

DSE210_Final_Exam

March 20, 2017

1 DSE 201 – Final Exam

1.1 Orysya Stus

1.2 3.20.2017

```
In [50]: import pandas as pd
import numpy as np
import sklearn
from sklearn.naive_bayes import BernoulliNB
from sklearn.cluster import KMeans
from sklearn.mixture import GMM
from sklearn import metrics
from sklearn import datasets
import seaborn as sns
import matplotlib.pyplot as plt
import re

%matplotlib inline
```

1.3 For a refresher on Bernoulli Native Bayes review:

http://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.BernoulliNB.html#sklearn.naive_bayes.BernoulliNB
http://scikit-learn.org/stable/auto_examples/text/document_classification_20newsgroups.html#sphx-gl-auto-examples-text-document-classification-20newsgroups-py

1.4 Q4

For this problem, you'll be using the 20 Newsgroups data set. There are several versions of it on the web. You should download "20news-bydate.tar.gz" from <http://qwone.com/~jason/20Newsgroups/>

The same website has a processed version of the data, "20news-bydate-matlab.tgz", that is particularly convenient to use. Download this and also the file "vocabulary.txt". Look at the first training document in the processed set and the corresponding original text document to understand the relation between the two. The words in the documents constitute an overall vocabulary V of size 61188. Build a Bernoulli Naive Bayes model using the training data. Write a routine that uses this naive Bayes model to classify a new document. To avoid underflow, work with logs rather than multiplying together probabilities.

1.4.1 Part (a)

Evaluate the performance of your model on the test data. What error rate do you achieve?

```
In [2]: train = sklearn.datasets.load_files('./20news-bydate-train/', load_content=True)
        test = sklearn.datasets.load_files('./20news-bydate-test/', load_content=True)
```

```
In [3]: vocab = {}
        reverse_vocab = {}
        count = 0
        a = open('./vocabulary.txt', 'r')
        for v in a:
            val = v.strip()
            vocab[val] = count
            reverse_vocab[count] = val
            count += 1
```

```
In [4]: categories = {}
        for i,e in enumerate(test.target_names):
            categories[i] = e
```

```
In [5]: test_map = dict(enumerate(test filenames))
        reverse_test_map = {v:k for k,v in test_map.items() }
```

```
In [6]: from sklearn.feature_extraction.text import CountVectorizer
        count_vect = CountVectorizer(strip_accents='unicode', decode_error = 'ignore')
        X_train_counts = count_vect.fit_transform(train.data)
        training_feature_names = count_vect.get_feature_names()
        X_train_counts.shape
```

```
Out[6]: (11314, 61188)
```

```
In [7]: from sklearn.feature_extraction.text import TfidfTransformer
        tf_transformer = TfidfTransformer()
        X_train_tfidf = tf_transformer.fit_transform(X_train_counts)
        X_train_tfidf.shape
```

```
Out[7]: (11314, 61188)
```

```
In [8]: count_vect = CountVectorizer(strip_accents='unicode', decode_error = 'ignore')
        X_test_counts = count_vect.fit_transform(test.data)
        test_feature_names = count_vect.get_feature_names()
        X_test_counts.shape

        tf_transformer = TfidfTransformer()
        X_test_tfidf = tf_transformer.fit_transform(X_test_counts)
        X_test_tfidf.shape
```

```
Out[8]: (7532, 61188)
```

```
In [9]: clf = BernoulliNB(alpha=1.0).fit(X_train_tfidf, train.target)
        predicted = clf.predict(X_test_tfidf)
```

```
In [10]: print 'The model is', metrics.accuracy_score(test.target, predicted), 'accurate'
         print 'The model has an error rate of', 1- metrics.accuracy_score(test.target, predicted)
```

