I Wish I Didn't Say That! Analyzing and Predicting Deleted Messages in Twitter

Saša Petrović

Miles Osborne

Victor Lavrenko

School of Informatics University of Edinburgh

School of Informatics University of Edinburgh School of Informatics University of Edinburgh

sasa.petrovic@ed.ac.uk miles@inf.ed.ac.uk vlavrenk@inf.ed.ac.uk

Abstract

Twitter has become a major source of data for social media researchers. One important aspect of Twitter not previously considered are *deletions* – removal of tweets from the stream. Deletions can be due to a multitude of reasons such as privacy concerns, rashness or attempts to undo public statements. We show how deletions can be automatically predicted ahead of time and analyse which tweets are likely to be deleted and how.

1 Introduction

In recent years, research on Twitter has attracted a lot of interest, primarily due to its open API that enables easy collection of data. The belief that tweets contain useful information has lead to them being used to predict many real-world quantities. For example, tweets have been used to predict elections (Tumasjan et al., 2010; O'Connor et al., 2010), stock market movement (Bollen et al., 2011), and even flu outbreaks (Ritterman et al., 2009). Twitter forbids distribution of raw tweets and their terms of service insist that any tweet collection must honor post-hoc deletion requests. That is, at any point in the future a user can issue a request to Twitter to delete a tweet. Predicting when a tweet is likely to be retracted by a user has important applications:

- Security. Twitter has become so ubiquitous that users often do not consider the potential confidentiality implications before they tweet.
- *Regret*. Users might post an inappropriate or offensive tweet in the heat of the moment, only to regret it later.
- *Public scrutiny*. High profile politicians at times tweet content that they later withdraw.

Here we report on the first results of automatically predicting if tweets will be deleted in the future. We also analyse why tweets are deleted.

2 Related Work

Predicting deleted messages has been previously addressed in the context of

emails (Dabbish et al., 2003; Dabbish et al., 2005). For example, (Dabbish et al., 2003) found that the most important factors affecting the chances of an email being deleted are the past communication between the two parties and the number of recipients of the email. However, it should be clear that people use tweets in very different ways to using email. The most similar work to ours is the recent analysis of censorship in Chinese social media (Bamman et al., 2012). The problem examined there is that of the government deleting posts in the Chinese social media site Sina Weibo (Chinese equivalent of Twitter). The authors analyze different terms that are indicative of a tweet being deleted and the difference between appearance of certain political terms on Twitter and on Sina Weibo. However, they make no attempt to predict what will be deleted and only briefly touch upon deleted messages in Twitter. While the main reason for deletion in Sina Weibo seems to be government censorship, there is no known censorship on Twitter, and thus the reasons for deletion will be quite different. To the best of our knowledge, we present the first analysis of deleted messages on Twitter.

3 Task Description

There are several ways in which a tweet can be deleted. The most obvious way is when its author explicitly deletes it (this is usually done by clicking on a Delete button available in most Twitter clients). Another way that a tweet becomes effectively deleted is when a user decides to make his tweets protected. Although the user's tweets are still available to read for his friends, no one else has access to them any more (unless the user decides to make them public again). Finally, the user's whole account might be deleted (either by their own choice or by Twitter), meaning that all of his tweets are also deleted. In the public streaming API, Twitter does not differentiate between these different scenarios, so we collapse them all into a single task: for each tweet predict if it will be deleted, by either of the aforementioned ways.

¹These results were also confirmed in (Tschang, 2012).

3.1 Example Deleted Tweets

Table 1 shows some examples of the various types of deleted tweets that we have discussed (identifiable information has been replaced by ***). Although we can never be sure of the true reason behind someone deleting a tweet, a lot of the time the reason is fairly obvious. For example, it is very likely that tweet 1 was deleted because the author regretted posting it due to its somewhat inappropriate content. On the other hand, tweet 2 was most likely posted by a spammer and got deleted when the author's account was deleted. Tweet 3 is probably an example of deleting a tweet out of privacy concerns - the author posted his email publicly which makes him an easy target for spammers. The fourth tweet is an example of a deleted tweet authored by a Canadian politician (obtained from the website politwitter.ca/page/deleted). Finally, tweet 5 is an example of a false rumour on Twitter. This tweet was retweeted many times right after it was posted, but once it became clear that the news was not true, many users deleted their retweets.

4 Predicting when Tweets will be Deleted

We now show the extent to which tweet deletion can be automatically predicted.

