# Forecasting Grocery Sales with Regression Models and Random Forests

Victor de Fontnouvelle

Justin Weltz

Thomas Thornton

# Introduction



- Finite supply
- Perishable

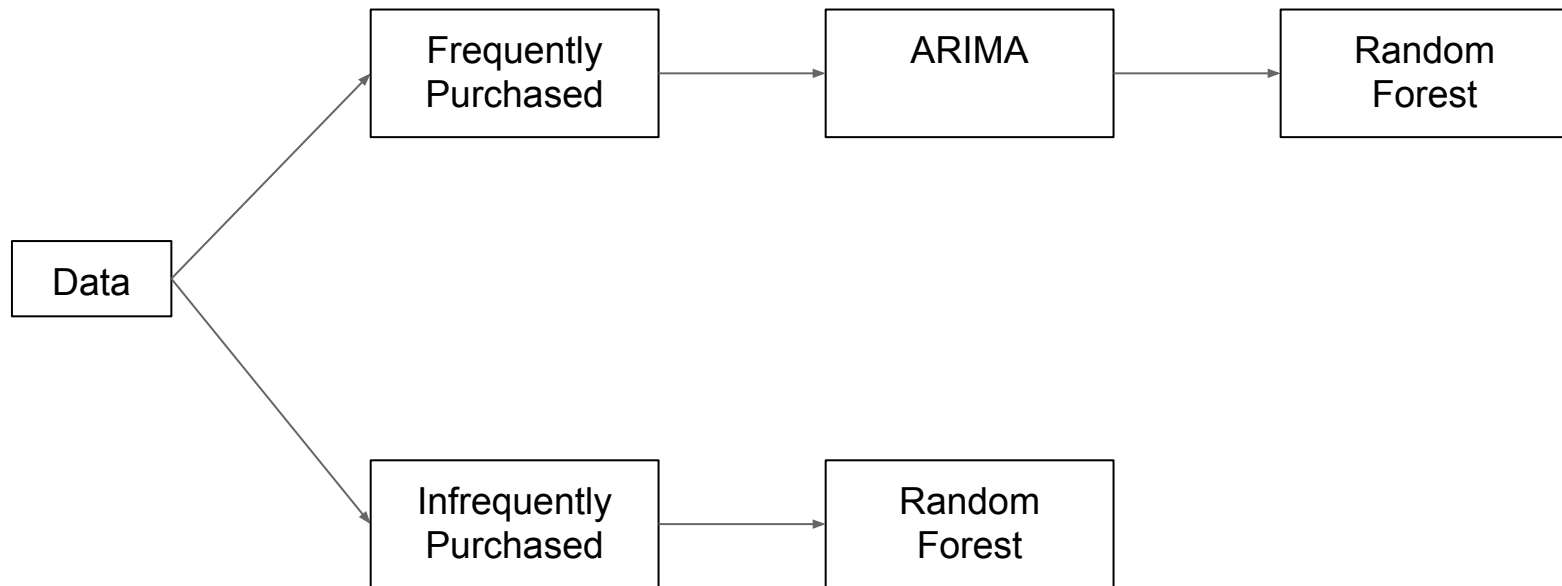How to predict demand?

