

## Learning Query Bias for Improved On-Line Document Classification

Ron Papka

Center for Intelligent Information Retrieval

Computer Science Department

University of Massachusetts

Amherst, MA 01003

papka@cs.umass.edu

<http://www.cs.umass.edu/~papka>

The following work describes two threshold selection algorithms for the document tracking task. The tracking task was defined by the Topic Detection and Tracking (TDT) research initiative, which is a DARPA sponsored effort comprising research groups from several commercial and academic sites. The task is a form of supervised learning and it is analyzed based on a system's ability to classify a test set of documents as relevant or non-relevant to a set of pre-specified topics derived from news events. This task is similar in nature to the document filtering track for the Text Retrieval Conferences (TREC) with three salient exceptions: 1) the information domain is a temporally ordered stream of broadcast news ; 2) solutions to tracking must make classification decisions on-line; and 3) classifiers are trained using only 1,2,4,8, or 16 relevant training documents from each topic. Much of the related research in this area has been explored using TREC's filtering methodology which focuses on training classifiers for general topics with one to a few thousand relevant training documents from heterogeneous sources including news, web, email, and other types of documents. What has emerged as the classifier of choice in the research is a query represented as a vector of features comprising stemmed words and associated weights (Figure 1). Query words and weights are determined through statistical and learning techniques utilizing inter- and intra-document word occurrences in the training data. The most common document representation is a vector of  $tf*idf$  weights co-occurring with the words in the query. Queries are compared to documents using a similarity function, which in many instances, is an extension of the inner-product of two vectors. This representation has the nice property of being extensible to linear classification using both supervised and unsupervised learning approaches. In what follows, we define and illustrate threshold estimator bias for document classification based on this representation. We discuss two methods that produce bias-reduced estimators resulting in improved classification accuracy.

q24= #WSUM( 1 0.69 crash 0.49 usair 0.49 plane 0.40 pittsburgh 0.38 flight 0.38 investigate 0.34 area 0.31 737 0.30 427 0.29 airport 0.26 wood 0.26 board 0.25 cause 0.24 safety 0.24 scene 0.24 site 0.23 piece 0.22 aircraft 0.22 pennsylvania 0.21 even);
--

Figure 1: 20 feature query formulated from the first 16 relevant documents from "Crash of US Air Flight 427"

Recently, tracking systems from various research sites were tested on the TDT2 evaluation corpus (~22K documents, 30 topics) [1]. Our system compared favorably to the best systems which used similar information retrieval representations to the one described above. An important component of our tracking system is a threshold estimator for queries formulated for the Inquiry retrieval engine. We used an estimator  $\hat{u}$  such that  $\hat{u} = 0.4 + \theta (b_{\text{optimized}} - 0.4)$ , where  $\theta$  is a global system parameter, 0.4 is the minimum belief value produced by Inquiry, and  $b_{\text{optimized}}$  is the belief value resulting from the query that, when applied to the topic's labeled training documents, optimizes a pre-specified utility function. The utility function used for this evaluation was defined by TDT and is a cost function based on error rates for false positives and false negatives produced by a system.

Using data from several experiments on the TDT1 corpus (~16K docs, 25 topics) and the TDT2 train and development corpora (~40K docs, 66 topics), it was determined that when fewer relevant training documents were used,  $b_{\text{optimized}}$ , our estimator  $\hat{u}$  when  $\theta = 1.0$ , was consistently below the parameter  $u$  it was trying to estimate i.e. the optimal threshold for the testing data of a particular topic, or simply  $b_{\text{optimal}}$ . In what follows, we assume the bias in our estimator  $\hat{u}$  is the quantity  $b(\hat{u}) = E[\hat{u}] - u$ .

We observed, as we tested various values for  $\theta$ , that bias decreased when more relevant training documents were used to formulate queries. We also observed similar but less significant increases in bias when more features were used. The observation that increasing training instances reduces the bias of an estimator, in general, is not surprising. James, for example, proves that estimates move toward the true population values when training instances are increased for data assumed to have multivariate-normal distributions [2]. He also proposes that once found, it may be possible to reduce the bias using a linear transformation. We test James's theory in the following experiments. We define a threshold estimator  $\hat{o}$ , such that  $\hat{o} = b_{\text{optimized}}$ . We then define a new bias-reduced threshold estimator  $\hat{e}$ , such that  $\hat{e} = m\hat{o} + b$ . We then compare the effectiveness of the query thresholds produced by our original estimator  $\hat{u}$ , which is bias-reduced using near optimal values for  $\theta$ , to those produced by  $\hat{e}$ , with parameters  $m$  and  $b$  learned through linear regression. In what follows, we explain these two methods that learn query bias for estimator  $\hat{o}$ .

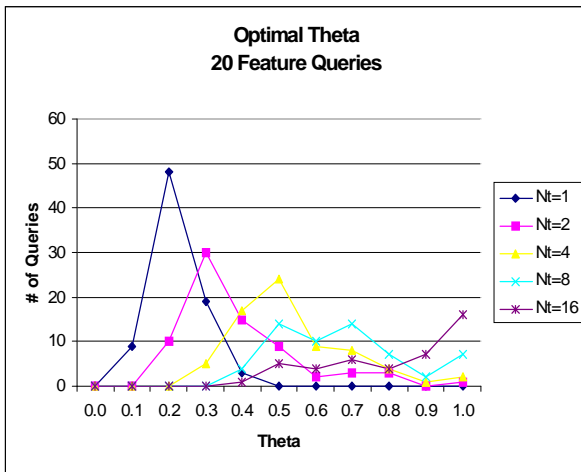


Figure 2: Histogram approach for estimator  $\hat{u}$ .

