# CS-734/834 Introduction to Information Retrieval: Assignment #5: 10.3, 10.6, 10.9, $SVM^{light}$

*Dr. Michael L. Nelson*

**Plinio Vargas**

pvargas@cs.odu.edu

# Contents

# List of Figures

# Listings

# 1 Exercise 10.3

Compute five iterations of HITS (see Algorithm 3) and PageRank (see Figure 4.11) on the graph in Figure 10.3. Discuss how the PageRank scores compare to the hub and authority scores produced by HITS.

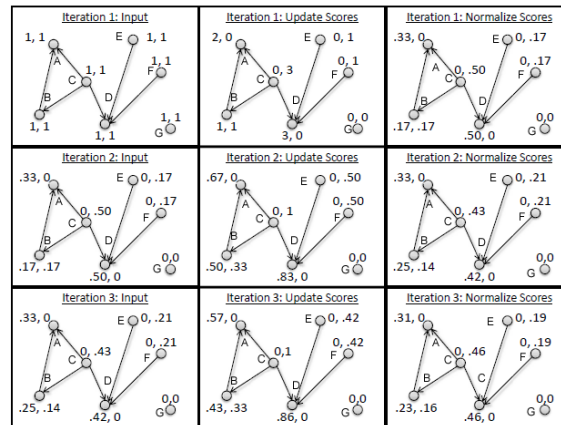Figure 1: HITS vs PageRank 5 Iteration Calculation



Fig. 10.3. Illustration of the HITS algorithm. Each row corresponds to a single iteration of the algorithm and each column corresponds to a specific step of the algorithm.

## 1.1 Approach

Vertices identification was added to the graph in order to better reference the nodes for further analysis. Graph $G$ on Figure 1 is composed of the following vertices: {A, B, C, D, E, F, G}.

### 1.1.1 HITS

To calculate 5 iterations of the HITS algorithm for the graph shown on Figure 1, the *Python* scripts *hits.py* was developed (See Listing 1). The data structure to construct the graph is a dictionary, shown on lines 30 through 36. Each *vertex* in the graph has another dictionary whose value contains the authority and hub score calculated during the iteration. The *vertex* has an outbound connection to another on. A dictionary with the key value 'o' has an array containing all outbound *vertices* connections. An additional dictionary, **node_value** on line 37, has the authority and hub tuple initial iteration values. The remaining lines are the implementation of the algorithm provided in textbook [3]

Listing 1: Script Generate $n$ Number of HITS Iterations

```
28    # creating the graph
29    graph = {'A': {'A':[0, 0]},
30              'B': {'o':['A'], 'B':[0, 0]},
31              'C': {'o':['A', 'B', 'D'], 'C':[0, 0]},
32              'D': {'D':[0, 0]},
33              'E': {'o':['D'], 'E': [0, 0]},
34              'F': {'o':['D'], 'F': [0, 0]},
```

```
35                'G': {'G':[0, 0]}}
36
37      node_value = {'A': [1, 1], 'B': [1, 1],'C': [1, 1],'D': [1, 1],'E': [1, 1],'F'
           : [1, 1],'G': [1, 1]}
38
39      for k in range(5):
40          hub_total = 0
41          auth_total = 0
42          for node in graph:
43              if 'o' in graph[node]:
44                  for vertex in graph[node]['o']:
45                      graph[node][node][1] += node_value[vertex][0]
46                      graph[vertex][vertex][0] += node_value[node][1]
47                      hub_total += node_value[vertex][0]
48                      auth_total += node_value[node][1]
49
50          for node in graph:
51              node_value[node][0] = graph[node][node][0] / auth_total
52              node_value[node][1] = graph[node][node][1] / hub_total
53              graph[node][node][0] = 0
54              graph[node][node][1] = 0
55
56          print('\nfor k=%d HITS Calculation is:' % (k + 1))
57          for node in sorted(node_value):
58              print('%c = (%.2f, %.2f)' % (node, node_value[node][0], node_value[
                  node][1]))
59
60      print('\nEnd Time:  %s' % strftime("%a,  %b %d, %Y at %H:%M:%S", localtime()))
61      print('Execution Time: %.4f seconds' % (time()-start))
62      return
```

### 1.1.2   PageRnk

To calculate five iterations of the PageRank algorithm for the graph shown on Figure 1, we developed the *Python* scripts *pagerank.py*. It has a difference the data structure as the one implemented in **HITS**. Here the *vertices* contain the values of the inward connection with them. Also, the dictionary contains a key *size* indicating the number of outbound connections, and the key *pr* indicating the calculated PR scored during the iteration. Additionally, the PR score at the beginning of the iteration is captured in the dictionary **pr** on line 45.

