

OLD DOMINION UNIVERSITY

Assignment 8

David Bayard

April 13, 2019

QUESTION 1.

Create two datasets; the first called Testing, the second called Training.

Solution:

Two datasets were created by copying and pasting the text from 20 emails from the Spam folder, and another 20 from the inbox. Specific emails that contained certain query terms such as "ODU", or "Monarchs" were selected to increase the odds of the classifier selecting the correct class per each document. Moreover, most of the emails shared a similar topic, being school related. This was purposefully done to test the classifier to a fuller extent.

The emails are split into two directories, being Spam and Not Spam. Within these directories, another two directories exist, called Test and Train. These folders represent the testing and training sets, containing 20 emails each. First, the training sets of Spam and Not Spam text files are run through the classifier, then the text files in the test directories are fed to the classifier.

QUESTION 2

Using the PCI book modified docclass.py code and test.py (see Slack assignment-8 channel) Use your Training dataset to train the Naive Bayes classifier (e.g., docclass.spamTrain()) Use your Testing dataset to test (test.py) the Naive Bayes classifier and report the classification results.

Solution:

As specified above, the PCI book implements the Naive Bayes classifier functionality. Changes were made to the docclass.py file in order to read from the.txt files, instead of the SQLite database. Other than these changes, the original functions supplied by the PCI book were used as supplied.

In order to train the Naive Bays classifier, three lines of code needed to be executed. First, an instance of the naivebays is instantiated and holds a reference to the getwords function. This is important because this function will later be used to extract all of the text from each of the .txt files, returning a lowercase count of each term.

```
1 cl=docclass.naivebays(docclass.getwords)
2
3 docclass.spamTrain(cl)
4 docclass.testClassifier(cl)
```

Listing 1: test.py driver function

When the naivebays object is instantiated, a classifier object is created, and passed a reference the getwords function. In the function below, the spam files numbered from 1-11 are opened, and read into the InputStream, and then into strings. Every iteration, these strings are passed to the classifiers train method, which converts the words to lowercase, and keeps count of how many documents contain each word, and in which class they fall into. Sample output of this classification is shown below.

```

1
2 for r in range(1,11):
3     spamFileName = "spamFile" + str(r)
4     notSpamFileName = "notSpam" + str(r)
5
6     spamFile = open("./Spam/Train/spam" + str(r) + ".txt", "r")
7     notSpamFile = open("./Not_Spam/Train/NotSpam" + str(r) + ".txt", "r")
8     text = spamFile.read()
9     text2 = notSpamFile.read()
10
11     cl.train(text2, "NotSpam")
12     cl.train(text, 'spam')

