

Assignment 3: Bayesian Inference

EECS 492: Artificial Intelligence

Fall 2015

Now: Thursday, October 29, 9:00 am

Due: Tuesday, November 17, 11:00 pm

In this assignment, you will implement one flavor of Naive Bayes (NB) classifier – the Bernoulli Naive Bayes – to solve an interesting problem in Natural Language Processing (NLP). Naive Bayes classifiers very succinctly and beautifully illustrate the core principle of Bayesian Inference – gathering evidence from many small sources, and accumulating them using a probabilistic framework to yield a single large signal. You will use the classifier to solve Authorship Attribution – a long-standing problem in NLP.

We will provide you with the data (text documents) and tokenizer,¹ and you will write your own code to implement the classifier (both training and testing), extract features, and produce output in the submission format (explained below). The Naive Bayes classifier is explained below.

Let d be a document, c be a class label, C is the set of all class labels, and $P(c|d)$ is the conditional probability of the class label given the document. We can now define the best class label c^* for d as the one that maximizes $P(c|d)$:

$$c^* = \operatorname{argmax}_{c \in C} [P(c|d)] \quad (1)$$

We will call this a *maximum a posteriori* (MAP) classifier because we will use Bayes Rule to compute $P(c|d)$ as a posterior probability from the prior probability $P(c)$ and the likelihood $P(d|c)$.

By Bayes Rule, we have:

¹The data is available here: <http://www.mathcs.duq.edu/~juola/problems/problemset.tar.gz>. Please feel free to download it and get started.

$$P(c|d) \propto P(c)P(d|c) \quad (2)$$

and $P(d|c)$ is approximated as:

$$P(d|c) \approx \prod_{i=1}^n P(f_i|c) \quad (3)$$

where document d is *represented by* n features f_1, f_2, \dots, f_n . Each feature imparts a small amount of evidence into the model. Note that we assume f_i 's are *independent* given the class label c . This is called the *conditional independence assumption* of Naive Bayes classifiers, and is the reason why they are called “Naive”. Despite this assumption, however, Naive Bayes works surprisingly well in practice.

We usually estimate $P(c)$ and $P(f_i|c)$ on the training data. The estimated value for $P(c)$ is $\hat{P}(c)$, and the estimated value for $P(f_i|c)$ is $\hat{P}(f_i|c)$. The final form of the classifier looks as follows:

$$c^* = \operatorname{argmax}_{c \in C} [\hat{P}(c) \prod_{i=1}^n \hat{P}(f_i|c)] \quad (4)$$

Since probabilities are small, and multiplication makes them smaller, a straight-forward implementation of Equation (4) is very likely to underflow, and return a value of zero. However, we are actually interested in the value c^* of c that maximizes the expression, rather than the actual maximum value. Since log is a monotonic function, we can take the log of Equation (4) and get the same maximum in Equation (5) without the danger of underflow.²

$$c^* = \operatorname{argmax}_{c \in C} [\log \hat{P}(c) + \sum_{i=1}^n \log \hat{P}(f_i|c)] \quad (5)$$

One of the principal variants of Naive Bayes classifier is the Bernoulli Naive Bayes, where $\hat{P}(f_i|c)$ is computed as:

$$\hat{P}(f_i|c) = \frac{N_{ci} + 1}{N_c + 2} \quad (6)$$

where N_{ci} is the number of documents with class label c that contain feature f_i , and N_c is the number

²Please use log base 2 in all cases.

of documents with class label c .

A natural generalization that is useful for other problems is the Multinomial Naive Bayes classifier, with the following more general formula for $\hat{P}(f_i|c)$:

$$\hat{P}(f_i|c) = \frac{T_{ci} + 1}{\sum_{j=1}^{|V|} T_{cj} + |V|} \quad (7)$$

where T_{ci} is the number of times feature f_i appeared under class label c , and V is the set of unique features in the training data (also known as the *vocabulary*). The addition of one in the numerator and $|V|$ in the denominator is called *add-one smoothing*. It is done so that zero probabilities of unseen features cannot make the whole likelihood zero.

For both Bernoulli and Multinomial Naive Bayes, prior probability estimate $\hat{P}(c)$ is computed as follows:

$$\hat{P}(c) = \frac{N_c}{N} \quad (8)$$

where N is the total number of documents (in the training set).

