

Automatic Tweet Generation: An Extractive Summarization Problem?

First Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Second Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Abstract

Social media such as Twitter have become an important method of communication, with potential opportunities for NLG to facilitate the generation of social media content. We focus on the generation of *indicative tweets* that contain a link to an external web page. While it is natural and tempting to view the linked web page as the source text from which the tweet is generated in an extractive summarization setting, it is unclear to what extent actual indicative tweets behave like extractive summaries. We collect a corpus of indicative tweets with their associated articles and investigate whether they can actually be derived from the articles using extractive methods. We also consider the impact of formality and genre differences between the article and the tweet. Our results demonstrate the limits of viewing indicative tweet generation as extractive summarization, and point to the need for the development of a methodology for tweet generation that is sensitive to genre-specific issues.

1 Introduction

With the rise in popularity of social media, message broadcasting sites such as Twitter and other microblogging services have become an important means of communication, with an estimated 500 million tweets being written each day¹. In addition to individual users, various organizations and public figures such as newspapers, government officials and entertainers have established themselves on social media in order to disseminate information or promote their products.

While there has been recent progress in the development of Twitter-specific POS taggers,

parsers, and other tools (Owoputi et al., 2013; Kong et al., 2014), there has been little work on methods for generating tweets, despite the utility this would have for users and organizations.

In this paper, we consider the generation of the particular class of tweets that contain a link to an external web page that is composed primarily of text. This class of tweets, which we call *indicative tweets*, would appear to be the easiest to handle using current methods, because there is a clear source of input from which a tweet could be generated. In effect, the tweet would be acting as an indicative summary of the article being linked to, and it would seem that existing methods in summarization can be applied.

There has in fact been some work along these lines, within the framework of extractive summarization. Lofi and Krestel (2012) describe a system to generate tweets from local government records through keyphrase extraction. Lloret and Palomar (2013) compares various extractive summarization algorithms applied on Twitter data to generate tweets from documents.

Lofi and Krestel do not provide a formal evaluation of their model, while Lloret and Palomar compared overlap between system-generated and user-generated tweets using ROUGE (Lin, 2004a). Unfortunately, they also show that there is little correlation between ROUGE scores and the perceived quality of the tweets when rated by human users for indicativeness and interest. More scrutiny is required to determine whether the wholesale adoption of methods and evaluation from extractive summarization is justified.

Beyond issues of evaluation measure, it is also unclear whether extraction is what people do. One of the original motivations behind extractive summarization was the observation that human summary writers tended to extract snippets of key phrases from the source text (Mani, 2001). In Twitter data, an additional issue arises in that the

¹<https://about.twitter.com/company>

genre of the source text, often a news article or other formal text, may be vastly different from the text of the tweet itself. Thus, a genre-appropriate extract may not be available.

We begin to address the above issues through a study that examines to what extent tweet generation can be viewed as an extractive summarization problem. We extracted a data set of indicative tweets containing a link to an external article, including the documents linked to through the tweets. We used this data and applied unigram, bigram and LCS (longest common subsequence) matching techniques inspired by ROUGE to show that we need a more involved approach than directly applying existing extractive summarization algorithms developed for news text. We also use stylistic analysis on the articles to examine the role of genre differences between the source text and the target tweet.

Our results point to the need for the development of a methodology for indicative tweet generation that is sensitive to stylistic factors.

2 Background and Related Work

With the increase in the number of users on Twitter, there has also been an increase in the number of studies on Twitter data, towards classifying and summarizing text, identifying intents of tweets, event summarization, and so on. Ghosh et al. (2011) classified the retweeting activity of users based on entropy. The study considered the occurrence of the same URL in a different tweet as a retweet, and was able to separate the tweets as automatic or robotic retweeting, campaigns, news, blogs and so on. The study shows some interesting trends of retweeting activity for each of these cases. In another study, Chen et al. (2012), were able to extract sentiment expressions based on their corpus of tweets, that resulted in extraction of both formal and slang sentiment bearing words.

