

Dynamic Language Model Adaptation Using Latent Topical Information and Automatic Transcripts

Berlin Chen

Graduate Institute of Computer Science & Information Engineering,
National Taiwan Normal University, Taipei, Taiwan
berlin@csie.ntnu.edu.tw

Abstract

This paper considers dynamic language model adaptation for Mandarin broadcast news recognition. Both contemporary newswire texts and in-domain automatic transcripts were exploited in language model adaptation. A topical mixture model was presented to dynamically explore the long-span latent topical information for language model adaptation. The underlying characteristics and different kinds of model structures were extensively investigated, while their performance was analyzed and verified by comparison with the conventional MAP-based adaptation approaches, which are devoted to extracting the short-span n -gram information. The fusion of global topical and local contextual information was investigated as well. The speech recognition experiments were conducted on the broadcast news collected in Taiwan. Very promising results in perplexity as well as character error rate reductions were initially obtained.

1. Introduction

Statistical language modeling, which aims to capture the regularities in human natural language and quantify the acceptance of a given word sequence, has continuously been a focus of active research in speech and language processing over the past three decades. The n -gram modeling (especially the bigram and trigram modeling) approach, which determines the probability of a word given the previous $n-1$ word history, is most prominently used [1][2]. The n -gram probabilities are normally estimated based on the maximum likelihood (ML) principle. However, in order to tackle the inevitable data sparseness problems when estimating the n -gram probabilities from a specific text corpus, a variety of smoothing or interpolation techniques have been proposed in the past several years [3]. Meanwhile, statistical language modeling also was introduced to information retrieval (IR) problems in the late 1990s, and research at a number of sites has confirmed that such a statistical modeling paradigm does provide an effective and theoretically attractive probabilistic framework for building IR systems [4].

On the other hand, with the rapid growth of accessible multimedia information over the Internet, large volumes of real-world speech information, such as broadcast radio and television programs, digital libraries and so on, are now being accumulated and made available to the public. Substantial efforts and very encouraging results on transcribing broadcast news speech have been reported, e.g., [5]. However, for complicated speech recognition tasks such as broadcast news transcription, it is still extremely difficult to build well-estimated language models, because the subject matters and

lexical characteristics for the linguistic contents of news articles are very diverse and are often changing with time. Various attempts have been made to adapt the language model by making use of either the contemporary corpus or the recognition hypotheses observed so far [6]. Two of the most widespread approaches to language model adaptation are count merging and model interpolation, which can be respectively viewed as a maximum *a posteriori* (MAP) language model adaptation with a different parameterization of the prior distribution and can be easily integrated into the n -gram modeling framework to capture the local regularities of word usage in the new task domain [7]. In contrast, the latent semantic analysis (LSI) approach originally formulated for relevance measures in various IR tasks also has been proposed to explore the latent topical factors for language model adaptation, e.g. [2]. LSI transforms the high-dimensional vector representations of a word and a document (or a search history) into a lower dimensional space (the so-called latent semantic space). The relevance measure can be estimated in the reduced space and then be transformed into an approximate probability measure. Though LSI has been demonstrated effective in a variety of speech recognition tasks, however, its derivation is based on linear algebra operations and lacks for a probabilistic framework for automatic model refinement or optimization.

Based on these observations, in this paper a topical mixture model (TMM) previously presented for information retrieval is investigated to dynamically explore the long-span latent topical information for language model adaptation [8][9]. The speech recognition experiments were carried out on the broadcast news collected in Taiwan. Both contemporary newswire texts and in-domain automatic transcripts were exploited in language model adaptation. The TMM model was first trained beforehand on a set of contemporary text articles or in-domain automatic transcripts, and then can be gradually optimized when being applied to speech recognition. Structures similar to the presented approach have also been investigated recently, e.g., [10]-[12]. The main differences between the presented approach and the previous ones are that we explicitly interpret the document (or the search history) as a generative mixture model used to predict the newly occurring word, and both the perplexity and word error rate experiments are simultaneously investigated with very good potential indicated. Besides, various kinds of model structures are extensively tested and their performance is compared with the conventional MAP-based adaptation approaches, which are devoted to extracting the short-span n -gram information. The fusion of global topical and local contextual information was investigated as well.

2. The NTNU Broadcast News System

The major constituent parts of the broadcast news system developed at National Taiwan Normal University (NTNU) as well as the speech and language data used in this paper will be briefly described in this section [13].

