
Banishing the Bid Bot: Detecting Online Bidding Fraud in Facebook Auctions

Team Random Clique

Daniel Geng
504588536

Department of Computer Science
University of California, Los Angeles
danielgeng@cs.ucla.edu

Meghana Ginpalli
804588573

Department of Computer Science
University of California, Los Angeles
mginpalli@ucla.edu

Sara Melvin
004590481

Department of Computer Science
University of California, Los Angeles
sara.melvin@gmail.edu

Sonu Mishra
604590548

Department of Computer Science
University of California, Los Angeles
mishrasonul993@gmail.com

Andrew Wong
104036702

Department of Computer Science
University of California, Los Angeles
anjuwong@ucla.edu

Abstract

stuff

1 Motivation

Facebook is a popular online social media platform connecting users from all around the world. Because of its high traffic, Facebook is an ideal form of commercial advertisement. Businesses are able to reach out to their existing customers as well as advertise to new customers by specifying which group of users they would like to advertise to. Although Facebook can provide space on both their website and mobile site, there is limited webpage “real estate”. Therefore, Facebook sets up auctions for businesses to bid on designated spaces for advertisement purposes.

Online shopping and bidding are quickly becoming a major proportion of retail transactions. [1] Consumers are increasingly enjoying the convenience and selection provided by the e-commerce platforms available today. Unfortunately, legitimate users are being frustrated by bots as they bid at the last possible second. Automated bids from bots are unfairly able to submit bids almost instantaneously, making it difficult for legitimate buyers to win auctions. According to a paper released by eBay, 65% of sessions each day are initiated by bots. [2]

“Facebook Recruiting IV: Human or Robot?” was a past Kaggle Competition that ended in June 2015. The competition tackles the problem of bot detection in Facebook Ads bidding. Our goal is to detect bot-user anomalies in this e-commerce environment using various anomaly detection methods.

2 Background

Existing approaches have been implemented in order to determine pattern detection for identifying anomalous users. There are many platforms relying on crowdsourcing for promoting products and services. However, there are bots pretending to be humans which lowers the consumers' trust in these sites. Thus, there is a lot of motivation for these companies to have some type of anomaly detector for removing these bots. These approaches are group-based, temporal-based, and edge attribute graph-based.

2.1 Group-based Approach

This approach aims to detect groups of anomalies based on features across multiple users. Previous methods have attempted to classify based on the similarity of feature groups. [3] Other methods detect group based anomalies by determining groups based on closeness of certain features, and then determining which groups are anomalous. [4] The benefits of using these group-based approaches is that it uses the sheer amount of data provided to its advantage. Generally, however, these approaches assume a static model for the data. Additionally, due to the large amount of data processed by these approaches, the methods are computationally expensive.

2.2 Temporal-based Approach

Unlike some of the other approaches, this method uses only one feature to detect bot-users' suspicious activity. The temporal-based approach operates under the assumption that the time between user events for bots is significantly different behavior than normal users. [5] [6] [7] [8] This assumption holds in most situations. However, bots could be compromised users, which camouflages the anomalous temporal feature(s) with their human host. Therefore, other features may need to be used for improved detection of bot-users.

2.3 Edge-attributed Graph-based Approach

Conventional graph-based approaches observe topological patterns of relations between entities (i.e. users, products, etc), but do not leverage the edge-attributes. The algorithm proposed in EdgeCentric utilizes more features by representing the relationship between entities as an edge-attributed multi-graph. [9] Given the distribution of various clusters of benign user behavior, the information required to describe the behavior of a new user is found in terms of its minimum description length (MDL) from the clusters. The MDL is used to compute how suspicious a certain user is. The strength of this approach, however, relies on the rating relationship between users, products, and buyers. In the project's bidding data, a seller and rating information is not disclosed, which potentially weakens the applicability of this approach.