Other studies using Twitter data include O'Connor et al. (2010), who use topic summarization for a given search for better browsing. Chakrabarti and Punera (2011) generate an event summary by learning about the event using a Hidden Markov Model over the tweets describing it. Wang et al. (2014) generate a coherent event summary by treating summarization as an optimization problem for topic cohesion. Inouye and Kalita (2011) compare multiple summariza-

tion techniques to generate a summary of multi-post blogs on Twitter.

Studies on classifying user intents in tweets are interpreted in different ways. Banerjee et al. (2012) analyze real time data to detect presence of intents in tweets. Wang et al. (2015) classify intents as food and drink, travel, career and so on, ones that can directly be used as intents for purchasing and can be utilized for advertisements. They also focus on finding tweets with intent and then classifying those. Gómez-Adorno et al. (2014) use features from text and stylistics to determine user intentions, which are classified as news report, opinion, publicity and so on. Mohammad et al. (2013) study the classification of user intents specifically for tweets related to elections. They study one election and classify tweets as ones that agree or disagree with the candidate, ones that are meant for humor, support and so on.

As described in Section 1, we analyze tweet generation using extractive summarization techniques. There has been one such study comparing different text summarization techniques for tweet generation by Lloret and Palomar (2013). Summarization systems were used to summarize texts to sentences and then were compared against each other, evaluated using the ROUGE metric for evaluation. The ROUGE-1, ROUGE-2 and ROUGE-L metrics were used and the tweets were compared against an ideal summary. ROUGE (Lin, 2004b) is a recall based n-gram counting evaluation metric developed for summarization (Nenkova, 2006). However, it reflects the summarization quality better when used with multiple reference texts and is not meant to be used at the sentence level. However, since extractive summarization algorithms are being compared, ROUGE is used for the evaluation.

The limits of extractive summarization have been studied by He et al. (2000) by comparing user preferences for multiple types of summaries for an audio-visual presentation. They demonstrate that the most preferred method of summarization is highlights and notes provided by the author, rather than transcripts or slides from the presentation. Conroy et al. (2006) have defined an oracle score towards the same aim. The oracle score is based on the maximum likelihood probability of words occurring in model summaries and is in turn used to generate summaries that perform better than any extracted and also human-generated

Mention intent of the tweets?

summaries. These studies show that extractive summarization algorithms may not generate good quality summaries even after giving high ROUGE evaluation scores.

3 Data Extraction and Preprocessing

3.1 Using Twitter for Data Extraction

As mentioned earlier, there have been numerous studies that used data from the public Twitter feeds. However, since none of the datasets used in these studies focused on tweets and related articles linked to these tweets separated into categories as required for this study, we extracted data directly from the site. This section describes extraction, cleaning and other preprocessing of the data.

3.2 Extracting Data

Data was extracted from Twitter using the Twitter REST API using 51 search terms, or hashtags. These hashtags were chosen from a range of topics including pop culture, international summit meetings discussing political issues, lawsuits and trials, social issues and health care issues. All these hashtags were trending (being tweeted about at a high rate) at the time of extraction of the data. To get a broader sample, the data was extracted over the course of 15 days in November, which gave us multiple news stories to choose from for the search terms. The search terms were chosen so that there would be equal representation in terms of various stylistic properties of text like formality, subjectivity, etc. For example, searches related to politics would be more formal, while those related to films would be informal, and would also have a lot more opinion pieces about them. A few examples of the search terms and their distribution in genre are shown in Table 1.

Only English tweets were extracted since the study is limited to English. In the beginning, about 30,000 tweets were extracted, and more than half of these tweets, around 16,000 contained URLs referencing some news articles, photos on photo sharing sites, and videos. The hashtags were chosen to maximise the number of articles related to the tweets. Hence, a lot of topics that were chosen were being tweeted about by news agencies and other popular news sources.

