Matthew Letter

Project 2

CS529 Machine Learning

March 19, 2015

## Introduction:

**A high-level description on how your code works. :**

The overall goal of the code was to use Naïve Bayes for classifying documents. The first part of which was training a Bayesian system. The data.txt and label.txt file were first parsed and turned into document objects and stored in a list. This list was used to train the Bayesian classifier. The document list was passed on and the probability of each word for each topic was calculated. Beta was set to 1/|vocabulary|. The reverted value for any word not found in a category was set to beta/|words in category|. $P(x_i|y_j)$ was set to |word in category|+beta/|words in category| and $P(y)$ was set to |documents in category/|number of documents|, and $P(y|x)$ was obtained by taking the max value of the log of $P(y_j)$ and adding the sum of all log $P(x_i)$. Which ever had the highest value was used as the classifier for the category. After the classifier was trained the test documents were loaded and parsed using the same approach as the training data. The test values were then classified using the trained Bayesian classifier. The number of correctly and incorrectly classified documents was recorded down and used to determine the % accuracy. The accuracy was the number of correct over |documents|.

**The accuracies you obtain under various settings:**

See figure 2

**Explain which options work well and why:**

When beta is set to between $10^{-2}$ and $10^{0}$ the accuracy was maximized. It was also noted that if the top 100 most common words are removed the accuracy is greatly increased. When beta is too small outliers will not get smoothed out. Conversely if beta is too large important information starts to become background noise.

## Questions:

**Question 1:**

It would be difficult to accurately estimate parameters because each document would be considered a unique index yet there are more than the 1000 words in each document (50,000). This makes it hard to get an accurate probability for each word at a unique index.

**Question 2:**

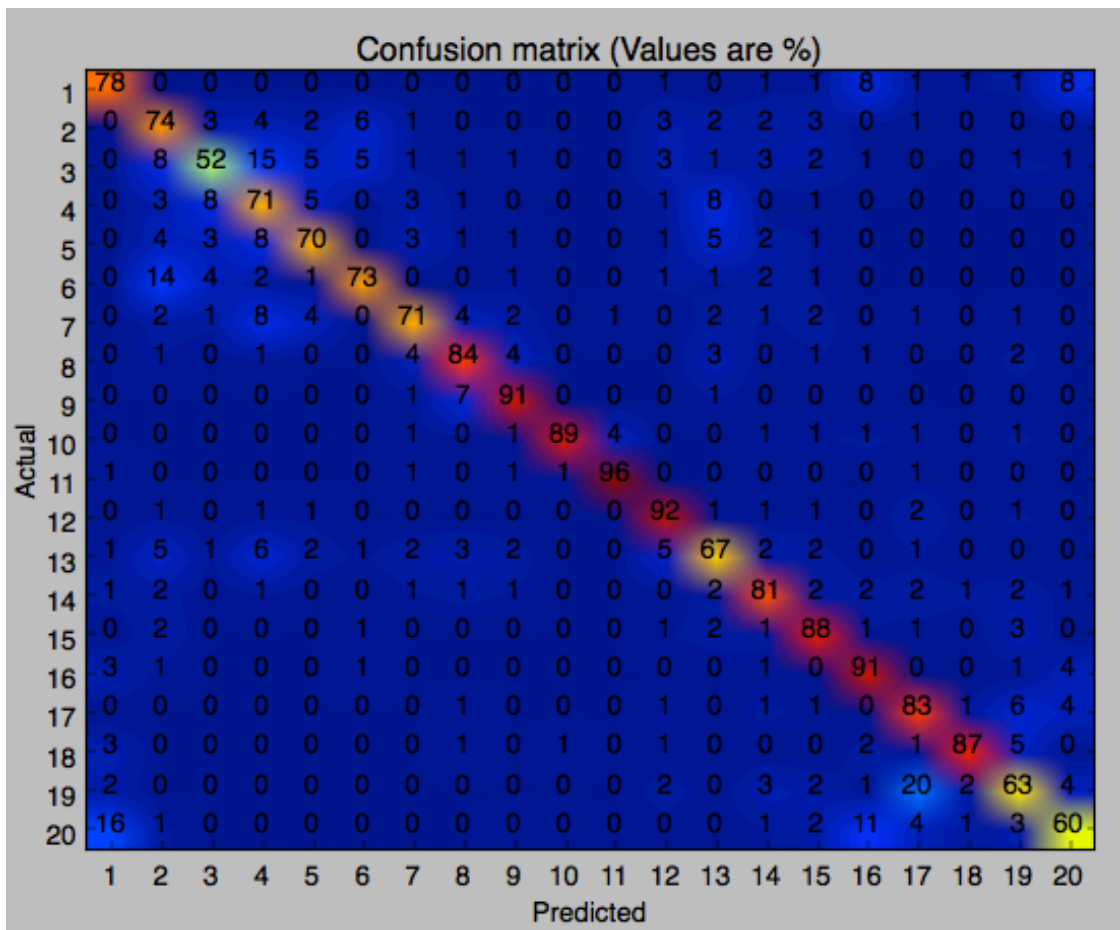Overall testing accuracy was 78.521% when beta was set to 1/|Vocabulary|.



