# Curbside Database Statistical Analysis

Tyler Campbell

2025-07-21

## Introduction

Welcome to my next project, a statistical analysis of my database using R! In this analysis, I will show basic visualizations and summary statistics of different entities in my database. I will also run a couple of regressions to test how different entities significantly impact each other. I hope that you enjoy!
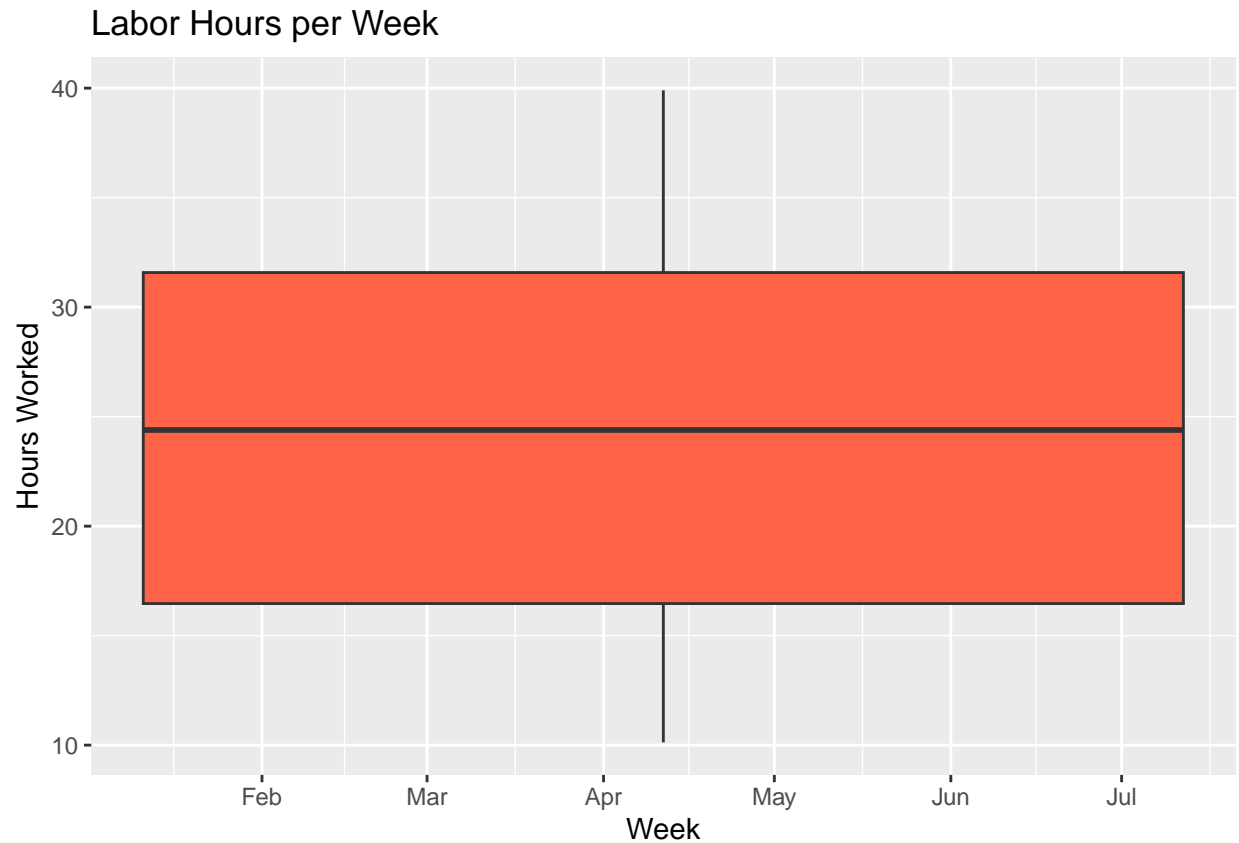
## Basic Visuals and Summary Stats