```

Listing 2: Pass file input stream as string to classifier

```
david@david-Aspire-E5-575:~/Documents/CS532/CS532_Assignment8ToPush
File Edit View Search Terminal Tabs Help

david@david-Aspire-E5-575:~/Documents/CS532/CS532_Assignment8ToPush
david@david-Aspire-E5-575:~/Documents/CS530/Project/Avocado2

{ dear: {'NotSpan': 7, 'span': 1}, colleagues: {'NotSpan': 1}, write: {'NotSpan': 1}, first: {'NotSpan': 1, 'span': 4}, invite: {'NotSpan': 2}, your: {'NotSpan': 2, 'span': 10}, participation: {'NotSpan': 1}, creating: {'NotSpan': 1, 'span': 1}, exhibit: {'NotSpan': 1}, featuring: {'NotSpan': 1}, the: {'NotSpan': 1, 'span': 10}, people: {'NotSpan': 2}, odu: {'NotSpan': 6}, photographer: {'NotSpan': 1}, glen: {'NotSpan': 1}, mclure: {'NotSpan': 1}, will: {'NotSpan': 1, 'span': 7}, span: 3}, campus: {'NotSpan': 4}, span: 1}, during: {'NotSpan': 3}, three: {'NotSpan': 2}, days: {'NotSpan': 1, 'span': 1}, april: {'NotSpan': 5}, span: 1}, photograph: {'NotSpan': 1}, students: {'NotSpan': 1}, span: 3}, staff: {'NotSpan': 1}, and: {'NotSpan': 10}, span: 10}, faculty: {'NotSpan': 1}, span: 2}, outdoors: {'NotSpan': 1}, random: {'NotSpan': 1, 'span': 1}, his: {'NotSpan': 1, 'span': 2}, distinctive: {'NotSpan': 1}, street: {'NotSpan': 2}, portrait: {'NotSpan': 1}, style: {'NotSpan': 1}, resulting: {'NotSpan': 1}, reflect: {'NotSpan': 2}, rich: {'NotSpan': 1}, diversity: {'NotSpan': 2}, shown: {'NotSpan': 1}, gordon: {'NotSpan': 1}, art: {'NotSpan': 1}, galleries: {'NotSpan': 1}, coming: {'NotSpan': 1, 'span': 2}, academic: {'NotSpan': 1}, span: 1}, year: {'NotSpan': 4}, span: 1}, please: {'NotSpan': 6}, span: 2}, look: {'NotSpan': 1}, for: {'NotSpan': 9}, span: 8}, camera: {'NotSpan': 1}, these: {'NotSpan': 3}, span: 1}, windows: {'NotSpan': 1}, time: {'NotSpan': 6}, span: 5}, kaufman: {'NotSpan': 1}, mall: {'NotSpan': 1}, activity: {'NotSpan': 1}, hour: {'NotSpan': 1, 'span': 1}, rain: {'NotSpan': 1}, date: {'NotSpan': 3}, needed: {'NotSpan': 1, 'span': 1}, follow: {'NotSpan': 1, 'span': 1}, arts: {'NotSpan': 1}, facebook: {'NotSpan': 1}, twitter: {'NotSpan': 1}, any: {'NotSpan': 1, 'span': 3}, scheduling: {'NotSpan': 1}, changes: {'NotSpan': 1}, secondly: {'NotSpan': 1}, virginia: {'NotSpan': 1}, festival: {'NotSpan': 1}, are: {'NotSpan': 8}, span: 6}, pleased: {'NotSpan': 2}, offer: {'NotSpan': 1, 'span': 2}, limited: {'NotSpan': 1}, number: {'NotSpan': 1}, free: {'NotSpan': 1, 'span': 3}, tickets: {'NotSpan': 2}, select: {'NotSpan': 2}, performance: {'NotSpan': 1}, attached: {'NotSpan': 1}, find: {'NotSpan': 2}, flyer: {'NotSpan': 2}, containing: {'NotSpan': 1}, details: {'NotSpan': 2}, about: {'NotSpan': 5}, span: 5}, visit: {'NotSpan': 2}, www: {'NotSpan': 1, 'span': 3}, vafest: {'NotSpan': 1}, org: {'NotSpan': 1, 'span': 2}, information: {'NotSpan': 4}, span: 1}, events: {'NotSpan': 2}, enter: {'NotSpan': 2}, win: {'NotSpan': 2}, 100: {'NotSpan': 3}, amazon: {'NotSpan': 1}, gift: {'NotSpan': 1}, card: {'NotSpan': 1}, answering: {'NotSpan': 1}, quick: {'NotSpan': 1}, survey: {'NotSpan': 1}, thank: {'NotSpan': 2}, NotSpan: 1}, you: {'NotSpan': 10}, NotSpan: 6}, interest: {'NotSpan': 3}, foxit: {'NotSpan': 2}, value: {'NotSpan': 2}, opinion: {'NotSpan': 1}, want: {'NotSpan': 2}, NotSpan: 2}, provide: {'NotSpan': 1}, with: {'NotSpan': 9}, NotSpan: 6}, best: {'NotSpan': 5}, service: {'NotSpan': 1}, NotSpan: 1}, help: {'NotSpan': 4}, NotSpan: 1}, taking: {'NotSpan': 3}, part: {'NotSpan': 2}, NotSpan: 2}, our: {'NotSpan': 5}, NotSpan: 4}, minute: {'NotSpan': 2}, answer: {'NotSpan': 1}, some: {'NotSpan': 2}, NotSpan: 2}, questions: {'NotSpan': 2}, NotSpan: 1}, experience: {'NotSpan': 5}, NotSpan: 1}, phantompdf: {'NotSpan': 2}, token: {'NotSpan': 1}, gratitude: {'NotSpan': 1}, one: {'NotSpan': 3}, NotSpan: 1}, ten: {'NotSpan': 1}, cards: {'NotSpan': 1}, advance: {'NotSpan': 2}, feedback: {'NotSpan': 2}, sincerely: {'NotSpan': 1}, NotSpan: 3}, team: {'NotSpan': 3}, NotSpan: 2}, give: {'NotSpan': 1}, business: {'NotSpan': 4}, email: {'NotSpan': 1, 'span': 1}, NotSpan: 2}, old: {'NotSpan': 2}, dominion: {'NotSpan': 2}, community: {'NotSpan': 5}, span: 2}, norfolk: {'NotSpan': 1}, police: {'NotSpan': 1}, department: {'NotSpan': 3}, investigating: {'NotSpan': 1}, shooting: {'NotSpan': 1}, reported: {'NotSpan': 1}, have: {'NotSpan': 6}, span: 