# SCALABILITY

Sankarshan and Paul
Airbnb Prediction Project

# MEANING OF SCALABILITY FOR OUR PROJECT

- Built a model pipeline in Spark that instantiates every possible regression model based on parameter inputs and combinations of input features

- Pipeline: Model Instantiation, Feature Importance, Evaluation Metrics

- Goal: find the best model features with brute force - or at least get an idea of where to look!

# PERFORMANCE

- We have only run our code in Spark on a personal computer - AWS coming soon!

- Works for as many as 9 features (1.5 hours), breaks with 10 or more.

# Executors

## Summary

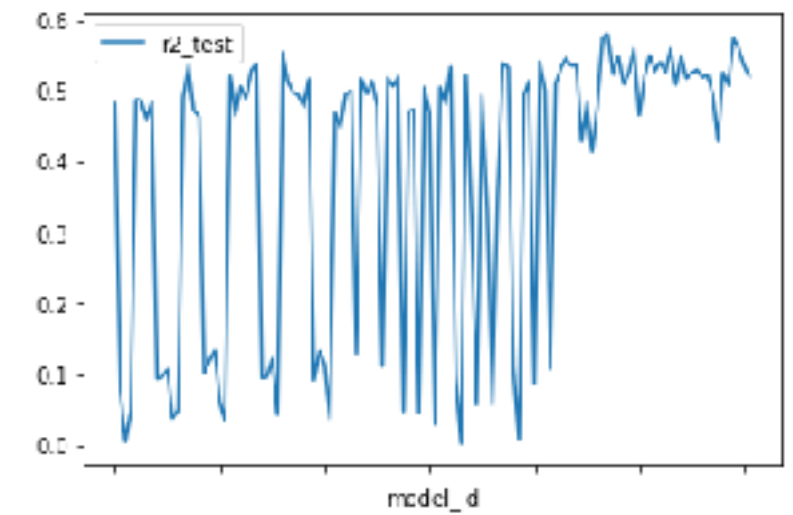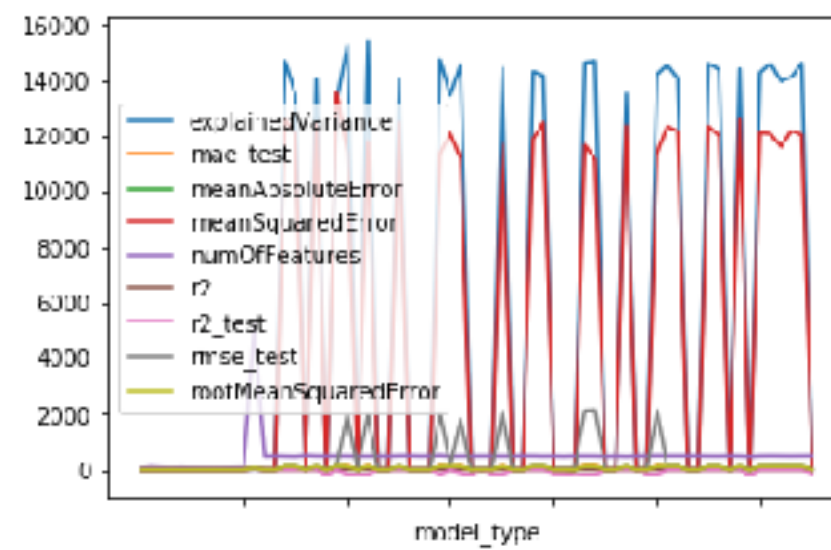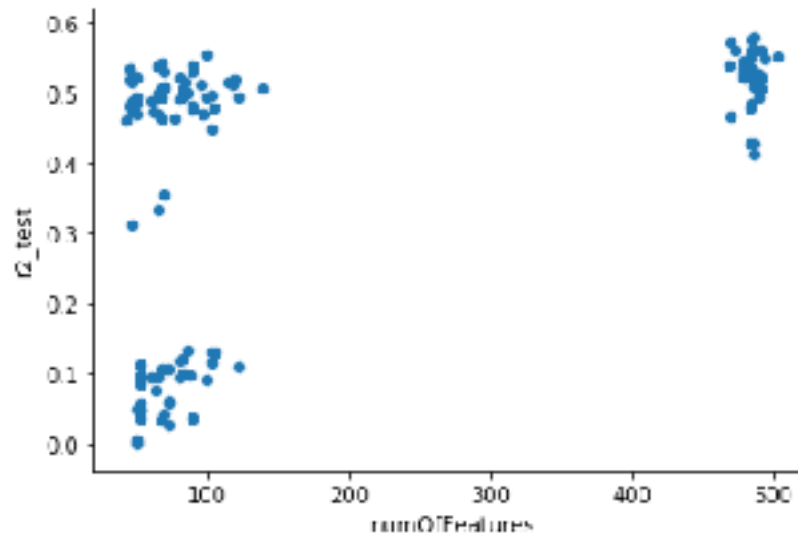| | RDD Blocks | Storage Memory | Disk Used | Cores | Active Tasks | Failed Tasks | Complete Tasks | Total Tasks | Task Time (GC Time) | Input | Shuffle Read | Shuffle Write | Blacklisted |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Active(1) | 134 | 134.5 MB / 384.1 MB | 0.0 B | 1 | 1 | 0 | 2517 | 2518 | 16 min (45 s) | 401.7 MB | 8.2 MB | 8.2 MB | 0 |
| Dead(0) | 0 | 0.0 B / 0.0 B | 0.0 B | 0 | 0 | 0 | 0 | 0 | 0 ms (0 ms) | 0.0 B | 0.0 B | 0.0 B | 0 |
| Total(1) | 134 | 134.5 MB / 384.1 MB | 0.0 B | 1 | 1 | 0 | 2517 | 2518 | 16 min (45 s) | 401.7 MB | 8.2 MB | 8.2 MB | 0 |

## Executors

Show 20 entries

Search:

| Executor ID | Address | Status | RDD Blocks | Storage Memory | Disk Used | Cores | Active Tasks | Failed Tasks | Complete Tasks | Total Tasks | Task Time (GC Time) | Input | Shuffle Read | Shuffle Write | Thread Dump |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| driver | 10.0.1.2:55816 | Active | 134 | 134.5 MB / 384.1 MB | 0.0 B | 1 | 1 | 0 | 2517 | 2518 | 16 min (45 s) | 401.7 MB | 8.2 MB | 8.2 MB | Thread Dump |

Showing 1 to 1 of 1 entries
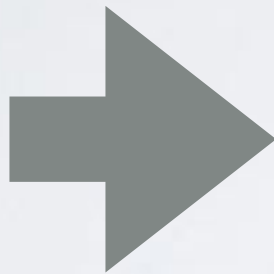
Previous  1  Next

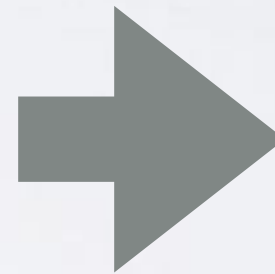# RESULTS

# DATA PIPELINE

**Data Acquisition**

- insideairbnb.com

- Scraper?

**Data Transformations**

- Compute Price Statistics

- Text: LDA, NLTK, Clustering

- Unstructured Features

- Geo Features: Distance from Ocean, Park/ Recreation Site, Active Businesses, Local Events

**Modeling**

- Model Exploration Phases

- Scaled Modeling

- Final Evaluation & Model Choice

# NEXT STEPS

- Run in AWS

- Expand capabilities of scalable pipeline

- Conclude work on scraper