

Naive Bayes Document Classifier

In this project, you will implement the Naive Bayes document classifier and apply it to the classic 20 newsgroups dataset (For a brief description and history of this dataset, but do NOT download the data, see: <http://qwone.com/~jason/20Newsgroups/>) In this dataset, each document is a posting that was made to one of 20 different usenet newsgroups.

Your goal is to write a program which can predict which newsgroup a given document was posted to.

Model 1 (do not implement this)

Say we have a document D containing d words; call the words $\{X_1, \dots, X_d\}$. The value of random variable X_i is the word found in position i in the document. We wish to predict the label Y of the document, which can be one of m categories. We could use the model:

$$P(Y | X_1 \dots X_d) \propto P(X_1 \dots X_d | Y) P(Y) = P(Y) \prod_{i=1}^d P(X_i | Y)$$

That is, each X_i is sampled from some distribution that depends on its position X_i and the document category Y . As usual with discrete data, we assume that $P(X_i | Y)$ is a multinomial distribution over some vocabulary V ; that is, each X_i can take one of $|V|$ possible values corresponding to the words in the vocabulary. Therefore, in this model, we are assuming (roughly) that for any pair of document positions i and j , $P(X_i | Y)$ may be completely different from $P(X_j | Y)$.

Question 1: In your answer sheet, explain in a sentence or two why it would be difficult to accurately estimate the parameters of this model on a reasonable set of documents (e.g. 1000 documents, each 1000 words long, where each word comes from a 50,000 word vocabulary). **[5 points]**

Model 2 (implement this)

To improve the model, we will make the additional assumption that:

$$\forall i, j \quad P(X_i | Y) = P(X_j | Y)$$

Thus, in addition to estimating $P(Y)$, you must estimate the parameters for the single distribution $P(X | Y)$, which we define to be equal to $P(X_i | Y)$ for all X_i . Each word in a document is assumed to be an iid draw from this distribution.

Data

The data file (available on UNM Learn) contains six files:

1. vocabulary.txt is a list of the words that may appear in documents. The line number is word's id in other files. That is, the first word ('archive') has wordId 1, the second ('name') has wordId 2, etc.
2. newsgrouplabels.txt is a list of newsgroups from which a document may have come. Again, the line number corresponds to the label's id, which is used in the .label files. The first line ('alt.atheism') has id 1, etc.
3. train.label Each line corresponds to the label for one document from the training set. Again, the document's id (docId) is the line number.
4. test.label The same as train.label, except that the labels are for the test documents.
5. train.data Specifies the counts for each of the words used in each of the documents. Each line is of the form "docId wordId count", where count specifies the number of times the word with id wordId in the training document with id docId. All word/document pairs that do not appear in the file have count 0.
6. test.data Same as train.data, except that it specifies counts for test documents.

If you are using matlab, the functions textread and sparse will be useful in reading these files

Implementation

Your first task is to implement the Naive Bayes classifier specified above. You should estimate $P(Y)$ using the MLE, and estimate $P(X|Y)$ using a MAP estimate with the prior distribution $\text{Dirichlet}(1 + \beta, \dots, 1 + \beta)$, where $\beta = 1/|V|$ and V is vocabulary.

Question 2: In your answer sheet, report your overall testing accuracy (Number of correctly classified documents in the test set over the total number of test documents), and print out the confusion matrix (the matrix C , where c_{ij} is the number of times a document with ground truth category j was classified as category i). **[10 points]**

Question 3: Are there any newsgroups that the algorithm confuses more often than others? Why do you think this is? **[5 points]**

3.4 Priors and Overfitting

In your initial implementation, you used a prior $\text{Dirichlet}(1 + \beta, \dots, 1 + \beta)$ to estimate $P(X|Y)$, and I told you set $\beta = 1/|V|$. Hopefully you wondered where this value came from. In practice, the choice of prior is a difficult question in Bayesian learning: either we must use domain knowledge, or we must look at the performance of different values on some validation set. Here we will use the performance on the testing set to gauge the effect of β .

Question 4: Re-train your Naive Bayes classifier for values of α between .00001 and 1 and report the accuracy over the test set for each value of β . Create a plot with values of β on the x-axis and accuracy on the y-axis. Use a logarithmic scale for the x-axis (in Matlab, the `semilogx` command). Explain in a few sentences why accuracy drops for both small and large values of β . **[5 points]**

Identifying Important Features

One useful property of Naive Bayes is that its simplicity makes it easy to understand why the classifier behaves the way it does. This can be useful both while debugging your algorithm and for understanding your dataset in general. For example, it is possible to identify which words are strong indicators of the category labels we're interested in.

Question 5: Propose a method for ranking the words in the dataset based on how much the classifier 'relies on' them when performing its classification (hint: information theory will help). Your metric should use only the classifier's estimates of $P(Y)$ and $P(X|Y)$. It should give high scores to those words that appear frequently in one or a few of the newsgroups but not in other ones. Words that are used frequently in general English ('the', 'of', etc.) should have lower scores, as well as words that only appear extremely rarely throughout the whole dataset. Finally, in your method this should be an overall ranking for the words, not a per-category ranking. **[5 points]**

Question 6: Implement your method, set β back to $1/|V|$, and print out the 100 words with the highest measure. **[5 points]**

Question 7: If the points in the training dataset were not sampled independently at random from the same distribution of data we plan to classify in the future, we might call that training set biased. Dataset bias is a problem because the performance of a classifier on a biased dataset will not accurately reflect its future performance in the real world. Look again at the words your classifier is 'relying on'. Do you see any signs of dataset bias? **[5 points]**

Turn in the following:

Your code. Submit your file through UNM Learn. **The due date is Midnight (+ 3 hrs buffer) contingent to the late policy stated in the syllabus.** Your code should contain appropriate comments to facilitate understanding. If needed, your code must contain a Makefile or an executable script that receives the paths to the training and testing files

A report of about 4 to 6 pages that describes:

- A high-level description on how your code works.
- The accuracies you obtain under various settings.
- Explain which options work well and why.
- Answers to questions 1 to 7

Rubric:

- Your code is thoroughly commented (10 pts)
- You provided a README file (5 pts)
- Your code executes correctly and is platform independent and you provided a Makefile or script for execution and clear instructions for execution are provided in the README file (20 pts)
- Implementation of NB is correct (15 pts)
- Your report is clear, concise and well organized (10 pts)
- Your answers to questions 1 to 7 (40 pts)

TOTAL: 100 pts (10 pts of your final grade)