**Data Glacier**

Your Deep Learning Partner

# G2M insight for Cab Investment firm

LISUM15

**2022.11.19**

# Introduction

**Problem Statement:**

XYZ is a private firm in US. Due to remarkable growth in the Cab Industry in last few years and multiple key players in the market, it is planning for an investment in Cab industry and as per their Go-to-Market(G2M) strategy they want to understand the market before taking final decision.

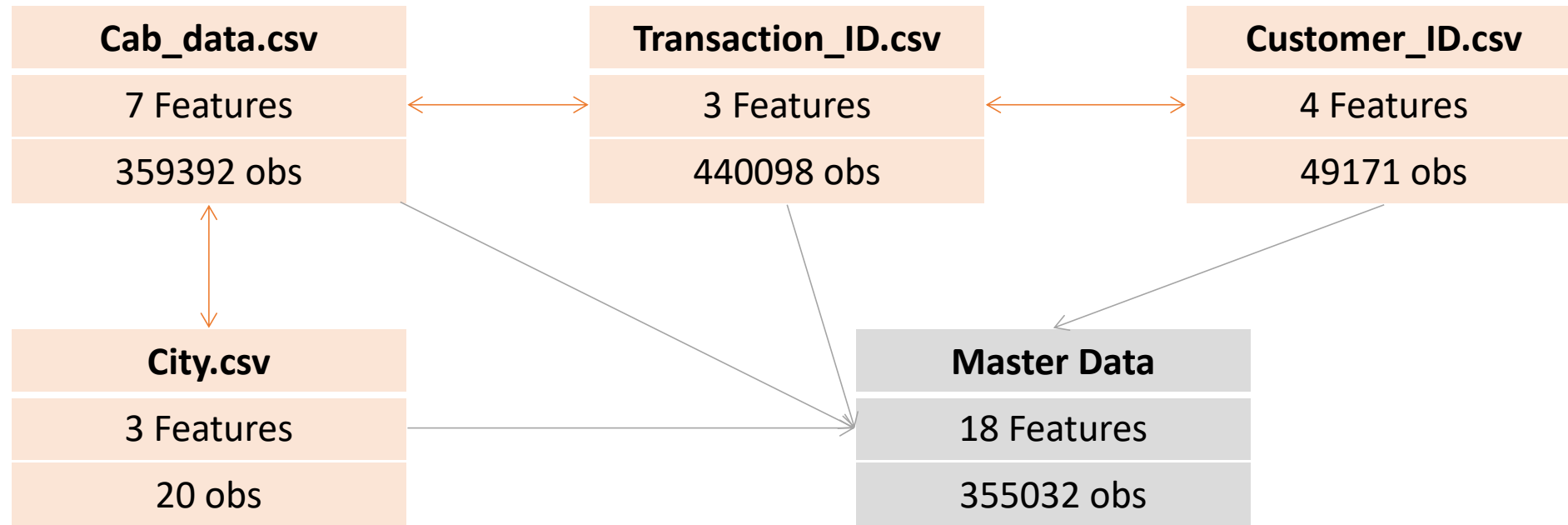**Analysis Outline:**

Understand data sets

Descriptive analysis of user groups

Explore potential factors for profit

Forecast the profitability of each cab

Recommendations

# Data Source

| Cab_data.csv | Transaction_ID.csv | Customer_ID.csv |
|:---:|:---:|:---:|
| 7 Features | 3 Features | 4 Features |
| 359392 obs | 440098 obs | 49171 obs |

| City.csv | Master Data |
|:---:|:---:|
| 3 Features | 18 Features |
| 20 obs | 355032 obs |

- The time period of this analysis is from 01/31/2016 to 12/31/2018.
- When creating master data, four new features are generated:
  Profit, Unit Price, Unit Cost and Number of active users

