

Tweet Acts: A Speech Act Classifier for Twitter

Soroush Vosoughi, Deb Roy

Massachusetts Institute of Technology
Cambridge, MA 02139

soroush@mit.edu, dkroy@media.mit.edu

Abstract

Speech acts are a way to conceptualize speech as action. This holds true for communication on any platform, including social media platforms such as Twitter. In this paper, we explored speech act recognition on Twitter by treating it as a multi-class classification problem. We created a taxonomy of six speech acts for Twitter and proposed a set of semantic and syntactic features. We trained and tested a logistic regression classifier using a data set of manually labelled tweets. Our method achieved a state-of-the-art performance with an average F1 score of more than 0.70. We also explored classifiers with three different granularities (Twitter-wide, type-specific and topic-specific) in order to find the right balance between generalization and overfitting for our task.

Introduction

In recent years, the micro-blogging platform Twitter has become a major social media platform with hundreds of millions of users. People turn to Twitter for a variety of purposes, from everyday chatter to reading about breaking news. The volume plus the public nature of Twitter (less than 10% of Twitter accounts are private (Moore 2009)) have made Twitter a great source of data for social and behavioural studies. These studies often require an understanding of what people are tweeting about. Though this can be coded manually, in order to take advantage of the volume of tweets available automatic analytic methods have to be used. There has been extensive work done on computational methods for analysing the linguistic content of tweets. However, there has been very little work done on classifying the pragmatics of tweets. Pragmatics looks beyond the literal meaning of an utterance and considers how context and intention contribute to meaning. A major element of pragmatics is the intended communicative act of an utterance, or what the utterance was meant to achieve.

It is essential to study pragmatics in any linguistic system because at the core of linguistic analysis is studying what language is used for or what we do with language. Linguistic communication and meaning can not truly be studied without studying pragmatics. Proposed by Austin (Austin 1962)

Copyright © 2016, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

and refined by Searle (Searle 1969), speech act theory can be used to study pragmatics. Amongst other things, the theory provides a formalized taxonomy (Searle 1976) of a set of communicative acts, more commonly known as speech acts.

There has been extensive research done on speech act (also known as dialogue act) classification in computational linguistics, e.g., (Stolcke et al. 2000). Unfortunately, these methods do not map well to Twitter, given the noisy and unconventional nature of the language used on the platform. In this work, we created a supervised speech act classifier for Twitter, using a manually annotated dataset of a few thousand tweets, in order to be better understand the meaning and intention behind tweets and uncover the rich interactions between the users of Twitter. Knowing the speech acts behind a tweet can help improve analysis of tweets and give us a better understanding of the state of mind of the users. Moreover, as we have shown in our previous works (Vosoughi and Roy 2015; Vosoughi 2015), speech act classification is essential for detection of rumors on Twitter. Finally, knowing the distribution of speech acts of tweets about a particular topic can reveal a lot about the general attitude of users about that topic (e.g., are they confused and are asking a lot of questions? Are they outraged and demanding action? Etc).

Problem Statement

Speech act recognition is a multi-class classification problem. As with any other *supervised* classification problem, a large labelled dataset is needed. In order to create such a dataset we first created a taxonomy of speech acts for Twitter by identifying and defining a set of commonly occurring speech acts. Next, we manually annotated a large collection of tweets using our taxonomy. Our primary task was to use the expertly annotated dataset to analyse and select various syntactic and semantic features derived from tweets that are predictive of their corresponding speech acts. Using our labelled dataset and robust features we trained standard, off-the-shelf classifiers (such as SVMs, Naive Bayes, etc) for our speech act recognition task.

