

Digital Speech Processing

Final Project Report

Paper Survey About Language Model and
Chinese Lyrics Decoding Using Trigram Language Model

電機四 林孟瑾 B04901009

1 Paper survey about language model

I really interested in the concepts of language model so I search for renown papers to understand more about the developments of language model.

1. A neural probabilistic language model [1]

According to Google Scholar website, there are 4685 citations about this paper.

(1) Introduction of this work

The most difficult issue in language modeling is curse dimensionality. Statistical language modeling utilizes word order to address this problem by constructing conditional probability between the neighboring words and using back-off and smoothing methods to deal with the words absent in training corpus. However, the author thinks that there are two issues that can be improved.

- i. Contexts farther than 2 words
- ii. Semantic and grammar similarity

To achieve these two goals, this work uses neural network model to deal with high-dimensional discrete distribution. It aims to train continuous distributed feature vectors and express the joint probability function of word sequences based on word feature vectors.

(2) Neural network model

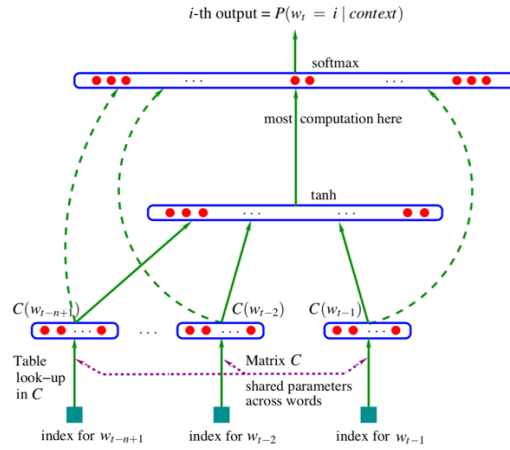
The training set is a sequence of words of vocabulary V which is large but finite set. The objective is to get the model $f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$ which gives high out-of-

sample likelihood. There are two parts of $f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$:

- i. Distributed feature vectors relating to each word in the vocabulary, expressed by a mapping C from any element i of V to a real vector $C(i) \in \mathbb{R}^m$
- ii. Probability function over words, expressed by a function g mapping input sequence of feature vectors for words to a conditional probability distribution over words in V for the next word w_t .

Therefore, the neural architecture as the following graph is a composition of two mappings C and g . C can be shared across all words in the context.

$$f(i, w_{t-1}, \dots, w_{t-n+1}) = g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$$



Training process includes looking for overall parameter set $\theta = (C, \omega)$ which maximizes L :

$$L = \frac{1}{T} \sum_t \log f(w_t, w_{t-1}, \dots, w_{t-n+1}; \theta) + R(\theta)$$

(3) Parallel Implementation

The computations of neural network model are far larger than traditional n-gram model. Nowadays, people use GPU to train neural network models to reduce computation time. So I think that GPU might not be common at the time this work published.

The parts requiring most computations are the activations of the output layers. Therefore, this work tries the methods of data parallel processing and parameter parallel processing to deal with immense computations.

- i. Data parallel processing
Each processor processes different subset of data and store data in the shared memory. The asynchronous implementation is used to allow each processor write data at any time without waiting others to release the lock. The noise that some processors overwrite data is negligible. However, shared memory parallel computer is too expensive.
- ii. Parameter parallel processing
This method is about parallelizing across the parameters of the output units. Each CPU computes the un-normalized probability for a subset of output data and updates the corresponding parameters. In parallelized stochastic gradient ascent, CPUs need to communicate the normalization factor of the output softmax and the gradients on the hidden layer and word feature layer. Still, there will be some computation duplicates among CPUs which are negligible.

The computation for processor is summarized as follows.

