

# Predicting NYC Taxi ETA

Liyan Nie  
Estelle Danilo

# Data

- 2016 NYC yellow taxi data (BigQuery)
- 2016 NYC weather data (BigQuery)

# Data Cleaning

# Missing values

Taxi data:

- Remove rows

|                          | missing_ratio |
|--------------------------|---------------|
| <b>dropoff_latitude</b>  | 47.11865      |
| <b>dropoff_longitude</b> | 47.11865      |
| <b>pickup_latitude</b>   | 47.11865      |
| <b>pickup_longitude</b>  | 47.11865      |

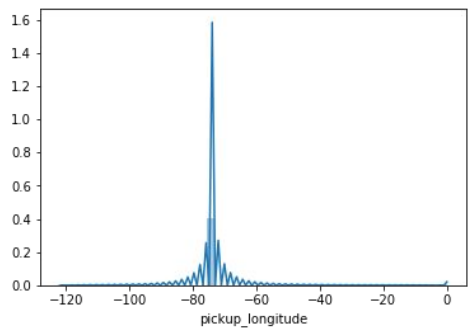
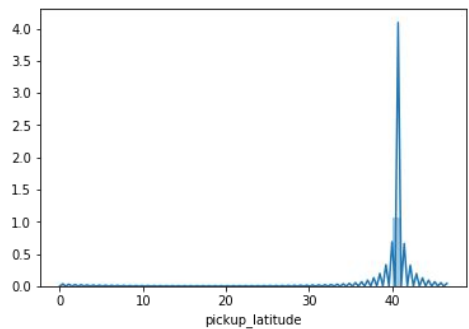
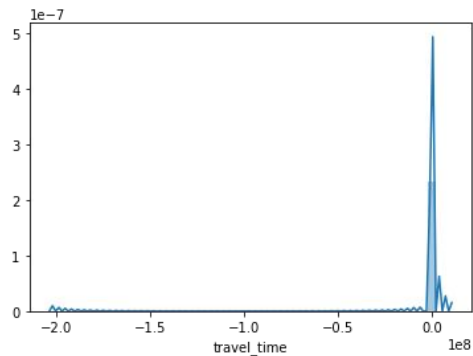
Weather data

- Remove missing rows
- Remove “sndp”

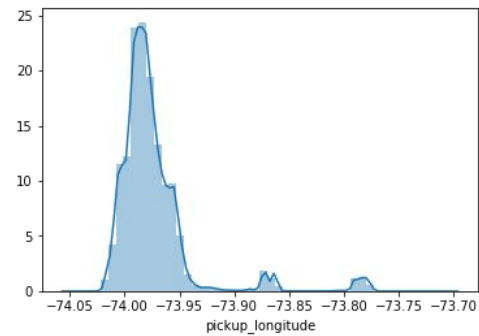
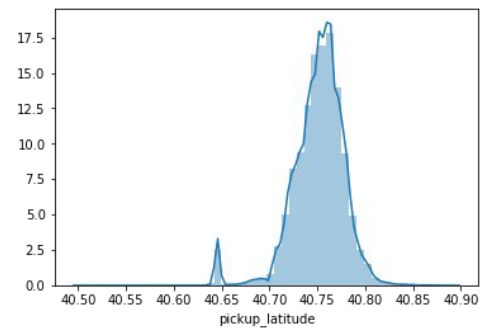
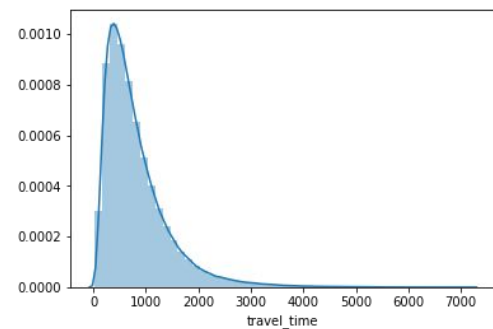
|              | missing_ratio |
|--------------|---------------|
| <b>visib</b> | 2.30176       |
| <b>wdsp</b>  | 2.86897       |
| <b>gust</b>  | 28.1718       |
| <b>sndp</b>  | 94.4242       |

# Outliers

- Remove `passenger_count == 0`
- Remove `rate_code == 99` (scale is 1-6 only)
- Only keep `lat` in `[40.6, 40.9]`
- Only keep `long` in `[-74.05, -73.7]`
- Only keep `travel_time` between 30s and 7200s ( 2 hours)



Remove outliers



# EDA & Identifying Features

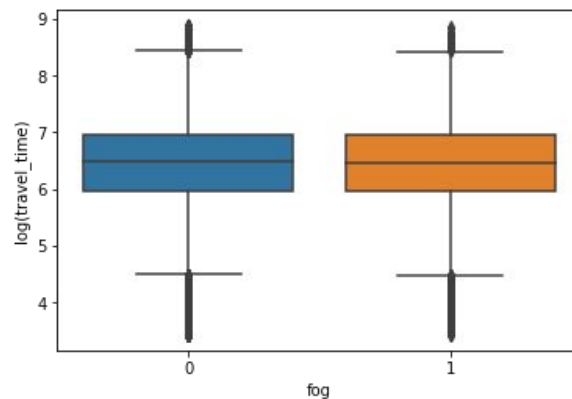
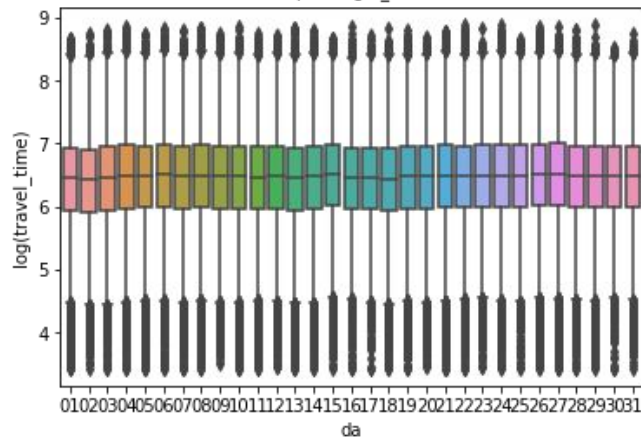
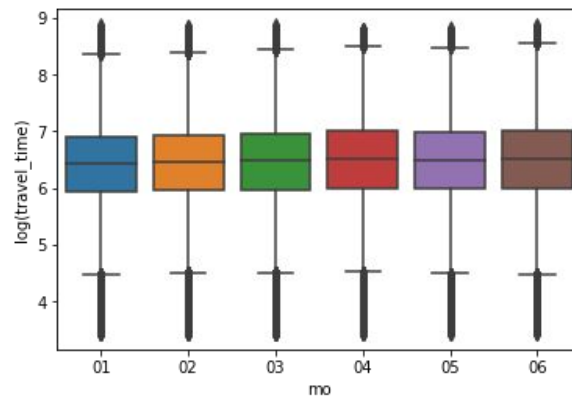
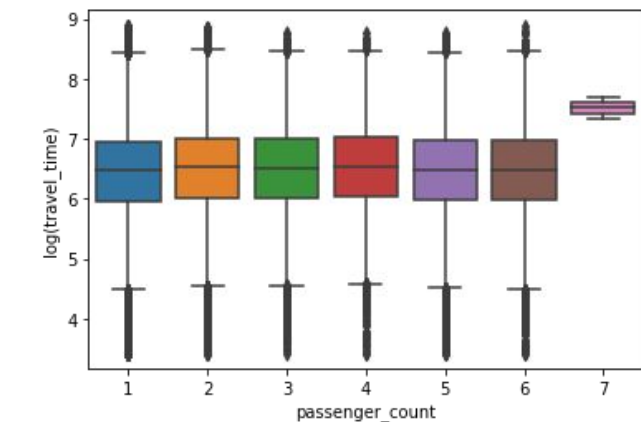
# Information leakage

Remove:

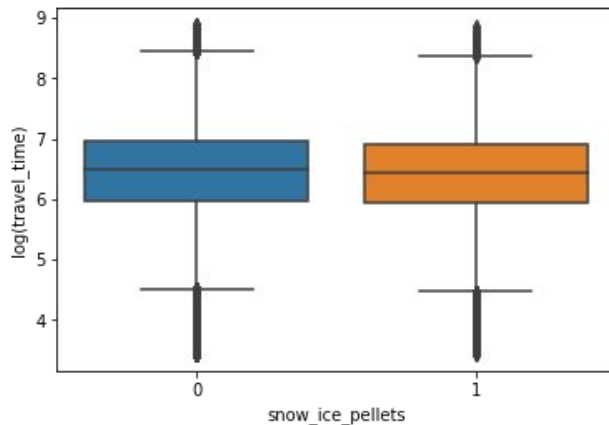
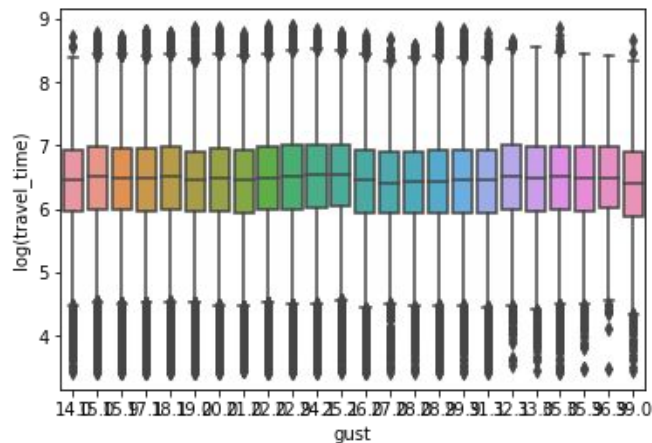
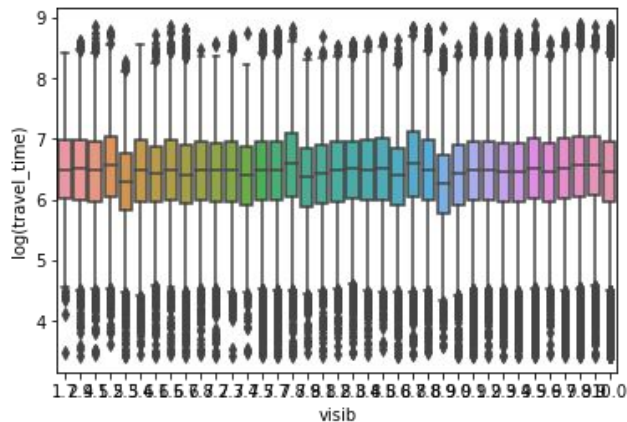
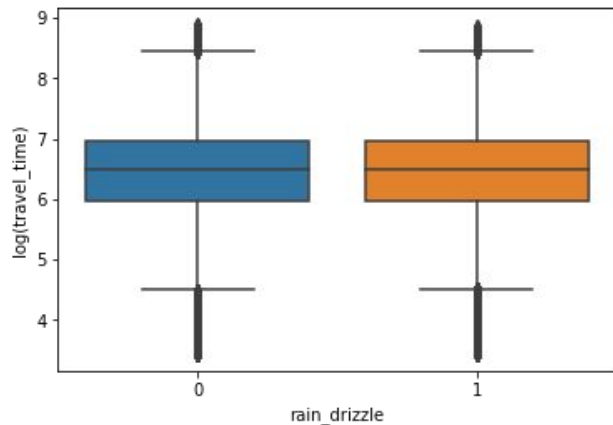
- Trip\_distance
- Fare\_amount
- Total\_amount
- Payment\_type
- Extra
- Mta\_tax
- Tip\_amount
- Tolls\_amount
- imp\_surcharge