Figure 1: Confusion matrix for accuracy of 78.521% and beta = 1/|vocabulary|

**Question 3:**

Newsgroups 3, 19, 20 are newsgroups that the algorithm appears confused with more than others. This could be because the words in these groups have a high frequency in multiple groups. For example in news group 20 the words used for this newsgroup could also be used frequently in newsgroups 1 and 16 and a little bit with 17 and 19 (based off the confusion matrix)
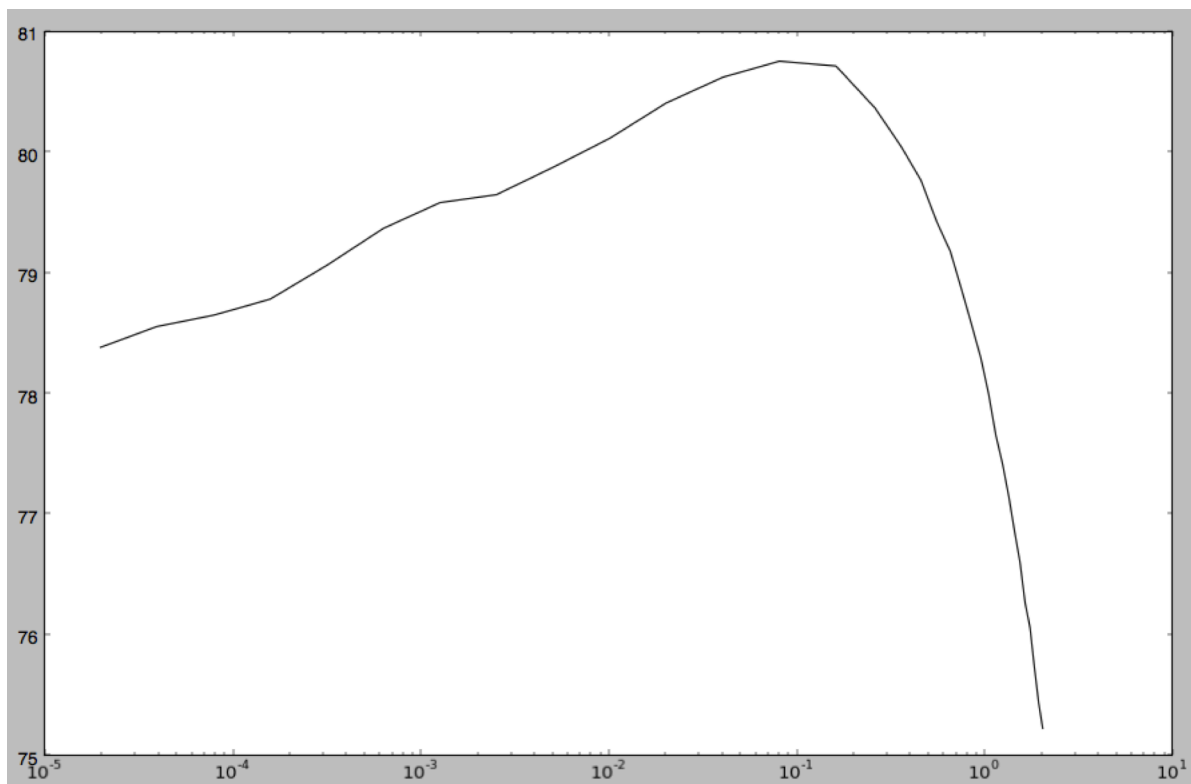
**Question 4:**



Figure 2:  Values of beta set from 0.00001 to 2

The accuracy drops for low beta because inaccuracies are no longer being reduced to the point that they begin to affect the actual classification accuracy. When beta is large actual useful information begins to get degraded.

**Question 5:**

To choose the top 100 words I ranked the words by finding the words with the highest information gained by choosing the words with the highest probability standard deviations. As the words with the most information gained would be those that had the highest standard deviation. Furthermore I removed the top 100 most frequently used words in the English language from the list of possible words to use for the top 100.

