

# Θεωρία Αποφάσεων

## Εργαστηριακή Άσκηση Χειμερινού Εξαμήνου 2021-2022

Τοτσίλα Διόνις  
ΑΜ : 1067532

Δουρούκας Γεώργιος  
ΑΜ : 1064872

### Περιεχόμενα

<b>1 Προ-επεξεργασία Δεδομένων, αναπαράστασή κειμένων στον διανυσματικό χώρο</b>	<b>2</b>
1.0.1 Προεπεξεργασία κειμένων	3
1.0.2 Stemming	3
1.0.3 TF-IDF	3
1.0.4 Κώδικας python3 για την προ-επεξεργασία των δεδομένων και την αναπαράσταση στον διανυσματικό χώρο	4
<b>2 Σύγκριση Διάφορων μοντέλων Μηχανικής Μάθησης</b>	<b>6</b>
2.1 Multinomial Naive Bayes	6
2.1.1 Κώδικας python3 Naive Bayes	7
2.2 Decision Trees	8
2.2.1 Κώδικας python3 Decision Trees	9
2.3 Extra Trees	10
2.3.1 Κώδικας python3 Extra Trees	10
2.4 Logistic Regression	11
2.4.1 Κώδικας python3 Logistic Regression	12
2.5 Neural Network	13
2.5.1 Κώδικας python3 Neural Network	13
2.6 Συμπεράσματα - Εφαρμογή σε απλό διακοσμητή chat room	15
2.6.1 Κώδικας python3 Server του chat room	16
2.6.2 Κώδικας python3 Client του chat room	19
<b>3 Παράρτημα</b>	<b>22</b>
<b>4 Βιβλιογραφία</b>	<b>24</b>

### Εισαγωγή

Σκοπός της εργασίας είναι να αναπτυχθεί ένα μοντέλο, το οποίο βασιζόμενο στις τεχνικές μηχανικής μάθησης, θα μπορεί να διακρίνει τα spam μηνύματα που λαμβάνει ένας χρήστης και να τον ενημερώνει για αυτά. Ως spam χαρακτηρίζεται ένα email το οποίο ο χρήστης δεν έχει επιλέξει αλλά ούτε και επιθυμεί να λαμβάνει και αποτελεί παγίδα που μπορεί να συναντήσει ένας χρήστης χρησιμοποιώντας την διαδικτυακή αλληλογραφία. Εφόσον τα παραπλανητικά μηνύματα δημιουργούν εμπόδιο στην ομαλή λειτουργία της αλληλογραφίας είναι αναγκαίο να βρεθούν λύσεις για την αποκατάσταση του προβλήματος. Οι διαδικασίες αυτές θα αναφερθούν αναλυτικά παρακάτω. Υπάρχουν διάφορων ειδών τέτοια μηνύματα με τα πιο επικίνδυνα να χαρακτηρίζονται τα phishing τα οποία έχουν στόχο να αλιεύσουν σημαντικά προσωπικά στοιχεία από το χρήστη (όπως στοιχεία τραπεζικού λογαριασμού) με χαρακτηριστικό ότι εμφανίζονται πλήρως πειστικά στο χρήστη και είναι ιδιαίτερα απειλητικά για αυτόν. Στην επίλυση όλων των παραπάνω θα βοηθήσει καταλυτικά η μηχανική μάθηση η οποία είναι ιδιαίτερα πολύτιμη καθώς επιτρέπει στους υπολογιστές να χρησιμοποιούνται για την αυτοματοποίηση των διαδικασιών λήψης αποφάσεων.

## 1 Προ-επεξεργασία Δεδομένων, αναπαράστασή κειμένων στον διανυσματικό χώρο

Για να είναι εφικτό το training των μοντέλων θα πρέπει τα κείμενα να αναπαρασταθούν στον διανυσματικό χώρο. Προτού όμως γίνει αυτό θα πρέπει αρχικά να γίνει μια προ-επεξεργασία του δοθέντος dataset[1] αλλά και να γίνει χρήση τεχνικών γλωσσικής τεχνολογίας οι οποίες θα κάνουν την διαδικασία ταξινόμησης πιο αποτελεσματική. Αρχικά μέσω της βιβλιοθήκης pandas[2] γίνεται ανάγνωση του dataset σε ένα dataframe, στο οποίο παρατηρείται ότι υπάρχουν 3 στήλες που δεν χρειάζονται και πρέπει να αποκοπούν. Στη συνέχεια δημιουργούνται μερικά διαγράμματα τα οποία παρουσιάζουν κάποιες γενικές πληροφορίες για την κατανομή του dataset, και τα χαρακτηριστικά που διακρίνουν την μια κατηγορία από την άλλη.