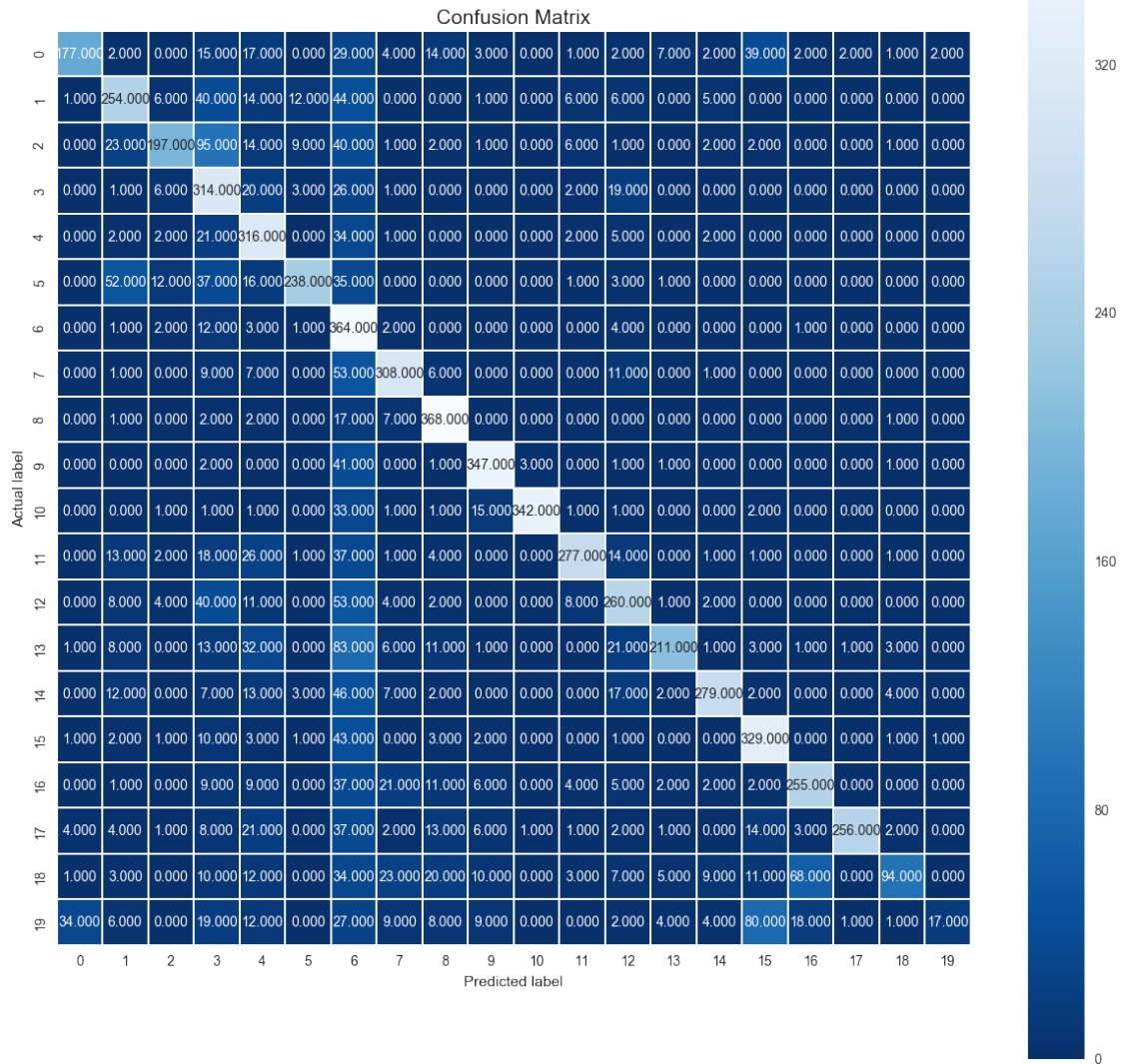
The model is 0.690785979819 accurate

The model has an error rate of 0.309214020181

1.4.2 Part (b)

Evaluate your final model on the test data. Construct a confusion matrix.

```
In [11]: cm = pd.DataFrame(metrics.confusion_matrix(test.target, predicted))
        plt.figure(figsize=(15, 15))
        sns.heatmap(cm, annot=True, fmt=".3f", linewidth=.5, square = True, cmap = 'magma')
        plt.ylabel('Actual label');
        plt.xlabel('Predicted label');
        plt.title('Confusion Matrix', size = 15);
```



1.5 Q5

```
In [21]: x = np.linspace(-10, 4, 10, endpoint=True)
         y = (1 * x + 10) / 2
         fig = plt.figure()
         ax = fig.add_subplot(111)
         ax.plot(x, y)
         ax.spines['right'].set_color('none')
         ax.spines['top'].set_color('none')
         ax.xaxis.set_ticks_position('bottom')
         ax.spines['bottom'].set_position(('data', 0))
```

```
ax.yaxis.set_ticks_position('left')
plt.annotate(r'positive', xy=(-6, 4), xycoords='data', xytext=(+50, +30),
```

Out [21]: <matplotlib.text.Annotation at 0x15df7278>



1.6 Q7

Consider the following observations:

$X = (-0.1, -0.2, 0.1, 0.2, 0, 0.1, -0.1, 0, -0.05, 0.1, 1.05, 1.1, 0.9, 0.8, 0.9, 1, 1.2, 1.1, 1.2, 0.9)$

Cluster this data into 2 classes using the K-means algorithm. What are the cluster center?

