# Feature Engineering
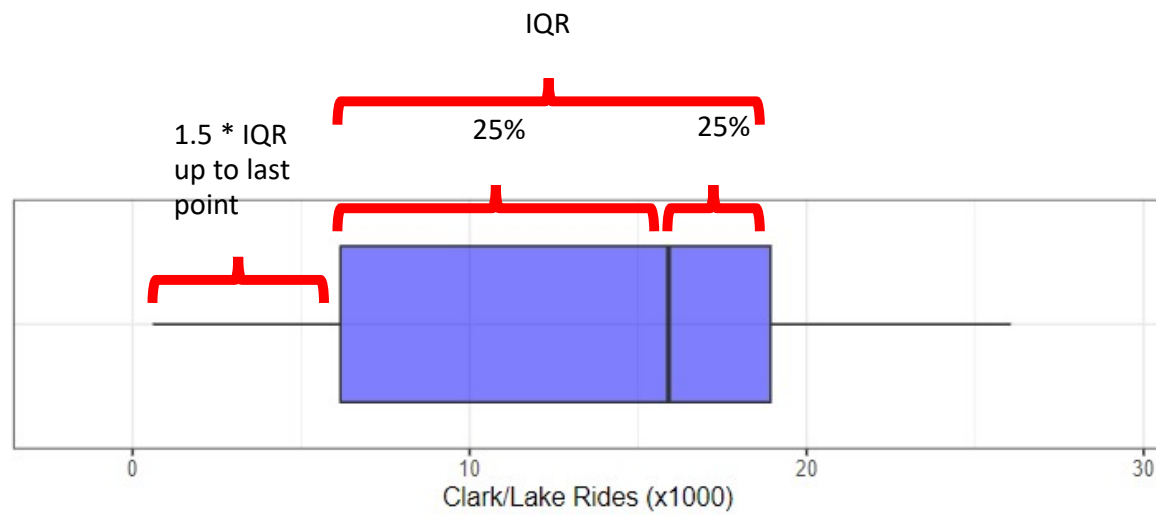
Chapter 4