1}, occurred: {'NotSpan': 1}, around: {'NotSpan': 1, 'span': 1}, yesterday: {'NotSpan': 1}, 800: {'NotSpan': 1, 'span': 1}, block: {'NotSpan': 3}, west: {'NotSpan': 1}, 45th: {'NotSpan': 1}, near: {'NotSpan': 1, 'span': 1}, intersection: {'NotSpan': 1}, colley: {'NotSpan': 1}, avenue: {'NotSpan': 1}, man: {'NotSpan': 1}, was: {'NotSpan': 3}, transported: {'NotSpan': 1}, sentara: {'NotSpan': 1}, general: {'NotSpan': 1}, hospital: {'NotSpan': 1}, non: {'NotSpan': 1}, life: {'NotSpan': 1}, threatening: {'NotSpan': 1}, gunshot: {'NotSpan': 1}, wounds: {'NotSpan': 1}, residents: {'NotSpan': 1}, who: {'NotSpan': 3}, span: 2}, live: {'NotSpan': 2}, span: 1}, this: {'NotSpan': 7}, span: 7}, area: {'NotSpan': 1}, can: {'NotSpan': 8}, span: 6}, expect: {'NotSpan': 1}, see: {'NotSpan': 6}, span: 2}, increased: {'NotSpan': 1}, presence: {'NotSpan': 2}, ask: {'NotSpan': 1, 'span': 1}, anyone: {'NotSpan': 1, 'span': 1}, incident: {'NotSpan': 1}, call: {'NotSpan': 2}, span: 1}, crime: {'NotSpan': 1}, line: {'NotSpan': 1}, 888: {'NotSpan': 1}, 562: {'NotSpan': 1}, 5887: {'NotSpan': 1}, lock: {'NotSpan': 1}, submit: {'NotSpan': 2}, tip: {'NotSpan': 2}, through: {'NotSpan': 3}, span: 1}, mobile: {'NotSpan': 1, 'span': 1}, app: {'NotSpan': 1, 'span': 1}, leading: {'NotSpan': 1, 'span': 1}, investigation: {'NotSpan': 1}, further: {'NotSpan': 1}, being: {'NotSpan': 1}, released: {'NotSpan': 1}, ann: {'NotSpan': 1}, hughes: {'NotSpan': 1}, public: {'NotSpan': 1, 'span': 1}, officer: {'NotSpan': 1}, buy: {'NotSpan': 1}, more: {'NotSpan': 8}, NotSpan: 1}, popular: {'NotSpan': 3}, higher: {'NotSpan': 2}, rated: {'NotSpan': 1}, lower: {'NotSpan': 1}, cost: {'NotSpan': 1}, pdf: {'NotSpan': 1, 'span': 1}, image: {'NotSpan': 1}, NotSpan: 1}, editors: {'NotSpan': 1}, today: {'NotSpan': 3}, off: {'NotSpan': 2}, coupon: {'NotSpan': 1}, works: {'NotSpan': 1}, NotSpan: 2}, studio: {'NotSpan': 2}, photo: {'NotSpan': 1}, mac: {'NotSpan': 1}, standard: {'NotSpan': 1}, sale: {'NotSpan': 1}, ends: {'NotSpan': 1}, act: {'NotSpan': 1}, now: {'NotSpan': 5}, NotSpan: 3}, here: {'NotSpan': 3}, NotSpan: 3}, code: {'NotSpan': 2}, mhq334: {'NotSpan': 1}, nlj8855x2t2szd2ul3: {'NotSpan': 1}, use: {'NotSpan': 4}, NotSpan: 3}, online: {'NotSpan': 6}, NotSpan: 2}, store: {'NotSpan': 1}, http: {'NotSpan': 2}, foxitsoftware: {'NotSpan': 1}, com: {'NotSpan': 1}, NotSpan: 1}, shipping: {'NotSpan': 1}, then: {'NotSpan': 2}, NotSpan: 2}, into: {'NotSpan': 5}, field: {'NotSpan': 1}, NotSpan: 1}, marked: {'NotSpan': 1}, promotional: {'NotSpan': 1}, click: {'NotSpan': 1}, NotSpan: 3}, apply: {'NotSpan': 2}, NotSpan: 1}, lets: {'NotSpan': 1}, edit: {'NotSpan': 2}, words: {'NotSpan': 1}, content: {'NotSpan': 1}, pages: {'NotSpan': 1}, images: {'NotSpan': 1}, NotSpan: 1}, file: {'NotSpan': 1}, convert: {'NotSpan': 1}, files: {'NotSpan': 1}, high: {'NotSpan': 2}, quality: {'NotSpan': 2}, word: {'NotSpan': 1}, powerpoint: {'NotSpan': 1}, excel: {'NotSpan': 2}, repurpose: {'NotSpan': 1}, create: {'NotSpan': 1}, NotSpan: 2}, fill: {'NotSpan': 1}, NotSpan: 1}, out: {'NotSpan': 2}, NotSpan: 1}, save: {'NotSpan': 3}, sign: {'NotSpan': 1}, NotSpan: 1}, send: {'NotSpan': 1}, dynamic: {'NotSpan': 3}, forms: {'NotSpan': 1}, NotSpan: 1}, publish: {'NotSpan': 1}, share: {'NotSpan': 3}, NotSpan: 1}, members: {'NotSpan': 1}, partners: {'NotSpan': 1}, get: {'NotSpan': 5}, NotSpan: 1}, work: {'NotSpan': 2}, done: {'NotSpan': 1}, faster: {'NotSpan': 1}, purchase: {'NotSpan': 2}, regards: {'NotSpan': 1}, software: {'NotSpan': 3}, solutions: {'NotSpan': 1}, fast: {'NotSpan': 2}, NotSpan: 1}, affordable: {'NotSpan': 1}, secure: {'NotSpan': 2}, discount: {'NotSpan': 1}, NotSpan: 1}, applies: {'NotSpan': 1}, licenses: {'NotSpan': 1}, only: {'NotSpan': 2}, NotSpan: 2}, order: {'NotSpan': 1}, resellers: {'NotSpan': 1}, distributors: {'NotSpan': 1}, not: {'NotSpan': 2}, NotSpan: 2}, eligible: {'NotSpan': 1}, cannot: {'NotSpan': 1}, combined: {'NotSpan': 1}, other: {'NotSpan': 5}, NotSpan: 1}, offers: {'NotSpan': 1}, terms: {'NotSpan': 1}, conditions: {'NotSpan': 1}, may: {'NotSpan': 2}, NotSpan: 1}, heart: {'NotSpan': 1}, career: {'NotSpan': 1, 'span': 5}, training: {'NotSpan': 1, 'span': 1}, center: {'NotSpan': 3}, experienced: {'NotSpan': 1}, dedicated: {'NotSpan': 1}, edu: {'NotSpan': 1, 'span': 5},
```