# Visualize response vs variables

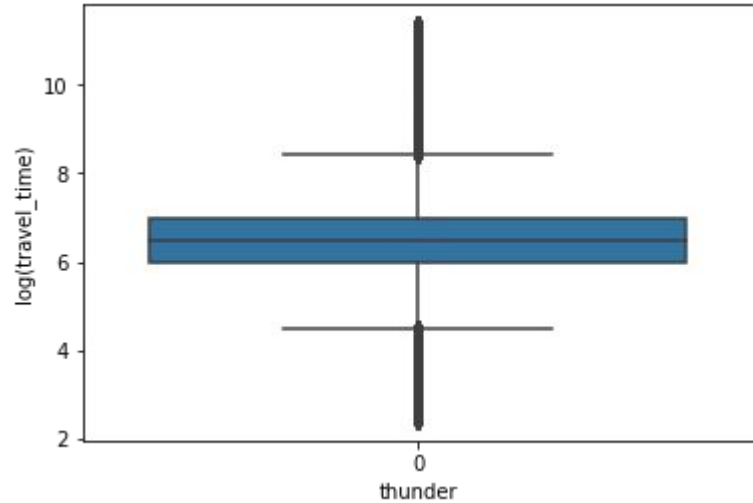
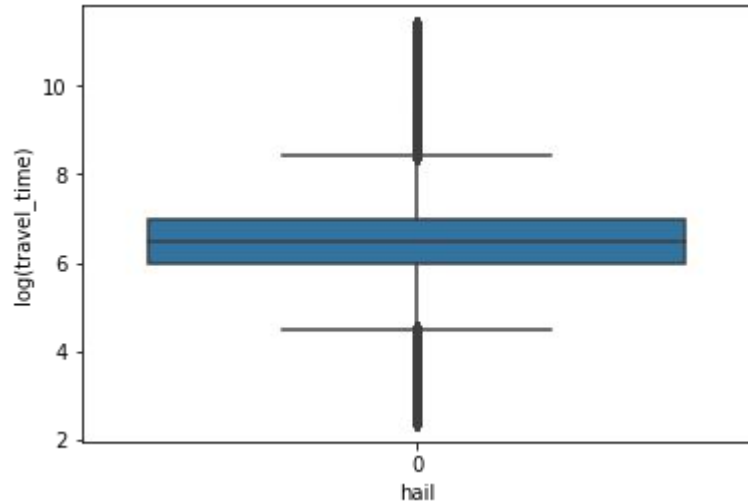


# Visualize response vs variables

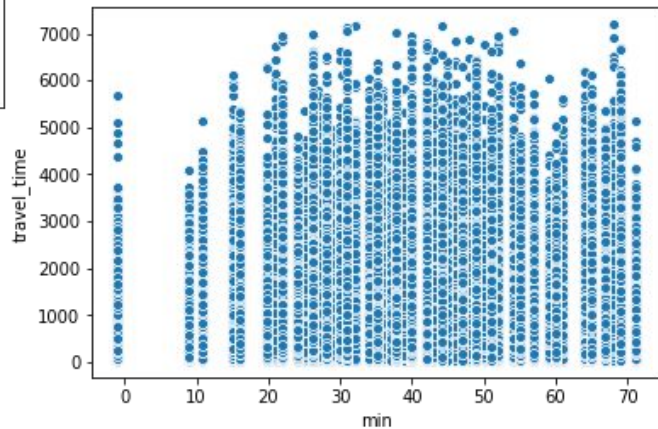
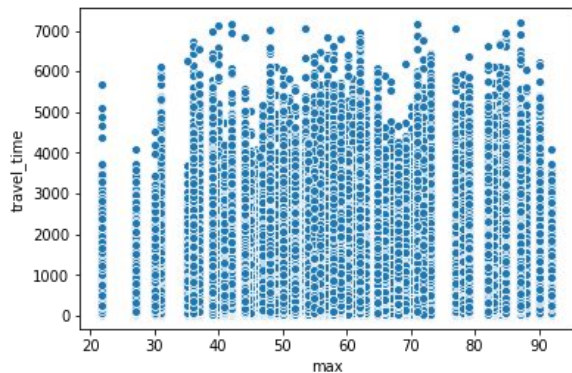
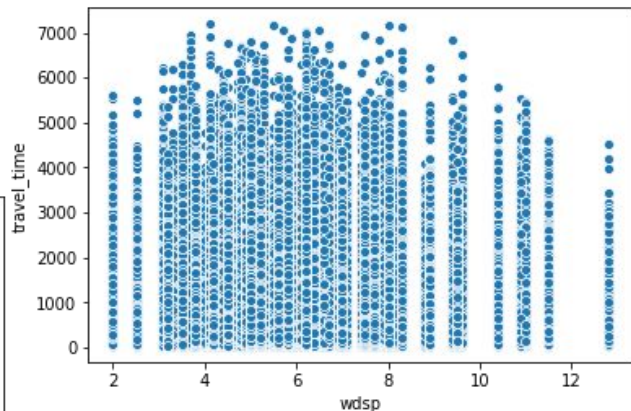
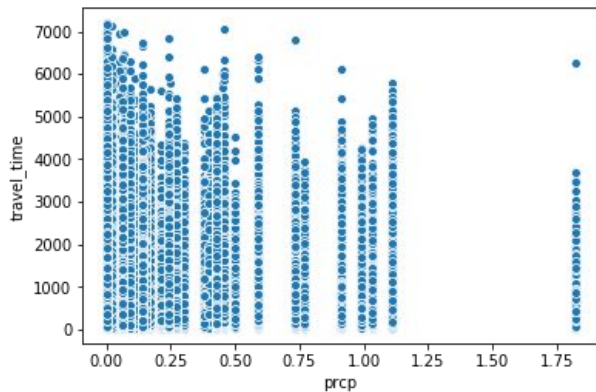
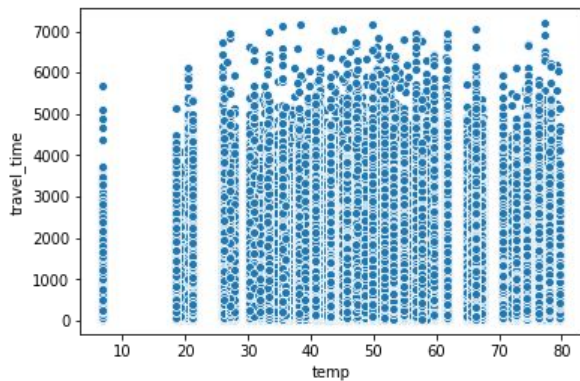


# Visualize response vs variables

Remove



# Visualize response vs variables



# EDA

Preliminary conclusion:

- None of these features alone seem to have strong prediction power. Need further feature engineering
- But first, run a few models with minimal preprocessing/feature engineering to get a baseline idea

# Baseline Modeling

# Identifying Features

Using all possible **non-leaking** features, given no strong signal in EDA

## Location

- pickup\_latitude
- pickup\_longitude
- dropoff\_longitude
- dropoff\_latitude

## Time

- day\_of\_year
- month\_of\_year

## Type of Taxi Trip

- vendor\_id
- passenger\_count
- rate\_code
- store\_and\_fwd\_flag
- payment\_type

# Features: Continuous vs. Categorical

## Categorical (12 levels or less)

- vendor\_id
- store\_and\_fwd\_flag
- payment\_type
- rate\_code (**continuous by default**)
- passenger\_count (continuous by default)
- month\_of\_year (continuous by default)



**OneHotEncode** (dropping first level)

## Continuous

- day\_of\_year
- pickup\_longitude
- pickup\_latitude
- dropoff\_longitude
- dropoff\_latitude



# Data Splitting

**Goal:** Find split by **pickup\_datetime** from January to June 2016 (where lat-long exist) that creates a 80-20% train/test split

**Method:** Training = pickup\_datetime > 2016-01-01 & < 2016-05-25  
Test set = pickup\_datetime >= 2016-05-25 & < 2016-07-01

**Result:** training and test splitted into **X\_train, X\_test** (features)  
and **y\_train, y\_test** (travel\_time)

# Baseline Models:

## 1) Decision Tree

To prevent overfitting →

## 2) Bagging

average out trees predictions

## 3) Random Forest

train on different samples of data

To split each node in tree:

**All features** considered

**only best split features**  
from a random subset of features

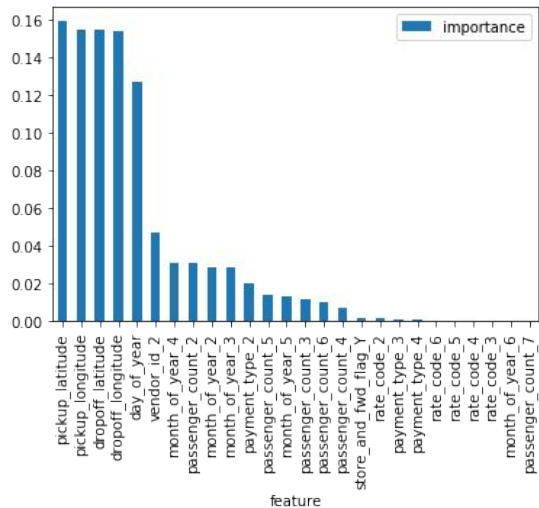
# Baseline Model Performance

- Low 5-fold cross validated test performance
- **Bagging** least underperforming

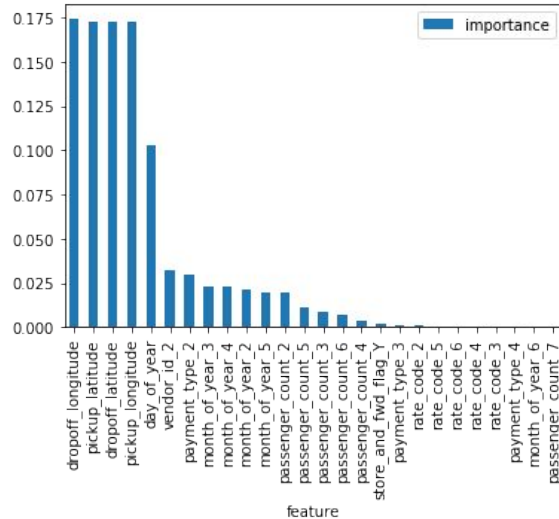
|                                | Decision Tree | Bagging Classifier  | Random Forest |
|--------------------------------|---------------|--|---------------|
| <b>R2</b>                      | -16           | -15  | -52           |
| <b>RMSE</b>                    | 4141          | 3301   | 3987          |
| <b>RMSLE</b>                   | 1.07          | 1.28   | 1.18          |
| <b>Training Subsample Size</b> | 200,000       | 50,000   | 50,000        |
| <b>Testing Subsample Size</b>  | 50,000        | 12,500   | 12,500        |

