Building Serverless ML Platform

# Data Analysis And Visualization Of Green Taxi Records

SRI HARSHA KORUKONDA BHATTAR (1953040)

AASHISH JOSHI (1886750)

ANKITH KANDALA (1894221)

# Contents:

➢**Introduction**

➢**Serverless Architecture**

➢**Methodology**

➢**Challenges**

➢**Future work**

# INTRODUCTION

**Building a serverless Machine Learning platform that analyses, visualizes Green-taxi user data and predicts future customer base.**
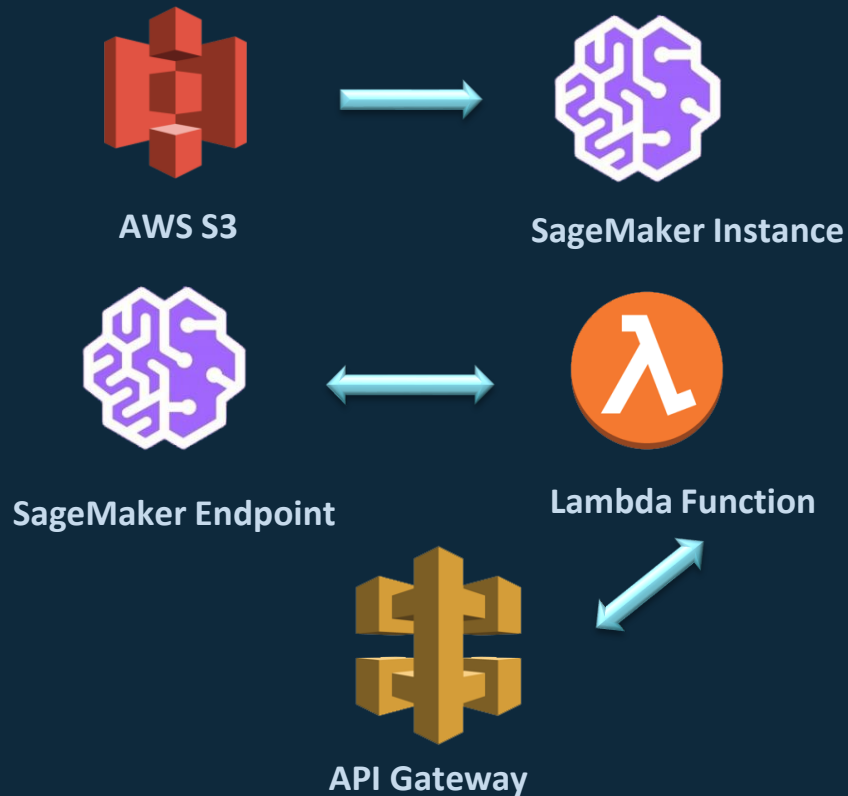
# SERVERLESS ARCHITECTURE

- ➢ Serverless is the native architecture of the cloud that enables you to shift more of your operational responsibilities to AWS.
- ➢ Serverless enables us to build modern applications with increased agility and lower total cost of ownership.

# METHODOLOGY

**AWS Tools used**:

- ➢ Amazon S3
- ➢ Lambda function
- ➢ Sagemaker Instance
- ➢ Sagemaker Endpoint
- ➢ AWS API

**AWS S3** → **SageMaker Instance**

**SageMaker Endpoint** ↔ **Lambda Function**

**API Gateway**

- ◇ 6 years Green taxi Trip records for the state of New York have been collected.

- ◇ Monthly Raw data files are stored in S3.

- ◇ We have implemented Python to organize and analyze the Raw data.

