Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

¹ An Approach to Email Classification Using Bayesian Theorem

Denil Vira¹, Dr. Denil Vira² and Pradeep Raja³

¹ Mumbai University

Received: 13 December 2011 Accepted: 3 January 2012 Published: 15 January 2012

6 Abstract

7 Email Classifiers based on Bayesian theorem have been very effective in Spam filtering due to

 $_{\rm 8}$ $\,$ their strong categorization ability and high precision. This paper proposes an algorithm for

⁹ email classification based on Bayesian theorem. The purpose is to automatically classify mails

¹⁰ into predefined categories. The algorithm assigns an incoming mail to its appropriate category

¹¹ by checking its textual contents. The experimental results depict that the proposed algorithm

¹² is reasonable and effective method for email classification.

13

14

2

3

Index terms—Bayesian, Email classification, tokens, text, probability, keywords.

Introduction nternet e-mail is an essential communication method for most computer users and has been 15 treated as a powerful tool intended to idea and information exchange, as well as for users' commercial and 16 social lives. Globalization has resulted in an exponential increase in the volume of e-mails. Nowadays, a typical 17 user receives about 40-50 email messages every day. For some people hundreds of messages are usual. Thus, 18 19 users spend a significant part of their working time on processing email. As the popularity of this mean of communication is growing, the time spent on reading and answering emails will only increase. At the same 20 time, a large part of email traffic consists of nonpersonal, non time critical information that should be filtered. 21 Irrelevant emails greatly affect the efficiency and accuracy of the aimed processing work. As a result, there has 22 recently been a growing interest in creating automated systems to help users manage an extensive email flow. 23 Consider the following scenario -You have just returned from a relaxing two week vacation. There has been no 24 phone, no email for two wonderful weeks, and now you are back. You open your inbox and ... Wow! There are 25 347 new messages! How could you manage to read all of them? Probably, you will spend the whole day trying 26 to sort out all this mail. Having done this burdensome work, you feel like you need a vacation again. What is 27 worse is that most of those messages are out of your interest or out of date. 28

Here comes the need for automatic email classification system that would sort the important and the unimportant mails, thus saving a much precious time of the users.

The rest of the paper is organized as follows. The next section describes the generic Bayesian filtering logic. Section 3 introduces the proposed algorithm for email classification. Section 4 presents experiments and results. The final section consists of the conclusion.

34 **1 II.**

³⁵ 2 Bayesian filter

Bayesian filter has been used widely in building spam filters. The Naïve Bayes classifier is based on the Baye's rule of conditional probability. It makes use of all the attributes contained in the data, and analyses them individually as though they are equally important and independent of each other. The rule for conditional probability is as follows P(H | E) = P(E | H)P(H) / P(E) ... Eq.(1)

Where P(H|E) is the conditional probability that hypothesis H is true given an evidence E; P(E|H) the conditional probability of E given H, P(H) the prior probability of H, P(E) the prior probability of E. In the case of spam classification, hypothesis H can be defined as spam or legitimate given an email E.

In applying this theorem, an email needs to be tokenized, and the extracted n tokens (i.e. words or phrases) are then used as the evidences E{e1, e2,?,en}, and their probability of spam or non-spam is calculated from

4 IMPLEMENTATION AND RESULTS

⁴⁵ previous emails that have been classified. Using past classified emails to estimate the probabilities of the tokens

 $_{\rm 46}$ $\,$ belonging to spam or no spam is a learning process and then they are used to predict the spam probability of a

⁴⁷ new incoming email. Assume that n key-words are extracted from the content of an email as evidences, then the ⁴⁸ probability that the email is spam can be calculated by:P(S|E(e1,e2...en)) = P(E(e1,e2...en)|S)*P(S) / P(E(e1,e2...en)|S)

49 e1,e2.. en)? Eq. (2)

⁵⁰ Where, P(S|E) is the probability that an email E is spam S. In practice, an assumption is commonly made ⁵¹ naïvely for simplifying calculation, that is, all the evidences are independent from each other. For example, if ⁵² an email contains four evidences: e1, e2, e3 and e4, then the joint and conditional probabilities given S can be ⁵³ accelerated by $P(c_1, c_2, c_3, c_4) = P(c_1)P(c_2)P(c_2)P(c_3)P(c_4)S) = P(c_1)S)P(c_2)SP(c_3)P(c_4)S)$

saily calculated by, P(e1,e2,e3,e4) = P(e1)P(e2)P(e3)P(e,), P(e1,e2,e3,e4|S) = P(e1|S)P(e2|S)P(e3|S)P(e4|S)

Thus the probability that a given email is considered spam when evidences (e1, e2, e3 and e4) appear in the given email is calculated by Equation (2), and a decision can be made if it is higher than a pre-set decision threshold, 0.5, usually.

57 III.

⁵⁸ 3 Proposed algorithm for email classification

The entire email classification process is divided into two phases. a) First phase is the Training phase/Learning 59 phase During this phase, the classifier will be trained to recognize attributes for each category. So that later on 60 when a new mail arrives, it compares the attributes of the mail with attributes of each category and the mail is 61 classified into the category having most similar attributes as that of the mail. To build the attribute list (also 62 referred as keywords database) for each category, the emails will be classified manually in to different categories 63 64 by the user. For better understanding of the working, we create two categories for emails say, work and personal. 65 The user manually specifies whether each mail belongs to work or personal category based on its contents including the subject and the body. 66

For each email classified manually by the user, the algorithm extracts keywords from the mail and stores them

into the keyword database for that particular category along with count of the number of occurrences. For eg.
If the user classifies the mail as work, the keywords extracted from that mail will be stored in work_database.
Every category will have its own set of keywords database.

Every category will have its own set of keywords database.
Also, before the mail is processed for keywords, it is first filtered. Mail is divided into set of tokens separated

⁷² by blank space or any other punctuation marks of English language. Proverbs, articles, html tags, noise words ⁷³ and other unnecessary contents are removed and then the keywords are extracted. Thus the training phase is ⁷⁴ responsible for building the keyword database for each predefined category.

b) The second is Classification phase Once the learning is sufficiently done, the algorithm is ready to move to 75 next phase. The new mails that arrive should to be automatically assigned their categories. Basically, here we 76 compare the contents of the mail with those of all category keyword databases using Bayesian theorem and look 77 for best matching category for the mail. The new incoming mail (also referred as unclassified mail) is broken into 78 tokens and filtered. The tokens are then compared with keyword databases of each category. The probability 79 that the mail belongs to a category is found out for every category. The category for which the probability 80 outcome is highest, it is then compared with the threshold of same category and classified into that category if 81 following condition holds true else the mail stays unclassified. P n = 1, if P(category | E) > threshold P n = 0 82 , if $P(\text{category} \mid E) < \text{threshold Where E is the newly arrived mail.}$ 83

⁸⁴ 4 Implementation and results

The implementation of the algorithm was done in Java. Authors' working email accounts were used for fetching the mails. A total of 5175 e-mails are used for the purpose of experiment.

 $_{87}$ In classification task, the performance of algorithm is measured in terms of accuracy. Accuracy = N c / N

N c : Number of correctly classified mails N : Number of total mails classified The accuracy of classification is
dependent on several parameters such as number of mails considered during training phase, number of categories
defined, threshold value of each category, size of mail etc.

We initially created two categories, work & personal and then gradually added three more categories. The test was repeated for different values of parameters. Through repeated tests, we choose the best-performing configuration for our algorithm. 10-fold cross-validation was used in our test: the corpus was partitioned randomly into ten parts, and the experiment was repeated ten times, each time reserving a different part for testing and using the remaining nine parts for training. All the results were then averaged over the ten iterations. The size of training data is the number of mails manually classified during learning phase.

97 Note that the accuracy is also dependent on number of overlapping words among different categories. We 98 observed that when two different categories have many keywords in common, the classifier accuracy is low. With 99 distinct keywords for each category, with minimum overlapping, the classifier accuracy was above 0.9. The results 100 also show that larger the size of training data better is the accuracy. ¹⁰¹ **5 VI.**

102 6 Conclusion

103 Considering the requirements of improving efficiency in processing emails, this paper introduced an approach to

104 classify mails based on Bayesian theorem. The algorithm was implemented and the experimental results were 105 studied. The results show that our approach to classify mails is a reasonable and effective one. However, there is

lot of work to be done in future. The future research includes improving accuracy of the classifier ,working with

107 languages other than English, making the keyword extraction more efficient, classifying based on attachments, understanding the semantics of email text and so on.



108

Figure 1: I

¹© 2012 Global Journals Inc. (US)

 $^{^2 \}mathbbm{O}$ 2012 Global Journals Inc. (US) Global Journal of Computer Science and Technology



Figure 2: 1.



15

Figure 3: Figure 1:5.

on proposed algorithm		
Size of	No. of	Accuracy
training data	categories	
100	2	0.841
200	2	0.925
600	2	0.992
1000	2	0.991
200	3	0.936
500	3	0.973
1000	3	0.995
200	4	0.762
500	4	0.836
1000	4	0.861
200	5	0.846
500	5	0.921
1000	5	0.967

Figure 4: Table 1 :

6 CONCLUSION

109 .1 Acknowledgements

The authors would like to acknowledge the contribution of Mr.Dhruven Vora and Mr.Jenish Jain for helping us in coding of the algorithm and testing of results. We would also like to acknowledge Ms.Prof. Kavita Kelkar for valuable suggestions and guidance. Lastly, we would like to thank the entire Department of Computer Engineering, K.J.Somaiya College of Engineering, Mumbai for providing us with all the resources needed for the research.

- [Sahami et al. ()] 'A Bayesian approach to filtering junk email'. M Sahami , S Dumais , D Heckerman , E Horvitz
 Proceedings of AAAI-98 Workshop on Learning for Text Categorization, (AAAI-98 Workshop on Learning for Text Categorization) 1998. p. .
- [Mccallum and Nigam ()] 'A Comparison of Event Models for Naive Bayes Text Classification'. Andrew Mccallum
 , Kamal Nigam . WS-98-05. AAAI-98 Workshop on Learning for Text Categorization, 1998. p. . (Technical
 Report)
- $_{121}$ [David et al. ()] 'A sequential algorithm for training text classfiers'. D David , William A Lewis , Gale .
- Proceedings of SIGIR-94, 17th ACM International Conference on Research and Development in Information
 Retrieval, (SIGIR-94, 17th ACM International Conference on Research and Development in Information
 Retrieval) 1994.
- [Androutsopoulos et al. ()] 'An Evaluation of Naïve Bayesian Anti-Spam Filtering'. I Androutsopoulos , J
 Koutsias , K V Chandrinos , G Paliouras , C D Spyropoulos . Proc. of the Workshop on Machine Learning
 in the New Information Age, 11th European Conference on Machine Learning, (of the Workshop on Machine
 Looming in the New Information Age, 11th European Conference on Machine Learning, Sparse
- Learning in the New Information Age, 11th European Conference on Machine LearningBarcelona, Spain)
 2000. 2000. p. .
- [Peng et al. ()] 'Augmenting naive bayes classifiers with statistical language models'. F Peng , D Schuurmans ,
 S Wang . Information Retrieval 2004. 7 p. .
- Beckermann et al. ()] Automatic categorization of email into folders: benchmark experiments on Enron and
 SRI corpora, R Beckermann , A Mccallum , G Huang . IR-418. 2004. University of Massachusetts Amherst
 (Technical report)
- [Hershkop and Stolfo ()] 'Combining e-mail models for false positive reduction'. S Hershkop , J Stolfo . Proc of KDD'05 of ACM? Chicago?, (of KDD'05 of ACM? Chicago?) 2005. p. . (s.n.]? pp)
- [Wang ()] 'Diversity between neural networks and decision trees for building multiple classifier systems'. W Wang
 Multiple Classifier Systems, 2000. p. .
- [Youn and Mcleod ()] 'Efficient Spam e-mail Filtering using Adaptive Ontology'. S Youn , D Mcleod . International Conference on Information Technology (ITNG'07), 2007. p. .
- [Androutsopoulos et al. ()] 'Learning to Filter Unsolicited Commercial E-Mail'. I Androutsopoulos, G Paliouras
 , E Michelakis . Technical Report of National Centre for Sciential Research 2004. (Demokritos)