# Baseline Feature Importance

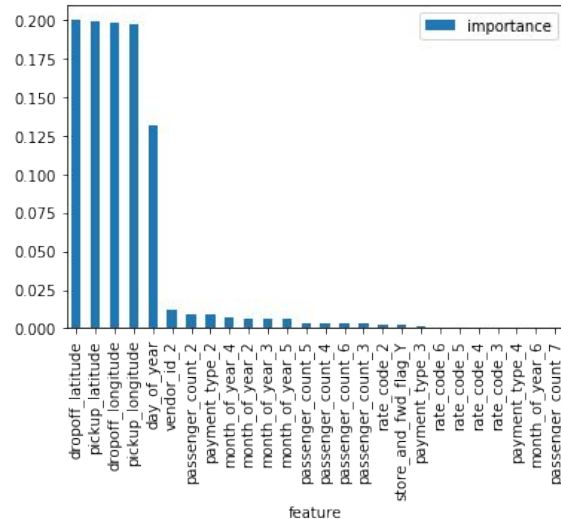
## Decision Tree



## Bagging Classifier



## Random Forest



- Location features most important
- Similar feature and magnitude importance across models

# GridSearched Baseline Models (5-fold)

# Hyperparameter Tuning

## 1) Decision Tree


max\_depth: [5,10],  
min\_samples\_leaf: [10, 100],  
min\_samples\_split : [50,200,500]

## 2) Bagging

n\_estimators: [10, 50],  
max\_samples' : [0.1, 0.5]

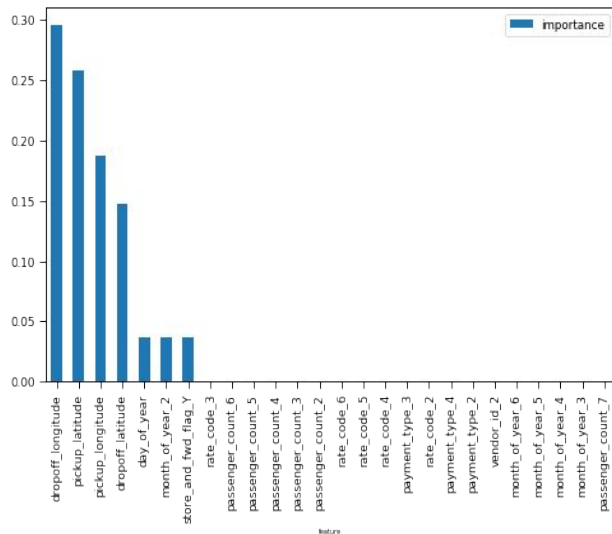
## 3) Random Forest

n\_estimators : [100, 200],  
max\_depth : [5,10]

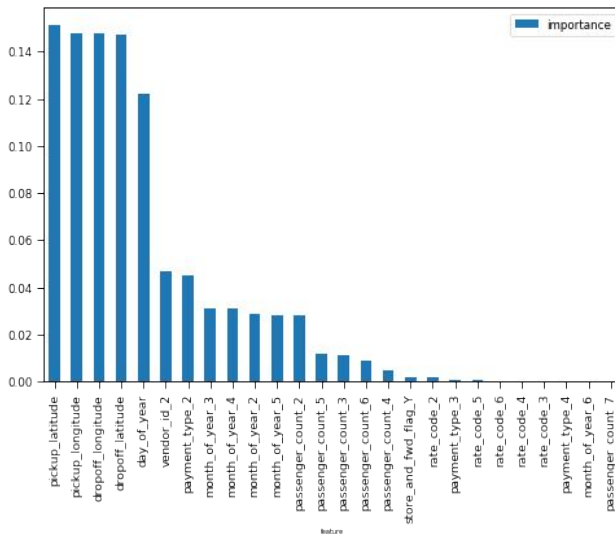
|              | Decision Tree | Bagging Classifier | Random Forest  |
|--------------|---------------|--------------------|---|
| <b>R2</b>    | -0.38         | -1.58              | -0.12   |
| <b>RMSE</b>  | 3053          | 3484               | 2390  |
| <b>RMSLE</b> | 1             | 1.32               | 0.87  |

# Feature Importance

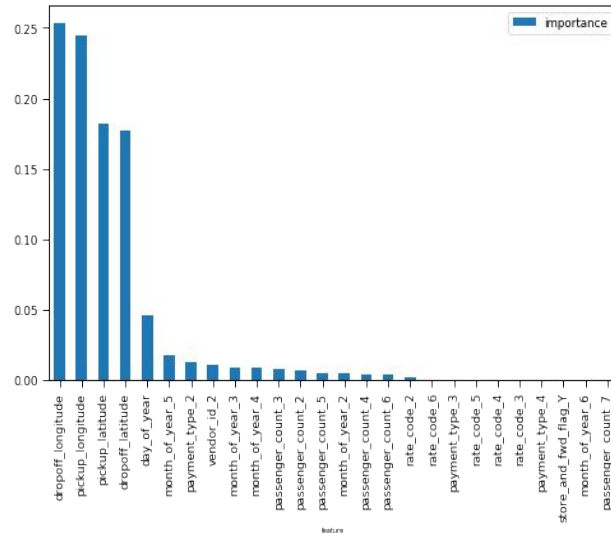
## 1) Decision Tree



## 2) Bagging



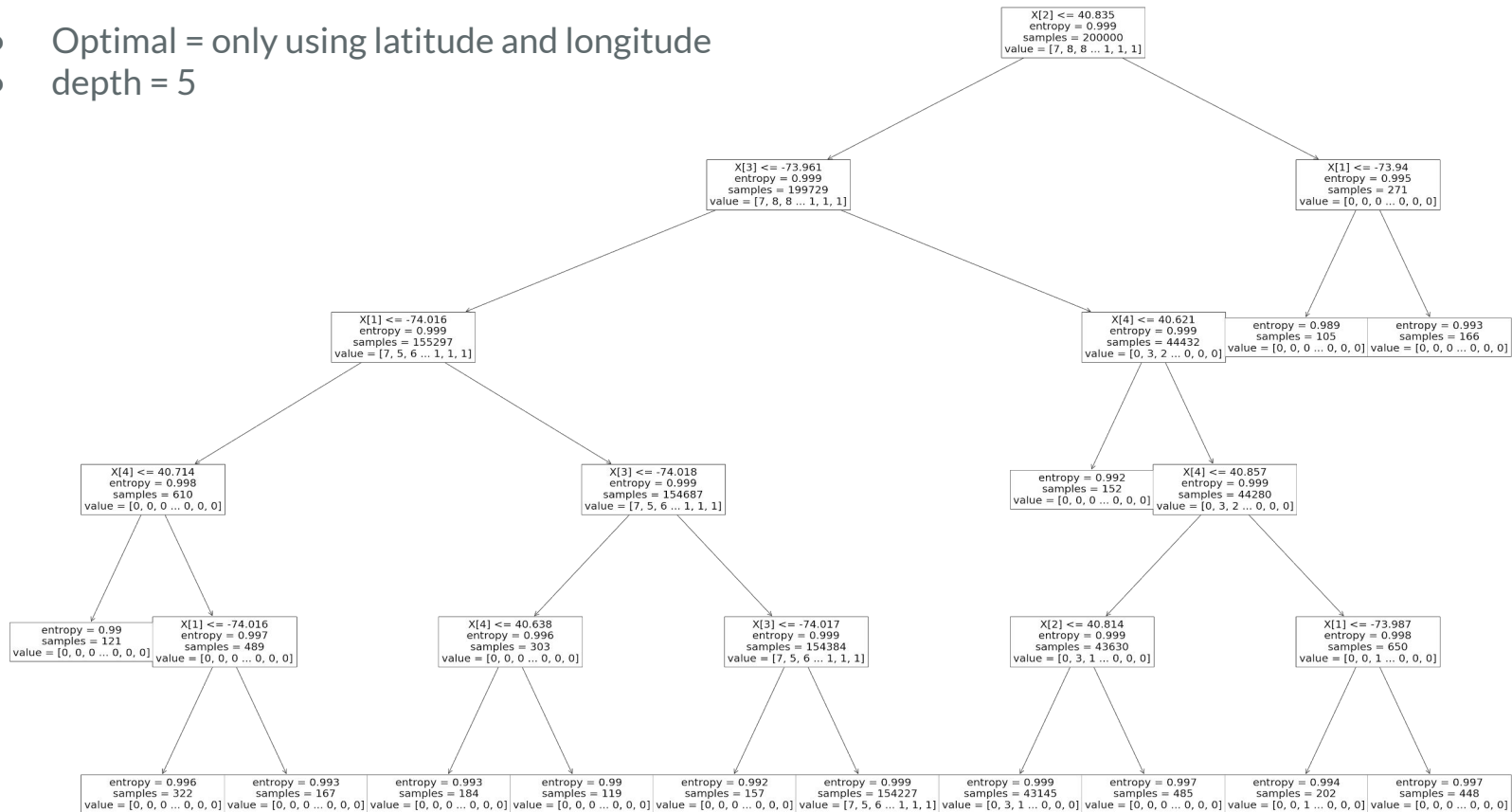
## 3) Random Forest



- Location features remain the most important
- More variance in magnitude in feature importances across models

# GridSearched Decision Tree

- Optimal = only using latitude and longitude
- depth = 5

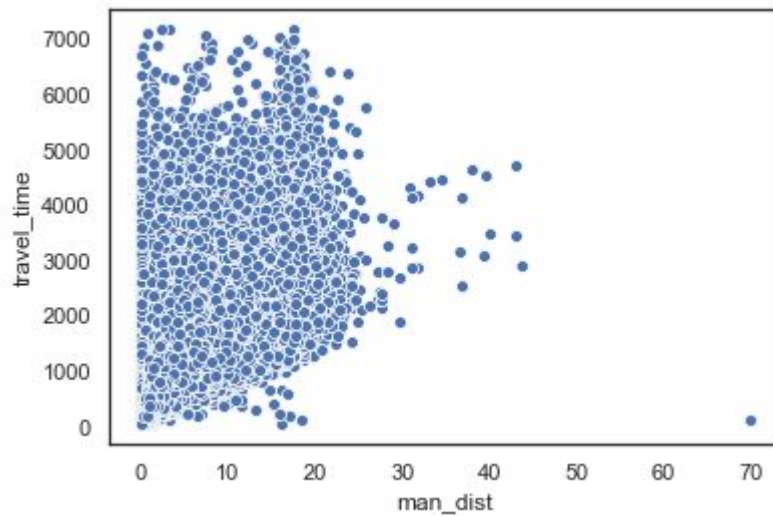
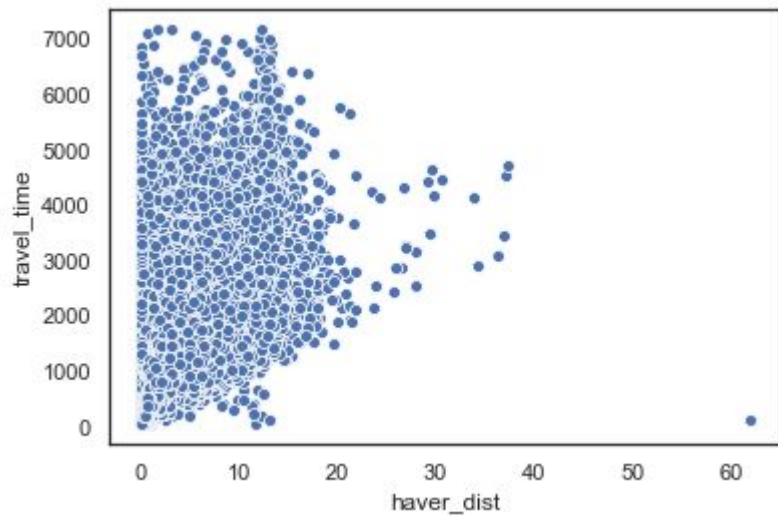




