

Lab 6, Information Retrieval and Text Mining

Due date: Monday, December 4, 2:00pm

Lab Assignment

This is a pair programming lab.

In this assignment you are asked to analyze a collection of text documents, the **Reuter 50-50** dataset that is often used in text mining and text classification tasks. You will create vector space representations of the documents from this dataset, and will use these representations to conduct a text mining study.

Determine if classification algorithms you implemented in Lab 3, namely, the K-Nearest Neighbors algorithm and Random Forests, can be used to reliably establish the authorship of the text documents in the dataset.

Notice that you get to reuse quite a bit of the Lab 3 code in this assignment.

The Data Collection: Reuter 50-50 Datasets

Reuter 50-50 collection of text documents¹ is a well-known and well-studied dataset often used in machine learning and information retrieval coursework, as well as research.

The dataset consists of a collection of news stories published by the **Reuters** news agencies. The dataset is broken into two parts called **training set** and **test set**. For our lab, the intents of these two parts of the dataset are not important and we will use both parts of the dataset together.

The dataset was constructed to study machine learning algorithms for authorship attribution. It consists of a selection of 50 authors who published news stories with **Reuters**. For each author, exactly 100 news stories they authored is placed in the dataset. 50 of the stories are placed in the **training set** part of the dataset and the remaining 50 – in the **test set**.

The dataset is available both from Lab 5 course web page, as well as from the University of California at Irvine (UCI) Machine Learning Datasets repository as a single zip file named **C50.zip**. When the file is unzipped, it creates the following directory structure:

```
dekhtyar@csc1nx11:~/.../C50 $ ls -al
total 8056
drwx-----.  4 dekhtyar domain^users    4096 Apr  7 23:50 .
drwx-----.  5 dekhtyar domain^users    4096 Apr  7 23:48 ..
drwx-----. 52 dekhtyar domain^users    4096 Apr  7 23:50 C50test
drwx-----. 52 dekhtyar domain^users    4096 Apr  7 23:50 C50train
-rw-----.  1 dekhtyar domain^users 8194031 Apr  7 22:15 C50.zip
```

```
dekhtyar@csc1nx11:~/.../C50 $ ls C50test/
AaronPressman    JaneMacartney    LydiaZajc        RobinSidel
AlanCrosby       JanLopatka       LynneO'Donnell   RogerFillion
AlexanderSmith   JimGilchrist     LynnleyBrowning  SamuelPerry
BenjaminKangLim  JoeOrtiz         MarcelMichelson  SarahDavison
BernardHickey    JohnMastrini     MarkBendeich     ScottHillis
BradDorfman      JonathanBirt     MartinWolk       SimonCowell
DarrenSchuettler JoWinterbottom   MatthewBunce     TanEeLyn
DavidLawder      KarlPenhaul      MichaelConnor    TheresePoletti
EdnaFernandes    KeithWeir        MureDickie       TimFarrand
EricAuchard      KevinDrawbaugh   NickLouth        ToddNissen
```

¹https://archive.ics.uci.edu/ml/datasets/Reuter_50_50

FumikoFujisaki	KevinMorrison	PatriciaCommins	WilliamKazer
GrahamEarnshaw	KirstinRidley	PeterHumphrey	
HeatherScoffield	KouroshKarimkhany	PierreTran	

The C50train directory has exactly the same structure.

Each directory inside C50train and C50test directories bears the name of one of the 50 authors. Inside each directory is 50 .txt files containing the stories written by that author. Each story is in a separate file whose name follows the following pattern:

<number>newsML.txt

where <number> is a 5- or 6-digit number representing the unique ID of the story².

The Task

Preliminary Steps

The first task for your assignment is to create vectorized tf-idf representations of each document and to store these representations.

Write a program `textVectorizer.java` or `textVectorizer.py`³ that takes as input the location of the root of the Reuters 50-50 dataset directory, and produces, as the result, a file containing the vectorized tf-idf representations for each of the 5000 documents in the dataset. The name of the output file can be an input parameter to your program (for the sake of flexibility). Additionally, your program shall produce (to make your life easier during the evaluation stages) a **ground truth** file that for each document (identified by its file name) stores the name of the author (name of the directory the file is in). The ground truth file can have any format you want, but a simple CSV file looking as follows will suffice:

```
421829newsML.txt,AaronPressman
424074newsML.txt,AaronPressman
...
236352newsML.txt,AlanCrosby
293241newsML.txt,AlanCrosby
...
```

Note: You must write the entire vectorizer from scratch without the use of text vectorization tools/APIs/functionality available in the specialized packages in the programming language of your choice. You are explicitly prohibited from using any `scikit learn` or `nltk` Python package functionality for this lab **with the exception of nltk's stemmers** (as well as their counterparts in other languages). The point of this assignment is to learn **how to build** this functionality.

