

Airbnb New User Bookings

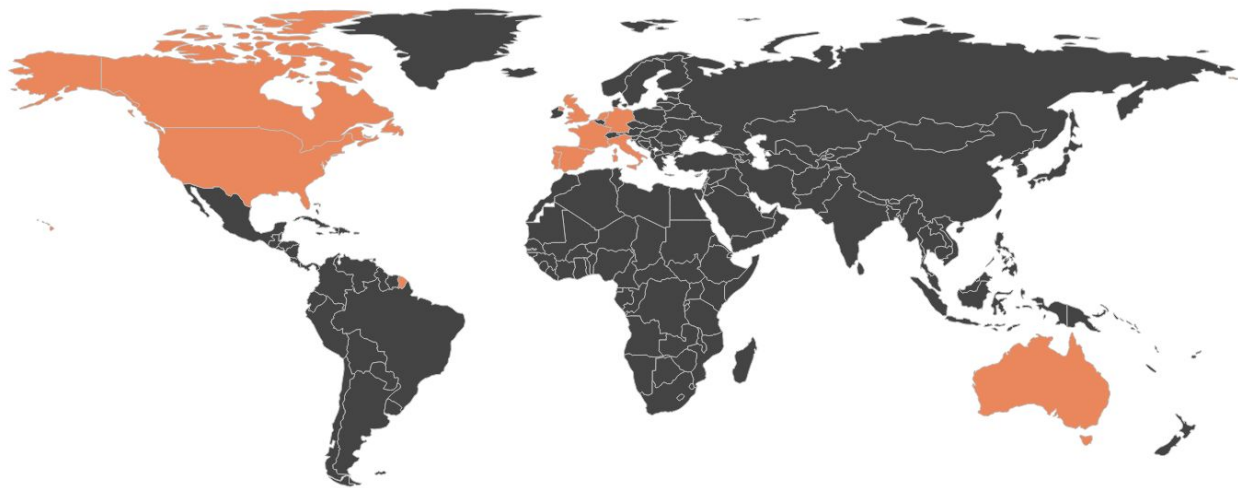


Robert Castellano

In collaboration with: Zi Jin, Yannick Kimmel, Michael Winfield

Introduction

- I. Goal of the project
- II. The data
- III. Insights into the data
- IV. Our strategy
 - A. Feature engineering
 - B. Stacking
 - C. Feature importance
- V. Results
- VI. Conclusions



Goals

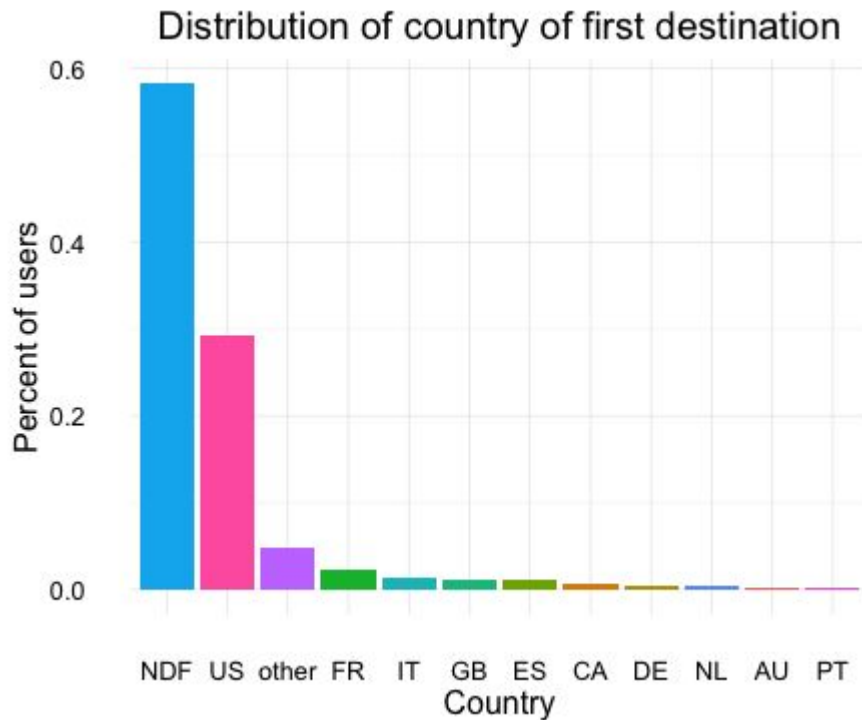
- Kaggle competition hosted by Airbnb, ending Feb 2016.
- Goal: Predict the country of a new user's first destination. This can include not booking (NDF).
- The competition allowed by the submission of five suggestions for each user.
- The competition was graded on normalized discounted cumulative gain (NDCG), which measures the performance of a recommendation system based on the relevance of the recommended entries.