# Feature Engineering

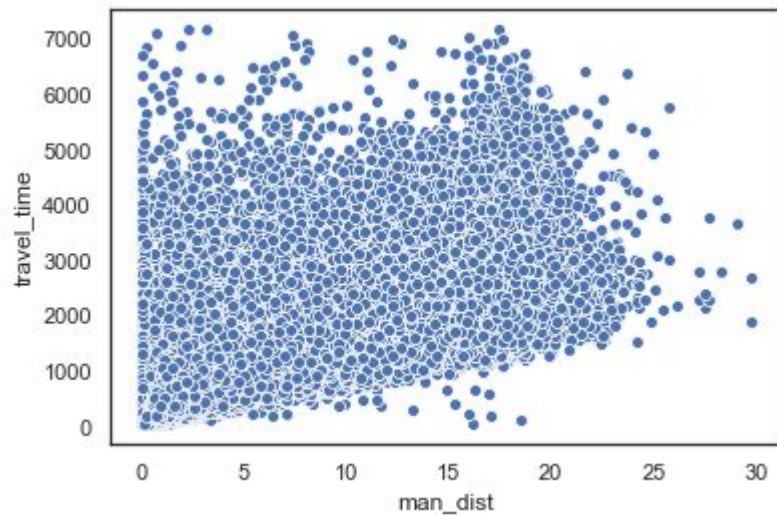
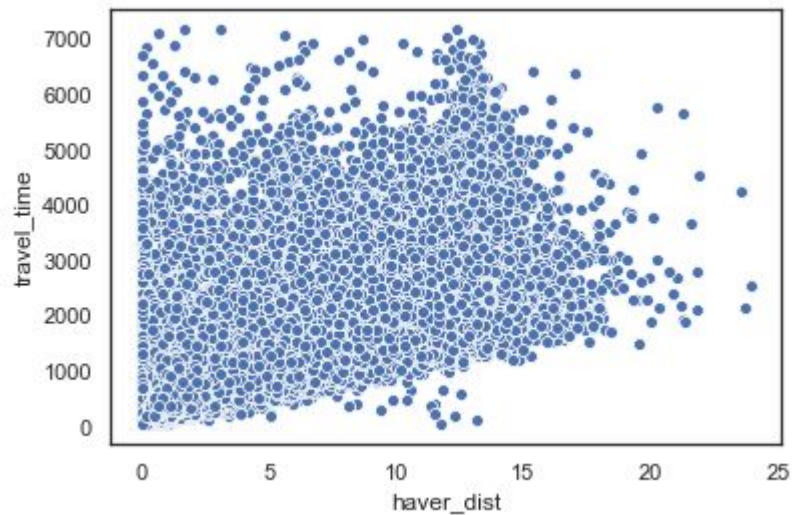
# Distance

Potentially strong predictors!



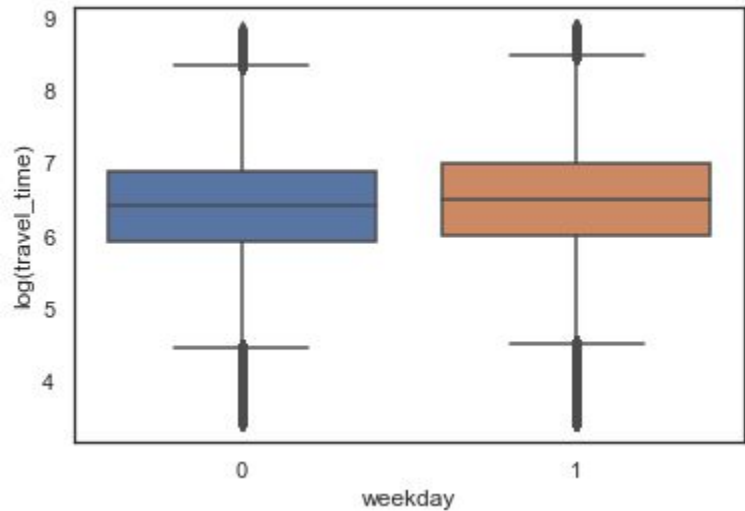
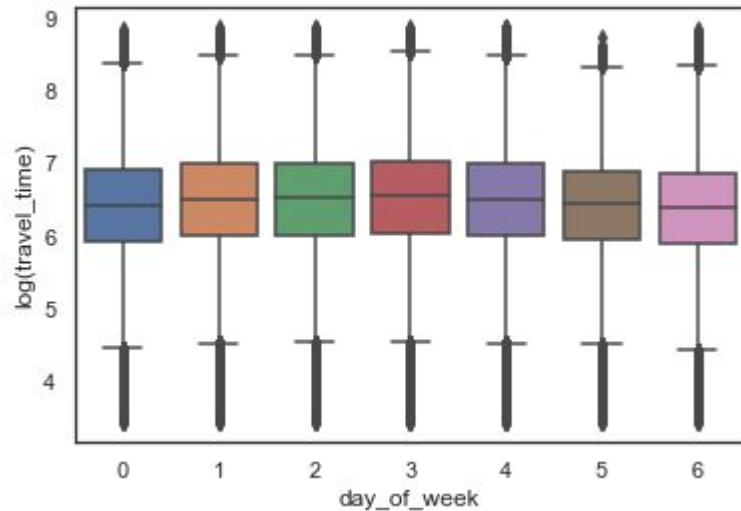
# Distance

Further remove outliers: Haversine > 25, and Manhattan > 30



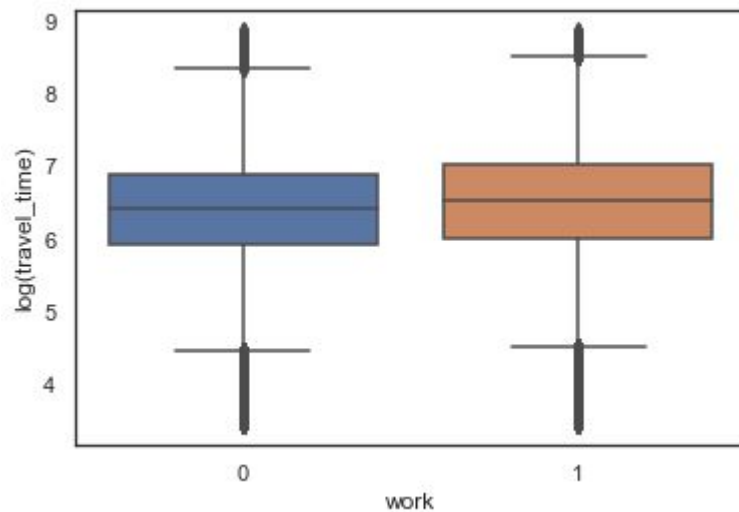
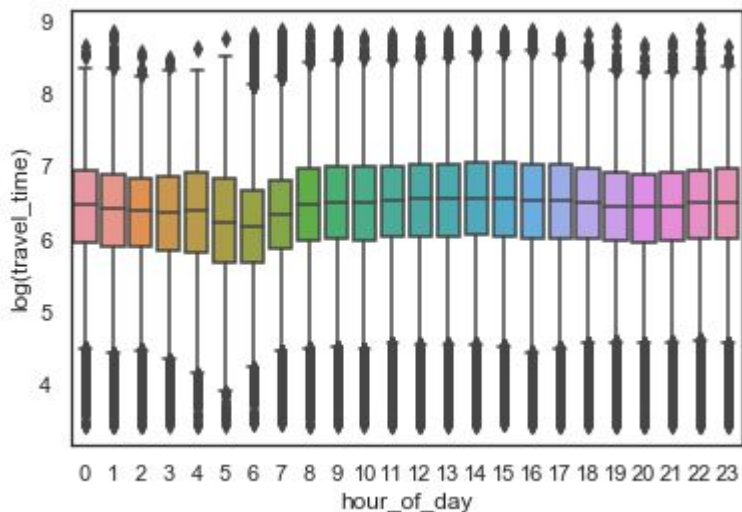
# Day of week and weekday

- Day\_of\_week: Mon to Fri seem to have longer travel time
- From day\_of\_week: weekday (1 or 0)



# Hour of day and work

- Hour\_of\_day: 8 am - 6 pm seem to have longer travel time
- From hour\_of\_day: work (1 or 0)