Παρατηρείται ότι το dataset είναι μη-ισορροπημένο και για να αποφευχθεί κάποιο overfit θα πρέπει να περιορίσουμε τα κείμενα της κατηγορίας ham, έτσι επιλέγοντα τυχαία κείμενα του ποσοστού της τάξεως 70% των αρχικών ham κειμένων

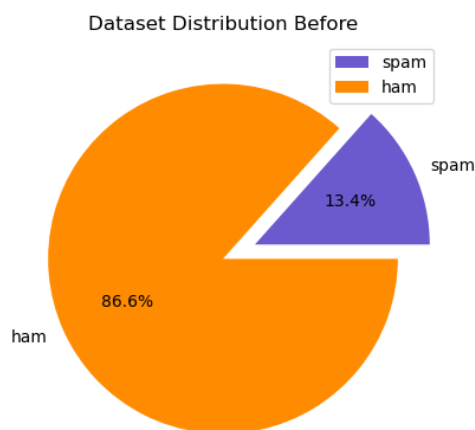


Figure 1: Κατανομή dataset πριν το undersampling

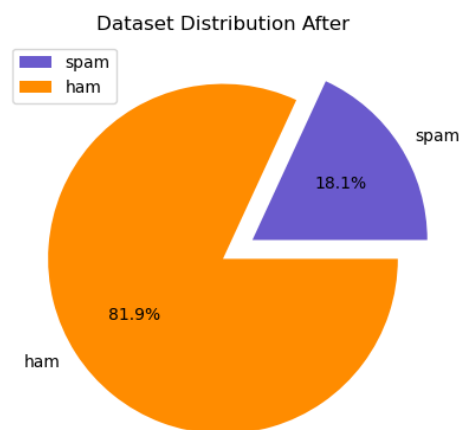


Figure 2: Κατανομή dataset μετά το undersampling

Παρατηρείτε επίσης ότι τα περισσότερα ham κείμενα περιέχουν από 0 έως 10 λέξεις ενώ τα περισσότερα spam από 20 έως 30, άρα το μήκος του κειμένου μπορεί να δώσει πληροφορία σχετικά με το αν ένα μήνυμα θα μπορούσε να είναι κακόβουλο ή όχι.

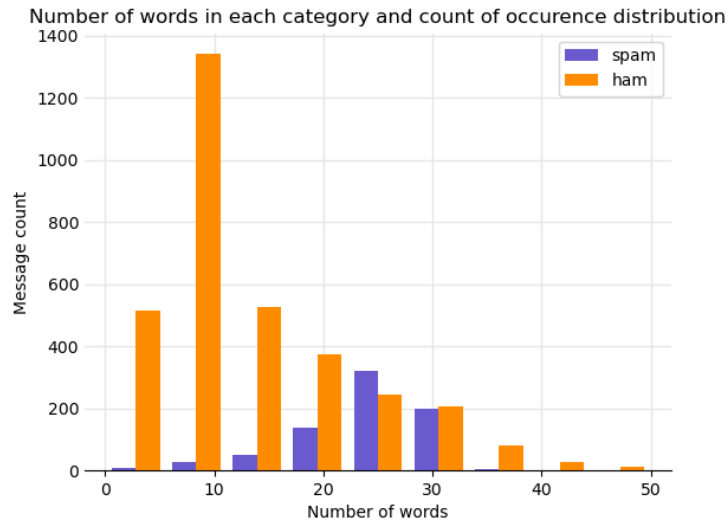


Figure 3: Κατανομή μήκους μηνύματος ανα κατηγορία

Είναι γνωστό πως τα περισσότερα spam μηνύματα έχουν άμεση σχέση με κάποιο χρηματικό ποσό, συνεπώς αφού παρατηρήθηκε ότι τα περισσότερα spam mail έχουν κάποιο χαρακτήρα συναλλάγματος όπως επίσης περιέχουν αριθμούς, κρίθηκε πως και αυτά τα 2 κριτήρια θα πρέπει προστεθούν στις παραμέτρους των διανυσμάτων.

### 1.0.1 Προεπεξεργασία κειμένων

Σειρά είχε η επεξεργασία των κειμένων, από κάθε κείμενο κρατήθηκαν μόνο οι λέξεις που αποτελούνται από αλφαριθμητικούς, καθώς τα σύμβολα και τα σημεία στίξης δεν μας δίνουν κάποια πληροφορία για το νοηματικό περιεχόμενο του κειμένου. Όλοι οι κεφαλαίοι χαρακτήρες μετατράπηκαν σε πεζούς και αφαιρέθηκαν τα stop words της αγγλικής (λέξεις που χρησιμοποιούνται συχνά σε κείμενα και κάνουν την διάκριση μη όμοιων νοηματικά προτάσεων).

### 1.0.2 Stemming

Γενικά για λόγους γραμματικής τα κείμενα έχουν διάφορες εκδοχές της ίδιας λέξης, όπως για παράδειγμα οι λέξεις organize, organizes και organizing, είναι προφανές ότι μια πρόταση που περιέχει τη λέξη organize και μια άλλη που περιέχει την λέξη organizes θα έχουν παρόμοιο νοηματικό περιεχόμενο. Στόχος λοιπόν του stemming είναι να "κανονικοποιήσει" παρόμοιες λέξεις δηλαδή να κάνει αναγωγή κλιτικών και παραγωγικών τύπων στο πρώτο κλιτικό. Για την επίτευξη του stemming χρησιμοποιήθηκε ο PorterStemmer της βιβλιοθήκης nltk[3].

