

# Troll Project

*Matt Cole*

*October 23, 2016*

## Introduction

Internet trolls are considered a menace in nearly all online communities as they strive to generate emotional responses from individuals. Motivations behind these users are unclear, but some social scientists suspect that a strange phenomenon known as the “disinhibition effect” may be to blame. Masked behind the apparent anonymity provided by many forms, sites, and the internet in general, social reservations that facilitate normal, face-to-face conversations can disappear, resulting in sometimes wild and rude behavior. This ‘troll behavior’ can stretch from simply posting the same status repeatedly to annoy and clog streams of information to violent threats and many, many ‘things’ in between. The effect of these trolls on Twitter, in particular, has been linked to many popular users leaving the platform entirely. The damage trolls may cause moves further than individuals, with Twitter’s troll problem partially responsible its recent loss of 35 billion in market capitalization according to some analysts (Hern, 2016). In this project, we focused on identifying political trolls on Twitter, an area considered to be particularly saturated with such users.

## Data

Usernamees were collected from Twitter using both the Twitter API and web app, and was limited to English tweets. Initial data collection consisted of searching twitter, via `twitteR` API wrapper for users with tweets containing one of the two major political party presidential candidate’s twitter handle within the tweet text (@hillaryclinton or @realdonaldtrump). This collection occurred on October 1st, 2016 at 4:14:19 PM EST. Twitter API mechanics resulted in a mix of most recent tweets (up to the second), as well as ‘hot’ tweets that were somewhat recent and popular in terms of favorites and retweets. In total, 142 usernames were collected through this method. Later, from October 7th through October 14th, because of an apparent lack of trolls in the API generated users, manual username harvesting was conducted. This manual harvest consisted of browsing political figures tweets and responses, in an attempt to find collect users with a higher proportion of trolls for modeling purposes. During this collection period, 132 additional usernames were added for a total of 274. Once a list of usernames was collected, each user’s 10 most recent tweets, as well as relevant tweet data (favorites, location, retweets, etc.) and data associated with the account (friends, description, etc.), were fetched.

Using only the text from each user’s 10 most recent tweets, including retweets and replies, users were classified as a troll or non-troll.

Trolls were defined as users with at least one of the previous 10 tweets appearing to be both:

- inflammatory: intended to arouse angry or violent feelings
- extraneous: irrelevant or unrelated to the subject being dealt with

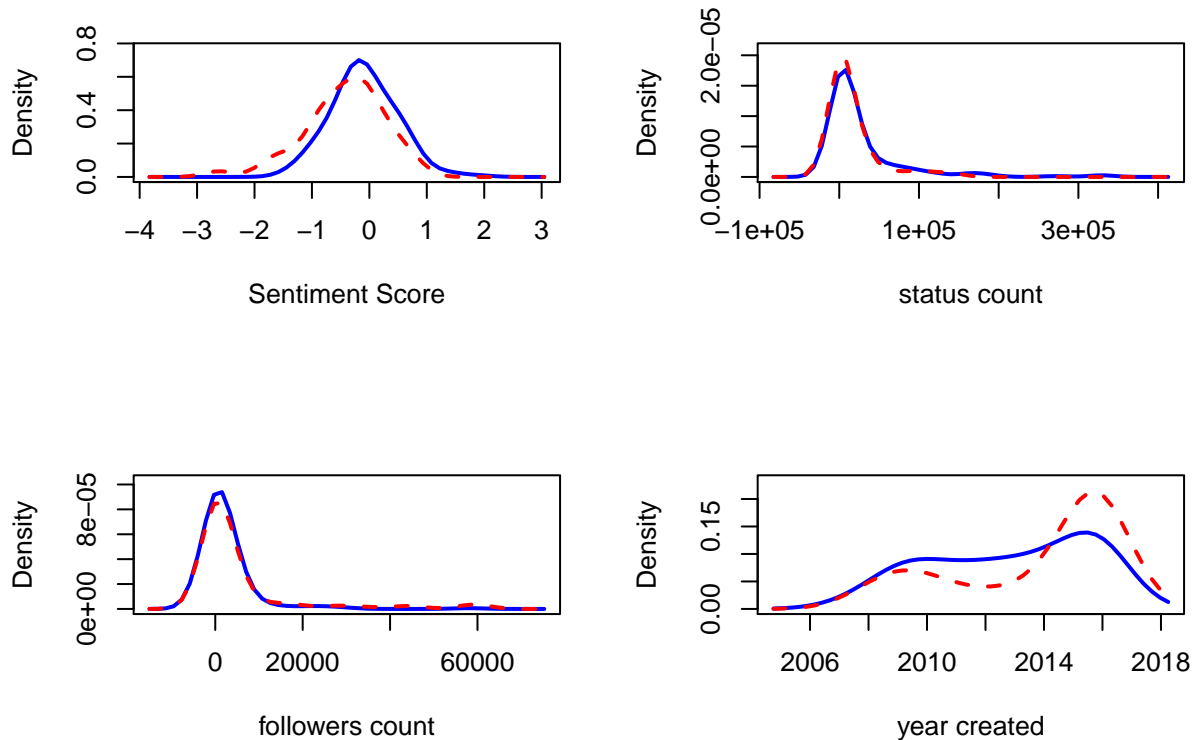
Ten tweets were chosen in an attempt to capture enough of each user’s tweeting tendencies and classifying users as trolls with at least one inflammatory and extraneous tweet allows us to capture more subtle trolls. In total, all 274 users were classified as troll/non-troll after examining roughly 2,740 individual tweets.

