

A First Investigation on Mongolian Information Retrieval

Guanglai Gao, Wei Jin, Fei Long, Hongxu Hou

School of Computer Science

Inner Mongolia University

Hohhot China 010021

Email: {csggl, csjinwei, csfeilong, cshhx}@imu.edu.cn

Abstract

In this paper we present an attempt to build a test collection for Mongolian IR as well as some preliminary tests about the key issues in Mongolian Information Retrieval: using a stoplist and using word stemming. Our preliminary tests will show that while these basic operations on Mongolian can bring slight improvements in retrieval effectiveness, many problems remain.

The results using stemming and stoplist show that the stemming can potentially lead to some gain in retrieval effectiveness; The stoplist slightly improve retrieval effectiveness, but it can reduce the index significantly.

Keywords: Mongolian IR, Language Model, Test Collection, Mongolian Information Processing

1、 Introduction

Traditional Mongolian (for short “Mongolian”) is the main language of Inner Mongolia Autonomous Region in China. It has a long history and is a widely used language [1]. Despite the fact that Mongolian is a language spoken by millions of people, no IR study has been carried out on it until now.

To our knowledge, no test collection for Mongolian IR exists. The current state of Mongolian processing is far behind other languages. This is the very reason that motivated our study: to construct a test collection in Mongolian and to promote research activities in Mongolian IR.

Our study has two main purposes. First, we aim to construct a test collection following the TREC methodology. Such a test collection can benefit other researchers. Second, we aim to perform some preliminary tests on the test collection to see if the traditional IR methods are effective in Mongolian IR. In particular, we will test several basic options for IR: the utilization of a stoplist, and word stemming.

In the following sections, we will first describe the Mongolian language and the problems for Mongolian IR. Then we describe our effort to construct a test collection. Several IR methods will be tested on the test collection.

2、 The Mongolian Language and Mongolian IR

Mongolian is a phonemic script language. It has 35 letters. A Mongolian word is formed by several letters. In the written language, a word is written top-down by letters joined cursively together. Figure 1 shows a Mongolian sentence. We can see that words are separated by spaces.

Generally, a Mongolian word can be divided into two parts: stem and affix. Mongolian language is an agglutinative language. It does not have pre-affix and mid-affix, but only has suffix. More than one suffix can be added one after another [2]. For example, the word (bolbasuragul – mature) is formed by the stem and three suffixes. Schematically, we can illustrate the formation of the word as follows:

(bol) – ripe (verb)+ (basun) →

model, we have that:

$$P(Q|\theta_D) = \prod_{w_i \in Q} P(w_i|\theta_D)$$

5. Experimentations

The Indri system is used for the experimentation [4].

5.1 Experimented methods

Here is an example to our queries.
 Өрнөд улс олимпын тоглоомын тухай (The information about Beijing Olympic Games). In this query, Өрнөд (about) and тоглоом (inquire) are stop words. The terms Өрнөд улс (Pekinese), олимпын тоглоом (Olympic's), спорт (locomotor), тоглоомчид (pageant's) are translated into Өрнөд (Beijing), олимпын (Olympic), спорт (athletics) and тоглоомчид (pageant) if word stemming is used.

Smoothing [5] is a method used to overcome both the 'zero probability' and data sparseness problem. Three kinds of tests have been carried out to compare Mongolian IR with/without word stemming and with/without stoplist.

- Dirichlet smoothing

$$P(w_i|\theta_D) = \frac{tf_{w_i,D} + \mu P(w_i|\theta_C)}{|D| + \mu} \quad (1)$$

Where $\mu=2500$ is the Dirichlet prior, θ_C is a language model of the collection and $|D|$ is the total count of words in document D.

- Jelinek-Mercer smoothing

$$P(w_i|\theta_D) = \frac{(1-\lambda)tf_{w_i,D} + \lambda P(w_i|\theta_C)}{|D|} \quad (2)$$

Where $\lambda=0.4$ is a smoothing parameter.

- Two-stage smoothing

$$P(w_i|\theta_D) = \frac{(1-\lambda)(tf_{w_i,D} + \mu P(w_i|\theta_C))}{|D| + \mu} + \lambda P(w_i|\theta_C) \quad (3)$$

Where $\lambda=0.4$ and $\mu=2500$ are the smoothing parameters.

5.2 Experimental results

The experimental results with the three smoothing methods are described in the following

tables.

- Smoothing methods

Figure 5 shows a comparison between different smoothing methods for document models. All the methods illustrated in the figure use stemming and stoplist. We can see that different smoothing methods lead to slightly different retrieval effectiveness. However, globally, the effectiveness is comparable. We can conclude that all these smoothing methods that have been successfully applied to other languages also apply to Mongolian IR.

- Using stemming

From Tables 1, we can see that when stemming is used, the effectiveness is improved. And stemming allowed us to greatly reduce the size of the index. The number of unique indexes is reduced from 125,796 (without stemming) to 74,001 (with stemming). As a consequence, the total size of the index is also reduced by 20%.

- Using stoplist

We can observe that using stoplist, we can also obtain slightly effectiveness than without using stoplist.

6. Concluding remarks and future work

This paper describes our first attempt to Mongolian IR. In this paper, we investigated the following fundamental problems in IR: word stemming and utilization of a stoplist.