Airbnb Kaggle Dataset

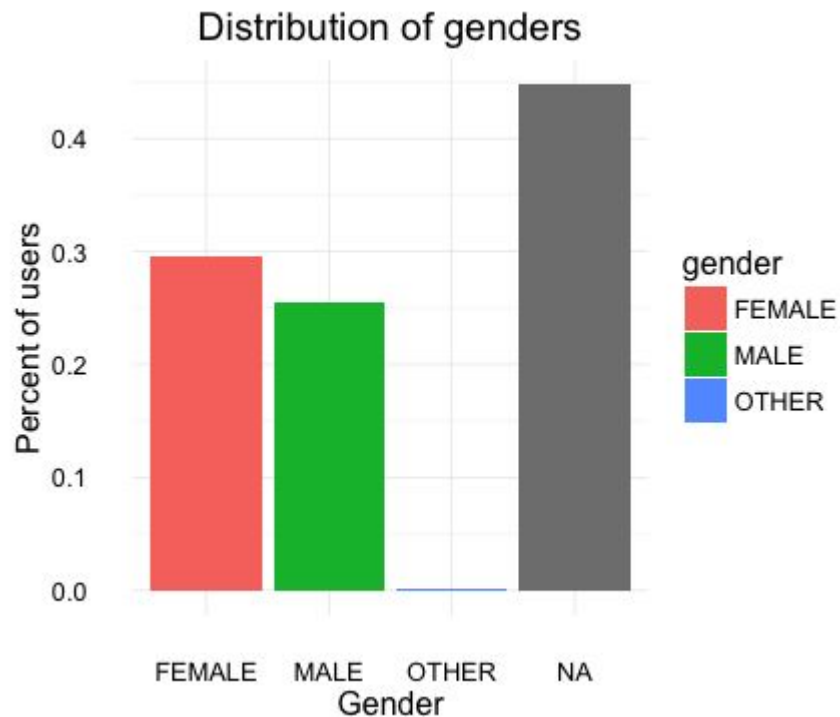
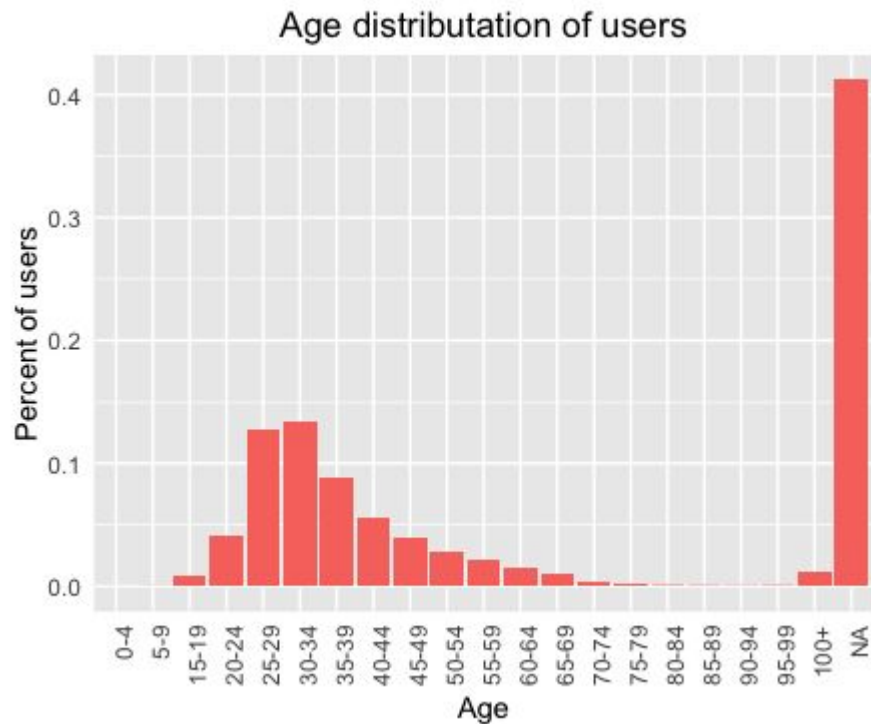
The Airbnb Kaggle dataset consisted of:

