
Factors Influencing Ride-Sharing Pricing in Boston: A Study on Lyft and Uber

6101 Intro to Data Science Midterm Presentation
Aaron Yang Jianjun Gao Luhuan Wang

Why this topic?

1. Ride-sharing Popularity
2. Easy Data Access
3. Dynamic Pricing Model
4. Comparative Analysis

The Uber logo is displayed within a black rectangular box. It consists of the word "Uber" in a white, sans-serif font.

Uber

The Lyft logo is displayed within a pink rectangular box. It features the word "lyft" in a white, lowercase, sans-serif font, with a stylized car icon integrated into the letter 't'.

lyft

Team Contributions

Aaron Yang	Jianjun Gao	Luhuan Wang
Outline & Proposal	Data Collection	Data Summary
Data Analysis(with Jianjun)	Data Analysis(with Aaron)	Slide Creation
Presentation	Presentation	Presentation
Paper Co-author	Paper Co-author	Paper Co-author

Key Concepts & Principles

Data Science and Statistical Methods	Programming	Visualization
Data Collection & Cleaning	R	Graphical Representation
Exploratory Data Analysis (EDA)	Python	
Regression Model		


Project Workflow Overview 1

Step No.	Phase	Tasks & Details
1.	Topic Decision	<ul style="list-style-type: none">- Team Meeting- Choose a topic based on available data
2.	Data Collection	<ul style="list-style-type: none">- Identify Data Sources- Gather data - Initial Data Cleaning
3.	Research & Background	<ul style="list-style-type: none">- Literature Review- Industry Knowledge
4.	Data Processing & EDA	<ul style="list-style-type: none">- Data Cleaning & Transformation- Descriptive Statistics- Visualization
5.	Data Analysis & Modeling	<ul style="list-style-type: none">- Hypothesis Testing- Regression Analysis- Advanced Modeling - Model Evaluation

Project Workflow Overview 2

Step No.	Phase	Tasks & Details
6.	Interpretation of Results	<ul style="list-style-type: none">- Discuss Findings- Statistical Significance
7.	Presentation Preparation	<ul style="list-style-type: none">- Slide Development- Rehearsal
8.	Presentation & Feedback	<ul style="list-style-type: none">- Deliver Presentation- Collect Feedback
9.	Final Paper	<ul style="list-style-type: none">- Write a paper



Data Collection & Filtering

 BM · UPDATED 4 YEARS AGO

▲ 123

New Notebook

Download (47 MB)

 vs 

Data Card

Code (23)

Discussion (3)

About Dataset

Uber vs Lyft

This is a very beginner-friendly dataset. It does contain a lot of NA values. It is a good dataset if you want to use a Linear Regression Model to see the pattern between different predictors such as `hour` and `price`.

A really amazing part of this dataset is that I have included the corresponding `weather_data` for that `hour` with a short summary of the `weather`. Other important factors are `temperature`, `wind`, and `sunset`.

Usability ⓘ
9.12

License
[CC0: Public Domain](#)

Expected update frequency
Quarterly

Tags
Travel

Kaggle DataBase

Data Collection & Filtering

[illegible]

More than
300,000 rows and
55 columns data



Original Dataset

[illegible]

Data Collection & Filtering

Drop null values, random sampling(random_state = 42),choose 6 columns

	hour	price	distance	surge_multiplier	cab_type	name
0	6	7.0	4.51	1.0	Lyft	Shared
1	0	10.5	2.80	1.0	Uber	WAV
2	12	7.0	1.09	1.0	Lyft	Lyft
3	3	15.5	0.92	1.0	Uber	Black
4	9	16.5	1.12	1.0	Lyft	Lux Black
...
5995	3	19.5	2.32	1.0	Uber	Black
5996	18	26.0	0.54	1.0	Lyft	Lux Black XL
5997	10	14.0	1.89	1.0	Uber	UberX
5998	13	9.0	2.08	1.0	Lyft	Lyft
5999	4	22.5	2.36	1.0	Uber	Black

6000 rows × 6 columns

Columns Adjustments

cab_type <chr>	name <chr>
Lyft	Shared
Uber	WAV
Lyft	Lyft
Uber	Black
Lyft	Lux Black
Uber	WAV



cab_company <chr>	cab_type <chr>
Lyft	Shared
Uber	WAV
Lyft	Lyft
Uber	Black
Lyft	Lux Black
Uber	WAV