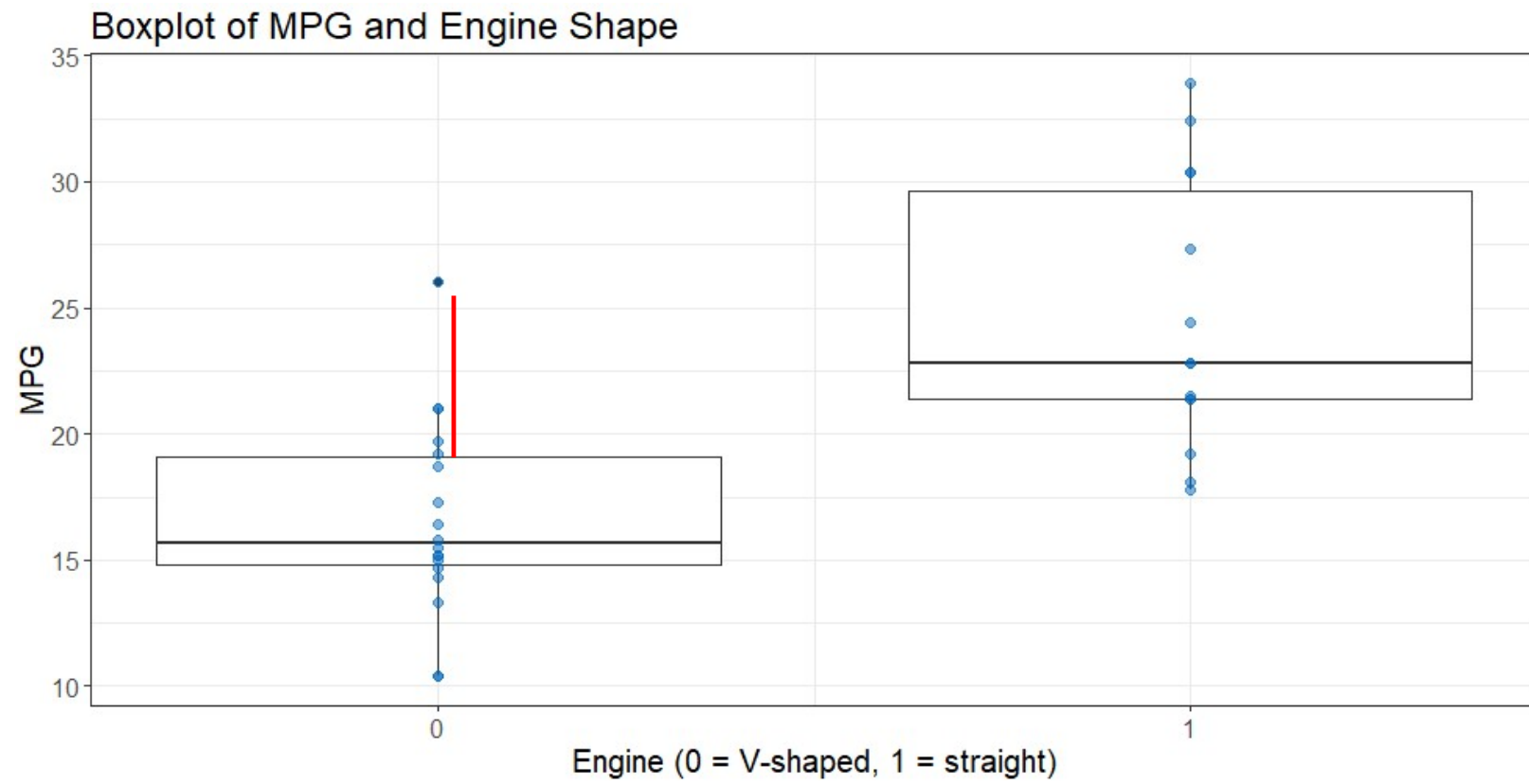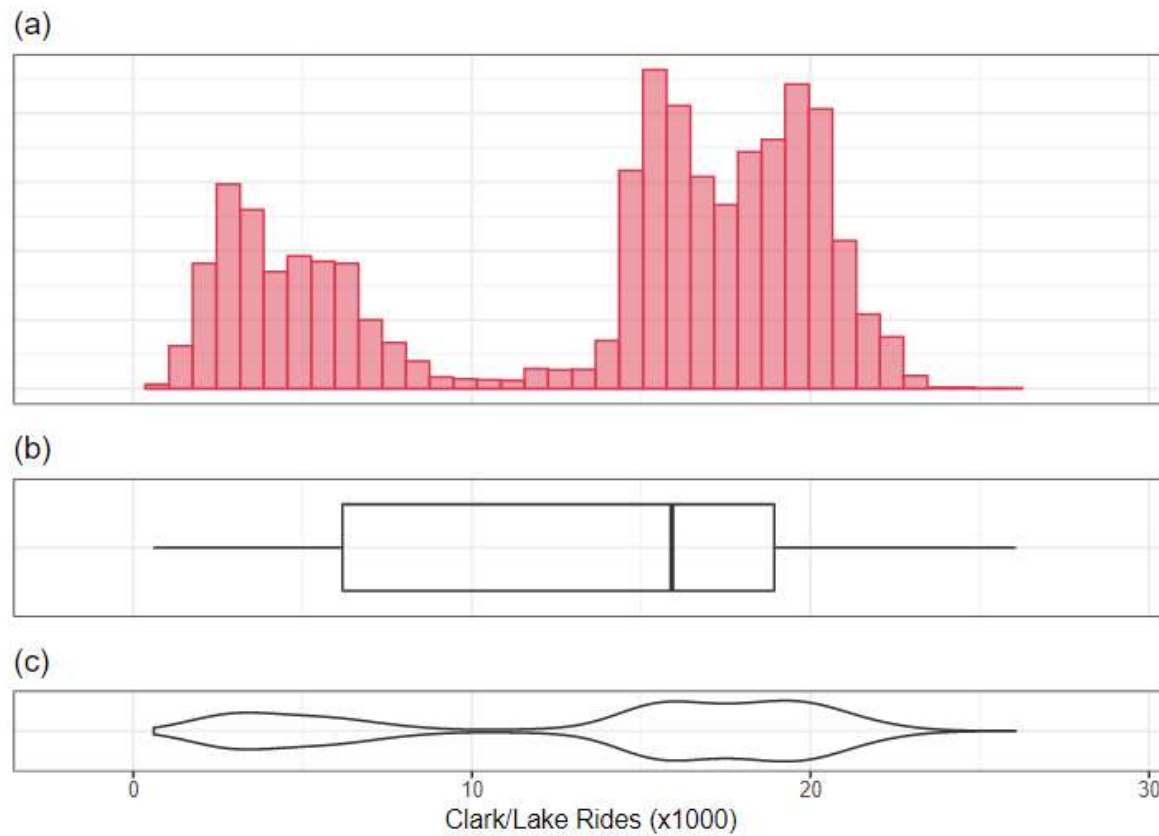
Stephen Kimel

