



# Open Research Online

---

The Open University's repository of research publications and other research outputs

## A latent variable model for query expansion using the hidden Markov model

Conference or Workshop Item

How to cite:

Huang, Qiang and Song, Dawei (2008). A latent variable model for query expansion using the hidden Markov model. In: ACM 17th Conference on Information and Knowledge Management (CIKM2008), 26-30 Oct 2008, Napa Valley, CA, USA.

For guidance on citations see [FAQs](#).

© 2008 The Authors

Version: Version of Record

Link(s) to article on publisher's website:

<http://dx.doi.org/doi:10.1145/1458082.1458310>

<http://dl.acm.org/citation.cfm?id=1458310>

---

Copyright and Moral Rights for the articles on this site are retained by the individual authors and/or other copyright owners. For more information on Open Research Online's data [policy](#) on reuse of materials please consult the policies page.

---

[oro.open.ac.uk](http://oro.open.ac.uk)

# A Latent Variable Model for Query Expansion Using the Hidden Markov Model

Qiang Huang  
Knowledge Media Institute  
The Open University  
Milton Keynes, UK  
q.huang@open.ac.uk

Dawei Song  
Knowledge Media Institute  
The Open University  
Milton Keynes, UK  
d.song@open.ac.uk

## ABSTRACT

We propose a novel probabilistic method based on the Hidden Markov Model (HMM) to learn the structure of a Latent Variable Model (LVM) for query language modeling. In the proposed LVM, the combinations of query terms are viewed as the latent variables and the segmented chunks from the feedback documents are used as the observations given these latent variables. Our extensive experiments shows that our method significantly outperforms a number of strong baselines in terms of both effectiveness and robustness.

## Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: Retrieval Models

## General Terms

Algorithm

## Keywords

Information retrieval, latent variable model, hidden Markov model

## 1. INTRODUCTION

In information retrieval, a prominent query language modeling approach is the Relevance Model (RM) [2]. Practically, variants of RM have shown encouraging performance in ad-hoc search [2, 4, 1] and topic detection and tracking [3], etc.

However, in these existing models, three issues, which in our belief will affect the retrieval effectiveness significantly, have not been addressed. Firstly, there often exist dependencies between query terms (i.e., intra-query term dependencies). Secondly, the distributions  $P(M_j)$ , i.e., the importance of different relevance feedback documents, should not be kept uniform. Instead, they should depend on the query terms and  $w$ . Thirdly, noise often exists within the feedback documents not any part of a relevant or pseudo-relevant document is necessarily relevant to the query.

The hypothesis of this paper is that tackling all the above identified three issues in a single query language modeling framework will further improve the retrieval effectiveness. In this paper, we first propose segmenting a document into

chunks by using an overlapped sliding window. We decompose the query into exhaustive combinations (subsets) of query terms and consider them as latent variables over the chunks. We then propose using a Latent Variable Model (LVM) which connects a chunk  $d$  and a word  $w$  through the latent variables. The dependencies between the latent variables are governed by an ergodic Hidden Markov Model (HMM), where the Viterbi algorithm is applied to optimize parameters involved in the HMM and the underlying LVM.

## 2. FRAMEWORK OF THEORY

In our model, the combinations of query terms are seen as the latent variables over the top-ranked documents. The reasons for selecting the query terms as the latent variables are that the word dependence to the query  $Q$  is to be estimated in the query language modeling framework, and the estimation of co-occurrence of words with query terms is key to query expansion. Eq. 1 shows an intuitive derivation of theory:

