# Predicting Taxi Driver Payments: A Big Data Pipeline

# Team Members

**01**   **Anastasia Pichugina** - Data Collection & Preparation, EDA

**02**   **Arthur Gubaidullin** - Feature Engineering, ML Modeling

**03**   **Israel Adewuyi** - Dashboard, Charts

**04**   **Anna Gromova** - Report writing, GitHub maintenance

# Project Goal

**01**  **Problem Statement**

Taxi companies need to balance driver earnings with customer pricing

**02**  **Objective**

Build a model to predict *driver_pay* based on trip features to:
- Detect unfair payments
- Optimize dynamic pricing
- Improve driver satisfaction

**03**  **Business Impact**

Driver churn reduction potential

# Datasets

▲ 45   ⟨⟩ Code   ⬇ Download

## Uber NYC for-hire vehicles trip data (2021)

NYC for-hire vehicle trip data, taxi zones zip code lookup, NYC weather 2021

Data Card   Code (8)   Discussion (2)   Suggestions (0)

**01   Trips Data (Parquet files)**
- 2.4M trips (January)
- Key fields: driver_pay, timestamps, locations

**02   Taxi Zones (CSV)**
- 265 zones
- Fields: borough, zone, service_zone

**03   Weather (CSV)**
- 356 days
- Key fields: temperature, humidity, snow, wind, uv, visibility, etc.

# Data Collection & Ingestion

- Converted collected from Parquet to CSV
- Designed PostgreSQL schema with optimized tables (*trips, taxi_zones, weather*)
- Ingested data to HDFS using Sqoop with Avro and Snappy compression

### taxi_zones

| | | |
|---|---|---|
| int | location_id | PK |
| varchar | borough | |
| varchar | zone | |
| varchar | service_zone | |

### weather

| | | |
|---|---|---|
| date | date_id | PK |
| varchar | station_name | |
| varchar | station_address | |
| varchar | resolved_address | |
| float | temperature | |
| float | feels_like | |
| float | dew_point | |
| float | humidity | |
| float | precipitation | |
| int | precipitation_prob | |
| varchar | precipitation_type | |
| float | snow | |
| float | snow_depth | |
| float | wind_gust | |
| float | wind_speed | |
| float | wind_direction | |
| float | sea_level_pressure | |
| float | cloud_cover | |
| float | visibility | |
| int | uv_index | |
| float | severe_risk | |

### trips

| | | |
|---|---|---|
| bigserial | trip_id | PK |
| varchar | hvfhs_license_num | |
| varchar | dispatching_base_num | |
| varchar | originating_base_num | |
| timestep | request_datetime | |
| timestep | on_scene_datetime | |
| timestep | pickup_datetime | |
| timestep | dropoff_datetime | |
| int | pu_location_id | FK |
| int | do_location_id | FK |
| float | trip_miles | |
| int | trip_time | |
| float | base_passenger_fare | |
| float | tolls | |
| float | bcf | |
| float | sales_tax | |
| float | congestion_surcharge | |
| float | airport_fee | |
| float | tips | |
| float | driver_pay | |
| char(1) | shared_request_flag | |
| char(1) | shared_match_flag | |
| char(1) | wav_request_flag | |
| char(1) | wav_match_flag | |
| date | request_date | FK |
| time | request_time | |

| | station_name character varying (50) | station_address character varying (100) | resolved_address character varying (100) | date_id [PK] date | temperature double precision | feels_like double precision | dew_point double precisi |
|---|---|---|---|---|---|---|---|
| 1 | nyc | nyc | New York, NY, United States | 2021-07-17 | 27.5 | 29.3 | |
| 2 | nyc | nyc | New York, NY, United States | 2021-09-16 | 23.3 | 23.3 | |
| 3 | nyc | nyc | New York, NY, United States | 2021-11-06 | 8.1 | 7.3 | |
| 4 | nyc | nyc | New York, NY, United States | 2021-10-11 | 19 | 19 | |
| 5 | nyc | nyc | New York, NY, United States | 2021-03-15 | -0.2 | -5.7 | |

| | location_id [PK] integer | borough character varying (50) | zone character varying (100) | service_zone character varying (50) |
|---|---|---|---|---|
| 1 | 81 | Bronx | Eastchester | Boro Zone |
| 2 | 220 | Bronx | Spuyten Duyvil/Kingsbridge | Boro Zone |
| 3 | 203 | Queens | Rosedale | Boro Zone |
| 4 | 159 | Bronx | Melrose South | Boro Zone |
| 5 | 178 | Brooklyn | Ocean Parkway South | Boro Zone |

# Data Storage & Preparation

- Created Hive tables with:
  - Partitioning by *dispatching_base_num*
  - Bucketing by location IDs (8 buckets)
- Implemented dynamic partitioning for future data

# Analysis Results

- HiveQL for aggregations (trip counts, avg. payments)
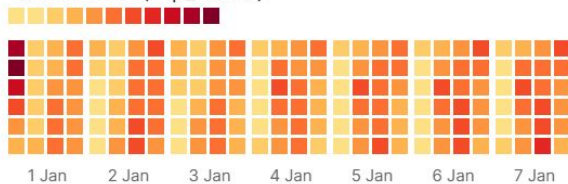- Spark SQL for complex joins (trips + weather)
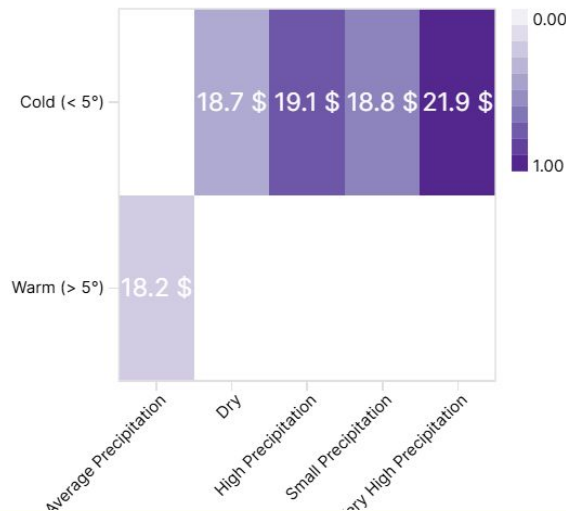
## Average number of trips per day

# 307k

Trips per day

## Heatmap: trips count by hours

Metric: SUM(trip_count)

1 Jan  2 Jan  3 Jan  4 Jan  5 Jan  6 Jan  7 Jan

## Heatmap with prices

| | Average Precipitation | Dry | High Precipitation | Small Precipitation | Very High Precipitation |
|---|---|---|---|---|---|
| Cold (< 5°) | | 18.7 $ | 19.1 $ | 18.8 $ | 21.9 $ |
| Warm (> 5°) | 18.2 $ | | | | |

0.00

1.00

# Analysis Results

## Companies distribution

HV0004: 0.95%

HV00...

HV00...

- HV0003
- HV0005
- HV0004
- (All) (Inv)

## Most popular destinations

- NA
- Crown Heig ◀ 1/20 ▶ (All) (Inv) ⊡ ⊏

350k
300k
250k
200k
150k
100k
50k
0

Bronx | Manhattan | Unknown

## Tips

SUM(trip_count)
120k
95.6k
75.7k

13.45
5.08
6.29

13.11
16.43
19.64

(avg_miles)

SUM(avg_

- No Tip
- Medium Tip
- High Tip
- Low Tip
- (All) (Inv)

# Analysis Results



Trip cost per hour

Trip duration by borough

Bronx: 14.19

EWR:...

St...

...

Qu...

Manhattan...

Brooklyn: 15.89

EWR  Queens  Brool  ◄ 1/4 ►  (All) (Inv)

# Feature Engineering

- Engineered features:
  - Cyclical encoding for *day/month* (sine/cosine)
  - One-Hot Encoding for *borough/service_zone/license_num*
- Dropped high-cardinality and irrelevant columns
- Handled missing values
- Joined tables

# ML Stage

- Trained/tested (70/30 split):
  - Train/test sizes: (1714935, 734381)
  - Linear Regression
  - Gradient Boosted Trees
- Optimized via 3–fold CV with grid search
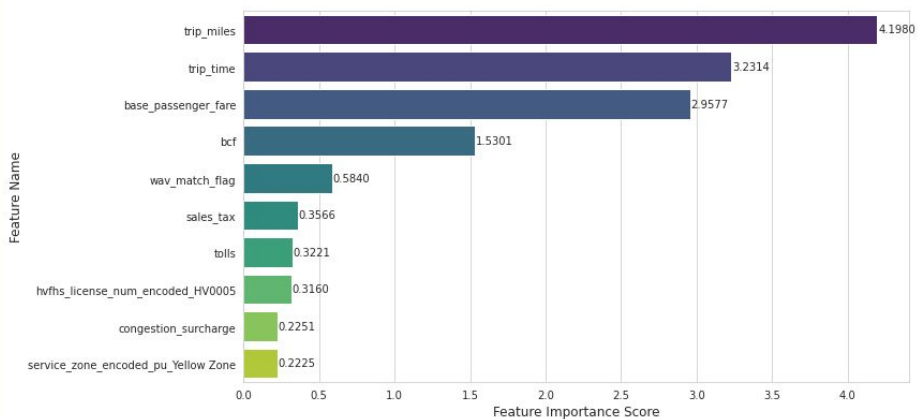
Key results:
- Best model: Linear Regression (R² 0.934, RMSE 2.84)
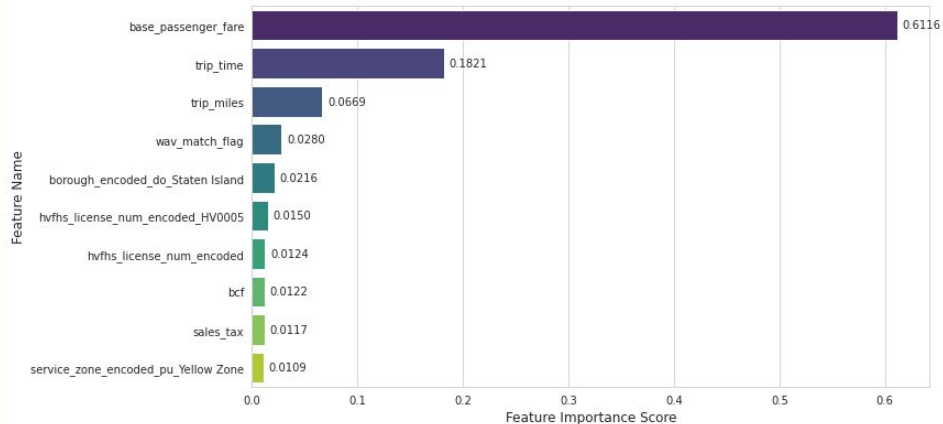- Outperformed GBT (R² 0.911, RMSE 3.37)

| Model | RMSE | MAE | R² |
|---|---|---|---|
| Linear Regression | 2.84 | 1.42 | 0.934 |
| GBT | 3.37 | 1.40 | 0.911 |

# Feature Importance Analysis



Top 10 Most Important Features (LR Model)

| Feature Name | Feature Importance Score |
|---|---|
| trip_miles | 4.1980 |
| trip_time | 3.2314 |
| base_passenger_fare | 2.9577 |
| bcf | 1.5301 |
| wav_match_flag | 0.5840 |
| sales_tax | 0.3566 |
| tolls | 0.3221 |
| hvfhs_license_num_encoded_HV0005 | 0.3160 |
| congestion_surcharge | 0.2251 |
| service_zone_encoded_pu_Yellow Zone | 0.2225 |

Top 10 Most Important Features (GBT Model)

| Feature Name | Feature Importance Score |
|---|---|
| base_passenger_fare | 0.6116 |
| trip_time | 0.1821 |
| trip_miles | 0.0669 |
| wav_match_flag | 0.0280 |
| borough_encoded_do_Staten Island | 0.0216 |
| hvfhs_license_num_encoded_HV0005 | 0.0150 |
| hvfhs_license_num_encoded | 0.0124 |
| bcf | 0.0122 |
| sales_tax | 0.0117 |
| service_zone_encoded_pu_Yellow Zone | 0.0109 |

# Challenges & Solutions

- Memory constraints during data preprocessing & fine-tuning:
  - Solved by using Spark distributed system;
  - Solved by using Gradient-based optimization model.
- Weather data alignment (dates format):
  - Addressed via date features preprocessing and careful joining of datasets.

Demo