# Corporacion Favorita

- 54 stores
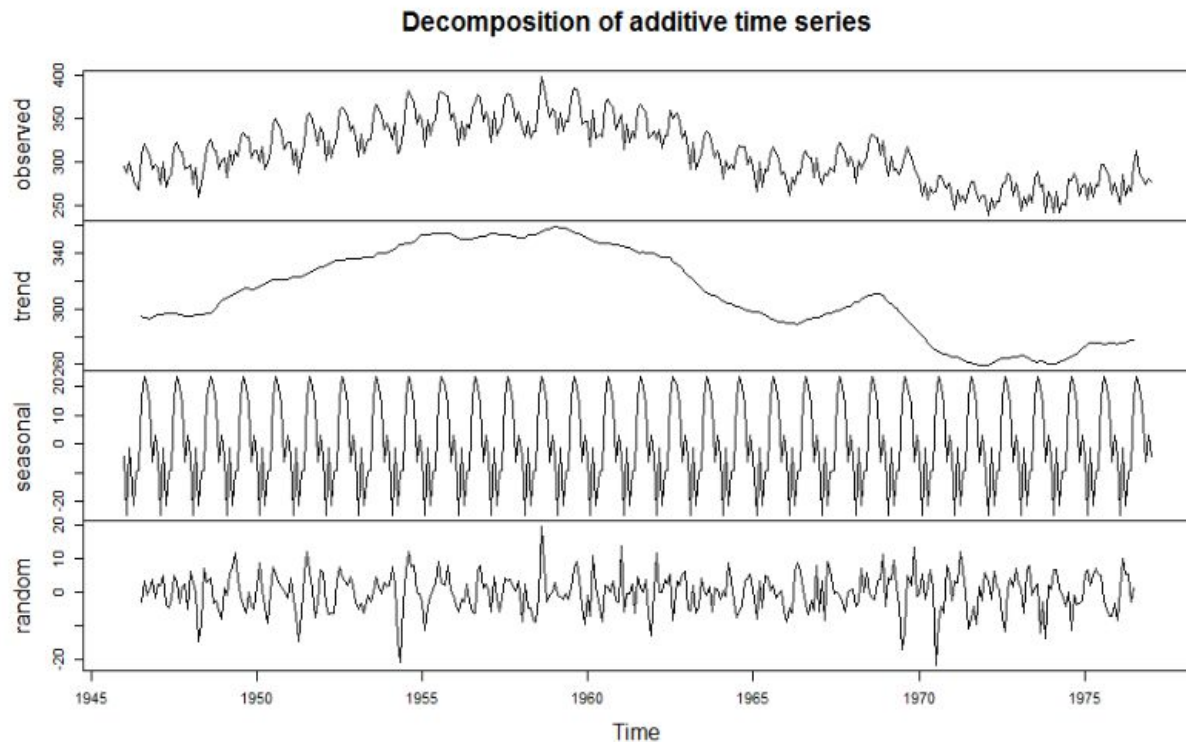- 200,000 items
- 125M transactions

# Datasets & Wrangling

- Items - type
- Stores - location, cluster
- Oil - date, price
- Train/Test Data - item number, store number, date ➔ **sales**
  - Joined tables
  - Converted categorical to binary
  - Added numerical variable: days since payday
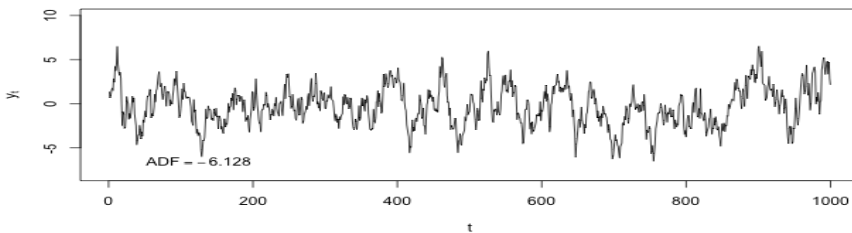
# Dual Model Pipeline

