



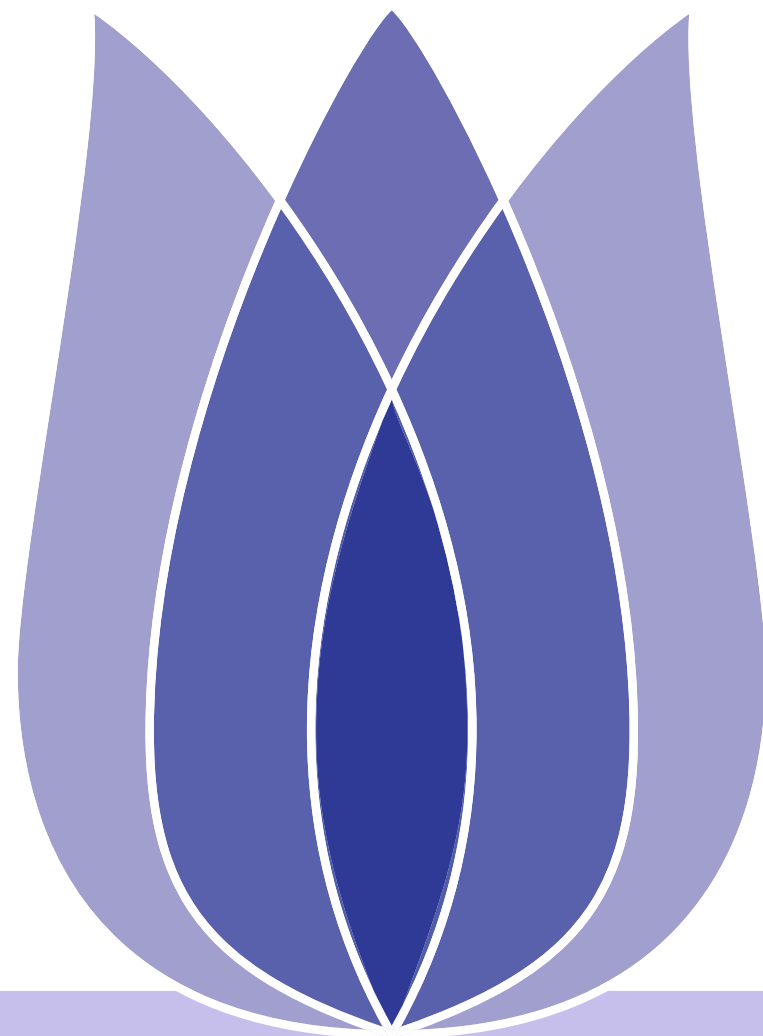
# New York City Taxi Trip Duration

Siyi Yu

Jilin University

Computer Science And Technology

April 25, 2021





# Overview

- [Project Introduction](#)
- [Data Analysis](#)
- [Model selection](#)
- [Result](#)

## Project Introduction

Project Introduction

## Data Analysis

Info

Spearman Correlation

Simple Attributes

Complex Attributes

Labels

## Model selection

Models

## Result

Models



Project Introduction

Project Introduction

Data Analysis

Model selection

Result

# Project Introduction

# Project Introduction

Project Introduction

Project Introduction

Data Analysis

Model selection

Result



## Data Fields for This Project

Intro

- **id, vendor\_id, store\_and\_fwd\_flag, passenger\_count**
- **pickup\_datetime, dropoff\_datetime**
- **pickup\_longitude, pickup\_latitude, dropoff\_longitude, dropoff\_latitude**
- **trip\_duration**



**TULIP**

Team for Universal Learning and Intelligent Processing



[Project Introduction](#)

**[Data Analysis](#)**

[Info](#)  
[Spearman Correlation](#)  
[Simple Attributes](#)  
[Complex Attributes](#)  
[Labels](#)

[Model selection](#)

[Result](#)

# Data Analysis





- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

train:

```
<class 'pandas.core.frame.DataFrame'>
Index: 1458644 entries, id2875421 to id1209952
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   vendor_id             1458644 non-null int64
1   pickup_datetime       1458644 non-null object
2   dropoff_datetime      1458644 non-null object
3   passenger_count        1458644 non-null int64
4   pickup_longitude      1458644 non-null float64
5   pickup_latitude       1458644 non-null float64
6   dropoff_longitude     1458644 non-null float64
7   dropoff_latitude      1458644 non-null float64
8   store_and_fwd_flag    1458644 non-null object
9   trip_duration         1458644 non-null int64
dtypes: float64(4), int64(3), object(3)
memory usage: 122.4+ MB
```