# Chicago Train Ridership

- "Two common mistakes of misunderstanding demand could be made. At one extreme, having too few cars on a line to meet weekday demand would delay riders from reaching their destination and would lead to overcrowding and tension. At the other extreme, having too many cars on the weekend would be inefficient leading to higher operational costs and lower profitability. Good forecasts of demand would help the CTA to get closer to optimally meeting demand."

- "Our illustration will narrow to predicting daily ridership at the Clark/Lake stop."
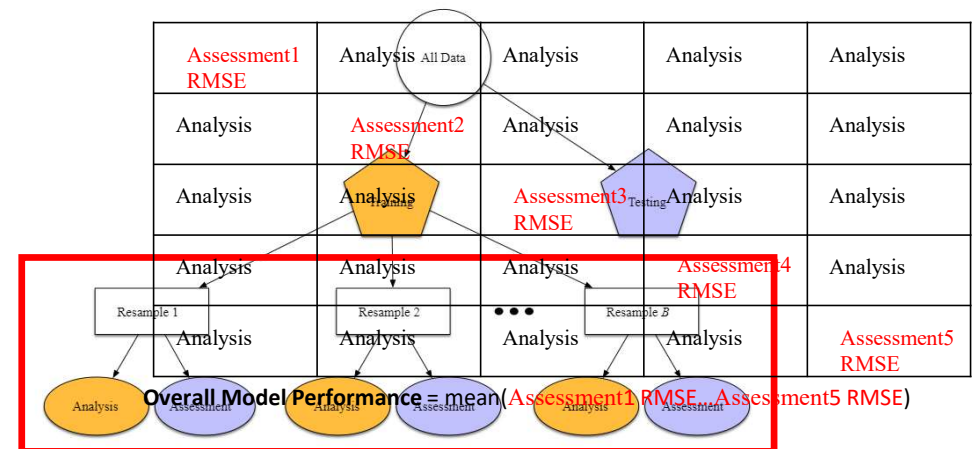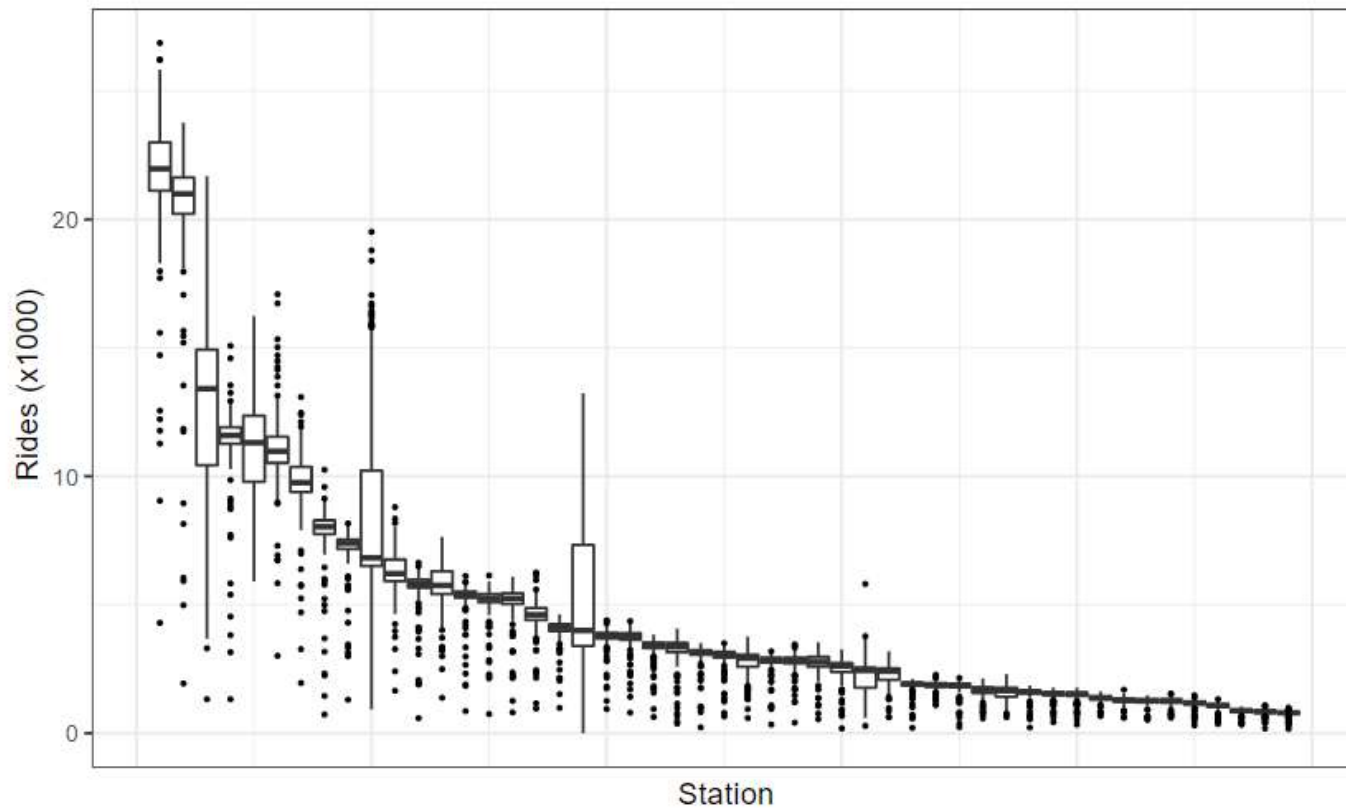
# Boxplot

Boxplot of MPG and Engine Shape

# Histogram vs. Bloxplot vs. Violinplot



(a)