4.1 Data

We use tweets collected from Twitter's streaming API during January 2012. This data consists of 75 million tweets, split into a training set of 68 million tweets and a test set of about 7.5 million more recent tweets (corresponding roughly to tweets written during the last three days of January 2012). A tweet is given the label 1, meaning it was deleted, if the notice about its deletion appeared in the streaming API at any time up to 29th February 2012. Otherwise we consider that the tweet was not deleted. In total, 2.4 million tweets in our dataset were deleted before the end of February.

4.2 Features

We use the following features for this task:

- Social features: user's number of friends, followers, statuses (total number of tweets written by a user), number of lists that include the user, is the user verified, is the tweet a retweet, is the tweet a reply. Additionally, we include the number of hashtags, mentions, and links in the tweet under social features, even though they are not strictly "social". We do this because these features are dense, and thus much more similar to other dense features (the "real" social features) than to sparse features like the author and text features.
- Author features: user IDs,
- *Text* features: all the words in the tweet.

Because of the user IDs and lexical features, the feature set we use is fairly large. In total, we have over 47 million features, where 18 million features are user IDs, and the rest are lexical features (social features account for only about a dozen of features). We do not use features like user's time zone or the hour when the tweet was written. This is because our preliminary experiments showed that these features did not have any effect on prediction performance, most likely because the author and text features that we use already account for these features (e.g., authors in different time zones will use different words, or tweets written late at night will contain different words from those written in the morning).

4.3 Learning Algorithm

In all our experiments we use a support vector machine (SVM) (Cortes and Vapnik, 1995) implemented in Liblinear (Fan et al., 2008). We note that while SVMs are generally found to be very effective for a wide range of problems, they are not well suited to large-scale streaming problems. A potential limitation is the fact that they require batch training, which can be prohibitive both in terms of space and time when dealing with large datasets. Because of this, we also explored the use of the passive-aggressive (PA) algorithm (Crammer et al., 2006), which is an efficient, online, max-margin method for training a linear classifier. Thus, we also present results for PA as an alternative for cases where the data is simply too big for an SVM to be trained.

4.4 Results

We formulate predicting deletions as a binary classification task – each tweet is assigned a label 0 (will not be deleted) or 1 (will be deleted). Because the two classes are not equally important, i.e., we are normally more interested in correctly predicting when something will be deleted than correctly predicting when something will not be deleted, we use the F_1 score to measure performance. F_1 score is standard, e.g., in information retrieval, where one class (relevant documents) is more important than the other.

Results are shown in Table 2. The random baseline randomly assigns one of the two labels to every tweet, while the majority baseline always assigns label 1 (will be deleted) to every tweet. We can see from the absolute numbers that this is a hard task, with the best F_1 score of only 27.0. This is not very surprising given that there are many different reasons why a tweet might be deleted. Additionally, we should keep in mind that we work on all of the crawled data, which contains tweets in nearly all major languages, making the problem even harder (we are trying to predict whether a tweet written in any language will be deleted). Still, we can see that the machine learning approach beats the baselines by a very large margin (this difference is statistically significant at p=0.01). Further improving perfor-

- 1 Another weekend without seeing my daughters-now if I'd shot my ex when we split I would of been out by now, missed opportunity:
- 2 Get more followers my best friends? I will follow you back if you follow me http://***
- 3 @*** yeah man email the contract to ***@gmail.com ... This has been dragged out too long big homie
- 4 Gov must enforce the Air Canada Act and save over 2,500 jobs. @*** http://*** #ndpldr
- 5 BREAKING: URGENT: News spreading like wildfire, BASHAR AL-ASSAD HAS ESCAPED #SYRIA! We're waiting for a confirmation

Table 1: Examples of tweets that have been deleted.

	F_1
Random baseline Majority baseline All features (SVM)	5.8 6.0 27.0
All features (PA)	22.8
Social features Lexical features User IDs	3.8 10.9 12.2

Table 2: Results for predicting deleted tweets.

mance in this task will be the focus of future work and this should enable researchers to distribute more stable Twitter datasets.

We mentioned before that using an SVM might be prohibitive when dealing with very large datasets. We therefore compared it to the PA algorithm and found that PA achieves an F_1 score of 22.8, which is 4.2 points lower than the SVM (this difference is significant at p=0.01) However, the SVM's gain in performance might be offset by its additional computational cost – PA took 3 minutes to converge, compared to SVM's 8 hours, and its memory footprint was two orders of magnitude smaller. Because efficiency is not our primary concern here, in the rest of the paper we will only present results obtained using SVM, but we note that the results for PA showed very similar patterns.