Descriptive Stats

```
unique(df$cab_company) #see unique values in company
```

```
## [1] "Lyft" "Uber"
```

```
summary(df) # descriptive stats
```

```
##      hour      price      distance  surge_multiplier
## Min.   : 0.0   Min.   : 2.5   Min.   :0.02   Min.    :1.00
## 1st Qu.: 6.0   1st Qu.: 9.0   1st Qu.:1.30  1st Qu.:1.00
## Median :12.0   Median :13.5   Median :2.19  Median :1.00
## Mean   :11.7   Mean    :16.5   Mean    :2.20  Mean    :1.02
## 3rd Qu.:18.0   3rd Qu.:22.5   3rd Qu.:2.94  3rd Qu.:1.00
## Max.   :23.0   Max.    :73.5   Max.    :7.46  Max.    :2.50
## cab_type      name
## Length:6000    Length:6000
## Class :character Class :character
## Mode :character  Mode :character
##
```

```
unique(df$cab_type) #see unique values in cab_type
```

```
## [1] "Shared"      "WAV"          "Lyft"          "Black"          "Lux Black"
## [6] "Lyft XL"     "Lux"          "UberX"         "Black SUV"      "UberPool"
## [11] "Lux Black XL" "UberXL"
```

```
## {r, results='markup', echo=TRUE}
# measure of variance
sd(df$price)
sd(df$distance)
sd(df$surge_multiplier)
```

```
[1] 9.43
[1] 1.14
[1] 0.103
```