Listing 2: Script Generate $n$ Number of PageRank Iterations

```
20  cachedStopWords = stopwords.words("english")
21
22
23  def main():
24      # record running time
25      start = time()
26      print('Starting Time: %s\n' % strftime("%a,  %b %d, %Y at %H:%M:%S", localtime
           ()))
27
28      # creating the graph
29
```

```python
30      graph = {'A': {'links': ['B', 'C'], 'size': 0, 'pr': 0},
31              'B': {'links': ['C'], 'size': 1, 'pr': 0},
32              'C': {'links': [], 'size': 3, 'pr': 0},
33              'D': {'links': ['C', 'E', 'F'], 'size': 0, 'pr': 0},
34              'E': {'links': [], 'size': 1, 'pr': 0},
35              'F': {'links': [], 'size': 1, 'pr': 0},
36              'G': {'links': [], 'size': 0, 'pr': 0}}
37      """
38      graph = {'A': {'links':['C'], 'size': 2, 'pr': 0},
39              'B': {'links': ['A'], 'size':1, 'pr': 0},
40              'C': {'links': ['A', 'B'], 'size': 1, 'pr': 0}}
41      """
42      landa = 0.15
43      n = len(graph)
44      # pr = {'A': 1/n, 'B': 1/n, 'C': 1/n}
45      pr = {'A': 1/n, 'B': 1/n, 'C': 1/n, 'D': 1/n, 'E': 1/n, 'F': 1/n, 'G': 1/n}
46      y = lambda x: graph[x]['links']
47      l = lambda x: 1 if graph[x]['size'] == 0 else graph[x]['size']
48      page_rank = lambda k: landa/n + (1 - landa) * sum([pr[x]/l(x) for x in y(k)])
49
50      for k in range(5):
51          print('Iteration k=%d:' % k)
52          for node in sorted(graph):
53              graph[node]['pr'] = page_rank(node)
54              print('PR(%c) = %.4f' % (node, graph[node]['pr']))
55
56          for node in graph:
57              pr[node] = graph[node]['pr']
58          print()
```

## 2  Solution

Looking at Figure 2, we can find some similarities and difference between **HITS** and **PageRank** algorithm results. In the HITS algorithm at the fifth iteration, *vertex* **D** shows the highest authority score, similarly the same *vertex* has the highest PR rank scored through the same number of iterations. The same was noticed for *vertex* B, which has the second highest authority rank scored on the **HITS** algorithm. It also has the second high PR score in the **PageRank** algorithm.

**PageRank** seemed to converged at the third iteration while **HITS** appeared to converged at the fifth iteration.

**HITS** providesa more descriptive information about why a particular node is important. It tells whether a node is providing information as an authority, or whether it is a hub pointing to the correct authority.

In **PageRank**, *vertex* G has a very small PR score even-though G is not connected to any of the other nodes in the graph. This is due to the $\lambda$ factor or the damping factor which, in our case, was set to 0.15. However, the same *vertex* shows no scores, either as an authority or as a hub.

Figure 2: Result After 5 Iterations for HITS and PageRank Algorithms



(a) HITS 5 Iteration Calculation Results



(b) PageRank 5 Iteration Calculation Results

# 3    Problem 2

Using http://www.cs.cornell.edu/People/tj/svm_light/ work through the "Inductive SVM" example, discuss in detail the steps and resulting output.

The "Inductive SVM" example from [2] used the Reuters-21578 collection. 1000 positive and 1000 negative examples were used for the machine-learning training, in order to determine if a particular document in the

collection should be classified as "corporate acquisition". A test file $< test.dat >$ representing 600 documents, it served as the input to evaluate the accuracy of the created model.