In order to test different retrieval methods, we constructed a small test collection. This is the first test collection for this language. Although it is still at its first stage, we have been able to perform some preliminary tests. These tests suggest the following conclusions:

- Word stemming in Mongolian IR is important. It can improve retrieval effectiveness.
- Using stoplist can slightly improve retrieval effectiveness, and it can reduce the index significantly.

	Dirichlet				Jelinek-Mercer				Two-Stage			
	No Stoplist No Stem	Stoplist No Stem	No Stoplist Stemming	Stoplist Stemming	No Stoplist No Stem	Stoplist No Stem	No Stoplist Stemming	Stoplist Stemming	No Stoplist No Stem	Stoplist No Stem	No Stoplist Stemming	Stoplist Stemming
0.0	0.80883	0.80883	0.83798	0.83084	0.83647	0.83647	0.89634	0.89634	0.80883	0.80883	0.84447	0.83798
0.1	0.81734	0.81189	0.82798	0.81957	0.79564	0.79564	0.80944	0.80422	0.81394	0.81038	0.82674	0.81647
0.2	0.80639	0.81167	0.82812	0.82300	0.74695	0.74695	0.76263	0.76263	0.80839	0.81334	0.83544	0.82994
0.3	0.78258	0.78547	0.80401	0.80363	0.71221	0.71221	0.73107	0.72839	0.79579	0.79143	0.79447	0.78843
0.4	0.69485	0.67688	0.73370	0.71963	0.67736	0.67736	0.71772	0.71593	0.68524	0.67067	0.73458	0.72597
0.5	0.61808	0.63787	0.66369	0.66584	0.60741	0.61489	0.65075	0.64300	0.59042	0.61291	0.64518	0.64376
0.6	0.53299	0.55743	0.58749	0.58731	0.52406	0.55099	0.54008	0.5642	0.50921	0.53320	0.56366	0.57028
0.7	0.39715	0.42766	0.47659	0.49094	0.42203	0.44711	0.45876	0.50105	0.35875	0.40444	0.43904	0.45464
0.8	0.27344	0.31301	0.26074	0.30571	0.25226	0.29133	0.27289	0.32652	0.26658	0.31192	0.2619	0.30918
0.9	0.11021	0.11768	0.10944	0.12929	0.10627	0.11604	0.10103	0.12024	0.11914	0.13279	0.11679	0.14784
1.0	0.03243	0.03284	0.03177	0.03253	0.03995	0.04078	0.03853	0.04030	0.03241	0.03288	0.03175	0.03258
Avg	0.53403	0.54375	0.56014	0.56439	0.52005	0.52998	0.54357	0.55480	0.52624	0.53843	0.55400	0.55973

Table1. Language Model smoothing

—x— Jelinek-Mercer —●— Dirichlet
—▲— Two-Stage

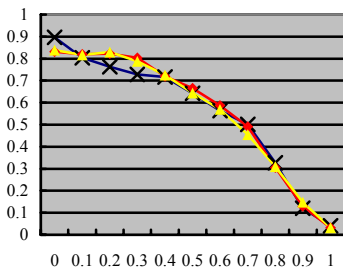


Figure 5. Recall – precision graphs for smoothing methods

—x— No Stoplist No Stem
—▲— No Stoplist Stemming

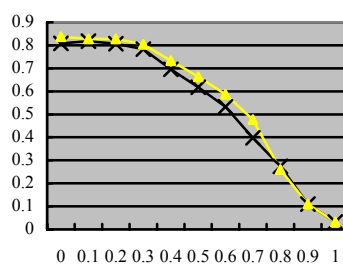


Figure 6. Recall – precision graphs for using stemming graphs

—x— No Stoplist No Stem
—▲— Stoplist No Stem

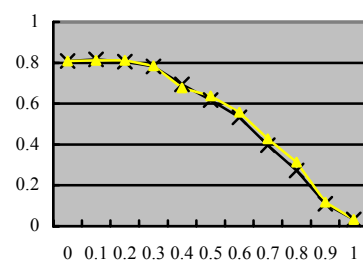


Figure 7. Recall – precision for using stoplist

The present study is still at its preliminary stage. We are still trying to construct a larger test collection for Mongolian IR. It is hopeful that other researchers can benefit from this collection in the future.

Among the future interesting investigations, we would like to investigate Mongolian morphology in more depth. In our current experimentations we only removed the last suffix. The stemming process should continue to remove other possible suffixes. Mongolian has seven vowels, and some of them have the same morphology but have different pronunciation. A possible solution is to perform query expansion by replacing letters in a query word that are confoundable with other letters.

Acknowledgement

We thank professor Jian-Yun Nie of Université de Montréal for helpful comments on this work.

This research is funded by State Key Basic Research and Development Program of China (973 Program) (2007CB316503).

References

- [1] Quejingzhabu. Mongolian Codes. Inner Mongolia University Publishing House. 2000.
- [2] Qingge’ertai. Mongolian Syntax. Inner Mongolia people publishing house. 1991.
- [3] Song, F., & Croft, B. A general language model for information retrieval. In Proceedings of the 1999 ACM SIGIR conference on research and development in information retrieval, 279-280, 1999.
- [4] Donald Metzler, Trevor Strohman, Howard Turtle and W. Bruce Croft. Indri at TREC 2004: Terabyte Track
- [5] C. Zhai and J. Lafferty. A study of smoothing methods for language models applied to information retrieval. ACM Trans. Inf. Syst., 22(2):179-214. 2004.