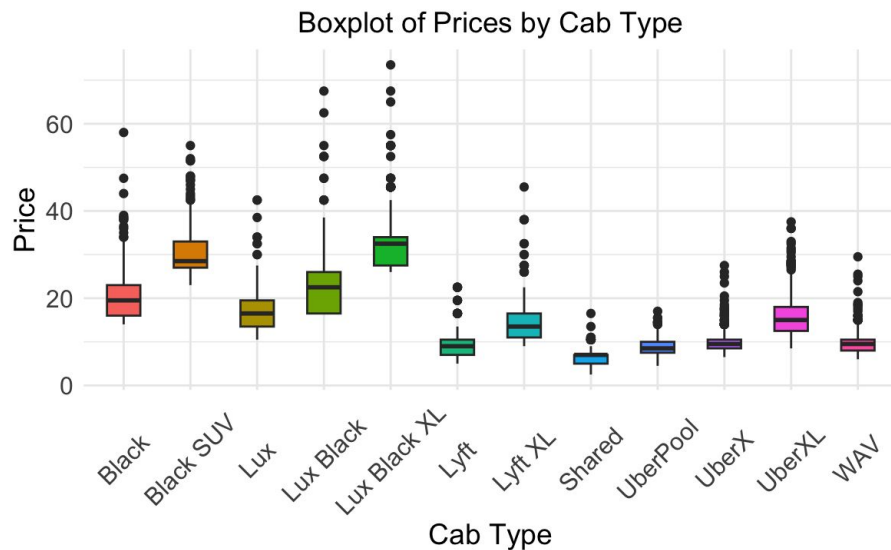
Data Visualization



0-8: Relatively High

15-23: Relatively Low

Data Visualization



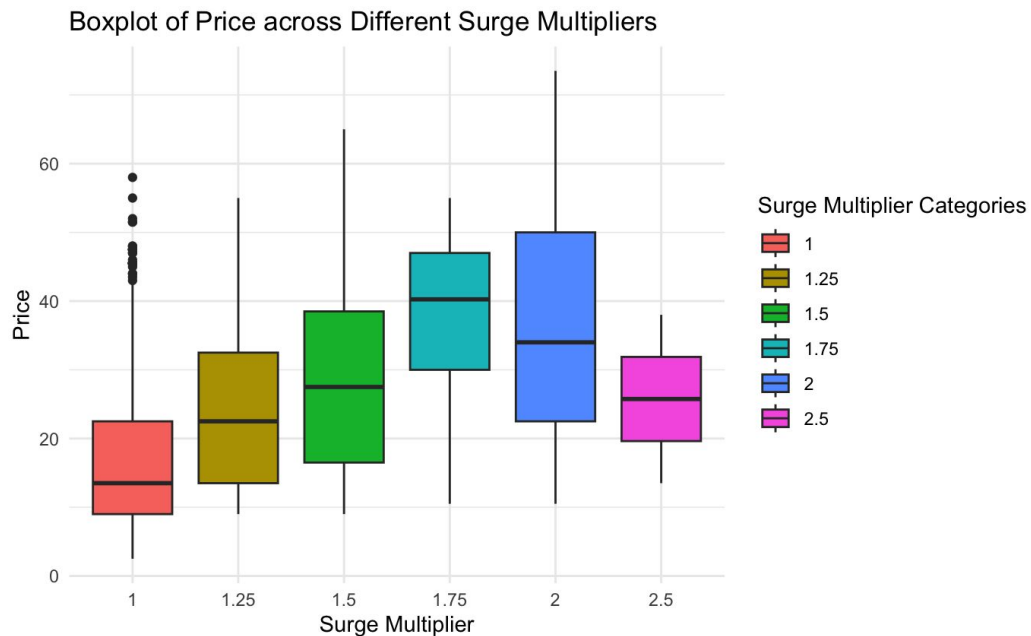
Highest Price:

Black SUV, Lux Black XL

Lowest Price:

Shared

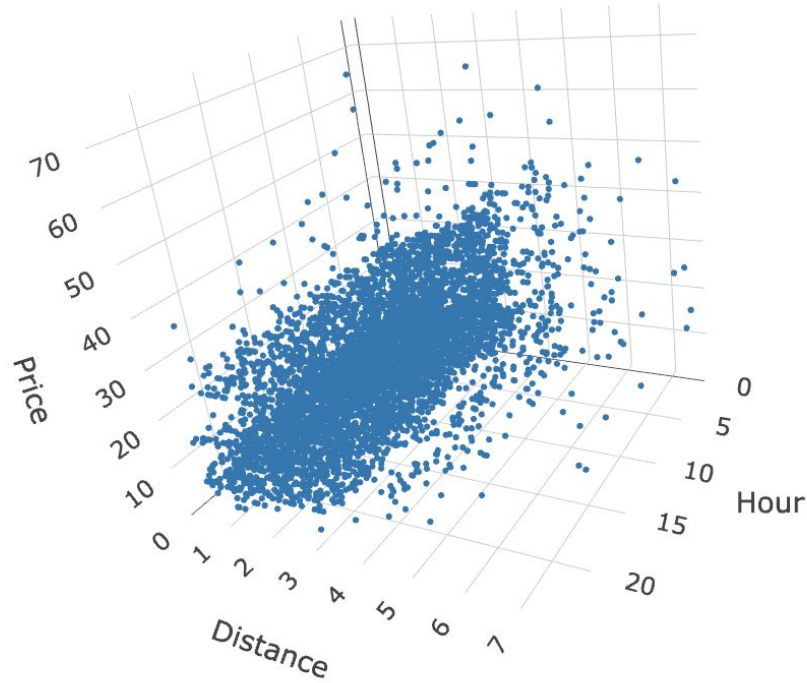
Data Visualization



**Highest Median
when Tipped 1.75**

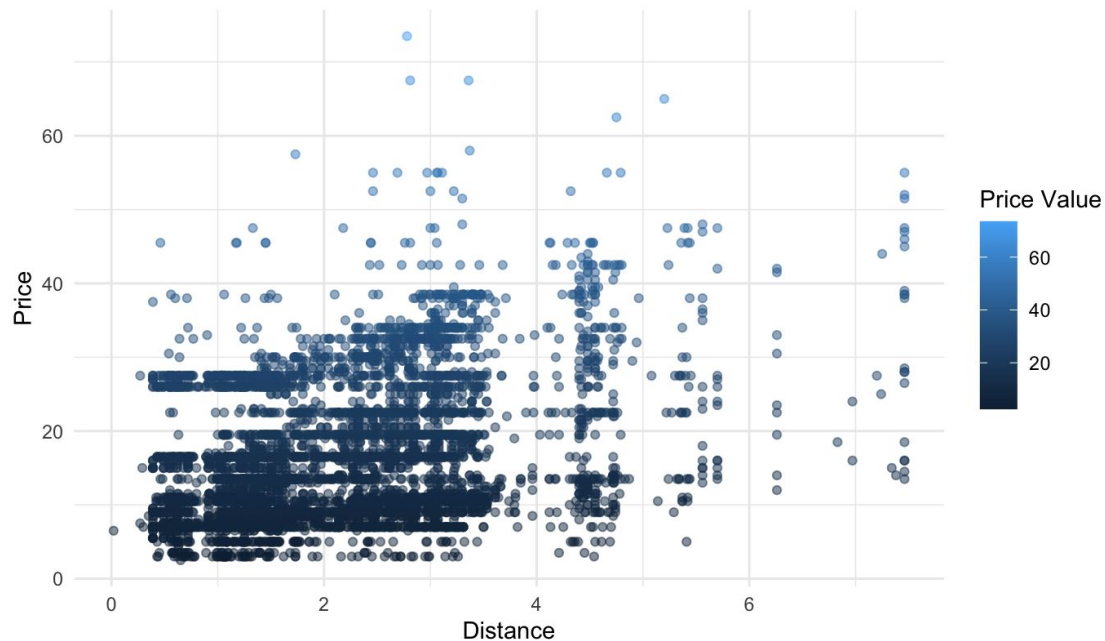
**Lowest Median when
Tipped 1**

Data Visualization



Data Visualization

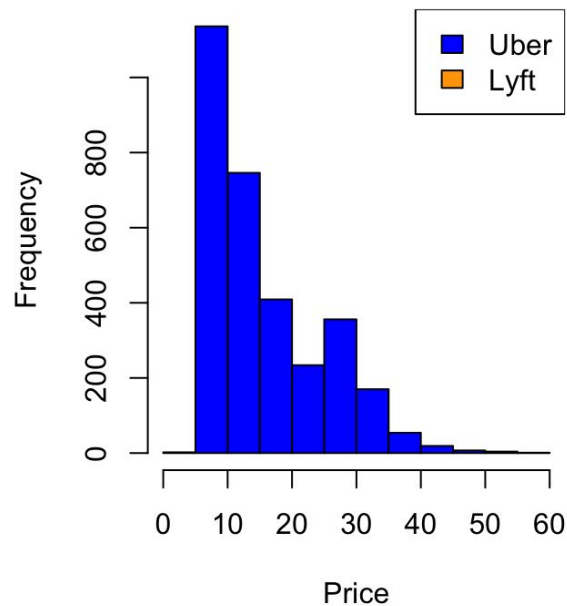
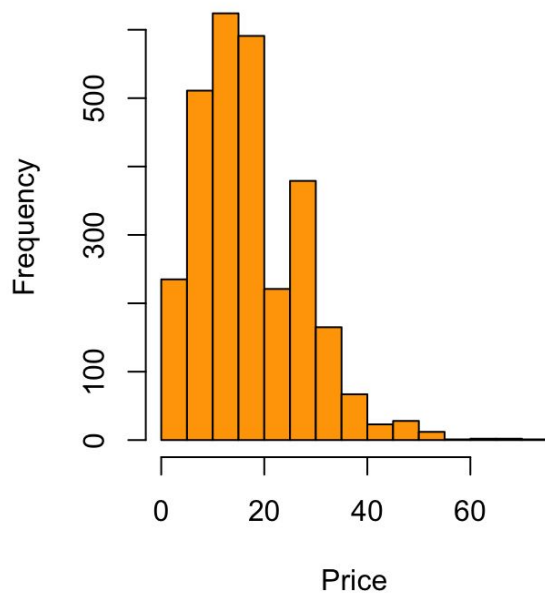
The Plot of Price vs. Distance