2.1. Front-End Processing

The front-end processing is conducted with the data-driven LDA-based (Linear Discriminant Analysis) feature extraction approach. The states of each HMM were taken as the unit for class assignment. The outputs of 18 Mel-frequency filter banks are chosen as the basic vector. The basic vectors from every nine successive frames were spliced together to form the supervectors for the construction of the LDA transformation matrix, which was then used to project the supervectors to a lower feature space. The dimension of the resultant vectors was set to 39.

2.2. Speech Corpus and Acoustic Training

The speech data set consists of about 176 hours of FM radio broadcast news, which were collected from several radio stations located at Taipei during November 1998 to April 2004. All the speech materials were manually segmented into separate stories, and each of them is a news abstract pronounced by one anchor speaker. Only 7.7 hours of speech data is equipped with corresponding orthographic transcripts, in which about 4.0 hours of data collected during 1998 to 1999 is used to bootstrap the acoustic training and the other 3.7 hours of data (506 stories) collected in September 2002 is for testing. An amount of 104.3 hours of the rest untranscribed speech data (about 18,600 stories) is reserved for unsupervised acoustic model training and unsupervised language model adaptation.

2.3. Lexicon and N -gram Language Modeling

The recognition lexicon initially consists of 67K words. A set of about 5K compound words was automatically derived using the forward and backward bigram statistics and was then added to the lexicon to form a new lexicon of 72K words. The background language models used in this paper consist of trigram and bigram models, which were estimated using a text corpus consisting of 170 million Chinese characters collected from Central News Agency (CNA) in 2000 and 2001 (the Chinese Gigaword Corpus released by LDC). On the other hand, another text corpus of about 39,000 news articles (20 million Chinese characters) collected from CNA during August to October 2002 is taken as the contemporary corpus for language model adaptation. The n -gram language models were trained with Katz backoff smoothing using the SRI Language Modeling Toolkit (SRILM) [14].

2.4. Speech Recognition

The speech recognizer was implemented with a left-to-right frame-synchronous Viterbi tree search as well as a lexical prefix tree organization of the lexicon. At each speech frame, a beam pruning technique, which considered the decoding scores of path hypotheses together with their corresponding unigram language model look-ahead and syllable-level acoustic look-ahead scores, was used to select the most promising path hypotheses. Moreover, if the word hypotheses ending at each speech frame had scores higher than a

predefined threshold, their associated decoding information, such as the word start and end frames, the identities of current and predecessor words, and the acoustic score, will be kept in order to build a word graph for further language model rescore. Once the word graph had been built, the Viterbi beam search with a more sophisticated language model was conducted on it to generate the most likely word sequence. In the baseline system, the word bigram language model was used in the tree search procedure while the trigram language model in the word graph rescore procedure.

3. Topical Mixture Model (TMM)

In information retrieval, the relevance measure between a query Q and a document D_j can be expressed as $P(D_j|Q)$; i.e., the probability that the document D_j is relevant given that the query Q was posed. Based on Bayes' theorem and some assumptions, this measure can be approximated by $P(Q|D_j)$, which stands for the probability of the query Q being posed, under the hypothesis that document D_j is relevant [8]. The query Q is treated as a sequence of input observations (terms or words), $Q = w_1 w_2 \dots w_n \dots w_N$, where the query terms are assumed to be conditionally independent given the document D_j . Therefore, the relevance measure $P(Q|D_j)$ can be decomposed as a product of the probabilities of the query terms generated by the document:

$$P(Q|D_j) = \prod_{n=1}^N P(w_n|D_j). \quad (1)$$

Each individual document D_j can be interpreted as a generative mixture model, which is just a special case of HMM. In the model, a set of K latent topical distributions characterized with unigram language modeling are used to predict the query terms, and each of the latent topics is associated with a document-specific weight. That is, each document can belong to many topics. The relevance measure therefore can be further expressed as:

$$P(Q|D_j) = \prod_{n=1}^N \sum_{k=1}^K P(w_n|T_k) P(T_k|D_j), \quad (2)$$

where $P(w_n|T_k)$ denotes the probability of the query term w_n occurring in a specific latent topic T_k , and $P(T_k|D_j)$ is the posterior probability (or weight) of topic T_k conditioned on the document D_j , with the constraint $\sum_{k=1}^K P(T_k|D_j) = 1$ imposed. More precisely, the topical unigram distributions, e.g. $P(w_n|T_k)$, are tied among the entire document collection, while each document D_j has its own probability distribution over the latent topics, e.g. $P(T_k|D_j)$. The key idea we wish to illustrate here is that the relevance measure of a query term w_n and a document D_j is not computed directly based on the frequency of w_n occurring in D_j , but instead based on the frequency of w_n in the latent topic T_k as well as the likelihood that D_j generates the respective topic T_k , which in fact exhibits some sort of concept matching. During training, the K-means algorithm is first used to partition the entire document collection into K topical classes. Hence, the initial topical unigram distribution for a cluster topic can be estimated according to the underlying statistical

characteristics of the documents being assigned to it, and the probabilities for each document generating the topics are measured according to its proximity to the centroid of each respective cluster as well. Then, given a training set of query exemplars with the corresponding query-document relevance information, each document mixture model can be optimized in a supervised manner by the expectation-maximization (EM) algorithm [8].

While the TMM retrieval model is applied to language model adaptation, a set of contemporary (or in-domain) articles are first collected and used to train their corresponding mixture models. However, because there is no any query exemplar provided for the document model to be trained, we simply use each individual document in the collection as a query to train its own mixture model in an unsupervised manner. In speech recognition, for each newly occurring word w_i , its corresponding search history H_{w_i} can be treated as a document. The corresponding document TMM model can be optimized using the EM algorithm, throughout the whole search process. In this work, we keep the topic factors $P(w_i|T_k)$ unchanged, but let the search history's probability distribution over the latent topics, $P(T_k|H_{w_i})$, be gradually updated as path extension is performed during the search process. Once the TMM model for a search history is estimated, it can thus be used to predict the occurrence probability of the newly occurring word w_i (acting here as a single-word query):

$$P_{TMM}(w_i|H_{w_i}) = \sum_{k=1}^K P(w_i|T_k)P(T_k|H_{w_i}) \quad (3)$$

Such a kind of language model probability to some extent can dynamically capture the underlying global topical information of the path hypothesis and can be further combined with the background n -gram (e.g. trigram) language probability, which provides the general constraint information of lexical regularities, to form an adapted language model for guiding the search process:

$$\tilde{P}_{Adapt1}(w_i|w_{i-2}w_{i-1}) = \lambda \cdot P_{TMM}(w_i|H_{w_i}) + (1-\lambda) \cdot P_{Back}(w_i|w_{i-2}w_{i-1}), \quad (4)$$

where $P_{TMM}(w_i|H_{w_i})$ and $P_{Back}(w_i|w_{i-2}w_{i-1})$ respectively are the TMM probability and background trigram probability, and λ is a tunable weighting parameter.

4. Experimental Setup

In this paper, the language model adaptation experiments were performed in the word graph rescoring procedure, as described in Section 2.4. A set of 506 broadcast news stories collected in September 2002 is used for testing. For each broadcast news story to be processed, its associated word graph was built beforehand by the tree search procedure and using the background bigram language model. As mentioned earlier, a set of about 39,000 text news stories collected from CNA during August to October 2002 is taken as the contemporary adaptation data, while another set of about 18,600 automatic transcripts (3.2 million Chinese characters) as the in-domain adaptation corpus. These two sets of corpora are postulated to be either temporally or stylistically consistent with the broadcast news speech to be tested, and therefore can be used to explore the latent topical and local contextual information which might be helpful for speech recognition. For the TMM approach, its training and special utilization for language model adaptation have been described

earlier in Section 3. As for the two conventional MAP-based adaptation approaches to be used here for comparison [7], i.e. count merging and model interpolation, the adaptation formulae (e.g. for trigram modeling) can be respectively written as:

$$\tilde{P}_{Adapt2}(w_i|w_{i-2}w_{i-1}) = \frac{\alpha \cdot C_{d,Cont}(w_{i-2}w_{i-1}w_i) + \beta \cdot C_{d,Back}(w_{i-2}w_{i-1}w_i)}{\alpha \cdot C_{Cont}(w_{i-2}w_{i-1}) + \beta \cdot C_{Back}(w_{i-2}w_{i-1})}, \quad (5)$$

and

$$\tilde{P}_{Adapt3}(w_i|w_{i-2}w_{i-1}) = \gamma \cdot P_{Cont}(w_i|w_{i-2}w_{i-1}) + (1-\gamma) \cdot P_{Back}(w_i|w_{i-2}w_{i-1}). \quad (6)$$