# Time Series Forecasting



Decomposition of additive time series

# The ARIMA Model Equations

$$Y_t = \alpha + \rho Y_{t-1} + e_t$$



$$Y_t = \alpha + \rho e_{t-1} + e_t$$

# Popularity

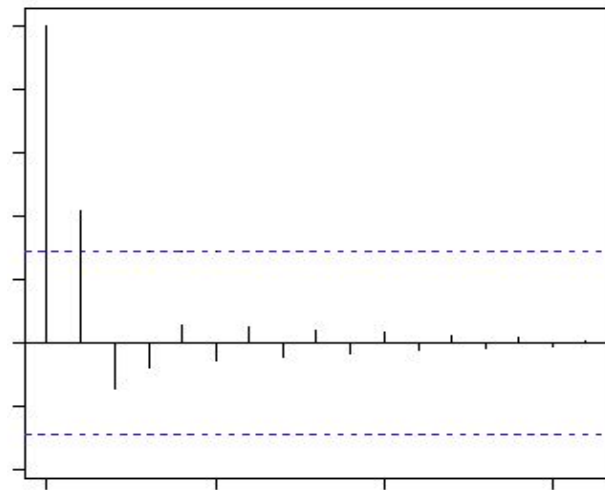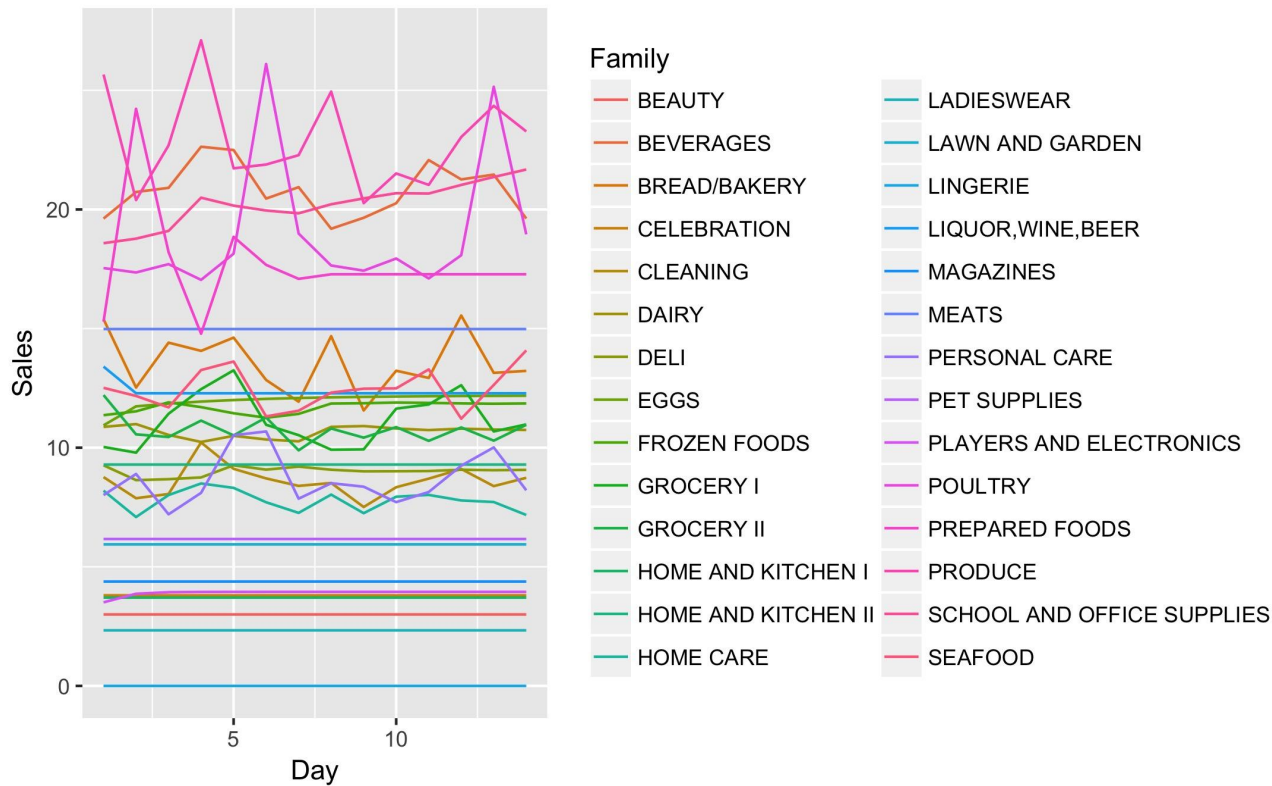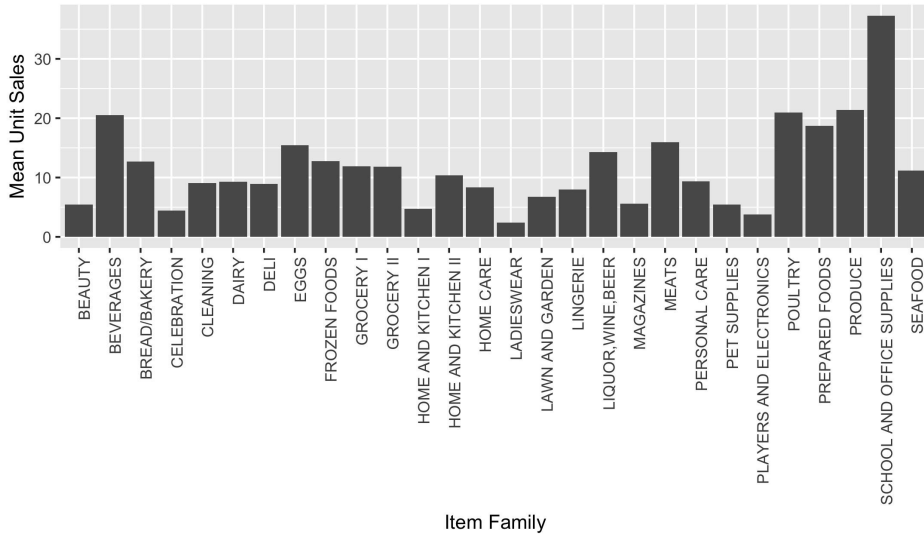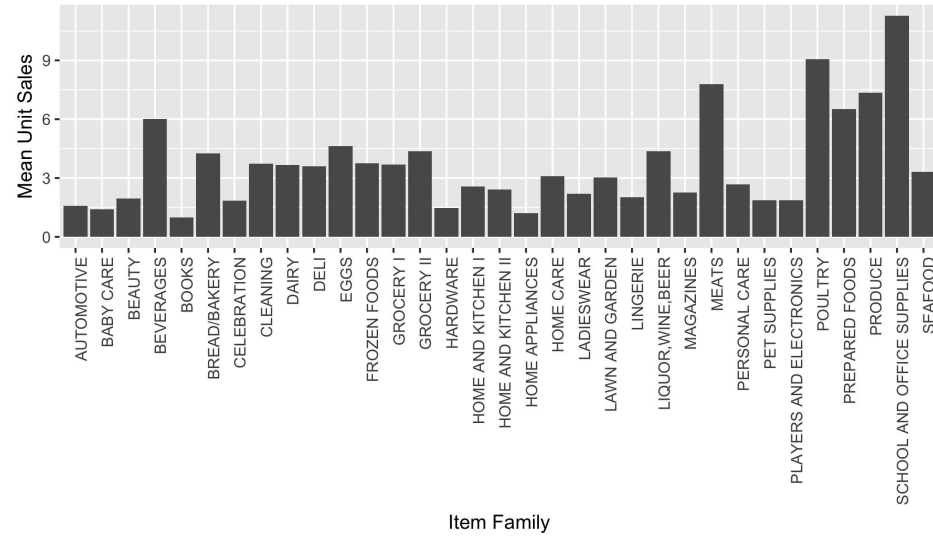# Random Forest Creation