- **User information:** Age, gender, web browser, avenue in which the user accessed AirBnB, country destination, timestamp of first activity, account created, and first booking.
- **Browser session data:** Action type, and time elapsed.
 - We included the number of actions a user took and the total time spent on the website in our model.
- Training set: 200,000 users--Jan 2010 to Jun 2014
Test set: 60,000 users--July 2014 to Sep 2014

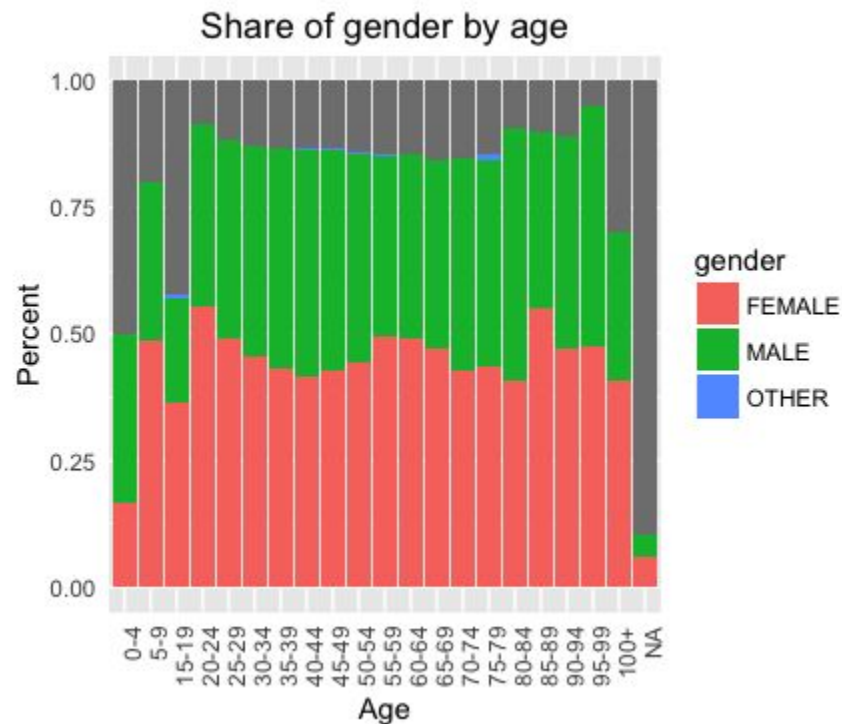
Airbnb User Booking Behavior



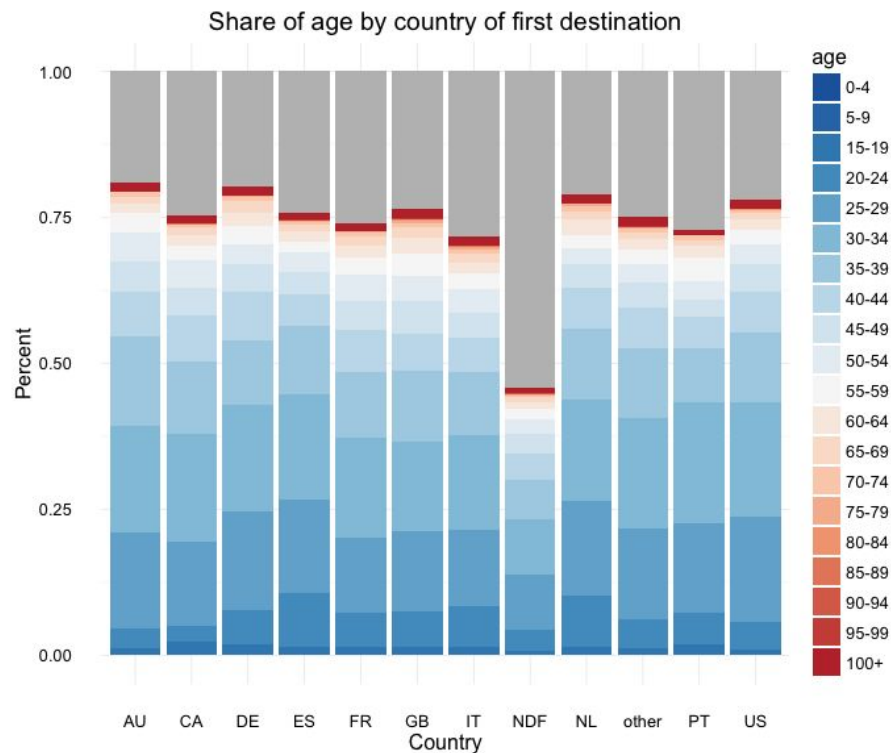
User demographics



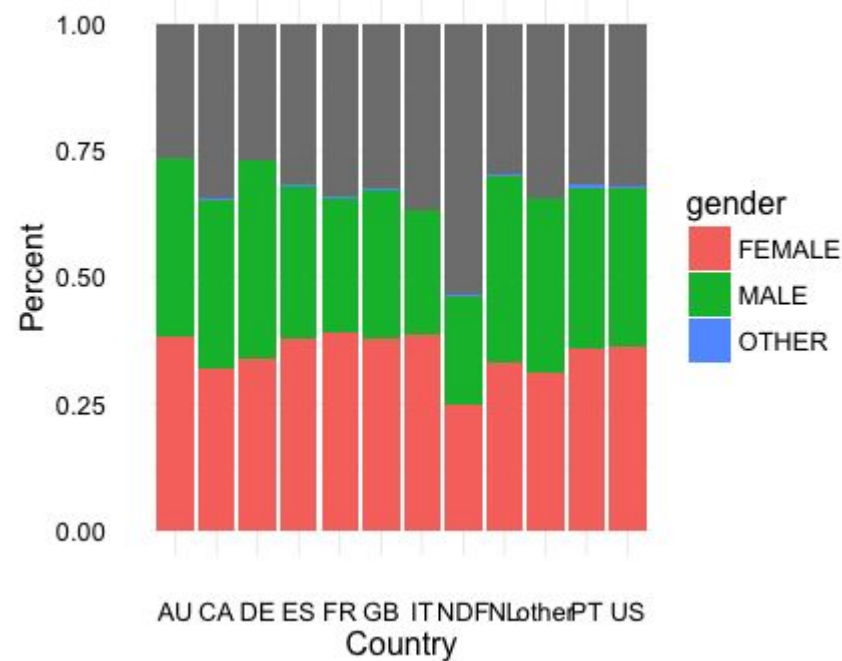
Age/Gender Missingness



Age & Gender on Country Destination

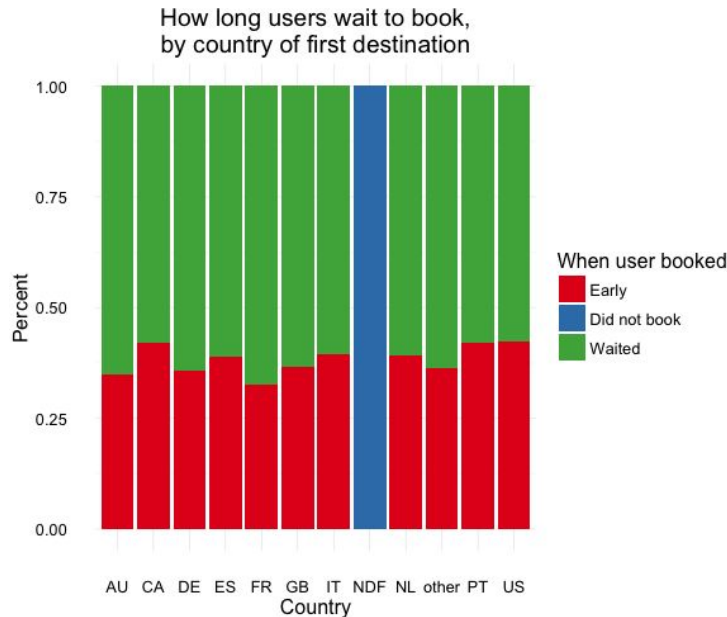
