

Clicks Activation Study

Pradeep Arkachar

Business Problem

Problem Statement

- A user is looking for loans on Google and enters a search query. Depending on how related the search query is to the business, we choose to bid for the query with the goal of maximizing the user impressions using Google Adwords (Google actually showing the ad to the user when they execute the search). This would lead to more clicks on ads, and potentially more signups, activations, and revenue. Here activation is defined as a user successfully contacted by the call center.

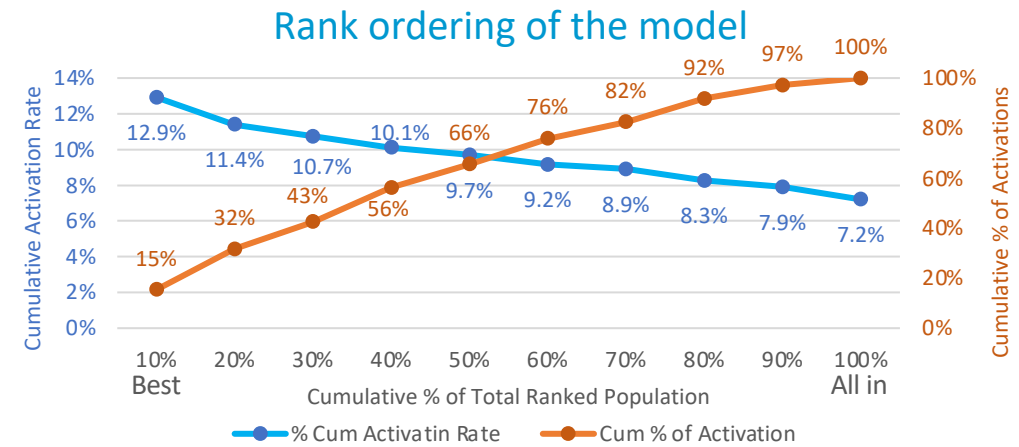
Executive Summary

Objectives

- Develop probability of activation per click model to optimize/maximize the user activation rate arising from the Google Adwords impressions
- Recommend strategy to improve marketing ROI
- Propose engineering architecture for model deployment and discuss limitations

Recommendations

- Developed a probability of activation model that can be used to rank order customers for an efficient marketing dialer strategy shown in the chart above
- This framework can be further enhanced by combining with profitability model to go after the most profitable and converting customers
- Proposed a cloud-based engineering architecture



The chart displays the ranking power of the model by showing the cumulative activation rate from the best 10% (left) to 100% of the population (right)

Insights

- 83% of the customers seek debt settlement - 61% via mobile and 21% via Desktop (higher activation)
- 98% of the customers are not prime - 50/50 split city vs non-city
- High converters in GA, PA, MD and CO; Philadelphia, Honolulu and Denver
- Cash loans have lowest conversion rate (3.8%)