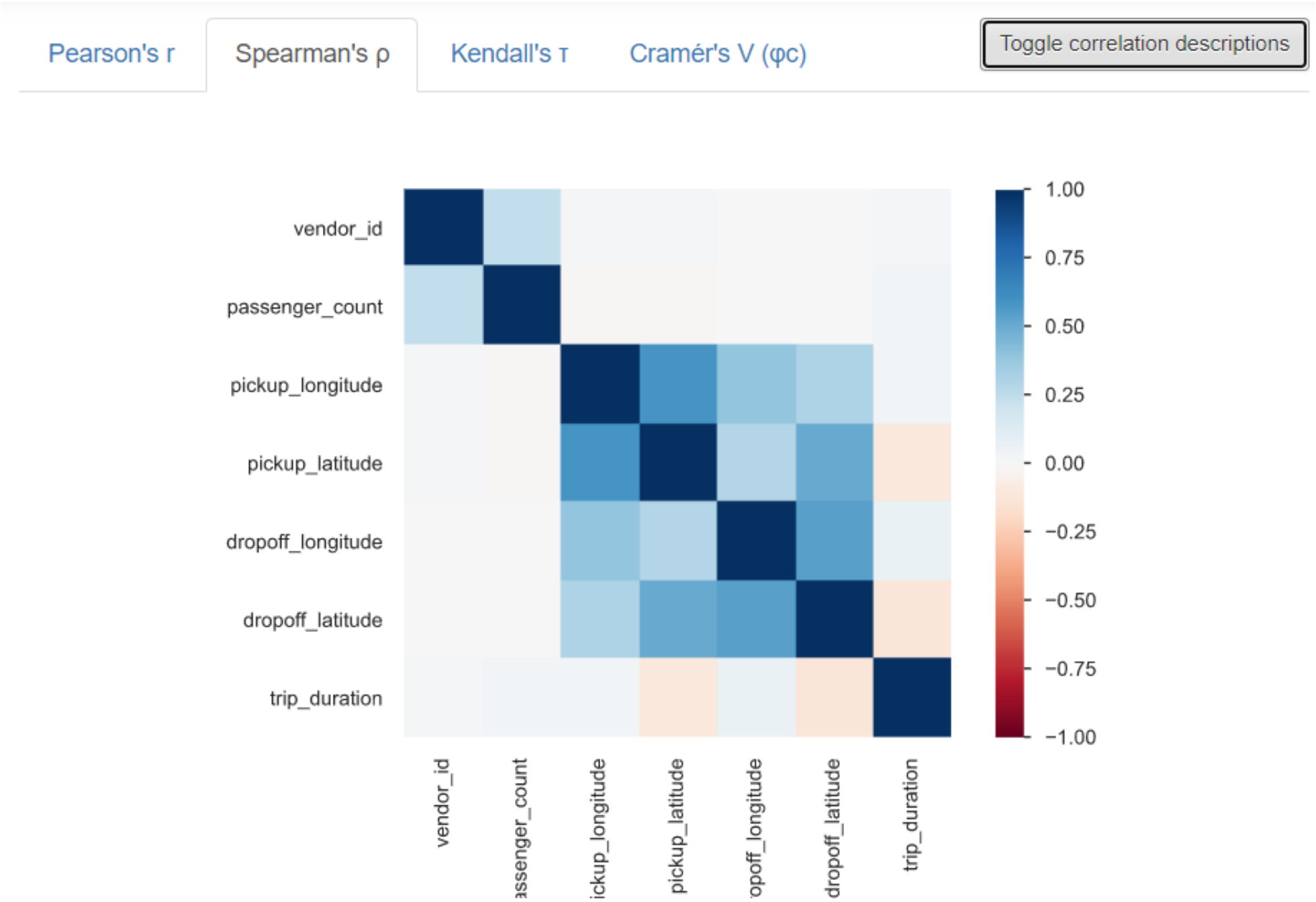
test:

```
<class 'pandas.core.frame.DataFrame'>
Index: 625134 entries, id3004672 to id0621643
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   vendor_id             625134 non-null int64
1   pickup_datetime       625134 non-null object
2   passenger_count        625134 non-null int64
3   pickup_longitude      625134 non-null float64
4   pickup_latitude       625134 non-null float64
5   dropoff_longitude     625134 non-null float64
6   dropoff_latitude      625134 non-null float64
7   store_and_fwd_flag    625134 non-null object
dtypes: float64(4), int64(2), object(2)
memory usage: 42.9+ MB
```



# Spearman Correlation

- [Project Introduction](#)
- [Data Analysis](#)
- [Info](#)
- [Spearman Correlation](#)**
- [Simple Attributes](#)
- [Complex Attributes](#)
- [Labels](#)
- [Model selection](#)
- [Result](#)



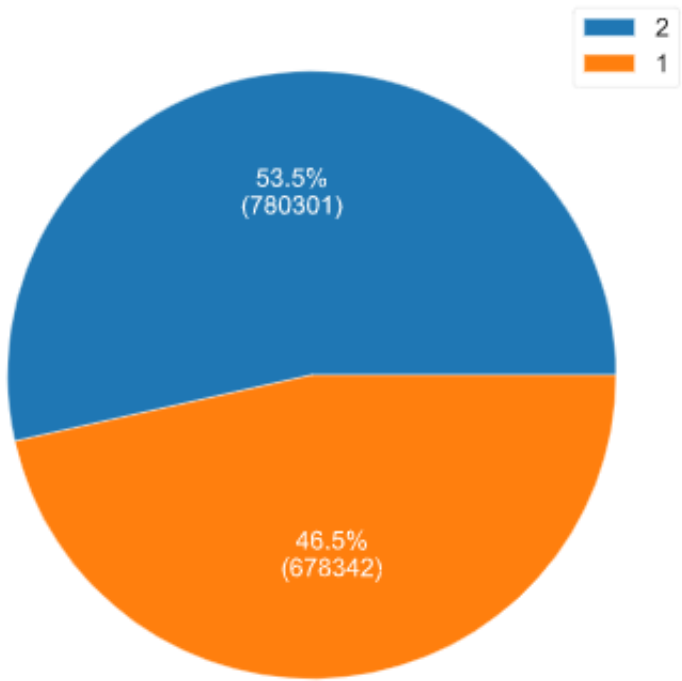




# Simple Attributes

- [Project Introduction](#)
- [Data Analysis](#)
- [Info](#)
- [Spearman Correlation](#)
- [Simple Attributes](#)**
- [Complex Attributes](#)
- [Labels](#)
- [Model selection](#)
- [Result](#)

