

# Report Homework1

MACHINE LEARNING

PAPA EDOARDO    MATRICOLA: 1962169



## Sommario

|  |    |
|--|----|
| 1. Project Overview .....                        | 3  |
| 1.1 Dataset .....                                | 3  |
| 2. Vectorizer .....                              | 4  |
| 3. Classification .....                          | 6  |
| 3.1 Accuracy .....                               | 6  |
| 3.2 Pecision.....                                | 6  |
| 3.3 Recall .....                                 | 6  |
| 3.4 F1-score .....                               | 7  |
| 4. Multiclass Classification inside project..... | 8  |
| 4.1 Split Data .....                             | 8  |
| 4.2 Model .....                                  | 8  |
| 4.3 Evaluations .....                            | 9  |
| 4.3.1 Bernoulli.....                             | 9  |
| 4.3.2 Logistic Regression.....                   | 11 |
| 4.3.3 Multinomial Naive Bayes .....              | 13 |
| 4.3.4 SVM Linear Kernel.....                     | 15 |
| 4.4 Conclusion of evaluation .....               | 16 |
| 5. Blind Test .....                              | 17 |

# 1. Project Overview

The goal of this project is solve one classification problem.

In particular the aim is a semantic prediction, more precisely, build a classifier that can predict which type of function we are facing, starting to specific list of assembly functions.

## 1.1 Dataset

If we choose a “*dataset.json*”, it contains 14397 functions, each function belongs to one of aformetioned classes (*Encryption, String Manipulation, Math, Sorting*):

- 2724 are Encryption Functions;
- 3104 are String Manipulation;
- 4064 are Sorting Functions;
- 4504 are Math Functions;

The structure of the dataset is a JSONL, so for each lines we have a JSON composed in the following way:

```
{
  "id": "828",
  "semantic": "string",
  "lista_asm": ["jmp qword ptr [rip + 0x220882]", "jmp qword ptr [rip + 0x220832]", ...],
  "cfg": {
    "directed": true,
    "graph": [
      ],
    "nodes": [
      {
        "id": 4276480,
        "asm": "554889e54883ec2048897df048837df0000f841d000000",
        "label": "0x414100:\\tpush\\trbp\\n0x414101:\\tmov\\trbp, rsp...",
        ...
      },
      ...
    ],
    "adjacency": [
      [
        {
          "id": 4276532
        },
        {
          "id": 4276503
        }
      ],
      ...
    ],
    "multigraph": false
  }
}
```

Figura 1 - Dataset Structure

The JSON file, that is composed of a multiple JSON lines, i decided to read it and go over using a pandas library, that allow read each lines through a function called “*read\_json*”. This function returns a DataFrame type that has a multi labels. In particular: *id*, *semantic*, *lista\_asm*, *cfg*.

So, each lines have:

- **ID** parameter that provide to give it a unique value inside the JSON;
- **SEMANTIC** specified the which classes belong;
- **LISTA\_ASM** specified the list of all assembly instructions;
- **CFG** represent the graph.

## 2. Vectorizer

The first question that someone can do are: “*why we should use a Vectorizer?*”

The answer of this question is very simply. In machine learning the algorithms operate on a numerical space of functionality. So, to improve machine learning on our data, we need to transform the data into a vector representation. Now we can apply a numerical machine learning.

Inside this project i decided to use 3 different type of vectorizer, where each one have different property:

- **CountVectorizer**: creates a matrix with each unique word is represented by a column of the matrix and each text sample from the document is a row in the matrix. The value of each cell is the count of the word in that particular text sample.
- **HashingVectorizer**: HashingVectorizer is very similar to CountVectorizer. It consist on hashing the words, then tokenize and encode document as needed. With this vectorizer, each token is mapped to a column position in a matrix.
- **TfidfVectorizer**: also this vectorizer is similar to CountVectorizer. In this case the vectorizer doesn't count the word but it calculate the word frequencies. This vectorizer have another important feature; infact it uses TF-IDF that provide to highlight words that are more interesting then the others. The TfidfVectorizer tokenize documents, learn the vocabulary and inverse document frequency weightings, and allow you to encode new documents.

