# Predicting Retention

WENQI CHENG

MARIO DIAZ DE LA ROSA

CLARK LI

# Cleaning up the Data - Target

Retained User = Active 30 days prior to data pull.

Data pulled on July 1$^{st}$, 2014.

**Consider user retained whenever**

$$\textbf{last\_trip\_date} \geq \textbf{'2014-06-01'}$$

This defines our classification labels → find factors that best predict value of **1**

**All features are arguably, intuitively relevant, so we decided to evaluate them all.**

# Cleaning up the Data - Numbers

Additional engineering:

- Dummified phone (Android = 1, iPhone = 0)

- Dummified city (King's Landing and Winterfell)

- Created an Account_Age metric by calculating number of days from signup_date to *data pull*
    - *NOT* to last_trip_date to prevent leakage

# Cleaning up the Data - Ratings

The Ratings Question:

- *what do the NaN represent?*

- *are we losing a valuable signal by removing NaN data?*

- *can the way customers use ratings inform our decision?*

Trade-off between potential information loss and additional data.

# Cleaning up the Data – Douchebag Matrix

# Cleaning up the Data - Ratings

# Naïve Bayes

Classifier applies Bayes Theorem assuming independence between features.

For the validation data set:

Multinomial NB accuracy = 68.9%

Gaussian NB accuracy = 73.8%

Not the best model to use due to low accuracy of MNB and difficult to interpret of justify GNB

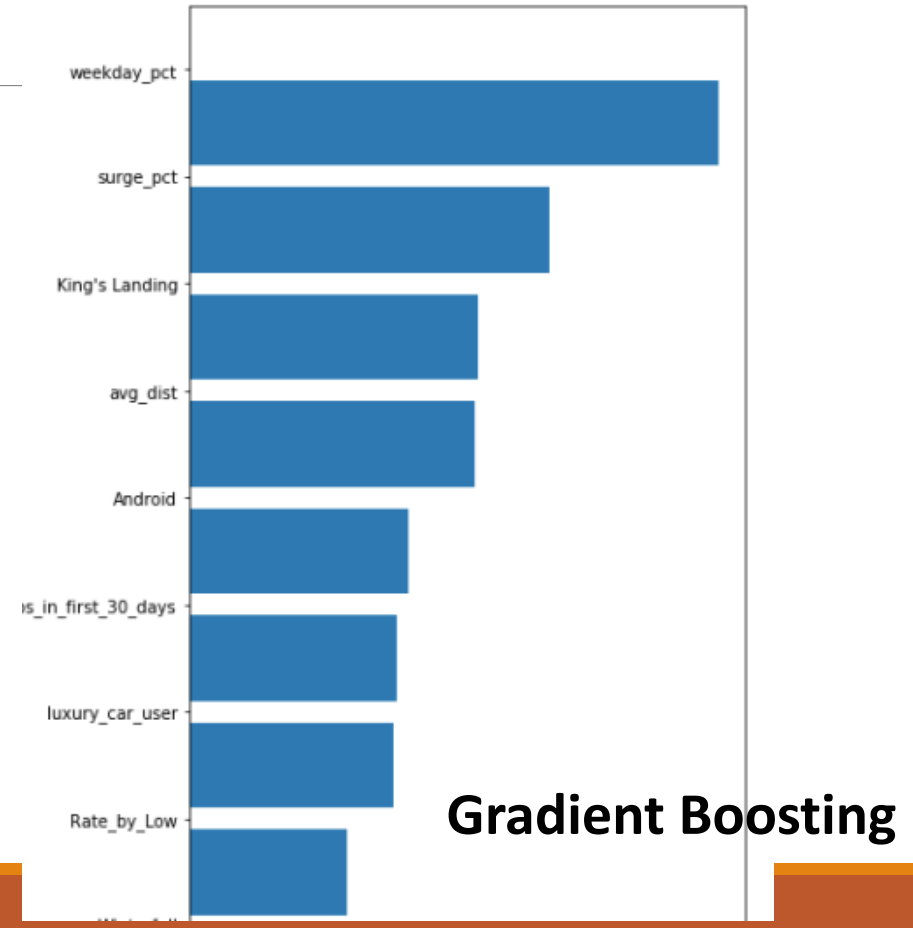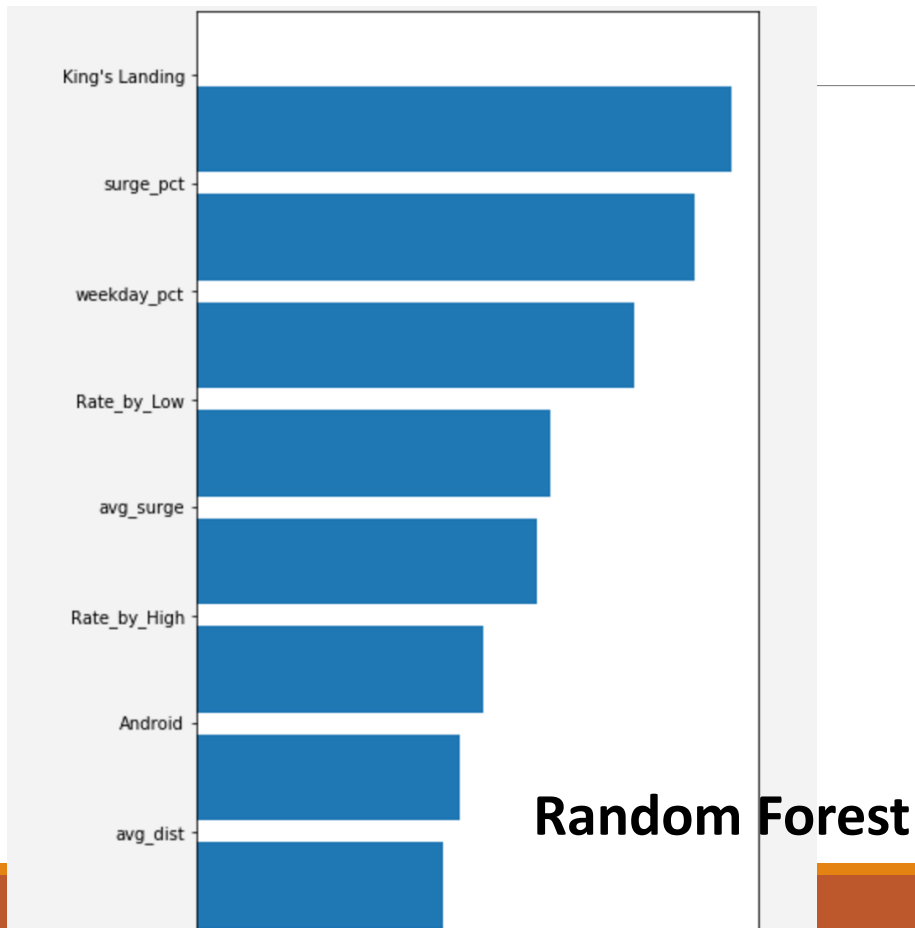# Logistic Regression vs. Random Forest vs. Gradient Boosting