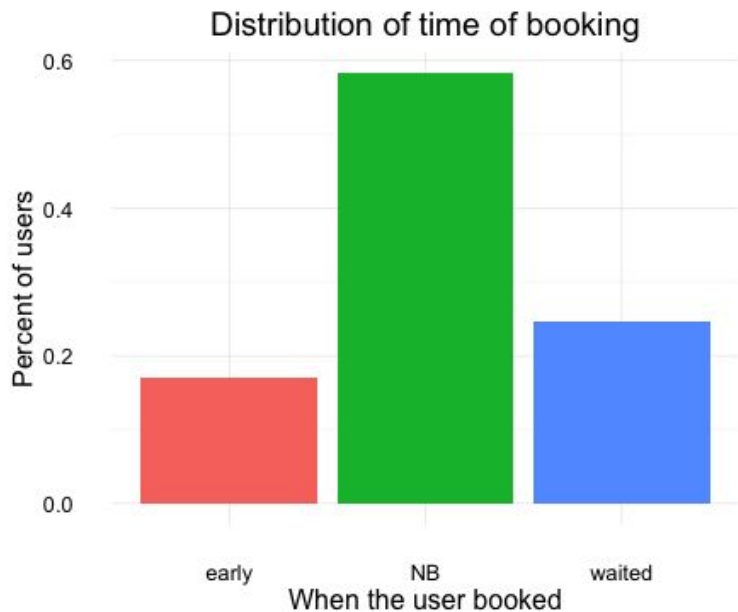


Share of gender by country of first destination



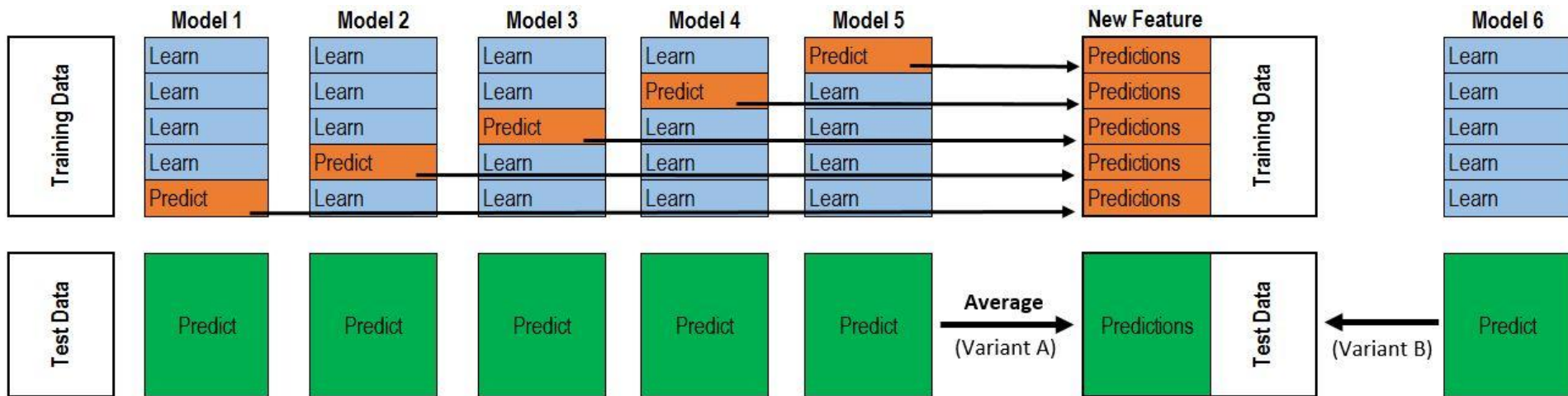
Time variable feature engineering

- *We decided to engineer 3 features based on user booking behavior, specifically the time between the creation of Airbnb accounts, a user's first activity on the website, and their date of first booking.*