The data from the tweets was cleaned by removing the tweets that were not in English as well as the ones that were retweeted, which is equivalent to re-publishing the same tweet from a different

Politics	Science & Technology
#apec2014	#rosetta
#G20	#lollipop
#oscarpistorius	#mangalayan
Events	Films and Pop culture
#haiyan	#TaylorSwift
#memorialday	#theforceawakens
#ottawashootings	#johnoliver
International	Sports
#berlinwall	#ausvssa
#ebola	#playingitmyway
#erdogan	#nycmarathon

Table 1: Table of Hashtags used for extraction. Table shows some examples of search terms chosen from various different categories.

user.

Unique URLs were first extracted from the 16,000 or so URLs in the data. Next, data from these unique URLs was extracted and then pre-processed. The newspaper package² was used to extract article text and the title from the web page. For the articles obtained from URLs, photos and video links for example, from Instagram and Youtube needed to be removed. For this, the data cleaning was achieved by removing articles by limiting word length of the extracted text to about 150 words. This ensured the removal of photos, videos, advertisements, incorrectly extracted articles from the data. After this preprocessing, the number of useful articles reduced from 6003 to 3066.

The final version of the data consists of all tweets along with all the information of the tweet itself, such as the text of the tweet, links to articles if any, hashtags, and so on. The article links from these tweets are stored as a separate file, with information about the articles themselves, along with some preprocessed data. This includes the URL itself and the text extracted from the article, as well as some extracted information such as sentence boundaries, POS tags for tokens, parse trees and dependency trees. This processing of the text was done using the CoreNLP toolkit developed at Stanford Manning et al. (2014). These were used later during analysis in Section 4

Tweets are linked to URLs through another file. A URL could have been tweeted through multiple tweets, all the ids of these tweets are linked to the

²<https://pypi.python.org/pypi/newspaper>

same URL. It should be noted that the tweet to article dataset contains only the articles that are significantly long texts about the subject with a title, and contain no advertisements, other languages, or links to images or videos. Table 2 shows an example of an entry in the dataset.

Tweet	'#RiggsReport: #CA as the #Election-Night exception. Voters rewarded #GOP nationally, but not in the #GoldenState. http://t.co/K542wvSNVz '
Title	'The Riggs Report: California as the Election Night exception'
Text	'When the dust settled on Election Night last week...'

Table 2: Example of a tweet, title of the article and the text.

4 Analysis

This section details the analyses performed on the data. The analyses mimic the ROUGE-1,2 and L methods of comparing documents, where we compare the tweet and article text using these methods. This is done to determine whether the tweets promoting the articles could be generated from the document text. The results of the comparison show that the tweet is not extracted from the article text.

We calculate the degree of common words - unigrams and bigrams, between the tweet and the text of the document. We also check least common subsequences between the tweet and the document. These are the ROUGE-1,2 and L style calculations. The hypothesis is that these results give an approximation of the degree to which the tweet is extracted from the document text.

For all these analyses, the stop words have been eliminated from the tweet as well as the document, so that only the significant words are taken into consideration. The hastags, references (@) and URLs from the tweets were also all removed.

4.1 Total match with text in article

To calculate the position of tweet text as a whole in the text, we checked for a complete substring match of the tweet in the text. Out of the 2471 unique instances of tweet text and the article text pairs, where a tweet text was checked against the text in the article, a complete match was found 23 times. 9 times out of these, the tweet text had been

matched against the title of the article extracted into the text. The rest of the results are significant, since the text of the tweet appears exactly as is inside the text of the article. For these cases, the user that wrote the tweet went through the article text, and the sentence that either seemed to be the most conclusive contribution of the article, or expressed the opinion of the user was extracted to be tweeted. We also checked to see if the tweet text matched with the article titles, and this was found not to be the case. This comparison shows that the tweet is extracted from the article very few times, and does not match with the title of the articles a lot of times either.