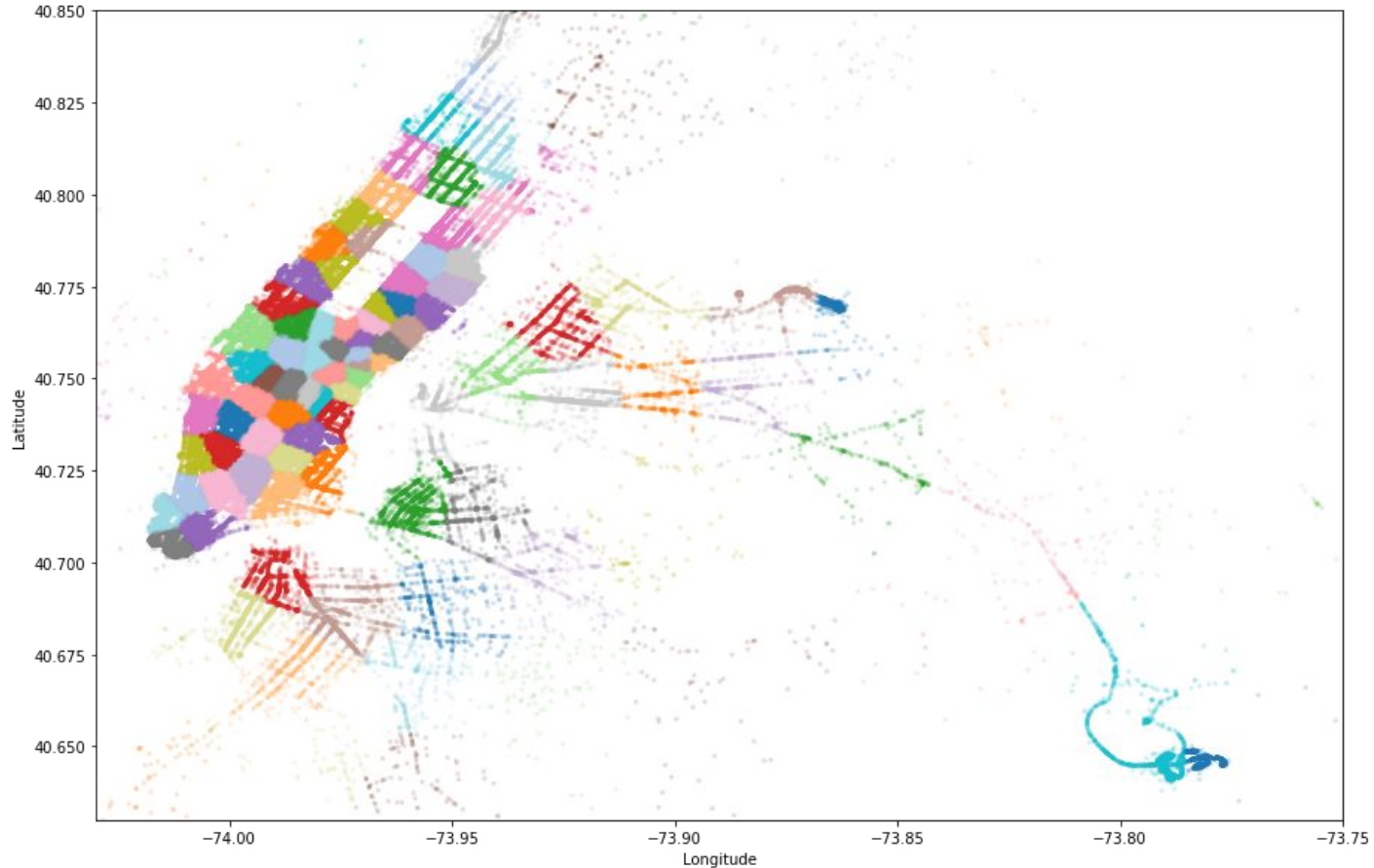
# Clustering

- Center\_latitude, center\_longitude: middle point between pickup and dropoff.  
Use for later

Clustering all lat/lon pairs:

- K-Means: into 100 small neighborhoods using euclidean distance
- Create pickup\_cluster, center\_cluster and dropoff\_cluster for each sample
- Can feed into linear models after one-hot-encoding

# Clustering visualization

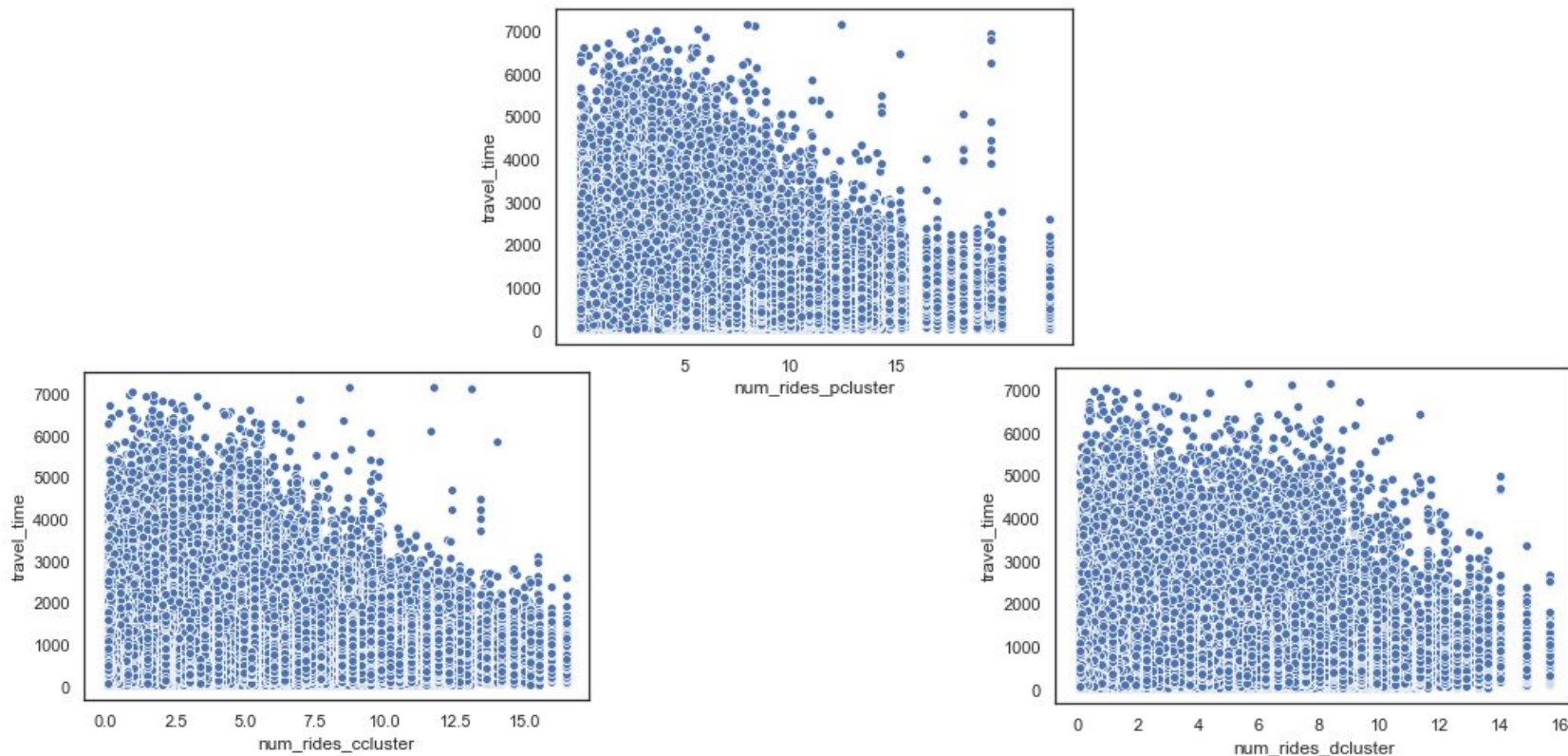


# Traffic

- For each sample, average number of rides in that hour of that week day in its pickup cluster, center cluster and dropoff cluster separately.
- Can show number of yellow cabs in beginning, middle and end of the trip
- Potentially can indicate traffic?



# Traffic

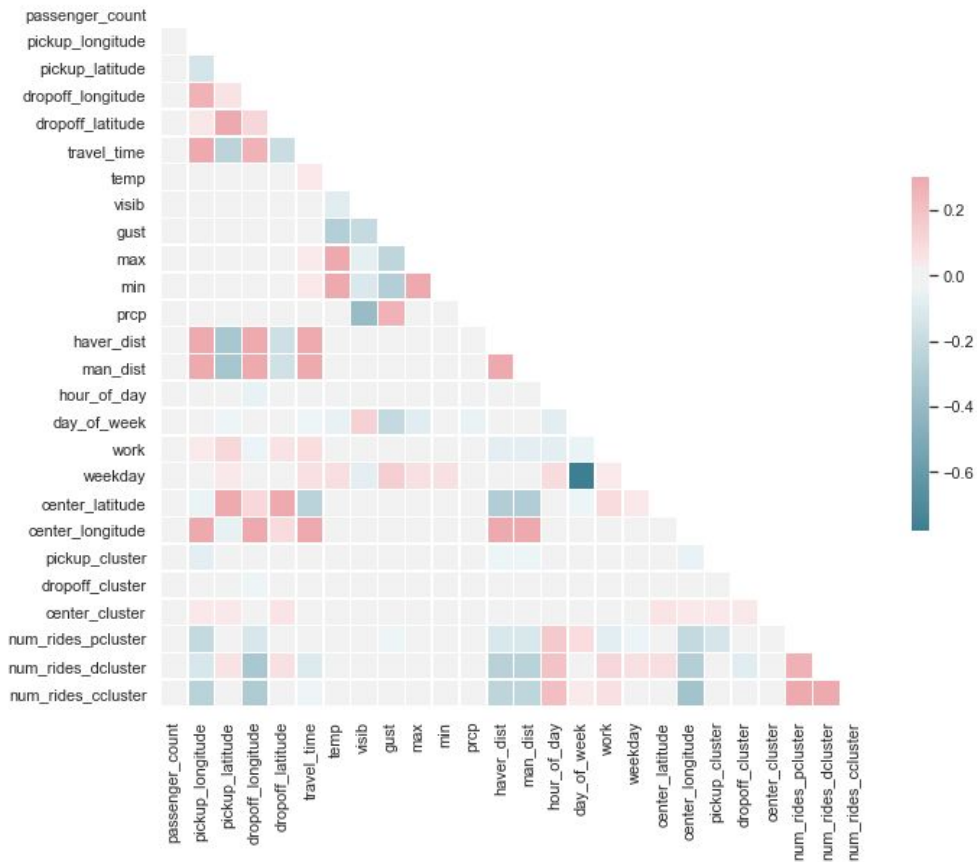


# Traffic

Why the downward trend?

- Maybe when there are a lot of taxis picking up and dropping off riders in an area at a certain time, it means the traffic is actually moving instead of congested. Therefore the shorter travel time?

# Correlations



# Model Training & Model Selection

# Preprocessing

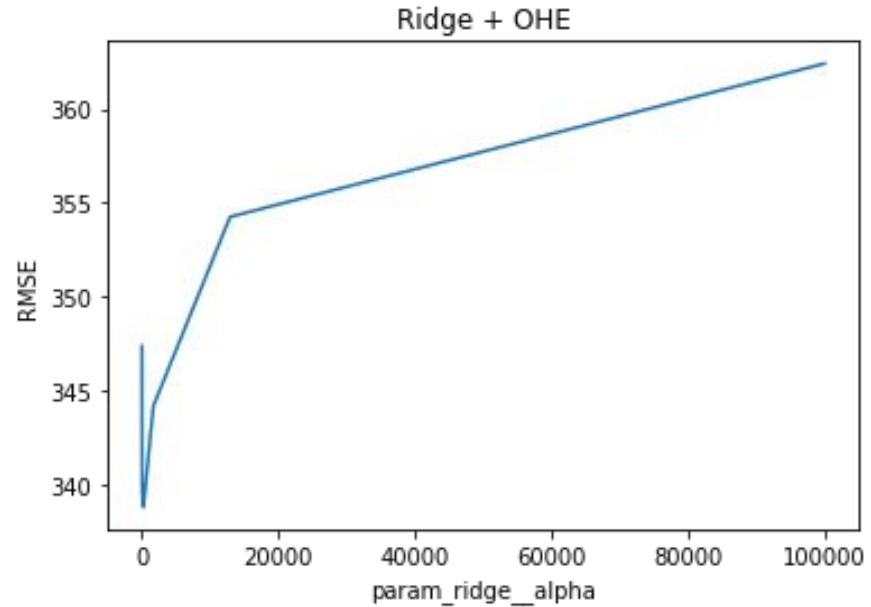
- **OneHotEncoder:** categorical features
- **TargetEncoder:** categorical features with > 15 categories , e.g. pickup\_cluster. (For each cluster, what is the average travel\_time? ) Instead of 100 columns using OHE, it's just one column
- **PowerTransformer:** make distributions of numerical features closer to normal
- **StandardScaler:** scale numerical features
- **ColumnTransformer:** do different preprocessing on different column types
- **Interaction features:** add interaction terms
- **Pipeline:** prevent test info leaking
- **GridSearchCV:** training and tune/select hyperparameters

