

# CMU Spoken Document Retrieval in Trec-8: Analysis of the role of Term Frequency TF

**M. Siegler, R. Jin and A. Hauptmann**

The participation of Carnegie Mellon University in the TREC-8 Spoken Document Retrieval Track used the basic same Sphinx speech recognition system as in TREC-7. Due to some unfortunate defaults in the parameter setup files, the speech recognizer did not perform in a reasonable manner. We will not analyze the results of the speech recognizer runs, as we believe the results contained abnormal types of errors, and insights or improvements on these errors would not generalize. A thorough examination of the speech recognition condition is given in [3]. However, we did evaluate a slightly modified weighting scheme in the reference (R1) and baseline (B1) conditions, which is described below.

## **1. Motivation for the Retrieval Formulas**

The formula we used for cmu-r1 and cmu-b1 runs is based on dtb described by Singhal[1]. In addition, to improve the accuracy, we use the standard pseudo-relevance feedback [2], conservative collection enrichment [1] and document expansion [1]. Since what we have done in that part is similar to what AT&T has done in TREC7 [1], we are not repeating the description of this approach.

Here we would like to discuss the modification we make in the dtb formula. As many people pointed out in previous TRECs, directly multiplying term frequency  $tf$  with inverse document frequency  $idf$  generally causes poor performance. The poor performance seems to be due to an overweighting of the  $tf$  term. To avoid this overestimation of  $tf$ , researchers have used  $\ln(tf+1)$  or even  $\ln(\ln(tf)+1)+1$ .

However, these approaches look more empirical than theoretical. Our modification on the  $tf$  term is a theoretically motivated attempt to resolve this problem.

## **2. Analysis**

Usually  $idf$  for a word  $A$  is written as  $\log((N+1)/M)$ . Here  $N$  is the total number of documents in the collection and  $M$  is number of documents having at least one occurrence of the word  $A$  within the collection. In other words,  $idf$  for word  $A$  can be thought of as  $-\log(p)$  where  $p$  is the probability that a document contains at least one occurrence of word  $A$  in the collection. Then the term  $tf*idf$  for a word  $A$  can be written as

$$Tf * idf = -tf * \log(p) = -\log(ptf). \quad (1)$$

We can easily extend the meaning of  $idf$  to interpret the term  $tf*idf$  for word  $A$  as  $-\log(p')$  and  $p'$  is the probability that a document contains at least  $tf$  occurrences of word  $A$ . So we have

$$Tf * idf = -\log(p) \quad (2)$$

Combing the two formulas together, we have  
 $p' = ptf$ ,

which means that the probability that a document contains at least  $tf$  occurrences of the word  $A$  is the probability that a document contains at least one occurrence of the word  $A$  to the power  $tf$ . Obviously this can be true only when the independent assumption that one occurrence of the word  $A$  has nothing to do with another occurrence of the word  $A$  is correct.

However, because of the complicated correlation within word occurrences, the independent assumption generally is wrong and will cause underestimation of probability  $p'$ . We think this is the reason why multiplying  $tf$  with  $idf$  directly generally causes overestimation and gives rise to poor performance.

### 3. Solution

To avoid this problem of overestimation by multiplying  $tf$  directly, we have come up with two solutions to replace  $tf*idf$ . Since  $tf*idf$  for word  $A$  can be interpreted as  $-\log(p')$  and  $p'$  is the probability that a document contains at least  $tf$  occurrences of word  $A$ , the key issue is how to estimate this probability  $p'$ .

One solution is using the word histogram directly. For each word  $A$ , we can build up a histogram function  $N(x, A)$  that tells the number of documents containing exact  $x$  occurrences of the word  $A$ . With this histogram function, we can compute the " $tf*idf$ " for word  $A$  as

$\log((N+1)/G(tf, A))$ .

Here  $G(tf, A)$  is defined as

$G(tf, A) = \text{Sum}(N(x, A))$  over  $x$  and  $x$  is integer from  $x$  to infinity.

The second method uses a fitted Gaussian distribution to estimate the probability  $p'$ . For each word  $A$ , we can compute the average occurrences of word  $A$   $avg\_A$  and standard deviation of occurrences of word  $A$   $std\_dev\_A$  from the histogram function  $N(x, A)$ . Now we can build the normalized Gaussian distribution as  $D(x, avg\_A, std\_dev\_A)$ .

Then the " $tf*idf$ " can be computed as

$-\log(I(tf, A))$ .

Here  $I(tf, A)$  is defined as

$I(tf, A) = \text{Integral of } D(x, avg\_A, std\_dev\_A)$  over  $x$  and  $x$  is from  $tf$  to infinity.

