

Data Science Capstone Project: NLP model

Miao YU

2014/10/10

Introduction

Portable office actually means the works done on the cellphones or tablets and we need a smart input system to saving our time on typing. So efficient keyboard is required. The core of this input system is a natural language processing model. This report is focused on this, covering from the very beginning, namely data collection, to the final model. The final model will out put some words as the next words of input sentences. The first few parts came from the milestone report with a major revision.

Data Acquisition and Cleaning

The data were downloaded from the course website (from [HC Corpora](#)) and unzipped to extract the English database as a corpus. Three text documents from the twitter, blog and news were found with each line standing for a document.

Data Pre-Summary

After check the three documents with **bash**, the basic summary of the data set is shown as follows:

	line counts	word counts	document size
twitter	2360148	30373603	166816544
news	1010242	34372530	205243643
blogs	899288	37334147	208623081

Table 1: Summary of the datasets

Twitter documents were short(of course less than 140 words) with a lot of informal characters and less grammar, which means more noise; news documents were written in a formal manner but the topics were only focused on news; blog's pattern is between twitter and news with less noise and more topics. The average length of each lines in the three database: blog > news > twitter, which means blog is the longest document class and longer document might help to build a better model for prediction in certain context.

So, the blog data will be good for us to build a model if those three document is too large to be loaded for exploring. However, using sampling will ease the burden on the calculation and finally I sampled 30,000 20,000 and 10,000 lines with seed from the blogs, news and twitter database for exploring and the left data will be sampled as test data sets.

Tokenization

The whole tokenization is aiming at removing meaningless characters and the words with low frequency to avoid overfit in the corpus. The final corpus will show the words or terms with a high frequency which will be helpful for exploring the relationship between the words and building a meaningful statistical model.

- each data set might overfit on the most common words while underfitted on the least common words
- using different sources of documents will be helpful for a supervised learning
- PCA analysis of the documents will also be helpful for a unsupervised learning to group the sources

Modeling

Pre-classification

Based on the exploratory analysis, I sampled the origin train data to get a corpus without capital characters, punctuation, numbers and stop words. I didn't stem the words because I found stemming would make the predicted words unclear. The most frequency words were used to classify the original data into different sources. Here I use news, blogs and twitters as three different sources and the most frequency 115 words occurred more than 1000 times in 60000 samples were used to class the sentences by a classification tree. Results were bad and it seemed all of the sentences will be grouped into the twitter group. This is a proof for that twitter will cover most kinds of words.

I finally dropped this idea but I think classification will be helpful to get many sub-models to speed up the whole models. Maybe some unsupervised learning models will be useful to get some topic groups for the classification and local prediction under a topic might be better for global prediction.

N-grams Extraction

The extraction of n-gram was the base of a n-gram model. After I got the n-gram corpus with **RWeka** package, the term-document matrix were used to get the counts for each n-grams(uni-gram, bi-gram and tri-gram in this study) by a row sum. It will really take a long time to get them and low frequency words would also cause overfit for certain corpus.

More data at certain frequency threshold would subset the meaningful terms from noisy terms. For example, "aaaaaaaaaa" and "capstone" might both occur two times in a small corpus. When we get a ten times bigger corpus, "aaaaaaaaaa" might still occur two or at most twenty times. But a much more meaningful word such as "capstone" might at least occur twenty times. Considering a bigger corpus will have more words, such threshold might set lower than ten.

So, I explored the data with 60,000 sentences but my final corpus is twenty times bigger, namely, 1,200,000 sentences. Here I only use a threshold of six to get a smaller but much meaningful corpus for next step. For different corpus, the threshold is actually a parameter to be optimized. On my PC, less than five will make the n-gram extraction very slow and six was my last choice and I got a uni-gram with 81,456 words, bi-gram with 413,605 words and tri-gram with 49,670 words. Using a higher probability of the last words in n-gram will predict next words for certain sentence.

Backoff and Smoothing

The n-gram model worked well if the terms were huge enough to cover any cases. However, building such model will cost a lot of time. Another way is just using a back-off model to change n-gram model into (n-1)-gram model for the unseen words in n-gram. The simplest back-off model will first get the probability of every (n-1) terms, order them and show the first few words as prediction. When no words were shown, a (n-1)-gram model will be used until uni-gram model, which will show the most common words in the corpus.

However, if there were only one terms in a tri-gram while many terms in a bi-gram for certain words. A simple back-off model can't distribute the probability to the candidates from both n and (n-1) gram. A common way is that smoothing the counts on the n-gram using interpolated smoothing models. The unseen words' probability could be saved by a discount on the observed terms' counts. I use an absolute discounting on each counts based on Ney et al.'s study:

$$D = \frac{n_1}{n_1 + 2 * n_2}$$

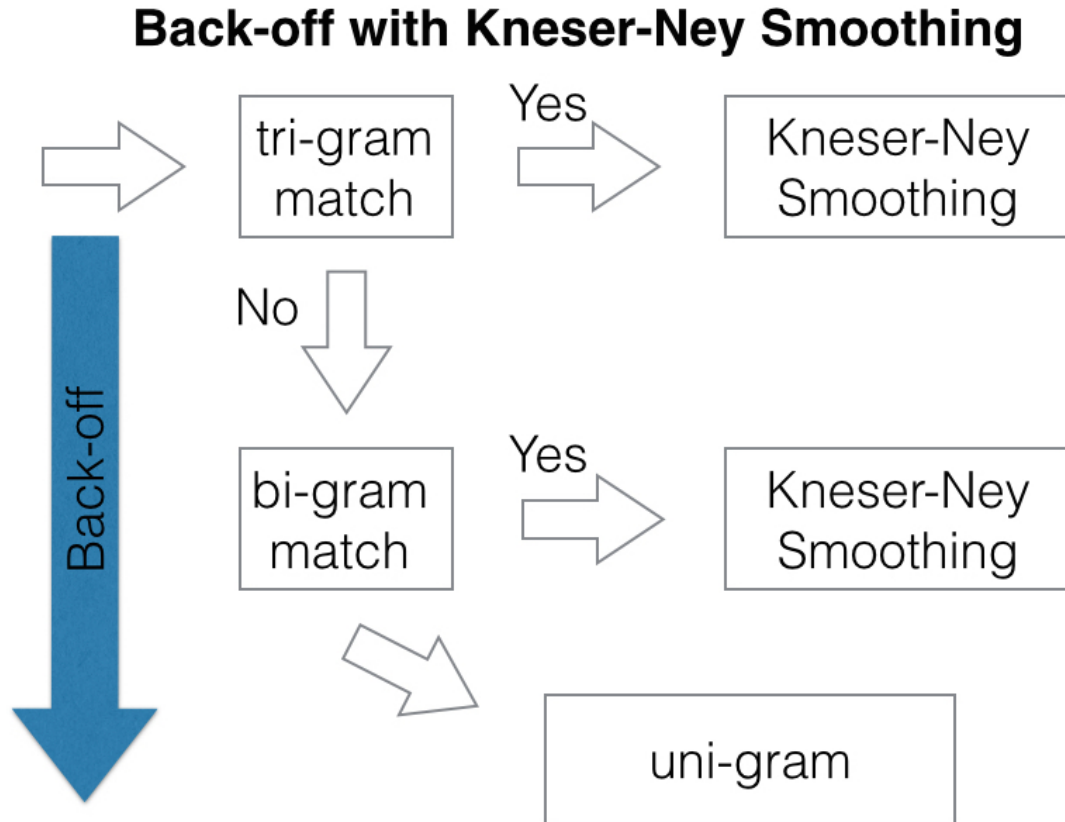
n_1 and n_2 were terms occurred exactly once and twice.

After I get the free space for interpolated smoothing of the (n-1)-gram model, the Kneser-Ney Smoothing were employed to get the probability with a combination of (n-1)-gram. The tri-gram formula is shown below:

$$P_{KN}(w_3|w_1, w_2) = \frac{\max(C(w_1, w_2, w_3) - D, 0)}{C(w_1, w_2)} + D * \frac{N(w_1, w_2, \cdot)}{C(w_1, w_2)} * (\frac{\max(N(\cdot, w_2, w_3) - D, 0)}{N(\cdot, w_2, \cdot)} + D * \frac{N(w_2, \cdot)}{N(\cdot, w_2, \cdot)} * \frac{N(\cdot, w_3)}{N(\cdot, \cdot)})$$

C stand for the counts on certain terms, N stand for the species of such terms and D is the discount. The basic idea of Kneser-Ney Smoothing was that when a word has many kinds of (n-1)-gram terms, its probability would be larger than another word has few kinds when they have the same counts in (n-1)-gram data sets.

I actually combined a Kneser-Ney Smoothing with a back-off model: When the model could find terms in the tri-gram, a tri-gram Kneser-Ney model will be used. While the model can't find a hit in tri-gram, a bi-gram Kneser-Ney Smoothing were run. Those two Kneser-Ney Smoothing have different discountings. The reason I choose a back-off model was that it would be faster. The detailed formula could be found in Körner's paper in the reference and the code for the whole process is shown in the RMD document's code chunk. A simple workflow is shown in the following picture.



Stupid Backoff Implementation

The model above were still slow to show predicted words and I speed up the model with less code and a relative small loss of prediction accuracy. The core of the code called Stupid Backoff implementation, which is often used in web-based corpus. The core of this backoff implementation is that using a fixed discount for (n-1)-gram's possibility of interpolated smoothing. With a huge corpus, the performance of Stupid Backoff implementation will show a similar prediction accuracy with the Kneser-Ney Smoothing. This method was much faster than Kneser-Ney Smoothing. The main function code is shown in the RMD document's code chunk.

Prediction

Well, though I spent a lot of time understanding the algorithms in the model, the performance of the final models were still poor in the course quiz. I think the main reason is that I only use a small part of the documents to build my model and the model is underfit. However, I suggest the reader to try it on the [web app](#) and a subjective feeling might be “objective” in this kind of model.

Summary

- the data were really BIG for PC
- exploratory analysis were very important for further modeling
- data clean should not remove the signals from noises
- N-grams could be get by a threshold using big data
- Back-off and interpolated smoothing models was useful for n-gram model
- Kneser-Ney Smoothing is slower than Stupid Backoff Implementation
- Next word prediction is a more ‘subjective’ model relying on the corpus
- Pre-classification might help to get a local best prediction
- Google is REALLY an important tool for data scientist
- this capstone project's topic is totally new to me and I am looking foreword to hearing from you

References

- Hornik, K. (2008). Journal of Statistical Software, 25(5).
- Körner, M. C. Implementation of Modified Kneser-Ney Smoothing on Top of Generalized Language Models for Next Word Prediction Bachelorarbeit, (September 2013).
- [Coursera course on NLP](#)
- [Coursera Discussion Board](#)