Table 1: Summary Statistics from Orders Data

| avg_order_price | min_order_price | max_order_price | avg_wages | total_labor_hours |
|---|---|---|---|---|
| 101.76 | 15.4 | 198.29 | 17.12 | 682.74 |

This is just some basic summary statistics. Everything seems prettu reasonable except for the min and max order prices. Realistically, there will be order prices that are more than 200 dollars and less than 15 dollars, but this is just synthetic data.
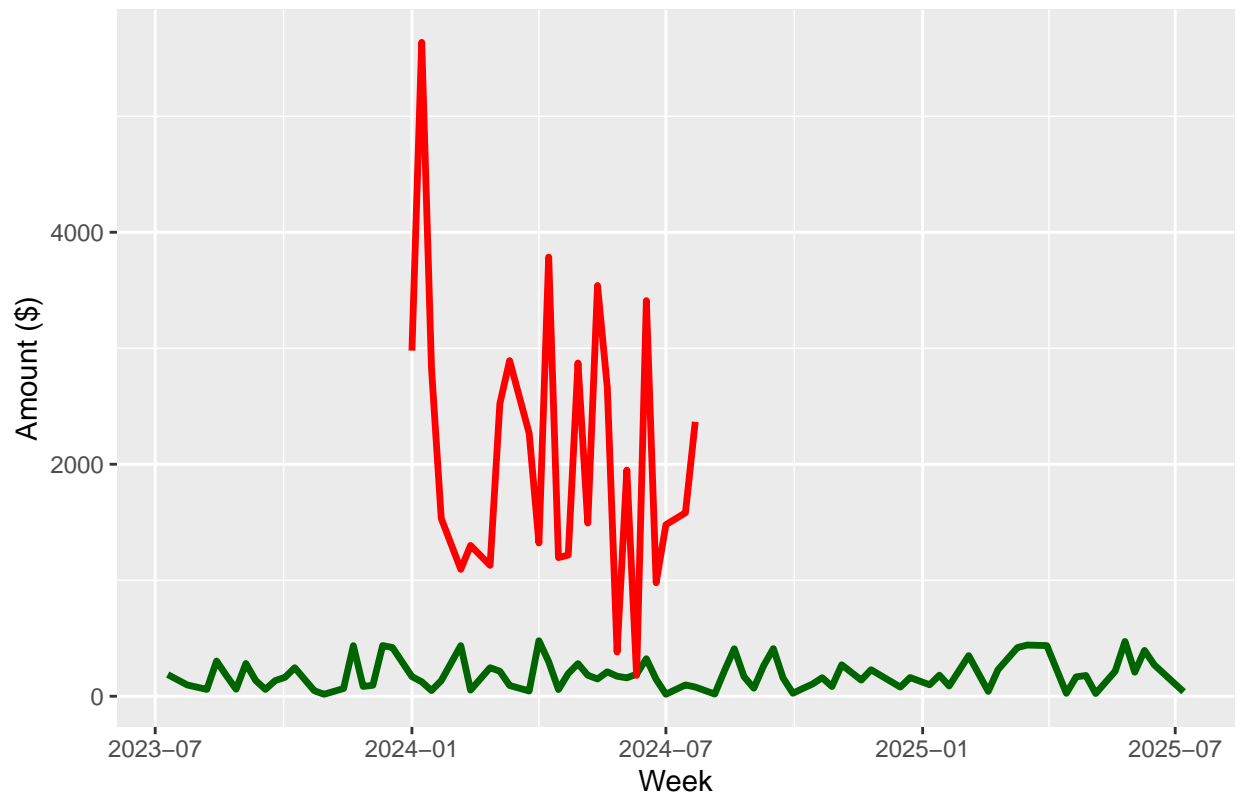
## Distribution of Order Prices



This is a bar graph of the price distribution. The majority of the orders are in the 25 and 60 dollar price range. This follows more of a uniform distribution and is not normally distributed. My intuition would think that it would either be more skewed to the right because less people probably purchase highly-priced orders.