Model Development

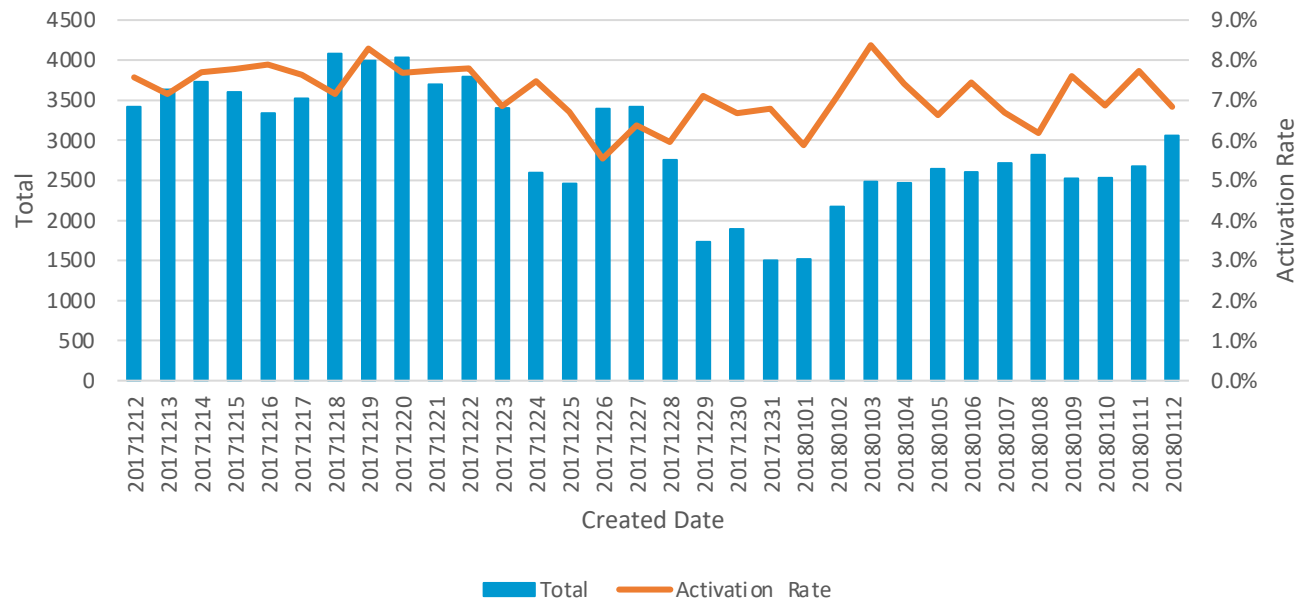
Data Summary

Total	Activated	Activation Rate
94,194	6,787	7.21%

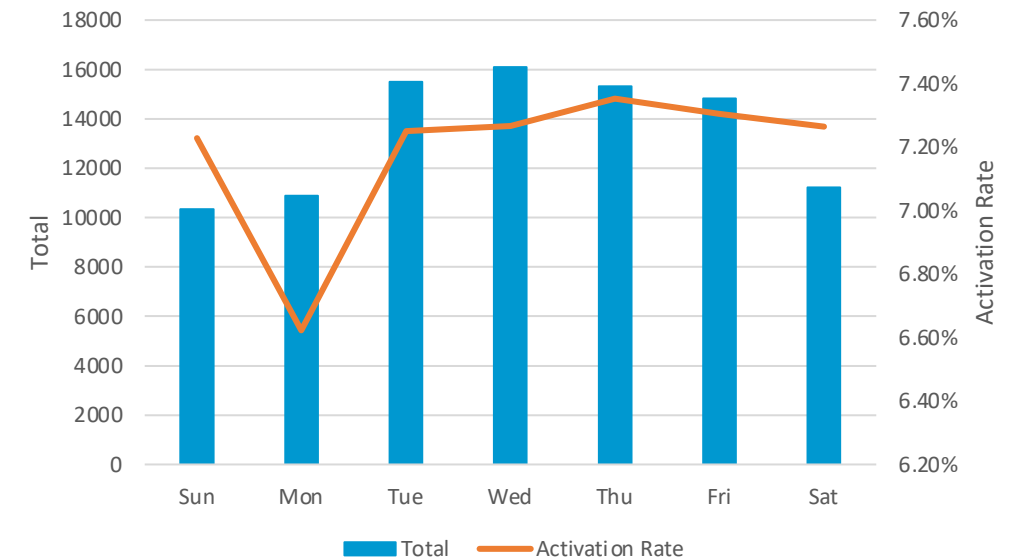
Variable name	Total Missing	% Missing
category_debt_type	9859	10.5%
location_in_query	15	0.0%
campaign_state	6	0.0%
is_activated (target)	0	0.0%
is_hardship	0	0.0%
is_prime	0	0.0%
in_city	0	0.0%
platform	0	0.0%
created_date	0	0.0%

Daily and Weekly Seasonality

Daily Activation Rate

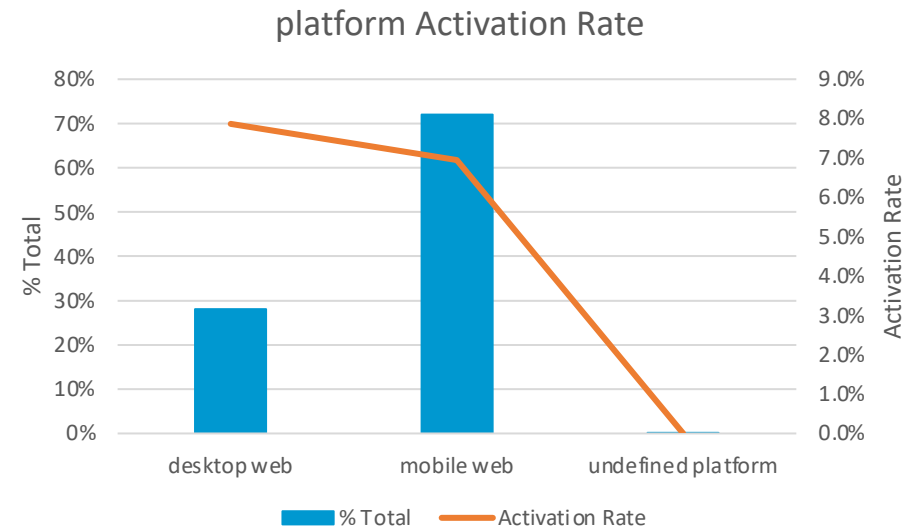
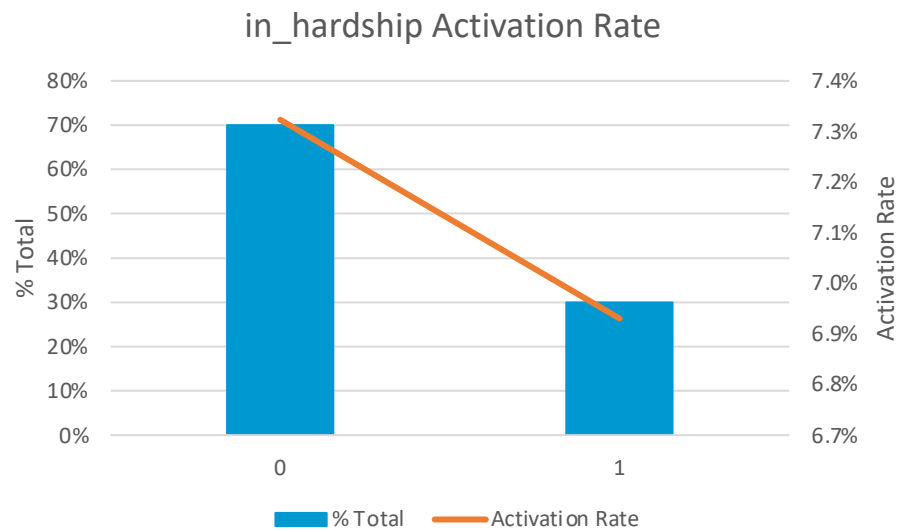
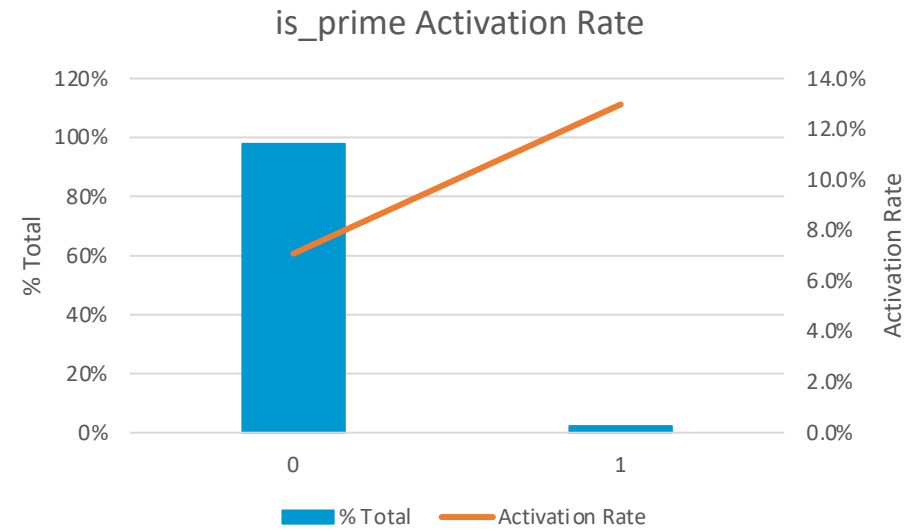
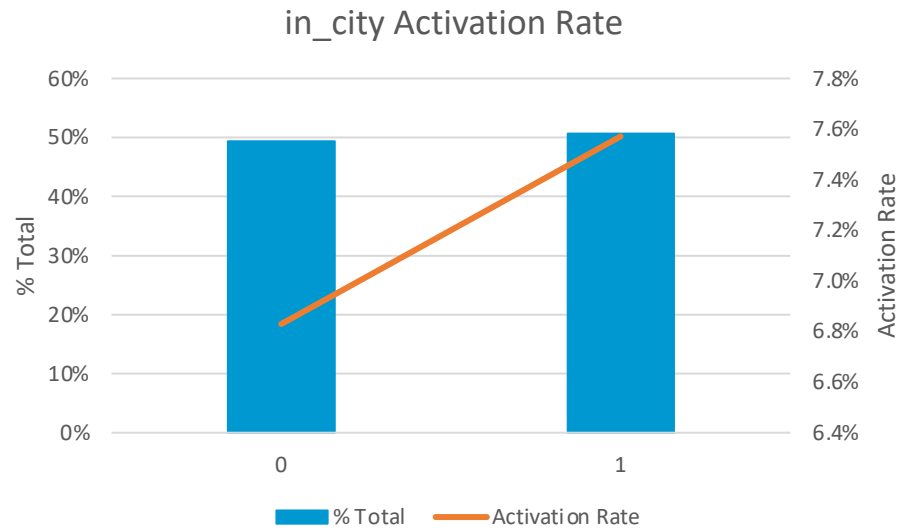


