# CS6502

# Applied Big Data & Visualization

## ML Project

**KAILASH MURALIDHARAN          19116608**

Your task here is to

1.  use the python scikit-learn to investigate the similarities and relationships that PCA can provide. It will not be possible for you to analyze the entire data set in this way so you should perform your analysis on *the first 1000 rows* of the table. Ideally the output of this phase should guide the direction you take in the second, ML, part of the assignment. You will need to decide which columns should be part of your analysis and which should be ignored. See some of the tutorial links we have posted in the lecture slides for help in deciding what are appropriate columns to consider.

2.  based on the information you gained from step 1, create a model[1] to predict the taxi fare ("fare" column in the dataset). Note that you may need to clean the data, pick a list of features (feature engineering), and then design your model.

Please email your zipped solution pack to the lecturer by the end of week 14 with subject "cs6502: ml proj".

Note that your solution pack should contain

- the query you use to selected the first 1000 rows
- key steps for your PCA (setup, command, etc)
- the sql for model creation
- the sql to evaluate the model
- the sql to predict using the model above
- link of the BigQuery commands you composed[2]
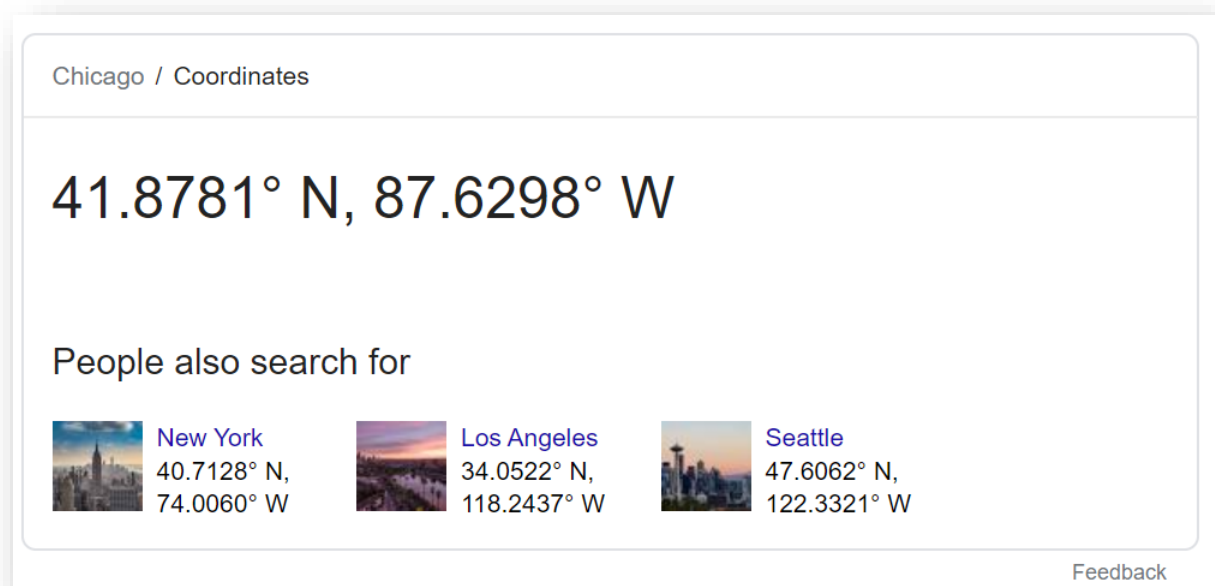- screenshot of the model evaluation report[3]

# Summary

| | Item | Links |
|---|---|---|
| 1. | Query used to select 1000 rows | https://console.cloud.google.com/bigquery?sq=3103539222:ed4b7d29551040a1a3441a2b2cf59031 |
| 2. | Results of the query containing 1000 rows | https://docs.google.com/spreadsheets/d/1xc51HUUrw51UBCwIfbpPT343E8GyODNaQc2khHywu30/edit?ts=5eb593ba#gid=2010390566 |
| 3. | Query used to create the Final Model | https://console.cloud.google.com/bigquery?sq=3103539222:88ef2a11074c451db05d9e1598f90d55 |
| 4. | Query used for Model Evaluation | https://console.cloud.google.com/bigquery?sq=3103539222:4c44667fd77d489e8a8164b3a02b843c |
| 5. | Query used for Model Prediction | https://console.cloud.google.com/bigquery?sq=3103539222:4be462936e3c40eba4f80b02cc69c4b4 |
| 6. | Final Model Evaluation Report |  |

For item 6:

Details    Training    Evaluation    Schema

| | |
|---|---|
| Mean absolute error | 1.8286 |
| Mean squared error | 11.1206 |
| Mean squared log error | 0.0143 |
| Median absolute error | 1.1597 |
| R squared | 0.9421 |

# 1. Query to select 1000 rows

I have used a query to check for the minimum, maximum and average of three numeric columns.



