

Tema 4

Clasificator Naive Bayes

Dat fiind setul de date spam_ham_dataset.csv, avem urmatoarele cerinte:

- De citit setul de date;
- De impartit in date de antrenament (3/4 din date) si date de test (restul de 1/4);
- De aplicat Naive Bayes din sklearn si de masurat acuratetea;
- De implementat propriul clasificador Naive Bayes si de masurat acuratetea; de comparat cu rezultatele din librarie;

Continutul fisierului **tema4.py** unde am aplicat Naive Bayes preimplementat si am comparat cu Naive Bayes-ul implementat de mine (care se afla in fisierul my_naive_bayes.py) este:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.naive_bayes import GaussianNB
from my_naive_bayes import MyNaiveBayes
from sklearn.metrics import accuracy_score

Data = pd.read_csv("spam_ham_dataset.csv")

# convert text to vector
count_vectorizer = CountVectorizer()
X = count_vectorizer.fit_transform(Data["text"]).toarray()

X_train, X_test, y_train, y_test = train_test_split(X, Data["label_num"].tolist(), test_size=0.25,
random_state=109)

gnb = GaussianNB()
gnb.fit(X_train, y_train)
y_pred = gnb.predict(X_test)
print("Accuracy of implemented Naive Bayes:", accuracy_score(y_test, y_pred))

# compare it with my implementation
mnb = MyNaiveBayes()
y_pred = mnb.implemented_classifier(X_train, y_train, X_test)
print("Accuracy of my implementation for Naive Bayes:", accuracy_score(y_test, y_pred))
```

Continutul fisierului **my_naive_bayes.py** care contine implementarea mea pentru Naive Bayes este:

```
import math

class MyNaiveBayes:
    def mean(self, nums):
        return sum(nums) / float(len(nums))

    def std_dev(self, nums):
        m = self.mean(nums)
        sigma = math.sqrt(sum([pow(x - m, 2) for x in nums]) / float(len(nums) - 1))
        return sigma

    def mean_and_std_dev(self, data):
        info = [(self.mean(attr), self.std_dev(attr)) for attr in zip(*data) if self.mean(attr) != 0 and
self.std_dev(attr)]
        return info

    def separate_by_class(self, X_train, y_train):
        dict = {}
        for i in range(len(X_train)):
            if y_train[i] not in dict:
                dict[y_train[i]] = []
            dict[y_train[i]].append(X_train[i])
        return dict

    def mean_and_std_dev_for_class(self, X_train, y_train):
        info = {}
        dict = self.separate_by_class(X_train, y_train)
        for class_val, instances in dict.items():
            info[class_val] = self.mean_and_std_dev(instances)
        return info

    def calculate_gaussian_probability(self, x, mean, stdev):
        expo = math.exp(-(math.pow(x - mean, 2) / (2 * math.pow(stdev, 2))))
        return (1 / (math.sqrt(2 * math.pi) * stdev)) * expo

    def calculate_class_probabilities(self, info, test):
        probabilities = {}
        for class_val, class_summaries in info.items():
            probabilities[class_val] = 1
            for i in range(len(class_summaries)):
                mean, std_dev = class_summaries[i]
```

```

        x = test[i]
        probabilities[class_val] *= self.calculate_gaussian_probability(x, mean, std_dev)
    return probabilities

```

```

def predict(self, info, test):
    probabilities = self.calculate_class_probabilities(info, test)
    best_label, best_prob = None, -1
    for class_val, probability in probabilities.items():
        if best_label is None or probability > best_prob:
            best_prob = probability
            best_label = class_val
    return best_label

```

```

def get_predictions(self, info, test):
    predictions = []
    for i in range(len(test)):
        result = self.predict(info, test[i])
        predictions.append(result)
    return predictions

```

```

def implemented_classifier(self, X_train, y_train, X_test):
    info = self.mean_and_std_dev_for_class(X_train, y_train)
    return self.get_predictions(info, X_test)

