

Duluth at SemEval-2017 Task 6: Language Models in Humor Detection

Xinru Yan

Department of Computer Science
University of Minnesota Duluth
Duluth, MN, 55812 USA
yanxx418@d.umn.edu

Ted Pedersen

Department of Computer Science
University of Minnesota Duluth
Duluth, MN, 55812 USA
tpederse@d.umn.edu

Abstract

This paper describes the Duluth system that participated in SemEval-2017 Task 6 #HashtagWars: Learning a Sense of Humor. The system completed Subtask A and Subtask B using N-gram language models, ranking well during evaluation. This paper includes the results of our system during development and evaluation stage, with two post-evaluation results.

1 Introduction

Since humor represents human uniqueness and intelligence to some extent, it has continuously drawn attention in different research areas such as linguistics, psychology, philosophy and computer science. In computer science, relevant theories derived from those fields have formed a relatively young area of study, *computational humor* (Zhang and Liu, 2014). Humor has not been addressed broadly in current computational research. Many studies have developed decent systems to produce humor (Ozbal and Strapparava, 2012). However, *humor detection* is essentially a more challenging and fun problem. For example, Mihalcea and Strapparava draw particular attention on automatic humor recognition in their work (Mihalcea and Strapparava, 2006). SemEval-2017 Task 6 focuses on *humor detection* by asking participants to develop systems that learn a sense of humor from the Comedy Central TV show, *@midnight with Chris Hardwick*. Our system, Duluth, applies the language model approach to detect humor by training N-gram language models on two sets of training data, the tweets data and the news data.

2 Background

Training **Language models** (LMs) is a straightforward way to collect set of rules by utilizing the

fact that words do not appear in an arbitrary order, which means we can gain useful information from a word and its neighbors (Jurafsky and Martin, 2009). A statistical language model is a model that computes the probability of a sequence of words or an upcoming word (Jurafsky and Martin, 2009). Below are two examples of language modeling:

To compute the probability of a sequence of words W given the sequence (w_1, w_2, w_3) , we have:

$$P(W) = P(w_1, w_2, w_3) \quad (1)$$

To compute the probability of an upcoming word w_3 given the sequence (w_1, w_2) , we have:

$$P(w_3|w_1, w_2) \quad (2)$$

The idea of word prediction with probabilistic models is called the N-gram model, which predicts the upcoming word from the previous N-1 words. An N-gram is a contiguous sequence of N words: a unigram is a single word, a bigram is a two-word sequence of words and a trigram is a three-word sequence of words. For example, in tweet “tears in Ramen #SingleLifeIn3Words”, “tears”, “in”, “Ramen” and “#SingleLifeIn3Words” are unigrams; “tears in”, “in Ramen” and “Ramen #SingleLifeIn3Words” are bigrams and “tears in Ramen” and “in Ramen #SingleLifeIn3Words” are trigrams.

When we use for example, a trigram LM, to predict the conditional probability of the next word, we are thus making the following approximation:

$$P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-2}, w_{n-1}) \quad (3)$$

The assumption that the probability of a word depends only on a small number of previous words is called the **Markov** assumption (Markov, 1954). According to the Markov assumption, here we

show the general equation for computing the probability of a complete word sequence using a tri-gram LM:

$$P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-2}, w_{k-1}) \quad (4)$$

In the study on how phrasing affects memorability, in order to analyze the characteristics of memorable quotes, researchers take the language model approach to investigate the distinctiveness feature (Danescu-Niculescu-Mizil et al., 2012). Specifically, to evaluate how distinctive a quote is, they evaluate its likelihood with the respect of the “common language” model which consists of the newswire sections of the Brown corpus. They employ LMs on the “common language” and come to the conclusion that movie quotes which are less like the “common language” are more memorable. The idea of using LMs to assess the memorability of a quote is suitable for our purpose of detecting how humorous a tweet is. Except for using tweets provided by the task to train N-gram LMs, our system also trained N-gram LMs on English news data in order to evaluate how distinctive, in this case, how funny, a tweet is comparing to the tweets model and the news model. Tweets that are more like the tweets model, or less like the news model, are ranked as being more funny. For our purpose, we trained bigram LMs and trigram LMs on both sets of training data.

KenLM, as a language modeling tool, is used in our system (Heafield et al., 2013). LMs are estimated from the corpus using modified Kneser-Ney smoothing without pruning. KenLM reads a text file and generates LMs in ARPA format. KenLM also implements back-off technique, which means if the N-gram is not found, it applies the lower order N-gram’s probability along with its back-off weights. Instead of using the real probability of the N-gram, KenLM applies the logarithm scheme.