# Results

| Model   | Test RMSE | Test R2 |
|---|-----------|---------|
| <b>Ridge</b> (OHE)                            | 458       | 0.56    |
| <b>Ridge</b> (Target + OHE + Interaction)     | 446       | 0.58    |
| <b>Ridge</b> (OHE + Target + Power + Scaling) | 487       | 0.50    |
| <b>Lasso</b> (OHE + Target + Power + Scaling) | 493       | 0.49    |
| <b>XGBoost</b> (OHE)                          | 393       | 0.67    |
| <b>XGBoost</b> (OHE + Target)                 | 374       | 0.70    |
| <b>RF</b> (OHE + Target)                      | 380       | 0.69    |

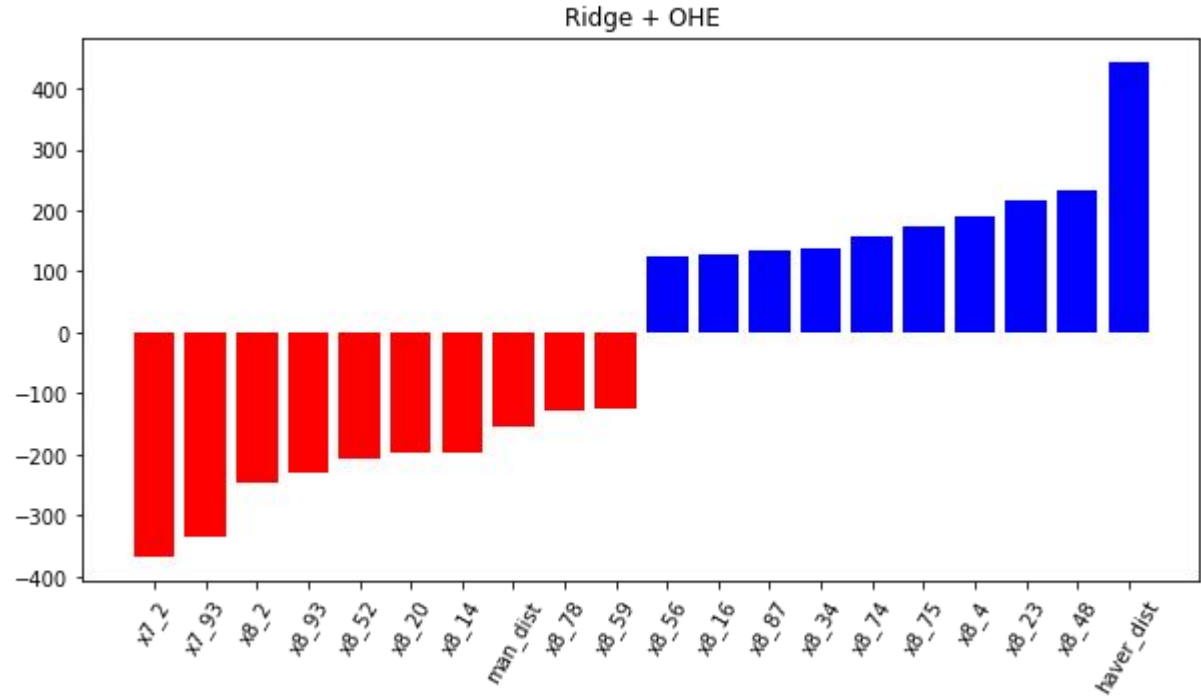
# Model 1: Ridge + OHE

- Test RMSE = 458
- Test R2 = 0.56
- Best alpha = 215.44



# Model 1: Ridge + OHE

- Haversine
- Manhattan
- Pickup cluster
- Dropoff cluster

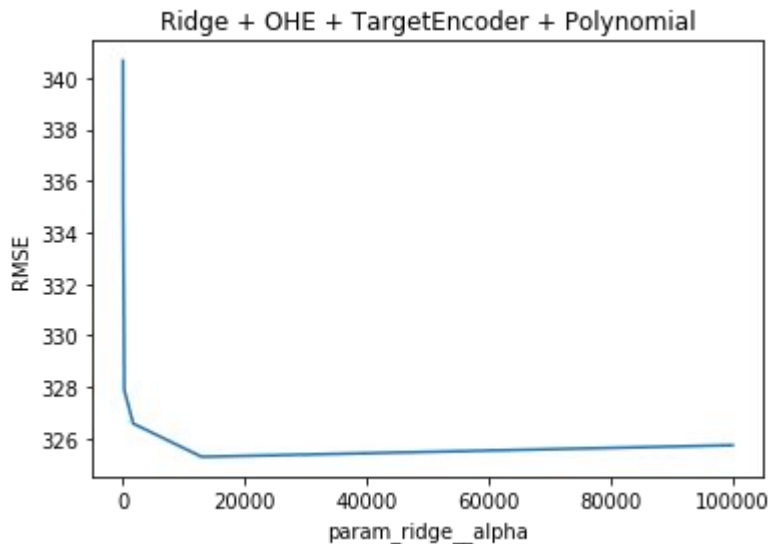




# Model 2:

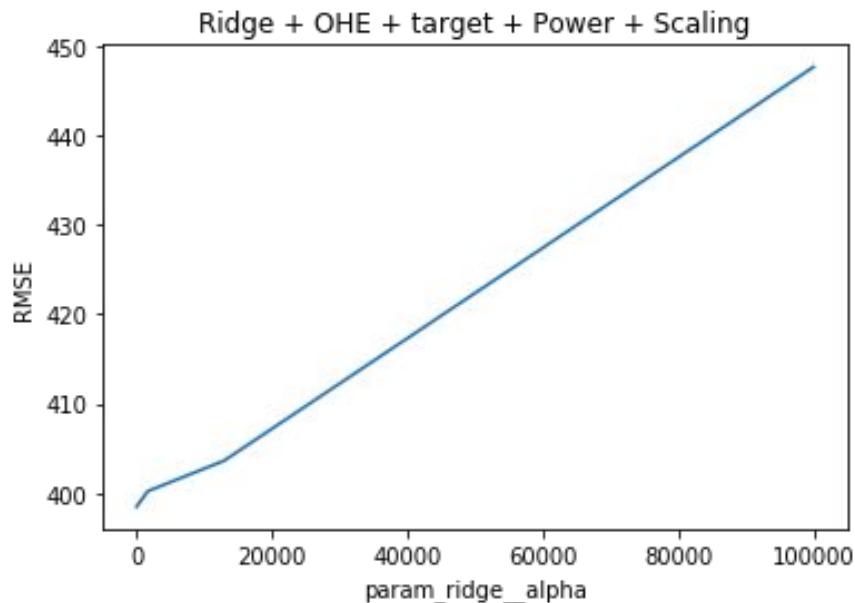
## Ridge (OHE, TargetEncode, Interaction)

- Test RMSE = 446
- Test R2 = 0.58
- Best alpha = 10000
- **Best among linear models**
- **Good for interpretability**



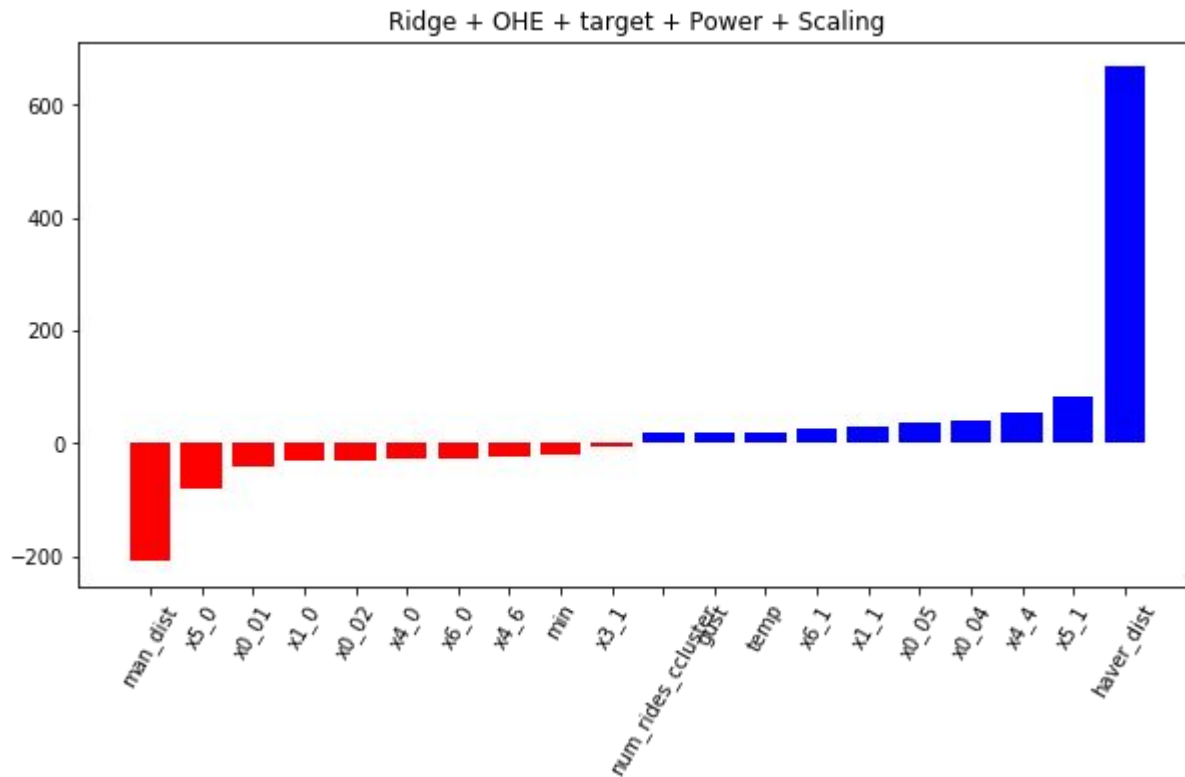
# Model 3: Ridge (OHE, TargetEncode, PowerTransform, Scaling)