3 Methods

Various anomaly detection algorithms are compared and contrasted, including BIRDNEST, Edge-Centric and classic machine learning algorithms such as K-means, logistic regression, kernel SVM, naive Bayes, and decision trees to discover bot-anomalous behavior. Additionally, these models are compared to the model implemented by the winner of this Kaggle Competition, Small Yellow Duck. Because the dataset from Kaggle (<https://www.kaggle.com/c/facebook-recruiting-iv-human-or-bot/data>) provides labeled data, we are focusing on using supervised algorithms for detection. We plan to use k-fold cross validation and tune our hyperparameters evaluation metrics.

To evaluate the efficacy model, we will apply the trained model to our testing data and compare the classifying algorithms based on performance. The metrics we will use to evaluate the performance of our algorithms is the area under the ROC curve (AUC).

4 Data

The data provided includes timestamps of bids, merchandise categories, bidding devices, countries of bids, and the reference URL for a bid. We plan to use these features and state-of-the art anomaly detection algorithms to determine which method will best identify bot-users.

The bidder dataset contained approximately 2,000 instances and had the following features: bidder_id, payment_account, address, outcome. The bid dataset contained approximately 7.6 million instances and had the following features: bid_id, bidder_id, auction, merchandise, device, time, country, ip, url.

The training data set provided by Kaggle had labels for each user, classifying each individual as a human or a bot. However, the testing data had no labels; therefore, only the training data set could be used for training, validation, and testing.

The training data set was joined with the bids data set, resulting in over 3 million bids. Overall, it contained 1,984 distinct bidders with 103 of the total users classified as bots.

5 Software Tools

Python was used for a majority of our data processing, given the large number of libraries that are publicly available. Sklearn has a number of built-in machine learning algorithms. Numpy is very useful for working with matrices and different feature vectors. Pandas has a high performance data analysis tool used to analyze large amounts of data. Matplotlib facilitates data visualization, which will yield insights into how our data are related. Additionally, Python has APIs for Apache Spark and Hadoop, which may prove useful for large-scale data processing.

6 Results

6.1 Classic Machine Learning Approaches

All machine learning models were implemented for binary classification, rather than a suspiciousness score. Initially, support vector machines (SVMs) and logistic regression were implemented. SVMs were tested with linear, polynomial, and radial kernels and logistic regression was tested using both primal and dual optimization functions.

However, because these models alone were not very effective, a myriad of boosting and ensemble methods were used in an attempt to improve performance. The following methods were tested:

- Adaboost – weights each classifier based on performance
- Gradient boosting – fits trees based on negative gradient of loss function
- Random forests
- Bagging – trains on subsets of data and uses them to vote
- Extra trees – trains on subsets of data and averages them

All of these models fit numerous classifiers—typically decision trees—to make up for the shortcomings of the earlier classifiers.

6.1.1 Features

The given data was pre-processed in various ways. A combination of the following features were tested in each model to varying degrees of performance.

- Log-bucketed IAT (base 2 and base 5) – each bucket is a feature
- Log-bucketed response time (base 2 and base 5) – each bucket is a feature
- Mean IAT
- Mean response time

- Total number of bids
- Mean bids per auction
- Number of devices
- Number of IPs
- Entropy of devices
- Entropy of IPs

6.2 BIRDNEST

BIRDNEST [[5]] is a Bayesian-based approach to anomaly behavior detection. This algorithm identifies anomalous users by implementing two phases. First the algorithm models user behavior distributions as a mixture model (BIRD model) to produce the posterior distributions of user's behavior. It then uses the global distribution of the mixture model to calculate a user's score by taking the negative log likelihood the user's posterior distribution would be generated by the mixture model (NEST metric).

Originally, the algorithm was built to identify rating-fraud by observing two features: rating behavior and inter-arrival times between user activity. In the context of auction bidding, the rating behavior was not used, however the temporal feature aspect of the algorithm seemed similar to the given situation. Therefore, this algorithm was chosen and adjusted to be a temporal-based approach.