**A Trend That
is not Obvious**

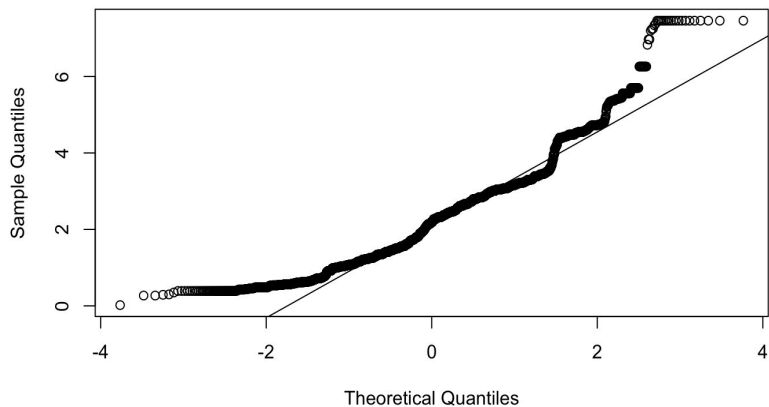
Data Visualization

Histogram of Prices by Cab Type (Ly) Histogram of Prices by Cab Type (Uk)

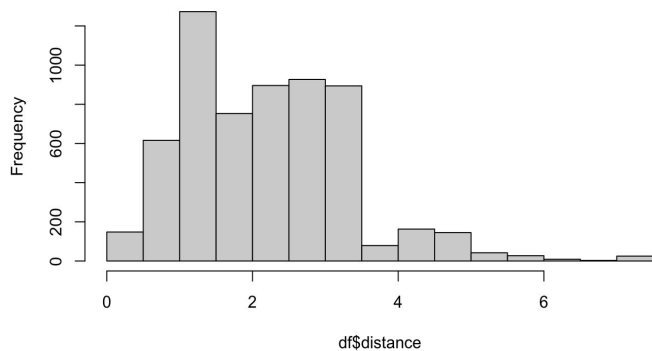


Normality Test For Distance

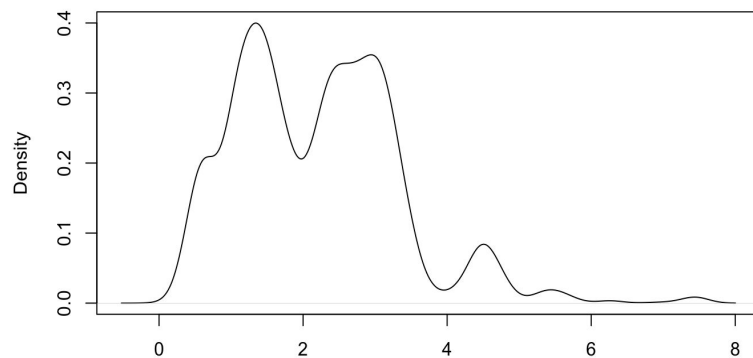
Normal Q-Q Plot



Histogram of df\$distance

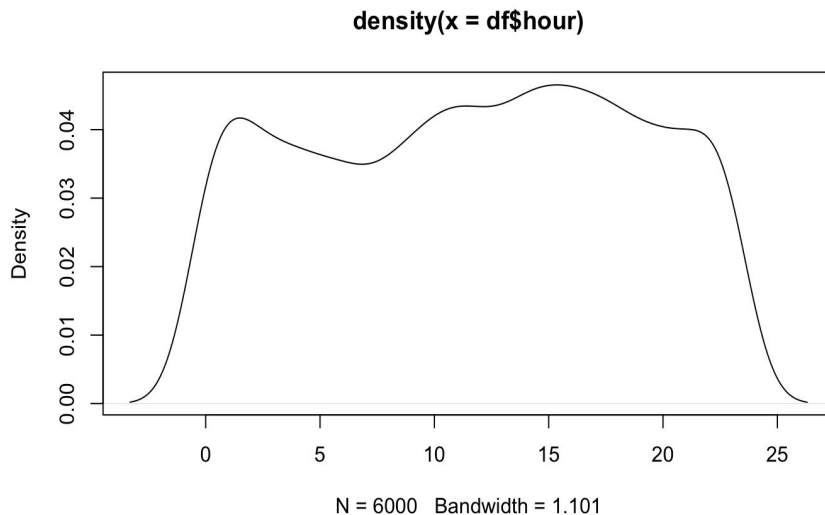
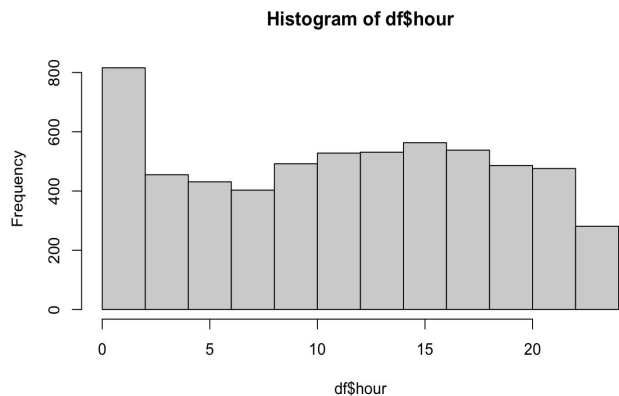
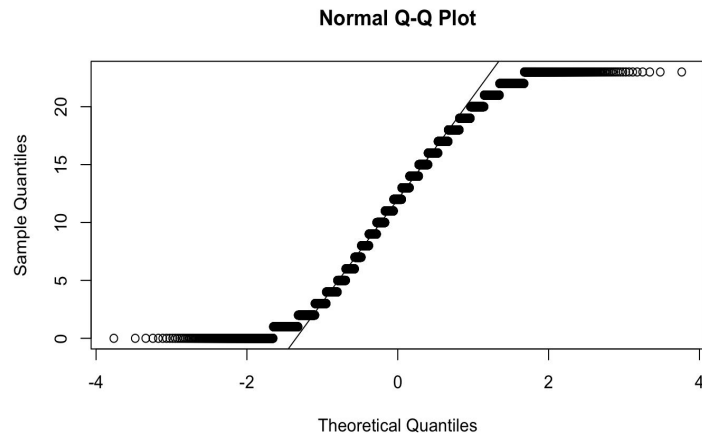


density(x = df\$distance)