# Data Manipulation

## Data Cleaning

➤ Missingness - Use pandas.merge() and drop the records with missing values

➤ Duplication - Use df.drop_duplicates() to drop duplicate records

➤ Outlier - Use boxplot to detect outliers

## New Features:

$$\text{profit} = \text{price\_charged} - \text{cost\_of\_trip}$$

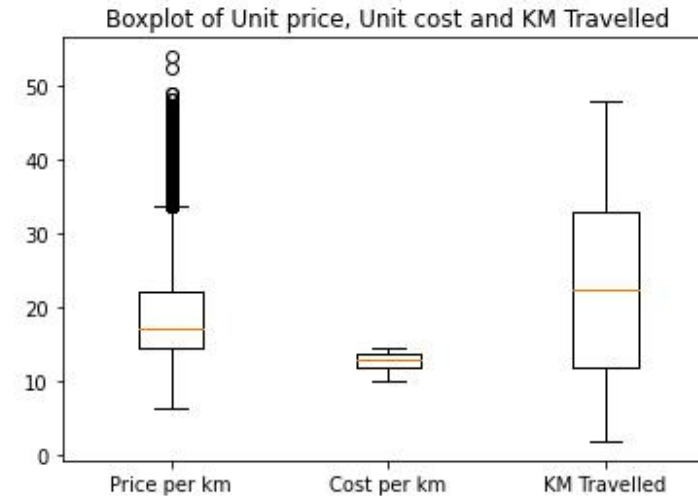$$\text{unit price} = \frac{\text{price\_charged}}{KM\_Travelled}$$

$$\text{unit cost} = \frac{\text{cost\_of\_trip}}{KM\_Travelled}$$

Active users: Users who made transactions during this time period

User_ratio: The proporation of cab users in the population of a city

$$User\_ratio = \frac{Number\_of\_Users}{Total\_population}$$

# Data Manipulation



Boxplot of Price & Cost

Boxplot of Unit price, Unit cost and KM Travelled

| Number of Outliers in Price Charged | |
| --- | --- |
| Yellow Cab | Pink Cab |
| 5861 | 18 |