$$P(w|\theta_Q) = \sum_{S_j \in \mathbf{S}} P(w|S_j)P(S_j) \quad (1)$$

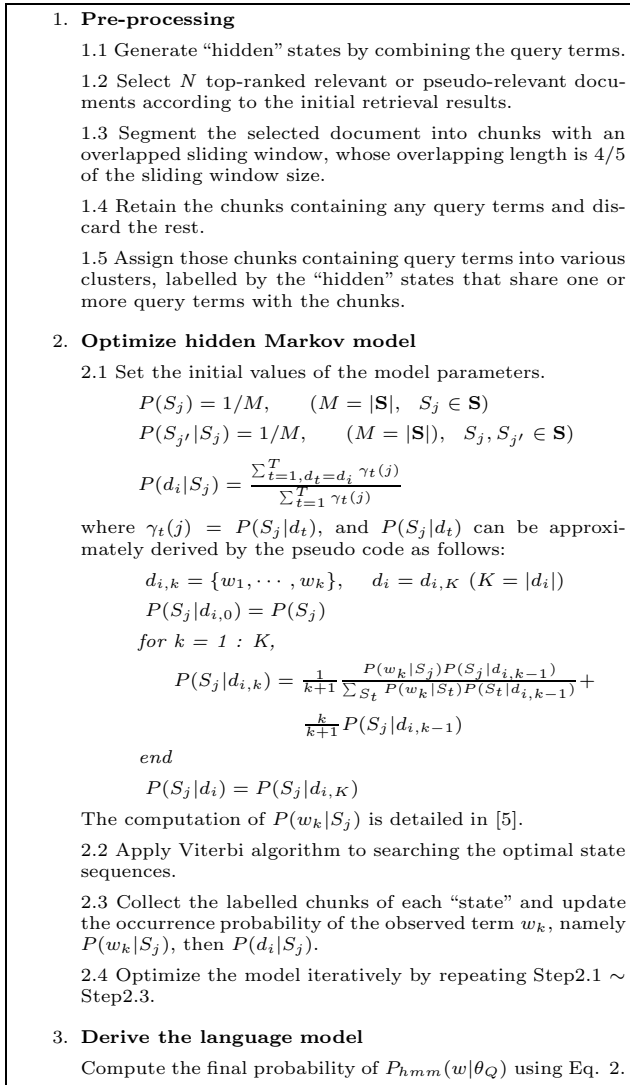
$S_j$  is a latent variable in the set  $\mathbf{S}$  ( $\mathbf{S} = \{S_1, \dots, S_M\}$ ) and  $w$  is the word whose occurrence probability in the expanded query model  $\theta_Q$  to be estimated. The latent variable  $S_j$  is generated from the query  $Q$ , where each  $S_j$  is defined as a combination of query terms. For example, given  $Q = \{q_1, q_2\}$ , the set of latent variables can be represented as  $\mathbf{S} = \{\{q_1\}, \{q_2\}, \{q_1, q_2\}\}$ . Compared with the traditional methods the use of all the combinations of query terms expands the observing space.

In Eq. 1, the relationship between the word  $w$  and the latent variable  $S_j$  is derived from the relevance feedback documents. As we have introduced in Section 1, the top-ranked documents in pseudo relevance feedback are not necessarily relevant to the query. Thus, we propose using the segmented chunks to connect  $S_j$  and  $w$ , and the relevance of each chunk to the query will also be considered. Then, a new equation is given as below:

$$P(w|\theta_Q) = \sum_{d_i \in \mathbf{d}, S_j \in \mathbf{S}} P(w|d_i, S_j)P(d_i|S_j)P(S_j) \quad (2)$$

$\mathbf{d}$  is the collection of chunks.  $P(S_j)$  is the prior distribution of latent variables,  $P(d_i|S_j)$  is the probability of an observed chunk  $d_i$  given a latent variable  $S_j$ , and  $P(w|d_i, S_j)$  is the probability of a word  $w$  in a chunk  $d_i$  given a specific latent variable  $S_j$ .

In order to calculate the three probability parameters in Eq. 2, we design a framework based on the Hidden Markov



**Figure 1: Outline of framework**

Model (HMM). The application of the HMM can not only estimate the prior distribution of each  $S_j$ , but also integrate the dependence between any two latent variables and their underlying observables through a state transition matrix. The details are presented in Figure 1.

Here, we regularize the model estimation with the original query model  $P(q_k|Q)$  to alleviate the data sparsity. In this paper, the original query model is defined as below:

$$P(q_k|Q) = \frac{\#q_k \cdot IDF(q_k)}{\sum_{j \in \{1 \dots |Q|\}} \#q_j \cdot IDF(q_j)} \quad (3)$$

where  $\#q_k$  is the frequency of query term  $q_k$  in  $Q$  and  $IDF(q_k)$  is the inverse document frequency (IDF) of  $q_k$ .