## Methods

Using each users previous 10 tweets, we were able to construct average tweet sentiment using the AFINN lexicon. This lexicon comprises of English words with integer values between minus five (negative) and plus

five (positive). Sentiment scores for each tweet were calculated and averaged together. Additionally, the NRC lexicon, which assigns scores to each of two sentiments (positive and negative) as well as eight emotions was used to generate an ‘angry’ score comprising of anger, fear, and disgust subscores as well as a ‘happy’ score comprising of joy and trust emotion subscores. As emojis have become increasingly popular in today’s web-based communication, an emoji dictionary was scrapped from <http://www.unicode.org> and emoji usage as well as sentiments were determined and averaged across tweets. Tweet source (how tweets were sent) was collected and divided into three categories, browser-based tweets, mobile tweets, and other tweet systems which included automated (bots) and unknown tweet sources.

Exploratory data analysis revealed that no single factor showed good separation, between trolls and non trolls (figure 1). Several covariates expected to have a strong relationship with troll status such as number of followers did not appear to show such in the 2 dimensional plots (figure 1). However, others such as AFFIN sentiment score did (figure 1).



**Figure 1** Smoothed density plots of common twitter metrics in both troll and non-troll groups. Trolls are denoted by the red line while non-trolls are denoted by the blue line. Distribution of AFFIN sentiment score and year of account creation appear to potentially be different in the two groups while there is no noticeable difference in follower or status count.

We focused on two models to predict troll status using logistic regression and classification trees. Classification and regression trees (CART) are a group of supervised machine learning methods which construct decision tree models utilizing observed data. Unlike other popular methods for classification or regression, regression trees are not global models, where a single predictive formula is employed to make decisions or predictions over all observations (James et al., 2013). Instead, CART methods break down the feature space containing the observations into small subspaces (called recursive partitioning) until the subspaces are able to be represented by relatively simple models that do not necessarily utilize all covariates presented in feature space. Benefits of CART methods include its high level of interpretability and simple models, automatic variable selection, and ease of visualization. The classification tree was grown using many features and pruned to minimize the cv-10 error rate. A major issue of classification trees is the potential for overfitting, with the cv-pruning to eliminate this issue. Logistic regression models the probability of a user being a troll for each observation  $i$  ( $p_i$ ) by treating the logit function,  $\log(\frac{p_i}{1-p_i})$  as a linear function of the covariates. In this study, a major limitation may be logistic regression’s inability to capture non-linear trends associated with the covariates.



## Discussion

### results

It was expected that the regression tree would be able to outperform the logistic regression as the relationship between covariates and troll status seem opaque at best and would likely be best modeled with a non-parametric method. As we saw, the classification tree was infact able to outperform the logistic regression models with respect to 10-fold cross-validation prediction error rate and recorded a sensitivity of 0.86 and a specificity of 0.78.