## 3.1   Input File Format

The training file has the form:

```
<line> .=. <target> <feature>:<value> <feature>:<value> ... <feature>:<value> # <info>
<target> .=. +1 | -1 | 0 | <float>
<feature> .=. <integer> | "qid"
<value> .=. <float>
<info> .=. <string>
```
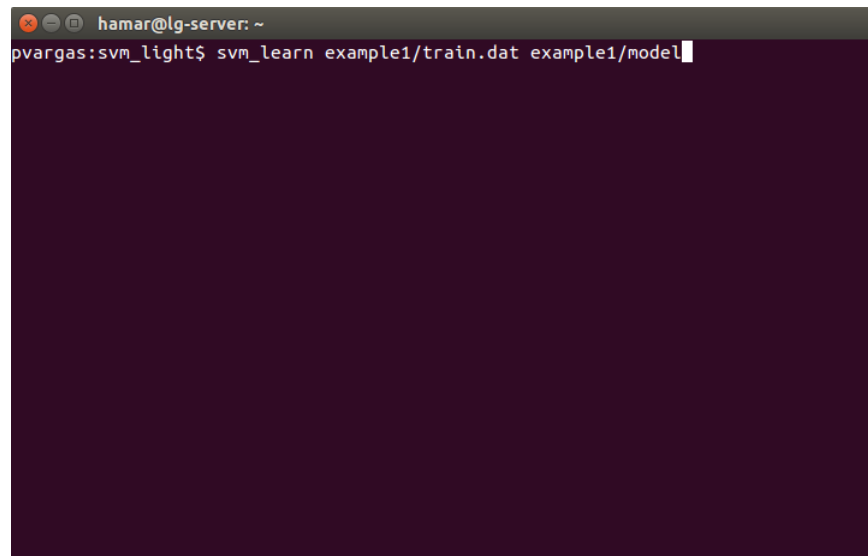
Considering that the **positive** and **negative** examples are included in a single file, in order to distinguish if the training document is for a positive example, the value of **1** is placed at the beginning of the record. On the contrary, a **-1** value indicates the training record is related to a negative example.

In order to create our prediction model, the following instruction should be placed in the command-line:

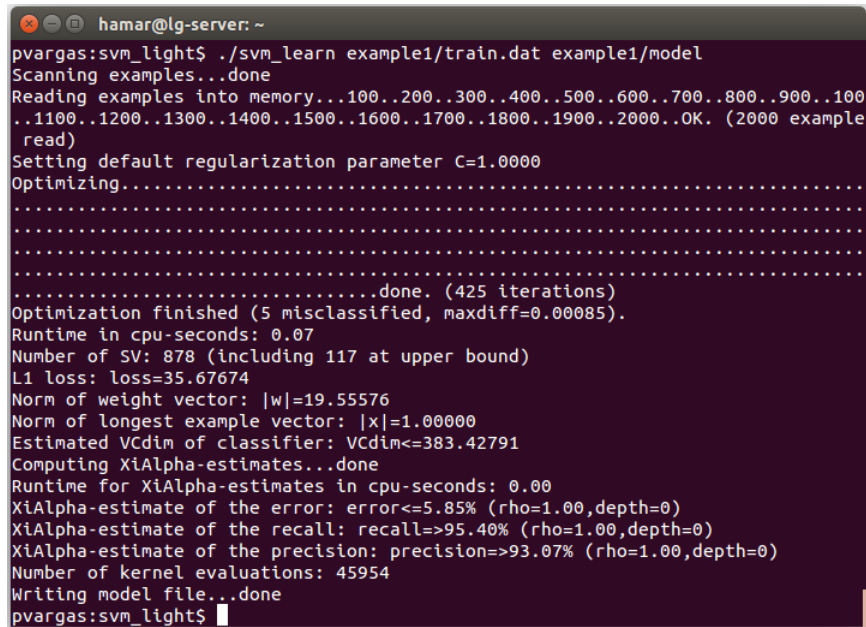**svm_learn example1/train.dat example1/model**

Figure 3: Training Module

## 3.2   Resulting Output

The accuracy of the report is shown on Figure 4

Figure 4: Data Training Result



### 3.2.1   Number of SV

Although the entire collection has **9,947** words, the number of features required to make the classification for any particular document in the collection as a "corporate acquisition" was calculated at: $SV = \mathbf{878}$. Therefore, the model only requires **878** features to obtain our classification.

### 3.2.2   Loss Function

Measurement of empirical error, define the error the classifier made on the training set.