Because outliers only exist in Price Charged, which may arise due to some abnormal operation and will interfere the correlation analysis between mileage and charge, so we drop these outliers, and the remaining master data is:
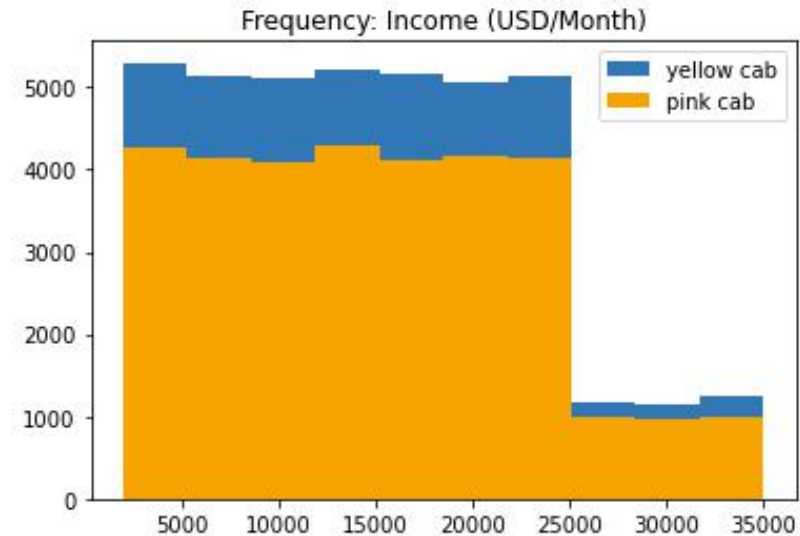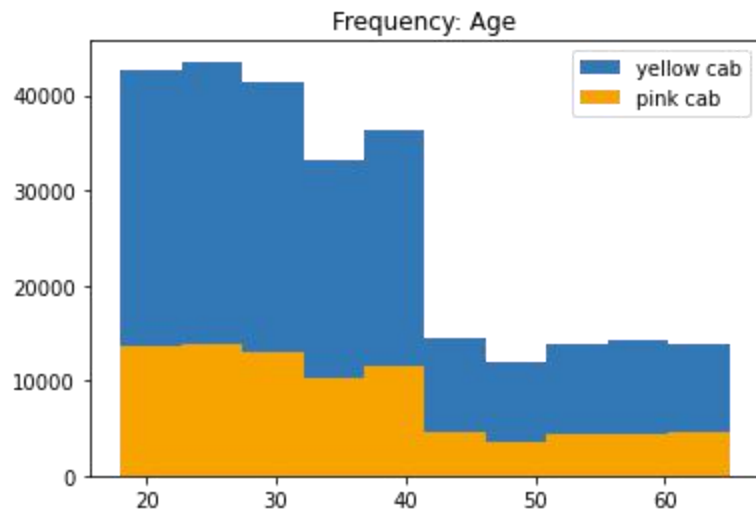