- i. Forward phase
 - a. Forward computation for the word features layer
 - b. Forward computation for the hidden layer
 - c. Forward computation for output units in i-th block
 - d. Computation and sharing the sum of the probabilities for output units among processors using MPI library
 - e. Probability normalization
 - f. Log-likelihood update

- ii. Backward/update phase
 - a. Backward gradient computation for output units in i-th block
 - b. Summing and sharing $\frac{\partial L}{\partial x}$ and $\frac{\partial L}{\partial a}$ across processors using MPI library
 - c. Hidden layer weights back-propagation and update.
 - d. Word feature vectors update for the input words

In addition, the mini-batch method is mentioned in the paper. By performing communications every K examples, it can reduce communication latency and save time.

(4) Result

The performance of neural network model is compared to n-gram models based on Brown corpus and AP News corpus using perplexity.

- i. Neural network can take advantage of more context
- ii. Hidden units are useful
- iii. Mixing output probabilities of neural network with interpolated trigram is of good performance

(5) Extension and future work

- i. An energy minimization network

In this model, output words also are mapped as its feature vector. The output of the model is as following:

$$E(w_{t-n+1}, \dots, w_t) = v \cdot \tanh(d + Hx) + \sum_{i=0}^{n-1} b_{w_{t-i}}$$

where b is the vector of biases, d is the vector of hidden units biases, v is the output weight vector and H is the hidden layer weight matrix. And that input and output words contribute to x :

$$x = (C(w_t), C(w_{t-1}), C(w_{t-2}), \dots, C(w_{t-n+1}))$$

And the conditional probability is computed as follows:

$$\hat{P}(w_t|w_{t-1}, \dots, w_{t-n+1}) = \frac{e^{-E(w_{t-n+1}, \dots, w_t)}}{\sum_i e^{-E(w_{t-n+1}, \dots, w_{t-1}, i)}}$$

In addition, this model can handle OOV easily than previous model. It takes weighted convex combination of feature vectors of other words with the weights proportional to their conditional probability. And the new word can be incorporated in V and re-compute the probability distribution in the new set.

- ii. Decomposing the network in sub-network to train the model more efficiently
- iii. Representing conditional probability using tree structure to reduce computation time by a factor of $\frac{|V|}{\log |V|}$
- iv. Propagating gradients only from a subset of the output words by selecting important words to accelerate training
- v. Introducing a-priori knowledge including semantic and grammar information
- vi. Interpreting the word feature representation learned by neural network model
- vii. Associating each word with multiple points in word feature space to handle polysemous words.

(6) Conclusion

The neural network model proposed in this paper uses learned distributed representation to inform model about the combinatorial number of other sentences. In addition, this smooth and compact representation can accommodate more conditioning variables to improve statistical language model.

2 Chinese lyrics decoding using trigram language model

1. Introduction

When I listen to the Chinese songs, I often misunderstand the exactly words in the lyrics because the tones of the words cannot be recognized in the songs. Human can figure out the possible lyrics using our language knowledge. How about the computer? While taking DSP course, I constantly think about how to solve this problem using computer program. I think this is a practical issue in life so I decide to handle this problem in the final project.

This computer program aims to decode purely ZhuYin lyrics sequences into Chinese characters using language model.

Input: 为 | 尤虫为么厂×为 | 尤虫为么厂×々么勿さ弓×所

Output: 兩隻老虎兩隻老虎跑得快

It contains two main steps:

1. Language model training

Train the trigram language model using srilm library [2] and lyric corpus collected from website [3]. While constructing the corpus, compute the perplexity to see the improvement of performance of language model.

2. ZhuYin lyrics sequences decoding

Using language model to decode the sequences made of solely ZhuYin. Compare the results with using trigram language model trained by corpus in hw3.

2. Execution

The commands corresponding to each task are as following.

1. Language model training

- Collect new lyrics data in “lyrics_utf.txt”
- Type “bash train_corpus.sh” to expand lyrics corpus and train trigram language model
- Type “make run_perplexity” to compute perplexity of the language model by running calculate_perplexity.cpp

2. ZhuYin lyrics sequences decoding

- Put the testing Chinese lyrics in directory `finaldata/` and name it `test_ori.txt`.
- Type “`bash convert_test_zhuyin.sh`” to convert `test ori.txt` to purely Zhuyin sequences `test zhuyin.txt`.