### 3.2.3   Norm of Weight Vector

According to [3], for all document pairs in the rank data, we would like the score for the document with the higher relevance rating (or rank) to be greater than the score for the document with the lower relevance rating. For our data $|x| = 19.55576$.
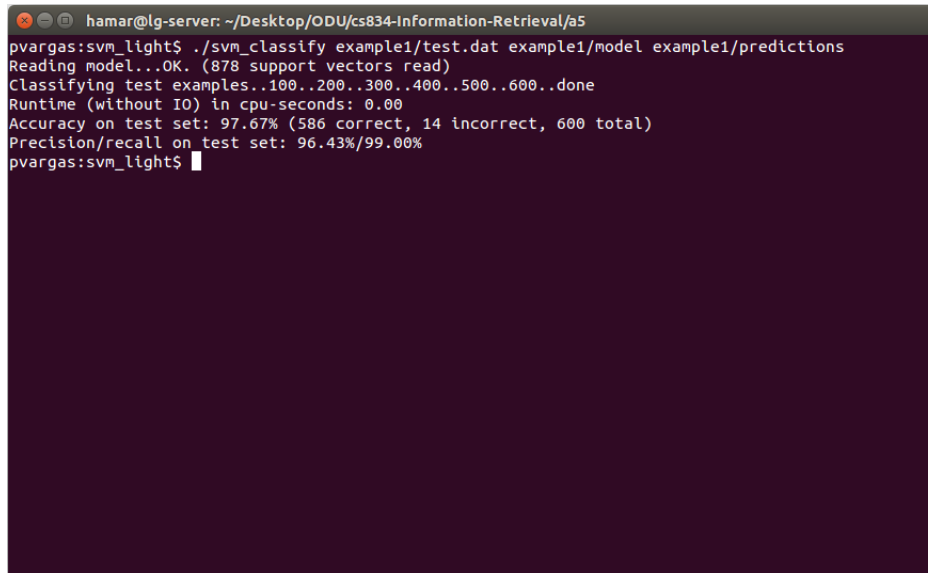
### 3.2.4   VC Dimension Result

The VC dimension or **VapnikChervonenkis dimension** measures the number of points required to shatter the classification model. For this example: $VCdim <= \mathbf{383.42791}$

### 3.2.5   Testing the Model

Looking at Figure 5, we can see that 14 out of 600 documents were incorrectly predicting, bringing the accuracy of this model to **97.6%** accuracy. This error falls into the predicted *XiAlpha-estimate* error of less than 5.85%.

Figure 5: Data Test Result

# 4   Problem 3

- create your own example modeled after the "Inductive SVM" example

- pick a topic (e.g., "Australia") and provide 100 positive and 100 negative examples for training data:

    a. using the Reuters-21578 collection (linked from the SVMlight page)

    b. or, create your own collection with crawled web pages

- pick 30 documents not in the training set for your test data

- stem the words in the collection, using TFIDF as the features (compute for the 230 documents)

- train, classify, and discuss the results

## 4.1   Approach

### 4.1.1   Topic Selection

The collection selected for this exercise was **Reuters-21578**. This collection was accessed via the python library **nltk.corpus**. The topic chosen was "**exports**." In order to obtain 100 positive examples for data training, a *Python* function was creating. The script *svm-light.py* contains this function (*find_features*) which is shown on Figure 3.

### 4.1.2   Selecting Positive Example Documents

The first line for all the documents in the corpus is the title, which is representative its content. The function *find_features* takes for argument a string, which is the term scanned in the document title (line 140). In order to obtain a document with plenty features, a minimum length of 500 characters was required for the document size to be considered for selection. The number of documents in the corpus under this criteria came out to be **136**.

The function returns an array containing the index for all the documents in the corpus where the term appears. The content of those documents were visually inspected for validity, and the first 100 indexes were stored in a file: *positive.dat*. Later on, this file was going to be used as a input for the 100 positive training data examples.

Listing 3: Scan for Feature in Document Title

```
132 def find_features(feature):
133     feature = tokenize(feature)[0]
134     documents = reuters.fileids()
135
136     id_array = []
137     for id in documents:
138         doc = reuters.raw(id)
139         first_line = doc.split('\n')[0]
140         if feature in tokenize(first_line) and len(doc) > 500:
141             print('<----------- %s -------------->' % id)
142             print(doc)
```

```
143            print(len(doc))
144            id_array.append(id)
145
146    print(id_array)
147    print(len(id_array))
148    categories = reuters.categories()
149    print(categories)
150
151    return id_array
```

### 4.1.3  Selecting Negative Example Documents

A similar approach was used to obtain the negative training examples. The term "**bank**" was entered as the feature into the function *find_features*. The result was an array of 229 document indexes. The first 100 documents were viewed for content accuracy. In this case, it was ensured the content was not related to the term **export**. The document indexes were saved into the file: *negative.dat*.