Variables Included:

- Item Family
- Store City
- Store Cluster
- Store Type
- Day of Week
- Days Since Payday
- Popularity

# Item Family

# Store City

# Store Cluster

# Store Type



Effect of Store Type for Highly-Transacted Items

Effect of Store Type for Non-Highly-Transacted Items

# Day of Week
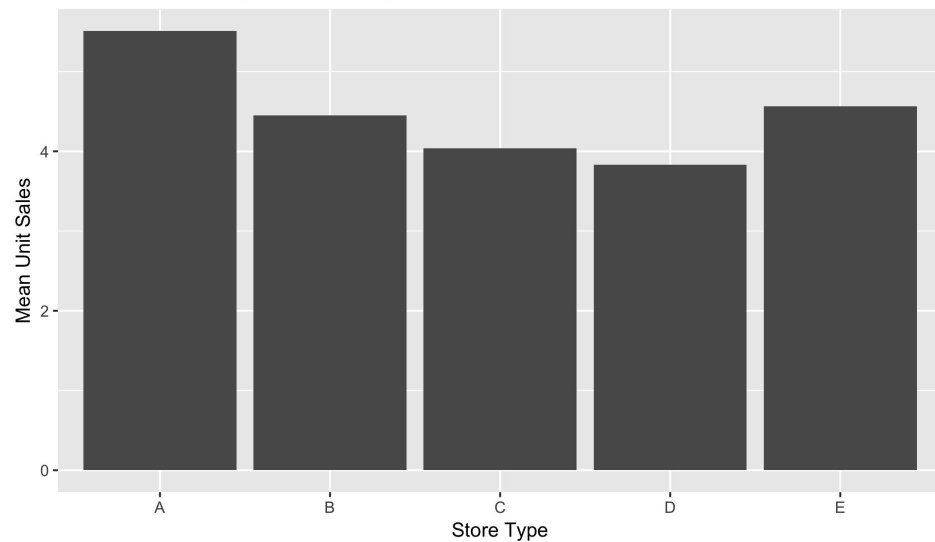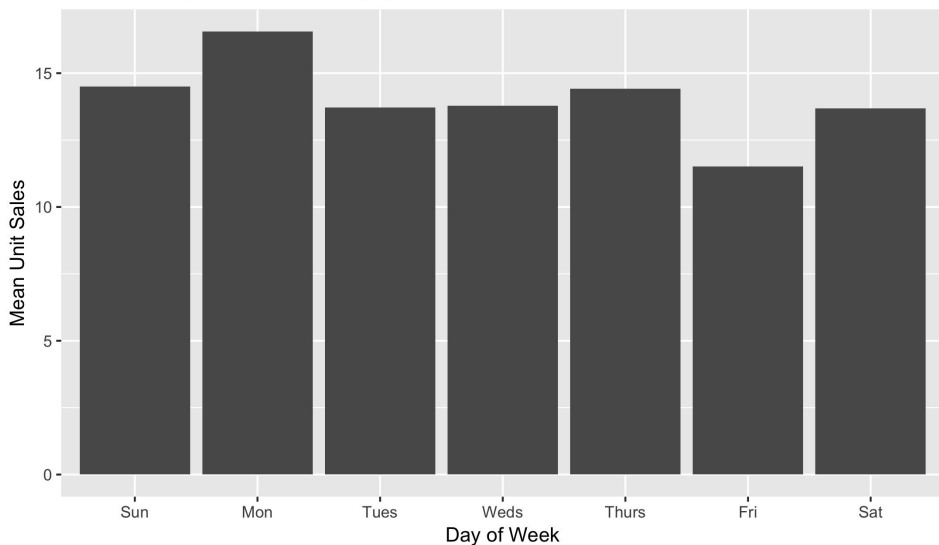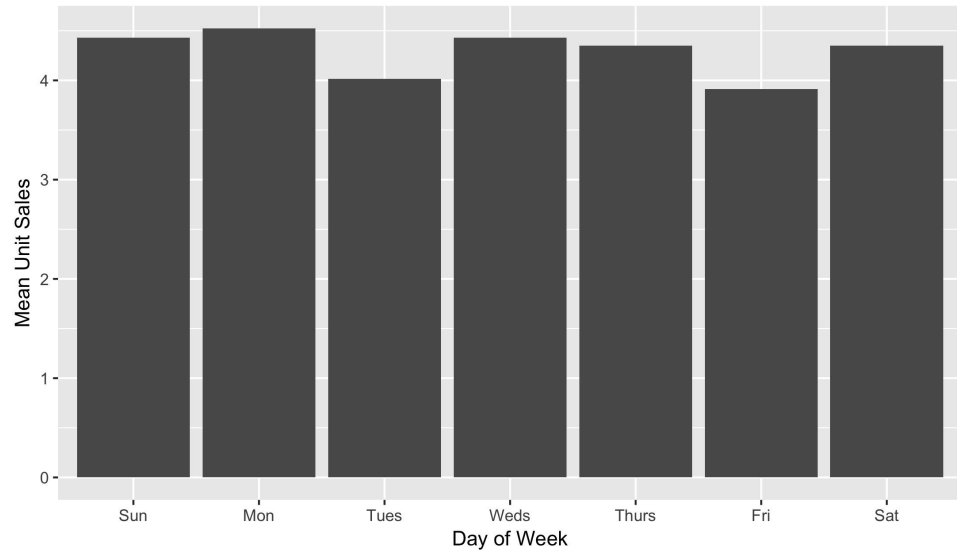