### 1.0.3 TF-IDF

Αφού λοιπόν ολοκληρώθηκε η επεξεργασία των κειμένων, το κάθε κείμενο αναπαραστάθηκε ως διάνυσμα μέσω της μετρικής TF-IDF (Term Frequency - Inverse Document Frequency). Στόχος της μετρικής αυτής είναι η ανάθεση  $tf \cdot idf$  βάρους σε κάθε όρο του κάθε κειμένου και υπολογίζεται από τον τύπο:  $w_{ik} = tf_{ik} * \log(\frac{N}{df_i})$

Όπου:

$tf_{ik}$  = η συχνότητα εμφάνισης του όρου i στο κείμενο k

$df_k$  = Το πλήθος των κειμένων που περιέχουν τον i

N = το συνολικό πλήθος των κειμένων.

Για τον υπολογισμό του μητρώου tf-idf για την συλλογή κειμένων χρησιμοποιήθηκε η TfidfVectorizer της βιβλιοθήκης sklearn[4].

### 1.0.4 Κώδικας python3 για την προ-επεξεργασία των δεδομένων και την αναπαράσταση στον διανυσματικό χώρο

```

1  #!/usr/bin/env python
2  # coding: utf-8
3  import numpy as np
4  import pandas as pd
5
6  # NLP stuff
7  from sklearn.feature_extraction.text import TfidfVectorizer
8  from utils import processing, has_numbers, has_currency, sms_len, decorate_axis
9  import matplotlib.pyplot as plt
10
11 # set to True if you want to plot stuff
12 diag = False
13
14 # set to True if you want to show plots instead of saving them
15 show = False
16
17
18 # Read dataset into Pandas-Data Frame beware of the encoding
19 df = pd.read_csv('data/dataset.csv', encoding="ISO-8859-1")
20
21 # since columns Unnamed: # are empty we need to drop them and keep only the usefull data (v1,
22   v2) tag and sms respectively
23 df.drop(['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], axis=1, inplace=True)
24
25 # Replace the column names so that they are more identifiable and easily accessed
26 df = df.rename(columns={'v1': "category", "v2": "text"})
27
28 # we will balance the dataset a little bit
29 ham = df[df['category'] == 'ham']
30 spam = df[df['category'] == 'spam']
31
32 # plot dataset distribution before undersampling
33 if diag == True:
34     colors=['slateblue', 'darkorange']
35     plt.pie([len(spam), len(ham)], labels = ['spam', 'ham'], colors = colors, explode = [0.2,
36     0], autopct='%1.1f%%')
37     plt.legend()
38     plt.title('Dataset Distribution Before')
39     if show:
40         plt.show()
41     else:
42         plt.savefig("plots/data_dist_bef.png")
43         plt.close()
44
45 # get 70% of ham demos
46 ham = ham.sample(frac=0.7)
47
48 # plot dataset distribution after undersampling
49 if diag == True:
50     colors=['slateblue', 'darkorange']
51     plt.pie([len(spam), len(ham)], labels = ['spam', 'ham'], colors = colors, explode = [0.2,
52     0], autopct='%1.1f%%')
53     plt.legend()
54     plt.title('Dataset Distribution After')
55     if show:
56         plt.show()
57     else:
58         plt.savefig("plots/data_dist_aft.png")
59         plt.close()
60
61 # create and plot a list that contains number of words in messages and their count of
62   occurence in the dataset
63 if diag == True:
64     dataset_ham_count = ham['text'].str.split().str.len()
65     dataset_ham_count.index = dataset_ham_count.index.astype(str) + ' words:'
66     dataset_ham_count.sort_index(inplace=True)
67
68     dataset_spam_count = spam['text'].str.split().str.len()
69     dataset_spam_count.index = dataset_spam_count.index.astype(str) + ' words:'
70     dataset_spam_count.sort_index(inplace=True)
71     bins = np.linspace(0, 50, 10)
72     fig = plt.figure()

```

```

70     ax0 = fig.add_subplot()
71
72     plt.title('Number of words in each category and count of occurrence distribution')
73     ax0.hist([dataset_spam_count, dataset_ham_count], bins, label=['spam', 'ham'], color=
74     colors)
75     decorate_axis(ax0)
76     plt.legend(loc='upper right')
77     plt.xlabel("Number of words")
78     plt.ylabel("Message count")
79     if show:
80         plt.show()
81     else:
82         plt.savefig("plots/words_per_cat.png")
83         plt.close()
84
85 # concatenate spam and ham demos
86 df = pd.concat([spam,ham], ignore_index=True)
87
88 # shuffle dataset
89 df = df.sample(frac=1).reset_index(drop=True)
90
91 # apply processing to every SMS and store results to a new column
92 df['processed'] = df.apply(processing, axis=1)
93
94 # add a new column that indicates if the message has numerical values
95 df['has_num'] = df.apply(has_numbers, axis=1)
96
97 # add a new column that indicates if the message has currency related stuff
98 df['has_money'] = df.apply(has_currency, axis=1)
99
100 # add a new column that contains the length of each message
101 df['length'] = df.apply(sms_len, axis=1).values.reshape(-1,1)
102
103 # the original message is no longer needed
104 df = df.drop(['text'], axis=1)
105
106 # tfidf vector representation of texts
107 Transformer = TfidfVectorizer(max_features=2500, max_df=0.8)
108 tfidf = Transformer.fit_transform(df.processed.values.astype('U'))
109 tfidfDF = pd.DataFrame(tfidf.todense())
110 tfidfDF.columns = sorted(Transformer.vocabulary_)
111
112 # store idf of corpus needed for generatoin of tf idf for new (unseen text)
113 idf = pd.DataFrame(Transformer.idf_.transpose())
114 idf = pd.DataFrame(data=idf.values, columns=tfidfDF.columns)
115
116 final = pd.concat([df, tfidfDF], axis=1)
117 final = final.drop(['processed'], axis=1)
118
119 # split and store train-test because every method will be trained and tested on the same data,
120 # so that the comparisons are as fair as possible
121 train = final.head(int(np.floor(len(final) * 0.8)))
122 test = final.tail(int(np.ceil(len(final) - len(final) * 0.8)))
123
124 train.to_csv('data/train.csv')
125 test.to_csv('data/test.csv')