Weekday Activation Rate

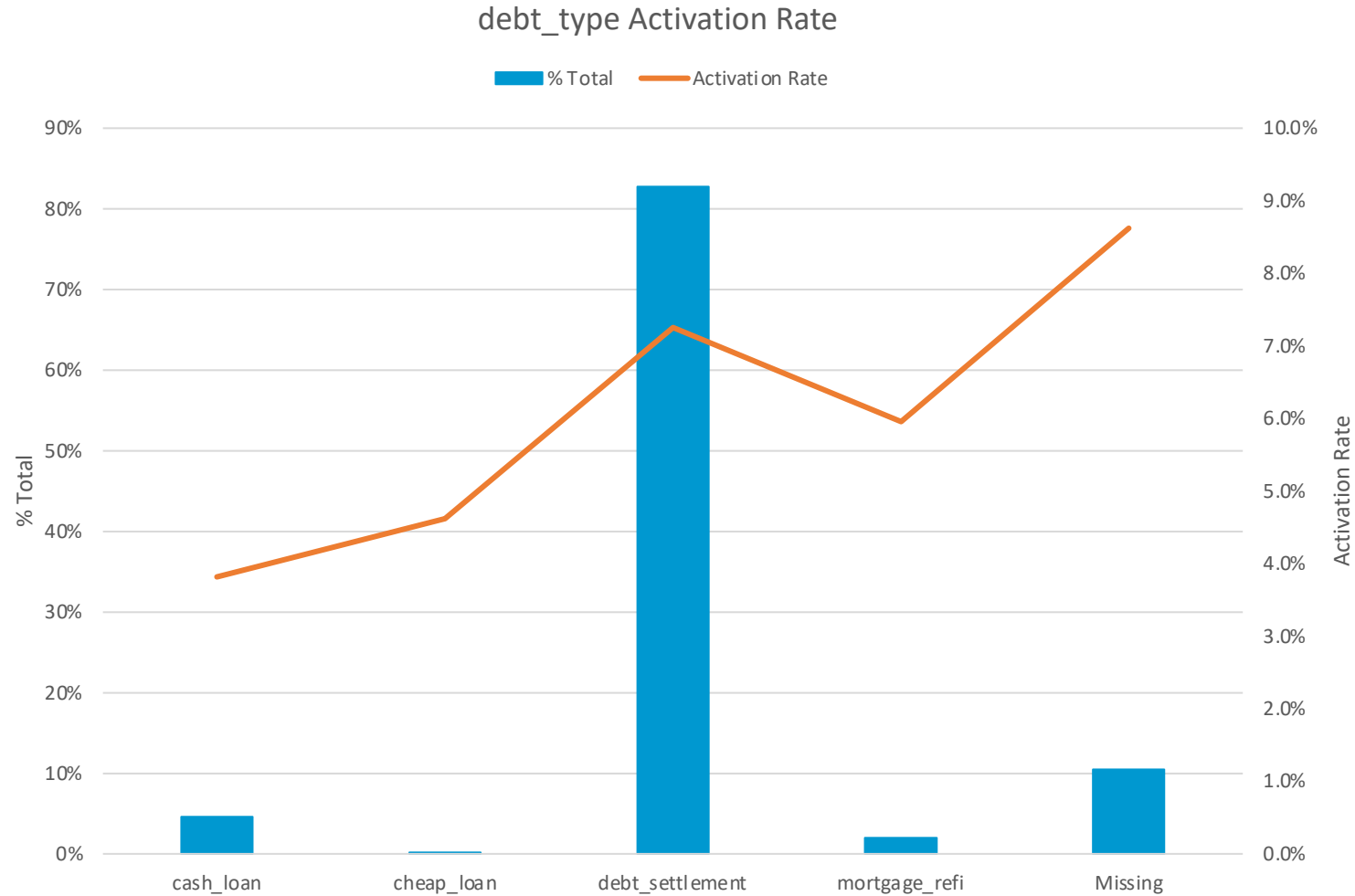


- Mondays showed lower activation rate than other weekdays
- However, given that the data only spans for a month (4 data points per weekday) and the lower activity during the last week of the year, this feature was dropped from consideration
- With more data points, assessment can be made if there is weekly and monthly seasonality

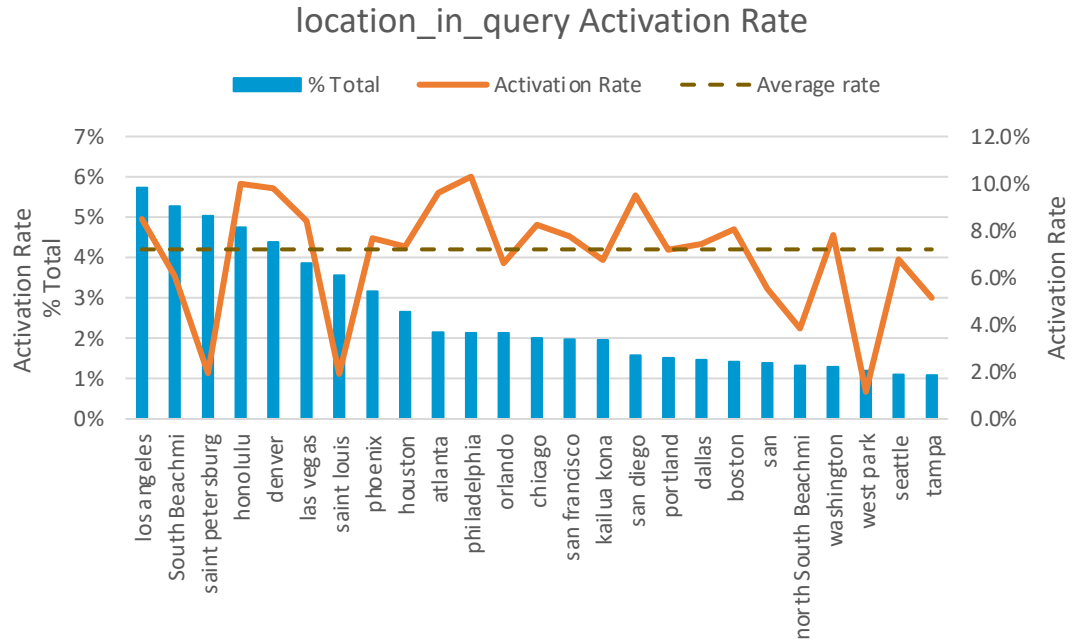
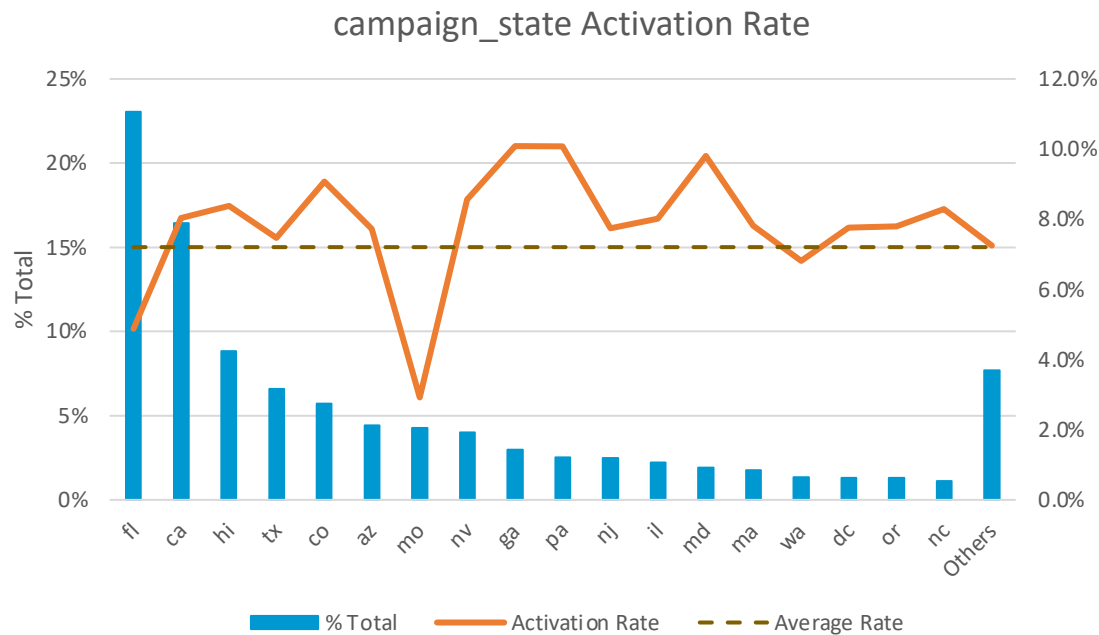
Activation Rates by features