- ◇ ML Algorithms – Multiple Linear Regression.

```python
In [158]: data = pd.read_csv('https://fianlprojectcc.s3-us-west-1.amazonaws.com/2014-01.csv')
          data.columns = ['lpep_pickup_datetime', 'Lpep_dropoff_datetime', 'Store_and_fwd_flag', 'RateCodeID', 'Pickup_longitude', 'Pickup_
          data.to_csv("data.csv", sep=',', index=False)
          columns = ['lpep_pickup_datetime', 'Lpep_dropoff_datetime', 'Total_amount', 'Trip_distance']
          df = pd.read_csv("data.csv", usecols = columns)
          df=df.dropna()
          #print(df)

          df['lpep_pickup_datetime'] = pd.to_datetime(df['lpep_pickup_datetime'],format='%m/%d/%Y %H:%M')
          df['Lpep_dropoff_datetime'] = pd.to_datetime(df['Lpep_dropoff_datetime'],format='%m/%d/%Y %H:%M')
          df['trip_duration']= df['Lpep_dropoff_datetime'] - df['lpep_pickup_datetime']
          data = df[~(df['lpep_pickup_datetime'] < '2014-01-01')]
          data = data[~(df['lpep_pickup_datetime'] > '2014-02-01')]
          data['count']=1
          #print(data)
          #print(data)

          reqcolumns = ['Trip_distance', 'Total_amount']
          reqcolumns2 = ['trip_duration', 'Total_amount']
          data_distance=data[reqcolumns]
          data_time = data[reqcolumns2]
          #print(data_time)

          by_hour=data.groupby(Grouper(key='lpep_pickup_datetime', freq='H')).sum()
          by_days=data.groupby(Grouper(key='lpep_pickup_datetime', freq='d')).sum()

          print(by_days)

          # data atribute has non group data
```
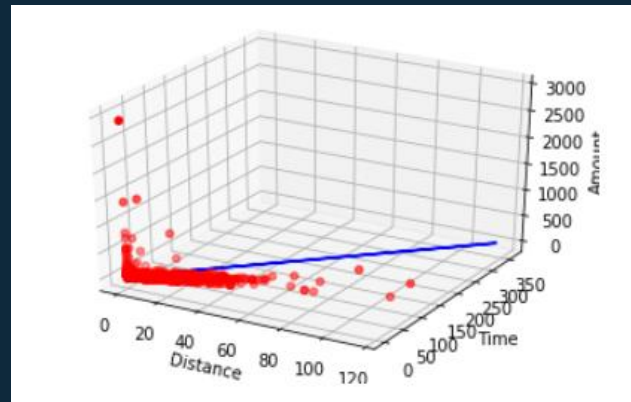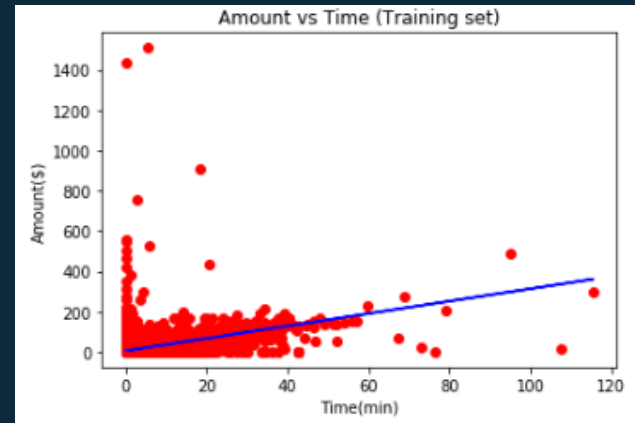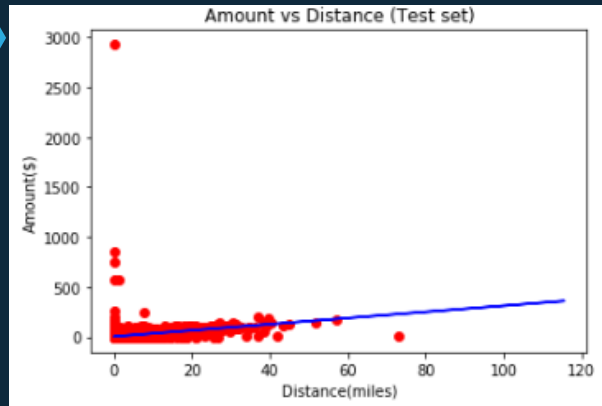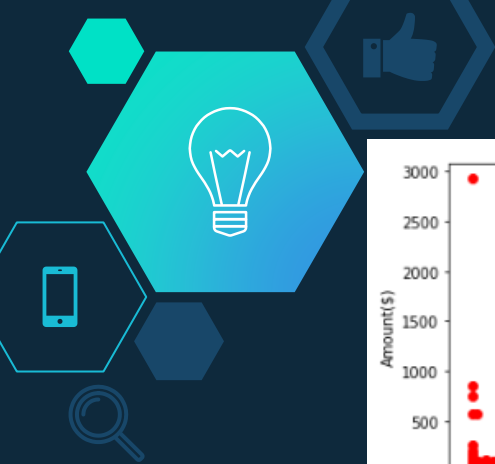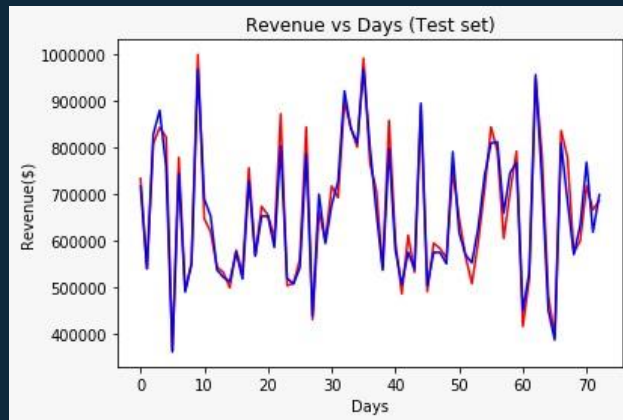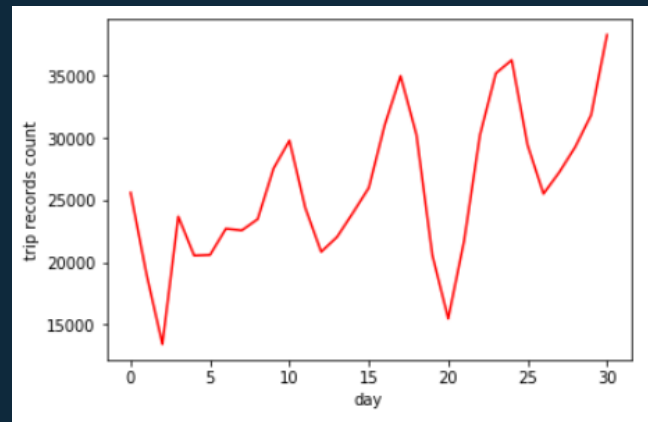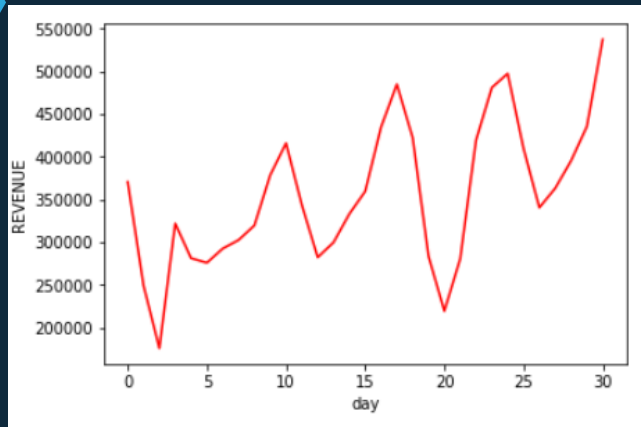
Amount vs Distance (Test set)

