

DalTREC 2005 Spam Track: Spam Filtering using N-gram-based Techniques

Vlado Kešelj, Evangelos Milios, Andrew Tuttle, Singer Wang, and Roger Zhang

Faculty of Computer Science
Dalhousie University, Halifax, Canada
{vlado,eem,swang,rogerz}@cs.dal.ca

Abstract

Note: This is a very rough draft of the report that will be submitted later for the proceedings.

1 Introduction

2005 is the second year that Dalhousie University participated actively in TREC. The project DalTREC¹ was initiated in 2003 and serves as an umbrella for different TREC-related activities at Dalhousie. The tracks of interest were Question Answering, Spam, HARD, Genomics, and Enterprise. We submitted runs for the Question Answering (QA) track and the Spam track. This paper discusses the Spam track.

2 Background Review

Spam, sometimes more formally known as "Unsolicited Commercial Email" (UCE) or "Unsolicited Bulk Email" (UBE), has grown from being a minor nuisance into a global menace. Millions of email users world-wide waste untold thousands of hours every day sorting through incoming mailboxes cluttered with unwanted or even offensive messages. Meanwhile, the growth in network resources consumed by spam continues at an unsustainable rate. The spiralling direct and indirect costs of spam threaten the reliability and utility of the email infrastructure that has become so vital to the global economy.

At the core of the spam problem is the fact that spam allows vendors to transfer the vast majority of their marketing costs to the people they are marketing to. Even though the response rate to spam is believed to be vanishingly small, it is still large enough to overcome the even smaller costs of sending spam. Therefore, it becomes a volume-driven business - the more spam you send, the more money you make. But the real cost of the ensuing flood of junk email is not small, and it is borne by the owners of the networks and servers that must deliver the spam, and by the individual recipients who must waste valuable time sorting through it.

A variety of legal, structural, and technical solutions have been proposed; most are controversial and many suffer from glaring weaknesses. Legal solutions face daunting jurisdiction problems. Structural solutions, such as replacing SMTP, face staggering infrastructure costs. Many technical solutions are extremely high-maintenance, resulting in a never-ending "arms race" between the elusive spammers and the network administrators attempting to find new ways to block them.

While the world searches for a long-term solution, there is an urgent need for something businesses can use to protect their networks from the escalating costs of spam. This research is motivated by that need.

¹<http://www.cs.dal.ca/~trecacct> — DalTREC URL

3 The Naïve Bayes Classifier

The Naïve Bayes classifier is a simple statistical algorithm with a long history of providing surprisingly accurate results. It has been used in several spam classification studies [1, 2, 3, 4], and has become somewhat of a benchmark. It gets its name from being based on Bayes's rule of conditional probability, combined with the "naïve" assumption that all conditional probabilities are independent (given the class) [5].

During training, the Naïve Bayes classifier examines all of the instance vectors from both classes. It calculates the prior class probabilities as the proportion of all instances that are spam ($\text{Pr}[\text{spam}]$), and not-spam ($\text{Pr}[\text{notspam}]$). Then (assuming binary attributes) it estimates four conditional probabilities for each attribute: $\text{Pr}[\text{true}|\text{spam}]$, $\text{Pr}[\text{false}|\text{spam}]$, $\text{Pr}[\text{true}|\text{notspam}]$, and $\text{Pr}[\text{false}|\text{notspam}]$. These estimates are calculated based on the proportion of instances of the matching class that have the matching value for that attribute.

To classify an instance of unknown class, the "naïve" version of Bayes's rule is used to estimate first the probability of the instance belonging to the spam class, and then the probability of it belonging to the not-spam class. Then it normalizes the first to the sum of both to produce a spam confidence score between 0.0 and 1.0. Note that the denominator of Bayes's rule can be omitted because it is cancelled out in the normalization step. In terms of implementation, the numerator tends to get quite small as the number of attributes grows, because so many tiny probabilities are being multiplied with each other. This can become a problem for finite precision floating point numbers. The solution is to convert all probabilities to logs, and perform addition instead of multiplication. Note also that conditional probabilities of zero must be avoided; instead a "Laplace estimator" (a very small probability) is used.