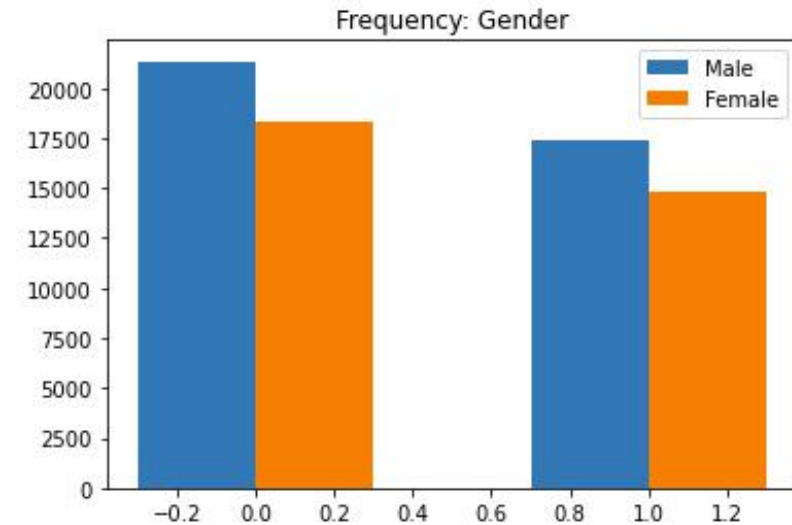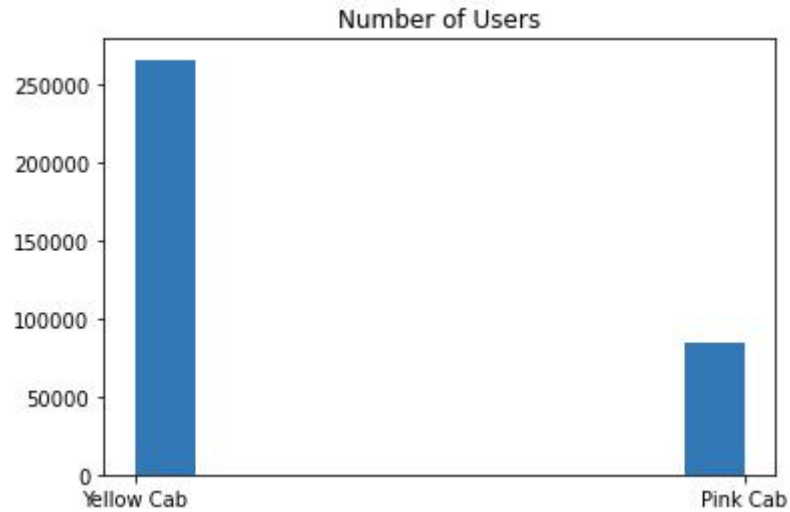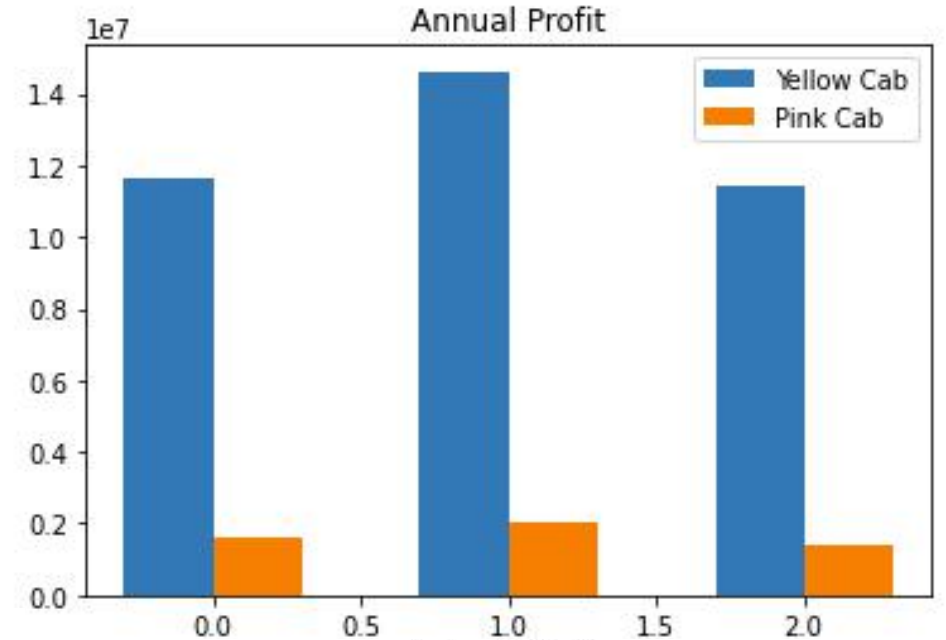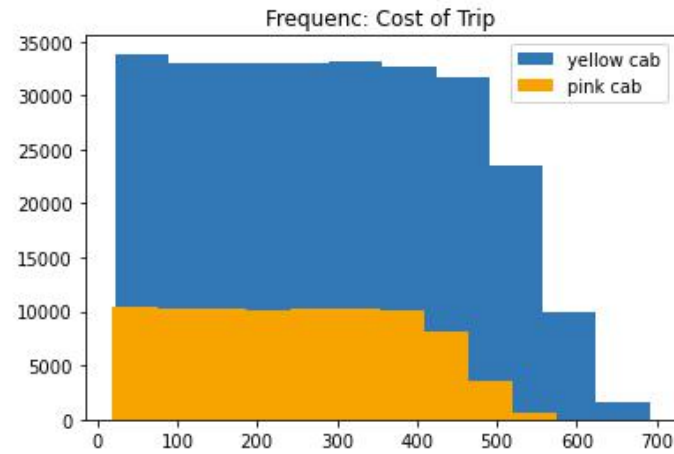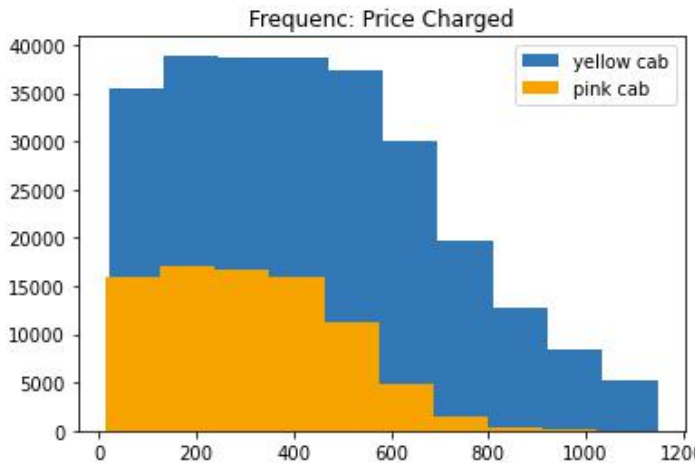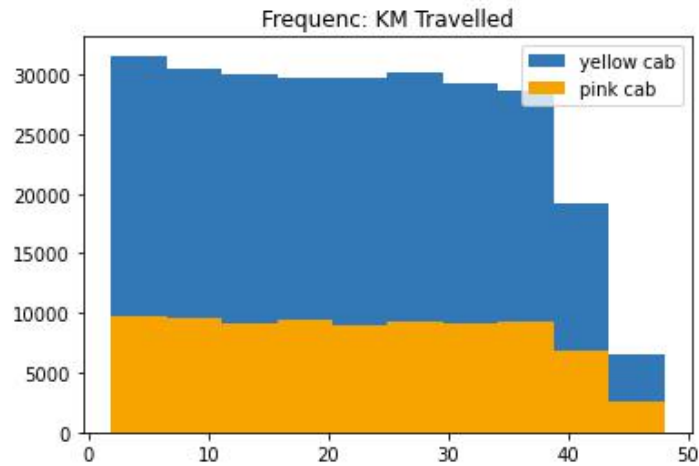
| Master Data |
| --- |
| 18 Features |
| 349153 obs |