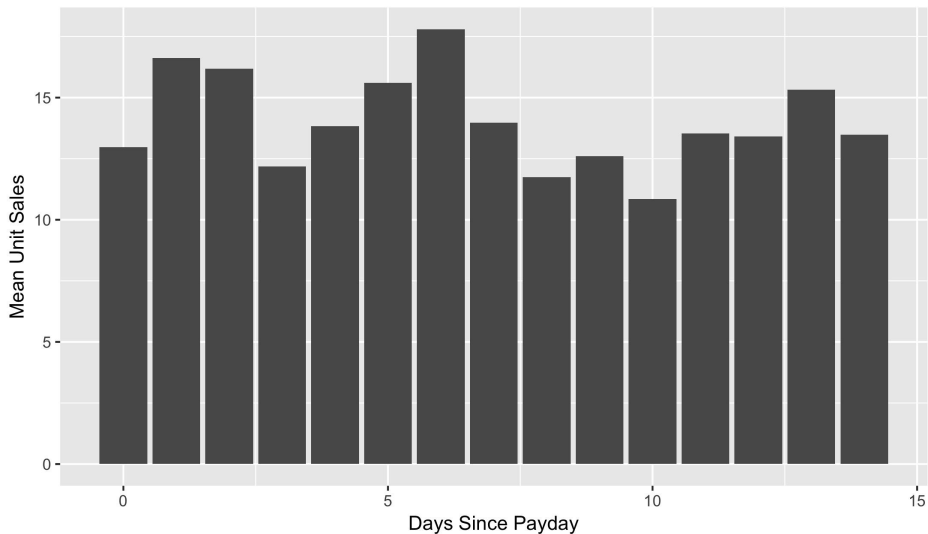Effect of Day of Week for Highly-Transacted Items