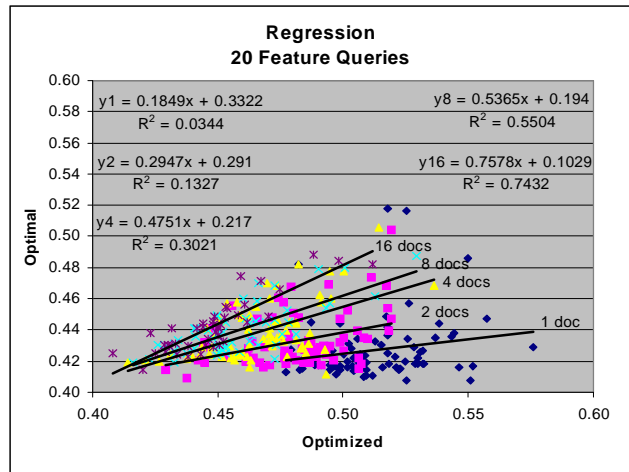


Figure 3: Regression approach for estimator  $\hat{e}$ .

Both approaches trained on 91 topics and associated labeled documents from the TDT1 and TDT2 train and development corpora. Training and testing are conducted independently over varying numbers of relevant training documents ( $N_t$ ) and varying numbers of query features. The first method, which utilizes threshold estimator  $\hat{u}$ , is illustrated using 20 word queries in Figure 2. Histograms of optimal values for  $\theta$  are collected for each value of  $N_t$  and for queries formulated with 10, 20, 50, 100, 200, 600, and 10000 features. Using this method on the training data determined that for 1 relevant document ( $N_t=1$ ) and 20 feature queries, 49 out of 91 queries had optimal cost when  $\theta = 0.2$ . For each pair of  $N_t$  and number of features, we calculate  $E[\theta]$  from the corresponding histogram, and use that value as the parameter for estimator  $\hat{u}$  when conducting evaluation experiments using the same pair. From the data in Figure 2 it is evident that as  $N_t$  increases  $E[\theta]$  increases. This observation coupled with the definition of estimator  $\hat{u}$  implies that  $b_{\text{optimized}}$  is closer to  $b_{\text{optimal}}$  when more relevant training examples are used.

The linear regression method used to train estimator  $\hat{e}$  is illustrated in Figure 3. Instead of collecting histograms, points represents by  $b_{\text{optimized}}$  and  $b_{\text{optimal}}$  are fitted using a line for each value of  $N_t$  and number of features. The slopes and intercepts of lines produced by the various regressions are subsequently used as parameters  $m$  and  $b$  for estimator  $\hat{e}$ . This method learns the same tendencies in threshold estimator bias as the approach for estimator  $\hat{u}$ . As the number of relevant training documents increases, the slope of the resulting regression line approaches 1. As the slope approaches 1,  $b_{\text{optimized}}$  approaches  $b_{\text{optimal}}$  which implies that higher values of  $N_t$  give rise to less estimator bias in the training data.

TDT2 Evaluation corpus ( 30 events )					
# of Relevant Training Documents ( $N_t$ )					
# of Features	1	2	4	8	16
10	8.2%	-20.4%	-13.2%	-21.9%	-2.9%
20	27.5%	-17.0%	-32.1%	-23.4%	-1.6%
50	5.0%	53.2%	0.0%	-4.0%	-8.6%
100	1.1%	12.7%	4.5%	-6.0%	-20.3%
200	3.3%	-7.5%	-8.2%	2.0%	-16.4%
600	27.8%	-6.0%	-7.3%	-8.0%	-1.8%
10000	27.8%	-6.0%	-7.3%	-8.0%	0.0%
Average	14.4%	1.3%	-9.1%	-9.9%	-7.4%

Figure 4: % Cost increase caused by replacing estimator  $\hat{u}$  with estimator  $\hat{e}$ .

From the recent TDT experiments evaluated by NIST, we knew estimator  $\hat{u}$  generalized and provided relatively low cost on the TDT2 evaluation corpus, so in the following experiment, we use it as a baseline to compare the effectiveness of estimator  $\hat{e}$ . We used the same query formulation and tracking processes and measured the percent increase in cost resulting from tracking with estimator  $\hat{e}$  instead of estimator  $\hat{u}$ . Since we are trying to minimize average cost across a set of queries, an increase in cost implies a decrease in classification accuracy. The results are listed in figure 4. They suggest that the histogram method works better for  $N_t=1$  and marginally so for  $N_t=2$ . However, for  $N_t > 2$  the regression method appears to reduce cost consistently for most of the query sizes tested. This relationship suggests a mixture of approaches where estimator  $\hat{u}$  is used for  $N_t \leq 2$  and estimator  $\hat{e}$  otherwise. We are, however, still evaluating these data, but we anticipate that estimator  $\hat{u}$  is not significantly more effective than estimator  $\hat{e}$  and vice versa, which would suggest that both methods are likely to do well on the currently available TDT corpora. We plan to continue investigating these approaches as other data sources with document assessments become available.

## References:

- [1] To appear in *Proceedings of the 4<sup>th</sup> Annual DARPA Broadcast News Workshop*, February 1999.
- [2] Mike James, *Classification Algorithms*, John Wiley & Sons, New York, 1985.