(a) Extracting Text

As seen in the image above, the words from all of the .txt files read in the function above are being extracted, transformed to lowercase, and counted. After the loop iterating over the files finishes running, the classifier object will contain each word in each document, a count of that word per document, and the class to which the word belongs. This data will be used when running the test files through the classifier to determine the correct class.

With the training data loaded into the classifier object, the testing data may now be loaded using the classify method in order to determine which class the test document belongs to. The code below loops over the remaining 20 documents, and predicts whether each document is spam or not spam.

```
1 for r in range(11,21):
2
3     spamFileName = "spamFile" + str(r)
4     notSpamFileName = "notSpam" + str(r)
5
6     spamFile = open("./Spam/Test/spam" + str(r) + ".txt", "r")
7     notSpamFile = open("./Not_Spam/Test/NotSpam" + str(r) + ".txt", "r")
8     text = spamFile.read()
9     text2 = notSpamFile.read()
10
11     predictNotSpam = cl.classify(text2, default='unknown')
12     predictSpam = cl.classify(text, default='unknown')
```

Listing 3: Clustering into 5 groups

The results of the classify function, as well as a title line from each document are recorded in the documentation.txt file. From reading the documentation.txt file, it is visible that the classifier issued one false positive, classifying notSpam15.txt as spam.