Using Searle's speech act taxonomy (Searle 1976), we established a list of six speech act categories that are commonly seen on Twitter: *Assertion*, *Recommendation Expression*, *Question*, *Request*, and *Miscellaneous*. Table 1 shows an example tweet for each of these categories.

| Act | Example Tweet |
|-----|---|
| Asr | authorities say that the 2 boston bomb suspects are brothers are legal permanent residents of chechen origin - @nbcnews |
| Rec | If you follow this man for updates and his opinions on #Ferguson I recommend you unfollow him immediately. |
| Exp | Mila Kunis and Ashton Kutcher are so adorable |
| Que | Anybody hear if @gehrig38 is well enough to attend tonight? #redsox |
| Req | rt @craigyh999: 3 days until i run the london marathon in aid of the childrens hopsice @sschospices . please please sponsor me here |
| Mis | We'll continue to post information from #Ferguson throughout the day on our live-blog |

Table 1: Example tweets for each speech act type.

Data Collection and Datasets

Given the diversity of topics talked about on Twitter, we wanted to explore topic and type dependent speech act classifiers. We used Zhao et al.’s (Zhao and Jiang 2011) definitions for topic and type. A *topic* is a subject discussed in one or more tweets (e.g., Boston Marathon bombings, Red Sox, etc). The *type* characterizes the nature of the topic, these are: *Entity-oriented*, *Event-oriented topics*, and *Long-standing topics* (topics about subjects that are commonly discussed).

We selected two topics for each of the three topic types described in the last section for a total of six topics (see Figure 1 for list of topics). We collected a few thousand tweets from the Twitter public API for each of these topics using topic-specific queries (e.g., #fergusonriots, #redsox, etc). We then asked three undergraduate annotators to independently annotate each of the tweets with one of the speech act categories described earlier. The *kappa* for the annotators was 0.78. For training, we used the label that the majority of annotators agreed upon (7,563 total tweets).

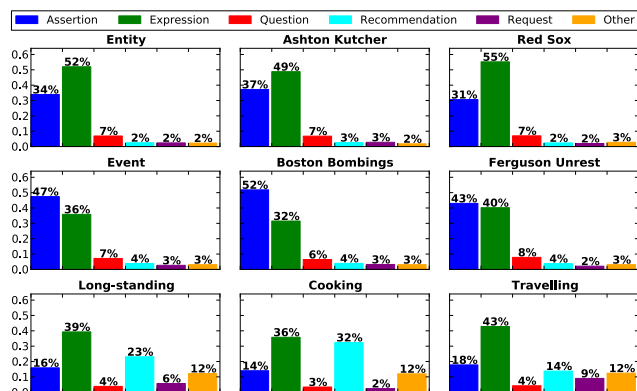


Figure 1: Distribution of speech acts for all six topics and three types.

The distribution of speech acts for each of the six topics and three types is shown in Figure 1. There is much

greater similarity between the distribution of speech acts of topics of the same type (e.g. Ashton Kutcher and Red Sox) compared to topics of different types. Though each topic type seems to have its own distinct distribution, *Entity* and *Event* types have much closer resemblance to each other than *Long-standing*. Assertions and expressions dominate in *Entity* and *Event* types with questions being a distant third, while in *Long-standing*, recommendations are much more dominant with assertions being less so. This agrees with Zhao et al.’s (Zhao and Jiang 2011) findings that tweets about *Long-standing* topics tend to be more opinionated which would result in more recommendations and expressions and fewer assertions.

The great variance across types and the small variance within types suggests that a type-specific classifier might be the correct granularity for Twitter speech act classification (with topic-specific being too narrow and Twitter-wide being too general). We will explore this in greater detail in the next sections of this paper.

Features

We studied many features before settling on the features below. Our features can be divided into two general categories: *Semantic* and *Syntactic*. Some of these features were motivated by various works on speech act classification, while others are novel features. Overall we selected 3313 binary features, composed of 1647 semantic and 1666 syntactic features.

Semantic Features

Opinion Words: We used the "Harvard General Inquirer" lexicon (Stone et al. 1968), which is a dataset used commonly in sentiment classification tasks, to identify 2442 strong, negative and positive opinion words (such as *robust*, *terrible*, *untrustworthy*, etc). The intuition here is that these opinion words tend to signal certain speech acts such as expressions and recommendations. One binary feature indicates whether any of these words appear in a tweet.