Before choose the Vectorizer i evaluated all the ones mentioned above. After that i decide to choose the TfidfVectorizer because it come out the better performant solution for this problem. If instead i used the CountVectorizer i could have some words that will appear many times and their aren't very significant in the coded vectors. I also decided to don't use HashingVectorized because this vectorizer require large vectors for encoding data. Moreover the hash is one-way function so there isn't way to convert the decoding data to word.

Is important to consider in this vectorizer the parameter **ngram\_range** that define the lower and upper limit range of values for different n-gram. This is useful for the vocabulary. In my case i consider a **ngram\_range**

= (1,4) because in this way i have different combinations of the word in vocabulary. It improve the result of evaluation.

After choose a vectorizer i execute the *fit\_transform* function. This function allow to learn the vocabulary and do the idf. This function is equivalent to call *fit* and after that *transform*.

One thing that is very important to think it's when evaluating a model we must use two different dataset, one for train and another one for test. This is very important because if use the same dataset this won't evaluate well and the model not works in realistic conditions. So we use two dataset: *train dataset* and *test dataset*.

## 3. Classification

Inside this project we'll facing to a classification report that it contains different type of informations and is about key metrics in a classification problem.

In particular, the information that we'll see are mentioned below.

### 3.1 Accuracy

This parameter allow immediately if a model is being trained correctly or not. However it doesn't give more detail to the application problem.

I cannot use only this parameter for the metric, because it not works good when we have a imbalance class dataset. In this case, we also consider the f1-score.

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Number of total predictions}}$$

In confusion matrix we have

$$\text{Accuracy} = \frac{\text{True positive} + \text{True negative}}{\text{Total predictions}}$$

### 3.2 Precision

The precision is the accuracy of the system predict a positive classes. In other words, helps when the costs of false positive are high.

In confusion matrix we have

$$\text{Precision} = \frac{\text{True positive}}{\text{True positive} + \text{False positive}}$$

### 3.3 Recall

The recall it's also called "sensitive" or "true positive rate" and indicates the ratio of positive instances correctly identified. In other words, helps when the cost of false negative is high.

$$\text{Recall} = \frac{\text{True positive}}{\text{True positive} + \text{False negative}}$$

## 3.4 F1-score

F1-score is an overall measure of a model's accuracy that combines precision and recall. For have a good f1-score, we could have a low false positives and low false negatives. An f1-score is considered perfect when it's 1, while the model is a failure when it's 0.

$$F1_{\text{score}} = 2 * \frac{\textit{Precision} * \textit{Recall}}{\textit{Precision} + \textit{Recall}}$$



## 4. Multiclass Classification inside project

The first think that we need to do are: split the dataset into two different dataset called by “*X\_train*” and “*X\_test*” with relative labels “*y\_train*” and “*y\_test*”. After that we need to create a model for prediction and at the end we need to predict the value and plot the result with *classification\_report* and *confusion\_matrix*.

### 4.1 Split Data

We want to considera the output as a semantic, so we take a column of the dataset labelled with “*semantic*”. The *train\_test\_split* is a function defined in Sklearn model section for split in a esy way the dataset into two different ones. The first is a training set and the second is a testing set.

We pass to this function 4 parmameters. The first parameter is a sparse matrix of the dataset passed before inside vectorizer (for this project i choose TfidfVectorizer). The second parameter is a list of all label of the dataset. The third parameter is a test dataset size, that decides the size of the data tha has to be split. For example if we use 0.2 as the value, automatically we’ll have a train set with 80% of the original dataset and the remain 20% is test set of the original dataset. The last parameter is a random state that will act as a seed for random number generator during the split.

### 4.2 Model

In this project different models are tested, such as:

- Logistic Regression;
- Multinomial Naive Bayes;
- Bernoulli
- SVM Linear Kernel

Each models will be evaluated using the classification report and confusion matrix.

## 4.3 Evaluations

### 4.3.1 Bernoulli

Bernoulli is a Naive Bayes classifier for multivariate Bernoulli models. This classifier is suitable for discrete data. I don't consider Random Forest because we have high data inside dataset that almost certainly generates overfitting, compromising the evaluation of the model.