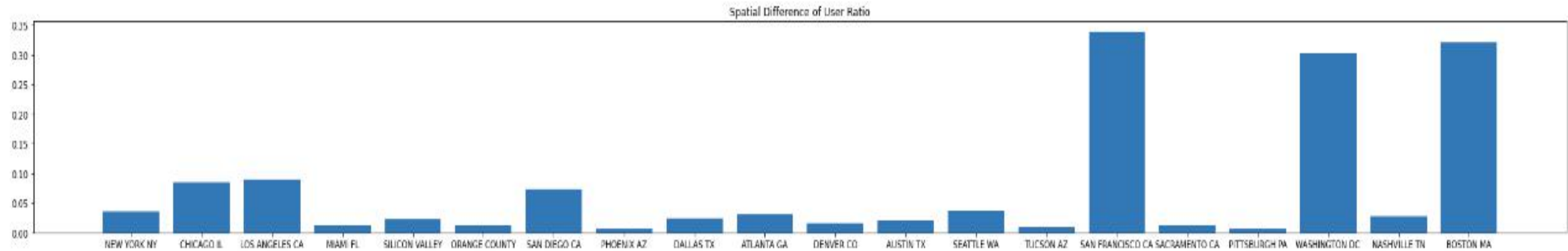There is no duplication and missingness existing in this master data.