4 Feature Selection

One of the difficulties in any text classification task is caused by the unstructured nature of natural language which, as opposed to highly structured database items, requires a very large number of attributes to model effectively. Among the various models developed over the last few decades, the Vector Space Model (VSM) [7] is the most popular today. In VSM, each document D is represented by a vector in m dimensional term space. The underlying term T of each dimension is often referred to as a feature or attribute and represented by the weight of T in D . Earlier models use binary term weights [8], which are either 1 when the term is present in the document, or 0 when the term is absent. This weighting scheme is conceptually simple, and computationally cheap, however it lacks the ability to count on term importance - terms that are more important to a document often tend to occur more frequently in that document. This leads to the development of TF-IDF (term frequency and inverse document frequency) [9] based models, which assumes that term importance is directly proportional to the number of times a term occurs in a document, and inversely proportional to the number of documents in which it occurs. Given a document D and a term T , $\text{TF}(D, T)$ is the number of times T occurred in D , and $\text{IDF}(T)$ is the number of documents in which T occurred divide into the total number of documents in the dataset. Since not all documents are of the same length, $\text{TF}(D, T)$ may turn out to be bigger in longer documents. In order to avoid this bias, we usually normalize it against the most frequent term in the document to get $\text{TFN}(D, T)$, meanwhile, in order to avoid big number arithmetic, the base-2 logarithm of $\text{IDF}(T)$ is usually used instead.

Since we incorporate two alternative feature representations in our filter, a "term" here can be either a naturally bounded word or a character N -gram. A character N -gram is simply a sequence of N letters taken from a text. For trigrams as an example, think of a width 3 sliding window moving from the beginning to the end of a text, one letter at a time, then the content of the window at each step is a trigram. N -gram representation has been shown competing word representation in text classification tasks [10]. For our submission, `dalSPAM1` incorporates words as features and `dalSPAM4` uses N -gram representation with N being 5.

One problem with both word representation and N -gram re representation is the excessive number of features, which has the effect of slowing down the classification of the email as it requires more I/O to load the features and their weights into memory. In the two submissions, `dalSPAM1` and

dalSPAM4, one of the goals is to see if we can still achieve good accuracy while using a smaller number of features. As a result we use mutual information (MI) [6, 2, 1, 4] to determine the best 1000 features for the classifier to work with. Our parser includes the values of all header fields, as well as the entire body text (although the Document class strips out MIME attachments before parsing, leaving only the text and HTML portions of a MIME-formatted message, along with the MIME headers for each attachment). When using word representation, the characters ",!?:-;|ç*;()[]" as well as the space, tab and new-line characters were used as token delimiters. Carriage returns were stripped out by the Document class before parsing. Tokens were truncated to a maximum length of 30 characters. For efficiency requirements of the TREC run, we do not use word stemming or stop word list. When using character N-gram representation, only alphabetical letters are considered significant, therefore all other characters are delimiters and are treated as spaces. Two or more consecutive spaces are treated as one.

5 The White-list

We incorporate an automatic white-list in our filter. Each time a miss-classified ham message is corrected during training, the sender's address will be automatically added to the white-list together with a time stamp. The length limit of the white-list is controlled by a parameter, which can be customized according to different user needs. When the length limit is exceeded, oldest entries will be truncated. Users may also add entries manually if they want. An entry can be either an individual address - which takes the form "user@domain", or a domain - which takes the form "@domain". Manually added entries are treated as permanent (with a special time stamp), and will never get truncated.

6 Major Software Components

Document - represents an RFC 822 compliant email message (RFC 822 is the basic format specification for Internet email messages);

Instance - represents an instance vector mapped from a Document;

Corpus - represents a collection of Documents;

Dictionary - stores the tokens represented by each attribute of an Instance;

Parser - parses a Document into a list of tokens;

FeatureSelector - given a set of token lists, produces a Dictionary;

Vectorizer - given a list of tokens and a Dictionary, produces an Instance;

Trainer - given a Corpus, produces a trained Classifier;

Classifier - given an Instance, returns a "spamminess" score between 0.0 and 1.0;

Evaluator - given a set of Documents, and a set of test configurations, runs a series of cross validation experiments.

