

# A Note on an Approximate Learning Algorithm with Limited Parameters for Language Models

Yasunari Maeda, Fumito Masui, Hideki Yoshida and Masakiyo Suzuki

Kitami Institute of Technology, 165 Koen-cho, Kitami-shi, Hokkaido 090-8507 Japan

## Summary

A lot of research in the field of NLP(natural language processing) for AI(artificial intelligence) has the goal of learning language models. In general, the aim is to minimize the divergence between the approximate model and the true model, but most learning algorithms are based on the maximum likelihood method. The existence of finite samples with high likelihood doesn't mean that the divergence between the approximate model and the true model is small. This paper proposes a new learning algorithm, the measure of which is divergence. The proposed algorithm is compared to previous algorithms using simulations.

## Key words:

language model, approximate learning algorithm, Markov chain, divergence

## 1. Introduction

Many researchers are now tackling language models in the field of NLP(natural language processing) for AI(artificial intelligence). Most papers used n-gram models as the language model. n-gram models can be represented as  $(n-1)$ -th order Markov chains, so training an n-gram model can be regarded as training a Markov chain with unknown parameters.

Some researchers have tried training Markov chains with an extremely large number of parameters; since the number of parameters of the true model is unknown, overcoverage is common. This requires basically unlimited memory for the parameters and so is rather impractical. Another research approach, which assumes that the number of parameters of the true model is known, is to train the Markov chain with enough parameters. Another research approach, which also assumes that the number of parameters of the true model is known, is to first train the Markov chain with enough parameters, and then develop approximate models with limited parameters to allow their implementation on computers that have limited memory. This paper adopts this last approach but introduces a better learning method.

Though the purpose of the approximate learning algorithm with limited parameters is to minimize the divergence between the approximate model and the true model, Brown's algorithm[1] uses the maximum

likelihood method. Its performance given finite samples is weak, since a high likelihood doesn't mean that the divergence is small.

This paper introduces an approximate learning algorithm in which divergence is used as a measure. I show that the algorithm can develop an approximate model that is closer to the true model than those yielded by Brown's algorithm.

I describe Markov chains in section 2 and previous research in section 3. The proposed algorithm is introduced in section 4. It is compared to Brown's algorithm in section 5. Section 6 concludes the paper.

## 2. Markov Chains

Almost all previous research studied n-gram models. These are language models that can be represented by  $(n-1)$ -th order Markov chains as follows:

$$p(w_i | w_0 w_1 \cdots w_{i-1}, \theta) = p(w_i | w_{i-n+1} w_{i-n+2} \cdots w_{i-1}, \theta), \quad (1)$$

where  $w_i$  is a word,  $w_i \in W$ ,  $W$  is a finite word set,  $\theta$  is a  $|W|^{n-1}(|W|-1)$ -dimensional vector of real-valued parameters,  $\theta \in \Theta$ ,  $\Theta$  is the set of parameters. This paper also adopts the Markov chain approach.

## 3. Previous Research

Many previous papers addressed the training of  $(n-1)$ -th order Markov chains. Prior work is described below.

### 3.1 Bayes coding

When the order of the true Markov model  $(n-1)$  is known, the true parameter vector  $\theta^*$ ,  $\theta^* \in \Theta$  is unknown and using  $|W|^{n-1}(|W|-1)$ -dimensional vector of real-valued parameters, divergence is minimized, in terms of the Bayes criterion, by Bayes coding[2]. The divergence is as follows:

$$D(p(\cdot|\cdot, \theta^*) \| p(\cdot|\cdot, \theta)) = \sum_{w^{n-1}} \pi(w^{n-1}, \theta^*) \sum_w p(w|w^{n-1}, \theta^*) \log \frac{p(w|w^{n-1}, \theta^*)}{p(w|w^{n-1}, \theta)}, \quad (2)$$

where  $\pi(w^{n-1}, \theta^*)$  is the stationary state probability of the true model. The approximate model in Bayes coding is shown as follows:

$$p_{Bayes}(w|w^{n-1}, \hat{\theta}) = \frac{h(w|w^{n-1}, x^M) + \beta(w|w^{n-1})}{\sum_{w_i} h(w_i|w^{n-1}, x^M) + \beta(w_i|w^{n-1})}, \quad (3)$$

where  $x^M$  is a string of length  $M$  for learning,  $x \in W$ ,  $h(w|w^{n-1}, x^M)$  is the number of times  $w$  appears next to  $w^{n-1}$  in  $x^M$ , and  $\beta(w|w^{n-1})$  is a parameter of Dirichlet distribution for the prior distribution of  $\theta$ .