Using the latitude and longitude values of Chicago, to consider only those trips that started and ended inside the city range.



**Conditions applied:**

- Picked values to apply for where conditions based on the min, max, avg summary and Chicago's location coordinates.
- Years filtered between 2016 and 2018, and months between January and September.
- Rest of the columns – NOT NULL condition
- In addition to the existing columns, following additional columns have been created
    - **total_fare_without_tips** (the value to be predicted), since the tips amount varies from customer to customer (created a column by subtracting tips from the trip_total)

- **hour_of_day**, from the timestamp
- **month_of_trip**, also from the timestamp
- **euclidean_distance**, the distance between the pickup and dropoff points
- **longitude,** the distance between the pickup and dropoff points longitude
- **latitude,** the distance between the pickup and dropoff points latitude
- **taxi_company**, the name of the taxi company

## 2. Base Model before PCA



## Evaluation Report:



- MSE= **14.3731**
- **R squared value** = **0.9258**, about **92.58**% of the variability in the dependent variable is explained by our model.

# 3. PCA using R statistical package

- Load the dataset into the dataframe df

```
1  df<-read.csv("/Users/kailashm/Downloads/BQML_ChicagoTaxiFares - 1000 row
2  str(df)
3
```

| ▶ df | 1000 obs. of 30 variables | ⊞ |
|------|---------------------------|---|

- The taxi_company column is changed to be of numeric data type (to be included in creating the model)

```
df$taxi_company = as.numeric(df$taxi_company)
str(df)
```

```
$ trip_miles          : num  17 16.5 7.5 4.2 6.5 5.4 5.9 4.9 4.6 4.7 ...
$ pickup_community_area  : int  76 76 32 6 6 6 6 6 6 6 ...
$ dropoff_community_area : int  28 8 6 8 32 28 28 8 8 8 ...
$ pickup_latitude     : num  42 42 41.9 41.9 41.9 ...
$ pickup_longitude    : num  -87.9 -87.9 -87.6 -87.7 -87.7 ...
$ dropoff_latitude    : num  41.9 41.9 42 41.9 41.9 ...
$ dropoff_longitude   : num  -87.6 -87.6 -87.7 -87.6 -87.6 ...
$ taxi_company        : num  18 18 18 21 21 20 21 16 16 20 ...
```

- Correlation Matrix

```
install.packages("ggcorrplot")
library(ggcorrplot)
correlation <- round(cor(df[,1:13]), 2)
ggcorrplot(correlation, lab=TRUE, colors = c("red", "white", "oran
```

The below matrix shows the correlation between the numeric variables in the dataset.



- Computing PCA using prcomp()

```r
install.packages("factoextra")
library(factoextra)
df.pca <- prcomp(df[,2:13], scale = TRUE)
fviz_eig(df.pca,addlabels = TRUE)
```

**This plot shows the percent contribution of the principle components.** We can see that the first two principle components explain about **57.5%** variation.



- We visualize the first two principle components to check the **% contribution of each variable** (features from the original dataset) **in both the principle components**.

```
fviz_contrib(df.pca, choice = "var", axes = 1:2, top = 12)
```

**Output Screenshot:**



Contribution of variables to Dim-1-2

```
> names(df[,2:13])
 [1] "hourofday"                "monthoftrip"
 [3] "euclidean_dist"           "trip_seconds"
 [5] "trip_miles"               "pickup_community_area"
 [7] "dropoff_community_area"   "pickup_latitude"
 [9] "pickup_longitude"         "dropoff_latitude"
[11] "dropoff_longitude"        "taxi_company"
```

# 4. Final Model after PCA and feature selection



**Final Prediction Model - Evaluation Report:**



- The MSE of the prediction model **after PCA,** is **11.1206**, **which is less than** that of the **base model MSE = 14.3731**
- The **R squared value** of the prediction model has **increased** from **0.9258** to **0.9421**, which suggests that the final model captures more variation in the **total_fare_without_tips** than the base model

# 5. Final Model - Evaluation Query



## Output Screenshot:

# 6. Final Model - Prediction Query

The prediction model uses training data from months 1 to 9 and the test data from months 10, 11, 12.



Below screenshot shows total_fare_without_tips and the predicted_total_fare_without_tips using the final prediction model.