- id:Determine the uniqueness of id
- vendor\_id:
  - ◆



- ◆ Convert int64 into string



- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

■ store\_and\_fwd\_flag:

■

Value	Count	Frequency (%)
False	1450599	99.4%
True	8045	0.6%

■ passenger\_count:

■

```
1    1033540
2     210318
5      78088
3      59896
6      48333
4      28404
0           60
7           3
9            1
8            1
Name: passenger_count, dtype: int64
```



# Complex Attributes

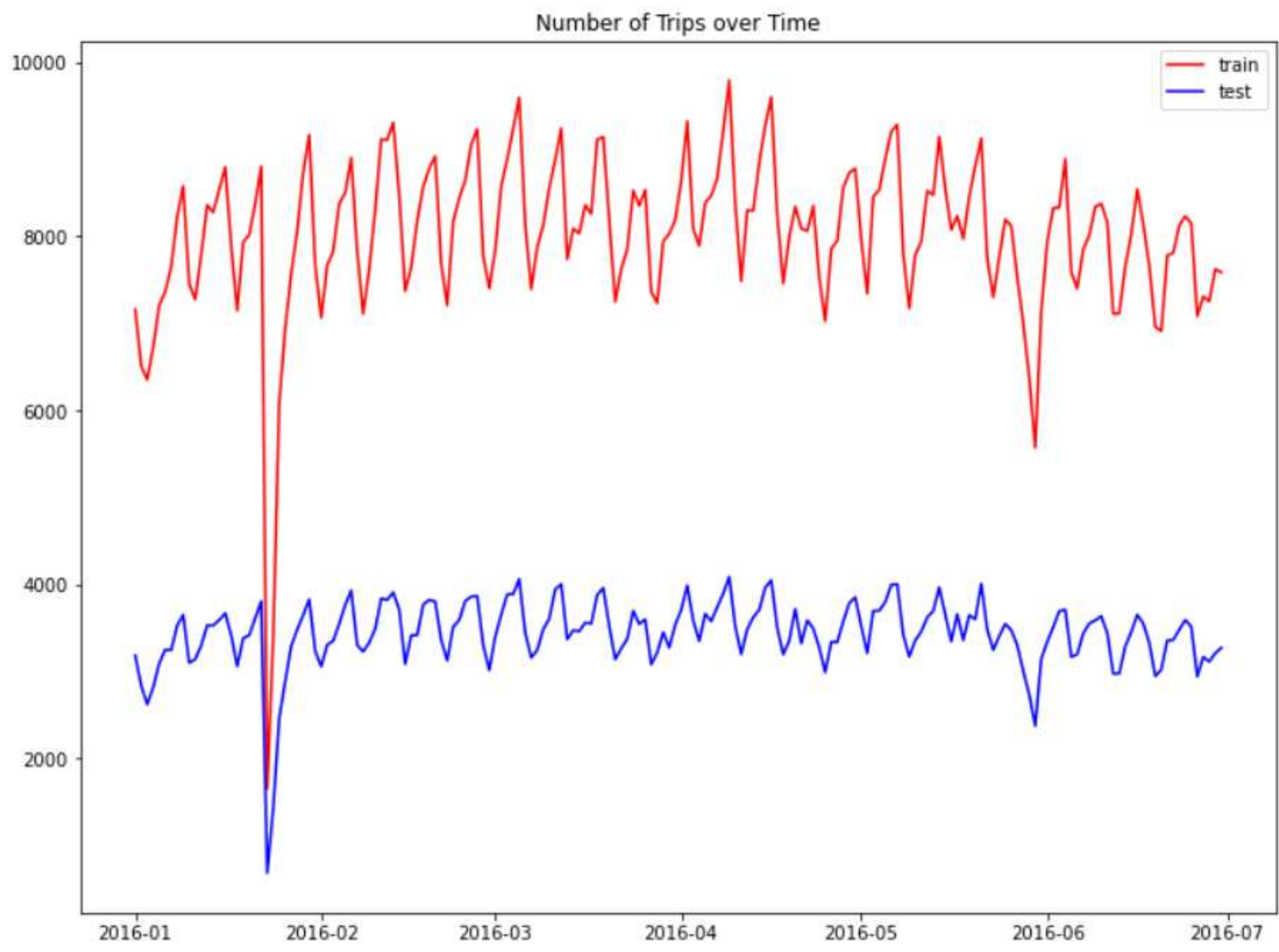
- [Project Introduction](#)
- [Data Analysis](#)
- [Info](#)
- [Spearman Correlation](#)
- [Simple Attributes](#)
- [Complex Attributes](#)
- [Labels](#)
- [Model selection](#)
- [Result](#)

- dropoff\_datetime:
  - ◆ Only exist in test.csv
  - ◆ check dropoff\_datetime-pickup\_datetime == trip\_duration
  - ◆ drop



- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

- pickup\_datetime:
  - ◆





[Project Introduction](#)

[Data Analysis](#)

[Info](#)

[Spearman Correlation](#)

[Simple Attributes](#)