```

## 2 Σύγκριση Διάφορων μοντέλων Μηχανικής Μάθησης

Παρακάτω ακολουθεί ο σχολιασμός του κάθε μοντέλου, για τα περισσότερα μοντέλα έγινε η χρήση της βιβλιοθήκης sklearn, για την λογιστική παρεμβολή και τα νευρωνικά δίκτυα έγινε χρήση της βιβλιοθήκης pytorch[5], καθώς θέλαμε να δούμε μια πιο πρακτική ανάπτυξη κώδικα για μοντέλα μηχανικής μάθησης

### 2.1 Multinomial Naive Bayes

Ο Naive Bayes[6] είναι ένας από τους απλούς και πιο αποτελεσματικούς αλγόριθμους επιβλεπόμενης μάθησης που είναι γνωστός και ως πιθανοτικός ταξινομητής, καθώς προβλέπει με βάση τη πιθανότητα ενός αντικειμένου. Είναι κατάλληλος για δυαδική ταξινόμηση και έχει καλή απόδοση σε περιπτώσεις κατηγορικών μεταβλητών εισόδου σε σύγκριση με αριθμητικές μεταβλητές. Ακόμη είναι χρήσιμος για την πραγματοποίηση προβλέψεων και την πρόβλεψη δεδομένων με βάση τα ιστορικά αποτελέσματα. Ο ταξινομητής αυτός λειτουργεί με βάση την αρχή της υπό όρους πιθανότητας, όπως δίνεται από το θεώρημα bayes  $P(\theta|\mathbf{D}) = P(\theta) \frac{P(\mathbf{D}|\theta)}{P(\mathbf{D})}$

Ο πολυωνυμικός αλγόριθμος Naive Bayes[7] είναι μια πιθανολογική μέθοδος εκμάθησης που χρησιμοποιείται κυρίως στην επεξεργασία φυσικής γλώσσας(NLP). Ο αλγόριθμος βασίζεται στο θεώρημα Bayes και προβλέπει την ετικέτα ενός κειμένου όπως ένα κομμάτι email ή ένα άρθρο εφημερίδας. Υπολογίζει την πιθανότητα κάθε ετικέτας για ένα δεδομένο δείγμα και στη συνέχεια δίνει την ετικέτα με την υψηλότερη πιθανότητα ως έξοδο. Πρόκειται για μια συλλογή πολλών αλγορίθμων όπου οι αλγόριθμοι μοιράζονται μια κοινή αρχή και αυτή είναι ότι κάθε χαρακτηριστικό που ταξινομείται δεν σχετίζεται με κανένα άλλο χαρακτηριστικό. Η παρουσία ή η απουσία ενός χαρακτηριστικού δεν επηρεάζει την παρουσία ή απουσία του άλλου χαρακτηριστικού και για αυτό το λόγο χρησιμοποιήθηκε ο Multinomial Naive Bayes, η επιτυχία του classifier ανέρχεται στο 97%. Ακολουθεί το confusion matrix του classification.

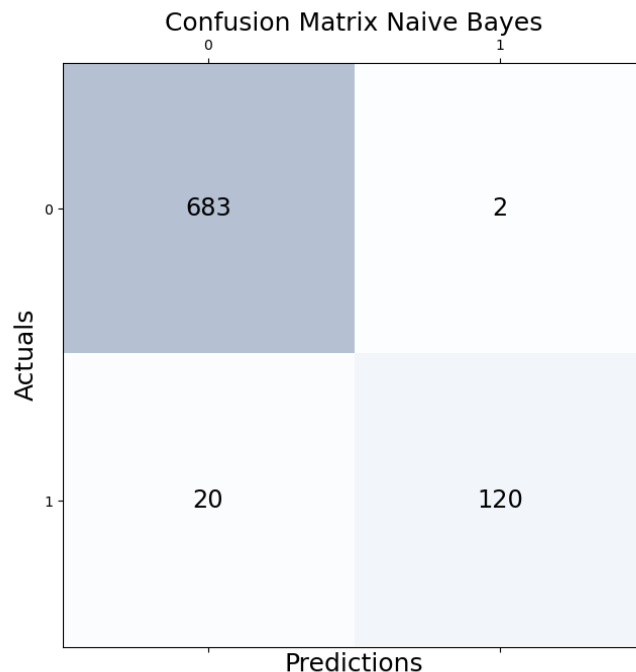


Figure 4: Polynomial Naive Bayes Classification Confusion Matrix

### 2.1.1 Κώδικας python3 Naive Bayes

```
1 from utils import read_data, plot_cf
2 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
3 from sklearn.naive_bayes import MultinomialNB
4 import pickle
5
6 # Read Train Test data
7 X_train, y_train = read_data('train')
8 X_test, y_test = read_data('test')
9
10 # init classifier and train
11 nb = MultinomialNB().fit(X_train,y_train)
12
13 # test
14 y_pred = nb.predict(X_test)
15
16 # Classification Evaluation
17 plot_cf(confusion_matrix(y_test,y_pred),"Naive Bayes")
18 print(classification_report(y_test,y_pred))
19 print(accuracy_score(y_test,y_pred))
20
21 # Store trained model
22 pickle.dump(nb, open('models/nb.pkl', 'wb'))
```

## 2.2 Decision Trees

Τα δέντρα αποφάσεων[8] είναι ένας τύπος εποπτευόμενης μηχανικής μάθησης όπου τα δεδομένα διαχωρίζονται συνεχώς σύμφωνα με μια συγκεκριμένη παράμετρο. Ο στόχος ενός δέντρου αποφάσεων είναι να δημιουργηθεί ένα μοντέλο εκπαίδευσης που μπορεί να χρησιμοποιηθεί για να προβλέψει την κλάση ή την τιμή της μεταβλητής στόχου μαθαίνοντας απλούς κανόνες απόφασης που συνάγονται από προηγούμενα δεδομένα. Αποτελείται από τον εσωτερικό κόμβο, τον κλάδο, και το κόμβο φύλλων. Κάθε εσωτερικός κόμβος αντιπροσωπεύει μια δοκιμή σε ένα χαρακτηριστικό. Κάθε κλάση αντιπροσωπεύει το αποτέλεσμα της δοκιμής και κάθε κόμβος φύλλων αντιπροσωπεύει μια ετικέτα κλάσης.

Βήμα 1: Εκκίνηση του δέντρου από τον ριζικό κόμβο, ο οποίος περιέχει το πλήρες σύνολο δεδομένων.

Βήμα 2: Εντοπίζεται το καλύτερο σύνολο δεδομένων χρησιμοποιώντας το Attribute Selection Measure(ASM).

Βήμα 3: Διαίρεται ο ριζικός κόμβος σε υποσύνολα που περιέχουν πιθανές τιμές.

Βήμα 4: Δημιουργία κόμβου αποφάσεων ο οποίος περιέχει το καλύτερο χαρακτηριστικό.

Βήμα 5: Δημιουργία αναδρομικών νέων δέντρων αποφάσεων χρησιμοποιώντας τα υποσύνολα του συνόλου δεδομένων που δημιουργήθηκαν στο προηγούμενο βήμα.

Η αποδοτικότητα των decision trees ανέρχεται στο 95%. Ακολουθεί το confusion matrix του classification.