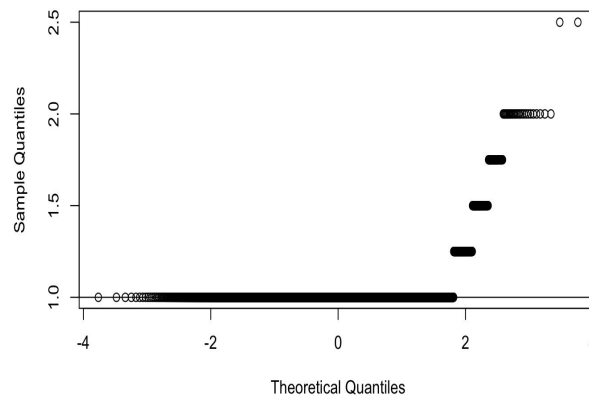
N = 6000 Bandwidth = 0.1808

Normality Test For Hour

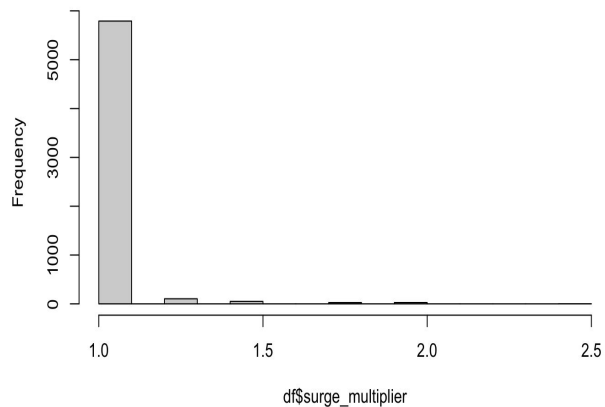


Normality Test For Surge_multiplier

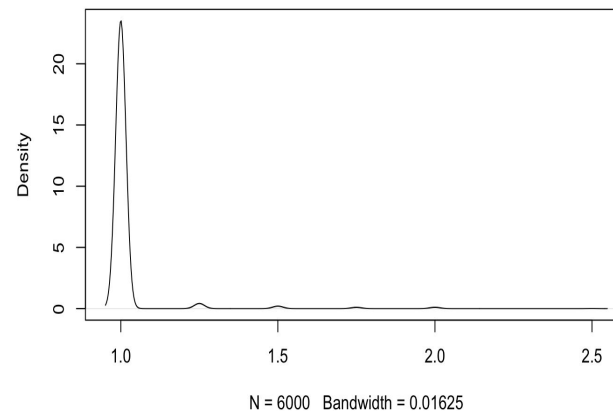
Normal Q-Q Plot



Histogram of df\$surge_multiplier

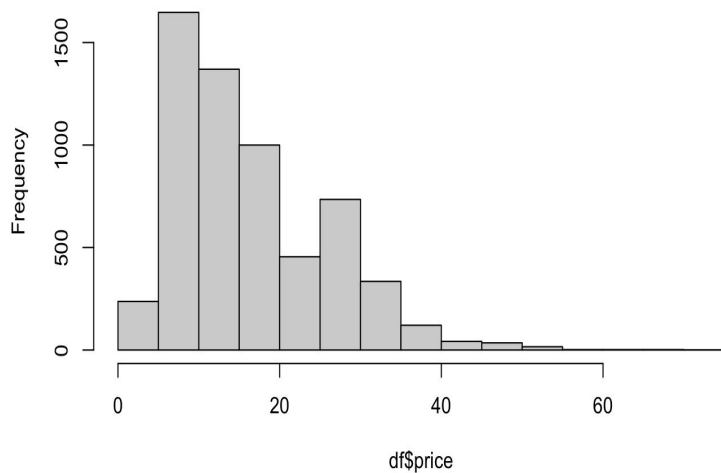


density(x = df\$surge_multiplier)