Naive Bayes classifier has been briefly touched upon in AIMA pages 499 and 808. For more details on the two flavors of Naive Bayes, please see Chapter 13 of the book “Introduction to Information Retrieval” by Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze (available on CTools under Resources/readings/13bayes.pdf).

Authorship Attribution

Now that we know about the Naive Bayes classifier, the next step is to use it to solve a practical problem in NLP. In Authorship Attribution, you are given documents with known authors. Your task is to train a model on these documents, and then test the model on documents of unknown authorship. For example, say we would like to know if a particular play of Shakespeare was written by Shakespeare or Marlowe. One way to answer that question will be to train a classifier (e.g., Naive Bayes) on the plays of Shakespeare and Marlowe, and then use the trained model on the disputed play to see what our model predicts. The success of computational stylometry in this regard stems partly from the fact that Bayesian reasoning can be very

effective at teasing apart the *authorial fingerprints* of individual people. It appears that stop words³ are good candidate features for this problem, since authors have little conscious control over the choice of stop words. Several lists of English stop words are available online. We have provided you with the one available at <http://www.lextek.com/manuals/onix/stopwords1.html>. This is a standard list, and has been used in several research papers.

The dataset we will use for this assignment, is a subset of the Ad-hoc Authorship Attribution Competition (AAAC) Dataset, publicly available at http://www.mathcs.duq.edu/~juola/authorship_materials2.html.⁴ The dataset is divided into 13 problems, the first eight being in English, and the rest in four other languages (French, Latin, Dutch, and Serbian-Slavonic). We will use Problems A, B, C, G, and H for this assignment. Each problem has a separate directory: “problemA”, “problemB”, “problemC”, and so on. Within a particular directory, say “problemA”, we have training examples with known authors:

- Atrain01-1.txt (example 1 from author 01)
- Atrain01-2.txt (example 2 from author 01)
- Atrain01-3.txt (example 3 from author 01)
- ...
- Atrain13-1.txt (example 1 from author 13)
- Atrain13-2.txt (example 2 from author 13)
- Atrain13-3.txt (example 3 from author 13)

and test examples with unknown authors:

- Asample01.txt
- Asample02.txt
- Asample03.txt
- and so on.

“problemA” directory has 13 test examples.

³Very frequently used English words, such as “a”, “an”, “the”, “am”, “is”, “are”, etc. Stop words are also often called *function words*, as they carry very little semantic content.

⁴Please make sure to download the complete dataset as a tarball: <http://www.mathcs.duq.edu/~juola/problems/problemset.tar.gz>.

For the assignment, we have provided you with a ground truth file (“test_ground_truth.txt”) that contains author information for the test examples of each AAAC problem. You will use this file to obtain the test accuracy of your classifier.

Assignment

Please implement the Bernoulli Naive Bayes classifier, as described above, for the Authorship Attribution problem. Tokenize the text using the given code.⁵ For Authorship Attribution, please use the given stop words as features. Performance of your classifier will be evaluated on the AAAC test examples. The performance measure we will use, is called *accuracy*. It is the percentage of correctly classified test examples. Please make sure your code takes one argument – the AAAC directory containing training and test documents for a particular problem. You may assume that the stop words file and ground truth file (“test_ground_truth.txt”) are present in the same directory as your code.

Output the following for each problem (A, B, C, G, and H):

1. **Test Accuracy:** What is the accuracy on the test data?
2. **Confusion Matrix:** What is the confusion matrix?⁶
3. **Feature Ranking:** Rank top 20 features by their *class-conditional entropy* (on training data). The class-conditional entropy (CCE) for a particular feature f_i is estimated as:

$$CCE_i \approx - \sum_{c \in C} [\hat{P}(c) \hat{P}(f_i|c) (\log \hat{P}(f_i|c))]^7 \quad (9)$$

For each feature, compute its class-conditional entropy. Then sort the features in descending order of their entropy, and present the top 20 features along with their entropy values.

4. **Feature Curve:** Sort the features in the descending order of their frequency (in the training data). Then take top 10 features, and train your model using a vocabulary of these top 10 features. Save the test accuracy of this model. Then take top 20 features, and save the test accuracy of the corresponding model. Then take top 30, top 40, and so on, and for each – save the corresponding test accuracy. In the end, plot the test accuracy against number of features.⁸ What conclusion can be drawn from this curve?

⁵Tokenization involves lowercasing the text, removing all punctuation, and splitting the text into words.

⁶It is a matrix where the rows are “ground truth” class labels, and columns are predicted classes. An example is given here: https://en.wikipedia.org/wiki/Confusion_matrix. We have provided you with a function to print out the confusion matrix.