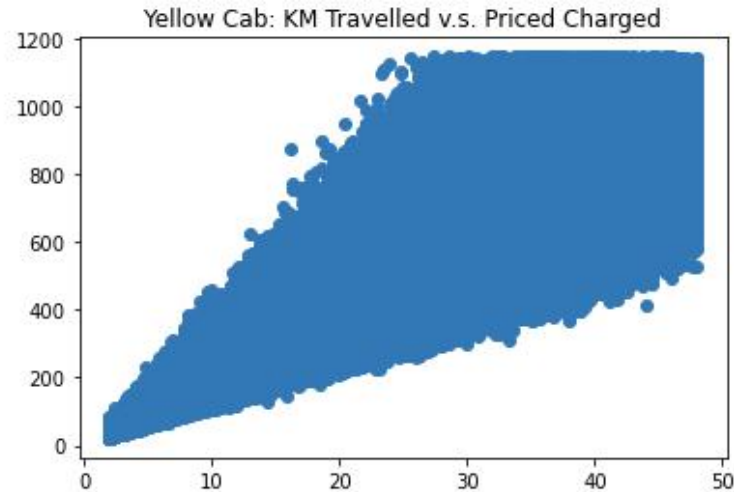# Descriptive analysis - Profitability

# Descriptive analysis - Spatial difference

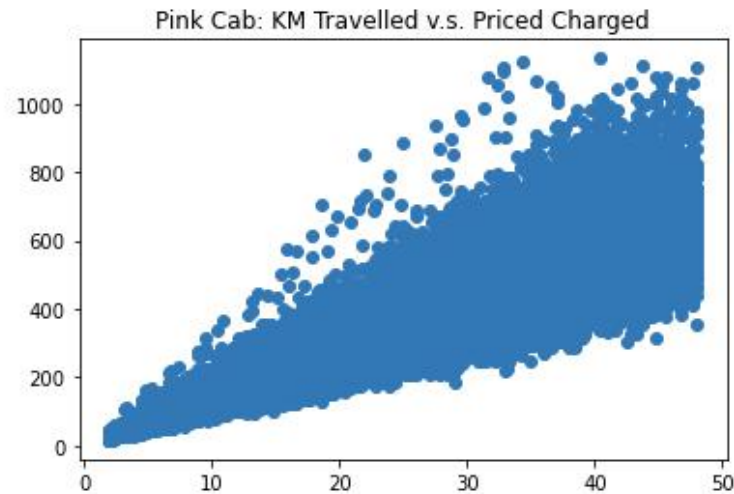We use User_ratio to represent the cab market scale of each city.



Spatial Difference of User Ratio

Data Glacier

# Regression analysis - Price



Yellow Cab: KM Travelled v.s. Priced Charged



Pink Cab: KM Travelled v.s. Priced Charged

**Yellow Cab:**

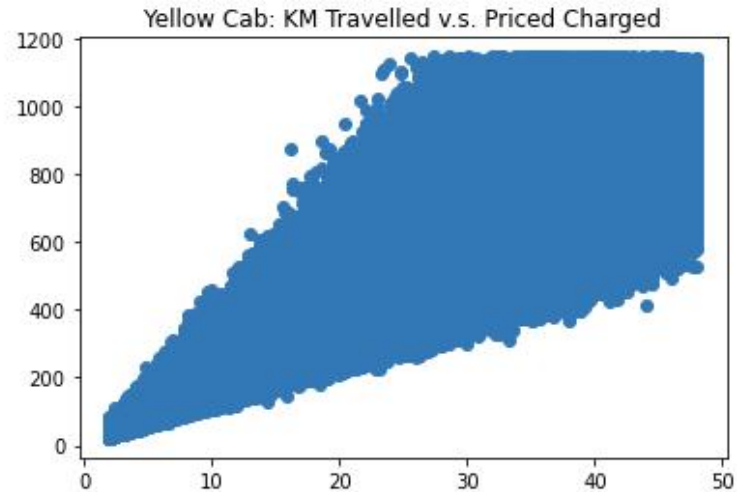$$Price\_Charged = 20.23 * KM\_Travelled + 0.88$$

$$R^2 = 0.74$$

**Pink Cab:**

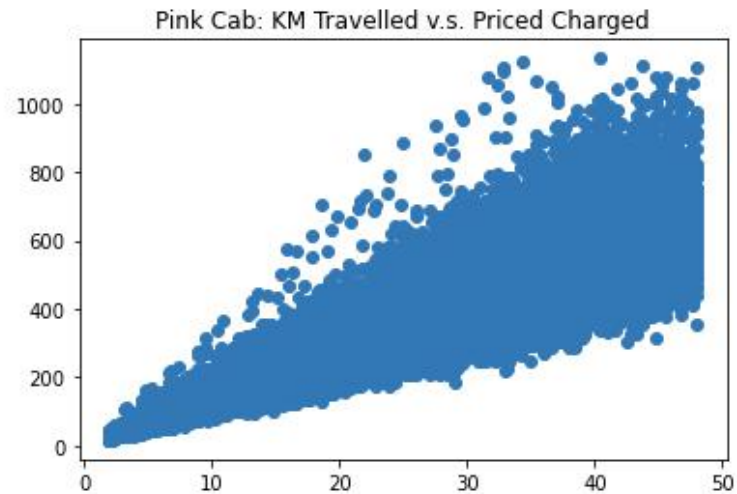$$Price\_Charged = 13.79 * KM\_Travelled - 0.55$$