#### Classification Report and Confusion Matrix

*Bernoulli with alpha value = 1.0*

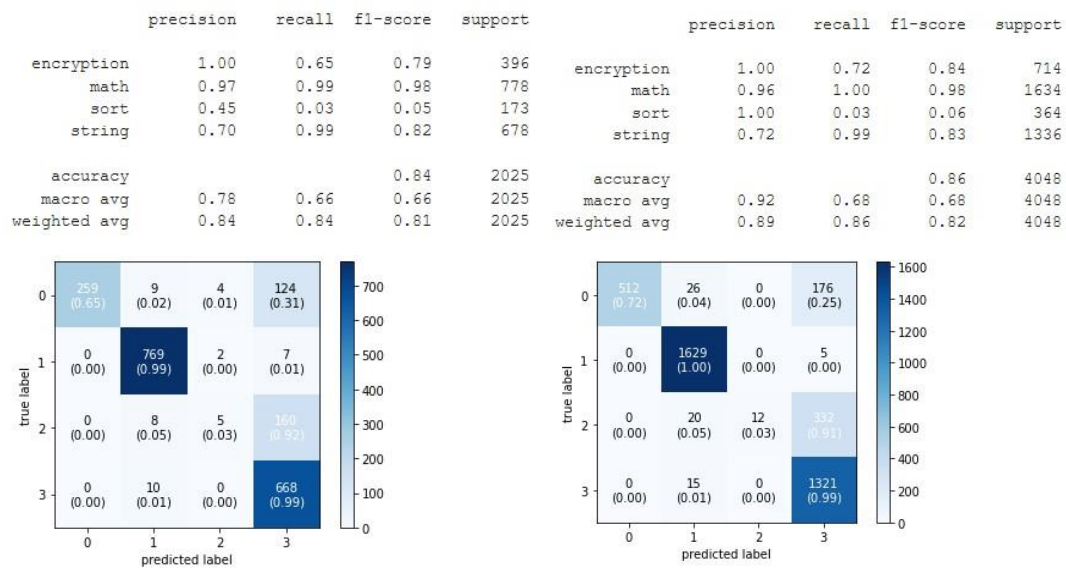


Figura 2 - Bernoulli test set (alpha = 1.0)

Figura 3 - Bernoulli train set (alpha = 1.0)

*Bernoulli with alpha value = 0.01*

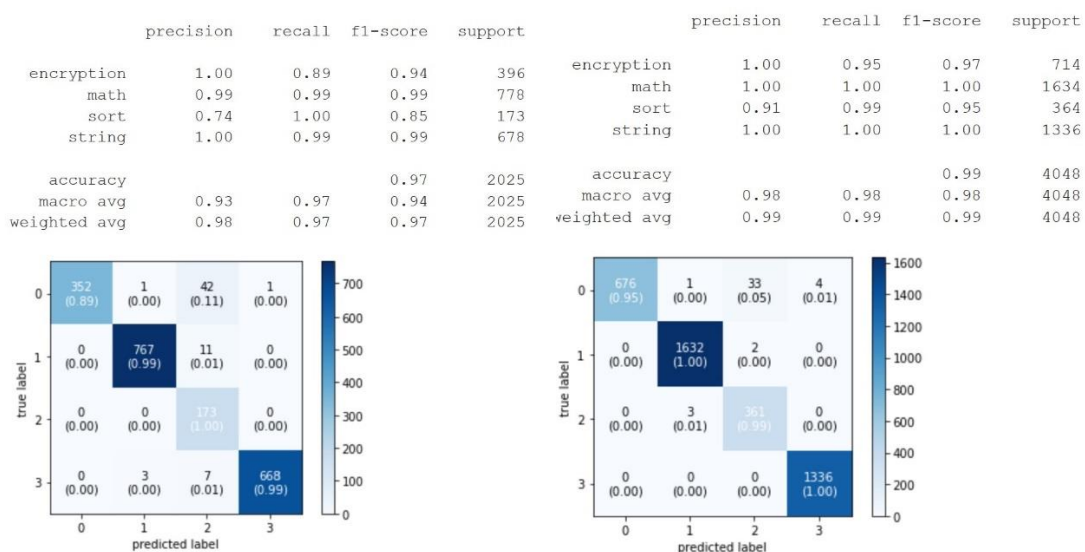


Figura 4 - Bernoulli test set (alpha = 0.01)