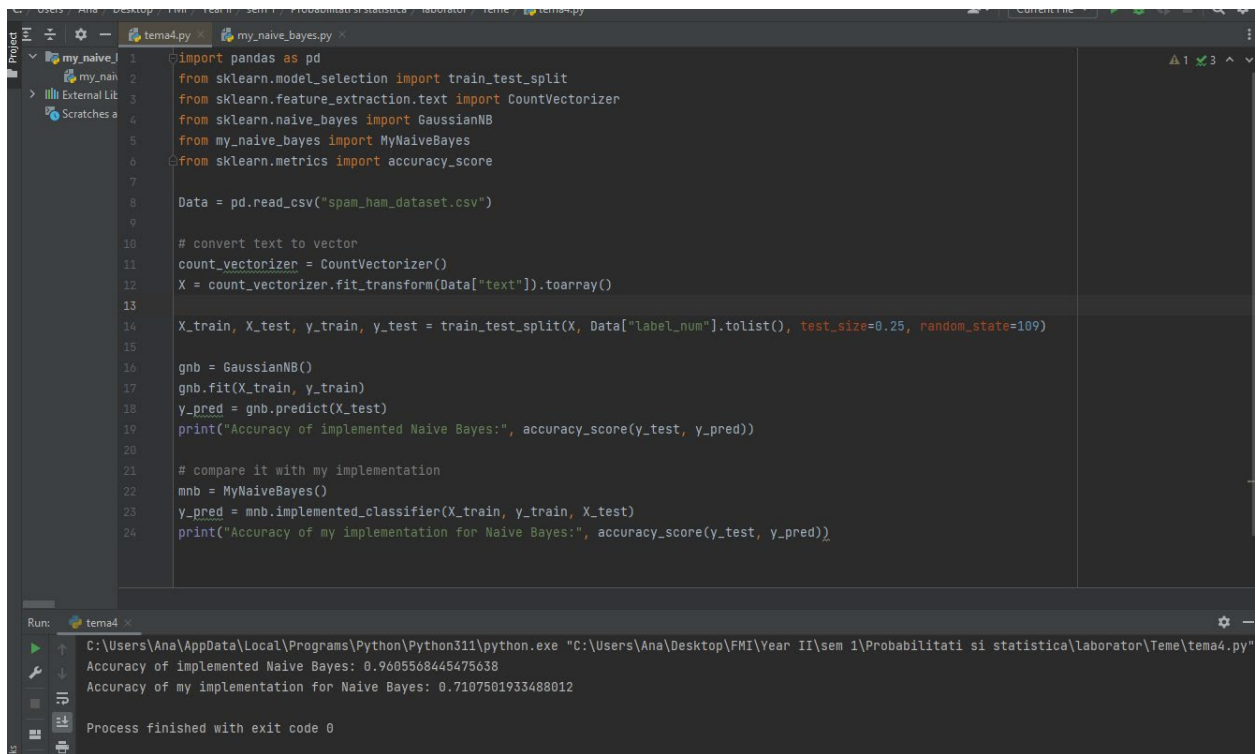
```

```

1  import math
2
3
4  class MyNaiveBayes:
5      def mean(self, nums):
6          return sum(nums) / float(len(nums))
7
8      def std_dev(self, nums):
9          m = self.mean(nums)
10         sigma = math.sqrt(sum([pow(x - m, 2) for x in nums]) / float(len(nums) - 1))
11         return sigma
12
13     def mean_and_std_dev(self, data):
14         info = [(self.mean(attr), self.std_dev(attr)) for attr in zip(*data) if self.mean(attr) != 0 and self.std_dev(attr)]
15         return info
16
17     def separate_by_class(self, X_train, y_train):
18         dict = {}
19         for i in range(len(X_train)):
20             if y_train[i] not in dict:
21                 dict[y_train[i]] = []
22             dict[y_train[i]].append(X_train[i])
23         return dict
24
25     def mean_and_std_dev_for_class(self, X_train, y_train):
26         info = {}
27         dict = self.separate_by_class(X_train, y_train)
28         for class_val, instances in dict.items():
29             info[class_val] = self.mean_and_std_dev(instances)
30         return info
31
32     def calculate_gaussian_probability(self, x, mean, stdev):
33         expo = math.exp(-(math.pow(x - mean, 2) / (2 * math.pow(stdev, 2))))
34         return (1 / (math.sqrt(2 * math.pi) * stdev)) * expo
35
36     def calculate_class_probabilities(self, info, test):
37         probabilities = {}
38         for class_val, class_summaries in info.items():
39             probabilities[class_val] = 1
40
41     def separate_by_class(self, X_train, y_train):

```

Cand rulez fisierul **tema4.py** se vor afisa urmatoarele rezultate:

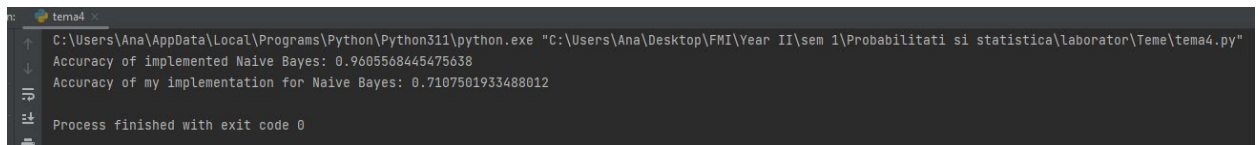


The screenshot shows a code editor with the file `tema4.py` open. The code imports `pandas`, `train_test_split`, `CountVectorizer`, `GaussianNB`, `MyNaiveBayes`, and `accuracy_score`. It reads the `spam_ham_dataset.csv` file, converts text to vectors, splits the data into training and testing sets, and compares the accuracy of a built-in `GaussianNB` classifier with a custom `MyNaiveBayes` implementation. The output shows that the custom implementation has a slightly higher accuracy.