| Logistic Regression | | Predicted | |
|---|---|---|---|
| | | Not Churn | Churn |
| Actual | Not churn | 2217 | 1555 |
| | Churn | 1030 | 5198 |

| | Accuracy | | Precision | Recall |
|---|---|---|---|---|
| Train | 0.75 | | | |
| Test | 0.74 | | 0.68 | 0.59 |

| Random Forest | | Predicted | |
|---|---|---|---|
| | | Not Churn | Churn |
| Actual | Not churn | 2370 | 1402 |
| | Churn | 867 | 5361 |

| | | | Precision | Recall |
|---|---|---|---|---|
| Train | 0.81 | | | |
| Test | 0.77 | | 0.73 | 0.63 |

| Gradient Boosting | | Predicted | |
|---|---|---|---|
| | | Not Churn | Churn |
| Actual | Not churn | 2408 | 1364 |
| | Churn | 883 | 5345 |

| | | | Precision | Recall |
|---|---|---|---|---|
| Train | 0.79 | | | |
| Test | 0.77 | | 0.73 | 0.64 |

# Important Features



**Random Forest**

**Gradient Boosting**

# Conclusion & Future Work

- Random Forest and Gradient Boosting gave better predictions

- Focus on:
  - Weekday Trip Percentage
  - Percentage of Trips with Surge Multiplier > 1

- Further validation of our models is necessary
- Need more data that reflects broader distribution of users

# Team Kuma



Questions?