<http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>

```
In [13]: X = [-0.1, -0.2, 0.1, 0.2, 0, 0.1, -0.1, 0, -0.05, 0.1, 1.05, 1.1, 0.9, 0.8, 0.9, 1, 1.2, 1.1, 1.2, 0.9]
X = np.array(X)
kmeans = KMeans(n_clusters=2, random_state=0)
kmeans.fit(X.reshape(-1,1))
```

```
Out [13]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
n_clusters=2, n_init=10, n_jobs=1, precompute_distances='auto',
random_state=0, tol=0.0001, verbose=0)
```

```
In [14]: kmeans.cluster_centers_
```

```
Out [14]: array([[ 1.015],
[ 0.005]])
```

1.7 Q8

Construct a Gaussian Mixture Model for the above data. What is your estimate for the number of mixtures?

http://www.astroml.org/book_figures/chapter4/fig_GMM_1D.html

<http://scikit-learn.org/stable/modules/generated/sklearn.mixture.GaussianMixture.html#sklearn.mixture>

```
In [56]: import warnings
          warnings.filterwarnings('ignore')

          for i in xrange(10):
              gmm = GMM(n_components=i+1)
              gmm.fit(X.reshape(-1,1))
              print 'Gaussian mixture model for', i+1, 'components, akaike information criterion (aic) for the'
              print 'Gaussian mixture model for', i+1, 'components, bayesian information criterion (bic) for the'
              print 'Gaussian mixture model for', i+1, 'components, means or estimates are:'

Gaussian mixture model for 1 components, akaike information criterion (aic) for the
Gaussian mixture model for 1 components, bayesian information criterion (bic) for the
Gaussian mixture model for 1 components, means or estimates are:
[[ 0.51]]
Gaussian mixture model for 2 components, akaike information criterion (aic) for the
Gaussian mixture model for 2 components, bayesian information criterion (bic) for the
Gaussian mixture model for 2 components, means or estimates are:
[[ 1.015]
 [ 0.005]]
Gaussian mixture model for 3 components, akaike information criterion (aic) for the
Gaussian mixture model for 3 components, bayesian information criterion (bic) for the
Gaussian mixture model for 3 components, means or estimates are:
[[ 0.99375351]
 [ 0.005      ]
 [ 1.02049175]]
Gaussian mixture model for 4 components, akaike information criterion (aic) for the
Gaussian mixture model for 4 components, bayesian information criterion (bic) for the
Gaussian mixture model for 4 components, means or estimates are:
[[ 0.99765427]
 [ 0.00269455]
 [ 1.02144721]
 [ 0.01174379]]
Gaussian mixture model for 5 components, akaike information criterion (aic) for the
Gaussian mixture model for 5 components, bayesian information criterion (bic) for the
Gaussian mixture model for 5 components, means or estimates are:
[[ 0.0100255 ]
 [ 1.01030616]
 [ 0.00269286]
 [ 0.9902593  ]
 [ 1.02263586]]
Gaussian mixture model for 6 components, akaike information criterion (aic) for the
Gaussian mixture model for 6 components, bayesian information criterion (bic) for the
```

Gaussian mixture model for 6 components, means or estimates are:

```
[[ 1.02250557]
 [ 0.00576591]
 [ 0.99368412]
 [ 0.01246569]
 [ 1.01279667]
 [ 0.00228867]]
```

Gaussian mixture model for 7 components, akaike information criterion (aic) for the

Gaussian mixture model for 7 components, bayesian information criterion (bic) for the

Gaussian mixture model for 7 components, means or estimates are:

```
[[ 0.00189    ]
 [ 1.01307606]
 [ 0.99493135]
 [ 0.00571948]
 [ 0.0161848  ]
 [ 0.01059467]
 [ 1.02303845]]
```

Gaussian mixture model for 8 components, akaike information criterion (aic) for the

Gaussian mixture model for 8 components, bayesian information criterion (bic) for the

Gaussian mixture model for 8 components, means or estimates are:

```
[[ 0.99831262]
 [ 0.00780372]
 [ 1.01345556]
 [ 0.00193757]
 [ 0.01177073]
 [ 0.00493955]
 [ 1.02313641]
 [ 0.01784174]]
```

Gaussian mixture model for 9 components, akaike information criterion (aic) for the

Gaussian mixture model for 9 components, bayesian information criterion (bic) for the

Gaussian mixture model for 9 components, means or estimates are:

```
[[ 0.01164422]
 [ 1.01509937]
 [ 0.00778792]
 [ 1.00085625]
 [ 0.00459122]
 [ 1.02451423]
 [ 0.00179331]
 [ 0.01629944]
 [ 0.99233726]]
```

Gaussian mixture model for 10 components, akaike information criterion (aic) for the

Gaussian mixture model for 10 components, bayesian information criterion (bic) for the

Gaussian mixture model for 10 components, means or estimates are:

```
[[ 0.00447655]
 [ 1.01746192]
 [ 0.01119657]
 [ 0.99794265]
 [ 0.00720346]
```

```
[ 1.02433142]
[ 0.00130828]
[ 1.01137019]
[ 0.98993299]
[ 0.01605409]]
```

1.7.1 2 components were used like K-means, showing that the Gaussian mixture model means or estimates are the same as those found through K-means clustering. Also, **n_components=2** has the lowest bic and aic therefore providing the best means.

In []: