MSE
MASTER OF SCIENCE IN ENGINEERING

# New Taxi Agencies

HES-SO, Big Data Analytics

Romain Claret, Jämes Ménétrey and Damien Rochat

MASTER OF SCIENCE
IN ENGINEERING

# Dataset

Provided by **Chris Whong**, Urbanist, New York City

Trips and fares taxis in New York City (2013)

2x 12 CSV files (one per month), total of **173'179'759** records

Trips dataset: Taxi identification, start and end time, GPS coordinates

Fares dataset: Taxi identification, start time, amount, payment type, tip

# Goals

In order to open new taxi agencies and increase customer satisfaction (and profit):

(1)    Find the best spots to find customers
(2)    Identify the most profitable spots

# Preprocessing

Removed out of bound GPS coordinates (1.75% of the records)

Removed invalid fares (1% of the records)

# Machine Learning

PySpark

K-means, centroids for clusters of best locations for new agencies
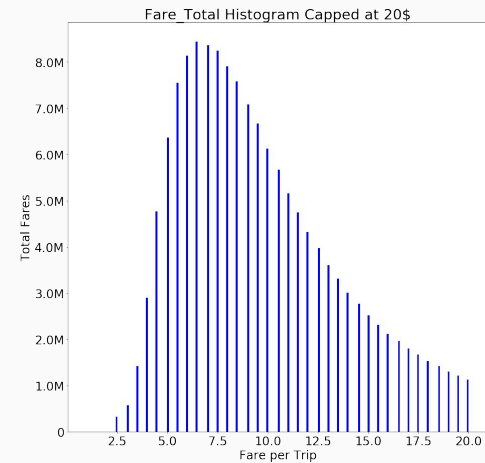
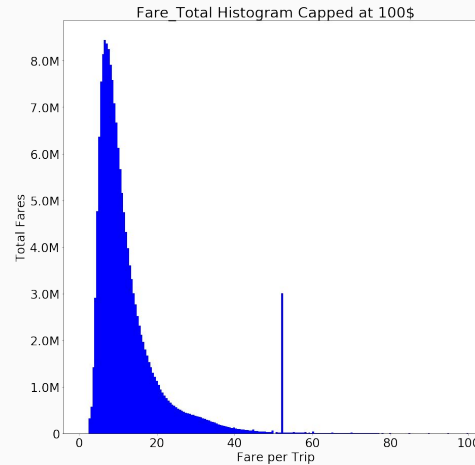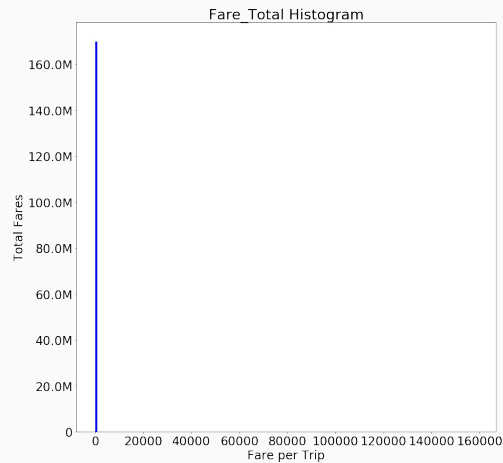Customers satisfaction:
- Using pickups coordinates

Profit:
- Using pickups coordinates
- Weighted with fare

# Machine Learning

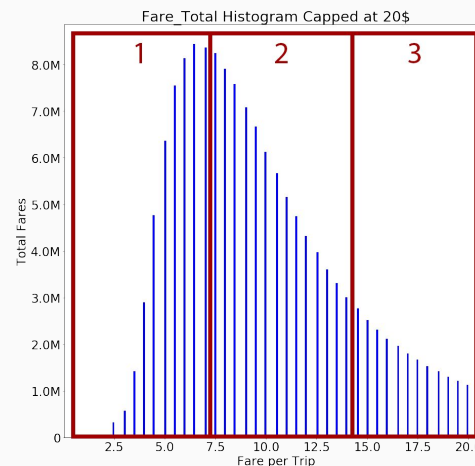total fare = amount of the fare + surcharge
maximum fare: 158'995.8125 $

# Optimization

Fare weight: proportionally duplicated records to add attraction
- Increased the dataset size of 2'575 %

Sampling fares in groups, example:

- 05.0$ → not cloned
- 11.5$ → cloned once
- 20.0$ → 3 cloned twice

MS≡  MASTER OF SCIENCE IN ENGINEERING

# Testing and evaluation

- Custom softwares for data transformations
- Our own Spark cluster
- Try, Fails and Meta Parameterizing

# Results

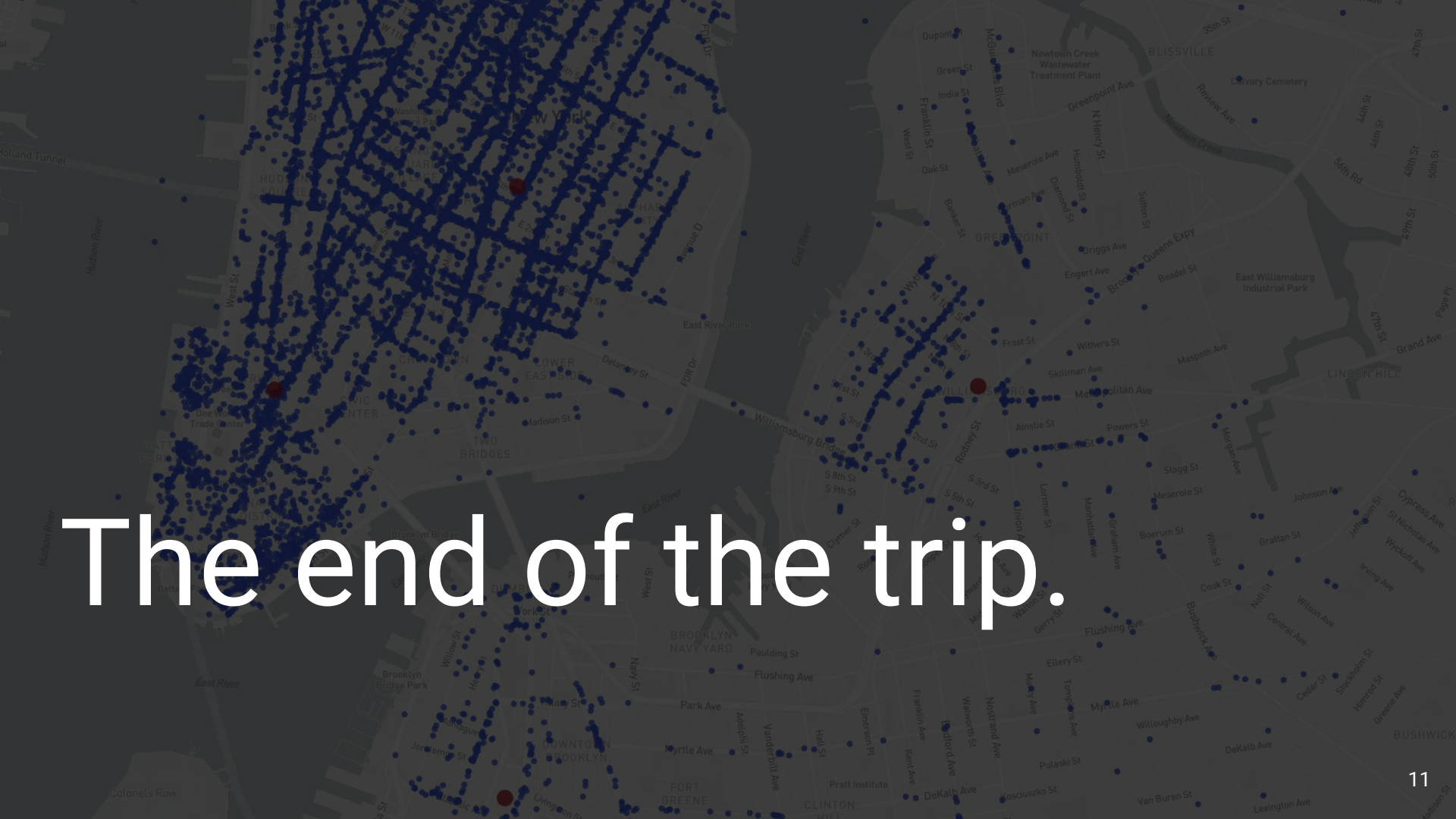Demo time



**http://bit.do/nyc-taxis**

# Conclusion

Best spots for customers are also the most profitable ones :)

Next steps:
- More processing time
- Zoom in data
- Prediction
- More data

# The end of the trip.

https://github.com/ZenLulz/hesso-bigdata-analytics