Figura 5 - Bernoulli train set (alpha = 0.01)

## Conclusion

The performance of this model isn't the best that i observed. I started trying the Bernoulli model with alpha *value* =  $1.0$ . There are a big gap between f1\_score of encryption, math, sort and string. I made a second step where i tried to change an alpha *value* =  $0.01$ . In this case i diminished the gap between f1\_score of encryption, math, sort and string. In conclusion the Bernoulli model generate an overfitting, so i completely reject this model for solve this problem.

### 4.3.2 Logistic Regression

The idea of Logistic Regression is find a relationship between specific feature and probability of particular comeout.

For this evaluation the Logistic Regression is created specifying some parameters. We need to find multiclass classifier, so we consider *multi\_class = multinomial* but for do this operation we need to choose different solver from the default one. In my case i select “*solver = saga*” and “*penalty = l1*”. The penalty is a regularization that allow to avoid the overfitting by discriminate the high-valued regression coefficients. In other words, it reduce the parameters for simplify the model.

#### Classification Report and Confusion Matrix

*Logistic Regression with penalty = l1*

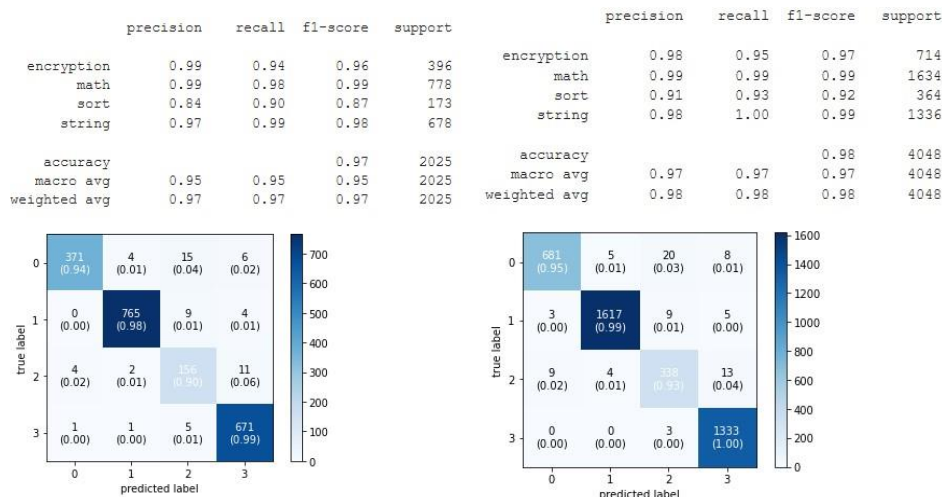


Figure 6 - Logistic Regression test set penalty l1

Figure 7 - Logistic Regression train set penalty l1

*Logistic Regression with penalty = l2*

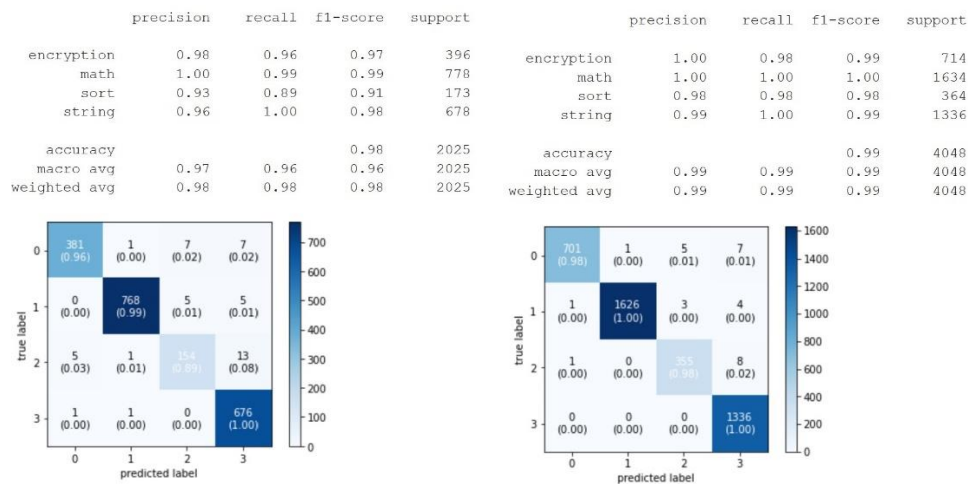


Figure 8 - Logistic Regression test set penalty l2

Figure 9 - Logistic Regression train set penalty l2

Consider the figure 6 we have 4x4 matrix, on the columns the prediction and on the rows the true values. The column 0 corresponding to “*encryption*”, 1 corresponding to “*math*”, 2 corresponding to “*sort*”, 3 corresponding to “*string*”. In this particular case, if we consider the *encryption* row, we can see that 371 of these are correctly and incorrectly predicted 4 of the *encryption* to be *math*, 15 of *encryption* to be *sort* and 6 of *encryption* to be *string*. Now, if we consider the *encryption* column, we can see there are 376 element that are predicted as *encryption* by the model but only 371 element are really *encryption*, 4 elements are *sort* and the last one is *string*. The same is for other three classes.

