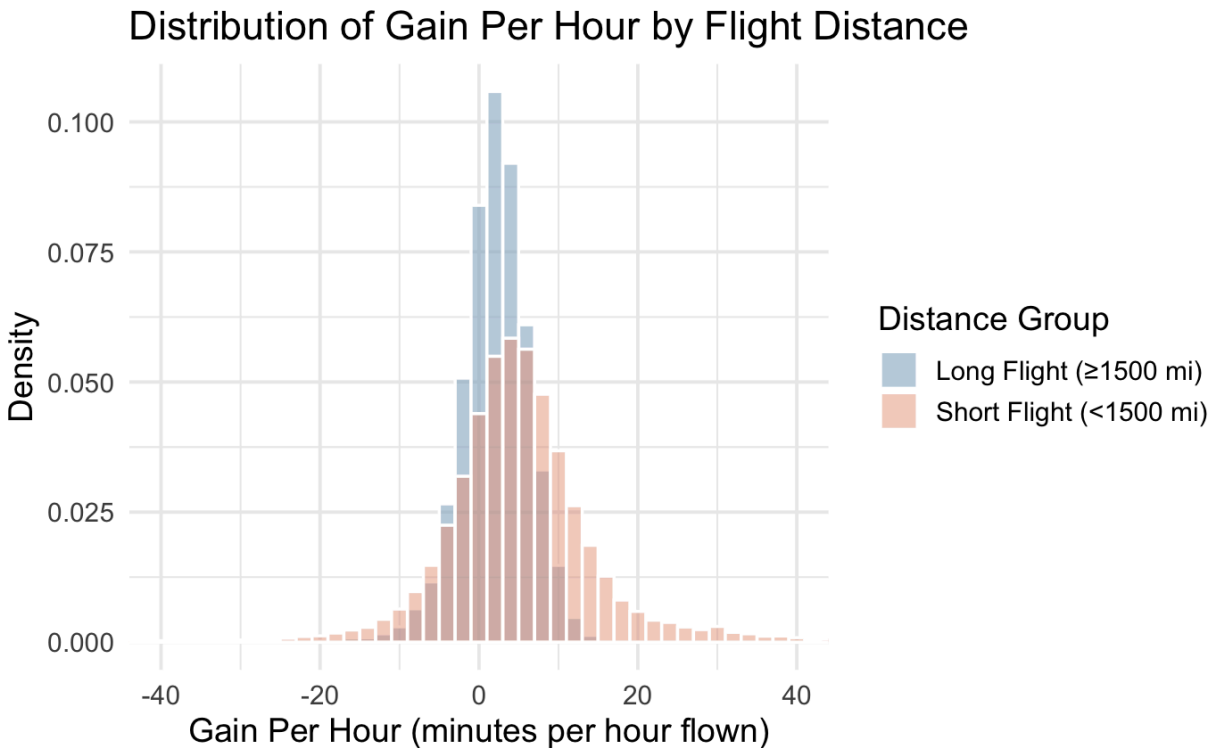
The observed difference in mean gain per hour (Very Late minus Not Very Late) was -0.66 minutes per hour. The 95 percent Bootstrap Confidence Interval was [-0.91, -0.40] minutes per hour. Because this interval is fully below zero, the data show that flights delayed more than 30 minutes are associated with a significantly reduced mean gain per hour when compared to less delayed flights. This means the more a flight is delayed at departure, the lower its ability to make up time during the flight. The effect is slightly stronger for the Very Late group (-0.66 minutes per hour) than for the Late group (-0.57 minutes per hour).

4. Does the average gain per hour differ for longer flights versus shorter flights?

To investigate whether longer flights gain more time per hour than short flights, I compared the average gain per hour across two distance groups: the long flights, which are greater than or equal to 1500 miles, and the short flights, which are less than 1500 miles. Long flights recovered an average of 2.02 minutes per hour, with a 95% confidence interval ranging from -7.35 to 9.95 minutes per hour. On the other hand, short flights gained an average of 4.88 minutes per hour, though their confidence interval was much wider, from -14.47 to 28.64 minutes per hour. The width and overlap of these intervals indicate substantial variability in both groups and provide little statistical evidence that either type of flight reliably recovers more time per hour. While the point estimate for short flights is higher, the imprecision reflected in the confidence interval means this difference should be interpreted with caution. Overall, the results do not strongly support the idea that flight distance has a meaningful impact on the amount of time a plane gains per hour while airborne.



This boxplot compares the gain per hour for United Airlines flights and shows that short flights (<1500 miles) tend to recover more time per hour in the air than long flights (>1500 miles). The median gain per hour is visibly higher for short flights, and their interquartile range is shifted upward, indicating that a typical short flight makes up more time while airborne. Long flights, by contrast, show a lower median gain and a tighter spread around zero, suggesting they recover less time per hour on average. Both groups contain extreme outliers, particularly negative values for time lost, but the overall pattern suggests that shorter flights generally gain more time per hour than longer ones.



The histogram comparing the gain per hour for long and short flights shows that both distributions are centered close to zero, indicating that most flights gain only a small amount of time per hour in the air. However, short flights display a wider spread, with more observations extending into both positive and negative gain values. This suggests greater variability in how much time shorter flights make up or lose while airborne. Long flights have a narrower distribution that is more tightly concentrated around zero, indicating more consistent performance with fewer extreme gains or losses. Although short flights reach higher positive gains per hour, the large overlap between the two distributions suggests that the overall pattern of time recovery is similar across both distance groups.

Discussion

We analyzed United Airlines flights using the nycflights13 dataset to see whether planes can make up for departure delays while in the air, and under what conditions that recovery might happen. Overall, there was little difference in total time gained between flights that left late and those that departed on schedule. But when we looked at time gained per hour, a clear pattern emerged: flights that departed late or very late actually made up less time per hour than flights that left on time. In other words, leaving late doesn't mean flying faster to catch up.

This trend held true regardless of distance. Neither short nor long flights showed a significant advantage, even though their median times varied. Destination, however, did matter. Flights to

Chicago (ORD) and Los Angeles (LAX) were most likely to recover time, while flights to Atlanta (ATL) were least likely. This suggests that certain routes may face operational constraints that limit recovery.

In short, while flights can sometimes make up time in the air, it's rare and only happens under very specific circumstances. This analysis also had limitations. We couldn't account for factors that likely influence whether a flight can speed up—such as airspace congestion, airline policies, and weather conditions. Future research should go beyond United Airlines and New York departures. It should focus on high-traffic routes and include variables like weather and airspace congestion to better understand when and why time recovery occurs.