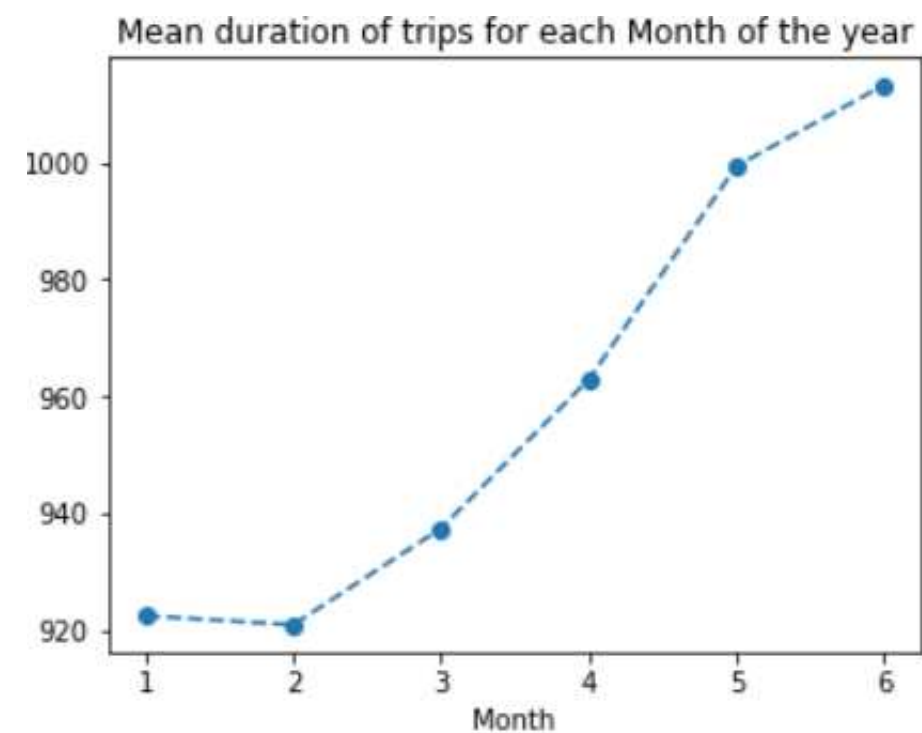
[Complex Attributes](#)

[Labels](#)

[Model selection](#)

[Result](#)

- pickup\_datetime:
  - ◆ Divide it into year, month, day, hour
  - ◆ The year is all 2016,so drop
  - ◆



- ◆ Since only six monthes in train and test(1,2,3,4,5,6),convert the type of month into string for one-hot encoding.



**TULIP**

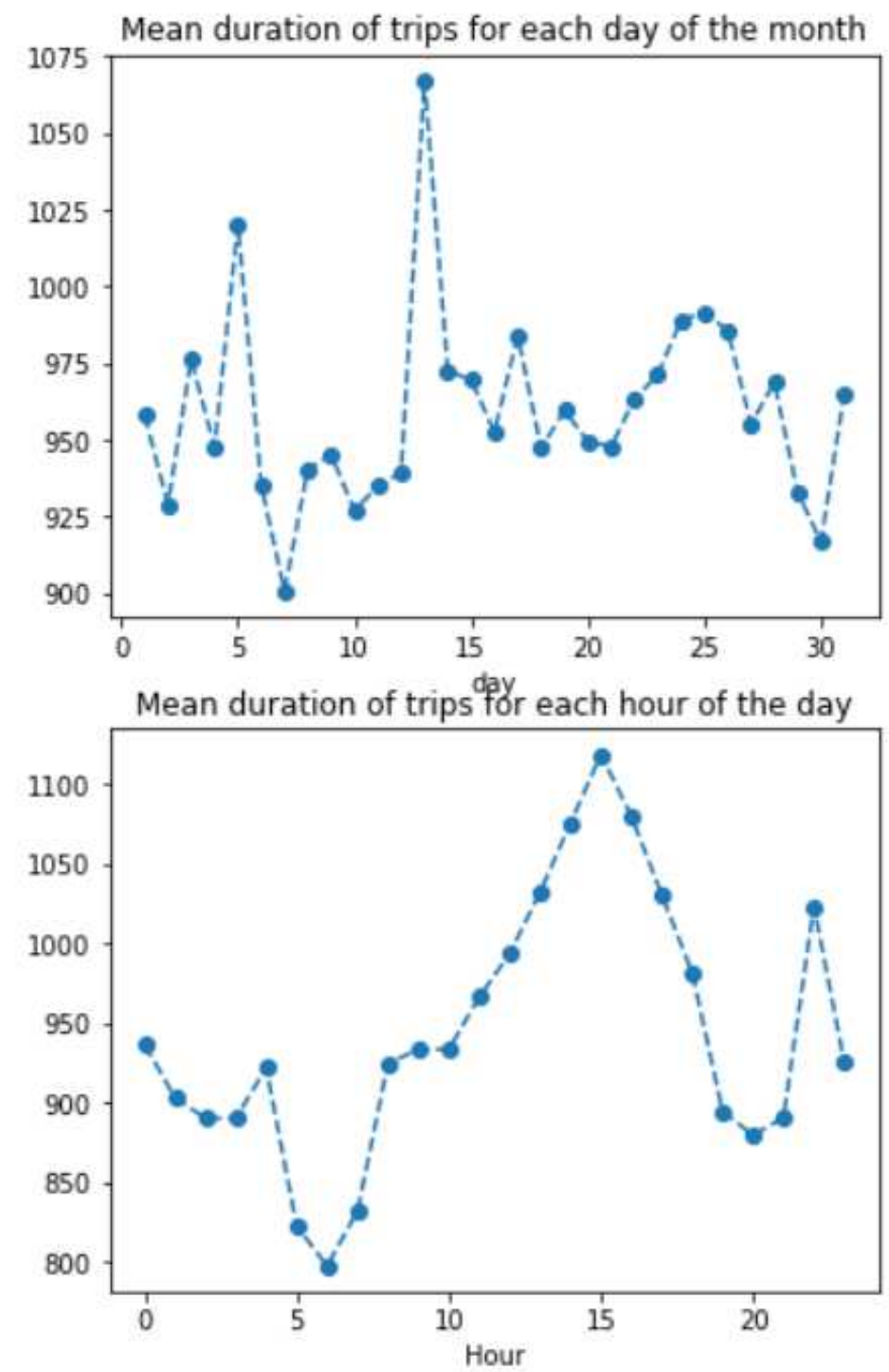
Team for Universal Learning and Intelligent Processing





- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

■ pickup\_datetime: ◆



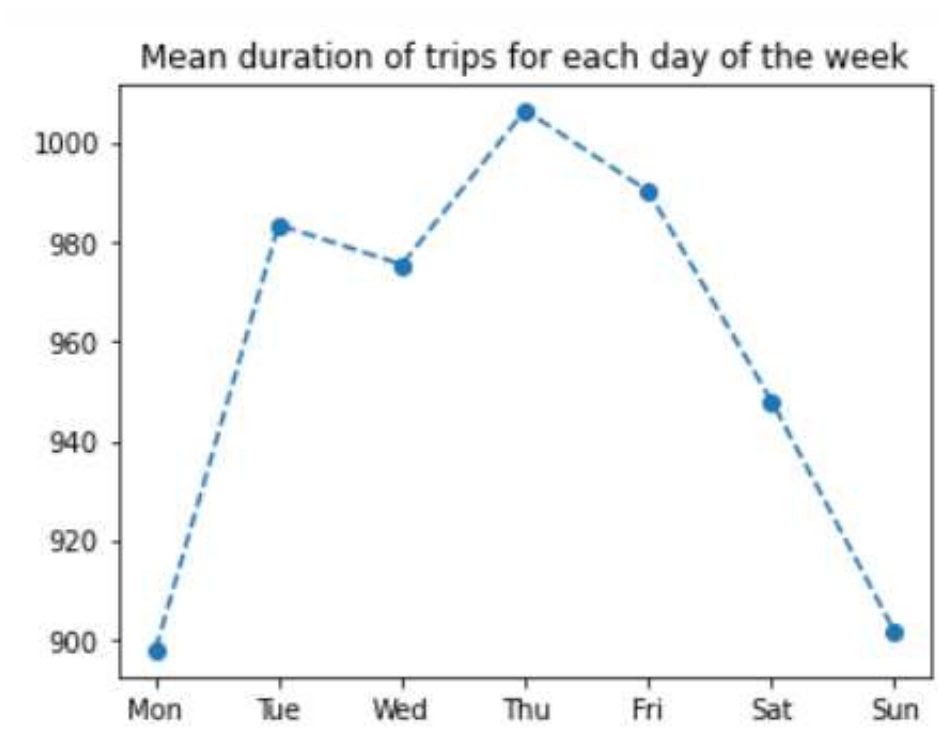
- It can be seen from the figure that the day of the month will affect the taxi trip turation.
- However,if we use one-hot enconding directly,it will cause the disaster of dimensionality.
- In term of the trip turation,we divide it into three categories.
- a:5,13      b:7,30      c:other
- The same process for the hour.
- x:0.1.2.3.4.8.9.10.19.20.21.23
- y:11.12.13.17.18.22
- z:14.15.16
- w:5,6,7





- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

- pickup\_datetime:
  - ◆ Derive weekday from the date.
  - ◆

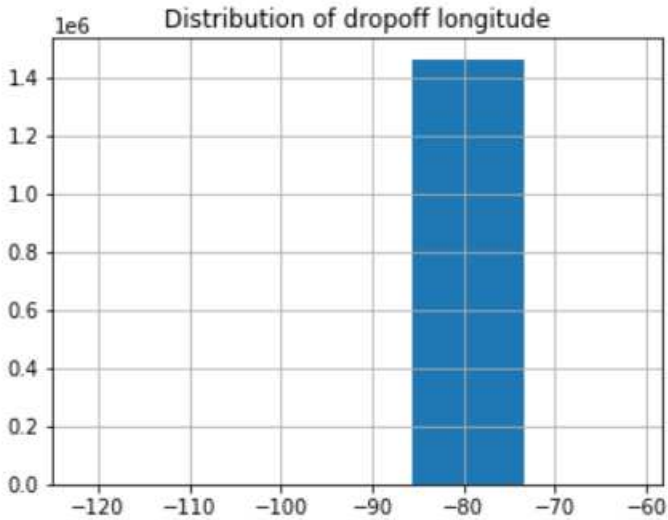
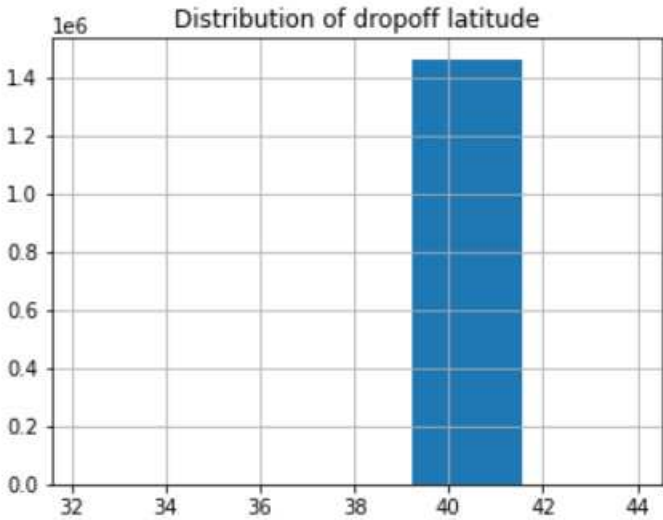
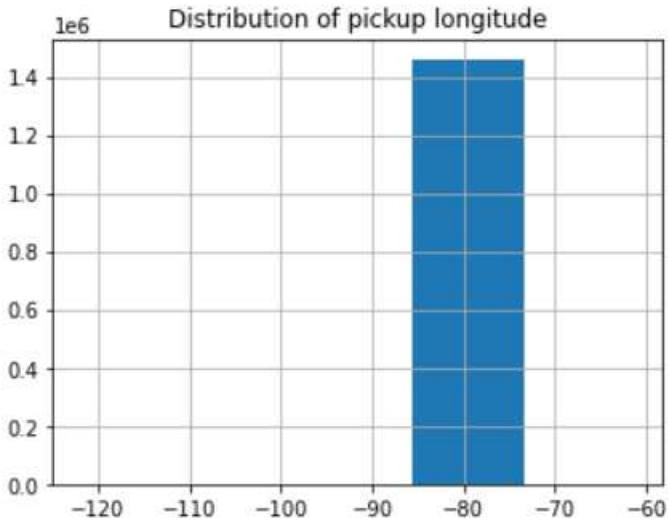
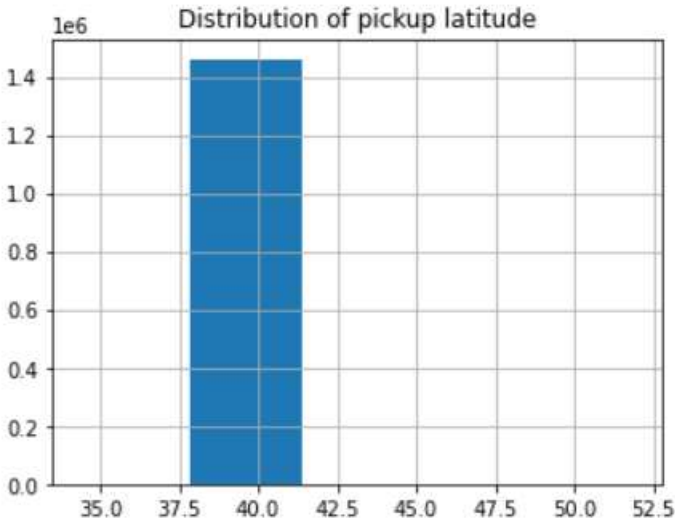


- ◆ Obviously,weekday affects the trip duration,so we add the attribute for the data.



- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

- latitude, longitude:
  - ◆ duration = distance / speed
  - ◆





[Project Introduction](#)

[Data Analysis](#)

[Info](#)

[Spearman Correlation](#)

[Simple Attributes](#)

[Complex Attributes](#)

[Labels](#)

[Model selection](#)

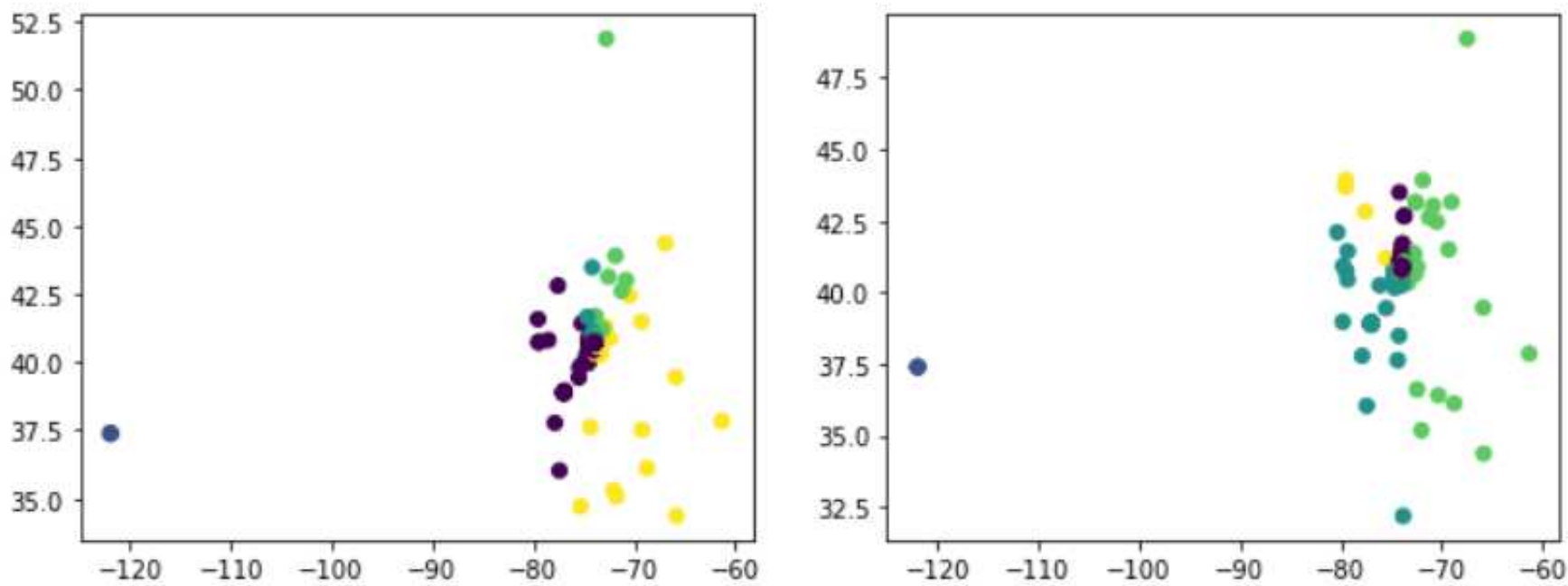
[Result](#)