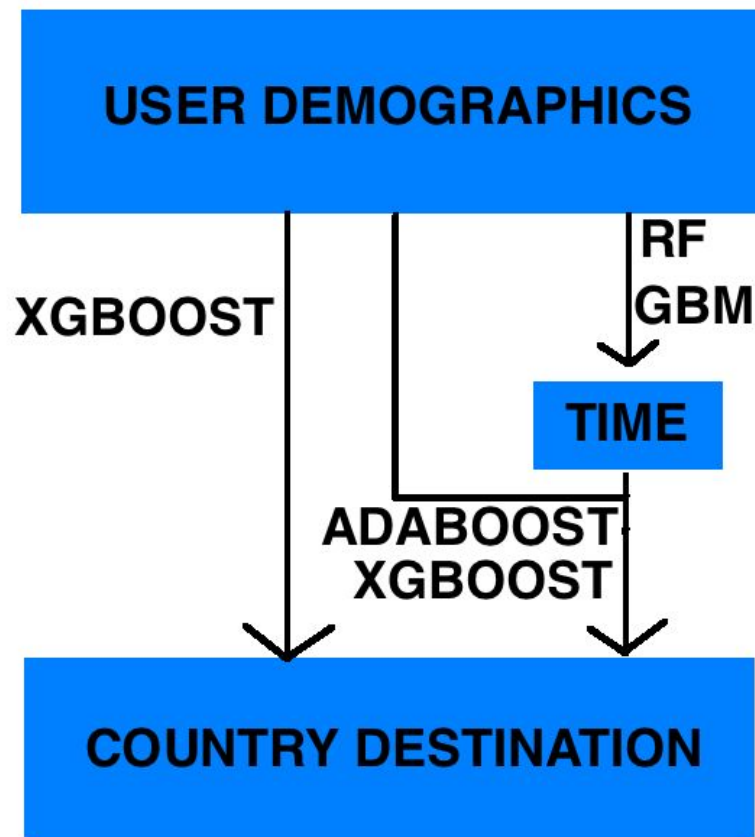


Stacking

- *Out-of-fold predictions of those three features were then added to the training dataset and test dataset through the process of stacking.*