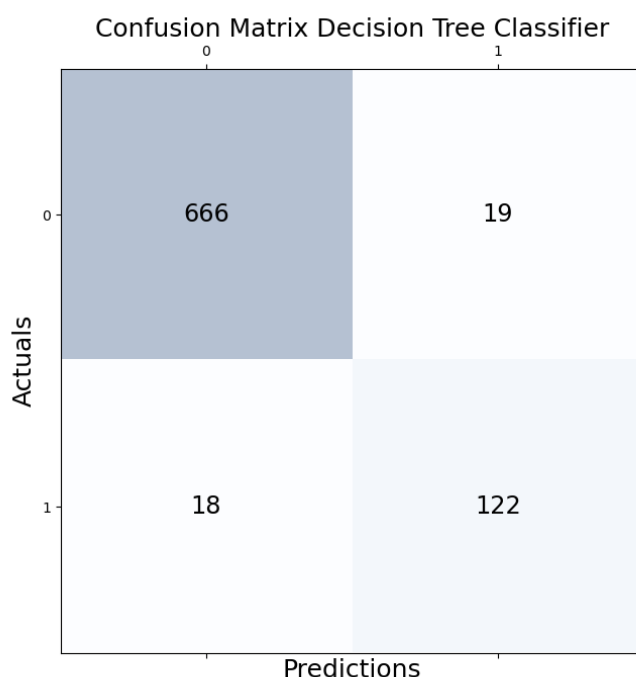


Figure 5: Decision Trees Classification Confusion Matrix



### 2.2.1 Κώδικας python3 Decision Trees

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
3 from utils import read_data, plot_cf
4 import pickle
5
6 # Read Train Test data
7 X_train, y_train = read_data('train')
8 X_test, y_test = read_data('test')
9
10 # init classifier
11 dtc = DecisionTreeClassifier(min_samples_split=7, random_state=252)
12
13 # train
14 dtc.fit(X_train, y_train)
15
16 # test
17 y_pred = dtc.predict(X_test)
18
19 # Classification Evaluation
20 plot_cf(confusion_matrix(y_test, y_pred), "Decision Tree Classifier")
21 print(classification_report(y_test, y_pred))
22 print(accuracy_score(y_test, y_pred))
23
24 # Store trained model
25 pickle.dump(dtc, open('models/dtc.pkl', 'wb'))
```

## 2.3 Extra Trees

Τα Extra Trees[9] είναι ένας αλγόριθμος μηχανικής μάθησης που συνδυάζει τις προβλέψεις από πολλά δέντρα αποφάσεων. Συχνά μπορεί να επιτύχει εξίσου καλή ή καλύτερη απόδοση από τον αλγόριθμο τυχαίων δασών. Επίσης έχουν το χαρακτηριστικό ότι ενώ προσθέτουν τυχαία δέντρα στο δάσος τους παράλληλα εξακολουθούν να έχουν βελτιστοποίηση. Οι προβλέψεις γίνονται με τη λήψη μέσου όρου πρόβλεψης των δέντρων απόφασης σε περίπτωση παλινδρόμησης ή χρησιμοποιώντας την πλειοψηφία στην περίπτωση ταξινόμησης. Η αποδοτικότητα των Extra trees ανέρχεται στο 97.6%. Ακολουθεί το confusion matrix του classification.

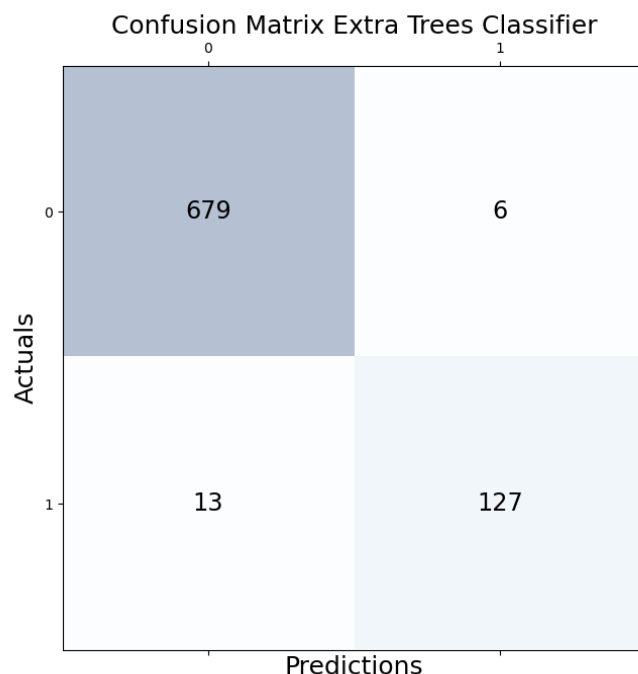


Figure 6: Extra Trees Classification Confusion Matrix

### 2.3.1 Κώδικας python3 Extra Trees

```

1 from sklearn.ensemble import ExtraTreesClassifier
2 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
3 import pickle
4 from utils import read_data, plot_cf
5
6 # Read Train Test Data
7 X_train, y_train = read_data('train')
8 X_test, y_test = read_data('test')
9
10 # init classifier
11 etc = ExtraTreesClassifier(n_estimators=37, random_state=252)
12
13 # Train
14 etc.fit(X_train, y_train)
15
16 # Test
17 y_pred = etc.predict(X_test)
18
19
20 # Classification Evaluation
21 plot_cf(confusion_matrix(y_test, y_pred), "Extra Trees Classifier")
22 print(classification_report(y_test, y_pred))
23 print(accuracy_score(y_test, y_pred))
24
25 # Store trained model
26 pickle.dump(etc, open('models/etc.pkl', 'wb'))

```

## 2.4 Logistic Regression

Η λογιστική παλινδρόμηση[10] είναι ένας από τους πιο δημοφιλείς αλγόριθμους μηχανικής μάθησης ο οποίος εμπίπτει στην τεχνική εποπτευόμενης μάθησης και αποτελεί έναν πολύτιμο αλγόριθμο που χρησιμεύει για την εκτίμηση διακριτών τιμών όπως 0/1, ναι/όχι, αληθές/ψευδές. Χρησιμοποιείται για τα προβλήματα κατηγοριοποίησης και είναι ένας αλγόριθμος προγνωστικής ανάλυσης και βασίζεται στην έννοια της πιθανότητας. Λειτουργεί με μία σύνθετη συνάρτηση κόστους η οποία μπορεί να οριστεί ως λογιστική συνάρτηση. Sigmoid function είναι μια μαθηματική συνάρτηση που παίρνει οποιοδήποτε πραγματικό αριθμό και τον αντιστοιχίζει σε μια πιθανότητα μεταξύ 1 και 0. Η αποδοτικότητα της Λογιστικής Παρεμβολής ανέρχεται στο 97.2%. Ακολουθεί το confusion matrix του classification.

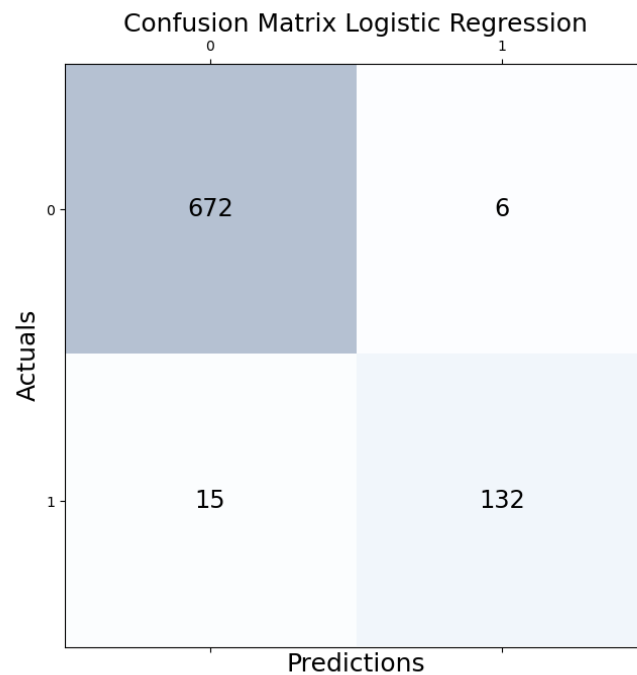


Figure 7: Logistic Regression Confusion Matrix

### 2.4.1 Κώδικας python3 Logistic Regression