- distance:
  - ◆ Euclidean distance
  - ◆ Manhattan distance



- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

- speed:
  - ◆ In term of pickup\_latitude,pickup\_longitude,dropoff\_longitude,drop\_latitude,divide data into 5 clusters by KMeans
  - ◆



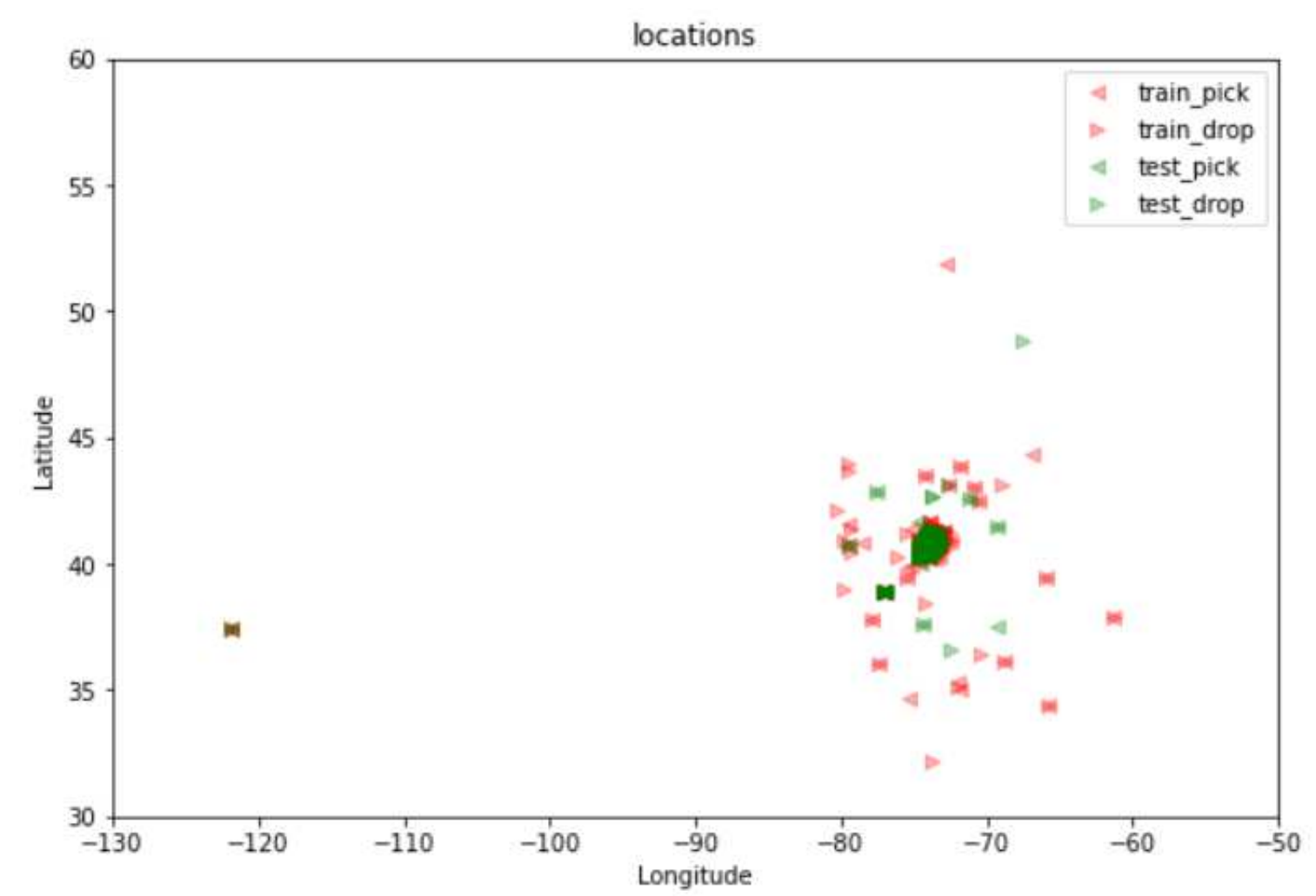
- ◆ Direction



- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

■ outliers:

◆



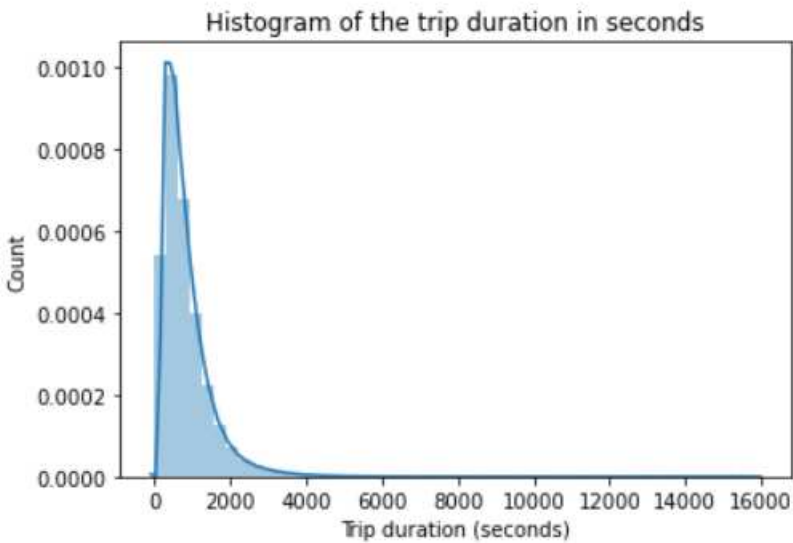
- ◆ From Google, we know New York City longitude vary from -74.03 to -73.75, and latitude vary from 40.63 to 40.85.
- ◆ In term of this, we drop outliers