(b)
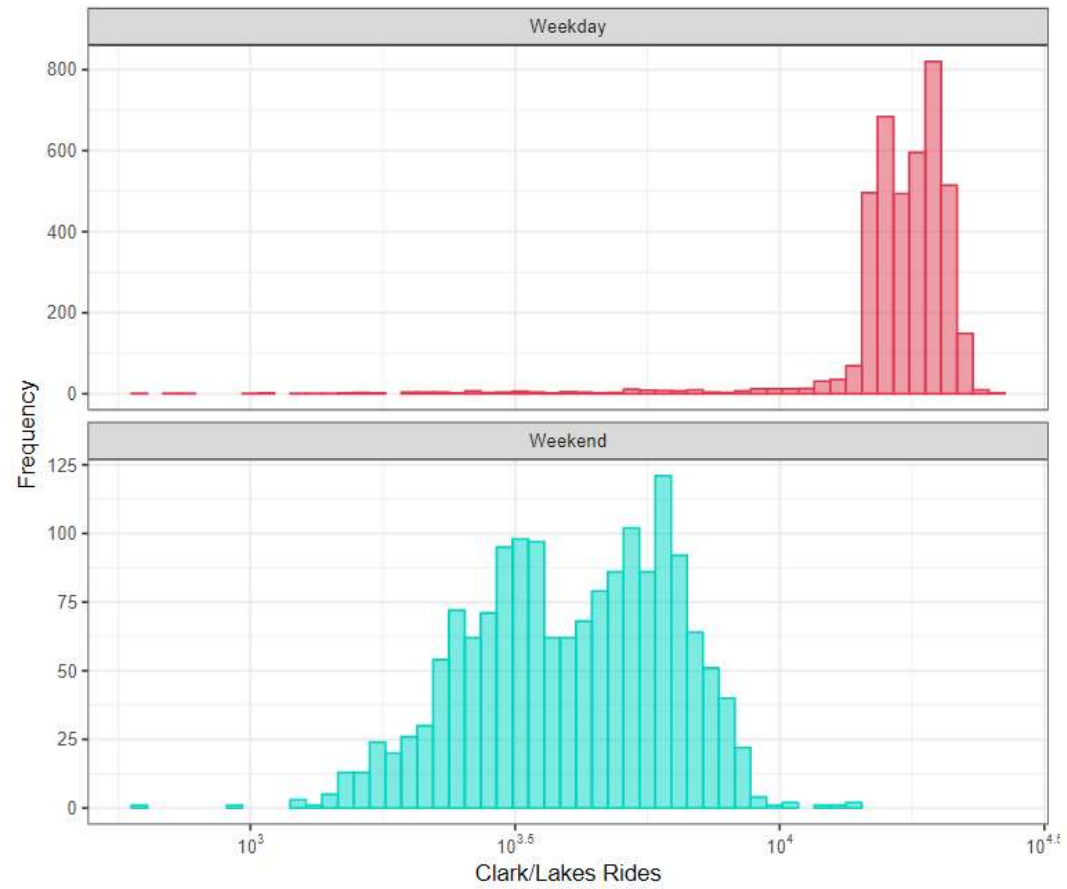
(c)

Clark/Lake Rides (x1000)