$$R^2 = 0.86$$

# Regression analysis - Cost



**Yellow Cab:**

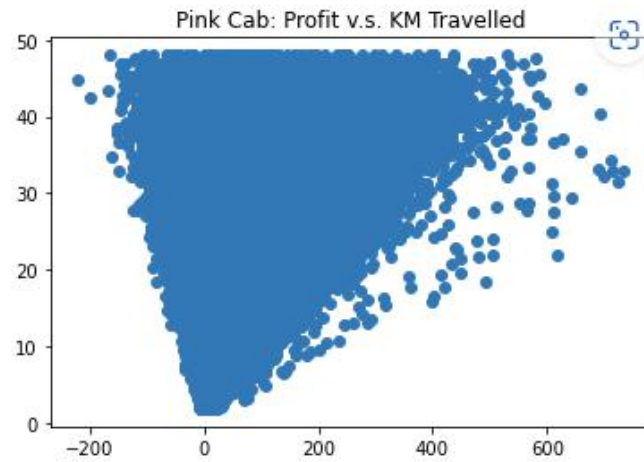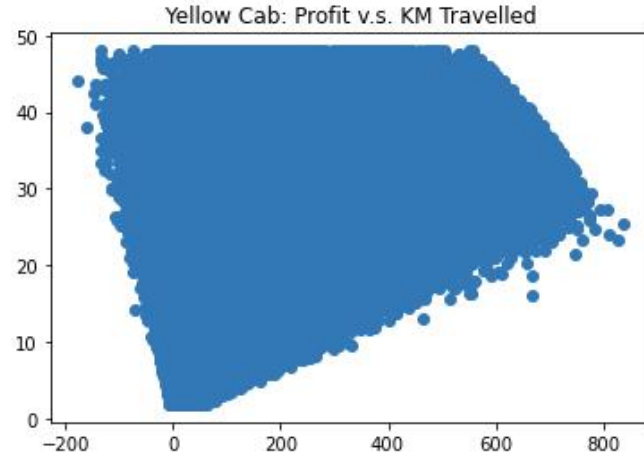$$Cost\_of\_Trip = 13.20 * KM\_Travelled + 0.03$$

$$R^2 = 0.987$$

**Pink Cab:**

$$Cost\_of\_Trip = 11.00 * KM\_Travelled + 0.05$$

$$R^2 = 0.987$$

# Causality Analysis - Profit


Yellow Cab: Profit v.s. KM Travelled


Pink Cab: Profit v.s. KM Travelled

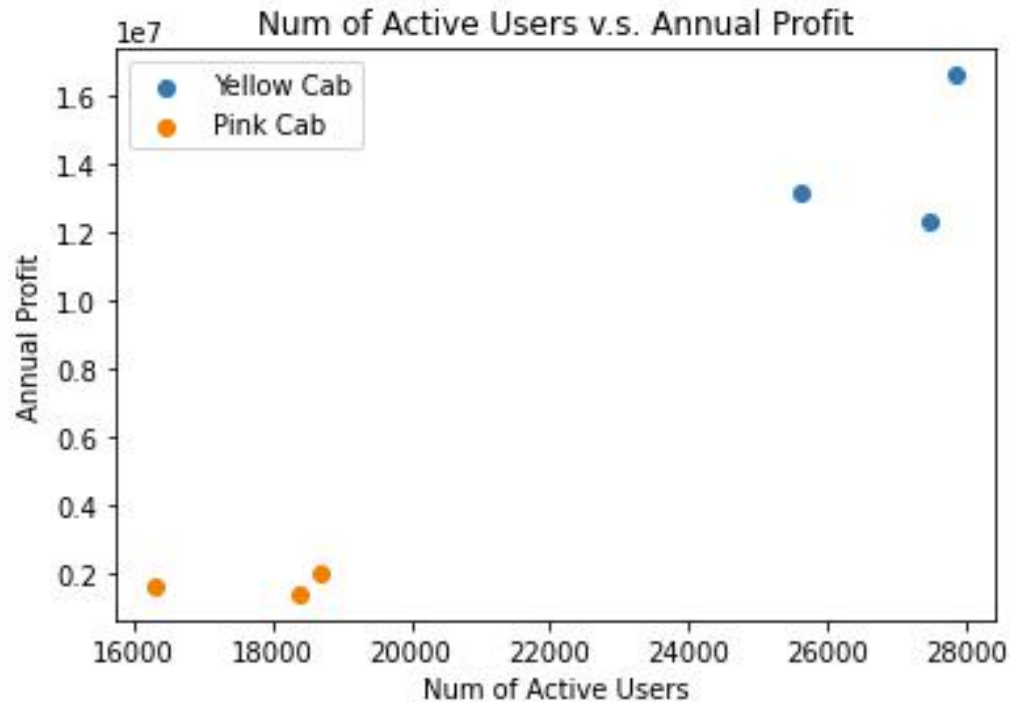| Pearson Correlation Coefficient Profit v.s. KM Travelled | |
| --- | --- |
| Yellow Cab | Pink Cab |
| 0.502 | 0.442 |