4.2 Percentage match for unigrams

Next, we did a percentage match with the text of the article. This was a bag-of-words check using unigrams from the tweet and the document. The order of the words in the tweet or the text did not matter. The results we got seem to suggest that a lot of significant words in the tweet are in fact present in the article. The minimum percentage match obtained was 60%. However, since the order of the words did not matter, this result can be traced back to the fact that tweet is based on the same topic as the document. Figure 1 shows the percentage of matches in the tweet and the article text as compared to the number of unigrams in the tweet. The mean of the match percentages is 29.53 and standard deviation is 20.2. Figure 2 shows the number of articles with same number of matching unigrams. The graph shows maximum number of articles with 2 unigrams matched. The number of articles with more matched goes on decreasing. The slight rise at the end - more than 10 matched unigrams - is accounted for by the completely matched tweets described above.

$$unigramMatch = \frac{|unigrams(tweet) \cap unigrams(text)|}{|unigrams(tweet)|} * 100 \quad (1)$$

4.3 Percentage match for bigrams

Similar to the unigram matching techniques, bigram percentage matching was also calculated. The text of the tweet was converted into bigrams and we then looked for those bigrams in the article text. The percentage was calculated similar to

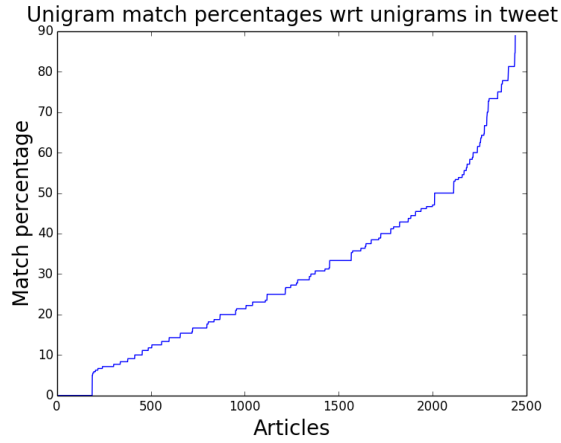


Figure 1: Unigram match percentage.

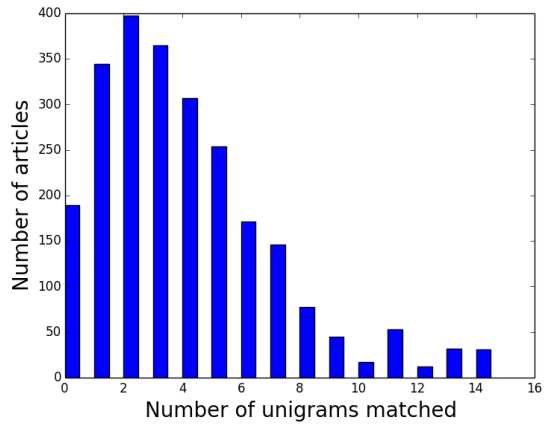


Figure 2: Bar plot of number of unique tweet-article pairs vs number of unigrams matched

the unigram matching done earlier. Figure 4 shows frequency of the number of articles for the bigrams matched for unique tweet-article pairs. There are no matched bigrams for most of the pairs. The number then decreases from one matched bigram till the end, where it increases a little at more than 10 matched bigrams, similar to the unigram frequency graph.

$$bigramMatch = \frac{|bigrams(tweet) \cap bigrams(text)|}{|bigrams(tweet)|} * 100 \quad (2)$$

Figure 3 shows the percentages of matches found in every article. Mean is 10.73 with a standard deviation of 18.5.