Workflow

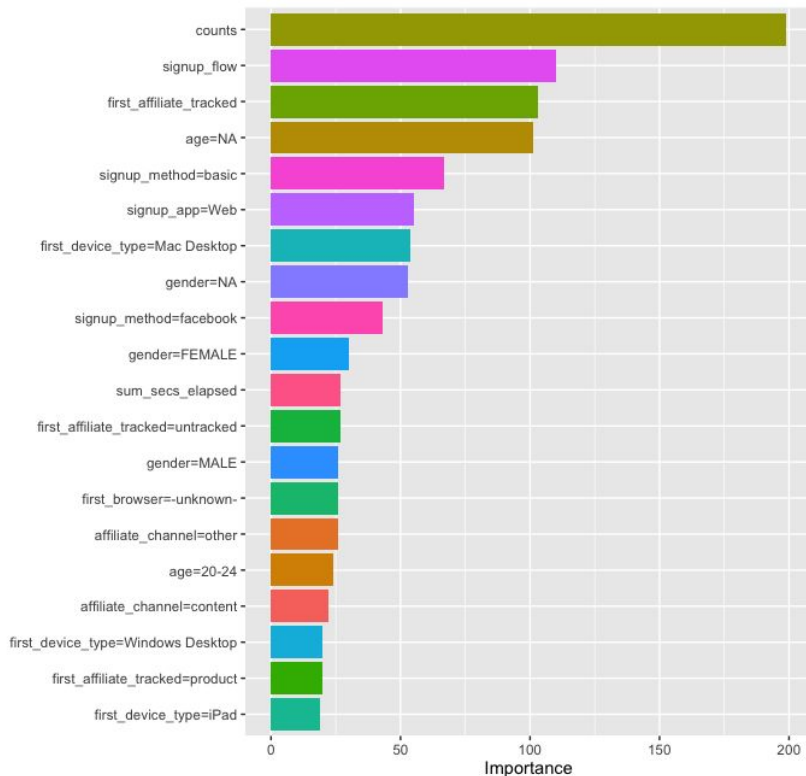


Predicting Country Destination

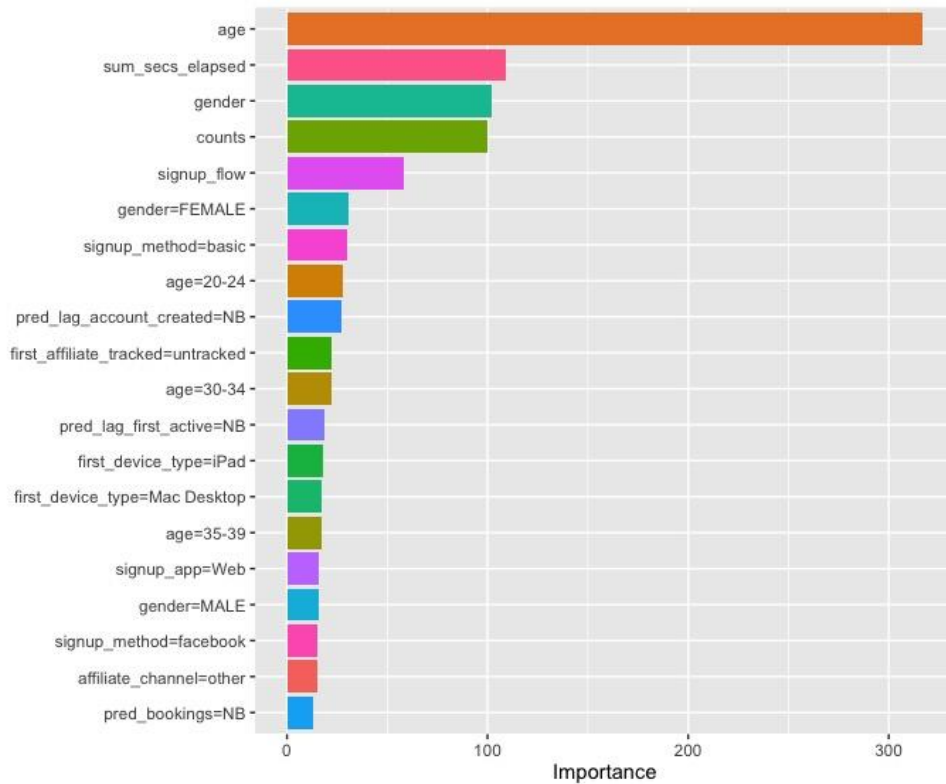
- First choice of either NDF or USA.
- Ran grid search cross-validation.
- Unstacked:
 - XGBoost -- Improved Kaggle ranking from #1165 to #374.
 - Best parameters: learning_rate = 0.1, max_depth = 4, n_estimators: 100
- Stacked:
 - XGBoost -- Kaggle ranking of #1030.
 - AdaBoost -- Kaggle ranking of #1028.

Variables of importance in XGBoost

Unstacked Model



Stacked Model



Summary

1. Performed exploratory data analysis on Airbnb new user information.
2. Data munging in Python and R.
3. Used R for visualization and the creation of a Shiny app.
4. Feature engineered time-lag-based variables using Python and R.
5. Fit models (XGBoost/Random Forest/AdaBoost) using Python.
6. Performed predictions on users using XGBoost--ranked at 374 on Kaggle.

Recommendations to Airbnb

- Invest in collecting more demographic data to differentiate country destinations.
- Flag users who decline to enter age and gender; such behavior correlates with users not booking.
- Continuously collect browser session activity; such data was helpful for predictions. This data was available only for newer users.

Further directions

Steps to improve our predictions:

- Run a multi-layer model to distinguish country of destination of those users who booked.
- Optimize tuning parameters for XGBoost on the stacked dataset.
- Stack country of destination predictions to dataset as features to improve predictions.
- Use multiple XGBoost models (stacked or unstacked) and ensemble them.