Activation Rates by Debt Type



Activation Rates by State and Location



- The states and locations with total counts less than 1000 were grouped together
- High conversion - GA, PA, MD and CO; Philadelphia, Honolulu and Denver
- Low conversion – MO; Saint Petersburg, West Park and Saint Louis

Modeling Consideration

Data

- 12/12/2017 to 01/12/2018 (~ 1 month)
- Model development data – first 3 weeks (~ 75%) train/test 80%/20%
- Out of time validation – last week (~ 25%)

Feature Engineering

- Location, State – Group categories < 25 Counts into 'Others'
- Missing values of Debt Category was its own group
- N-1 dummy variables for each categorical variable

Algorithm

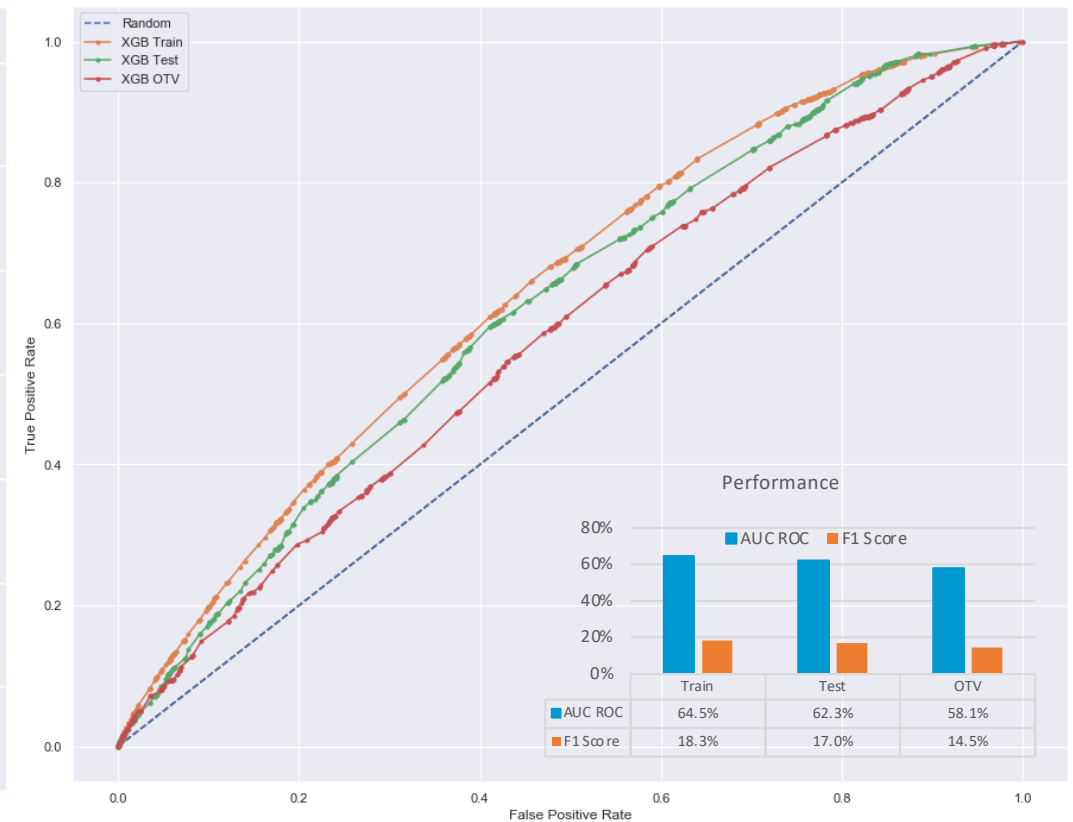
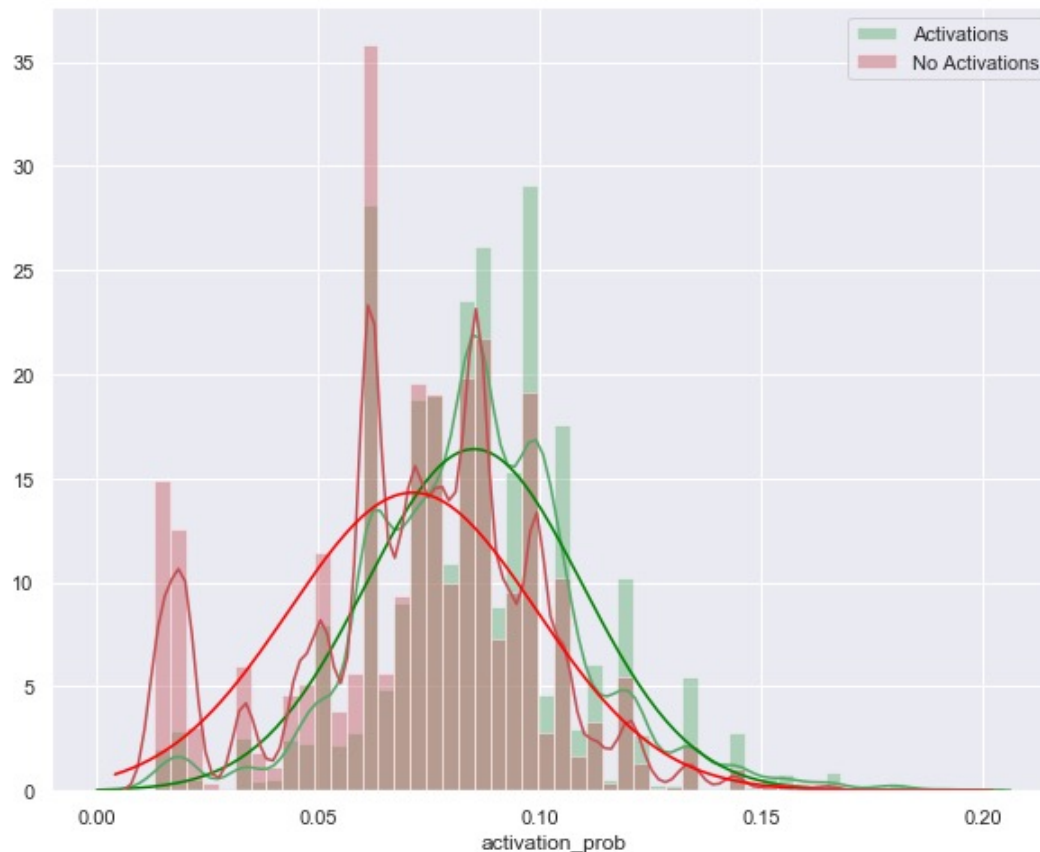
- XGBoost algorithm was chosen because of its fast development and outperformance, while offering some transparency