### 3.2 Brown's algorithm

There is a strong need to reduce the dimension of the parameter vector if we are to implement an approximate model on a computer with limited memory. I now describe Brown's algorithm for the case of 2-gram models. At first, an approximate model with a  $|W|(|W|-1)$ -dimensional parameter vector is learned based on the maximum likelihood method as follows:

$$p_{ml}(w_j|w_i, \hat{\theta}) = \frac{h(w_j|w_i, x^M)}{\sum_{w_j} h(w_j|w_i, x^M)}. \quad (4)$$

Next, the dimension of the approximate model is reduced as follows:

$$\hat{p}_{Brown}(w_i|w_{i-n+1} \cdots w_{i-1}) = p(w_i|c_i) p(c_i|c_{i-n+1} \cdots c_{i-1}), \quad (5)$$

where  $c_i$  is a subset of  $W$ ,  $c_i \cap c_j = \emptyset$ ,  $i \neq j$ ,  $\bigcup_i c_i = W$ ,

$C = \{c_0, c_1, \dots, c_{|C|-1}\}$ ,  $C$  is a partition of  $W$ ,  $|\cdot|$  is the cardinality of a set. The cardinality of  $C$  is reduced to  $|C|-1$  by merging  $c_i$  and  $c_j$ . The initial  $C$  is  $W$ ,

$C = \{c_0, c_1, \dots, c_{|W|-1}\} = \{w_0, w_1, \dots, w_{|W|-1}\}$ . The pair of  $(c_i, c_j)$  to be merged is determined as follows:

$$(c_i, c_j) = \arg \min_{c_i, c_j} (I(C; C) - I(C'(c_i, c_j); C'(c_i, c_j))), \quad (6)$$

where  $C'(c_i, c_j)$  is the new partition of  $W$  created by merging  $c_i$  and  $c_j$ ,  $I(C; C)$  is the mutual information of  $C$  and  $C$  as follows:

$$I(C; C) = \sum_{c_k, c_l} p(c_k c_l) \log \frac{p(c_k c_l)}{p(c_k)}, \quad (7)$$

where each probability distribution in the right hand side of formula (7) is calculated by the maximum likelihood

method by formula (4). When the merging step is repeated  $T$  times, the dimension of the parameter vector is reduced to  $(|W|-T)(2|W|-T-2)$ . Formula (6) is also based on the maximum likelihood method.

Unfortunately, the maximum likelihood method doesn't offer good performance if the sample number is finite. A high likelihood doesn't mean that the divergence between the approximate model and the true model is small.

## 4. Proposed Algorithm

This section proposes a new approximate learning algorithm for  $(n-1)$ -th order Markov chains. This algorithm uses divergence as a measure of fitness. At first, an approximate model with  $|W|^{n-1}(|W|-1)$ -dimensional parameter vector is calculated by formula (3), which minimizes the divergence between the approximate model and the true model using the Bayes criterion. Next, the dimension of the approximate model is reduced as follows:

$$\hat{p}_{Proposed}(w_i|w_{i-n+1} \cdots w_{i-1}) = p(w_i|s_i), \quad w_{i-n+1} \cdots w_{i-1} \in s_{i-1}, \quad (8)$$

where  $s_i$  is a subset of  $W^{n-1}$ ,  $s_i \cap s_j = \emptyset$ ,  $i \neq j$ ,  $\bigcup_i s_i = W^{n-1}$ ,

$S = \{s_0, s_1, \dots, s_{|S|-1}\}$ ,  $S$  is a partition of  $W^{n-1}$ . The cardinality of  $S$  is reduced to  $|S|-1$  by merging  $s_i$  and  $s_j$ .

The initial  $S$  is  $W^{n-1}$ ,

$$S = \{s_0, s_1, \dots, s_{|W|^{n-1}-1}\} = \{w_0 \cdots w_0, w_0 \cdots w_0 w_1, \dots, w_{|W|-1} \cdots w_{|W|-1} w_{|W|-1}\}.$$

The pair of  $(s_i, s_j)$  to be merged is determined as follows:

$$\begin{aligned} (s_i, s_j) &= \arg \min_{s_i, s_j} D(p(\cdot|\cdot, S) \| p(\cdot|\cdot, S'(s_i, s_j))) \\ &= \pi(s_i, S) \sum_{w_k} p(w_k|s_i, S) \log \frac{p(w_k|s_i, S)}{p(w_k|s_i \cup s_j, S'(s_i, s_j))} \\ &\quad + \pi(s_j, S) \sum_{w_k} p(w_k|s_j, S) \log \frac{p(w_k|s_j, S)}{p(w_k|s_i \cup s_j, S'(s_i, s_j))}, \end{aligned} \quad (9)$$

where  $S'(s_i, s_j)$  is the new partition of  $W^{n-1}$  created by merging  $s_i$  and  $s_j$ ,  $\pi(s_i, S)$  is the stationary state probability of the approximate model based on partition  $S$ , the initial  $\pi(s_i, S)$  is calculated based on the model of formula (3), the initial  $p(w_k|s_i, S)$  is equal to

$p_{Bayes}(w_k|w^{n-1}, \hat{\theta})$ , where  $w^{n-1} = s_i$ , calculated by formula

(3),  $p(w_k | s_i \cup s_j, S'(s_i, s_j))$  is calculated as follows:

$$p(w_k | s_i \cup s_j, S'(s_i, s_j)) = \frac{\pi(s_i, S)}{\pi(s_i, S) + \pi(s_j, S)} p(w_k | s_i, S) + \frac{\pi(s_j, S)}{\pi(s_i, S) + \pi(s_j, S)} p(w_k | s_j, S), \tag{10}$$

$$p(w_k | s_l, S'(s_i, s_j)) = p(w_k | s_l, S), \quad l \neq i, \quad l \neq j,$$

$\pi(s_i \cup s_j, S'(s_i, s_j))$  is calculated as follows:

$$\pi(s_i \cup s_j, S'(s_i, s_j)) = \pi(s_i, S) + \pi(s_j, S), \tag{11}$$

$$\pi(s_l, S'(s_i, s_j)) = \pi(s_l, S), \quad l \neq i, \quad l \neq j.$$

The increase in divergence caused by merging is minimized by formula (9). When merging is repeated  $T$  times, the dimension of the parameter vector is reduced to  $(|W|^{n-1} - T)(|W| - 1)$ .

Divergence is used as the convergence measure in the proposed algorithm. The idea of this algorithm is based on a part of Nomura's research[3]. The next section uses simulations to show that the approximate model calculated by the proposed algorithm offers better convergence to the true model than the approximate model calculated by Brown's algorithm.

### 5. Simulations

In this section I show a result of simulating the 2-gram model shown in Fig.1(at the end of this paper). Let  $|W|$  be equal to 10. The true model is as follows:

$$p(w_j | w_i, \theta^*) = \begin{pmatrix} .24556 & .00002 & .00102 & .46333 & .25257 & .01985 & .01691 & .00065 & .00007 & .00002 \\ .78176 & .01112 & .00885 & .10790 & .00157 & .08864 & .00007 & .00001 & .00007 & .00001 \\ .28659 & .05758 & .00001 & .00040 & .64797 & .00003 & .00542 & .00000 & .00197 & .00004 \\ .02189 & .22285 & .59295 & .00114 & .00032 & .00000 & .00023 & .16062 & .00000 & .00000 \\ .02851 & .03988 & .04005 & .04520 & .11245 & .08491 & .02554 & .36051 & .22218 & .04077 \\ .00051 & .05621 & .00621 & .35798 & .27244 & .02884 & .16469 & .10585 & .00050 & .00677 \\ .30394 & .03310 & .03919 & .00011 & .00010 & .00012 & .04755 & .57360 & .00226 & .00003 \\ .02786 & .30086 & .00016 & .00110 & .39047 & .25020 & .00246 & .00068 & .02603 & .00018 \\ .00498 & .00181 & .69491 & .11962 & .00284 & .00400 & .01403 & .09930 & .05549 & .00302 \\ .01825 & .45628 & .51917 & .00038 & .00056 & .00004 & .00508 & .00007 & .00006 & .00011 \end{pmatrix}, \tag{12}$$

where the value of an element located on the  $i$ -th row,  $j$ -th column is  $p(w_j | w_i, \theta^*)$ .

Divergence is calculated between the true model and each approximate model. Brown65 means an approximate model with a 65-dimensional parameter vector calculated by Brown's algorithm, Brown48 means an approximate model with a 48-dimensional parameter vector calculated by Brown's algorithm, and pro45 means an approximate model with a 45-dimensional parameter vector calculated by the proposed algorithm. Fig.1 shows the average results

of 100 runs. A uniform distribution was used as the prior distribution of Bayes coding.

The approximate model yielded by the proposed algorithm diverges less from the true model than those of Brown's algorithm, though the dimension of pro45 is less than those of Brown65 and Brown48. Similar results were gained in other simulations.

### 6. Conclusion

There are many learning algorithms for language models with limited parameters. Brown's algorithm uses likelihood as a measure, but likelihood is not appropriate for obtaining approximate models that are close to the true model. Because the goal is to minimize the divergence between the approximate model and the true model, divergence should be used instead of likelihood.

This paper proposed a new algorithm that uses divergence as a measure of fitness. Simulations showed that the approximate models calculated by the algorithm are closer to the true models than those calculated by Brown's algorithm.

Since the algorithm is a kind of greedy algorithm, the approximate models yielded by the algorithm are not the best in terms of any one criterion. Further work will be to create some theoretical guarantee of the algorithm.

Learning language models can be regarded as the clustering of words. Therefore, the proposed algorithm may be useful as a clustering algorithm.

### References

- [1] Peter F. Brown, Vincent J. Della Pietra, "Class-Based n-gram Models of Natural Language", Computational Linguistics, Vol.18, No.4, pp.467-479, 1992.
- [2] Matsushima, Hirasawa, "On universal codes for FSMX sources", The 18<sup>th</sup> Symposium on Information Theory and Its Applications (SITA95), pp.377-380, 1995.
- [3] Nomura, Matsushima, Hirasawa, "On Bayes Coding in the case of limited memories", TECHNICAL REPORT OF IEICE, IT96-15, pp.31-36, 1996.