6.2.1 Implementation

Features Extracted

1. Inter-arrival time (IAT) distribution - time between user's activity events
2. Response Time - time between the previous bidder's bid and the user

For selecting features, the first hypothesis was that bots will act similarly in a bidding environment as in a rating environment. Originally, the BIRDNEST paper showed an early spike for rating-bot IATs compared to humans. It was theorized that the bidding bot would be activated by the bot-creator once the bot-creator had a paying client and then would bid on as many auctions that fit the client's demographic as possible. This would be a sort of blanket-bidding which would lead to smaller IATs than humans.

For the IAT distribution, the IAT values were computed and binned in a log-based binning histogram with log-base of 2 resulting in 47 buckets. The IAT feature was calculated in a similar manner as the BIRDNEST authors performed in their source code [?]. There were adjustments made in the following: reading in the data and data structures for storing and tracking features.

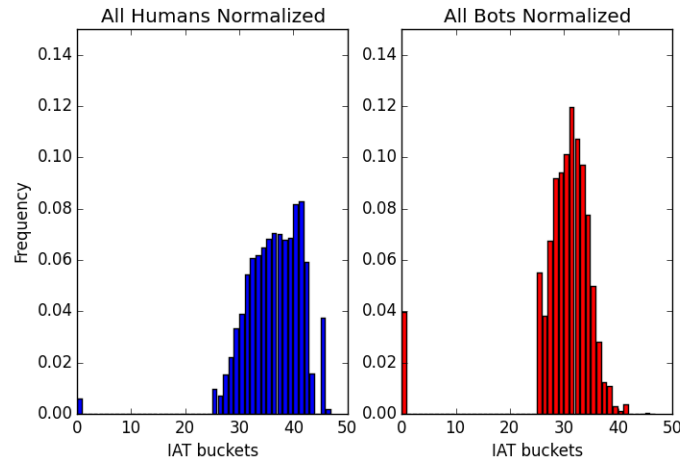
As seen from the figure below, we did not end up seeing the dramatic early spikes as we were hoping but on average the bots had a smaller IAT. We therefore adjusted the log-base to see more or less of a granularity of this distribution. Other log bases were attempted but did not dramatically alter the distributions. Therefore, we chose a log-base of 2 to give us the most detailed histogram in the most reasonable processing time.

After analyzing this feature we had another theory. Bot-creators activate their bots to be watchdogs on certain auctions paid for by their customers to out-bid any competitive bid. This response time would then be smaller than a normal human's reaction time to bids.

Upon searching for Kaggle winner blogs, the team came across the second place winner called the Small Yellow Duck (will be referred to later in this paper) who used, amongst other features, the median and mean response time of every bid as a feature in her solution. To reflect this feature for better comparison to the professional's approach, the Small Yellow Duck's response times were extracted and joined with a table consistent with the BIRDNEST format. The response time distribution was binned with the same log-base as IAT resulting which resulted in 44 buckets.

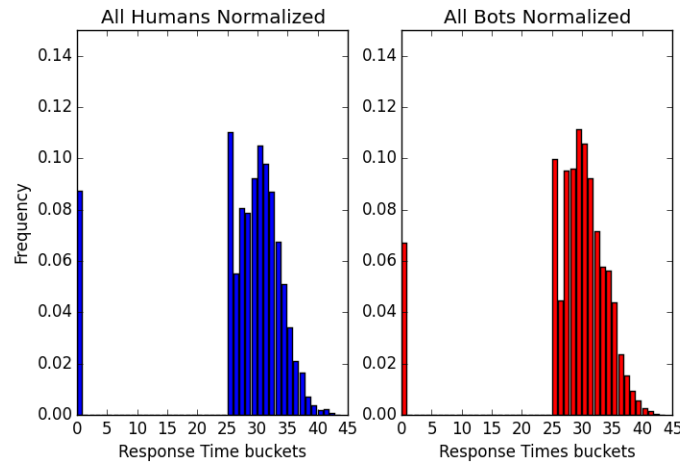
As the figure below will show, there is not a distinct distribution difference between these groups as originally hoped. This is because human users do not necessarily represent one person. Many

Figure 1: This figure shows the normalized true human IAT distribution (left) and true bot normalized IAT distribution (right)



human users in this data set have multiple people bidding with the same username on the same auctions. These multi-party users show very similar bot-like response times.

Figure 2: This figure shows the normalized true human response time distribution (left) and true bot normalized response time distribution (right)



6.2.2 Testing Methodology

To test which features contribute the most, this algorithm was tested with just IAT as the only feature and then with both IAT and response time as features.

The BIRD mixture model is calculated by taking a random subset of each cluster's user distributions to estimate hyperparameters until convergence. Therefore, BIRDNESST results had a slightly varying AUC each time the model was run. Consequently, the algorithm was run with different number of clusters parameter, K - one cluster to twenty clusters - to gain a true sense of the performance. Each user's score was averaged over these twenty runs and then re-ranked from highest to lowest. The true positive rate and false positive rate were calculated from this re-ranked list to produce the ROC curve and ultimately the AUC.

6.2.3 Assumptions

For the testing, the following assumptions were made:

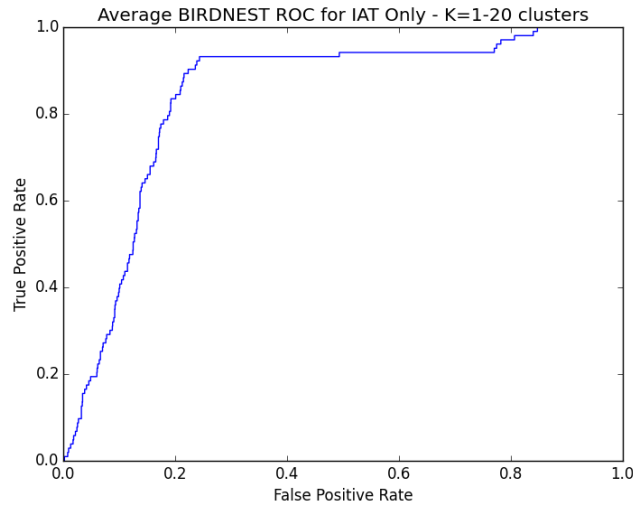
- All users that did not have more than one bid in the data set were given zero IAT distributions.
- All bids that initiated the auction and all bids that were made by the same bidder at the same auction at the same time were not binned in the response time buckets. At first these were considered zero, but upon testing it was found that ignoring them would lead to similar results.
- All instances where a user bid more than once at the same time on the same auction with the same username were considered duplicates and only one of these instances was kept.

6.2.4 Evaluation

Results: IAT only

Despite not seeing the “early” spike as expected in the bot IAT distribution, using IAT as the only feature produced decent results. The figure below shows the ROC with an AUC of only 0.8415.

Figure 3: BIRDNEST ROC curve of the average user scores over twenty different runs from one to twenty clusters, K, with only IAT as the feature.



As you seen from the figure below, the algorithm detects similar anomalous distributions as the true IAT-bot distribution and the true IAT-human distribution. This result shows that BIRDNEST mixture model detects distributions that are anomalous to the rest of the normal distributions. However, the algorithm's performance decreases because our data consists of this “noise” in the form of these multi-party users that exhibit bot-like behavior.

Results: IAT and Response Time

Using Response Time and IAT distributions as features, the performance was only marginally improved. The figure below shows the ROC with an AUC of 0.8456, which is only a 0.0041 improvement.

As mentioned previously, these multi-party users in the data set that exhibit bot-like behavior are a sort of noise in our data set. Therefore, we conclude that BIRDNEST is sensitive to noise and would work better in an environment where users are limited to have one person logged in at a time.

Figure 4: This figure shows the BIRDNEST identified human IAT distribution (left) and BIRDNEST detected bot IAT distribution (right)

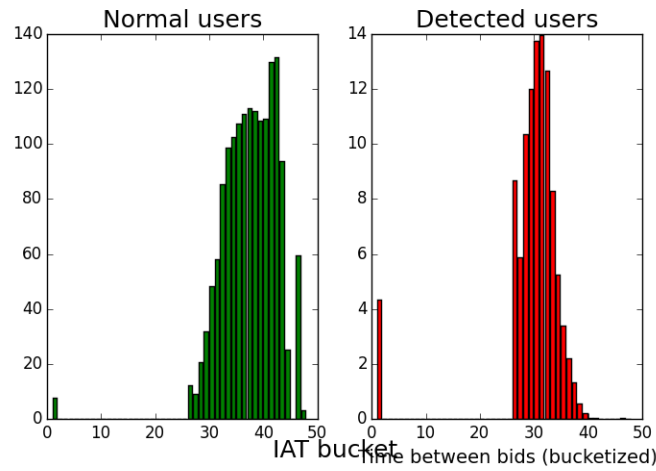
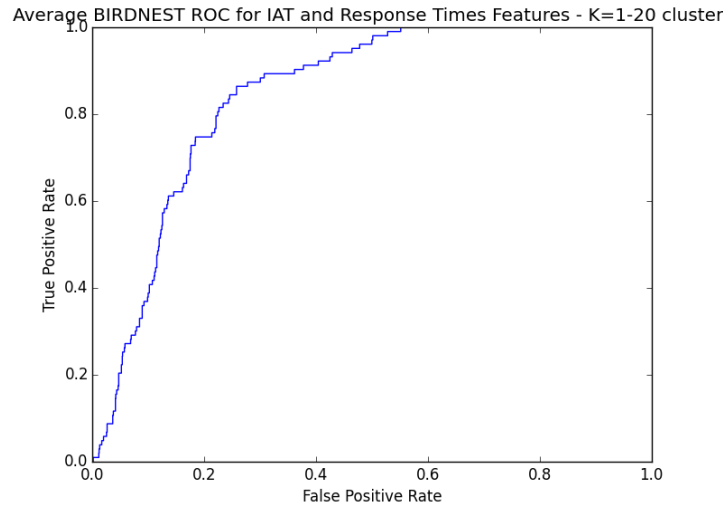


Figure 5: BIRDNEST ROC curve of the average user scores over twenty different runs from one to twenty clusters, K, for response time and IAT features.



6.3 Kaggle Winner: “Small Yellow Duck”

The second place winner of this Kaggle Competition is a professional data scientist and goes by the name “Small Yellow Duck”. Her classification model takes the average of the probabilities predicted by five instances of the Random Forest Classifier, which is an ensemble learning method that constructs multiple decision trees to create a stronger classification model.

Small Yellow Duck first makes a couple of observations before implementing feature extraction. First, human bidding activity tends to peak daily as seen in the figure below. Over the course of three days, there are three peaks where human bidding activity is higher. This may be because of auctions ending at the same time everyday.

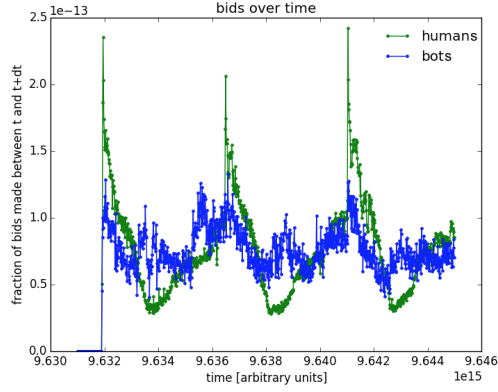
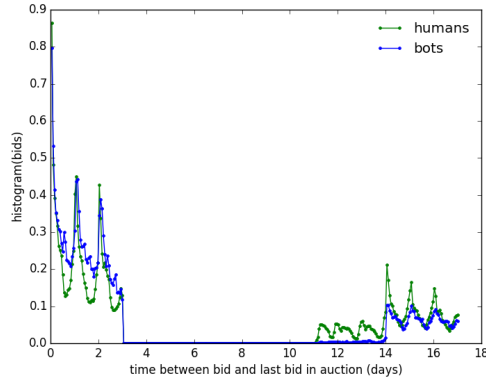


Figure 6: Daily peaks of human bidding activity

Another observation that she made is that auctions tend to last for longer than two weeks. In this timeframe, robots do not place any bids between 11 to 14 days before the auction ends. In the figure below, the auction runs for approximately 17 days. For some reason, between days 11 and 14, human bidding commences while bots tend to keep silent in this duration.



During feature extraction, Small Yellow Duck was able to extract the following features: the median time between user's bid and user's previous bid, the mean number of bids a user made per auction, entropy for how many bids a user placed on each day of the week, max number of bids in a 20 min span, total number of bids placed by user, average number of bids a user placed per URL, number of bids placed by the user on each of three weekdays in the data, and min and median times between a user's bid and previous bid by another user in the same auction.

In order to evaluate her results, Small Yellow Duck does cross validation with 100 different training and validation splits where 80% was training data and 20% was validation data. In terms of runtime, training and predicting takes around 3 minutes while 100-fold cross validation takes around 20 minutes. As a result, the AUC score tends to average at 0.94167. Based on her AUC score, our highest performing algorithm would allow us to be in the top 15% of the leaderboard.

7 Discussion

In conclusion, we started this project trying to compare and contrast different algorithms and features to see which one would perform better. In the end, our results were actually comparable to the Kaggle winner's approach. The optimal algorithm turned out to be Gradient Boost and other boosted machine learning approaches.

From this project, we gained a few takeaways and future direction. First of all, manipulating a lot of data can be difficult but using the right tools, finding the right features, organizing or sorting it before hand, or storing them as binaries make a huge difference. Secondly, temporal features are not nearly as powerful as we originally thought. Once we added entropy terms and different metrics for the number of resources each user was using to bid, we did see dramatic improvement. Although our results are comparable, we could improve these scores by extracting more interesting features that would lead to an improvement of accuracies. Furthermore, as 'state-of-the-art' as certain algorithms are, we should never underestimate the power of numbers. Combining multiple estimators through different boosting methods allows an estimator to fill in the gaps where another estimator might fail, making for a much more powerful classification algorithm. Although boosting the algorithms greatly improves performance, we would have liked to see if there was a singular algorithm that is just as powerful on its own without taking too much training time. Also, we did attempt a group-based and graph-based approach to apply to our data set but the scalability proved to be a huge issue and made it difficult to debug in the time we needed for this project and was not scalable in preprocessing data into a graph. In conclusion, we learned quite a lot of interesting ways to manipulate large amounts of data and get important features out of this. In addition, we were able to determine what types of algorithms performed better than others.

8 Contributions

Daniel: Most of report

Meghana: Presentation, Kaggle winner code/writeup, References

Sara: BIRDNEST, Presentation, References

Sonu: Presentation

Andrew: Machine learning code, Presentation

References

- [1] "15 Mind-Blowing Stats About Online Shopping." 15 Mind-Blowing Stats About Online Shopping.
- [2] Zhao, K. and Pengju, Y. "Machine Learning Method In eBay Bot Detection." Ebay Tech Blog. September 4 2014.
- [3] Giatsoglou, M. et. al. "ND-Sync: Detecting Synchronized Fraud Activities." Advances in Knowledge Discovery and Data Mining, 2015, 201-214t
- [4] Qi, Y., Xinran He, and Yan Liu. "GLAD: Group Anomaly Detection in Social Media Analysis- Extended Abstract." arXiv:1410.1940.
- [5] Hooi, B., et. al. "BIRDNEST: Bayesian Inference for Ratings-Fraud Detection." arXiv:1511.06030
- [6] Costa, A.F. et. al. "RSC: Mining and Modeling Temporal Activity in Social Media." Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2015, pp. 269-278
- [7] Vaz de Melo, P. et. al. "The Self-Feeding Process: A Unifying Model for Communication Dynamics in the Web." Proceedings of the International Conference on World Wide Web (WWW), pp. 1319-1330, 2013.
- [8] Malmgren, M.D. et. al. "Characterizing Individual Communication Patterns." Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2009, pp. 607-616.
- [9] Shah, N. et. al. "EdgeCentric: Anomaly Detection in Edge-Attributed Networks." arXiv:1510.05544.