You can use standard parsing techniques available to you in your programming languages such as `split()` and `strip()` methods. You can take advantage of the data structures provided to you by the NumPy package (or similar packages in other programming languages).

When implementing vectorized representations of the documents, please observe that the overall vocabulary of the Reuter 50-50 dataset is significantly larger than the number of unique words used in any specific single document, and therefore, the tf-idf vectors of individual documents will be rather sparse.

To support your vectorization efforts, and other tasks of this assignment, you will implement a **Vector** class (or a **Vector** ADT) in your host programming language. Your **Vector** class will store sparse tf-idf vector representations of text documents. It shall also implement at least two similarity score computations for a pair of vectors: **cosine** similarity and **okapi**.

²The Reuter 50-50 dataset was built out of a much larger data set of Reuters stories, so there is no rhyme or reason to the numbers contained in the filenames. All you need to know is that they are all unique.

³As usual, if you are using a different programming language, name your program accordingly.

Authorship Attribution Through Classification

You will use the K-Nearest Neighbors classifier and the Random Forest classifier to determine if they can help you establish the authorship of the articles.

K-Nearest Neighbors requires a distance or similarity metric to properly operate. You will use the cosine similarity and `okapi` metrics in this exercise.

Write two programs: `knnAuthorship.java/knnAuthorship.py` and `RFAuthorship.java/RFAuthorship.py` which take as input

- the file of vectorized document representations,
- a flag indicating the similarity metric to be used in the computation,
- the appropriate to each method parameters:
 - for KNN: a value of k (number of nearest neighbors to check)
 - for Random Forests: values for (a) number of decision trees to build, (b) number of attributes in a decision tree, and (c) number of data points to be used to construct a decision tree, as well as (d) any parameters you need for the C45 algorithm itself (e.g., the threshold).
- if needed, any other input parameters (document in your `README` file)

The output of each program shall be an authorship label predicted for each of the documents in the Reuters 50-50 dataset.

Your KNN implementation shall use an *all-but-one* validation (take a document, find k nearest neighbors among the remaining 4999 documents in the dataset). The output can be printed out directly to the screen, or placed in an output file.

Your Random Forest implementation does not need a cross-validation in this case (using subsets of the entire dataset to build individual trees in the forest is sufficient to prevent overfit), so, it should simply build the forest and make the predictions.

However, for each of the methods you can consider running the classifier in one of two ways:

- **Approach 1: 50-class classification problem.** In this case, you run a single classification task that has a class variable consisting of 50 labels (individual authors).
- **Approach 2: one-vs-the-world classification problem.** A one-vs-the-world approach is to create a binary classification problem for each of the fifty authors in the dataset. You will wind up with fifty binary classification problems where you will build a model to recognize a specific author vs simply recognizing that a document is NOT written by that author.

You are **not required** to implement both approaches, but you could - in a search for the most accurate classification results.

One-vs-the-world approach may prove to be more robust, but it comes with a cost of creating 50 classifiers where the 50-class classification approach needs to build only one model.

To evaluate the accuracy of the predictions, write an `classifierEvaluation.java` or `classifierEvaluation.py` program that takes as input the file generated by the `knnAuthorship` or the `RFAuthorship` program, as well as the ground truth file (remember you had to generate one!)⁴, and produces as the result the following output:

- For each author in the dataset, output the total number of *hits* (correctly predicted documents), *strikes* (false positives predicted) and *misses* (document written by the author, which were not attributed to the author).
- For each author in the dataset report the precision, the recall, and the f1-measure of the KNN procedure.
- Overall, report the total number of documents with correctly predicted authors, the total number of documents with incorrectly predicted authors, and the overall accuracy of the prediction.

⁴It may be helpful if both files are stored in exactly the same format.

- The full 50x50 confusion matrix in a CSV form (with the top row, and first column containing the author names/labels). This can be dumped directly into an output file, rather than printed to screen.