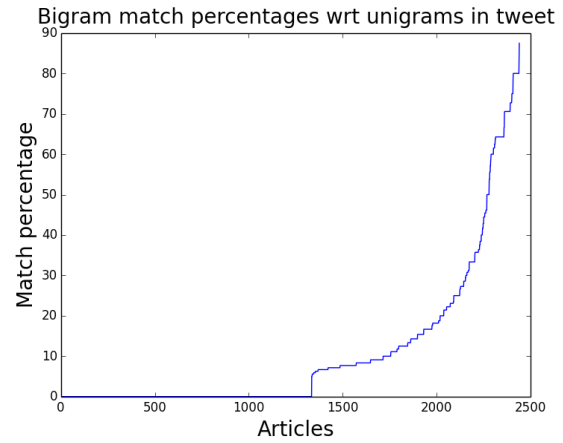


Figure 3: Bigram match percentage.

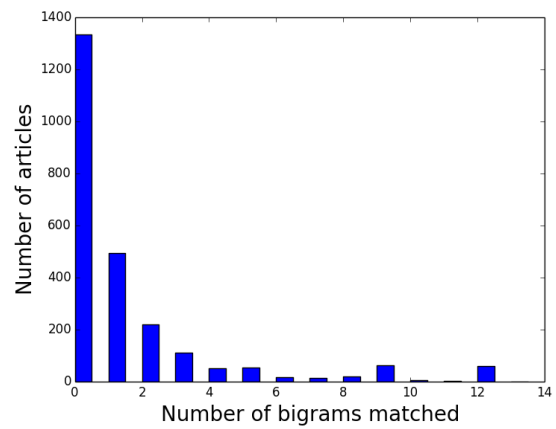


Figure 4: Bar plot of number of unique tweet-article pairs vs number of bigrams matched

4.4 Percentage matching inside a window in the article text

The next analysis was to check for a significant word matching inside a two or three sentence window inside the article text. We used a three sentence long window using the sentence boundary information obtained during preprocessing. After the text of the window was extracted, we performed a similar analysis as the last one, except on a smaller set of sentences. Again, the order of the unigrams didn't matter. Next, the matching percentages from all such windows in the articles were compared and the maximum out of these was considered for the highest match percentage and match position for the final results. The result from this experiment is shown in Figure 5. Here, the mean of the values is 26.6% and deviation 17%.

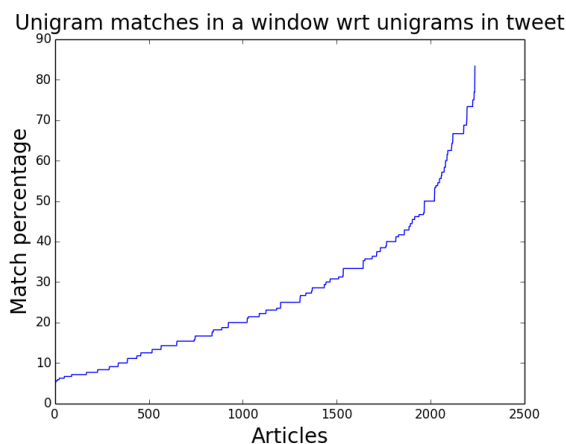


Figure 5: Percentages of common words in tweet and text.

4.5 Longest Common Subsequence match inside a window for the text

The percentage match analyses were a bag-of-words approach disregarding the order of the words inside the texts and tweets. To respect the order of the words in the sentence of the tweet, we also used the least common subsequence algorithm between the tweet text and the document text. This subsequence matching was done inside a sentence window of 5 sentences. Again, the final result for the article was the window in which the maximum percentage was recorded among all windows. The percentage match was calculated against the number of words in the tweet, as found

in the least common subsequence calculated between the two texts. These numbers are shown in Figure 6. The mean here is 44.6% and the standard deviation is 22.7%.