```

1 import torch
2 from torch.utils.data import DataLoader
3 from utils import CustomDataset, read_data, plot_cf
4 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
5 train = True
6
7 # Simple Logistic Regression Module (sigmoid activated)
8 class LogisticRegression(torch.nn.Module):
9     def __init__(self, input_dim, output_dim):
10         super(LogisticRegression, self).__init__()
11         self.linear = torch.nn.Linear(input_dim, output_dim)
12
13     def forward(self, x):
14         outputs = torch.sigmoid(self.linear(x))
15         return outputs
16
17 # Read Test Data
18 X_train, y_train = read_data('train')
19 X_test, y_test = read_data('test')
20 train_dataset = CustomDataset(X_train, y_train)
21
22 # Train in batches
23 dataloader = DataLoader(train_dataset, batch_size=128, shuffle=True)
24 batches_per_epoch = X_train.shape[0]//128
25
26 # Hyper Parameters
27 epochs = 1000
28 learning_rate = 1e-2
29
30 # 1 binary output (binary classification)
31 output_dim = 1
32 input_dim = X_test.shape[1]
33 model = LogisticRegression(input_dim, output_dim)
34 criterion = torch.nn.BCELoss()
35 optimizer = torch.optim.AdamW(model.parameters(), lr=learning_rate)
36
37 if train:
38     # train loop
39     for epoch in range(epochs):
40         running_loss = 0.0
41         for i, data in enumerate(dataloader, 0):
42             inputs, outputs = data
43             batch_loss = 0.0
44             # zero the parameter gradients
45             optimizer.zero_grad()
46
47             # forward + backward + optimize
48             output = model(inputs.view(-1, input_dim))
49             loss = criterion(output, outputs.view(-1, output.shape[1]))
50             running_loss += loss.item()
51             loss.backward()
52             optimizer.step()
53         train_mean_loss = running_loss/batches_per_epoch
54         print("Epoch: ", epoch+1, "Loss: ", train_mean_loss)
55     print('Finished Training')
56
57     # store trained model
58     torch.save(model.state_dict(), 'models/logistic_regression.pt')
59 else:
60     model.load_state_dict(torch.load('models/logistic_regression.pt'))
61 model.eval()
62 # test
63 # no derivative calculation so that the calculations are faster
64 with torch.no_grad():
65     y_pred = []
66     for i in range(X_test.shape[0]):
67         # round to 0 or 1 so that class is determined
68         output = torch.round(model(X_test[i])).detach().numpy()
69         y_pred.append(output)
70     plot_cf(confusion_matrix(y_test, y_pred), "Logistic Regression")
71     print(confusion_matrix(y_test, y_pred))
72     print(classification_report(y_test, y_pred))
73     print(accuracy_score(y_test, y_pred))

```

## 2.5 Neural Network

Τα νευρωνικά δίκτυα[11] , γνωστά και ως τεχνητά νευρωνικά δίκτυα ή προσομοιωμένα νευρωνικά δίκτυα, αποτελούν υποσύνολο της μηχανικής μάθησης και βρίσκονται στην καρδιά των αλγορίθμων βαθιάς μάθησης. Είναι μια σειρά αλγορίθμων που προσπαθεί να αναγνωρίσει τις υποκείμενες σχέσεις σε ένα σύνολο δεδομένων μέσω μιας διαδικασίας που μιμείται τον τρόπο λειτουργίας του ανθρώπινου εγκεφάλου. Επίσης μπορούν να προσαρμοστούν στις μεταβαλλόμενες εισόδους και έτσι το δίκτυο παράγει το καλύτερο δυνατό αποτέλεσμα χωρίς να χρειάζεται επανασχεδιασμός των κριτηρίων εξόδου. Οι πληροφορίες τροφοδοτούνται στο επίπεδο εισόδου το οποίο τις μεταφέρει στο κρυφό επίπεδο. Οι διασυνδέσεις μεταξύ των δύο επιπέδων εκχωρούν βάρη σε κάθε είσοδο τυχαία. Η loss function σε ένα νευρωνικό δίκτυο ποσοτικοποιεί τη διαφορά μεταξύ του αναμενόμενου αποτελέσματος και του αποτελέσματος που παράγεται από το μοντέλο μηχανικής μάθησης. Συγκεκριμένα η binary cross entropy συγκρίνει κάθε μια από τις προβλεπόμενες πιθανότητες με την έξοδο κλάσης που μπορεί να είναι είτε 0 είτε 1. Στη συνέχεια υπολογίζει τις πιθανότητες με βάση την απόσταση από την αναμενόμενη τιμή. Αυτό δείχνει πόσο κοντά ή μακριά είναι από την πραγματική τιμή. Ο Adam[12] είναι ένας προσαρμοστικός αλγόριθμος βελτιστοποίησης ρυθμού μάθησης που έχει σχεδιαστεί ειδικά για την εκπαίδευση, την κατάλληλη δηλαδή προσαρμογή των βαρών που προαναφέραμε, σε βαθιά νευρωνικά δίκτυα. Συνδυάζει τις καλύτερες ιδιότητες των αλγορίθμων AdaGrad[13] και RMSProp για να παρέχει έναν αλγόριθμο βελτιστοποίησης που μπορεί να χειριστεί αραιές κλίσεις σε θορυβώδη προβλήματα. Η αποδοτικότητα των νευρωνικών δικτύων ανέρχεται στο 97.45%. Ακολουθεί το confusion matrix του classification.

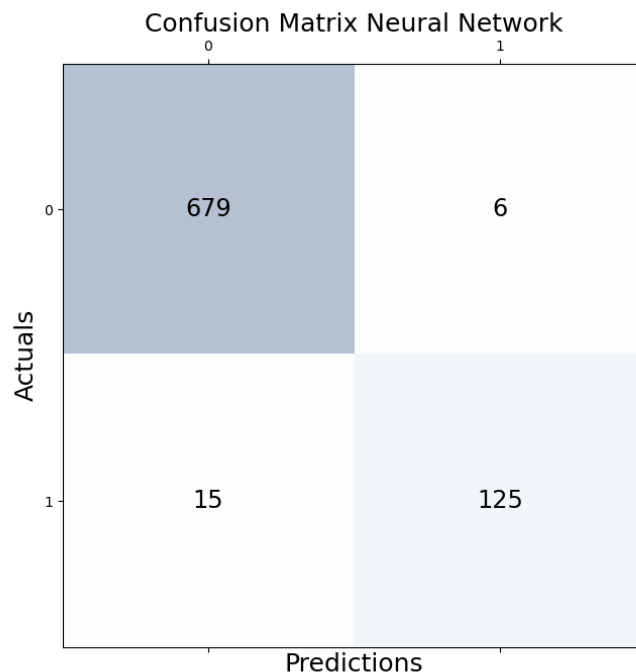


Figure 8: Neural Network Confusion Matrix

### 2.5.1 Κώδικας python3 Neural Network