- "Given the range of the daily ridership numbers, there was some question as to whether the **outcome** should be modeled in **the natural units or on the log scale**. On one hand, the **natural units makes interpretation of the results easier** since the RMSE would be in terms of riders. However, if the **outcome were transformed** prior to modeling, it would ensure that *negative* ridership **could not be predicted**. The **bimodal nature** of these data, as well as distributions of ridership for each year that have **a longer tail on the right** made this decision difficult. In the end, a **handful of models were fit both ways** to make the determination. The **models** computed in the **natural units** appeared to have slightly **better performance** and, for this reason, all models were analyzed in the natural units."
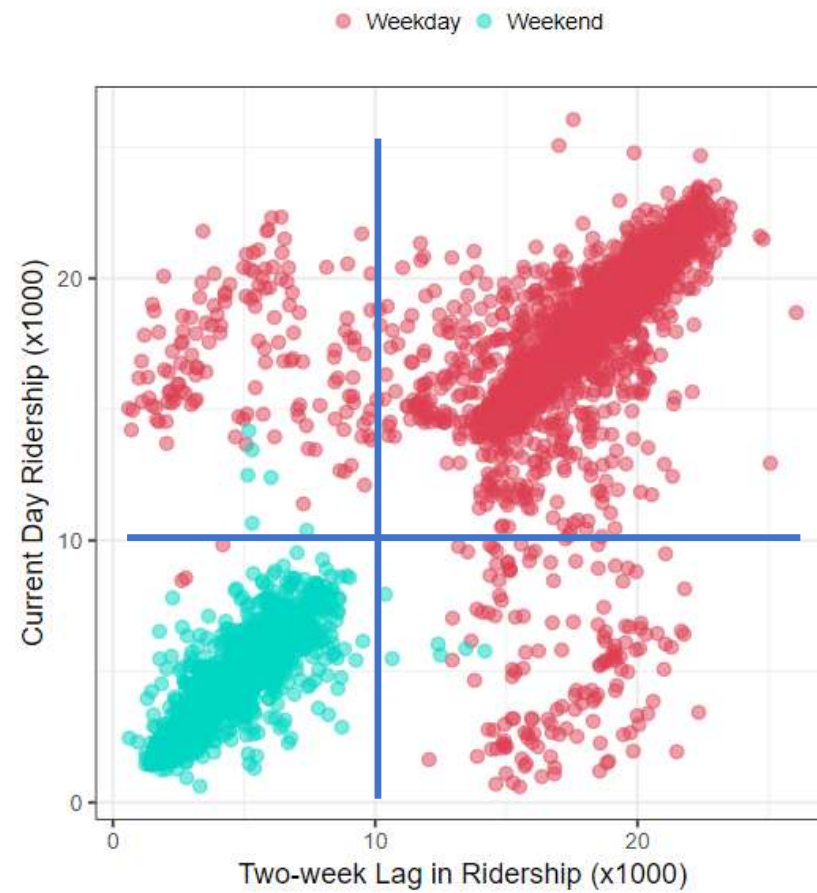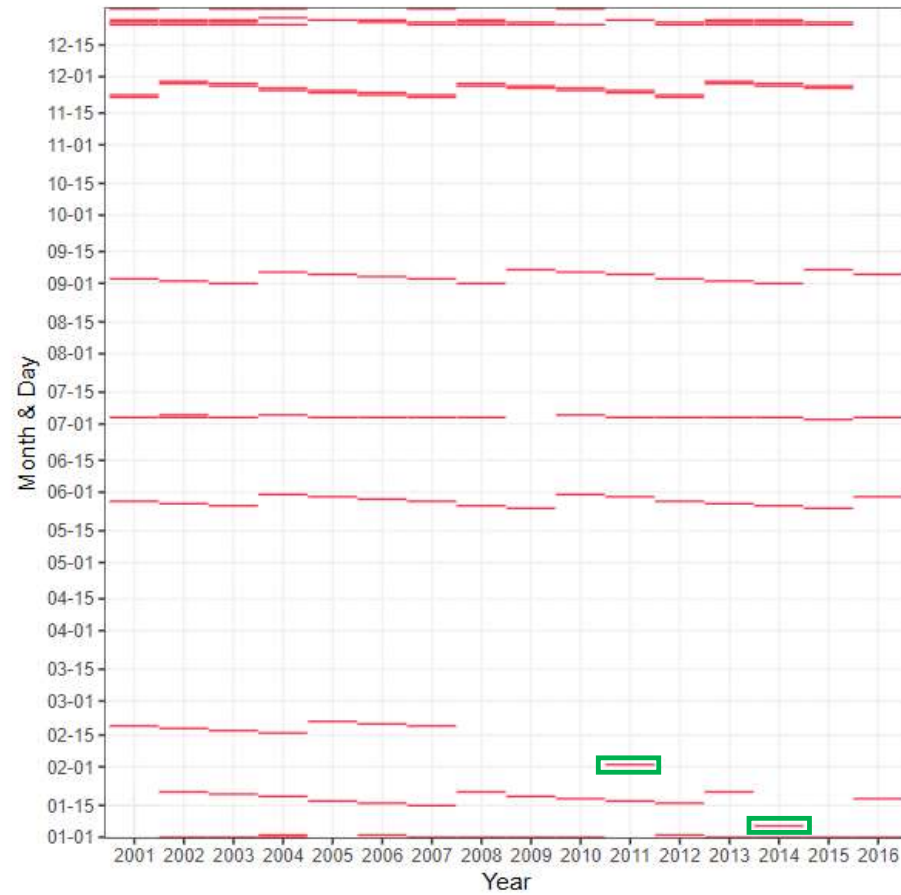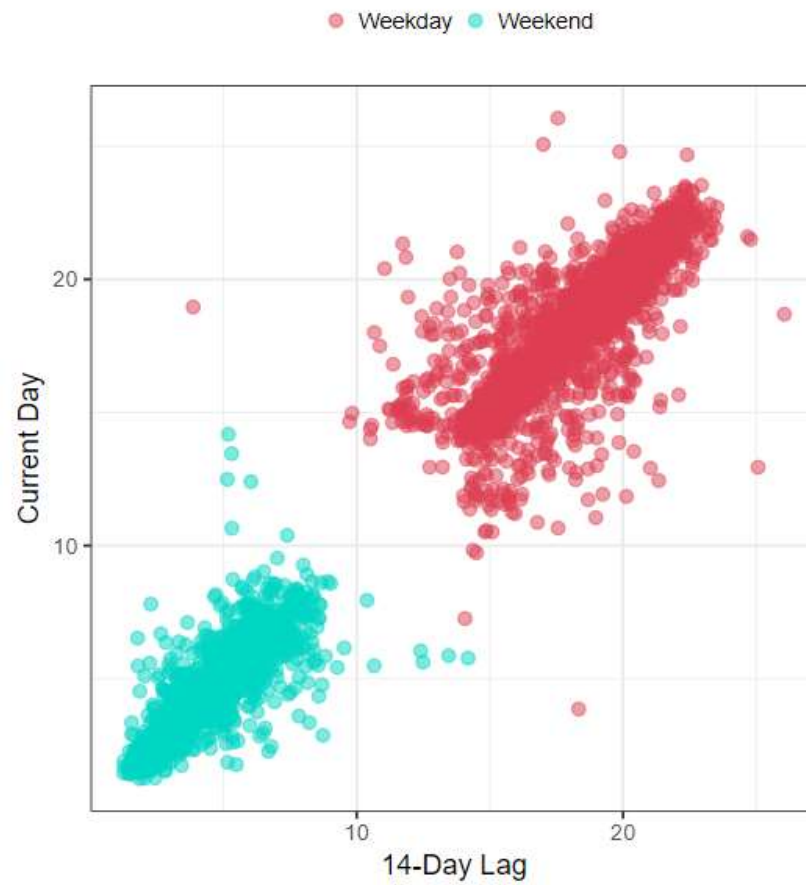
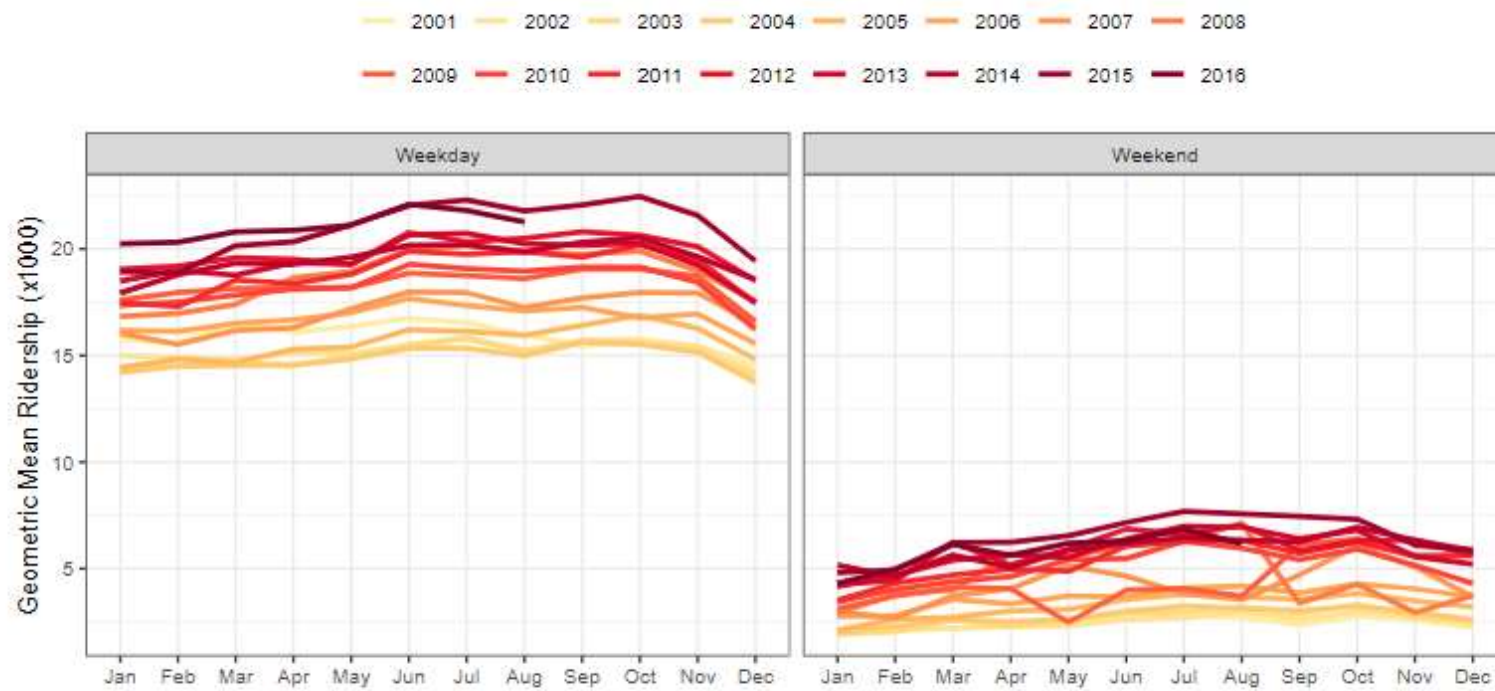# Observing Predictor Variables with Ridership