**Conclusion:**

Profit is correlated to KM_Travelled, but not strong, so there are some other factors affecting the profit of each trip.

Data Glacier

# Causality Analysis - Profit



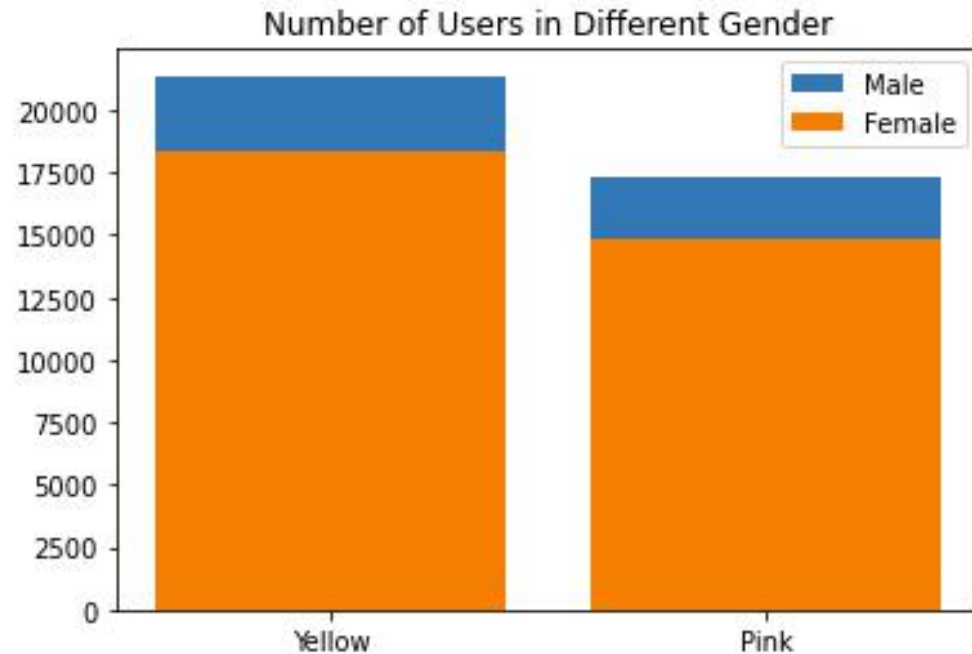| Pearson Correlation Coefficient Profit v.s. Num of Active Users | |
|:---:|:---:|
| Yellow Cab | Pink Cab |
| 0.46 | 0.23 |

**Conclusion:**

Profit is correlated to Number of Annual Active Users, but not strong.

**Data Glacier**
Your Deep Learning Partner

# Hypothesis Test - Gender preference



Number of Users in Different Gender

**Contingency Table:**

|        | Male  | Female |
|--------|-------|--------|
| Yellow | 21376 | 18379  |
| Pink   | 17363 | 14811  |

**Chi square test result:**

p-value = 1  >> 0.05

Conclusion: there is no gender preference of choosing cabs.

# Recommendation

# Thank You