Vulgar Words: Similar to opinion words, vulgar words can either signal great emotions or an informality mostly seen in expressions than any other kind of speech act (least seen in assertions). We used an online collection of vulgar words¹ to collect a total of 349 vulgar words. A binary feature indicates the appearance or lack thereof of any of these words.

Emoticons: Emoticons have become ubiquitous in online communication and so cannot be ignored. Like vulgar words, emoticons can also signal emotions or informality. We used an online collection of text-based emoticons² to collect a total of 362 emoticons. A binary feature indicates the appearance or lack thereof of any of these emoticons.

Speech Act Verbs: There are certain verbs (such as *ask*, *demand*, *promise*, *report*, etc) that typically signal certain speech acts. Wierzbicka (Wierzbicka 1987) has compiled a

¹<http://www.noswearing.com/dictionary>

²<http://pc.net/emoticons/>

total of 229 English speech act verbs divided into 37 groups. Since this is a collection of verbs, it is crucially important to only consider the verbs in a tweet and not any other word class (since some of these words can appear in multiple part-of-speech categories). In order to do this, we used Owoputi et al.’s (Owoputi et al. 2013) Twitter part-of-speech tagger to identify all the verbs in a tweet, which were then stemmed using *Porter Stemming* (Porter 1980). The stemmed verbs were then compared to the 229 speech act verbs (which were also stemmed using Porter Stemming). Thus, we have 229 binary features coding the appearance or lack thereof of each of these verbs.

N-grams: In addition to the verbs mentioned, there are certain phrases and non-verb words that can signal certain speech acts. For example, the phrase “*I think*” signals an expression, the phrase “*could you please*” signals a request and the phrase “*is it true*” signals a question. Similarly, the non-verb word “*should*” can signal a recommendation and “*why*” can signal a question.

These words and phrases are called n-grams (an n-gram is a contiguous sequence of n words). Given the relatively short sentences on Twitter, we decided to only consider unigram, bigram and trigram phrases. We generated a list of all of the unigrams, bigrams and trigrams that appear at least five times in our tweets for a total of 6,738 n-grams. From that list we selected a total of 1,415 n-grams that were most predictive of the speech act of their corresponding tweets but did not contain topic-specific terms (such as *Boston*, *Red Sox*, etc). There is a binary feature for each of these sub-trees indicating their appearance.

Syntactic Features

Punctuations: Certain punctuations can be predictive of the speech act in a tweet. Specifically, the punctuation *?* can signal a question or request while *!* can signal an expression or recommendation. We have two binary features indicating the appearance or lack thereof of these symbols.

Twitter-specific Characters: There are certain Twitter-specific characters that can signal speech acts. These characters are *#*, *@*, and *RT*. The position of these characters is also important to consider since Twitter-specific characters used in the initial position of a tweet is more predictive than in other positions. Therefore, we have three additional binary features indicating whether these symbols appear in the initial position.

Abbreviations: Abbreviations are seen with great frequency in online communication. The use of abbreviations (such as *b4* for *before*, *jk* for *just kidding* and *irl* for *in real life*) can signal informal speech which in turn can signal certain speech acts such as expression. We collected 944 such abbreviations from an online dictionary³ and Crystal’s book on language used on the internet (Crystal 2006). We have a binary feature indicating the presence of any of the 944 abbreviations.

³<http://www.netlingo.com/category/acronyms.php>

Dependency Sub-trees: Much can be gained from the inclusion of sophisticated syntactic features such as dependency sub-trees in our speech act classifier. We used Kong et al.’s (Kong et al. 2014) Twitter dependency parser for English (called the *TweeboParser*) to generate dependency trees for our tweets. Dependency trees capture the relationship between words in a sentence. Each node in a dependency tree is a word with edges between words capturing the relationship between the words (a word either modifies or is modified by other words). In contrast to other syntactic trees such as *constituency trees*, there is a one-to-one correspondence between words in a sentence and the nodes in the tree (so there are only as many nodes as there are words). Figure 2 shows the dependency tree of an example tweet.