This study did shed light on a strange phenomena of different classes of trolls. From our anecdotal experience, there seemed to be three types or classes of twitter trolls.

- ‘personal accounts’ comprising of students, professionals, etc. whom occasionally say mean/emotionally charged things on twitter.
- bots: usually newly created, low tweet/follower/friend count which seem to spew politically charged propaganda.
- troll/fan pages: these accounts seem to be semi-autonomous as they respond to policial figures (@therealdonaldtrmp, @hillaryclinton) almost instantaneously while also being able to respond to certain accounts in a thought out manner. These accounts have many tweets, were created long ago and have many followers. In fact, Fortune, POLITICO and the New York Magazine all have run stories of these ‘bot armies’.

### limitations

The main limitation of this study was a lack of available data, particularly trolls. Although there were 274 observations there were only 92 trolls in our dataset, more would likely allow us to identify additional relationships and interactions between variables and troll status. In addition, there is the possibility for human error as a single person classified all tweets. This error could come from keystroke error or implicit political bias. Political bias could be difficult to identify and overcome where beliefs of what actually ocured could alter whether a statement is either inflamitory or extraneous. A good example of this are statements such as “Bill Clinton is a rapist”, which were classified as both inflamitory or extraneous, but to an individual whom believes that the statement is true, would likely not view it as extraneous as it is simply stating the truth.

### Generalization

It is plausible that these models could generalize to other spaces within twitter or even outside of the social media site where trolling is a serious issue. However, the model relies heavily on non-text covariates such as friends, dates, and followers and applying this model to, say an anonymous online form would likely be less effective.

### future work

In the future, it would be beneficial to further optimize tuning parameters associated with these models and inspect how other algorithms perform, particularly convolutional neural networks as well as support vector machines which may be able to identify aditional hard to see relationships in the data resulting in better prediction accuracy. In addition, future work should include identifying word n-grams which may allow us to gain additional insight beyond lexicons. n-grams are popular in computational linguistics as two words in sequence contain additional information about the meaning of the sentence compared to lexicon analysis of the same two words. Future would could focus on classifny individual tweets as troll / non-troll tweets

or even non-binary degrees of ‘troll’, for instance, classifying three types of trolls mentioned earlier in the discussion.

## Thank you

Thanks to Finn Årup Nielsen manually labeled the words in the AFINN lexicon from 2009-2011, which was used in this analysis.

Thank you to <http://unicode.org/emoji/charts/full-emoji-list.html> for letting me borrow the emoji database without asking.

## References:

Hern, Alex. “Did Trolls Cost Twitter \$3.5bn and Its Sale?” The Guardian. Guardian News and Media, 18 Oct. 2016. Web. 21 Oct. 2016.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning (Vol. 6). New York: springer.

Martin, J. H., & Jurafsky, D. (2000). Speech and language processing. International Edition, 710.

## Packages

Hadley Wickham and Romain Francois (2016). dplyr: A Grammar of Data Manipulation. R package version 0.5.0. <https://CRAN.R-project.org/package=dplyr>

Terry Therneau, Beth Atkinson and Brian Ripley (2015). rpart: Recursive Partitioning and Regression Trees. R package version 4.1-10. <https://CRAN.R-project.org/package=rpart>

Hadley Wickham, Jim Hester and Romain Francois (2016). readr: Read Tabular Data. R package version 1.0.0. <https://CRAN.R-project.org/package=readr>

Bowman, A. W. and Azzalini, A. (2014). R package ‘sm’: nonparametric smoothing methods (version 2.2-5.4) URL <http://www.stats.gla.ac.uk/~adrian/sm>, [http://azzalini.stat.unipd.it/Book\\_sm](http://azzalini.stat.unipd.it/Book_sm)

Stephen Milborrow (2016). rpart.plot: Plot ‘rpart’ Models: An Enhanced Version of ‘plot.rpart’. R package version 2.1.0. <https://CRAN.R-project.org/package=rpart.plot>

Jeff Gentry (2015). twitterR: R Based Twitter Client. R package version 1.1.9. <https://CRAN.R-project.org/package=twitter>