Validation

- K-Fold Cross validation and hyperparameter tuning was performed to avoid overfitting
- Expected VS actuals plot showed that the model fit the actual data well

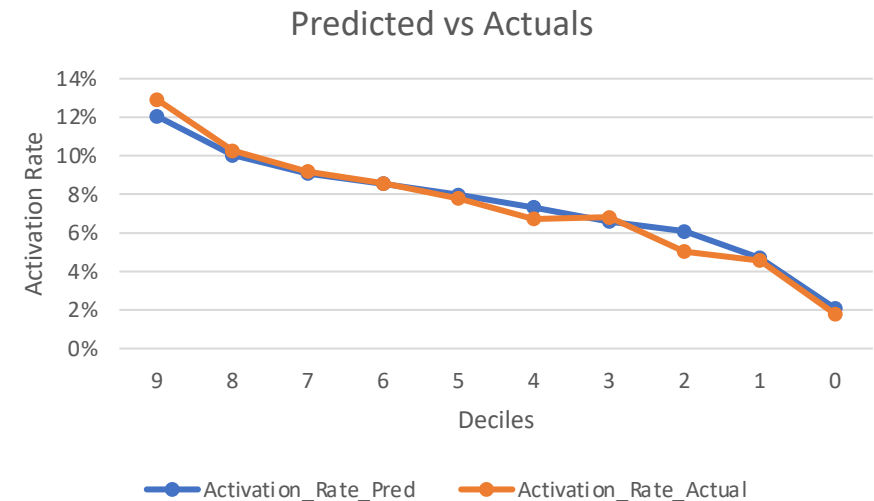
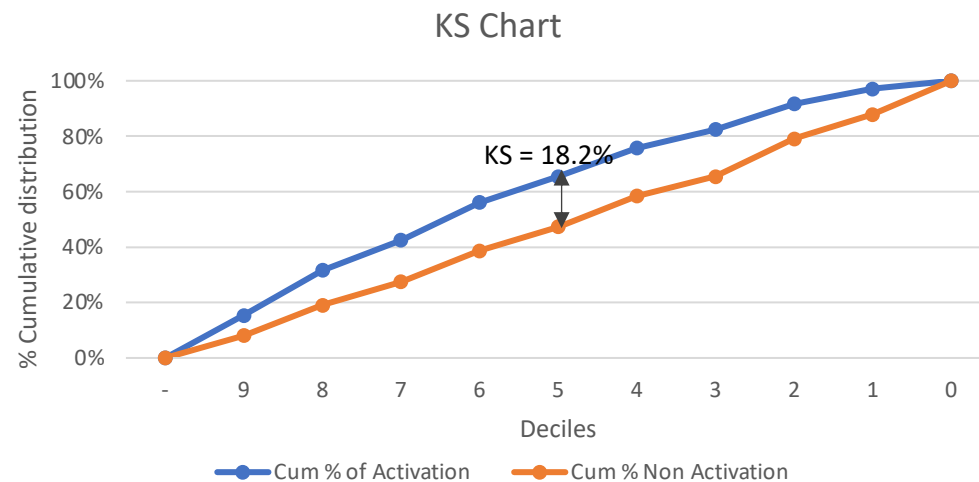
Modeling Evaluation

- As seen in the distribution chart below, as well as the AUC and F1 scores, the model is weak in terms of separating activations from non-activations, however, it rank orders properly as demonstrated by the rank ordering table on the next slide