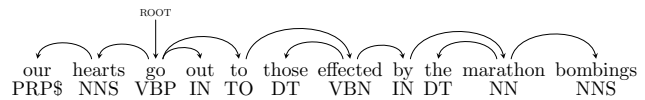


Figure 2: The dependency tree and the part of speech tags of a sample tweet.

We extracted sub-trees of length one and two (the length refers to the number of edges) from each dependency tree. Overall we collected 5,484 sub-trees that appeared at least five times. We then used a filtering process identical to the one used for n-grams, resulting in 1,655 sub-trees. There is a binary feature for each of these sub-trees indicating their appearance.

Part-of-speech: Finally, we used the part-of-speech tags generated by the dependency tree parser to identify the use of adjectives and interjections (such as *yikes*, *dang*, etc). Interjections are mostly used to convey emotion and thus can signal expressions. Similarly adjectives can signal expressions or recommendations. We have two binary features indicating the usage of these two parts-of-speech.

Supervised Speech Act Classifier

We trained four different classifiers on our 3,313 binary features using the following methods: *naive bayes (NB)*, *decision tree (DT)*, *logistic regression (LR)*, *SVM*, and a baseline max classifier *BL*. We trained classifiers across three granularities: *Twitter-wide*, *Type-specific*, and *Topic-specific*. All of our classifiers are evaluated using 20-fold cross validation.

Table 2 shows the performance of our five classifiers trained and evaluated on all of the data. We report the F1 score for each class. As shown in Table 2, the logistic regression was the performing classifier with a weighted average F1 score of .70. Thus we picked logistic regression as our classifier and the rest of the results reported will be for LR only. Table 3 shows the average performance of the LR classifier for Twitter-wide, type and topic specific classifiers.

The topic-specific classifiers’ average performance was better than that of the type-specific classifiers (.74 and .71 respectively) which was in turn marginally better than the

| | As | Ex | Qu | Rc | Rq | Mis | Avg |
|-----|------------|------------|------------|------------|------------|------------|------------|
| BL | 0. | .59 | 0. | 0. | 0. | 0. | .24 |
| DT | .57 | .68 | .79 | .32 | 0. | .29 | .58 |
| NB | .72 | .76 | .71 | .40 | 0. | .41 | .66 |
| SVM | .71 | .80 | .86 | .35 | .13 | .43 | .69 |
| LR | .73 | .80 | .87 | .30 | .16 | .45 | .70 |

Table 2: F1 scores for each speech act category. The best scores for each category are highlighted.

| | As | Ex | Qu | Rc | Rq | Mis | Avg |
|--------------|-----|-----|-----|-----|-----|-----|------------|
| Topic | .73 | .87 | .98 | .57 | .03 | .26 | .74 |
| Type | .71 | .84 | .98 | .37 | .11 | .25 | .71 |
| Twitter-wide | .73 | .80 | .87 | .30 | .16 | .45 | .70 |

Table 3: F1 scores for Twitter-wide, type-specific and topic-specific classifiers.)

performance of the Twitter-wide classifier (.70). This confirms our earlier hypothesis that the more granular type and topic specific classifiers would be superior to a more general Twitter-wide classifier.

Next, we wanted to measure the contributions of our semantic and syntactic features. To do so, we trained two versions of our Twitter-wide logistic regression classifier, one using only semantic features and the other using syntactic features. As shown in Table 4, the semantic and syntactic classifiers’ performance was fairly similar, both being on average significantly worse than the combined classifier. The combined classifier outperformed the semantic and syntactic classifiers on all other categories, which strongly suggests that both feature categories contribute to the classification of speech acts.

| | As | Ex | Qu | Rc | Rq | Mis | Avg |
|-----|------------|------------|------------|------------|------------|------------|------------|
| Sem | .71 | .80 | .62 | .22 | 0. | .23 | .64 |
| Syn | .59 | .81 | .94 | .12 | 0. | 0. | .62 |
| All | .73 | .80 | .87 | .30 | .16 | .45 | .70 |

Table 4: F1 scores for each speech act category for semantic and syntactic features.