For the count merging formula in (5), $C_{d,Cont}(w_{i-2}w_{i-1}w_i)$ and $C_{d,Back}(w_{i-2}w_{i-1}w_i)$ are respectively the discounted trigram counts [3] accumulated from the contemporary and background text corpora, $C_{Cont}(w_{i-2}w_{i-1})$ and $C_{Back}(w_{i-2}w_{i-1})$ are respectively the bigram counts accumulated from the contemporary and background text corpora as well, and α and β are tunable weighting parameters; while for the model interpolation formula in (6), $P_{Cont}(w_i|w_{i-2}w_{i-1})$ and $P_{Back}(w_i|w_{i-2}w_{i-1})$ are the trigram probabilities respectively estimated from the contemporary and background text corpora, and γ is a tunable weighting parameter. Language model adaptation using in-domain automatic transcripts can be expressed using the same formulae as well.

5. Experimental Results

The baseline result obtained by performing word graph rescoring with the background trigram language model alone is shown in Row 2 of Table 1. A character error rate (CER) of 15.22% and a perplexity (PP) of 752.49 were initially obtained. We first evaluate the performance level of the TMM adapted language model in (4) by varying the model complexities (the number of latent topics is ranged from 16 to 256) and using either the contemporary newswire texts (denoted as Texts) or in-domain automatic transcripts (denoted as Transcripts). The weighting parameter λ in (4) was initially set to 0.1 in this research. As can be seen from Table 1, both CER and PP are steadily reduced as the topic mixture number increases, when the contemporary texts were used for language model adaptation. A best result of CER of 14.47% (4.93% relative reduction) and PP of 457.74 (39.17% relative reduction) is obtained when the TMM topic number is set to 256. Though the performance seems not to be saturated yet, these results clearly demonstrate the effectiveness of the TMM approach for dynamic language model adaptation. Whereas, as the in-domain automatic transcripts were instead used for language model adaptation, a best result of CER of 14.82% (2.62% relative reduction) is achieved when the topic number is 64, and a best result of PP of 460.46 (38.81% relative reduction) is achieved when the topic number is 256. In summary, language model adaptation using the in-domain automatic transcripts is quite competitive with that using the contemporary newswire texts in PP reduction, but only reaches about half of CER reduction as that provided by using the contemporary newswire texts. The influence of the recognition errors of the automatic transcripts on language model adaptation is currently under extensive investigation.

Then, the experimental results as the count merging and model interpolation adaptation approaches are respectively applied are listed in Table 2 for comparison. The weighting

	CER (%)	PP
Baseline	15.22	752.49
+ TMM (16 Topics, Texts)	14.83	576.95
+ TMM (32 Topics, Texts)	14.73	555.93
+ TMM (64 Topics, Texts)	14.58	529.49
+ TMM (128 Topics, Texts)	14.53	492.40
+ TMM (256 Topics, Texts)	14.47	457.74
+ TMM (16 Topics, Transcripts)	14.99	577.46
+ TMM (32 Topics, Transcripts)	14.92	559.70
+ TMM (64 Topics, Transcripts)	14.82	534.98
+ TMM (128 Topics, Transcripts)	14.87	502.50
+ TMM (256 Topics, Transcripts)	14.92	460.46

Table 1: The baseline results and the results achieved by using the TMM-based adaptation approaches.

	CER (%)	PP
+ Count Merging (Texts)	13.70	458.79
+ Model Interpolation (Texts)	13.74	430.59
+ Count Merging (Transcripts)	14.87	635.03
+ Model Interpolation (Transcripts)	15.11	508.52

Table 2: The results achieved by two variants of the MAP-based adaptation approaches.

