

잠재의미색인화모형에 의한 패췌지검색방법

리 청 한

위대한 령도자 김정일동지께서는 다음과 같이 교시하시였다.

《현시대는 과학과 기술의 시대이며 이르는 곳마다에서 요구하는것은 기술입니다. 기술을 몰라가지고서는 경제조직사업과 생산지휘를 바로할수 없으며 사회주의건설에 적극 이바지할수 없습니다.》(《김정일전집》 제2권 499~500페이지)

질문응답체계에서 중요한 문제의 하나는 다량의 문서집합에서 질문에 대한 정답이 들어있는 패췌지를 검색하는것이다. 그것은 질문에 적합한 정답이 들어있는 패췌지를 먼저 검색함으로써 제한된 시간안에 패췌지속에 포함된 정답을 찾는것이 질문응답체계의 성능평가에서 매우 중요하기때문이다.[4] 이로부터 대규모의 문서집합에서 질문에 적합한 패췌지를 검색하는 방법들이 많이 제안되였다.[1-3]

패췌지검색에서 가장 많이 리용되는 방법으로서는 일반정보검색에서 흔히 리용하는 벡토르모형을 리용한 검색방법[1], BM25검색모형과 확률모형을 리용한 검색방법[2]이 있다. 그러나 이 방법들은 문서의 크기에 비해 질문이 짧은 경우 검색의 정확도가 떨어지는 결함으로 하여 질문응답체계의 패췌지검색에서는 적합하지 않다.

본문에서는 사용자의 질문이 짧은 경우 패췌지검색에 효과적인 방법인 잠재의미색인화모형에 의한 패췌지검색방법을 제기한다.

1. 잠재의미색인화모형에 의한 패췌지검색

패췌지집합 P 를 $P=(\hat{p}_1, \hat{p}_2, \dots, \hat{p}_n)$ 에 의하여 정의된 패췌지행렬이라고 하자. 여기서

$$\hat{p}_j = \frac{p_j}{\|p_j\|}$$

이다.

이때 P 를 특이값분해하면 P 는 3개의 특징을 가지는 행렬로 분해된다. 즉

$$P=USV^T.$$

여기서 U 와 V 는 각각 크기가 $m \times r$, $n \times r$ ($r=\text{rank}(P)$)인 행렬이다. 그리고 S 는 특이값 σ_i ($\sigma_i \geq \sigma_j$, $i \leq j$)를 가지는 대각선행렬로서 다음과 같이 표시된다.

$$S=\text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$$

한편 U 와 V 에서 대응하는 렬과 함께 S 에 들어있는 가장 큰 특이값 k ($k < r$)를 유지하면 P 는 $P_k=U_k S_k V_k^T$ 에 의하여 근사되는데 이 근사값은 용어들사이의 잠재의미관계를 나타낸다. 여기서 U_k , S_k , V_k 는 각각 크기가 $m \times k$, $k \times k$, $n \times k$ 인 행렬이다.

여기에 기초하여 패썬지와 질문사이의 유사성을 계산해보자. 이를 위해 $V_k = (v_{ji})$ 에서 $v_j = (v_{j1}, \dots, v_{jk})$ ($1 \leq j \leq n$, $1 \leq i \leq k$) 를 행벡토르라고 하자. 그러면 k 차원공간에서 패썬지 p_j 는 $p_j^* = S_k v_j^T$ 로 표현되고 초기질문 역시 $q^* = U_k^T q$ 로 표현된다.

결국 질문과 패썬지의 유사성은 $\text{sim}(p_j^*, q^*)$ 에 의하여 얻어진다.

2. 잠재의미색인화모형에 의한 패썬지검색실행

패썬지에서 $P_0 \subset P_1$ 이고 $Q(P_0) = Q(P_1)$ 이면 $PS_{\text{DIDF}}(P_0) > PS_{\text{DIDF}}(P_1)$ 로 된다.

정의 패썬지 P 에 대하여 $P' \subset P$ 이고 $Q(P') = Q(P)$ 를 만족시키는 P' 가 없다면 P 를 질문용어최소패썬지(Q -minimal passage)라고 한다.

이러한 정의에 기초하여 $[1, |P|]$ 에서 질문용어최소패썬지를 검색하는 알고리즘은 다음과 같다.

```

scanner=[];
size=0;
PS=0;
for each q in Q([1, N]) {
  Insert((PL[q][0], q), scanner);
  size=size+1;
}
while(size>0) {
  (L, q1)=scanner[0];
  idfsum=0;
  for (r=0; r<size; r=r+1) {
    (R, qr)=scanner[r];
    idfsum=idfsum+idf[qr];
    PS=idfsum*exp(-β*(R-L));
  }
  Remove(scanner);
  size=size-1;
  remove(PL[q1]);
  if (PL[q1] exist) {
    insert((PL[q][0], q1), scanner);
    size=size+1; } }

```

여기서 $\text{insert}((I, q), L)$ 은 i 의 증가순서로 정렬된 목록 L 에 (I, q) 를 삽입하는 함수, $\text{remove}(L)$ 은 L 의 첫 요소를 제거하는 함수이며 PS 는 득점값이 가장 높은 패썬지를 나타낸다. 그리고 매 질문용어 $q \in Q(1, N)$ 에 대하여 $PL[q]$ 는 위치목록, $PL[q][h]$ 는 D 에서 q 의 h 번째 위치를 나타낸다.

탐색기의 초기상태는 패썬지 P 에서 나타나는 모든 질문용어의 가장 왼쪽 위치목록에 의하여 주어진다.

한편 잠재의미색인화모형에 의한 패썬지검색알고리즘의 계산시간은 $O(kn)$ 이다. 여기서 n 은 P 에서 질문용어의 총빈도수이고 k 는 질문용어의 수이다.

3. 실험결과 및 분석

론문에서는 잠재의미색인화모형에 의한 패췌지검색의 성능을 평가하기 위하여 대상자료로서 《조선전사》(1~15권)에 기초하여 만든 360개의 표준질문과 응답패췌지들을 준비하였다. 여기에 기초하여 질문응답체계에서 패췌지검색성능을 평가할 때 흔히 리용되는 MRR(거췌순위평균)평가척도를 가지고 평가하였는데 그 결과는 표와 같다. 표에서 보는바와 같이 제안된 방법이 선행한 방법들보다 검색정확도가 높다는것을 알수 있다.

표. 검색성능평가결과	
검색모형	MRR
벡토르검색모형	0.408
BM25검색모형	0.504
확률검색모형	0.502
제안된 방법	0.523

맺 는 말

잠재의미색인화모형에 의한 패췌지검색방법은 사용자질문이 짧게 제기되는 질문응답체계의 특징을 반영한 패췌지검색방법으로서 질문응답체계의 패췌지검색에서 검색의 정확도와 체계의 응답시간을 단축할수 있다.

참 고 문 헌

- [1] Wei Xu et al.; Proceedings of the 5th International Joint Conference on Natural Language Processing, 1046, 2011.
- [2] Petr Knuth et al.; Proceedings of the 23rd International Conference on Computational Linguistics, 590, 2010.
- [3] Baoxun Wang et al.; Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, 1230, 2010.
- [4] Nicolas Foucault et al.; Proceedings of Recent Advances in Natural Language Processing, 716, 2011.

주췌105(2016)년 8월 5일 원고접수

A Method of Passage Retrieval by using Latent Semantic Indexing

Ri Chong Han

We propose a passage retrieval method by using latent semantic indexing model to improve the well-known VSM and experimentally show that the retrieval on latent semantic indexing is also advantageous for dealing with short queries on condition that documents are long.

Key words: passage retrieval, question answering system, model