$$LCSMatch = \frac{|lcs(tweet, text)|}{|unigrams(tweet)|} * 100 \quad (3)$$

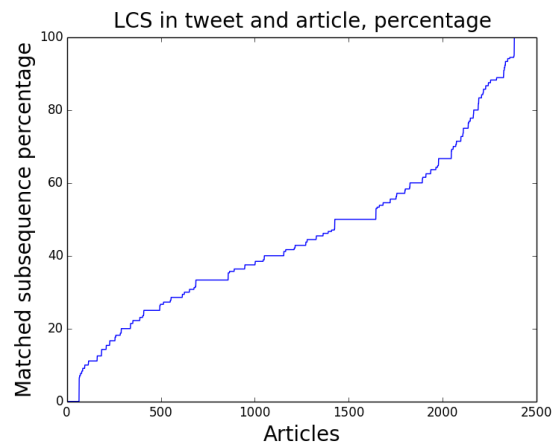


Figure 6: Percentages of words matching in tweet and document text using an LCS algorithm.

5 Results

As seen in the results of the analyses performed in Section 4, the tweets have little in common with the articles they are related to. The analyses are based on the ROUGE-1,2 and L like calculations. This shows that extractive summarization algorithms cannot be directly applied to articles to generate tweets.

To tie in the results of the findings above with some intuitive notions about the text and see how formality interacts with the results, we also calculated the formality of the articles. This formality score was correlated with the longest common subsequence. To achieve this, the degree of formality of the text was calculated with the help of Brooke and Hirst (2013). The formality lexicon was generated by analyzing the stylistics of text and can be used to measure formality of a given text. They calculate formality scores for words and sentences by training a model on a large corpus based on the appearance of words in specific documents. Their model represents words as vectors and the formal and informal seeds appear in opposite halves of the graphs, suggesting that we can use these seeds to determine if an article is

formal or informal. The lexicon consists of words and phrases and the degree of formality for their occurrence. Thus, more formal words are marked on a positive scale and informal words like those occurring in colloquial language are marked on a negative scale. The degree of formality was calculated using this lexicon.

$$formalityScore = \frac{|unigramsArticle \cap formalitySet|}{|unigramsArticle|} * 10 \quad (4)$$

The formality lexicon gave positive weights for formal expressions and negative for informal expressions. After calculating the formality weights for all articles, it was observed that they all had a total negative normalized weight, meaning a lot more informal expressions were getting matched. Hence, we used just the formal word occurrences for calculating the weight. Thus, above a certain cut-off weight, the article could be considered formal, else would be considered informal. To make sure these formality scores intuitively made sense, we calculated the average formality score for each hashtag used in the search during data extraction and ordered them, shown in Table 3

Lowest	Highest
#theforceawakens	#KevinVickers
#TaylorSwift	#erdogan
#winteriscoming	#apex

Table 3: Table of hashtags(broadly, topics) with highest and lowest formality according to the lexicon.

This formality score for each article was then correlated with the percentage of match obtained using the longest common subsequence algorithm. The Pearson correlation value was 0.41, with a p-value of 7.08e-66. The p-value justifies that we can reject the null hypothesis, and say with confidence that there is a correlation between the formality scores and the ROUGE-L scores of the tweets and articles. Hence, we can say that the more formal the subject or the article, there are higher chances of the tweet being extracted directly from the article.

6 Conclusion

We have described a study on investigating whether indicative tweet generation can be viewed as an extractive summarization problem. By ana-

lyzing a collection of indicative tweets that we collected according to measures inspired by extractive summarization evaluation measures, we find that most tweets cannot be recovered from the article that they link to, demonstrating a limit to the effectiveness of extractive methods.

We further performed an analysis to determine the role of formality differences between the source article and the Twitter genre. We find evidence that formality is an important factor, as the more formal the source article is, the less extractive the tweets seem to be. Future methods that can change the level of formality of a piece of text without changing the contents will be needed.