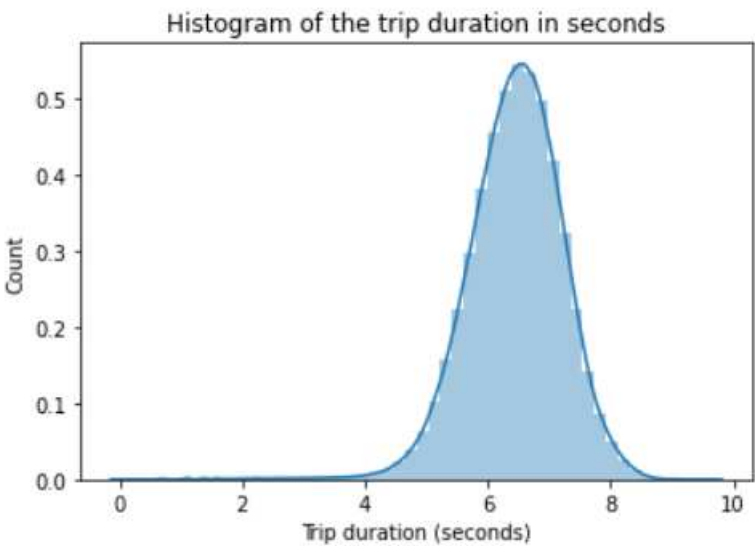
# Labels

- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

- trip\_duration:
  - ◆



- ◆ Through power conversion, it is more like Gaussian distribution.
- ◆







- Project Introduction
- Data Analysis
- Info
- Spearman Correlation
- Simple Attributes
- Complex Attributes
- Labels
- Model selection
- Result

- outliers:
  - ◆ Trip duration varies from 1 second to 3526282 second, and there are some outliers.
  - ◆ We keep the trip duration between  $\text{mean} - 3 \times \text{std}$  and  $\text{mean} + 3 \times \text{std}$ . (approximately 99% )



[Project Introduction](#)

[Data Analysis](#)

**[Model selection](#)**

[Models](#)

[Result](#)

# Model selection



# Models

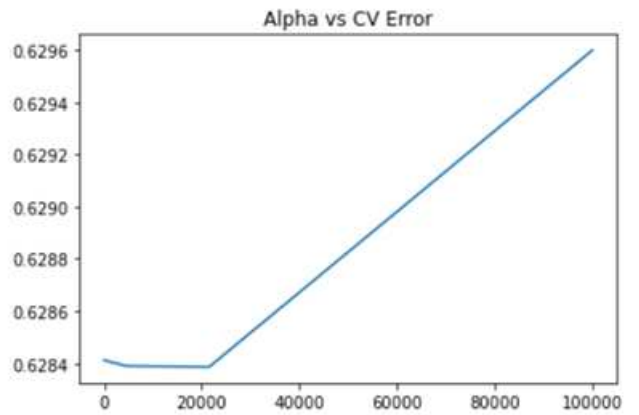
- [Project Introduction](#)
- [Data Analysis](#)
- [Model selection](#)
- [Models](#)
- [Result](#)

- Models:
  - ◆ Ridge
  - ◆ Bagging
  - ◆ Boosting
  - ◆ RandomForest
  - ◆ Lightgbm
  - ◆ Xgboost

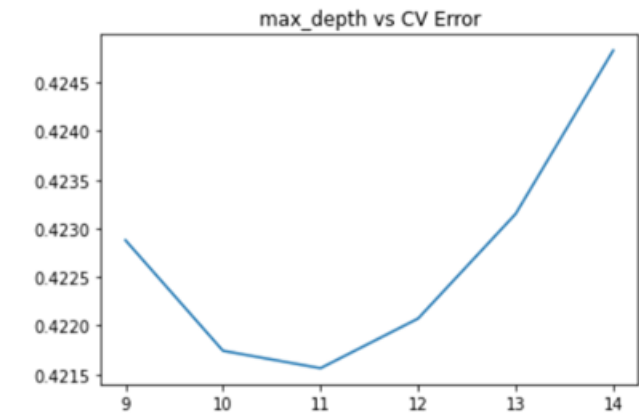


- [Project Introduction](#)
- [Data Analysis](#)
- [Model selection](#)
- [Models](#)**
- [Result](#)

- Ridge
- 



- Xgboost
- 





[Project Introduction](#)

[Data Analysis](#)

[Model selection](#)

**[Result](#)**

[Models](#)

# Results



# Results

- [Project Introduction](#)
- [Data Analysis](#)
- [Model selection](#)
- [Result](#)
- [Models](#)



<a href="#">results.csv</a>	0.59486	0.59596	<input type="checkbox"/>
9 days ago by Daylight Dream			
<a href="#">add submission details</a>			



<a href="#">results.csv</a>	0.42527	0.42716	<input type="checkbox"/>
a day ago by Daylight Dream			
<a href="#">add submission details</a>			



<a href="#">results.csv</a>	0.40978	0.41166	<input type="checkbox"/>
4 hours ago by Daylight Dream			
<a href="#">change XGB</a>			





# Contact Information

Undergraduate Siyi Yu  
School of Computer Science and Technology  
Jilin University, China

 YUSIYICSAT@163.COM