### 4.1.4  Selecting Testing Documents

Considering that only 100 indexes from 136 available were taken to be trained as positive examples, and the next 15 indexes were selected to serve as input for the testing data. The same was done to find negative documents for the testing stage. In total, 30 documents were used for testing: 15 positive and 15 negative. The indexes for those documents were saved under the files: *positive-test.dat* and *negative-test.dat* respectively.

### 4.1.5  Building Word Index

The online resource on [1] was very useful to stem the words in the collection. The PorterStemmer function from the **nltk** library facilitated the stemming action. The function *create_vocabulary* completed the stemming task. The function shown on Listing 4 reads all the documents in the corpus, then it tokenized each of the document and it placed the stemmed word into a set object(lines 305-306). Considering that the set does not have duplicate elements, the resulting set becomes the vocabulary. Its content was written to the file: *words*.

Listing 4: Stemming Words in the Collection

```
301 def create_vocabulary():
302    documents = reuters.fileids()
303    for document_id in documents:
304        text = tokenize(reuters.raw(document_id))
305        for word in text:
306            word_set.add(word)
307
308    with open('words', 'w') as f:
309        for record in word_set:
310            f.write('%s\n' % record)
```

### 4.1.6    Creating Training Feature Records

$SVM^{light}$ requires the training input file in a specific format, as it is shown in section 3.1. The resource cited on [1] contained a function that calculates **TFIDF** which is the value for the features in our training document. This function named *tf_idf* was incorporated into the script *svm-light.py*, and it takes a string as a parameter, which in our case, it is the entire document. This function generates a **TfidVectorizer** object available in the **sklearn.feature_extraction** library.

The function *tf_idf* calls the function *feature_values* prior to passing the calculated vector. This last function takes all elements in the vector and converts them into a set of tuples $(w, v)$, where $w$ is the tokenized word, and $v$ is its calculated **TFIDF** value.

Finally, the function *write_train_data*, shown on Listing 5, for each index marked as positive or negative in the training files, *positive.dat* and *negative.dat* respectively, is opened and passed as an argument to the *tf_idf* (line 113). The returned tuples $(w, v)$ are placed into a dictionary (lines 119-120), where the word $w$ is converted to the index in the vocabulary.

Since $SVM^{light}$ requires the features to be placed in ascending order, the dictionary is sorted by key value prior to writing the features into *train.dat*, which it is the file that it will be used for training data.

Listing 5: Creating Training Data

```
102  def write_train_data(train_file, mode, train_data_file, pos_neg):
103      # get training data
104      train_data = []
105      with open(train_file, mode) as w_file:
106          with open(train_data_file, 'r') as f:
107              for record in f:
108                  record = record.strip()
109                  # include index
110                  train_data.append(record)
111
112                  # get TFIDF features for document
113                  features = tf_idf(reuters.raw(record))
114
115                  unorder_list= {}
116
117                  print(pos_neg, end=' ')
118                  w_file.write('%s ' % pos_neg)
119                  for key, value in features:
120                      unorder_list[word_dic[key]] = value
121                      # print('%d:%f' % (word_dic[key], value), end=' ')
122
123                  for key in sorted(unorder_list):
124                      print('%d:%f' % (key, unorder_list[key]), end=' ')
125                      w_file.write(' %d:%f' % (key, unorder_list[key]))
126
127                  print()
128                  w_file.write('\n')
129      return
```
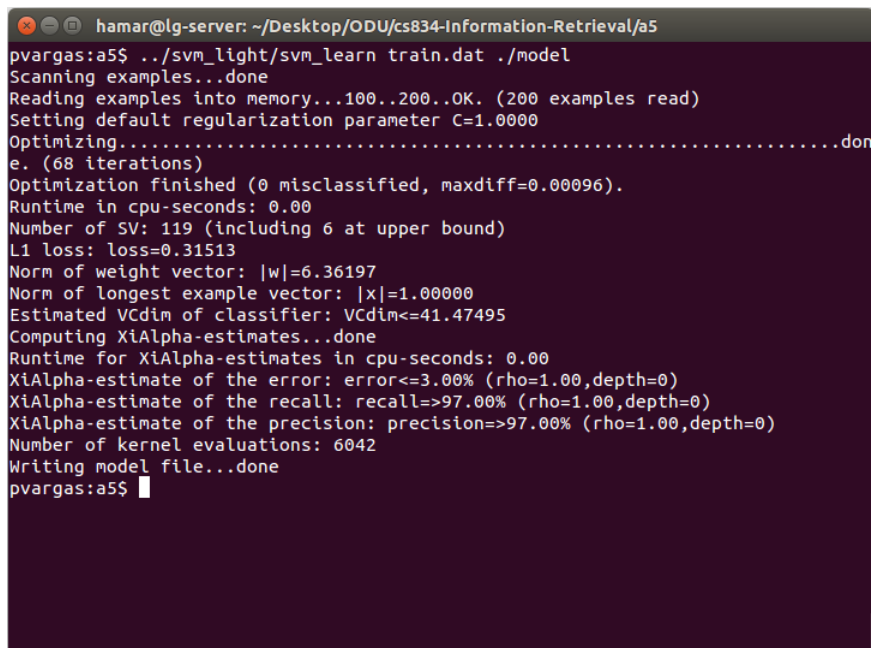
## 4.2    Solution