- Test RMSE = 487
- Test R2 = 0.50
- Best alpha = 27.83



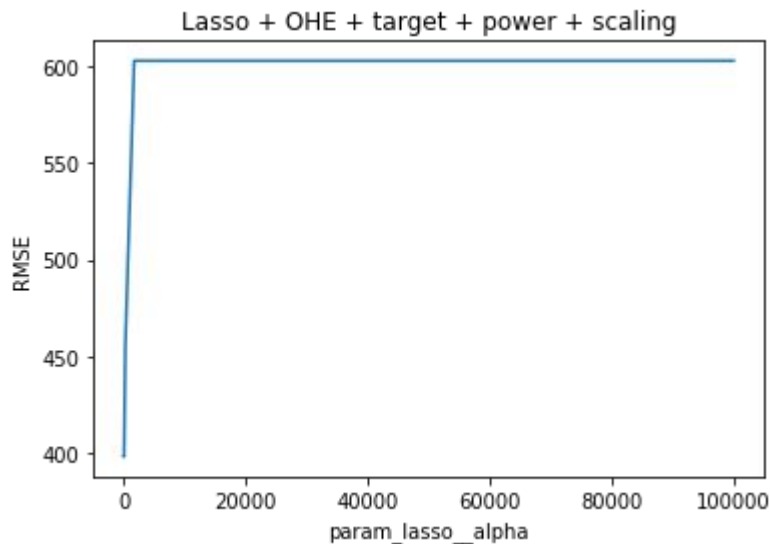
# Model 3: Ridge (OHE, TargetEncode, PowerTransform, Scaling)

- Haversine
- Manhattan
- Work
- Month
- Fog
- Day of week



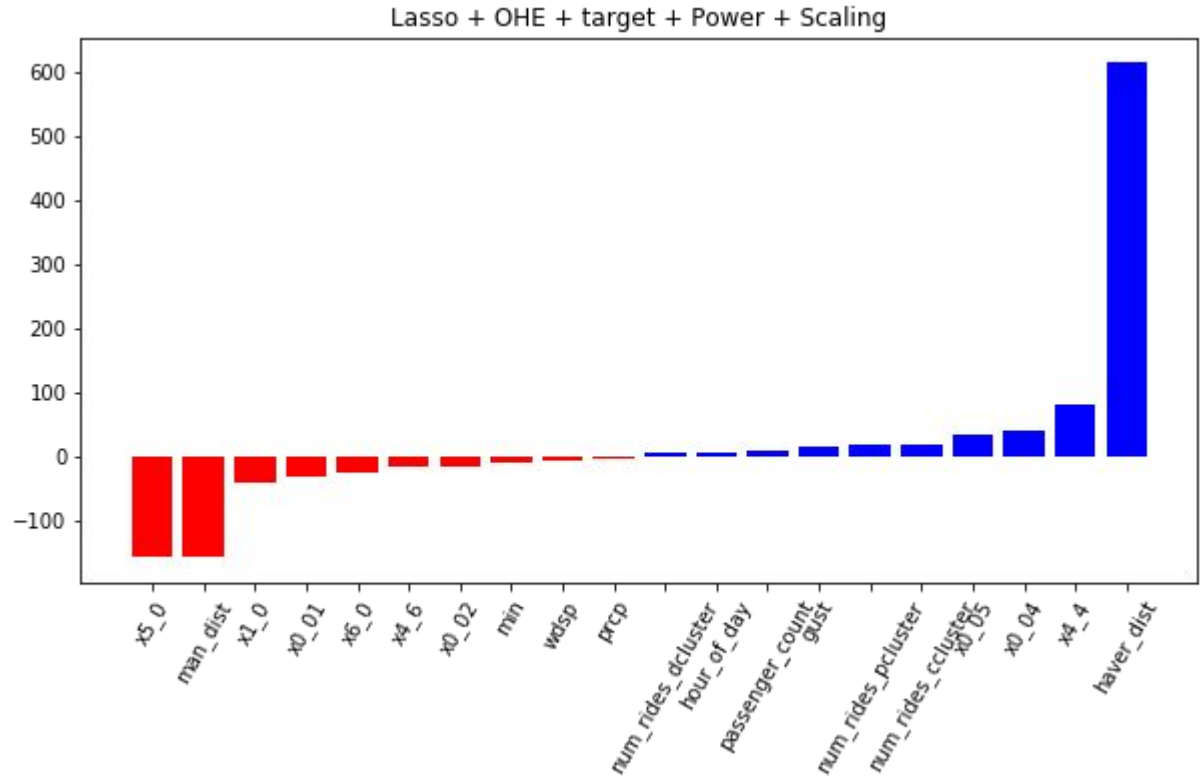
# Model 4: Lasso (OHE, TargetEncode, PowerTransform, Scaling)

- Test RMSE = 493
- Test R2 = 0.49
- Best alpha = 0.46



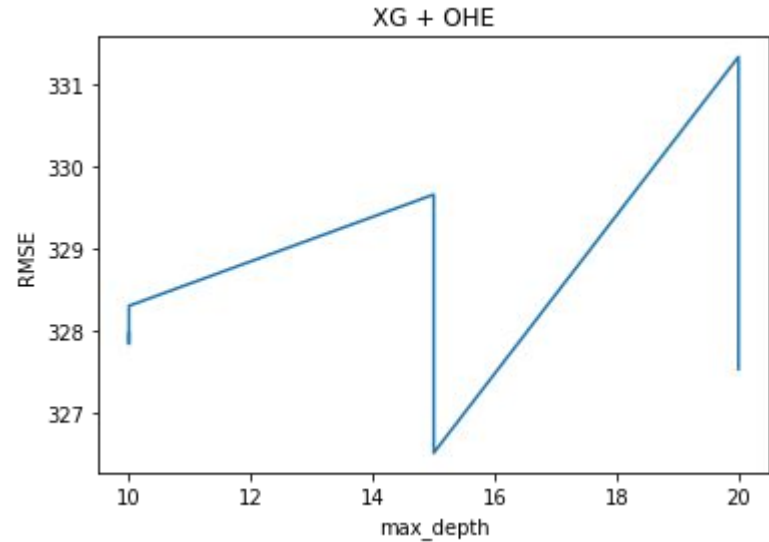
# Model 4: Lasso (OHE, TargetEncode, PowerTransform, Scaling)

- Haversine
- Manhattan
- Month
- Fog
- Day of week
- Work
- Number of rides in center cluster



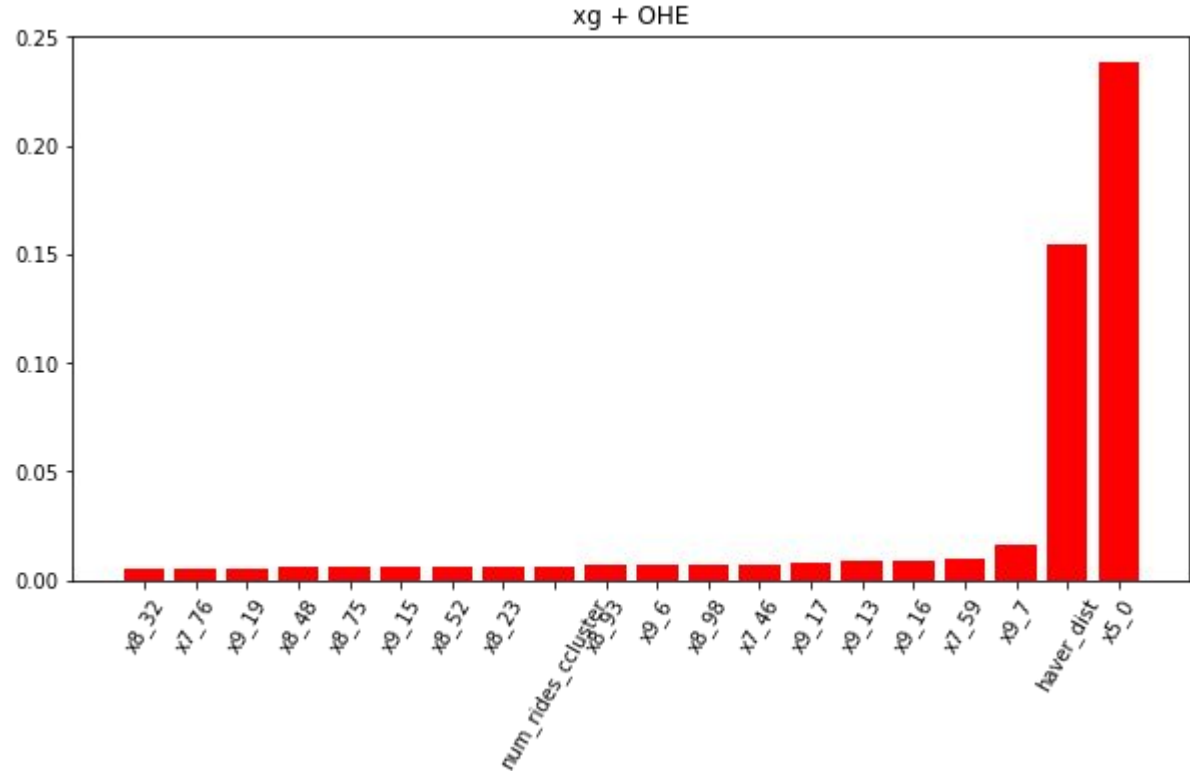
# Model 5: XGBoost (OHE)

- Test RMSE = 393
- Test R2 = 0.67
- Best max\_depth = 15



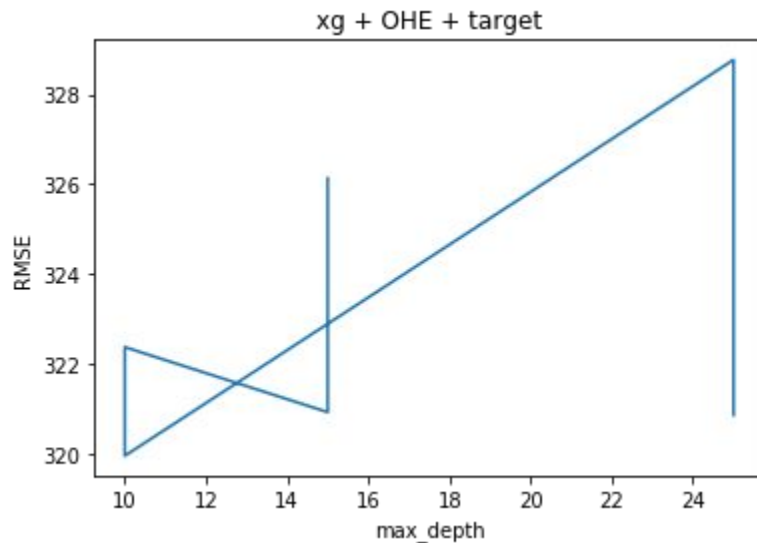
# Model 5: XGBoost (OHE)

