Naïve Bayes Classifiers Christina Wallin, Period 3 Computer Systems Research Lab 2008-2009

<u>Abstract</u>

One part of computational linguistics is the classification and comparison of texts into classes. This comparison could be author classification, spam filtering, or in the case of this project, the classification of texts from different genres of news stories. The goal of the project is to be able to classify a text as belonging to one of two classes based on the actual words in the text. This method of classification is called a naive Bayes classifier. In this age of huge amounts of data available online, a classifier which could discriminate between two types of news could be increasingly useful.

Background

The naive Bayes classification is a relatively simple method for classifying texts based on the false assumption that all of the variables, in this case words in the documents, are independent of each other. Even though this assumption is false, this project is done to achieve fundamental understanding concerning the effectiveness of the naive Bayes as compared to other methods, and to find a way of improving upon the performance of this classifier. However, it is also done to try to provide a way for news stories from different genres to be classified. There have been many studies which have used the naive Bayesian classification method. One paper which puts together information from previous studies is "Idiot's Bayes--Not So Stupid After All?" by David Hand and Keming Yu, which using theoretical and real data situations shows that the Naive Bayes is not a horrible method because of its false assumption that all of the variables (in my case, occurrences of words) are independent. In "Improving the Performance of Naive Bayes for Text Classification," by Yirong Shen and Jing Jiang, the authors explain a way to improve the naive Bayes by combining it with logistic regression, a mathematical method. The naive Bayes method itself is a bit overconfident in classification, and so by combinging it with logistic regression and using the Porter stemmer, the classification is improved.

<u>Methodology</u>

The programming language used for these manipulations is Python. The database which I will be using to classify texts is the 20 Newsgroups database. There are 20 different genres of news stories, and it is divided into training and testing sets. The training and testing sets are separated in time by a little, and so it is more realistic.

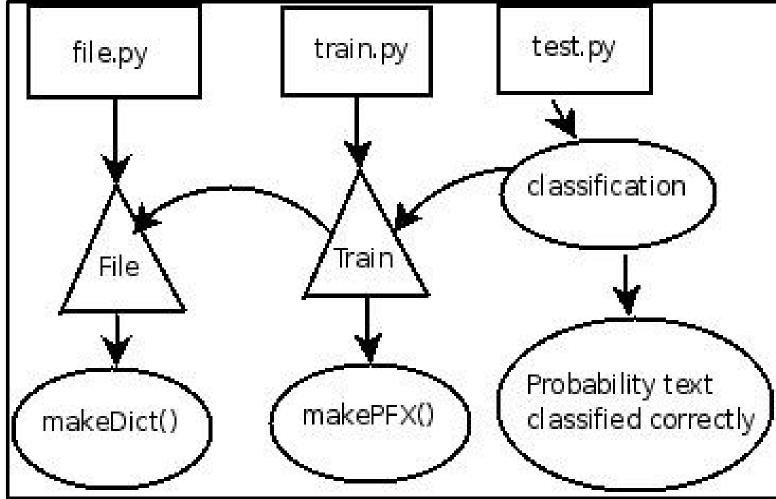
The first step in classifying a document is to read and parse the file using the file.py program, removing all punctuation and case. Then, I make a dictionary with the words occurring and their frequencies. In this step, it is possible to use a Porter stemmer to stem words to their roots —for example, "running" and "runs" would both go to "run."

Next, for each class/genre, train.py trains the program as to what characteristics are most prevalent. It does this with the words themselves, creating an array containing the PFX, or probability that each word occurs in the class. As the program goes through in turn all of the files in a class, making a dictionary for each one, the PFX is calculated as the number of texts in a class which contain at least one instance of a word over the total number of texts in a class.

Then with this probability vector for each class, I can calculate the probability that a text is of a specific class in test.py by generating the probability vector for that specific text and comparing it to the PFX for each class. For this, each variable (i.e. the occurrence or non-occurrence of each word) has to be calculated in order to form the probability that it is in a particular class. Then, the probabilities of each word are multiplied in order to determine what the probability that the file is in a specific class is. The program has been tested by using small data sets and checking manually whether or not the frequencies were correct, and using my testing program. I also tested the program by making sample data sets based on a programmed-in probability for each word. With this perfect data, I was able to check my PFX vector and found that it figured out the same probability as was programmed in. Thus, the probability calculation is correct. I was also able to check whether the testing part classified the perfect data correctly, which it did.

<u>Results</u>

I now have a working naïve Bayes classifier, and so I have conducted an experiment to determine whether or not using the Porter stemmer to stem words improves the percentage of files classified correctly. I expected that it would help because words similar to each other would be counted together. I found that it in fact does not help, and was about a percentage point less effective. For example, with the stemmer, 82.6 percent correctly classified vs 83.6 percent without in alt.atheism and rec.autos. I will repeat this test when I have the multivariate model working, for perhaps it helps to stem in that case.



Program Methodology