



SpamDetector

Project for Data Mining and Machine Learning

Francesco Campilongo

Contents

1	Introduction	2
2	Datasets	3
2.1	Data Preprocessing	3
3	Design	4
3.1	Implementation	4
3.1.1	Libraries	4
3.1.2	Functions	4
3.1.3	Main code	5
4	Results	6

Chapter 1

Introduction

The purpose of this documentation is to prove the machine learning algorithms utility for the detection of spam, useless messages, into dataset of emails.

This is a basic text classification problem, and so the algorithm used are very known in the lecture.

The project is carried out on python.

Chapter 2

Datasets

The datasets used for this project are two:

- E-mail dataset
- SMS dataset

Those two datasets were found on Kaggle.com both of them in a .csv format.

2.1 Data Preprocessing

The datasets found were already pretty good for the type of result this project wants to achieve. The E-mail dataset is characterized from three columns: "label", "text" and "label_num", the second column contains the actual message, the first and third columns are the same, but in the first one there are the actual words "ham", for the messages which are not classified as spam, and "spam", instead in the third column there are "0" for the ham messages and "1" for the spam messages;

The first column has been discharged since not useful for the execution of the classification algorithm.

The SMS file had need a little more work to obtain a dataset usable, because it has five columns named "v1", "v2", "", "" and "", the first contains the words "ham", for the messages which are not classified as spam and "spam"; the second column contains the actual messages and the last three columns are basically blank, so in order to obtain a good dataset to work with the last three columns have been discharged the first two are taken in consideration, but in the first one the "ham" word is changed with a "0" and the "spam" word with a "1", in order to have the same kind of datasets between the E-mail and the SMSs.

Chapter 3

Design

The application design is pretty basic, it is written in python (version 3) does not involve any classes, just the creation of some functions in order to keep the code clean and more understandable.

3.1 Implementation

3.1.1 Libraries

The libraries involved into this application are two:

- Pandas
- Sklearn

Pandas gives the possibility to use a dataset in .csv (and many others) from the disk, implementing it on a dataframe into the main memory to use all the machine learning algorithm that are imported from Sklearn.

3.1.2 Functions

The functions are meant to implement all the algorithm used, which are the following:

- Support Vector Machine
- Multinomial N  ive Bayes
- K Nearest Neighbor
- Random Forest
- AdaBoost

These functions are written to use them with different datasets, the SMS dataset and the E-Mail dataset as already mention earlier.

There is also another function that is called to show the result of all the algorithms.

3.1.3 Main code

This section is used to show some insight on the main code.

Here will be shown a bit of the dataset cleaning and processing that has been made, and described into the Data Preprocessing chapter.

```
1 df2 = pd.read_csv("sms_spam.csv", encoding = "ISO-8859-1")
2 dfsp = df2[['v1', 'v2']]
3 dfsp.loc[dfsp["v1"] == 'ham', "Category"] = 0
4 dfsp.loc[dfsp["v1"] == 'spam', "Category"] = 1
5 dfsp = dfsp.rename(columns={'v2': 'Content'})
6 dfs = dfsp[['Content', 'Category']]
```

The following code will give a look at one of the function, the one that implement the Support Vector Machine. Since these functions are quite similar there is no need to show them all here.

```
1 def SVM(x_train, y_trainFeat, y_testFeat):
2     classifier = LinearSVC()
3     classifier.fit(y_trainFeat, x_train)
4     predRes = classifier.predict(y_testFeat)
5     return predRes
```

The following code lines will show the functions calls, also in this case in order to avoid repetitions only few functions call will be shown.

```
1 # KNN
2 predResMailKNN = KNN(x_train, y_trainFeat, y_testFeat)
3 # Metrics and results
4 print("\tK Nearest Neighbors RESULTS")
5 print("Neighbors Number: 1")
6 show_res(actual_x, predResMailKNN)
7
8 # RF
9 predResMailRF = RF(x_train, y_trainFeat, y_testFeat)
10 # Metrics and results
11 print("\tRandom Forest RESULTS")
12 show_res(actual_x, predResMailRF)
```

Chapter 4

Results

The results are as expected very goodscrivere roba riguardante il motivo per il quale i risultati sono expected magari elencando alcune peculiarità degli algoritmi

The metrics that shows the results goodness are: Accuracy, F Score and Confusion Matrix; Chosen in order to give a deep look at the algorithms behavior.

Here the results for what concern the E-mail dataset:

```
1 /-----SpamDetector for e-Mail-----/
2 Support Vector Machine results
3 Accuracy Score: 98.6473
4 F Score: 98.3941
5 Confusion Matrix:
6 [[716 13]
7 [ 1 305]]
8
9 Multinomial Naive Bayes results
10 Accuracy Score: 91.3043
11 F Score: 88.5304
12 Confusion Matrix:
13 [[727 2]
14 [ 88 218]]
15
16 K Nearest Neighbors results
17 Neighbors Number: 1
18 Accuracy Score: 95.8454
19 F Score: 94.9084
20 Confusion Matrix:
21 [[718 11]
22 [ 32 274]]
23
24 Random Forest results
25 Accuracy Score: 95.4589
26 F Score: 94.5432
27 Confusion Matrix:
28 [[706 23]
29 [ 24 282]]
30
31 Adaboost results
32 Estimators Number: 100
33 Accuracy Score: 96.9082
34 F Score: 96.3159
35 Confusion Matrix:
36 [[709 20]
37 [ 12 294]]
```

Here the results for what concern the SMS dataset:

```
1 /-----SpamDetector for SMS-----/
2 Support Vector Machine results
3 Accuracy Score: 98.5650
4 F Score: 96.8657
5 Confusion Matrix:
6 [[960  0]
7  [ 16 139]]
8
9 Multinomial Naive Bayes results
10 Accuracy Score: 97.1300
11 F Score: 93.4249
12 Confusion Matrix:
13 [[960  0]
14  [ 32 123]]
15
16 K Nearest Neighbors results
17 Neighbors Number: 1
18 Accuracy Score: 94.7982
19 F Score: 87.0259
20 Confusion Matrix:
21 [[960  0]
22  [ 58  97]]
23
24 Random Forest results
25 Accuracy Score: 97.4888
26 F Score: 94.3510
27 Confusion Matrix:
28 [[959  1]
29  [ 27 128]]
30
31 Adaboost results
32 Estimators Number: 100
33 Accuracy Score: 97.7578
34 F Score: 95.0885
35 Confusion Matrix:
36 [[956  4]
37  [ 21 134]]
```