To get more insight into the task, we look at how different feature types affect performance. We can see from the last three rows of Table 2 that social features alone achieve very poor performance. This is in contrast to other tasks on Twitter, where social features are usually found to be very helpful (e.g., (Petrović et al., 2011) report F_1 score of 39.6 for retweet prediction using only social features). Lexical features alone achieved reasonable performance, and the best performance was achieved using user ID features. This suggests that some users delete their tweets very frequently and some users almost never delete their tweets, and knowing this alone is very helpful. Overall, it is clear that there is benefit in using all three types of features, as the final performance is much higher than performance using any single feature group.

We performed ablation experiments where we re-

moved social features from the full set of features one at a time and measured the change in performance. We found that the only two features that had an impact greater than 0.1 in F_1 were the number of tweets that the user has posted so far (removing this feature decreased F_1 by 0.2), and is the tweet a retweet (removing this feature decreased F_1 by 0.16). This is interesting because the number of statuses is usually not found to be helpful for other prediction tasks on Twitter, while the followers number is usually a very strong feature, and removing it here only decreased F_1 by 0.07.

The number of followers a user has is often considered one of the measures of her popularity. While it is certainly not the only one or the "best" one (Cha et al., 2010), it is still fairly indicative of the user's popularity/influence and much easier to collect than other ones (e.g., number of mentions). In the next experiment, we are interested in seeing how well our system predicts what popular users (those with at least a certain number of followers) will delete. In addition, we look at how well our system works for verified users (celebrities). Arguably, predicting whether a celebrity or a user with 10,000 followers will delete a tweet is a much more interesting task than predicting if a user with 3 followers will do so. To do this, we run experiments where we only train and test on those users with the number of followers in a certain range, or only on those users that are verified. We can see from Table 3 that the situation between groups is very different. While for users with less than 1,000 followers the performance goes down, our system does much better on users that have lots of followers (it is also interesting to note that the baseline is much higher for users with more followers, which means that they are more likely to delete tweets in the first place). In fact, for users with more than 10,000 followers our system achieves very good performance that it could actually be applied in a real scenario. For celebrities, results are somewhat lower, but still much higher than for the whole training set.

5 Why are tweets deleted?

One of the fundamental questions concerning deleted tweets is why are they deleted in the first place. Is it the case that most of the deletion notices that we see in the stream are there because users deleted their accounts? Or is it the case that most of the deleted tweets

User group	Here	Baseline	# in test set	Deletion type	% of tweets in test set	Accuracy
Followers < 1,000	17.8	5.8	6.8M	Manual deletion	85.2	18.8
Followers $\in [1k, 10k]$	33.7	6.6	640k	Protected	12.2	17.5
Followers $\in [10k, 100k]$	66.0	17.7	50k	Account deleted	2.6	29.5
Followers $> 100,000$	86.4	41.5	5.5k			
Celebrities	39.5	6.0	3.5k	Table 4: Proportion	n of different types of de	letions and

performance of our algorithm across these types. Table 3: F_1 score for different groups of users. The third column shows our results for named groups. The

last column shows the number of users in the test set that fall into each category.

come from active Twitter users who change their mind about posting a tweet (for one reason or another)? Or are tweets deleted by Twitter because they are sent by spammers? Here we try to answer these questions by looking at profiles of users who deleted tweets.

We take the 200000 deleted tweets from the test set and query Twitter's API to retrieve the account status of their author. There are three possible outcomes: the account still exists, the account exists but it is protected, or the account does not exist any more. Deleted tweets from the first type of user are tweets that users manually delete and are probably the most interesting case here. Deleted tweets from users who have made their accounts protected are probably not really deleted, but are only available to read for a very small group of users. The third case involves users who have had their entire accounts deleted and thus none of their tweets are available any more. While it is possible for a user to delete his account himself, it is much more likely that these users are spammers and have had their accounts deleted by Twitter. Statistics about these three types of deletions are shown in Table 4. Most of the deleted tweets are genuine deletions rather than a consequence of deleting spammers, showing that there is much more to predicting deletions than simply predicting spam tweets.