References

- Nilanjan Banerjee, Dipanjan Chakraborty, Anupam Joshi, Sumit Mittal, Angshu Rai, and Balaraman Ravindran. 2012. Towards analyzing micro-blogs for detection and classification of real-time intentions. In *ICWSM*.
- Julian Brooke and Graeme Hirst. 2013. A multi-dimensional bayesian approach to lexical style. In *HLT-NAACL*, pages 673–679.
- Deepayan Chakrabarti and Kunal Punera. 2011. Event summarization using tweets. *ICWSM*, 11:66–73.
- Lu Chen, Wenbo Wang, Meenakshi Nagarajan, Shaojun Wang, and Amit P Sheth. 2012. Extracting diverse sentiment expressions with target-dependent polarity from twitter. In *ICWSM*.
- John M Conroy, Judith D Schlesinger, and Dianne P O’Leary. 2006. Topic-focused multi-document summarization using an approximate oracle score. In *Proceedings of the COLING/ACL on Main conference poster sessions*, pages 152–159. Association for Computational Linguistics.
- Rumi Ghosh, Tawan Surachawala, and Kristina Lerman. 2011. Entropy-based classification of ‘retweeting’ activity on twitter. *arXiv preprint arXiv:1106.0346*.
- Helena Gómez-Adorno, David Pinto, Manuel Montes, Grigori Sidorov, and Rodrigo Alfaro. 2014. Content and style features for automatic detection of users intentions in tweets. In *Advances in Artificial Intelligence-IBERAMIA 2014*, pages 120–128. Springer.
- Liwei He, Elizabeth Sanocki, Anoop Gupta, and Jonathan Grudin. 2000. Comparing presentation summaries: slides vs. reading vs. listening. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*, pages 177–184. ACM.

- David Inouye and Jugal K Kalita. 2011. Comparing twitter summarization algorithms for multiple post summaries. In *Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (Social-Com), 2011 IEEE Third International Conference on*, pages 298–306. IEEE.
- Lingpeng Kong, Nathan Schneider, Swabha Swayamdipta, Archana Bhatia, Chris Dyer, and Noah A. Smith. 2014. A dependency parser for tweets. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1001–1012, Doha, Qatar, October. Association for Computational Linguistics.
- Chin-Yew Lin. 2004a. Rouge: A package for automatic evaluation of summaries. In *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*, pages 74–81.
- Chin-Yew Lin. 2004b. Rouge: A package for automatic evaluation of summaries. In *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*, pages 74–81.
- Elena Lloret and Manuel Palomar. 2013. Towards automatic tweet generation: A comparative study from the text summarization perspective in the journalism genre. *Expert Systems with Applications*, 40(16):6624–6630.
- Christoph Lofi and Ralf Krestel. 2012. iparticipate: Automatic tweet generation from local government data. In *Database Systems for Advanced Applications*, pages 295–298. Springer.
- Inderjeet Mani. 2001. *Automatic summarization*, volume 3. John Benjamins Publishing.
- Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J Bethard, and David McClosky. 2014. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60.
- Saif M Mohammad, Svetlana Kiritchenko, and Joel Martin. 2013. Identifying purpose behind electoral tweets. In *Proceedings of the Second International Workshop on Issues of Sentiment Discovery and Opinion Mining*, page 1. ACM.
- Ani Nenkova. 2006. Summarization evaluation for text and speech: issues and approaches. In *INTER-SPEECH*.
- Brendan O’Connor, Michel Krieger, and David Ahn. 2010. Tweetmotif: Exploratory search and topic summarization for twitter. In *ICWSM*.
- Olutobi Owoputi, Brendan O’Connor, Chris Dyer, Kevin Gimpel, Nathan Schneider, and Noah A. Smith. 2013. Improved part-of-speech tagging for online conversational text with word clusters. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 380–390, Atlanta, Georgia, June. Association for Computational Linguistics.
- Lu Wang, Claire Cardie, and Galen Marchetti. 2014. Socially-informed timeline generation for complex events. *constitution*.
- Jinpeng Wang, Gao Cong, Xin Wayne Zhao, and Xiaoming Li. 2015. Mining user intents in twitter: A semi-supervised approach to inferring intent categories for tweets. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.