- Work
- Haversine
- Pickup cluster
- Dropoff cluster
- Hour of day



# Model 6: XGBoost (OHE + TargetEncode)

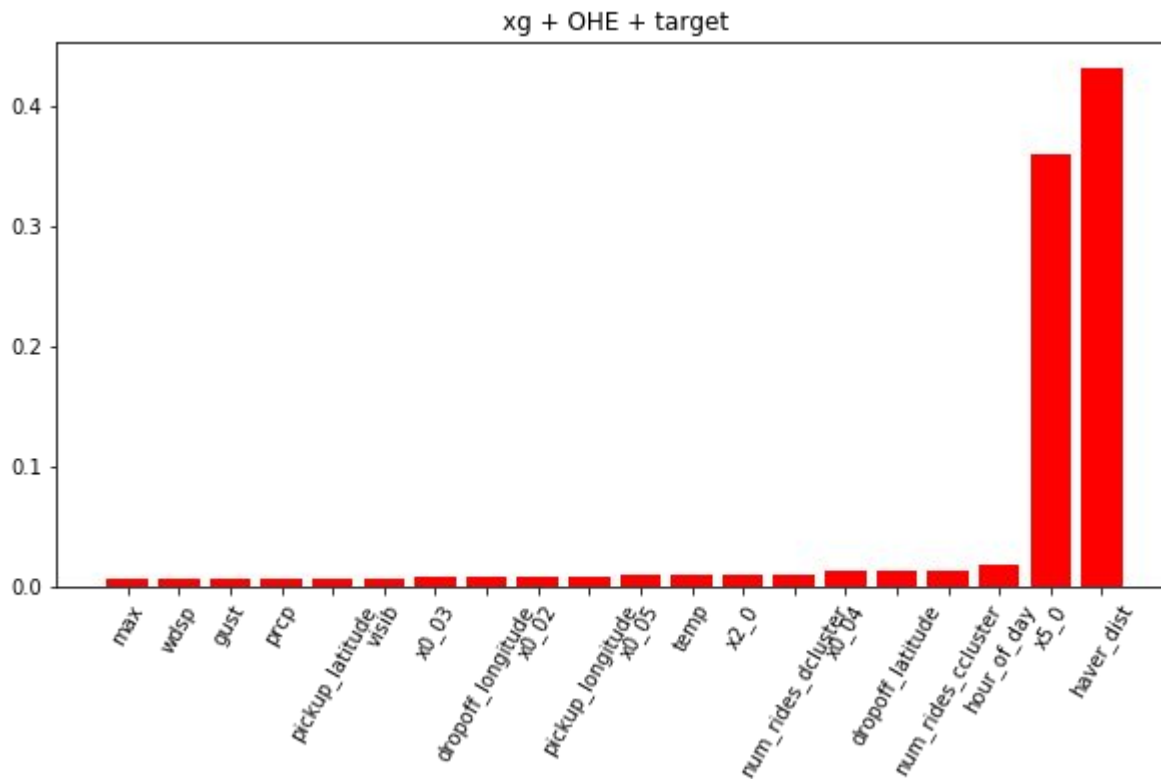
- **Best for prediction**
- Test RMSE = 374
- Test R2 = 0.70
- Best max\_depth = 10





# Model 6: XGBoost (OHE + TargetEncode)

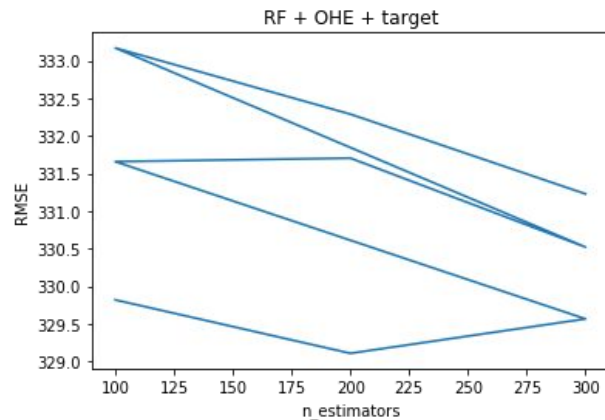
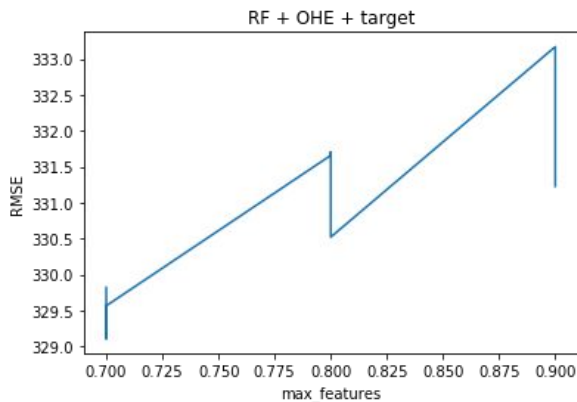
- Haversine
- Work
- Hour of day
- Number of rides in center cluster
- Dropoff latitude



# Model 7:

## Random Forest (OHE + TargetEncode)

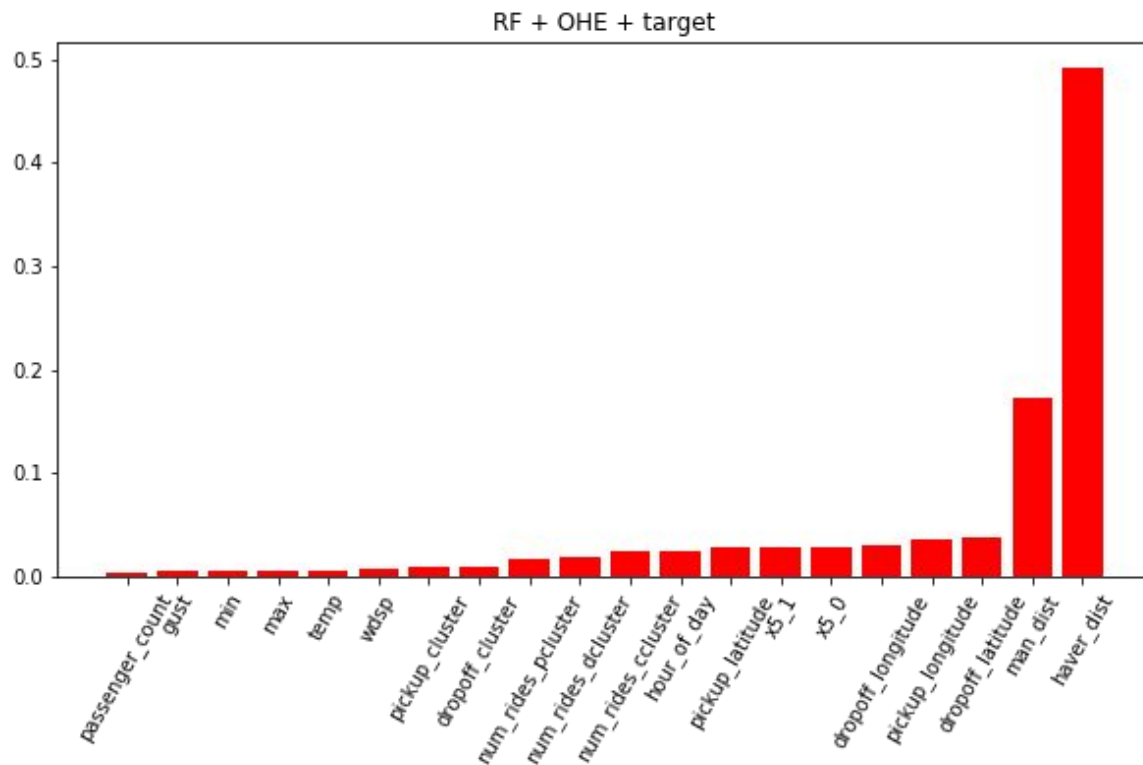
- Test RMSE = 380
- Test R2 = 0.69
- Best max\_features = 0.7
- Best n\_estimators = 200



# Model 7:

## Random Forest (OHE + TargetEncode)

- Haversine
- Manhattan
- Dropoff latitude
- Pickup latitude
- Pickup longitude
- Dropoff longitude



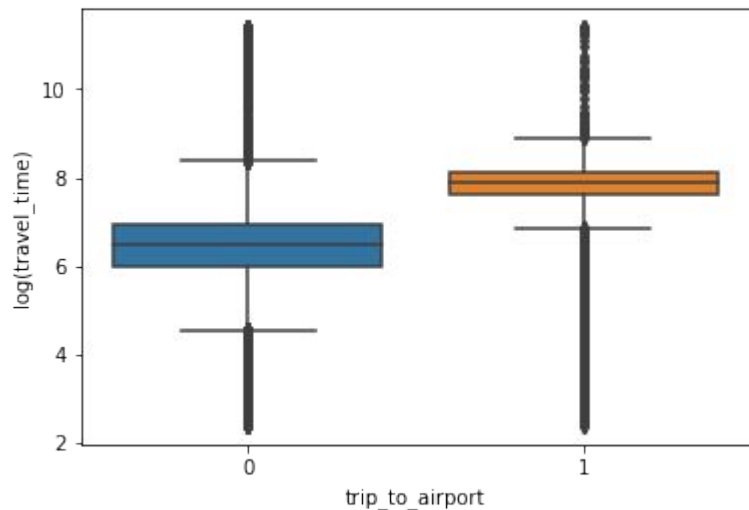
# Conclusion

- Most important features: Haversine dist, Manhattan dist, work, dropoff cluster, pickup cluster
- Ensemble tree-based models > linear models
- Linear models and XGBoost are fast to train

# Important lessons

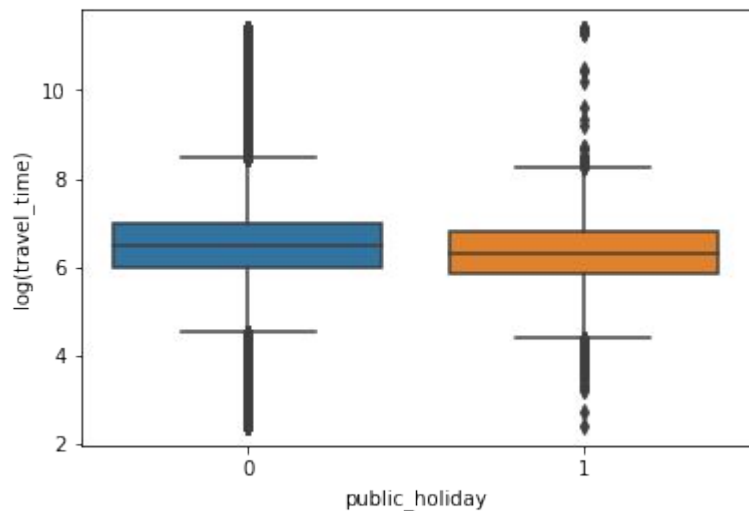
- Removing outliers: very important
- Even tree-based model that are supposedly robust to outliers were performing badly
- Always start with a simple model with simple preprocessing and variable transformations
- Then move on gradually to more complex models with more sophisticated preprocessing
- This way can have an idea of a baseline idea on how complex your model needs to be

# Next Steps: Trip to Airport



Trip with rate\_code == 2 or 3,  
to either JFK or Newark

# Next Steps: NY Public Holiday



→ only from January to June

| Federal & NY Public Holiday | Number of Trips in Dataset |
|-----------------------------|----------------------------|
| New Year                    | 5,023                      |
| Martin Luther King Jr. Day  | 5,087                      |
| President's Day             | 5,119                      |
| Memorial Day                | 3,817                      |
| Independence Day            | 0                          |
| Labor day                   | 0                          |
| Columbus Day                | 0                          |
| Veterans Day                | 0                          |
| Thanksgiving Day            | 0                          |
| Christmas                   | 0                          |
| <b>All Public Holidays</b>  | <b>19,046</b>              |