3 Method

Our system estimated tweet probability using N-gram LMs. Specifically, it solved the comparison (Subtask A) and semi-ranking (Subtask B) subtasks in four steps:

1. Corpus preparation and pre-processing: Collected all training data files to form one training corpus. Pre-processing included filtering and tokenization.
2. Language model training: Built N-gram LMs by feeding the corpus to the KenLM Language Model Toolkit.
3. Tweet scoring: Computed log probability for each tweet based on the trained N-gram LM.
4. Tweet prediction: Based on the log probability
 - Subtask A – Given two tweets, compared them and predicted which one is funnier.
 - Subtask B – Given a set of tweets associated with one hashtag, ranked tweets from the funniest to the least funny.

3.1 Corpus Preparation and Pre-processing

The Duluth system uses two distinct sets of training data: the tweets data and the news data. The tweets data was provided by the SemEval task. It consisted of 106 hashtag files with about 21,000 tokens. We collected in total of 6.2 GB of English news data with about 2,000,000 tokens from the News Commentary Corpus and the News Crawl Corpus from 2008, 2010 and 2011¹.

3.1.1 Preparation

To prepare the tweets training corpus, the system took each tweet from the hashtag files, which included tweets from both *train_dir* and *trial_dir* from the task, and created a text file with each tweet on its own line. During the development stage of the system we trained LMs solely on the *train_dir* data, which included 100 hashtag files; we tested the system on the *trial_dir* data consisting of 6 hashtag files. For the news data, the system read each sentence from the news files and created a text file with each sentence per line to form the news training corpus.

3.1.2 Pre-processing

The pre-processing consisted of two steps: filtering and tokenization. The filtering step was only for the tweet training corpus. We experimented with various filtering and tokenization combinations during the development stage to determine the best setting.

- Filtering: the filtering process included removing the following elements from the tweets:

¹<http://www.statmt.org/wmt11/translation-task.html#download>

- URLs
- Twitter user names: Tokens starting with the “@” symbol
- Hashtags: Tokens starting with the “#” symbol
- Tokenization: For both training data sets we split text by spaces and punctuation marks

3.2 Language Model Training

Once we had the corpora ready, we used the KenLM Toolkit to train the N-gram LMs on each corpus. We trained two different LMs, bigrams and trigrams, on both tweets and news training data sets. Table 1 shows an example ARPA file of the trigram LM we trained on the tweets data:

N-gram 1 = 21580
N-gram 2 = 60624
N-gram 3 = 73837
unigram:
-4.8225346 <unk> 0
0 <s> -0.47505832
-1.4503417 </s> 0
-4.415446 Donner -0.12937292
...
bigrams:
...
-0.9799023 Drilling Gulf -0.024524588
...
trigrams:
...
-1.171928 I'll start thinking
...

Table 1: An example ARPA file of a trigram LM trained on the tweets data

The ARPA file starts with a header, listing the number of each N-gram. Each N-gram line starts with the logarithm probability of that N-gram, followed by the N-gram which consists of N words. The logarithm of the back-off weight for the N-gram optionally follows after. Notice that there are three “special” words in a language model: the beginning of a sentence denoted by <s>, the end of a sentence denoted by </s> and the out of vocabulary word denoted by <unk>. In order to be able to handle the unknown words to estimate the probability of a tweet more accurately, in all our experiments we kept the <unk> word in our LMs. To derive the best setting of the LMs for both tasks,

we experimented using the language model with and without sentence boundaries.

3.3 Tweet Scoring

After training the N-gram LMs, the next step was scoring. For each hashtag file that needed to be evaluated, based on the trained N-gram LM, a logarithm score was assigned by our system for each tweet in the hashtag file. The larger the score, the more likely that the tweet appeared with respect to that LM. Table 2 shows an example of two scored tweets from hashtag file *Bad_Job_In_5_Words.tsv* based on the trigram LM trained on the tweets data.

3.4 Tweet Prediction

The system sorted tweets for each hashtag file based on their score, meaning the funniest one was listed on the top i.e. if the system used a tweets LM, the tweets would be sorted in descending order. In the case that it used a news LM, the tweets would be sorted in ascending order. For Subtask A, given a hashtag file, the system went through the sorted list of tweets, compared each pair of tweets and produced a tsv format file. For each tweet pair, if the first tweet was funnier than the second one, the system would output the tweet_ids for the pair followed by “1”. Otherwise it output the tweet_ids followed by “0”. For Subtask B, given a hashtag file, the system output the tweet_ids starting from the funniest.

4 Experiments and Results

In this section we present the results from the development stage (Table 3), the evaluation stage (Table 4), as well as two post-evaluation results (Table 4). Since we implemented both bigram and trigram LMs during the development stage but only results from trigram LMs were submitted to the task, we evaluated bigram LMs in the post-evaluation stage. Note that the accuracy and distance measurements listed in Table 3 and Table 4 are provided by the task.

Table 3 shows results from the development stage. From this table we can estimate the best setting to train LMs for both data sets: for the tweets data we decided to use trigrams and omit sentence boundaries; for the news data we chose to train trigram LMs on a tokenized news corpus.