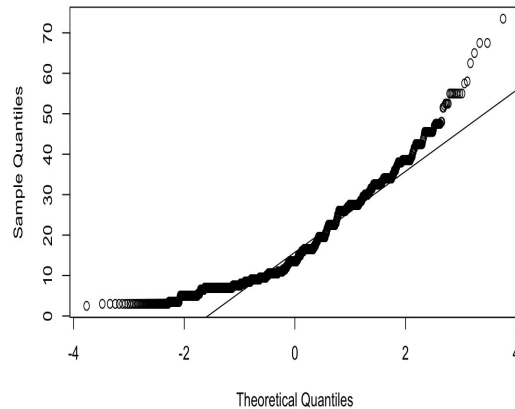


Normality Test For Price

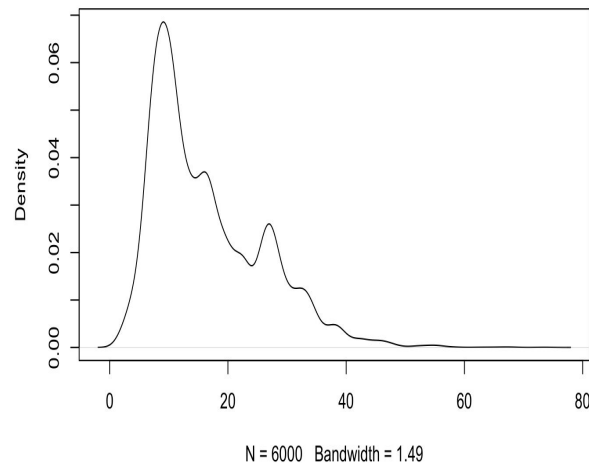
Histogram of df\$price



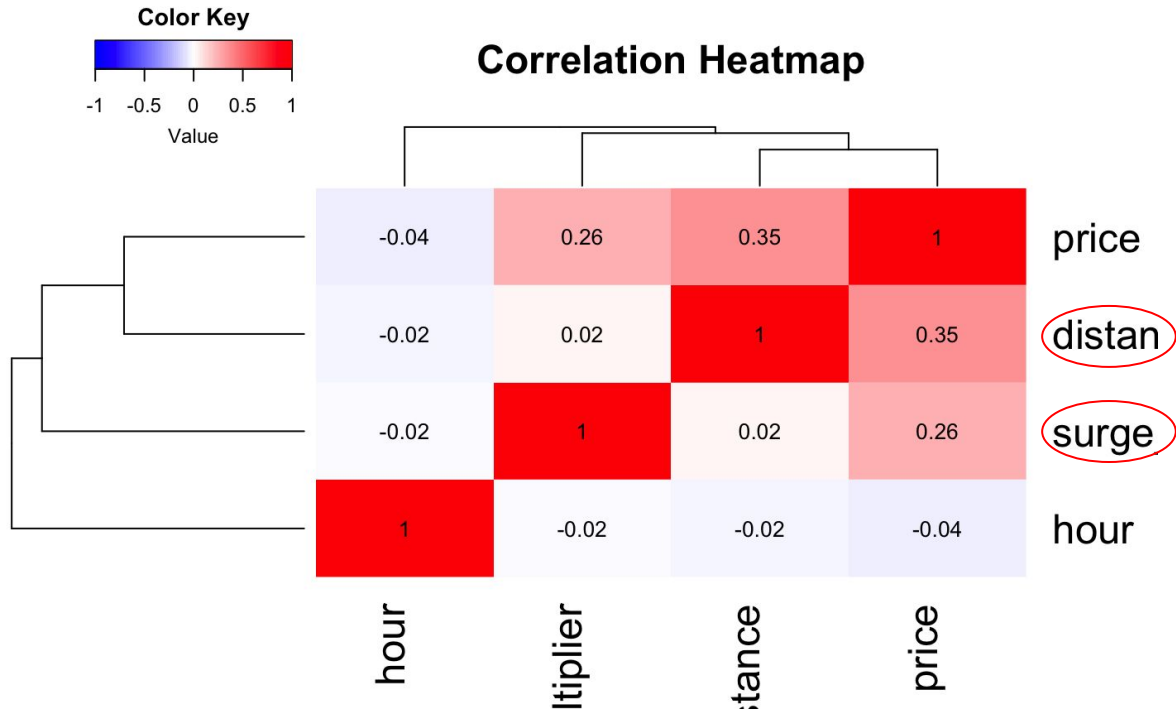
Normal Q-Q Plot



density(x = df\$price)



Correlation Heatmap for Numeric Columns



Chi-Squared Test for Cab Company & Cab Type

```
# Chi-squared test
chisq_result <- chisq.test(df$cab_company, df$cab_type)

chisq_result
```

```
##
##  Pearson's Chi-squared test
##
## data:  df$cab_company and df$cab_type
## X-squared = 6000, df = 11, p-value <2e-16
```

Cab_company and cab_type are significantly related.

Linear Regression Analysis

```
Call:
lm(formula = price ~ distance + surge_multiplier, data = df)

Residuals:
    Min       1Q   Median       3Q      Max
-32.53  -6.75  -1.61   4.98  38.56

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   -12.9730     1.1094   -11.7  <2e-16 ***
distance         2.8232     0.0962    29.4  <2e-16 ***
surge_multiplier 22.8997     1.0703    21.4  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.52 on 5997 degrees of freedom
Multiple R-squared:  0.184,    Adjusted R-squared:  0.183
F-statistic: 674 on 2 and 5997 DF,  p-value: <2e-16
```

Distance, surge_multiplier are good predictors for price.

In terms of 0.184 R-squared, the model is not so good.

Summary

-the reasons we chose this topic

With the development of social technology and fast-paced life, online shared-vehicle-services have attracted more and more attention.



Selected graphics

- Used python to do the initial data cutting

- Normality Test

 - Boxplot

 - Heatmap

Final result (according to EDA)

- **Weak relation to: Price and Hours**
- **Price Fluctuations Throughout the Day**
- **Possible Peak Hours:** the first peak close to the 5th hour might /the subsequent decrease might signify a lull during mid-morning.
- **Midday Drop & Evening Rise**

Final result (according to EDA)

Strong relation to: Surge_multiplier and Distance with Price

Hours is not a Key-point

Final result (according to EDA)

Strong relation to: Cab_company and Cab_type

Cab_company: This may result in software design and vehicle design being different from each other:

- Uber is the cheaper company.

- Uber is available all over the world, while Lyft is only available in the United States and Canada.



Reference Page

1. Brllrb. (2019). Uber and Lyft Dataset: Boston, MA. Kaggle.
<https://www.kaggle.com/datasets/brllrb/uber-and-lyft-dataset-boston-ma/data>
2. Peng, R. D., & Matsui, E. (2016). *Art of Data Science*.

Thank you for your attention!