Rank Ordering Table

group_10	Activation Probability Range	Total	% Cum Dist	Total Predicted Activation	Total Actual Activation	Activation_Rate_Pred	Activation_Rate_Actual	% Cum Activation Rate	% of Activation	% of Non-activation	Cum % of Activation	Cum % Non Activation	Difference
(Best) 9	(0.105, 1]	8082	8.6%	974	1044	12.05%	12.92%	12.92%	15.4%	8.1%	15.4%	8.1%	7.3%
8	(0.0956, 0.105]	10748	20.0%	1078	1103	10.03%	10.26%	11.40%	16.3%	11.0%	31.6%	19.1%	12.5%
7	(0.0866, 0.0956]	8064	28.6%	733	740	9.08%	9.18%	10.73%	10.9%	8.4%	42.5%	27.5%	15.1%
6	(0.0826, 0.0866]	10781	40.0%	923	924	8.56%	8.57%	10.12%	13.6%	11.3%	56.2%	38.7%	17.4%
5	(0.0769, 0.0826]	8151	48.7%	649	635	7.97%	7.79%	9.70%	9.4%	8.6%	65.5%	47.3%	18.2%
4	(0.0694, 0.0769]	10383	59.7%	761	698	7.33%	6.72%	9.15%	10.3%	11.1%	75.8%	58.4%	17.4%
3	(0.0621, 0.0694]	6613	66.7%	436	450	6.60%	6.80%	8.90%	6.6%	7.1%	82.4%	65.5%	16.9%
2	(0.0523, 0.0621]	12516	80.0%	760	631	6.07%	5.04%	8.26%	9.3%	13.6%	91.7%	79.1%	12.6%
1	(0.0331, 0.0523]	8111	88.6%	381	370	4.70%	4.56%	7.90%	5.5%	8.9%	97.2%	87.9%	9.2%
(Worst) 0	(0.0, 0.0331]	10745	100.0%	224	192	2.08%	1.79%	7.21%	2.8%	12.1%	100.0%	100.0%	0.0%
Total		94194		6920	6787	7.35%	7.21%					KS	18.2%



Comparing Algorithms and Sampling Methods

- Given the imbalance between the minority and majority class, the following algorithms and sampling methods were considered to improve performance

Algorithms

- XGBoost
- Random Forest
- LGBM
- CatBoost

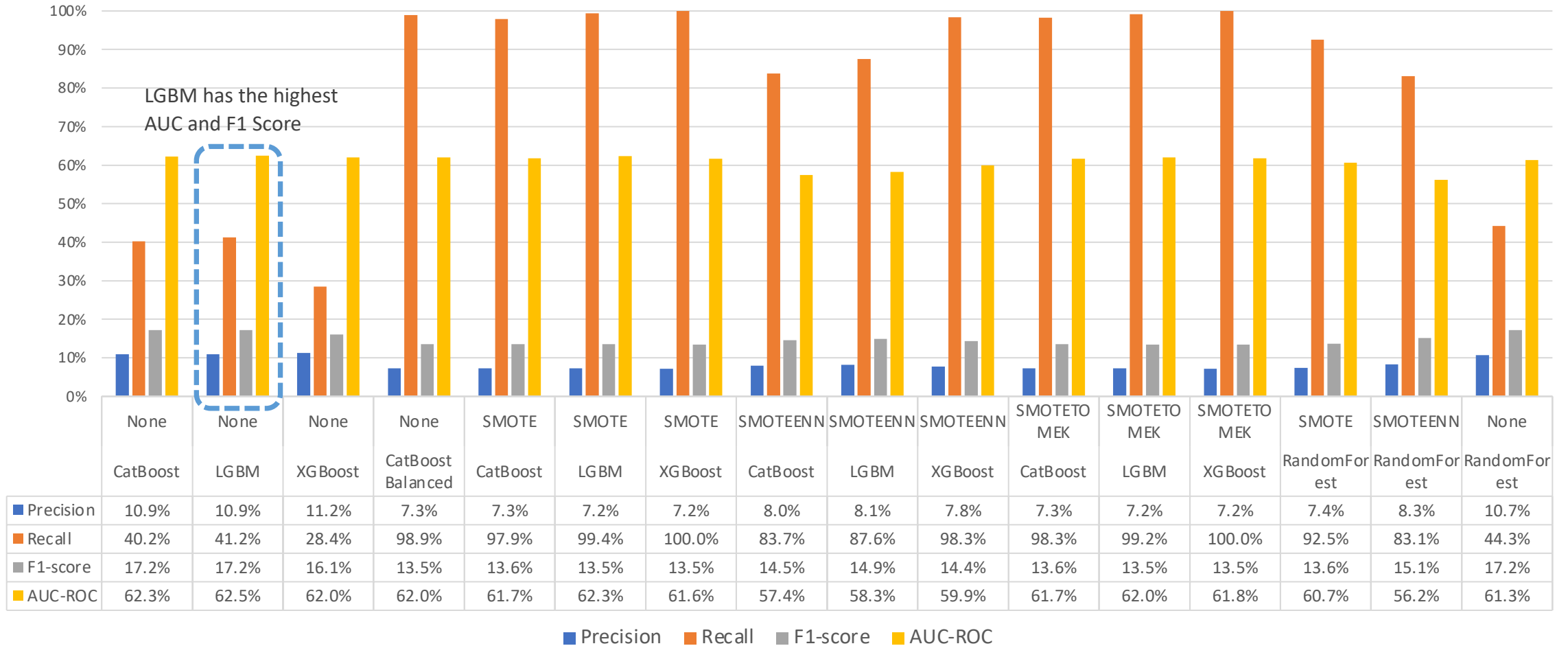
Sampling Methods

- No Sampling
- SMOTE (Oversampling Minority)
- SMOTE ENN (Hybrid)
- SMOTE Tomek (Hybrid)

