

# An approach of the Naive Bayes classifier for the document classification <sup>1</sup>

Ioan Pop

## Abstract

To perform the ranking document or the Web Mining tasks we have considered an approach based on the Naive Bayesian algorithm. The implementation of Naive Bayes algorithm was made in the PMML language.

**2000 Mathematics Subject Classification:** 60-04, 65C60

**Keywords:** Naive Bayesian classifier, document classification, Web Mining

## 1 Introduction

This article approaches the implementation of the Naive Bayes (shortly NB) classifier. It shows that the algorithm NB improves the tasks of the Web Mining by the accuracy documents classification. Its applications are important in the following areas: e-mail spamming; filtering spam results out of search queries; mining log files for computing system management; machine learning for Semantic Web; document ranking by text classification; hierarchical text categorization; managing content with automatic classification and other areas from Web Mining.

---

<sup>1</sup>*Received 19 August, 2007*

*Accepted for publication (in revised form) 17 October, 2007*

## 2 The NB probabilistic model

Abstractly, the probability model for a classifier is a conditional model  $p(C|F_1, F_2, \dots, F_n)$  over a dependent class variable  $C$  with a small number of outcomes or classes, conditional on several feature variables  $F_1$  through  $F_n$ . The problem is that if the number of features  $n$  is large or when a feature can take on a large number of values, then basing such a model on probability tables is infeasible. We therefore reformulate the model to make it more tractable. The Bayes' theorem relates the conditional and marginal probabilities of stochastic events  $C$ , and  $F$ :

$$(1) \quad Pr(C|F) = \frac{Pr(F|C)Pr(C)}{Pr(F)}$$

where:  $P(C)$  is the prior probability of hypothesis  $C$ ;  $P(F)$  is the prior probability of training data  $F$ ;  $P(C|F)$  is the probability of given  $F$  and;  $P(F|C)$  is the probability of  $F$  given  $C$ . Using Bayes' theorem for several feature variables  $F_n$ , we can rewrite this as:

$$(2) \quad p(C|F_1, \dots, F_n) = \frac{p(C)p(F_1, \dots, F_n|C)}{p(F_1, \dots, F_n)}.$$

In practice we are only interested in the numerator of that fraction, since the denominator does not depend on  $C$  and the values of the features  $F_i$  are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model (1) which can be rewritten using repeated applications of the definition of conditional probability as:

$$(3) \quad \begin{aligned} p(C, F_1, \dots, F_n) &= \\ &= p(C)p(F_1|C)p(F_2|C, F_1)p(F_3|C, F_1, F_2)\dots p(F_n|C, F_1, F_2, \dots, F_{n-1}) \end{aligned}$$

This means: assuming that each feature  $F_i$  is conditionally independent of every other feature  $F_j$  for  $j \neq i$  and  $p(F_i|C, F_j) = p(F_i|C)$  the model (1) can be expressed as:

$$(4) \quad p(C, F_1, \dots, F_n) = p(C)p(F_1|C)p(F_2|C), \dots = p(C) \prod_i p(F_i|C).$$

This means that under the above independence assumptions, the conditional distribution over the class variable  $C$  can be expressed like this:

$$(5) \quad p(C, F_1, \dots, F_n) = \frac{1}{Z} p(C) \prod_{i=1}^n p(F_i|C),$$

where  $Z$  is a scaling factor dependent only on  $F_1, \dots, F_n$ , i.e., a constant if the values of the feature variables are known. The corresponding classifier for this model is the classify function defined as follows:

$$(6) \quad \text{classify}(f_1, \dots, f_n) = \operatorname{argmax}_c P(C = c) \prod_{i=1}^n p(F_i = f_i|C = c).$$

### 3 The NB model for the document classification

Consider the problem of classifying documents by their content, for example spam and non-spam E-mails. Imagine that documents are drawn from a number of classes of documents which can be modelled as sets of words where the (independent) probability that the  $i$ -th word of a given document occurs in a document from class  $C$  can be written as  $p(w_i|C)$ . (Simply, we assume that the probability of a word in a document is independent of the length of a document, or that all documents are of the same length.) Then the probability of a given document  $D$ , given a class  $C$ , is

$$(7) \quad p(D|C) = \prod_i p(w_i|C).$$

Using the Bayesian result above, and the assuming that there are only two classes,  $S$  and  $\neg S$  (e.g. spam and not spam) we can write:

$$(8) \quad \frac{p(S|D)}{p(\neg S|D)} = \frac{p(S)}{p(\neg S)} \prod_i \frac{p(w_i|S)}{p(w_i|\neg S)}.$$

Thus, the probability ratio  $p(S|D)/p(\neg S|D)$  can be expressed in terms of a series of likelihood ratios. The actual probability  $p(S|D)$  can be easily computed from  $\log(p(S|D)/p(\neg S|D))$  based on the observation that

$p(S|D) + p(\neg S|D) = 1$ . Taking the logarithm of all these ratios, we have:

$$(9) \quad \ln \frac{p(S|D)}{p(\neg S|D)} = \ln \frac{p(S)}{p(\neg S)} + \sum_i \ln \frac{p(w_i|S)}{p(w_i|\neg S)}.$$

Finally, the document can be classified as follows. It is spam if  $\ln \frac{p(S|D)}{p(\neg S|D)} > 0$ , otherwise it is not spam.

## References

- [1] I. Rish, *An empirical study of the naive Bayes classifier*, IJCAI 2001 Workshop on Empirical Methods in Artificial Intelligence. (available online: PDF, PostScript).
- [2] M. Mozina, J. Demsar, M. Kattan, and B. Zupan, *Nomograms for Visualization of Naive Bayesian Classifier*, In Proc. of PKDD-2004, pages 337-348. (available online: PDF), 2004.
- [3] O. R. Duda, P. E. Hart and D. G. Stork, *Pattern classification (2nd edition)*, Section 9.6.5, p. 487-489, Wiley, ISBN 0471056693,2000.
- [4] PMML 3.0, <http://www.dmg.org/pmml-v3-0.html>
- [5] PMML 3.1, <http://www.dmg.org/pmml-v3-1.html>
- [6] PMML, <http://www.sourceforge.net/projects/pmml>

Department of Informatics, Faculty of Sciences

“Lucian Blaga” University of Sibiu

Str. Dr. I. Rațiu nr. 5-7, 550012 - Sibiu, România,

E-mail: [ioan.pop@ulbsibiu.ro](mailto:ioan.pop@ulbsibiu.ro)