7 N-gram Based Techniques (Dal2 and Dal3 submissions)

The Dal2 and Dal3 runs used a simple raw byte n-gram-based technique. The difference is only in the n-gram size: Dal2 uses n-grams of size 4, and Dal3 uses n-grams of size 5. The programs will maintain two files named `ham.profile` and `spam.profile` containing a list of up to 3000 byte n-grams with associated count frequencies. The Perl package was used `Text::Ngrams` to collect the n-grams. If the input file is longer than 10,000, only the first 10,000 bytes are used in order not spend too much time on n-gram collection.

The program classify collects n-grams from the email file, and measure the dissimilarities of the file profile with the ham and spam profiles using the formula:

$$d = \sum_{n \in \text{profile}} \left(\frac{f_1(n) - f_2(n)}{\frac{f_1(n) + f_2(n)}{2}} \right)^2 = \sum_{n \in \text{profile}} \left(\frac{2 \cdot (f_1(n) - f_2(n))}{f_1(n) + f_2(n)} \right)^2 \quad (1)$$

If the dissimilarities of the e-mail file from the ham and spam profiles are s_{ham} and s_{spam} , respectively, then the score is computed using the formula:

$$score = \frac{s_{ham}}{s_{ham} + s_{spam}}.$$

The message is classified as spam if $score > 0.5$, otherwise it is classified as ham. The program classify will print the *score* and the classification result.

The program `train` updates the ham or spam profile. It reads the appropriate profile, multiply all frequencies with 0.9 (which we call *profile aging*), collect up to 3000 most frequent n-grams from the email file with count frequencies, and add this message profile to the existing profile. It will keep only up to 3000 most frequent n-grams in the profile and save it to the class profile file.

8 Evaluation

References

2004. *Proceedings of the Thirteenth Text REtrieval Conference (TREC-13)*, Gaithersburg, MD, USA. NIST.

[1] I. Androutsopoulos, J. Koutsias, K. V. Chandrinos, G. Paliouras, and C. D. Spyropoulos. An evaluation of naive bayesian anti-spam filtering. In *Proceedings of the Workshop on Machine Learning in the New Information Age, 11th European Conference on Machine Learning (ECML 2000)*, pages 9-17, Barcelona, Spain, 2000.

[2] I. Androutsopoulos, G. Paliouras, V. Karkaletsis, G. Sakkis, C. D. Spyropoulos, and P. Stamatopoulos. Learning to filter spam E-mail: A comparison of a naive bayesian and a memory-based approach. In *Proceedings of the Workshop on Machine Learning and Textual Information Access, 4th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD 2000)*, pages 1-13, Lyon, France, 2000.

[3] Yanlei Diao, Hongjun Lu, and Dekai Wu. A comparative study of classificationbased personal e-mail filtering. In *Proceedings of PAKDD-00, 4th Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 408-419, Kyoto, JP, 2000. Springer Verlag, Heidelberg, DE.

[4] Karl-Michael Schneider. A comparison of event models for naive bayes anti-spam e-mail filtering. In *Proceedings of the 10th Conference of the European Chapter of the Association for Computational Linguistics*, 2003.

[5] Ian H. Witten and Eibe Frank. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann, 2000.

[6] Mehran Sahami, Susan Dumais, David Heckerman, and Eric Horvitz. A Bayesian approach to filtering junk E-mail. In *Learning for Text Categorization: Papers from the 1998 Workshop*, Madison, Wisconsin, 1998. AAAI Technical Report WS-98-05.

[7] G. Salton et al, A Vector Space Model for Automatic Indexing, *Commun. ACM*, 18(11):613-620, 1975.

[8] S.O. Kimbrough and J.R. Oliver, On Relevance and Two Aspects of the Organizational Memory Problem, *4th Annual Workshop on Information Technology and Systems*, 302-311, 1994.

[9] T. Joachims, Text categorization with support vector machines: Learning with many relevant features, *European Conference on Machine Learning*, 10:137-142, 1999.

[10] Y. Miao, Comparing Document Clustering using N-grams, Terms, and Words, Master Thesis, Faculty of Computer Science, Dalhousie University, 2004.