- However, the performance of the algorithms were very similar in terms of AUC ROC
- The sampling methods also didn't have any significant improvement in the performance
- The chart on the next slide shows the comparison between different algorithms
- A probability threshold of 0.09+ was set for activation

Algorithms and Sampling Methods

Comparison between Algorithms/Sampling methods



Proposed Next Steps To Improve Performance

- Explore additional data that can be obtained systematically
 - Credit bureau data
 - Social media data from third party data aggregators
- Assess trade off between model execution speed/ease of execution and performance
- Side by side with Call Center associates can provide additional insights into the key features they look for in a prospect; Bringing those features into modeling dataset will improve the prediction

Model Performance Tracking Post Implementation

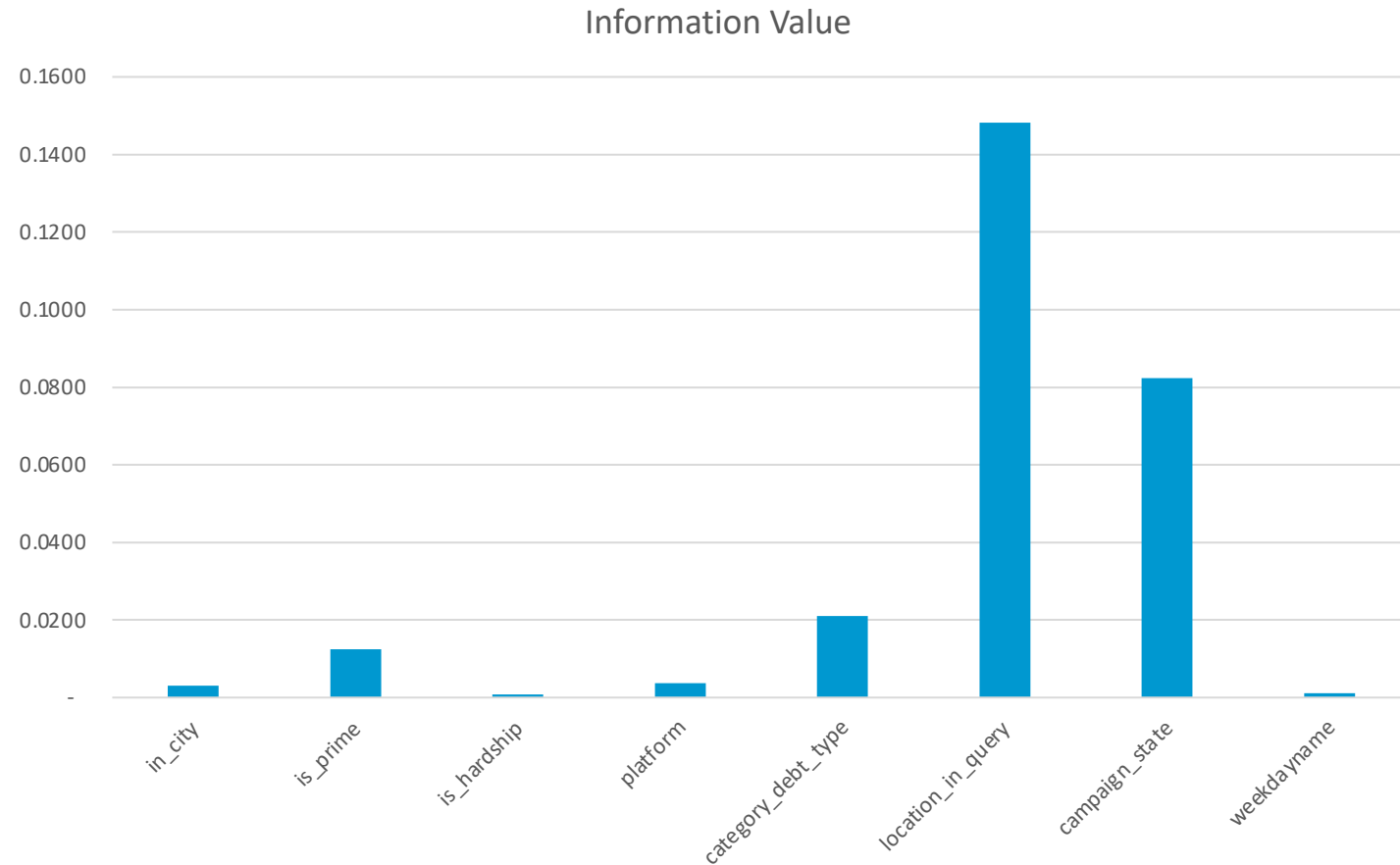
- Model performance monitoring and tracking
 - Periodic cadence (weekly, monthly)
 - Check for population stability (PSI), rank ordering reversals, AUC, F Score
- Define and monitor KPIs – In this case activation rate
 - A model developed with rigor will significantly improve activation rate thus increasing marketing ROI
- Evaluate additional data/features to improve model prediction on an ongoing basis

Appendix

Feature Importance (Partial list)



Information Values



- State and Location have IV indicating that targeting certain geographies can lead to high conversion rate