**Question 6:**

After implementing the proposed solution for question 5 here is the list of top 100 words

libraries , library , remained , english , realizing , ny , wallpapers , brad , for , attributed , knows , nationally , shift , corel , tapscott , cdt , cynical , kevin , blanks , package , printing , jk , iwll , kent , born , grabbing , science , anthology , teenage , keeps , unfortunately , ba , vanheyningen , materials , dominant , ames , information , games , outta , davidr , pattern , attest , server , signature , pretend , opossum , combining , but , correct , light , shoulder , votes , enlighten , hamish , throws , lynn , evaluation , per , tired , noyes , damaging , mantis , deal , people , accept , announced , nfs , guy , leon , stay , socket , nestorius , broken , the , cap , lists , project , questions , virginia , modified , kept , users , theocracy , vogle , mmwang , gif , produce , pretty , discriminate , axes , shoudl , cetera , justin , start , someday , murderers , write , mormons , atheist , nm

**Question 7:**

After review of the top 100 words that the classifier is relying on I made a confusion matrix using only the top 100 words for classifying the documents. I ended up getting a classification rate of 44.144% accuracy. Based on the below confusion matrix we see that some documents were actually really well classified based on the top 100 words leading me to believe that there is a biased in the data.
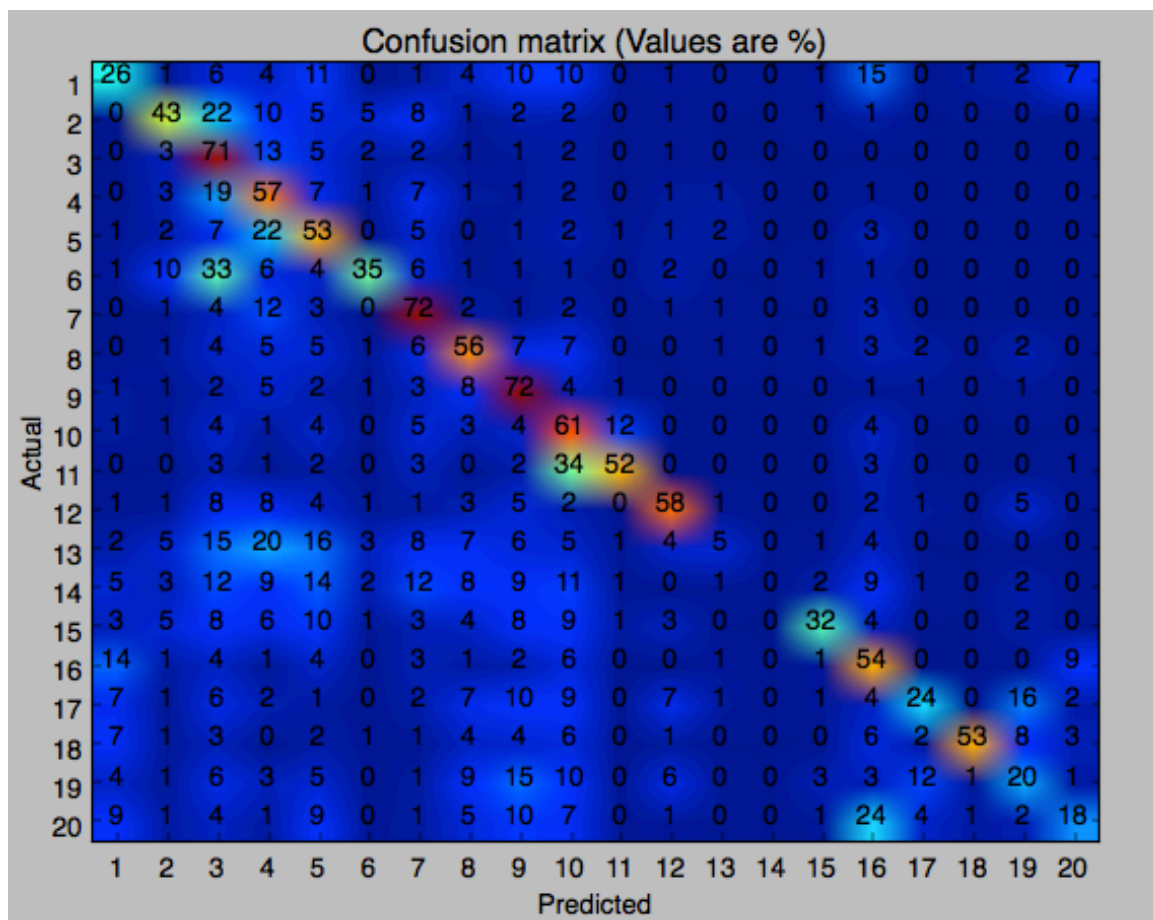


Figure 3: classification with the top 100 words with accuracy of 44.144%