	CER (%)	PP
+ TMM (256 Topics, Texts)	13.27	306.75
+ Count Merging (Texts)		
+ TMM (256 Topics, Texts)	13.37	311.28
+ Model Interpolation (Texts)		
+ TMM (64 Topics, Transcripts)	14.65	477.74
+ Count Merging (Transcripts)		
+ TMM (64 Topics, Transcripts)	15.00	452.33
+ Model Interpolation (Transcripts)		
+ TMM (64 Topics, Texts & Transcripts)	13.35	319.93
+ Count Merging (Texts & Transcripts)		
+ TMM (256 Topics, Texts & Transcripts)	13.23	283.33
+ Count Merging (Texts & Transcripts)		

Table 3: The results achieved by combing the TMM-based and MAP-based adaptation approaches.

parameters, i.e. α and β in (5) and γ in (6), are set at optimum values for each experimental condition. As can be seen, these two approaches are quite comparable to each other. Count merging is slightly better than model interpolation in CER reduction, while model interpolation is better in PP reduction. If we further compare the TMM-based approach with the MAP-based ones, it can be found that the TMM approach is competitive with the MAP approaches in PP reduction, but only reaches half of CER reduction as that provided by the MAP approaches when the contemporary texts were used for language model adaptation, which also implies that the local word regularity information inherent in the adaptation corpus is still vital for speech recognition and should be taken into account when performing language model adaptation.

Finally, we combine the TMM-based approach with the MAP-based approaches in an attempt to explore both the subject domain and word regularity information for language model adaptation. The contemporary newswire texts and in-domain automatic transcripts, as well as their combination, were exploited. The results are shown in Table 3. As can be observed, the fusion of these two kinds of information sources does provide additional gains. For example, the combination

of the TMM approach (with 256 topics) with the count merging approach, as shown in the last row of Table 3, achieves the best result of CER of 13.23% (13.07% relative reduction) and PP of 283.33 (62.35% relative reduction).

6. Conclusions

In this paper we have presented a topical mixture model and made use of the automatic transcripts for dynamic language model adaptation. The underlying characteristics and different kinds of the model structures were extensively investigated and tested. We compared it with two conventional MAP-based approaches. The fusion of global topical and local contextual information has been investigated as well. Very promising results in both perplexity and character error rate reductions were initially obtained. More in-depth investigation and analysis of the TMM-based approach as well as its possible application to spoken document understanding and organization are currently undertaken [15].

7. References

- [1] R. Rosenfeld, "Two Decades of Statistical Language Modeling: Where Do We Go from Here," *Proc. IEEE*, 88 (8), 2000.
- [2] J. R. Bellegarda, "Statistical Language Model Adaptation: Review and Perspectives," *Speech Communication* 42, 2004.
- [3] S. F. Chen, J. Goodman, "An Empirical Study of Smoothing Techniques for Language Modeling," *Computer Speech and Language* 13, 1999.
- [4] W. B. Croft (editor), J. Lafferty (editor). *Language Modeling for Information Retrieval*. Kluwer-Academic Publishers, 2003.
- [5] P. Beyerlein et al., "Large Vocabulary Continuous Speech Recognition of Broadcast News – The Philips/RWTH approach," *Speech Communication* 37, 2002.
- [6] M. Federico, N. Bertoldi, "Broadcast News LM adaptation Using Cotemporary Texts," in *Proc. Eurospeech 2001*.
- [7] M. Bacchiani, B. Roark, "Unsupervised Language Model Adaptation" in *Proc. ICASSP 2003*.
- [8] B. Chen et al., "Statistical Chinese Spoken Document Retrieval Using Latent Topical Information," in *Proc. ICSLP 2004*.
- [9] B. Chen, W. H. Tsai, J. W. Kuo, "Statistical Language Model Adaptation for Mandarin Broadcast News Transcription," in *Proc. ISCSLP 2004*.
- [10] D. Gildea, T. Hoffmann, "Topic-based Language Models Using EM," in *Proc. Eurospeech 1999*.
- [11] S. Wang et al., "Semantic N-gram Language Modeling with the Latent Maximum Entropy Principle," in *Proc. ICASSP 2003*.
- [12] D. Mrva and P. C. Woodland, "A PLSA-based Language Model for Conversational Telephone Speech," in *Proc. ICSLP 2004*.
- [13] B. Chen et al., "Lightly Supervised and Data-Driven Approaches to Mandarin Broadcast News Transcription," in *Proc. ICASSP 2004*.
- [14] A. Stolcke, "SRI language Modeling Toolkit," version 1.3.3, <http://www.speech.sri.com/projects/srlm/>.
- [15] L. S. Lee and B. Chen, "Spoken Document Understanding and Organization for Efficient Retrieval/Browsing Applications," to appear in *IEEE Signal Processing Magazine*, 2005.