- c. Type “make run” to produce result file in
result_final/run/run_test.txt
- d. Type “make run_compare” to produce result file using
trigram language model trained by hw3 corpus in
result_final/run_compare/run_compare_test.txt

4. Details

I collect lyrics data of Chinese pop songs from website [3] and use srilm library [2] to train trigram language model. During data collection and training, I examine the change of perplexity of language model with respect to testing data. The sentences in testing data are completely different from the training corpus.

1. Data collection

Since that some characteristics of Chinese lyrics are different from the normal documents such as news reports or articles, it is necessary to use lyrics data set to train language model. There are some traits of Chinese lyrics which are different from normal documents:

- a. Repetition

Ex. 兩隻老虎兩隻老虎 跑得快跑得快

- b. Inversion

Ex. 當有你的溫熱我想我很快樂

- c. Rhyme

Ex. 一閃一閃亮晶晶 滿天都是小星星 掛在天空放光明

Therefore, using the data set of Chinese lyrics, the language model can “learn” various characteristics of this kind of data. I collect lyrics data of several Chinese pop music singers’ songs from website [3] to construct lyrics corpus.

Besides, it is important to handle big5 format correctly when collecting data. I convert the original format when copy data from website (utf-8) into big5 by the following command in train_corpus.sh.

```
cat corpus_utf_train.txt | iconv -c -f UTF-8 -t BIG5 > big_corpus_train_big5.txt
```

2. Perplexity computation

At first, I refer to class materials for the computation of perplexity and I use the total number of words in the testing corpus for the N_D in the formula. However, I find that the result of perplexity is strange. So I search for some information about the computation of perplexity. And I found that N_D is related to the number of counts of the n-grams.

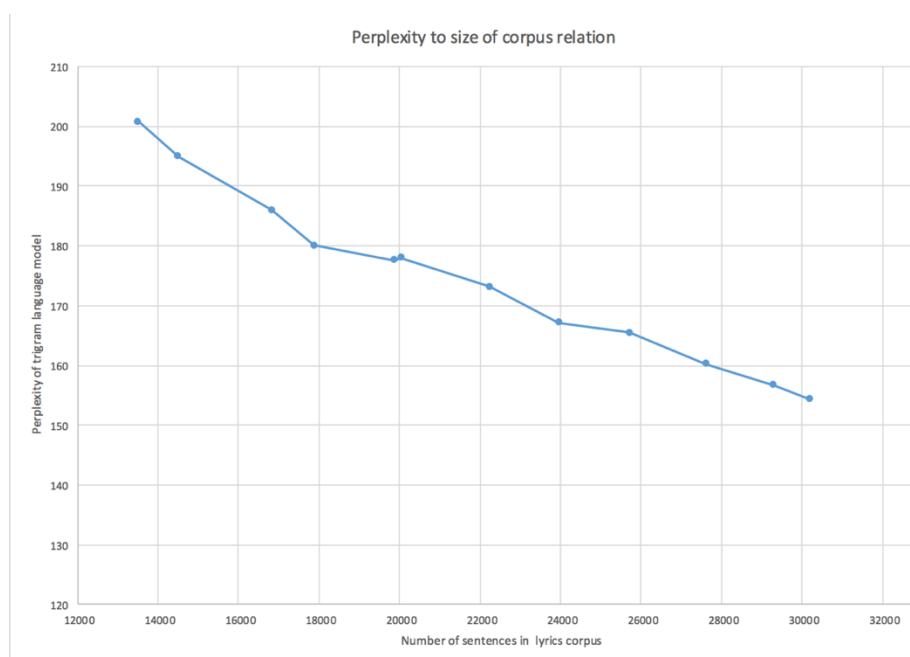
To compute the perplexity and verify performance of language model, I collect some random lyrics without repeating with the data in the corpus to construct testing corpus which contains 733 lines of lyrics.