Furthermore we can use the standard error function to represent  $I(tf, A)$  as the following:

$0.5 * \text{err}((tf - avg\_A)/\text{sqrt}(2)/std\_dev\_A)$   
if  $tf \geq avg\_A$   
 $I(tf, A) =$   
 $0.5 * \text{err}((avg\_A - tf)/\text{sqrt}(2)/std\_dev\_A) + 0.5$   
if  $tf < avg\_A$

Both these two approaches have their advantages and disadvantages. The good side of first approach is that it uses the exact data and makes no assumption or approximation. However it may be misled by the local fluctuation. As for the second approach, it complements the down side of the first approach by using fitted Gaussian distribution. However it may cause disaster if the data doesn't fit in Gaussian distribution or when a small data set doesn't reliably estimate the true average and standard deviation.

To obtain the good properties of both approaches, we created a new method. We modify the histogram function  $N(x, A)$  for word  $A$  as follows: Instead of using natural granularity 1 for  $x$ , we use standard deviation  $\text{std\_dev\_}A$  as the granularity. We define the granularity function for each word  $A$  as

$\text{grand}(A) = \text{MAX}(\text{std\_dev\_}A/3, 1)$ .

Now we can define a new histogram function  $N'(x, A)$  as

$N'(x, A) = \text{Sum}(N(y, A))$  over  $y$  and  $y$  is from  $\text{floor}(x / \text{grand}(A)) * \text{grand}(A)$  to  $\text{ceiling}(x / \text{grand}(A)) * \text{grand}(A)$ .

From the definition of  $N'(x, A)$ , it is easily seen that  $N'(x, A)$  is the same for all  $x$  in the range  $[\text{floor}(x / \text{grand}(A)) * \text{grand}(A) .. \text{ceiling}(x / \text{grand}(A)) * \text{grand}(A)]$ .

Now we can use the first approach to compute the "tf\*idf" except that this time the histogram function is  $N'(x, A)$  instead of  $N(x, A)$ . Since the new histogram function is defined as a sum of the old histogram function over an interval on the order of standard deviation, it will be more stable and avoids some risks of the first approach.

#### **4. Experiment**

From experiments we have performed on the TREC data, we find out that the approach described above for computing factor "tf\*idf" is better than  $\text{tf} * \text{idf}$  or  $(\ln(\text{tf}) + 1) * \text{idf}$ . However, we did not see any significant improvement in performance of our formula over  $(\ln(\ln(\text{tf}) + 1) + 1) * \text{idf}$ . Instead our formula is generally slightly worse than the  $(\ln(\ln(\text{tf}) + 1) + 1) * \text{idf}$  factor.

By comparing documents weighted by our schema and weighted by  $(\ln(\ln(\text{tf}) + 1) + 1) * \text{idf}$ , we find out that our schema still has the problem of overestimating  $\text{tf}$  especially when the  $\text{tf}$  is larger. We think it is due to the fact that when  $\text{tf}$  is close to the largest  $\text{tf}$ , the  $G(\text{tf}, A)$  is very inaccurate because of the lack of histogram data. In the future we will introduce special treatment for this "ending effect".

#### **5. Conclusion**

In this paper we attempted  $\text{tf} * \text{idf}$  a more theoretic interpretation and point out a possible reason why multiplying  $\text{tf}$  directly with  $\text{idf}$  causes the poor performance with the given interpretation. We came up with a word histogram based method of integrating  $\text{tf}$  and  $\text{idf}$  into one factor which is  $\log(1/p')$  and  $p'$  is the probability that a document has at least  $\text{tf}$  occurrences of a particular word. We have several success over  $\text{tf} * \text{idf}$  and  $(\ln(\text{tf})+1)*\text{idf}$  and fail to compete with  $(\ln(\ln(\text{tf})+1)+1)*\text{idf}$ . We think the failure is due to the "ending effect" and we will pursue the problem further in the future.

#### **6. References**

[1] Singhal.A. AT&T at TREC7. In E. M. Voorhees and D. K. Harman, editors, TEXT RETRIEVAL CONFERENCE (TREC-7), page 141-151, 1998.

[2] J.J Rocchio. Relevance feedback in information retrieval. In The SMART Retrieval System-Experiments in Automatic Document Processing, page 313-323, Englewood Cliffs, NJ, 1971. Prentice Hall, Inc.

[3] Siegler,M. Integration of Continuous Speech Recognition and Information Retrieval for Mutually Optimal Performance, PhD. Thesis, Electrical and Computer Engineering, Carnegie Mellon University, 1999

<http://www.cs.cmu.edu/~msiegler/publish/PhD/thesis.ps.gz>