Amount vs Time (Training set)

**Prediction Accuracy for trip pricing – 92.4%**
**Prediction Accuracy for revenue and customer base – 86.8%**

# Greentaxiroject

Throttle  |  Qualifiers ▼  |  Actions ▼  |  Select a test event ▼  |  Test  |  Save

▼ 📁 Greentaxiroject  ⚙▼
  📄 lambda_function.py

lambda_function ✕  ⊕

```python
1  import os
2  import io
3  import boto3
4  import json
5  import csv
6  # grab environment variables
7  ENDPOINT_NAME = os.environ['ENDPOINT_NAME']
8  runtime= boto3.client('runtime.sagemaker')
9  def lambda_handler(event, context):
10     print("Received event: " + json.dumps(event, indent=2))
11
12     data = json.loads(json.dumps(event))
13     payload = data['data']
14     print(payload)
15
16     response = runtime.invoke_endpoint(EndpointName=ENDPOINT_NAME,
17                                        ContentType='text/csv',
18                                        Body=payload)
19     print('res is !!!!')
20     print(response)
21     result = json.loads(response['Body'].read().decode())
22     print('result is!!!!')
23     print(result)
24     pred = int(result['predictions'][0]['score'])
25     #pred = int(result['predictions'][0])
26     predicted_label = 'delay' if pred == 1 else 'no delay'
27
28     return predicted_label
```

22:27  Python  Spaces: 4  ⚙

## Environment variables (1)

The environment variables below are encrypted at rest with the default Lambda service key.

Edit

| Key | Value |
| --- | --- |
| ENDPOINT_NAME | greentaxiproject |

# CHALLENGES

◇ Column names and data types were different for different monthly reports so preprocessing the data for analysis was a challenging job

◇ We were unable to integrate Location specific data in our predictions.

# FUTURE WORK

◇ We would like to use API collaboration platforms like Postman and provide two user interfaces. We would also integrate location data into our analysis.

◇ User-side Application.

◇ Company Interface.

# Thank You!

**Any questions?**