Overall, the classifier performed well, labeling 19 out of 20 documents in the proper class, but ranking the notSpam document as spam is a big mistake. This would be fixed by increasing the threshold for a document to be classified as spam. While this may increase the odds of labeling spam documents as not spam, it helps prevent documents that are not spam from being labeled as spam.

QUESTION 3

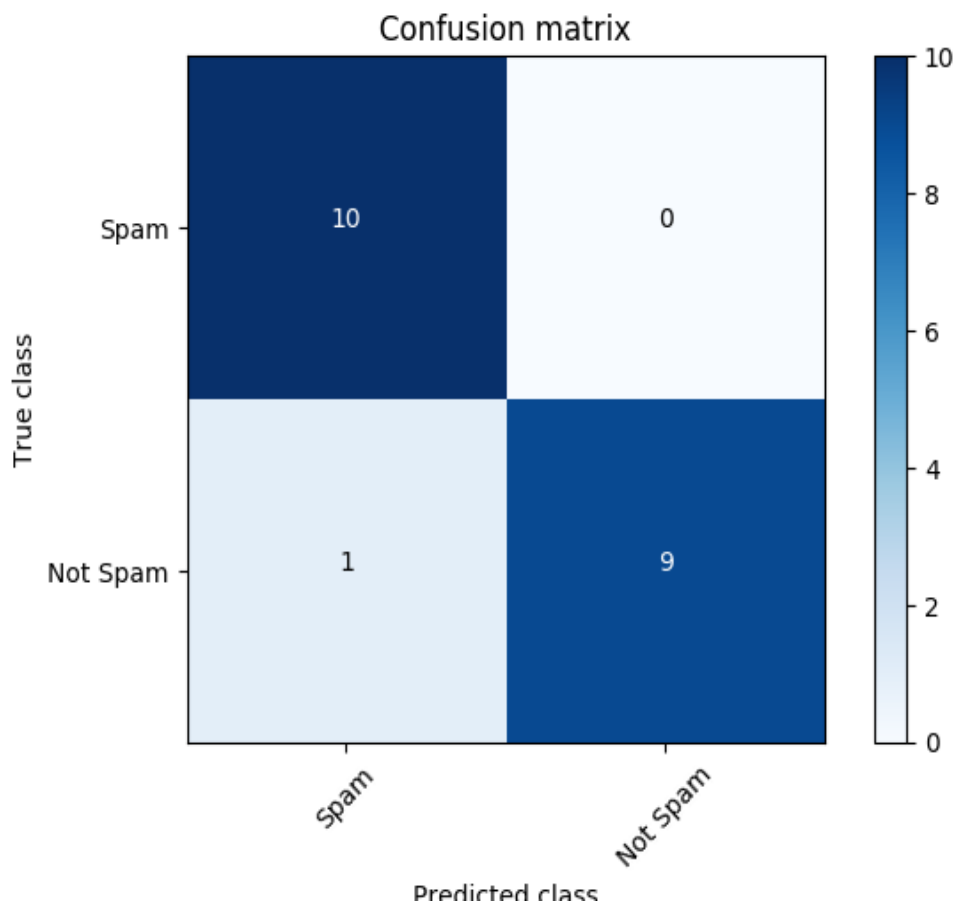
Draw a confusion matrix for your classification results

Solution:

The sklearn library is used here to draw the confusion matrix. It takes the true class values along with the predicted class values and draws the confusion matrix.

In the confusion matrix shown below, the top left corner represents the properly selected spam class, and the bottom left row represents false positives of spam. Thus out of the 10 spam documents, 10 were properly predicted as spam, but 1 document was classified as spam when it is not spam.

The bottom right corner represents the documents that are not spam. The classifier predicted only 9 documents as not spam, and the top right row displaying false negatives remains empty.



QUESTION 4

Report the precision and accuracy scores of your classification results

Solution:

The formulas necessary to calculate the accuracy and precision are displayed below.

$$Precision = \frac{tp}{tp + fp}$$
$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn}$$

The equations above represent tp as true positive, fp = false positive, tn as true negative, and fn as false negative. The equations are filled in with the values from the confusion matrix.

$$Precision = \frac{10}{10 + 1} = \frac{10}{11} = .91$$
$$Accuracy = \frac{10 + 9}{9 + 10 + 1 + 0} = \frac{19}{20} = .95$$

As seen above, the classifier has a fairly high precision and accuracy. These values would undoubtedly change given a larger sample size, but for a sample of only 20 files, the values are reasonable.