## Labor Hours per Week



This is a boxplot that shows the number of hours worked by partners during different weeks. It stays the same throughout the year with the majority of hours worked being between 17 and 33 primarily. Is this realistic? Probably not because the number of hours worked might change during different times of the year. For example. At the H-E-B that I work at in Lubbock, a lot of the partners are college students, so they will go home during the summer which causes them to work less.

Revenue vs Labor Cost Over Time

This shows the comparison of labor and revenue over time. As you can see, this H-E-B is NOT doing too hot. This is obviously not realistic of what real data would show.

## Distribution of Customer Ratings



This shows how many ratings each score got on the 1-5 survey. 4 has the most ratings while 5 has the least. This is probably not very realistic because this average will be close to 3, and it is typically closer to 5.

# Linear Regressions

In this section, I will be running regression analysis in order to test the correlation between two different entities and the strength of the models. There will be two tests that I will be running. The first test is whether or not number of items ordered affects the total revenue. The second test is whether or not the order revenue (or order size) generally affects the ratings that customers give on their surveys. Hope you enjoy!

### order revenue vs number of items ordered
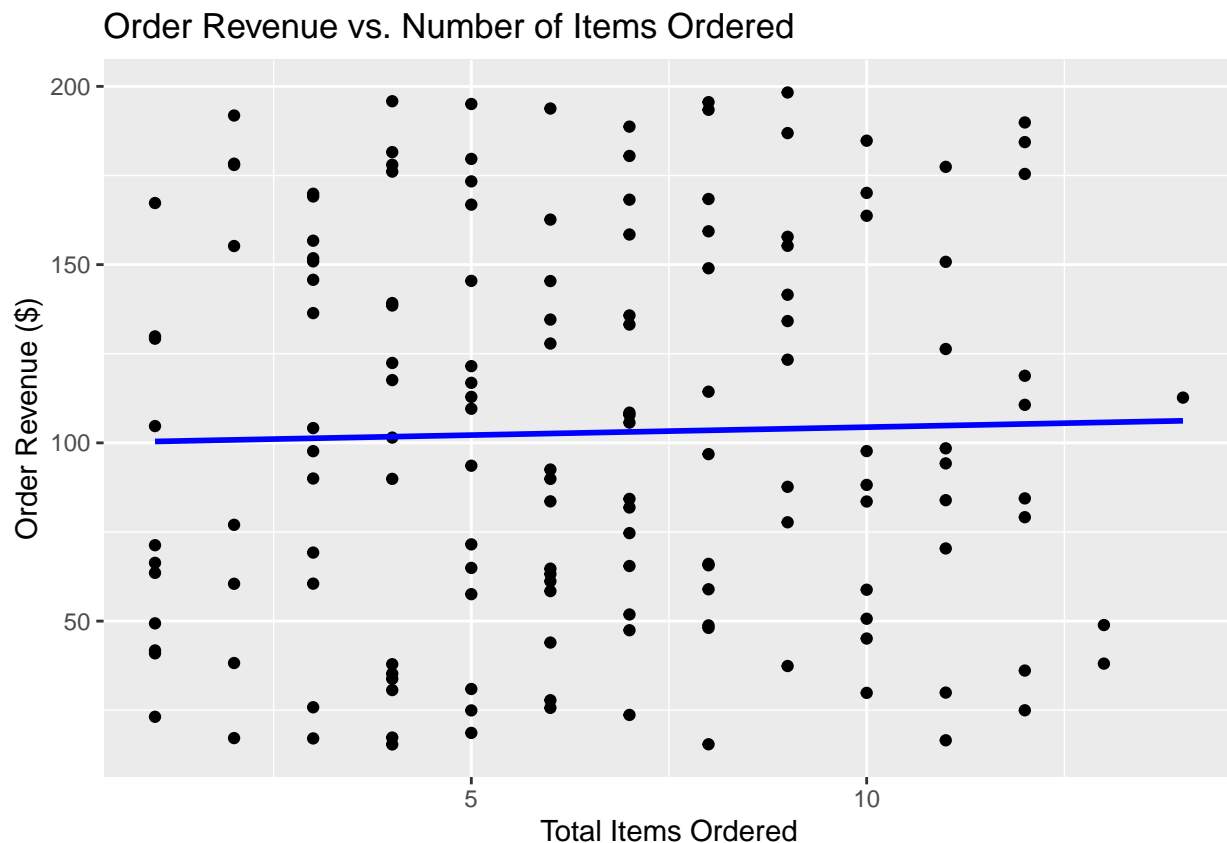
```
##
## Call:
## lm(formula = order_price ~ total_items, data = order_qty_rev)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -88.282 -44.617  -6.528  50.286  94.330
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  99.9490     9.8285  10.169   <2e-16 ***
```

```
## total_items    0.4457      1.3709    0.325       0.746
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 55.13 on 148 degrees of freedom
## Multiple R-squared:  0.0007137,  Adjusted R-squared:  -0.006038
## F-statistic: 0.1057 on 1 and 148 DF,  p-value: 0.7455
```

Let's look at the summary statistics for this linear regression. The coefficient for the intercept is 99 which means that when 0 items are ordered, the model predicts a revenue of approximately 99. This provides little statistical significance, and it serves primarily as a baseline. The slope for the line is .45. This means that for every 1 item ordered, the model predicts an increase of about 45 cents in revenue which is a very small amount statistically.

Moving on to the statistical significance section. The p-value of this test is .745. This means that there is very little significance between number of items ordered and order revenue. This means we fail to reject the null. The R-squared is .0007. This means that less than 1 percent of the variation in order revenue is explained by the number of items purchased There is not enough evidence to say that number of items ordered affects the revenue. This, however, seems contrary to my initial intuition. Because this is synthetic data and far from realistic, take this analysis with a grain of salt.

```
## `geom_smooth()` using formula = 'y ~ x'
```



Order Revenue vs. Number of Items Ordered

Looking at this plot, the line is very flat with a very slight upwards path. This supports what I have previously said with that the revenue does not visibly increase with the number of items ordered. The slight upwards path might indicate there is barely any positive correlation, but it is definitely not enough. The data points are scattered all over the place vertically. This just helps reinforce the lack of correlation.
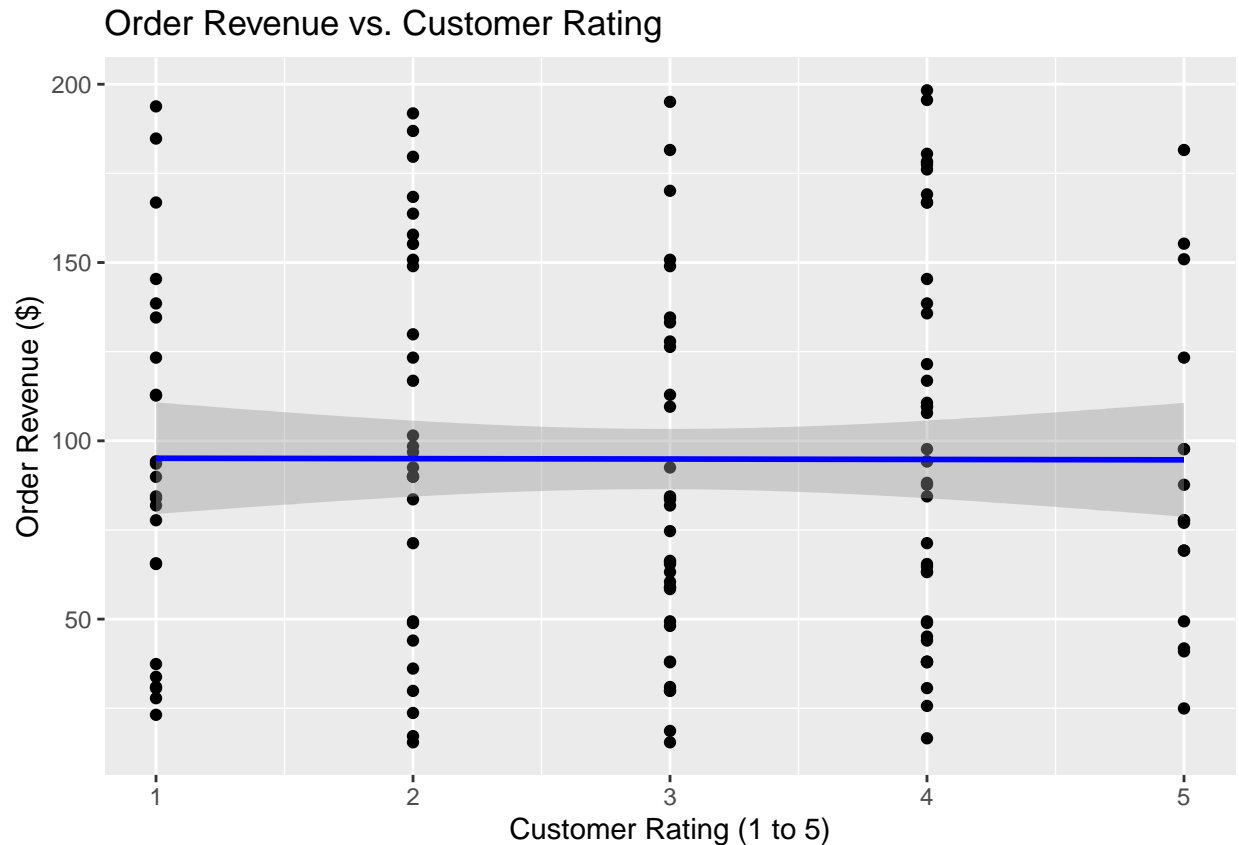
# Customer Rating vs Order Revenue

```
##
## Call:
## lm(formula = order_price ~ one_to_five_rating, data = order_rating_rev)
##
## Residuals:
##     Min     1Q  Median     3Q     Max
## -79.552 -45.435  -6.781  39.532 103.513
##
## Coefficients:
##                    Estimate Std. Error t value Pr(>|t|)
## (Intercept)          95.2271    10.9191   8.721 6.64e-15 ***
## one_to_five_rating   -0.1124     3.3809  -0.033    0.974
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 51.26 on 142 degrees of freedom
## Multiple R-squared:  7.789e-06,  Adjusted R-squared:  -0.007034
## F-statistic: 0.001106 on 1 and 142 DF,  p-value: 0.9735
```

Looking at the summary statistics for this next regression, the intercept is 95 which is the prediction of order revenue when the rating is 0. This is not significant since the survey is on a 1-5 scale. The rating slope is -.11 which means that as customer ratings increase, order revenue slightly decreases. However this is not statistically significant either.

The p-value is .97 which is way above .05. This means that we reject the null hypothesis. The rating does not have a significant affect on the revenue. The R-squared is essentially 0. This means that the ratings explain virtually 0% of the variation in order revenue.

```
## `geom_smooth()` using formula = 'y ~ x'
```

**Order Revenue vs. Customer Rating**

The plot is the same as the previous regression that I did. The flat line and scattered plots further prove the lack of relation between survey ratings and order revenue.

# Conclusion

In conclusion, in this statistical analysis that I conducted on my H-E-B curbside database, I was able to visualize multiple different entity relationships, such as the distribution of order prices, the comparison of revenue and the number of labor hours among the workers.

Using synthetic data has skewed my intuition solely because of how unrealistic this data is. It is comical to see how poorly the store is doing when looking at revenue and labor costs. However, this has made the project more fun!

Looking at the last portion, the linear regressions, this was the most boring part solely because there was not significant correlation between any of the regressions that I ran.

All in all, I really enjoyed this project using R, and I hope you have gathered some valuable insights of my intuition!