The following chart shows part of the process of collecting different singer's lyrics set from website [3] and the trend that perplexity decreases while constructing the lyrics corpus. For each singer, I collect all of their songs they published before. It shows that the performance respective to testing corpus is gradually improved. In “吳亦凡” ’ s case, though, adding lyrics of 吳亦凡’ s songs results in the increase of perplexity. It might because that lyrics data in this case is not useful for the testing corpus. The lyrics corpus after removing empty lines are saved in “corpus_utf.txt” .

	Perplexity	Number of sentences	Average reduced perplexity per sentence
The corpus before adding the following singers' lyrics	200.808	13510	
1. 蔡健雅	194.92	14507	-0.00590572
2. 王心凌	186.011	16828	-0.00383843
3. 李宗盛	180.097	17882	-0.00561101
4. 方大同	177.639	19865	-0.0012395

5. 吳亦凡	178.018	20050	0.0020487
6. 蕭敬騰	173.14	22241	-0.0022264
7. 許茹芸	167.175	23957	-0.0034761
8. 張震嶽	165.481	25714	-0.0009641
9. 謝和弦	160.232	27622	-0.0027511
10. 李榮浩	156.687	29298	-0.0021152
11. 田馥甄	154.322	30211	-0.0025904

The following graph shows the degree of reducing perplexity of each singer's data. The data of “李宗盛” has the greatest performance of reducing perplexity. It might be due to that his songs seldom repeat lyrics compared to other singers. Or that his songs have more variable words and structures of sentences to improve language model.



5. Results

1. Children song lyrics test data

I train the trigram language model using the corpus in hw3 to compare the results of small test data containing some children song lyrics. Since the corpus in hw3 contains mainly news documents, the performance of resulting trigram_ref model is worse than the trigram_train which use lyrics corpus to train. As a result, data collection can extremely influence the performance of certain tasks.

answer 兩隻老虎兩隻老虎 跑得快跑得快 一閃一閃亮晶晶 滿天都是小星星 掛在天上放光明 好像許多小眼睛 妹妹背著洋娃娃 走到花園來看花 娃娃哭著叫媽媽	result of trigram_train	<s> 兩隻老虎兩隻老虎 </s> <s> 跑的快跑的快 </s> <s> 一閃一閃亮晶晶 </s> <s> 滿天都是小星星 </s> <s> 掛在天上放光明 </s> <s> 好像許多小眼睛 </s> <s> 沒沒背著洋娃娃 </s> <s> 走到花園來看花 </s> <s> 娃娃哭著叫媽媽 </s>
	result of trigram_ref	<s> 兩隻老虎兩隻老虎 </s> <s> 跑得快的快 </s> <s> 以善意閃亮晶晶 </s> <s> 滿天都是小行星 </s> <s> 掛在天上放光明 </s> <s> 好像許多小眼睛 </s> <s> 妹妹被這樣娃娃 </s> <s> 走到花園來看花 </s> <s> 娃娃哭著叫媽媽 </s>

2. Pop song lyrics test data

After converting Chinese lyrics data into purely ZhuYin sequences, I compare the results between two trigram language model using famous pop songs. Some songs are new publications while some are the old singers' works. Generally, the performance of trigram language model trained by lyrics corpus to decode ZhuYin sequences is much better.

	accuracy of trigram_train	accuracy trigram_ref	accuracy improvement
周杰倫 - 不愛我就拉倒	0.835735	0.723343	0.155379
周杰倫 - 等你下課	0.887701	0.794118	0.117845
周杰倫 - 千里香	0.679104	0.597015	0.137499
陳綺貞 - 沙發海	0.845411	0.763285	0.107595
陳綺貞 - 變色龍	0.915966	0.840336	0.089999
陳綺貞 - 蜉蝣	0.954198	0.732824	0.302084