To mix the original query model into the newly generated one, two methods are used. **HMM-I** is an automatic method to integrate the original query model directly into the HMM, in which the original query model as a hidden state  $S_Q$ . We therefore can obtain Equation 4.

$$P_{HMM-I}(w|\theta_Q) = \sum_{d_i \in \mathbf{d}} \sum_{S_j \in \{\mathbf{S}, S_Q\}} P(w|d_i, S_j)P(d_i|S_j)P(S_j) \quad (4)$$

The second method, **HMM-II**, uses the linear interpolation to combine the original query model.

$$P_{HMM-II}(w|\theta_Q) = \lambda P(w|\theta_Q) + (1 - \lambda)P(w|Q) \quad (5)$$

The latter method involves manual adjustment of the interpolation parameter to generate the optimal retrieval performance.

### 3. EXPERIMENTS AND CONCLUSION

We evaluate our methods by testing TREC topics (Topics151–200, 601–700, and 501–550) on three large collections (AP88–90, ROBUST, and WT10G), respectively. we apply our approach to two scenarios: **pseudo-relevance feedback** and **true relevance feedback**. As a comparison, we list the values of MAP obtained by using HMM-II and HMM-I in Table 1.

**Table 1: Comparison of Optimal MAPs using HMM-II and HMM-I**

Pseudo-relevance Feedback					
Collection	KL	RM1	RM2	HMM-I	HMM-II
AP88–90	0.2077	0.2603	0.2676	0.2814 <sup>†</sup>	0.2830 <sup>†</sup>
ROBUST	0.2920	0.3129	0.3143	0.3613 <sup>†</sup>	0.3660 <sup>†</sup>
WT10G	0.2032	0.2131	0.2134	0.2305 <sup>†</sup>	0.2370 <sup>†</sup>
Relevance Feedback					
Collection	KL	RM1	RM2	HMM-I	HMM-II
AP88–90	0.2077	0.3218	0.3427	0.4034 <sup>†</sup>	0.4168 <sup>†</sup>
ROBUST	0.2920	0.4126	0.4432	0.5014 <sup>†</sup>	0.5122 <sup>†</sup>
WT10G	0.2023	0.2571	0.2993	0.3764 <sup>†</sup>	0.3859 <sup>†</sup>

<sup>†</sup> The improvement is statistically significant at the level of 0.05 according to the Wilcoxon signed rank test

Table 1 shows the optimal value when using HMM-II. It is found that using HMM-II can only obtain slightly and insignificantly better performances than HMM-I. However, HMM-I integrates original query model neatly in the Hidden Markov Model in a fully automatic way. Therefore, the HMM-I has demonstrated its robustness and effectiveness from both theoretical and practical perspectives.

In this paper, we present a novel method to build a latent variable model using the HMM for query expansion. This paper tries to address some specific issues on the dependencies between query terms, the different degree of relevance of different chunks in the feedback documents to the query, and noise existing within the feedback documents. Our technique incorporates these key issues in a single comprehensive framework and apply the HMM to estimating and optimizing the structure of the LVM. Our experimental results therefore show that our method always obtains significant improvements in comparison with KL, RM1 and RM2.

### 4. REFERENCES

- [1] J. Bai, D. Song, P. Bruza, J. Nie, and G. Cao. Query expansion using term relationships in language models for information retrieval. In *Proceedings of CIKM'2005*.
- [2] V. Laverenko and W. B. Croft. Relevance-based language models. In *Proceedings of SIGIR'2001*.
- [3] V. Lavrenko, J. Allan, E. DeGuzman, D. LaFlamme, V. Pollard, and S. Thomas. Relevance models for topic detection and tracking. In *Proceedings of HLT'2002*.
- [4] D. Song and P. Bruza. Towards context sensitive information inference. *JASIST*, 54(3):321–334, 2003.
- [5] D. Song, Q. Huang, S. Rueger, and P. Bruza. Facilitating query decomposition in query language modeling by association rule mining using multiple sliding windows. In *Proceedings of ECIR'2008*.