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.naive_bayes import GaussianNB
5 from my_naive_bayes import MyNaiveBayes
6 from sklearn.metrics import accuracy_score
7
8 Data = pd.read_csv("spam_ham_dataset.csv")
9
10 # convert text to vector
11 count_vectorizer = CountVectorizer()
12 X = count_vectorizer.fit_transform(Data["text"]).toarray()
13
14 X_train, X_test, y_train, y_test = train_test_split(X, Data["label_num"].tolist(), test_size=0.25, random_state=109)
15
16 gnb = GaussianNB()
17 gnb.fit(X_train, y_train)
18 y_pred = gnb.predict(X_test)
19 print("Accuracy of implemented Naive Bayes:", accuracy_score(y_test, y_pred))
20
21 # compare it with my implementation
22 mnb = MyNaiveBayes()
23 y_pred = mnb.implemented_classifier(X_train, y_train, X_test)
24 print("Accuracy of my implementation for Naive Bayes:", accuracy_score(y_test, y_pred))
```

Run: tema4

```
C:\Users\Ana\AppData\Local\Programs\Python\Python311\python.exe "C:\Users\Ana\Desktop\FMI\Year II\sem 1\Probabilitati si statistica\laborator\Teme\tema4.py"
Accuracy of implemented Naive Bayes: 0.9605568445475638
Accuracy of my implementation for Naive Bayes: 0.7107501933488012
Process finished with exit code 0
```



The screenshot shows a terminal window with the same output as the IDE's Run console. It displays the path to the Python executable, the command to run `tema4.py`, and the resulting accuracy scores for both the implemented and custom Naive Bayes classifiers.

```
tema4
C:\Users\Ana\AppData\Local\Programs\Python\Python311\python.exe "C:\Users\Ana\Desktop\FMI\Year II\sem 1\Probabilitati si statistica\laborator\Teme\tema4.py"
Accuracy of implemented Naive Bayes: 0.9605568445475638
Accuracy of my implementation for Naive Bayes: 0.7107501933488012
Process finished with exit code 0
```