Table 4 shows the results of our system applying trigram LMs during evaluation along with bi-

The hashtag: #BadJobIn5Words		
tweet_id	tweet	score
705511149970726912	The host of Singled Out #Bad-JobIn5Words @midnight	-19.923433303833008
705538894415003648	Donut receipt maker and sorter #BadJobIn5Words @midnight	-27.67446517944336

Table 2: Scored tweet according to the trigram LM. The format follows .tsv file provided by the task. The first column shows tweets_id; the second column shows tweets; the third column shows the probability score computed based on the trigram LM.

DataSet	N-gram	# & @ re-moved	Sentence Bound-aries	Lowercase	Tokenization	Subtask A Accuracy	Subtask B Distance
tweets	trigram	False	False	False	False	0.543	0.887
tweets	trigram	False	True	True	False	0.522	0.900
tweets	bigram	False	False	False	False	0.548	0.900
news	trigram	NA	False	False	True	0.539	0.923
news	trigram	NA	False	False	False	0.460	0.923
news	bigram	NA	False	False	False	0.470	0.900

Table 3: Development results. The development results are based on data from *trial_dir*. In general, trigram LMs outperform bigram LMs.

DataSet	N-gram	# & @ re-moved	Sentence Bound-aries	Lowercase	Tokenization	Subtask A Accuracy	Subtask B Distance
tweets	trigram	False	False	False	False	0.397	0.967
tweets	bigram	False	False	False	False	0.406	0.944
news	trigram	NA	False	False	True	0.627	0.872
news	bigram	NA	False	False	True	0.624	0.853

Table 4: Evaluation and post-evaluation results. The evaluation and post-evaluation results are based on evaluation data from *gold_dir*. The trigram LM trained on the news data ranked 4th place for Subtask A and 1st place for Subtask B during evaluation.

gram LMs results from the post-evaluation runs. It demonstrates that trigram LMs work better than bigram LMs.

5 Discussion and Future Work

We focused on training bigram and trigram LMs because tweets are normally short and concise. Trigrams outperform bigrams because trigrams have relatively better coverage than bigrams.

After comparing the amount of tweets data and news data we used, we believe that the lack of tweets data could have caused the tweets LMs to perform worse. Therefore, one way to improve the system, especially the tweets data LM, is to collect more tweets that participate in the hashtag wars. We would also like to train news LMs using the

same amount of data we have for the tweets to see how the results compare. Additionally, we want to gather more news data and see if the quantity of news data would still make a difference.

Besides, we would like to try some machine learning techniques, specifically deep learning methods such as recurrent neural networks. Studies have shown that neural network based LMs work effectively and outperform standard back-off N-gram models (Mikolov et al., 2011). In addition, recurrent neural networks are capable of forming short term memory so it can better deal with problems associated with sequences. It would be interesting to see if some combination of these methods could enhance the system.

References

- Cristian Danescu-Niculescu-Mizil, Justin Cheng, Jon Kleinberg, and Lillian Lee. 2012. [You had me at hello: How phrasing affects memorability](#). In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers - Volume 1*. Association for Computational Linguistics, Stroudsburg, PA, USA, ACL '12, pages 892–901. <http://dl.acm.org/citation.cfm?id=2390524.2390647>.
- Kenneth Heafield, Ivan Pouzyrevsky, Jonathan H. Clark, and Philipp Koehn. 2013. Scalable modified Kneser-Ney language model estimation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*. Sofia, Bulgaria, pages 690–696.
- Daniel Jurafsky and James H. Martin. 2009. *Speech and Language Processing (2Nd Edition)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- AA Markov. 1954. Theory of algorithms [translated by jacques j. schorr-kon and pst staff] imprint moscow, academy of sciences of the ussr, 1954 [jerusalem, israel program for scientific translations, 1961; available from office of technical services, united states department of commerce] added tp in russian translation of works of the mathematical institute, academy of sciences of the ussr, v. 42. *Original title: Teoriya algorifmov*. [QA248. M2943 Dartmouth College library. US Dept. of Commerce, Office of Technical Services, number OTS 60-51085]
- Rada Mihalcea and Carlo Strapparava. 2006. Learning to laugh (automatically): Computational models for humor recognition. *Computational Intelligence* 22(2):126–142.
- Tomáš Mikolov, Stefan Kombrink, Lukáš Burget, Jan Černocký, and Sanjeev Khudanpur. 2011. Extensions of recurrent neural network language model. In *Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on*. IEEE, pages 5528–5531.
- Gözde Ozbal and Carlo Strapparava. 2012. Computational humour for creative naming. *Computational Humor 2012* page 15.
- Renxian Zhang and Naishi Liu. 2014. [Recognizing humor on twitter](#). In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*. ACM, New York, NY, USA, CIKM '14, pages 889–898. <https://doi.org/10.1145/2661829.2661997>.