Finally, we compared the performance of our classifier (called TweetAct) to a logistic regression classifier trained on features proposed by, as far as we know, the only other supervised Twitter speech act classifier by Zhang et al. (called Zhang). Table 5 shows the results. Not only did our classifier outperform the Zhang classifier for every class, both the semantic and syntactic classifiers (see Table 4) also generally outperformed the Zhang classifier.

| | As | Ex | Qu | Rc | Rq | Mis | Avg |
|----------|------------|------------|------------|------------|------------|------------|------------|
| Zhang | .67 | .60 | .73 | .18 | 0. | .19 | .59 |
| TweetAct | .73 | .80 | .87 | .30 | .16 | .45 | .70 |

Table 5: F1 scores for each speech act category for our classifier compared to the Zhang classifier.

Conclusions and Future Work

In this paper, we presented a supervised speech act classifier for Twitter. We treated speech act classification on Twitter as a multi-class classification problem and came up with a taxonomy of speech acts on Twitter with six distinct classes. We then proposed a set of semantic and syntactic features for supervised Twitter speech act classification. Using these features we were able to achieve state-of-the-art performance for Twitter speech act classification, with an average F1 score of .70. Speech act classification has many applications; for instance we have used our classifier to detect rumors on Twitter in a companion paper (Vosoughi and Roy 2016).

References

- Austin, J. L. 1962. *How to do things with words*. Clarendon Press, Oxford.
- Crystal, D. 2006. *Language and the internet*.
- Kong, L.; Schneider, N.; Swayamdipta, S.; Bhatia, A.; Dyer, C.; and Smith, N. A. 2014. A dependency parser for tweets. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, Doha, Qatar, to appear*.
- Moore, R. J. 2009. Twitter data analysis: An investor’s perspective. In <http://techcrunch.com/2009/10/05/twitter-data-analysis-an-investors-perspective-2/>. techcrunch.
- Owoputi, O.; O’Connor, B.; Dyer, C.; Gimpel, K.; Schneider, N.; and Smith, N. A. 2013. Improved part-of-speech tagging for online conversational text with word clusters. In *HLT-NAACL*, 380–390.
- Porter, M. F. 1980. An algorithm for suffix stripping. *Program: electronic library and information systems* 14(3):130–137.
- Searle, J. R. 1969. *Speech acts: An essay in the philosophy of language*, volume 626. Cambridge university press.
- Searle, J. R. 1976. *A taxonomy of illocutionary acts*. Linguistic Agency University of Trier.
- Stolcke, A.; Ries, K.; Coccaro, N.; Shriberg, E.; Bates, R.; Jurafsky, D.; Taylor, P.; Martin, R.; Van Ess-Dykema, C.; and Meteer, M. 2000. Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational linguistics* 26(3):339–373.
- Stone, P.; Dunphy, D. C.; Smith, M. S.; and Ogilvie, D. 1968. The general inquirer: A computer approach to content analysis. *Journal of Regional Science* 8(1):113–116.
- Vosoughi, S., and Roy, D. 2015. A human-machine collaborative system for identifying rumors on twitter. In *proceedings of ICDMW 2015*, 47–50.
- Vosoughi, S., and Roy, D. 2016. A semi-automatic method for efficient detection of stories on social media. In *proceedings of the 10th ICWSM*.
- Vosoughi, S. 2015. *Automatic Detection and Verification of Rumors on Twitter*. Ph.D. Dissertation, Massachusetts Institute of Technology.
- Wierzbicka, A. 1987. *English speech act verbs: A semantic dictionary*. Academic Press Sydney.
- Zhao, X., and Jiang, J. 2011. An empirical comparison of topics in twitter and traditional media. *Singapore Management University School of Information Systems Technical paper series*. Retrieved November 10:2011.