## Conclusion

This is the best result that i’ve obtained. I started trying the Linear Regression with *penalty = l1* (figure 6). We can see from this evaluations that Logistic Regression has high values with regard to *accuracy* and also the other metrics have high values, in particular the *f1\_score*. The most important thing is check if the model causes overfitting. Fortunately, in this case the model doesn’t cause a overfitting problem. I made a second step where i changed the value of penalty to l2. We can see the *f1\_score* and *accuracy* was improve, but we introduced a overfitting problem. So, the best result is using this model with *penalty=l1*.

### 4.3.3 Multinomial Naive Bayes

#### Classification Report and Confusion Matrix

*Multinomial NB with  $\alpha = 1.0$*

|              | precision | recall | f1-score | support |              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|--------------|-----------|--------|----------|---------|
| encryption   | 1.00      | 0.80   | 0.89     | 396     | encryption   | 1.00      | 0.87   | 0.93     | 714     |
| math         | 0.99      | 0.97   | 0.98     | 778     | math         | 1.00      | 0.98   | 0.99     | 1634    |
| sort         | 1.00      | 0.02   | 0.03     | 173     | sort         | 1.00      | 0.03   | 0.06     | 364     |
| string       | 0.72      | 1.00   | 0.83     | 678     | string       | 0.74      | 1.00   | 0.85     | 1336    |
| accuracy     |           |        | 0.87     | 2025    | accuracy     |           |        | 0.88     | 4048    |
| macro avg    | 0.93      | 0.70   | 0.69     | 2025    | macro avg    | 0.93      | 0.72   | 0.71     | 4048    |
| weighted avg | 0.90      | 0.87   | 0.83     | 2025    | weighted avg | 0.91      | 0.88   | 0.85     | 4048    |