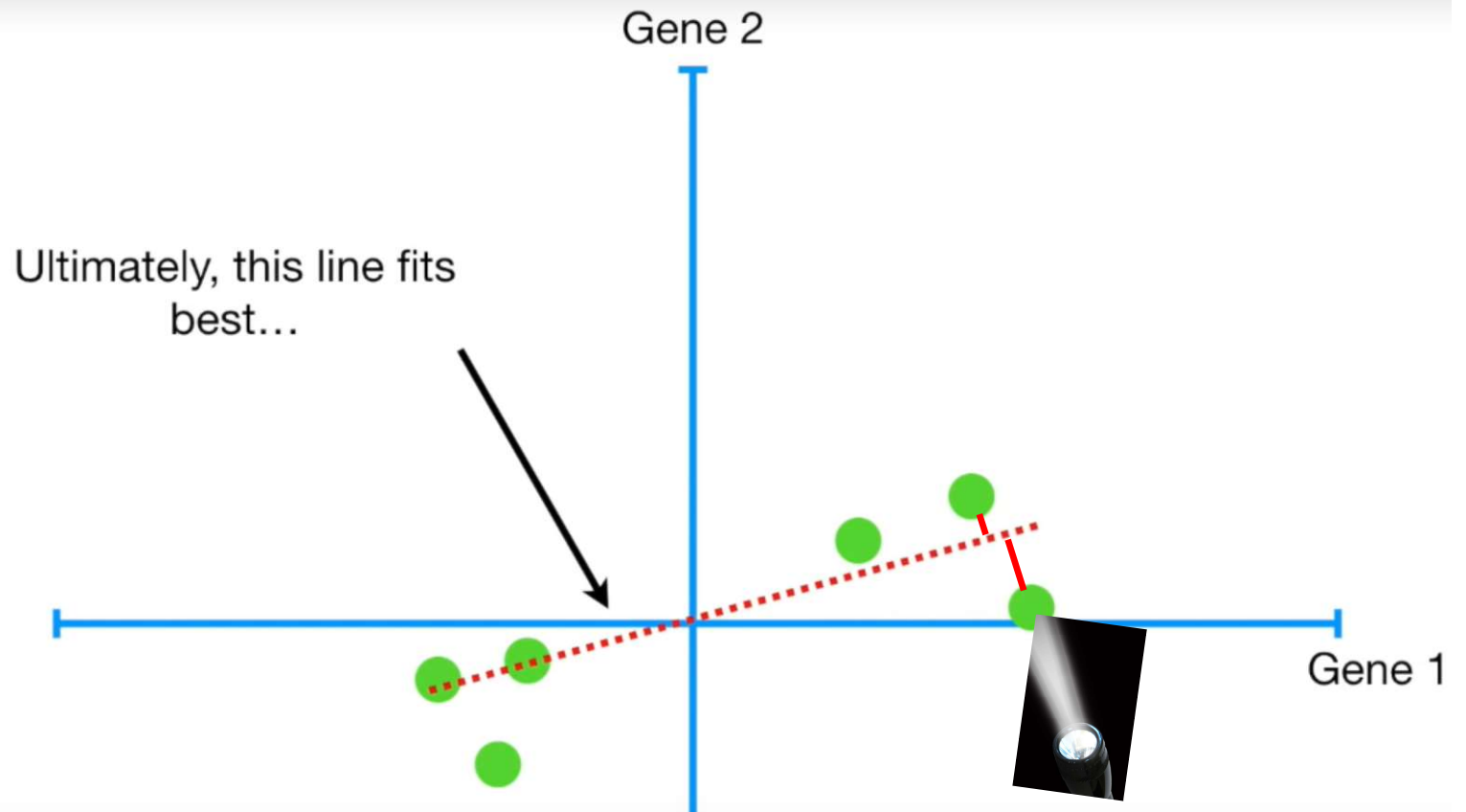# Faceting

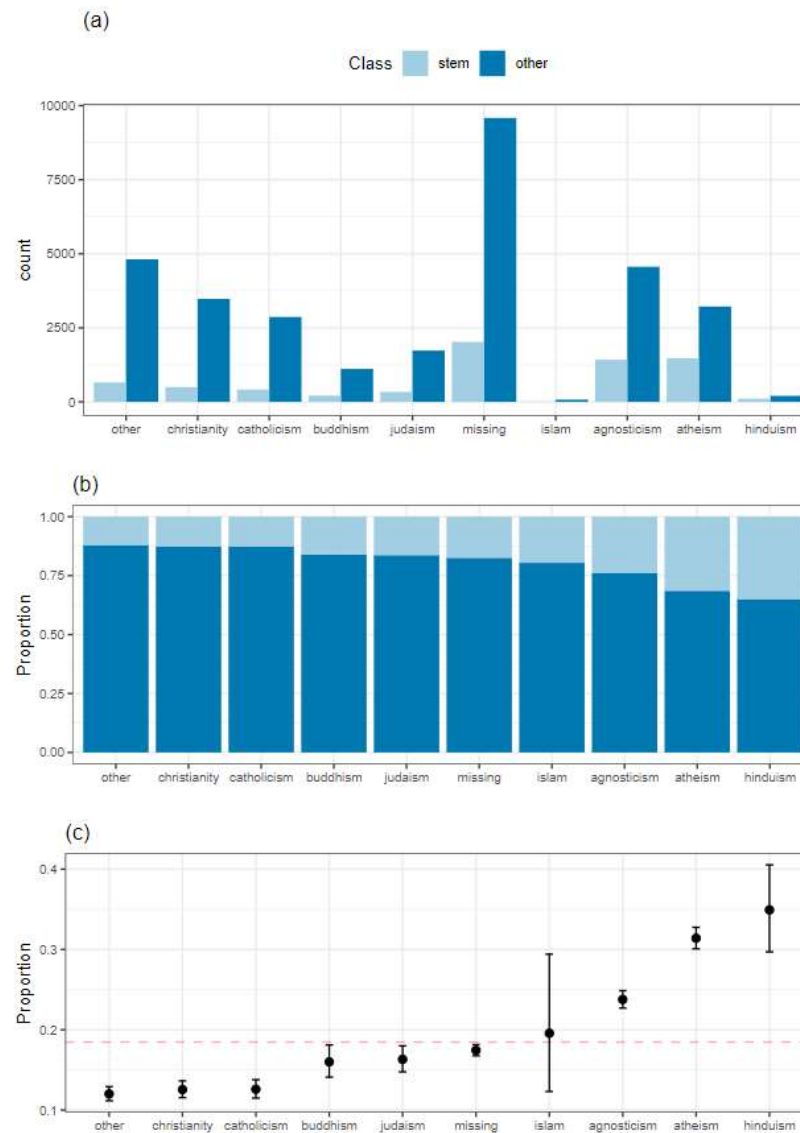# Scatterplots

# Heatmap

# After Some Feature Engineering

# Line Charts

# PCA

Goal: Retain the most information while reducing the number of dimensions
