

Big Data Mining Techniques (M118)

Winter Semester 2020-2021

Deadline: Last Day Before Exam Period
Assignment for teams of 2 students

Goal

The purpose of the project is to familiarize you with the basic steps of the process followed for applying data mining techniques, namely: collection, preprocessing / cleaning, conversion, application of data mining techniques and evaluation. Implementation will be done in the Python programming language using the SciKit Learn and Keras tool. The thesis consists of two (2) tasks related to categorization, duplication detection. Two (2) separate competitions have been created for the requirements of the job on the Kaggle platform. You will need to sign up in the Kaggle platform using your academic email (STUDENT_ID@di.uoa.gr) and upload the output files with the predictions. The Kaggle platform provides you with 42 hours of GPU usage if you want to speed up your calculations with neural networks. Pay special attention to the report because the work is first graded by the quality of the documentation.

Requirement 1: Text classification

Description

The requirement is related to text classification of news articles. The data are organized in CSV files whose fields are separated by the '|' character. There are two files:

1. train_set.csv (111795 items): This file will be used to train your algorithms and contains the following fields:
 - a. Id: A unique number for the article
 - b. Title: The title of the article
 - c. Content: The content of the article
 - d. Label: The category to which the article belongs
2. test_set.csv (47912 items): This file will be used to predict in new (unseen) data. It contains the same fields as in the training file except from the Label field. You will be asked to predict this field using classification algorithms.

There are 4 categories of articles and they are presented in the table below.

Beat the Benchmark

You should select and experiment with any classification algorithm, preprocess steps in order to beat the performance of the best performing model of the previous question. you should justify the methodology that you choose to follow.

Evaluation Results

The report should include the following table with the evaluation of your techniques using the train-set and 5-Fold Cross Validation.

Statistic Measure	SVM (BoW)	Random Forest (BoW)	SVM (SVD)	Random Forest (SVD)	My Method
Accuracy					
Precision					
Recall					
F-Measure					

A description of the above results should be included in the report.

Output File

Your code should create the output file "testSet_categories.csv" which will contain the predictions for articles in the test set dataset (the ones where the Label field is not given). You should use your best model. The format of the testSet_categories.csv file, which will contain the categories of articles given in the Test set, is shown below:

Id	Predicted
1	Business
2	Technology
...	

For the file "testSet_categories.csv" the above formatting should be used *strictly* separating the two fields with the comma (",") character and should also have the first line with the two field names (Id and Predicted) followed by your model predictions in the following lines specifying the article Id from the test set and the predicted label.

You will need to upload your file to the Kaggle contest at the address <http://www.kaggle.com/c/bigdata2021classification>.

Hint:

1. Because the text files are large see:
https://scikit-learn.org/0.15/modules/scaling_strategies.html
2. Use the Kaggle computing resources if you want to try out complex neural networks. (This is outside the scope of the course).
3. For computing the evaluation metrics: Precision/Recall/F-Measure you will not use the number of examples in each class (Macro).

Requirement 2: Nearest Neighbor Search and Duplicate Detection

Question 2a: De-Duplication with Locality Sensitive Hashing

Description

In this question you will be given a train set file with small texts. Every text is a Quora question. You will also be given a test set in the same format. The purpose of the question 2a is to find how many of the documents in the test-set already exist in the train-set. As duplicate we define document pairs with similarity more than a threshold $\tau=0.8$. You have to search for duplicate documents using the appropriate LSH family in order to reduce the time required for the detection.

You have to do that considering the metrics:

1. Cosine Similarity: Random projection LSH family. Set the parameter K from 1 to 10.
2. Jaccard Similarity: Min-Hash LSH family. Set number of permutations to {16,32,64}

Evaluation Results

You need to evaluate the performance of the LSH algorithm and you should report:

1. The total LSH Index Creation Time (BuildTime)
2. The total time it took to answer all the test set questions. (QueryTime)
3. TotalTime: BuildTime + QueryTime
4. The number of duplicates in the test-set.

In your report should include a table as follows:

Type	BuildTime	QueryTime	TotalTime	#Duplicates	Parameters
Exact-Cosine	0	600	600	1000	-
Exact-Jaccard	0	300	300	1500	-
LSH-Cosine	30	200	230	800	K=2
LSH-Cosine	50	150	200	600	K=3
LSH-Jaccard	100	50	150	900	-

Things to Consider:

1. Try to use vectorized operations.
2. You can use available implementations of the LSH families
3. Use <http://ekzhu.com/datasketch/lsh.html>

Question 2b: Same Question Detection

Description

In this question you should find if questions with similar format ask the same thing. Consider the example:

1. What restaurants should I visit during my holidays trip in **Dublin**?
2. What restaurants should I visit during my holidays in **Athens**?

Clearly, the above two questions share (9) words but the question is not the same. Specifically you will be given pairs of Quora questions and you need to find out if the pair contains questions that ask the same. In other words, you want to create a model that can answer if two questions are ultimately identical. In this question you should experiment with heuristic features that could solve the above problem. **In your report you should describe in detail the similarity features you selected and how you trained your algorithm.**

For this question, there are two files:

1. train_set.csv (283014 pairs): This file will be used to train your algorithms and it contains the following fields:
 - a. Id: A unique number for the pair
 - b. Question1: The first question
 - c. Question2: The second question
 - d. IsDuplicate: Column describing whether the pair is a duplicate or not
2. test_set.csv (47912 items): This file will be used to make predictions for new data. this file contains all fields of the training file except from the IsDuplicate field. You will be asked to predict this field using classification algorithms.

Evaluation Results

You should evaluate the technique using 5-fold Cross Validation and you should report the following metrics: Accuracy, Precision, Recall and F1. **Your report should include a table as the following:**

Method	Precision	Recall	F-Measure	Accuracy
Method-1	0,9	0,9	0,9	0,9

Method-2	0,8	0,8	0,8	0,8
----------	-----	-----	-----	-----

A description of the above results should be included in the report.

Output Files

You should use your best model and create the file `duplicate_predictions.csv` containing the test set predictions. The file format is CSV and is shown below:

Id	Predicted
1	1
2	0
...	

For the file "duplicate_predictions.csv" the above formatting should be used *strictly* separating the two fields with the comma (",") character and should also have the first line with the two field names (Id and Predicted) followed by your model predictions in the following lines specifying the article Id from the test set and the predicted label to indicate whether the documents in the pair with the specified Id are similar or not.

You will need to upload your file to the Kaggle contest at the address <http://www.kaggle.com/c/bigdata2021duplicatedetection>.

Hint:

1. In this question it is important to experiment with different pre-processing techniques for the questions and with different heuristic features.
2. The question pairs were annotated by Quora, and because this task is to some extent subjective, there may be errors in both the train-set and the test-set.
3. For Precision / Recall / F-Measure you will not use the number of examples in each class (Macro-Precision, Macro-Recall, Macro-F-Measure).

Regarding the deliverables

The folder you deliver should have the name:

Ass1_name1_AM1_name2_AM2.

The folder should contain:

1. A text with detailed analysis on the experiments you did and the methods you tried in PDF format. Your report should also contain all the tables and plots requested and should not exceed 30 pages. In the report you should include a description of your experiments and everything you can think of to show what experiments you did, why

the specific results of the methods you selected, how these methods work, and commentary on your results. **All tasks will be evaluated on the basis of the detailed documentation and the extent to which the tasks are being implemented.**

2. The requested output files.
3. The source code files.

Datasets

Available at: <http://195.134.67.98/documents/BigData/datasets2020.tar.gz>

Username: bigdata

Password: d@t@s3t