⁷Please use log base 2 in all cases.

⁸X-axis is number of features, Y-axis is test accuracy.

What to submit

A zipped directory with one subdirectory: **code**, and one PDF file. Name your directory “A3-uniquename”, where “uniquename” is replaced by your own uniquename. The **code** subdirectory contains the code, the given stop words file, the given ground truth file, and a README with enough information to run your code. Specifically, your program should take one argument: the AAAC directory containing training and test files, and it should print out the test accuracy, the confusion matrix, the feature ranking, and the test accuracy for top k features (k varying from 10 to the maximum in steps of 10). We will go through the code to make sure they are present. Please submit on CTools, through the Assignments tool. If for any reason that doesn’t work, submit it to the Drop Box, and send email.

Here are two sample commands for C++ and python:⁹

```
./mycode AAAC_problems/problemA/
```

```
python mycode.py AAAC_problems/problemA/
```

Now, please answer the following questions (they are graded).¹⁰ Once you are done, save this PDF as “answers.pdf” under “A3-uniquename” directory, zip it, and upload on CTools.

AAAC Problem A

- What is the test accuracy?¹¹

The test accuracy was 0.384615384615

- What is the confusion matrix?

The confusion matrix was

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13
1 0 0 0 0 0 0 0 0 1 0 0 0 0
2 0 0 0 0 0 0 0 0 0 1 0 0
3 0 0 0 0 1 0 0 0 0 0 0 0 0
4 0 0 0 1 0 0 0 0 0 0 0 0 0
5 0 0 0 0 1 0 0 0 0 0 0 0 0
6 1 0 0 0 0 0 0 0 0 0 0 0 0
7 0 0 0 0 0 0 1 0 0 0 0 0 0
8 0 0 0 0 0 0 0 0 0 0 1 0 0
9 0 0 0 0 0 0 0 0 1 0 0 0 0
10 0 0 0 0 0 0 0 0 0 0 1 0 0
11 0 0 0 0 1 0 0 0 0 0 0 0 0
12 0 0 0 0 0 0 0 0 0 0 0 0 1
13 0 0 0 0 0 0 0 0 0 1 0 0 0
```

⁹For C++, you may want to use something like the following command: `g++ -O3 -std=c++11 mycode.cpp -o mycode.`

¹⁰Please fill in the PDF directly. To edit PDF, you can use an online utility like <https://www.pdfescape.com/>.

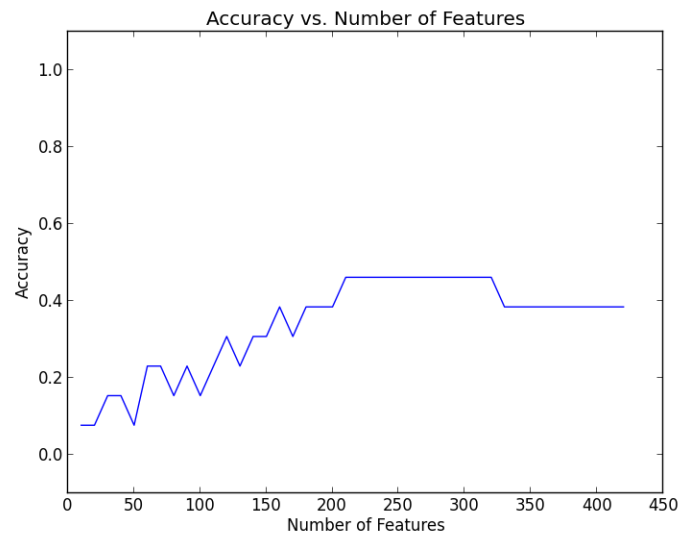
¹¹We got 38.4615% with our implementation.

- What is the feature ranking? Please list the top 20 features, along with their class-conditional entropy.

The feature ranking was:

| Feature | Entropy |
|----------|----------------|
| everyone | 0.353681544803 |
| need | 0.348943051238 |
| always | 0.348943051238 |
| others | 0.347727876904 |
| working | 0.344326761443 |
| men | 0.344326761443 |
| high | 0.344326761443 |
| off | 0.344326761443 |
| think | 0.343111587109 |
| around | 0.343111587109 |
| his | 0.341774209005 |
| work | 0.341147469829 |
| me | 0.340843032907 |
| really | 0.340803442211 |
| making | 0.340681238441 |
| been | 0.340681238441 |
| should | 0.339466064108 |
| long | 0.338373093544 |
| whole | 0.33728012298 |
| our | 0.336064948646 |

- Please give the feature curve as described above (should be legible). What conclusion can be drawn from this curve?



From this curve we can conclude that after around 300 features, the program begins to overfit the data

AAAC Problem B

- What is the test accuracy?

The test accuracy was 0.307692307692

- What is the confusion matrix?

The confusion matrix was

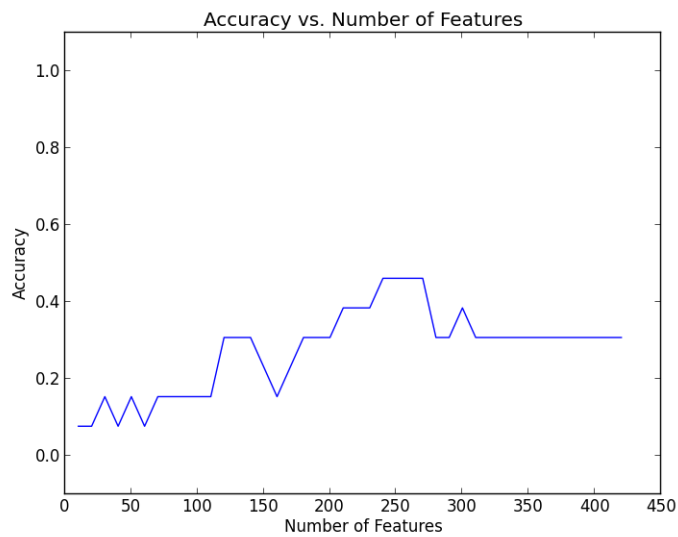
| | | | | | | | | | | | | | |
|----|---|---|---|---|---|---|---|---|---|----|----|----|----|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

- What is the feature ranking? Please list the top 20 features, along with their class-conditional entropy.

The feature ranking was

| Feature | Entropy |
|---------|----------------|
| going | 0.354896719136 |
| against | 0.353681544803 |
| find | 0.351251196135 |
| always | 0.348943051238 |
| having | 0.347850080674 |
| greater | 0.347850080674 |
| areas | 0.344326761443 |
| around | 0.343111587109 |
| again | 0.343111587109 |
| taken | 0.341896412775 |
| between | 0.341896412775 |
| do | 0.341896412775 |
| year | 0.340803442211 |
| her | 0.340803442211 |
| enough | 0.340803442211 |
| wanted | 0.340681238441 |
| put | 0.340681238441 |
| making | 0.340681238441 |
| great | 0.340681238441 |
| point | 0.339588267878 |

- Please give the feature curve as described above (should be legible). What conclusion can be drawn from this curve?



From this curve, we conclude that after around 250 features, the program begins to overfit the data.

AAAC Problem C

- What is the test accuracy?

The test accuracy was 0.444444444444

- What is the confusion matrix?

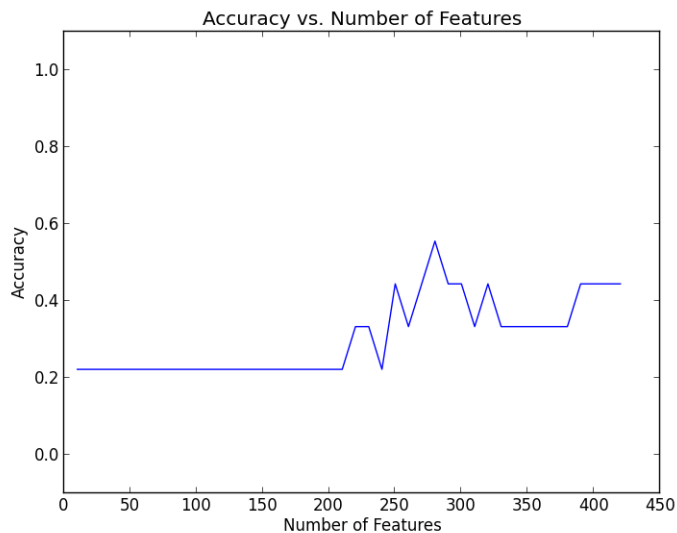
The confusion matrix was

| | | | | | |
|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 |
| 1 | 0 | 2 | 0 | 0 | 0 |
| 2 | 0 | 2 | 0 | 0 | 0 |
| 3 | 0 | 1 | 0 | 0 | 0 |
| 4 | 0 | 0 | 2 | 0 | 0 |
| 5 | 2 | 0 | 0 | 0 | 0 |

- What is the feature ranking? Please list the top 20 features, along with their class-conditional entropy.

The feature ranking was
 yours 0.36394971844
 gets 0.359330775865
 ways 0.354711833291
 showing 0.354711833291
 ends 0.351455121113
 everything 0.348049125965
 beings 0.346836178538
 interesting 0.343941124578
 pointed 0.34343018339
 orders 0.34138641864
 presented 0.339322182003
 o 0.333349197786
 says 0.332838256598
 want 0.33214853349
 sides 0.33214853349
 largely 0.33214853349
 younger 0.331446527251
 parted 0.330794491848
 number 0.328730255211
 fully 0.328219314023

- Please give the feature curve as described above (should be legible). What conclusion can be drawn from this curve?



From the curve, we conclude that after around 270 features, the program begins to overfit the data

AAAC Problem G

- What is the test accuracy?

The accuracy was 0.25

- What is the confusion matrix?

The confusion matrix was

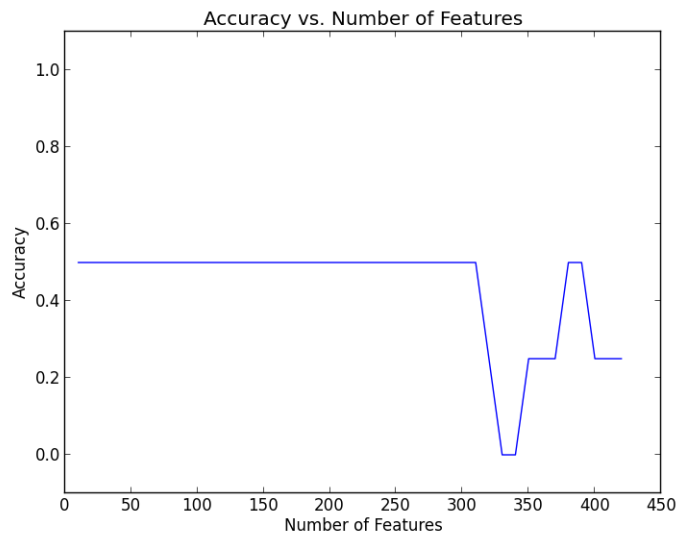
```
0 1 2
1 1 1
2 2 0
```

- What is the feature ranking? Please list the top 20 features, along with their class-conditional entropy.

The feature ranking was

| Feature | Entropy |
|-----------|----------------|
| works | 0.36651629275 |
| oldest | 0.36651629275 |
| lets | 0.36651629275 |
| facts | 0.36651629275 |
| differ | 0.36651629275 |
| x | 0.344201937618 |
| wants | 0.344201937618 |
| thinks | 0.344201937618 |
| states | 0.344201937618 |
| s | 0.344201937618 |
| generally | 0.344201937618 |
| evenly | 0.344201937618 |
| anywhere | 0.344201937618 |
| anyone | 0.344201937618 |
| shows | 0.336505833505 |
| sees | 0.336505833505 |
| presents | 0.336505833505 |
| per | 0.336505833505 |
| o | 0.336505833505 |
| mostly | 0.336505833505 |

- Please give the feature curve as described above (should be legible). What conclusion can be drawn from this curve?



From the curve, we conclude that after around 300 features, the program begins to overfit the data

AAAC Problem H

- What is the test accuracy?

The test accuracy was 0.666666666667

- What is the confusion matrix? The confusion matrix was

| | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 1 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 |
| 3 | 0 | 1 | 0 |

- What is the feature ranking? Please list the top 20 features, along with their class-conditional entropy.

The feature ranking was

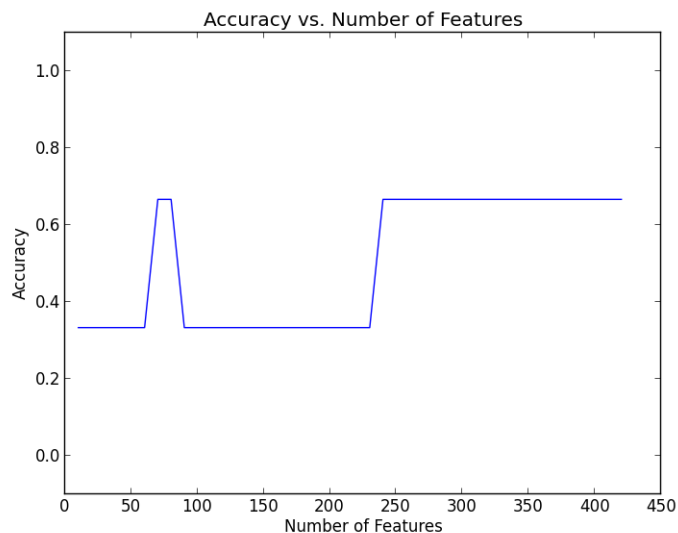
Feature | Entropy

```

-----
works 0.366204096223
within 0.366204096223
ways 0.366204096223
wants 0.366204096223
turn 0.366204096223
thus 0.366204096223
taken 0.366204096223
somewhere 0.366204096223
seemed 0.366204096223
rooms 0.366204096223
presenting 0.366204096223
presented 0.366204096223
points 0.366204096223
pointing 0.366204096223
others 0.366204096223
opens 0.366204096223
nobody 0.366204096223
later 0.366204096223
knew 0.366204096223
interests 0.366204096223

```

- Please give the feature curve as described above (should be legible). What conclusion can be drawn from this curve?



From the feature curve, it looks like the program still fits the data well after 400 features, but it usually under fits the data for less than 250 features.

Scoring Details

Please make sure your code runs on CAEN before you submit.

We will grade the following on both completion and correctness. That is, each part has to be complete **and** correct to receive full credit.

- Code compiles and runs on CAEN (10 points)
- Test accuracy (5 problems * 4 points each = 20 points)
- Confusion matrix (5 problems * 4 points each = 20 points)
- Feature ranking (5 problems * 4 points each = 20 points)
- Feature curve (5 problems * 4 points each = 20 points)
- Code prints the values for feature curve¹² (5 problems * 2 points each = 10 points)

Total: 100 pts.

Extra Credit

Please fill in the PDF directly. To edit PDF, you can use an online utility like <https://www.pdfescape.com/>.

- Implement the Multinomial Naive Bayes classifier to solve another interesting problem in NLP – Language Identification. In Language Identification, the goal is to identify the language of a given document. Train a classifier on documents with known language, and test it on other documents. Language identification is interesting because people often mix two or more languages while writing, esp. on social media like Facebook and Twitter. This phenomenon is variously known as *code switching* or *code mixing*. Given that we can automatically identify people’s language, it becomes easier to analyze cultural, social, political, and economic backgrounds of a large section of human population who do not necessarily speak English. For language identification problems, character bigrams (sequences of two characters) have been found to perform well.¹³ What dataset to use for Language Identification?

¹²Test accuracy for top k features, k varying from 10 to the maximum in steps of 10.

¹³Example of character bigrams: the phrase “hot dog.” has seven character bigrams – “ho”, “ot”, “t ”, “ d”, “do”, “og”, and “g.”. Space characters and punctuation symbols are considered valid while extracting character bigrams.

A good choice will be a subset of the Europarl dataset.¹⁴ You may use character bigrams (with punctuation and spaces) as features. What is the test accuracy, confusion matrix, feature ranking, and feature curve for Language Identification?

- Instead of character bigrams, use character trigrams (three-character sequences) for Language Identification. Does it make any difference? How about higher-order character n-grams?

- Do word n-grams have an impact on performance ($n > 1$)?¹⁵

¹⁴<http://www.statmt.org/europarl/>.

¹⁵The sentence “I go to UMich” has three word bigrams – “I go”, “go to”, and “to UMich”, and two word trigrams – “I go to”, and “go to UMich”.

- Instead of feature frequency, use a different measure of salience, e.g., *tfidf*.¹⁶ What difference does it make?
- *Stemming* refers to the practice of partial or complete removal of inflection from a word. For example, “run”, “runs”, and “running” are all stemmed to the base form “run”. Several English stemmers are available online. Please use the Porter Stemmer. Does stemming have any impact on performance of Authorship Attribution and/or Language Identification?
- Does the *length* of the text samples in number of words/characters (for training and/or test) have an impact on performance?

¹⁶*tfidf*, or *term frequency inverse document frequency*, is defined as the product of term frequency (tf) and inverse document frequency (idf). The formula to compute *tfidf* is explained here: <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>. Please note that the term frequency component needs to be normalized to prevent a blow-up of score due to high-frequency terms.