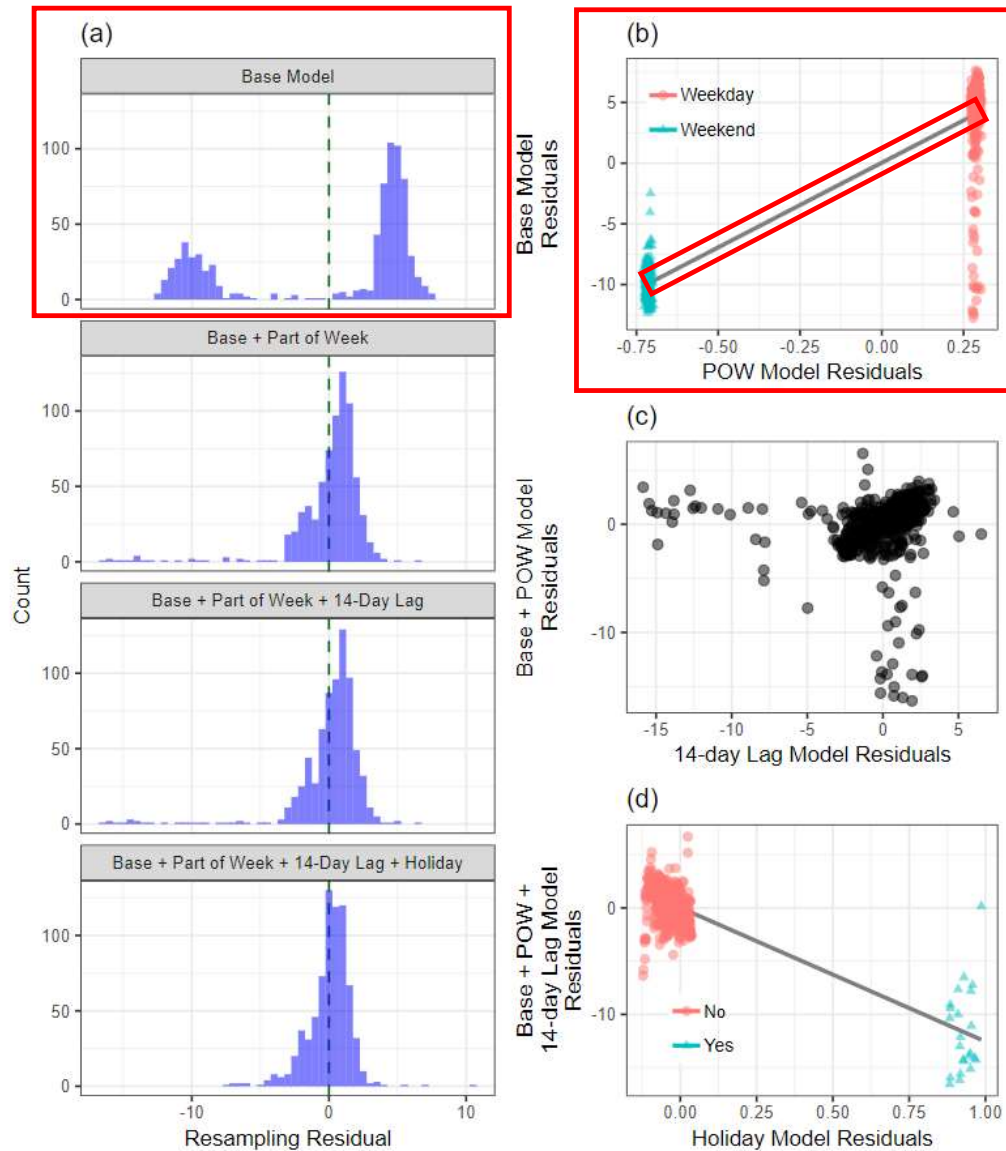


13

# Are Proportions Different?

# Post-Modeling Charts

## Partial Regression Plots

Base Model is:

ridership = month + year + week

weekday = month + year + week