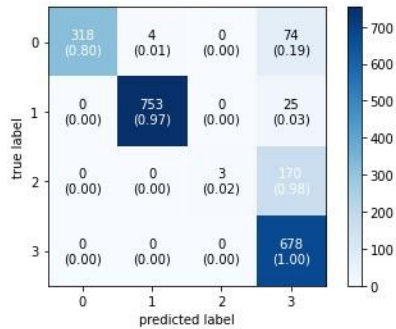


Figura 10 - Multinomial Naive Bayes test set  $\alpha=1.0$

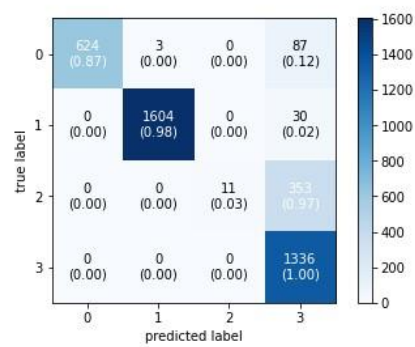


Figura 11 - Multinomial Naive Bayes train set  $\alpha=1.0$

*Multinomial NB with  $\alpha = 0.01$*

|              | precision | recall | f1-score | support |              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|--------------|-----------|--------|----------|---------|
| encryption   | 1.00      | 0.97   | 0.98     | 396     | encryption   | 1.00      | 0.99   | 1.00     | 714     |
| math         | 1.00      | 0.98   | 0.99     | 778     | math         | 1.00      | 1.00   | 1.00     | 1634    |
| sort         | 0.89      | 0.98   | 0.93     | 173     | sort         | 0.99      | 0.99   | 0.99     | 364     |
| string       | 0.99      | 1.00   | 0.99     | 678     | string       | 0.99      | 1.00   | 1.00     | 1336    |
| accuracy     |           |        | 0.99     | 2025    | accuracy     |           |        | 1.00     | 4048    |
| macro avg    | 0.97      | 0.98   | 0.98     | 2025    | macro avg    | 1.00      | 1.00   | 1.00     | 4048    |
| weighted avg | 0.99      | 0.99   | 0.99     | 2025    | weighted avg | 1.00      | 1.00   | 1.00     | 4048    |

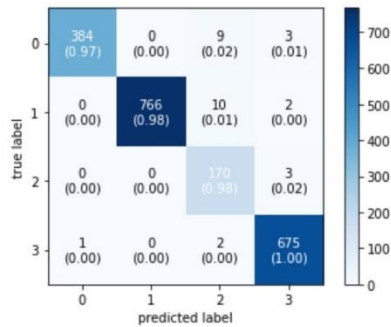


Figura 12 - Multinomial Naive Bayes test set  $\alpha=0.01$

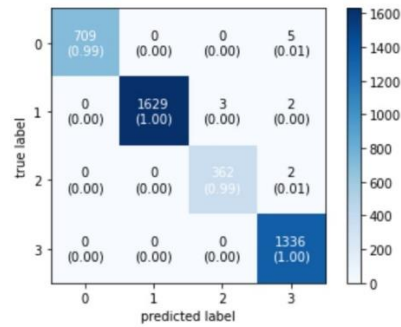


Figura 13 - Multinomial Naive Bayes train set  $\alpha=0.01$

## Conclusion

The performance of this model isn't the best that i observed. I started trying the Multinomial NB model with  $\alpha \text{ value} = 1.0$ . There are a big gap between f1\_score of sort and others values. I made a second step where i tried to change an  $\alpha \text{ value} = 0.01$ . In this case i diminished the gap between f1\_score of encryption, math, sort and string. In conclusion the Multinomial NB model generate an overfitting, so i completely reject this model for solve this problem.

#### 4.3.4 SVM Linear Kernel

Linear Kernel SVM is mainly used when the data is linear separable. With this model we can use a C parameter that it's useful for the SVM optimization for avoid misclassifying. For large values of C, the optimization will choose a smaller-margin hyperplane if that hyperplane. Instead a small value of C will cause the optimizer to look for a larger-margin separating hyperplane.