Given this classification of deletions, we are interested in finding out how our approach performs across these different groups. Is it the case that some deletions are easier to predict than others? In order to answer this question, we test the performance of our system on the deleted tweets from these three groups. Because each of the three test sets now contains only positive examples, we measure performance in terms of accuracy instead of F_1 score. Note also that in this case accuracy is the same as recall. The third column of Table 4 shows that i) predicting deletions that are a result of deleted accounts (i.e., spotting spammers) is much easier than predicting genuine deletions, and ii) predicting which tweets will become protected is the hardest task.

Our manual analysis of the tweets discovered that a lot of deleted tweets contained curse words, leading us to examine the relationship between cursing and deletion in more detail. Curse words are known to express negative emotions (Jay, 2009), which lead us to hypothesize that tweets which contain curse words are more likely to be deleted. In order to test this hypothesis, we calculate the probabilities of a tweet being deleted conditioned on whether it contains a curse word. We use a list of 68 English curse words, and only consider English tweets from the test set. We find that the probability of deletion given that the tweet contains a curse word is 3.73%, compared to 3.09% for tweets that do not contain curse words. We perform a two-sample z-test and find that the difference is statistically significant at p=0.0001, which supports our hypothesis.

6 Conclusion

We have proposed a new task: predicting which messages on Twitter will be deleted in the future. We presented an analysis of the deleted messages on Twitter, providing insight into the different reasons why people delete tweets. To the best of our knowledge, we are the first to conduct such an analysis. Our analysis showed, e.g., that tweets which contain swear words are more likely to be deleted. Finally, we presented a machine learning approach and showed that for certain groups of users it can predict deleted messages with very high accuracy.

References

[Bamman et al.2012] David Bamman, Brendan O'Connor, and Noah A. Smith. 2012. Censorship and Deletion Practices in Chinese Social Media. *First Monday*, 17(3).

[Bollen et al.2011] Johan Bollen, Huina Mao, and Xiaojun Zeng. 2011. Twitter mood predicts the stock market. *Journal of Computational Science*, 2(1):1–8.

[Cha et al.2010] Meeyoung Cha, Hamed Haddadi, Fabricio Benevenuto, and Krishna P. Gummadi. 2010. Measuring user influence in Twitter: The million follower fallacy. In 4th International AAAI Conference on Weblogs and Social Media (ICWSM), pages 10–17.

[Cortes and Vapnik1995] Corina Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20(3):273–297.

[Crammer et al.2006] Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. 2006. Online passive-aggressive algorithms.

- *The Journal of Machine Learning Research*, 7:551–585.
- [Dabbish et al.2003] Laura Dabbish, Gina Venolia, and JJ Cadiz. 2003. Marked for deletion: an analysis of email data. In *CHI '03 extended abstracts on Human factors in computing systems*, CHI EA '03, pages 924–925, New York, NY, USA. ACM.
- [Dabbish et al.2005] Laura Dabbish, Robert E. Kraut, Susan Fussell, and Sara Kiesler. 2005. Understanding email use: predicting action on a message. In Proceedings of the SIGCHI conference on Human factors in computing systems, pages 691–700. ACM.
- [Fan et al.2008] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. 2008. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874.
- [Jay2009] Timothy Jay. 2009. The utility and ubiquity of taboo words. *Perspectives on Psychological Science*, 4(2):153–161.
- [O'Connor et al.2010] Brendan O'Connor, Ramnath Balasubramanyan, Bryan R. Routledge, and Noah A. Smith. 2010. From tweets to polls: Linking text sentiment to public opinion time series. In *Proceedings* of the 4th International Conference on Weblogs and Social Media, pages 122–129. The AAAI Press.
- [Petrović et al.2011] Saša Petrović, Miles Osborne, and Victor Lavrenko. 2011. RT to Win! Predicting Message Propagation in Twitter. In *Proceedings of ICWSM*.
- [Ritterman et al. 2009] Joshua Ritterman, Miles Osborne, and Ewan Klein. 2009. Using prediction markets and twitter to predict a swine flu pandemic. In *1st International Workshop on Mining Social Media*.
- [Tschang2012] Chi-Chu Tschang. 2012. An analysis of sina weibo censorship using weiboscope search data. http://partnews.mit.edu/.
- [Tumasjan et al.2010] A. Tumasjan, T.O. Sprenger, P.G. Sandner, and I.M. Welpe. 2010. Predicting elections with Twitter: What 140 characters reveal about political sentiment. In *International AAAI* Conference on Weblogs and Social Media, Washington, DC.