### 4.2.1    Training $SMV^{light}$

To create our model, we typed the command **svm_learn train.dat model**, as it is shown on Figure 6. The file *train.dat* is the training data, and *model* is the file that resulted from the machine learning. The estimated error is within **3%**, therefore we are expected to have an accuracy greater or equal to **97%**.
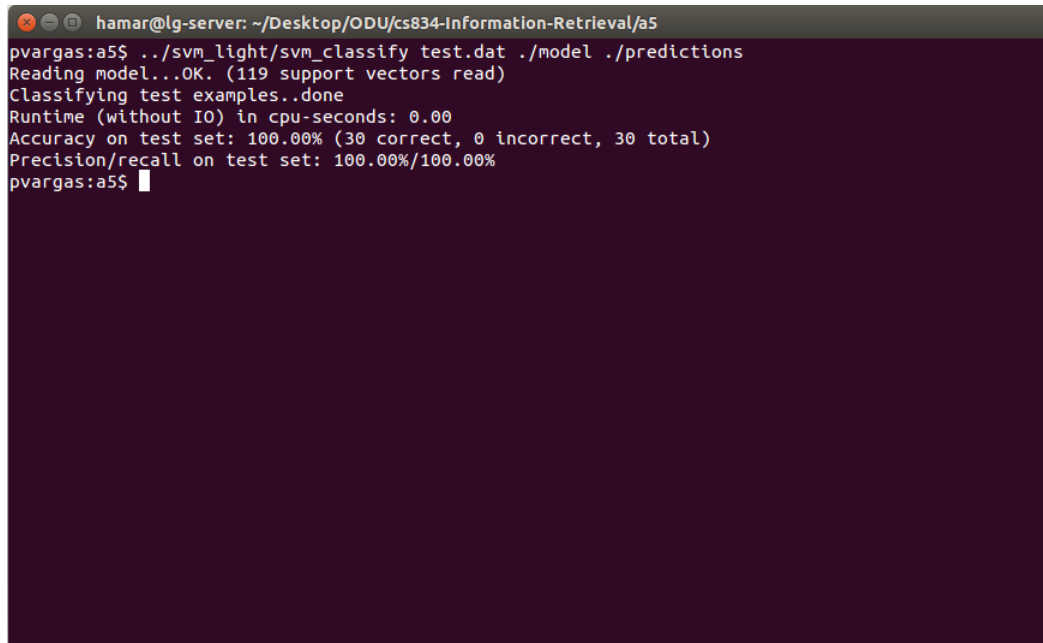
Figure 6: Data Training Result



### 4.2.2    Testing the Model

To run the test, we typed the command **svm_classify test.dat predictions**, as shown on Figure 7. The file *test.dat* contained the test data, and *predictions* is the file that resulted after running the classification application. The accuracy of our test came out to be **100%**, which agrees with the prediction on section 4.2.1, which predicted that the results should be greater than or equal to **97%**.

Figure 7: Testing Results

# References

[1] *Classifying reuters-21578 collection with python. (n.d.) retrieved december 14, 2016, from https://miguel-malvarez.com/2015/03/20/classifying-reuters-21578-collection-with-python-representing-the-data/.*

[2] *Svm light. (n.d.) retrieved december 13, 2016, from http://www.cs.cornell.edu/people/tj/svm_light/.*

[3] T. S. W.B. Croft, D. Metzler, *Search Engine Information Retrieval in Practice*, Pearson Education, 2015.