#### Classification Report and Confusion Matrix

*SVM Linear Kernel with C = 1.0*

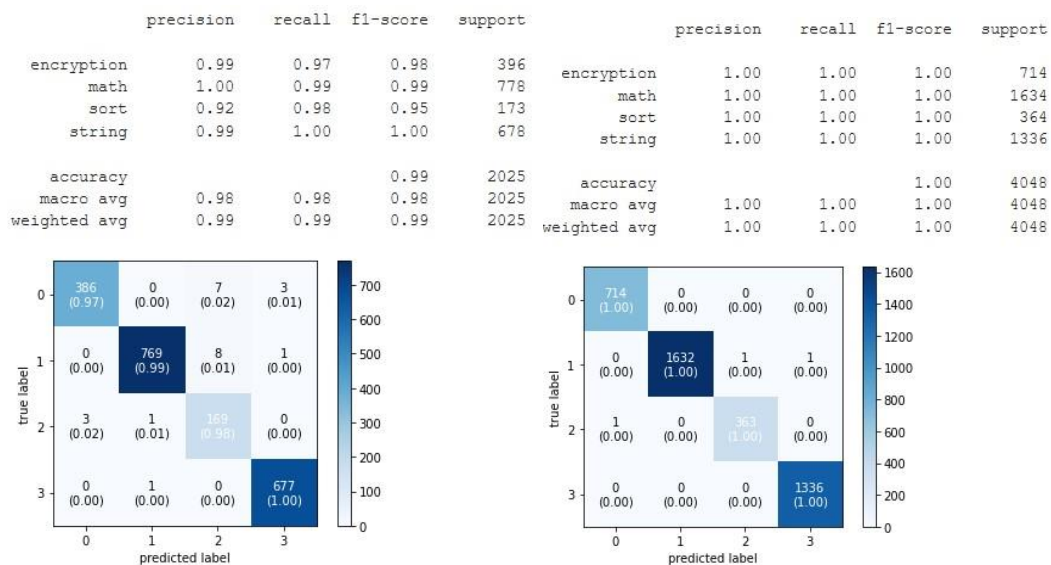


Figura 14 - SVM Linear Kernel test set C=1

Figura 15 - SVM Linear Kernel train set C=1

*SVM Linear Kernel with C = 10.0*

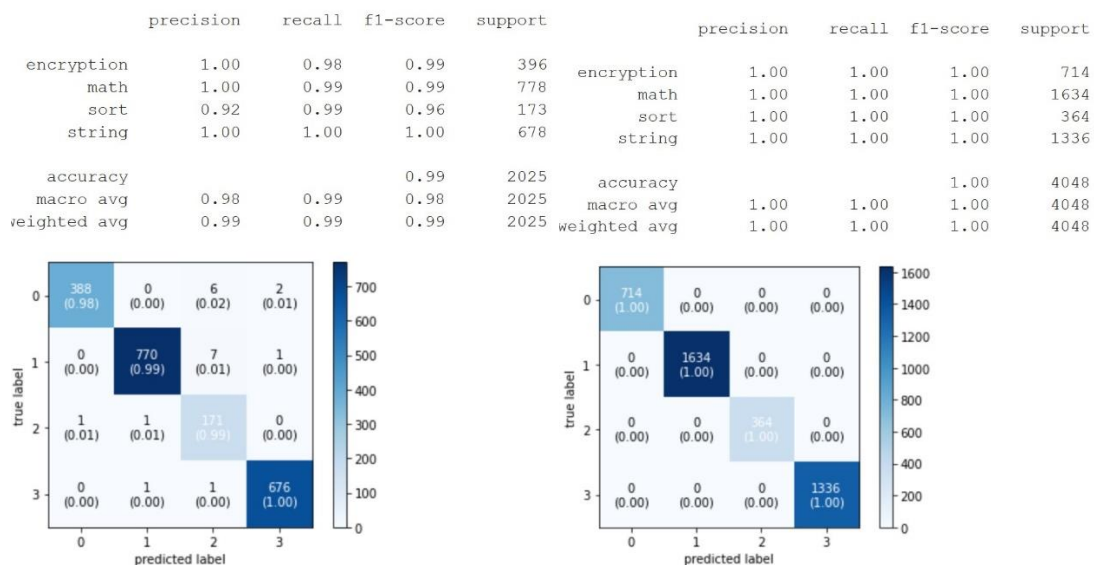


Figura 16- SVM Linear Kernel test set C=10

Figura 17 - SVM Linear Kernel train set C=10



## Conclusion

I started trying this model with parameter of  $C=1$  and how we can see, this model generates overfitting, so we cannot considerate it. Even if it has a good values of metrics, in particular *f1\_score*. I made a second step where i change the value of  $C$  to 10 and also in this case it generates overfitting.

### 4.4 Conclusion of evaluation

In conclusion, the best model for solve this problem is Logistic Regression because we want to avoid as much as possible the overfitting problem, and this model is the best for prediction of *semantic*. Another advantage point of this model it's that the *f1\_score* metric is very good.

## 5. Blind Test

In this project we have a file available called “*blindtest*” that contains a set of data without semantic information. So the aim is predict all the semantics for the input data and make another file with the result. In this phase i read the file and i pass it inside the vectorizer (TfidfVectorizer). After that i do the *fit* and *transformation*. I make prediction using as model Logistic Regression. After do this operations, i write the predictions inside a .txt file. For each row we can find a prediction that corresponding line of *blindtest* file.