Please note, if you are implementing a one-vs-the-world classification, then instead of the `ClassifierEvaluation` program described above, your analysis should be performed as follows (name the program performing it `OneWorldEval`:

- Program first runs binary classification (using either your KNN or Random Forest implementation) for each of the fifty class labels.
- The prediction for each binary classification model is stored separately (create a working directory and dump 50 output files in it).
- For each class label, the program evaluates the precision and recall of this class based on the binary classification output for the model built to recognize this class, and report the number of hits, strikes, and misses for that specific class label (author).
- The full 50x50 confusion matrix here is not needed.

You can reuse and repurpose your code from Lab 3 for this task. Your evaluator may be a bit more complex, but it can be built conveniently as a branch of your Lab 3 evaluator code.

Your analytical goal is to determine which authors are easier to predict authorship for, and which authors are hard (and who they tend to get confused with). You can do the final analytical step by hand - simply looking at the results of your evaluation program, or, alternatively, you can add code to your evaluation program that reports this information out automatically.

Tuning Parameters for Information Retrieval/Vectorization Part of the Assignment

When experimenting with both analytical methods you may wind up running your analyses multiple time with a different set of parameters. For both methods, you can construct different datasets based on what pre-processing techniques were used to create them. Generally speaking, for the preprocessing, you have four key choices to consider:

- **No stopword removal, no stemming.**
- **Stop word removal, but no stemming.**
- **Stemming, but no stopword removal.**
- both **Stemming** and **stopword removal**.

Within this set of possibilities, different lists of stop words can be used to create more possible inputs to both the classification and clustering authorship attribution methods.

Finally, your similarity score (cosine or okapi) is another parameter that can be used to build different authorship attribution models.

Tuning Parameters for Classification Part of The Assignment

For the KNN-based authorship attribution, you have an additional hyperparameter to tune: K , the number of neighbors to use. You have some natural boundaries on K as part of your procedure: clearly, $K \leq 99$, and quite possibly needs to be made significantly smaller.

For the Random Forest-based authorship attribution, you have your three RF parameters (number of trees, number of attributes in a tree, number of data points used to build a tree), plus any C45 parameters you are willing to expose (generally speaking you can set a standard threshold, and commit to using either information gain or information gain ratio, so there is no need for extra parameters here beyond the three above).

Your random forest should be fairly large (remember, there are many words in the document collection, the Random forest should give each a chance to participate in helping recognize authorship), the number of attributes in a tree probably also needs to be larger than what we are used to be (10-20 might be a good target, otherwise

the number of decision trees will really need to explode), while the number of data points can be selected with an idea that you want to have representative writings from all authors in the data. (you could simply use 50% as your target, but can probably get away with a smaller percentage to speed things up).

Please note: training your models for this assignment may **take a long time**. Please make sure to give yourself at least 4-5 days to run all the code/collect all the information.

Reporting your findings

Prepare a report that describes the following:

- The experiment you ran for each of the two analytical methods for determining authorship: what datasets were explored, which values of hyperparameters were investigated.
- The results of your best runs for each of the two analytical methods. You should present in tabular (or easy to read text) form the accuracy measures for determining each author, as well as the overall accuracy. Any raw results shall be included in the appendix to your report.
- A reflection on each of the methods w.r.t. the method's overall ability to properly attribute authorship of the articles, as well as any specific information that stands out: which authors are easy to predict? which are hard? which tend to be confused with each other?
- A comparison of two methods to each other and a final determination of which method proved to be more accurate.

Submission

For this lab we will use CSL **handin** tool.

You shall submit the following.

- All the code you wrote to fulfill the requirements of this lab in a single zip or tarred gz file. Resubmit any code from Labs 3 or 4 that you reused.
- A README file with instructions on how to run your code to perform different tasks of the assignment.
- Output files for best respective runs generated by your evaluation programs for both analytical methods for authorship detection (specify filenames in README file). (Note, some of this information may also be included in your report appendix, but having an appendix does not absolve you from submitting the raw output files).
- Your report in PDF format.

Place **all code** files and directories *except* for the report file, the README file and the output files, into a single archive named **lab05.tar.gz** or **lab05.zip** and submit it. Submit the PDF version of the report, the README file and the output files separately, outside of the archive.

Use **handin** to submit as follows:

```
$ handin dekhtyar 466-lab06 <FILES>
```

Good Luck!