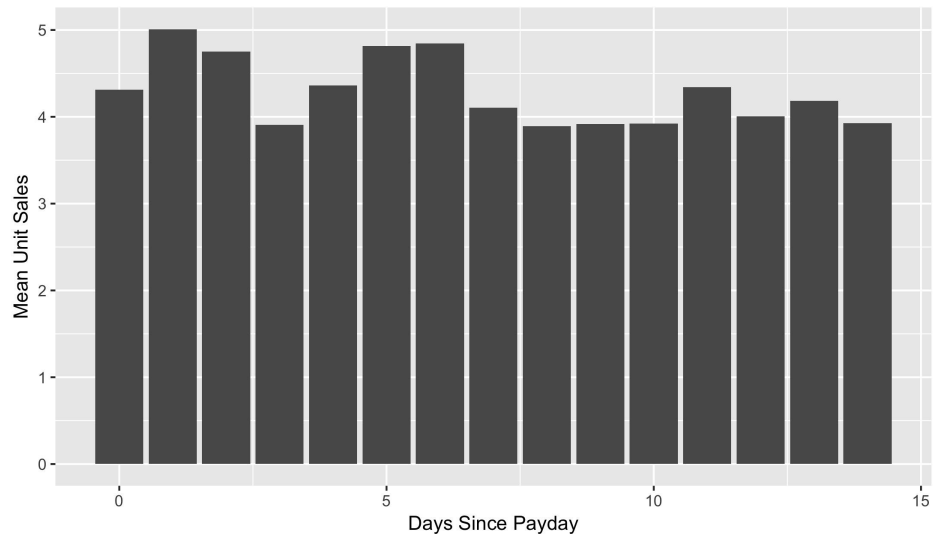Effect of Day of Week for Non-Highly-Transacted Items
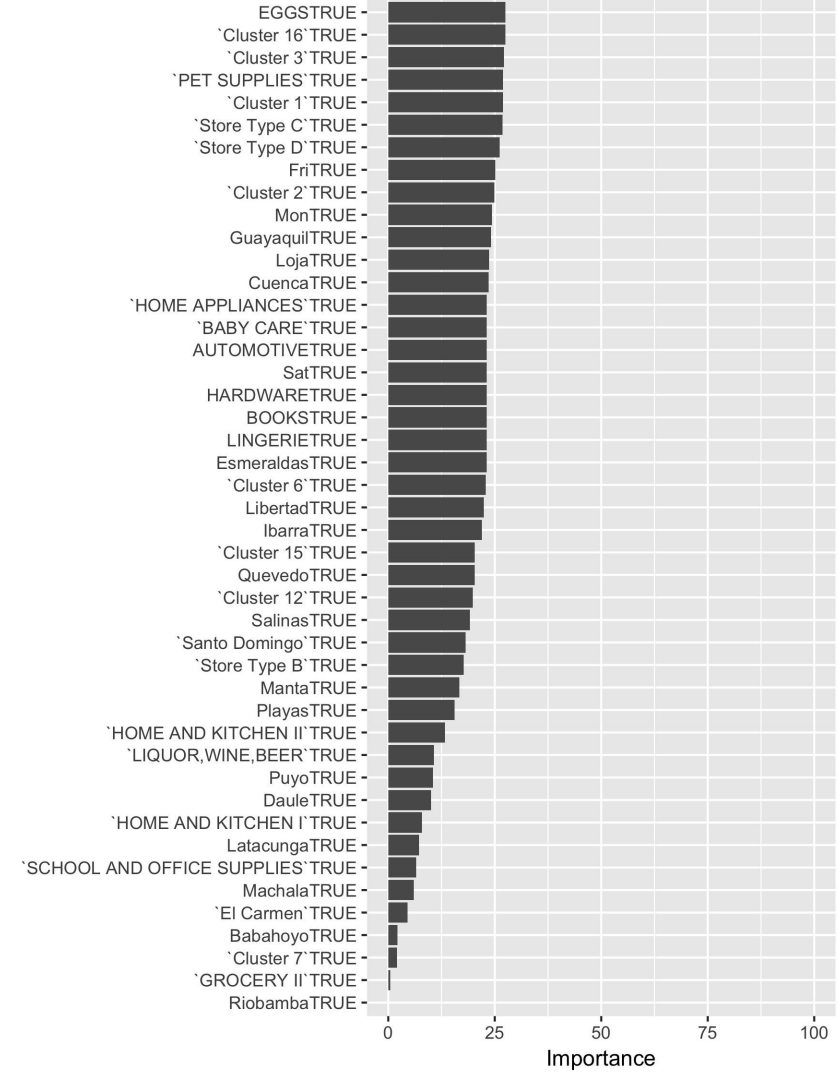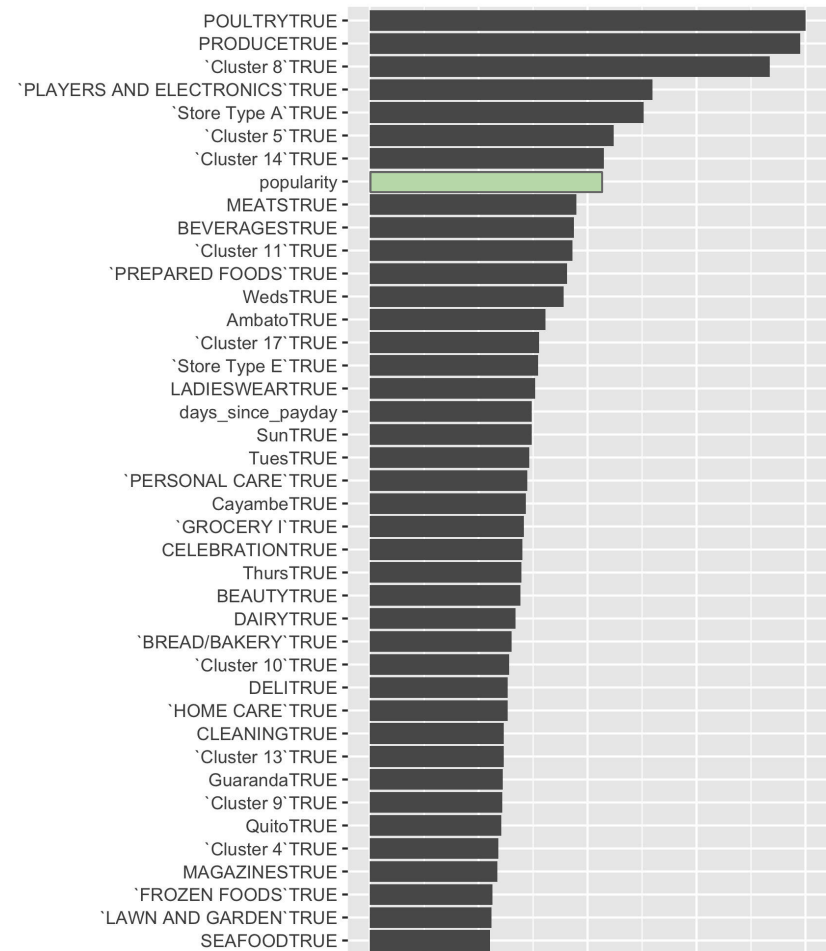
# Days Since Payday
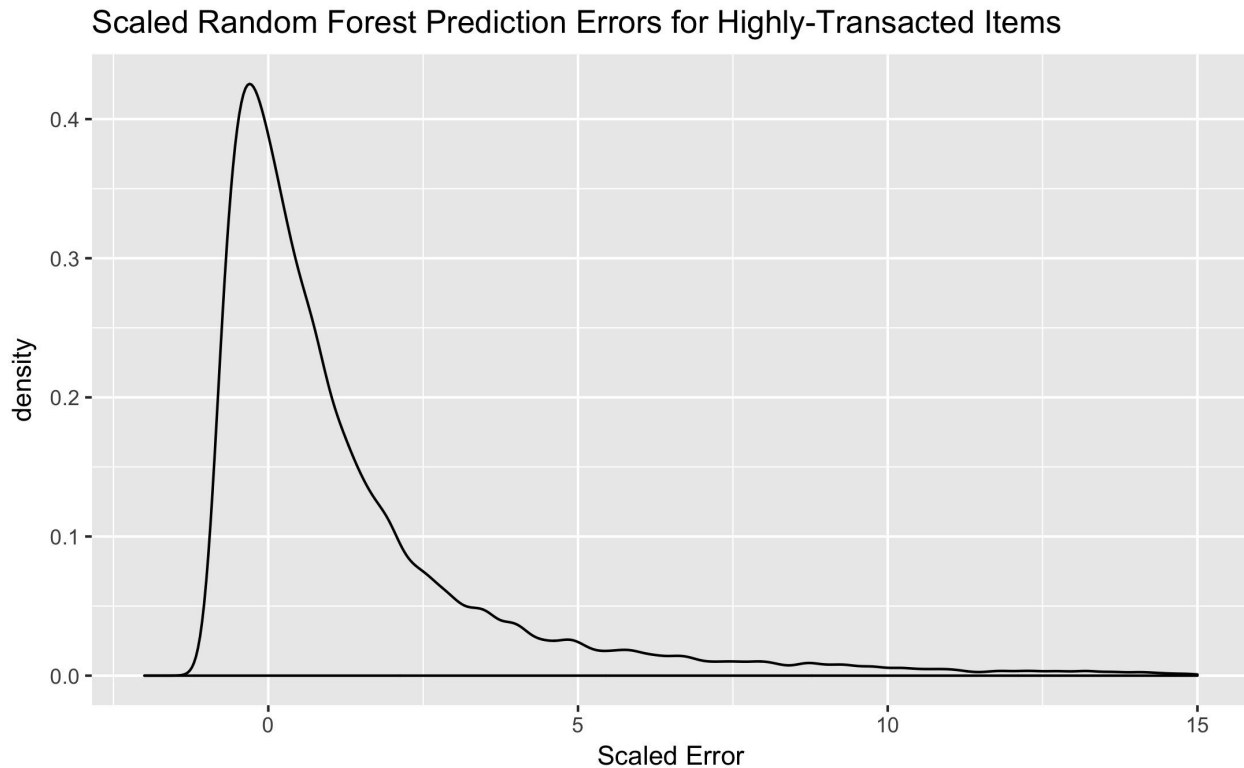


Effect of Days Since Payday for Highly-Transacted Items

Effect of Days Since Payday for Non-Highly-Transacted Items

# Results: Variable Importance

# Results - Prediction Errors



Scaled Random Forest Prediction Errors for Highly-Transacted Items

# Results: Kaggle

| | | | | |
|---|---|---|---|---|
| 933 | ▼ 153 | **Matheus Facure** | | 1.294 |
| 934 | ▼ 153 | **Yosuke Abe** | | 1.295 |
| 935 | ▼ 153 | **rjuer** | | 1.295 |
| 936 | ▼ 153 | **Jeffrie** | | 1.295 |
| 937 | ▼ 153 | **Anjukan Kathirgamanathan** | | 1.299 |
| 938 | ▼ 153 | **mhaulrich** | | 1.303 |
| 939 | new | **Victor de Fefontnouvelle** | | 1.307 |
| 940 | ▼ 106 | **Magic Logic** | | 1.309 |
| 941 | ▼ 154 | **tomgrek** | | 1.310 |

1036 Teams Total

# Potential Improvements

- **0.25% of training set used**
- Add data for no sales
- Incorporate other datasets
  - Holiday
  - Weather
  - Economy