```

1 import torch
2 from torch.utils.data import DataLoader
3 from utils import CustomDataset, read_data, plot_cf
4 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
5 import torch.nn as nn
6 import torch.nn.functional as F
7
8 train = False
9 # Read Test Data
10 X_train, y_train = read_data('train')
11 X_test, y_test = read_data('test')
12 train_dataset = CustomDataset(X_train, y_train)
13 print(X_train.shape)
14 print(X_test.shape)
15

```

```

16 dataloader = DataLoader(train_dataset, batch_size=128, shuffle=True)
17
18
19 # Hyper Parameters
20 epochs = 500
21
22 # input_dim = tfidf vector length
23 input_dim = X_train.shape[1]
24
25 # 1 binary output (binary classification)
26 output_dim = 1
27 learning_rate = 1e-3
28 class BinaryClassification(nn.Module):
29     def __init__(self, input_dim):
30         super(BinaryClassification, self).__init__()
31         self.layer_1 = nn.Linear(input_dim, 512)
32         self.layer_2 = nn.Linear(512, 128)
33         self.layer_3 = nn.Linear(128, 32)
34         self.layer_out = nn.Linear(32, 1)
35         self.relu = nn.ReLU()
36
37
38     def forward(self, inputs):
39         x = self.relu(self.layer_1(inputs))
40         x = self.relu(self.layer_2(x))
41         x = self.relu(self.layer_3(x))
42         x = torch.sigmoid(self.layer_out(x))
43
44
45         return x
46
47 model = BinaryClassification(input_dim)
48 criterion = torch.nn.BCELoss()
49 optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
50 batches_per_epoch = X_train.shape[0]//128
51
52 if train:
53     for epoch in range(epochs):
54         running_loss = 0.0
55         for i, data in enumerate(dataloader, 0):
56             # get the inputs; data is a list of [inputs, output]
57             inputs, outputs = data
58             # print(inputs)
59             batch_loss = 0.0
60             # zero the parameter gradients
61             optimizer.zero_grad()
62             # forward + backward + optimize
63             output = model(inputs.view(-1, input_dim))
64             # print(output)
65             loss = criterion(output, outputs.view(-1, output.shape[1]))
66             running_loss += loss.item()
67             loss.backward()
68             optimizer.step()
69             # scheduler.step(loss.item())
70         # scheduler.step()
71         train_mean_loss = running_loss/batches_per_epoch
72
73         print("Epoch: ", epoch+1, "Loss: ", train_mean_loss)
74     print('Finished Training')
75     torch.save(model.state_dict(), 'models/nn.pt')
76 else:
77     model.load_state_dict(torch.load('models/nn.pt'))
78 with torch.no_grad():
79     model.eval()
80     y_pred = []
81     for i in range(X_test.shape[0]):
82         # round to 0 or 1 so that class is determined
83         output = torch.round(model(X_test[i])).detach().numpy()
84         y_pred.append(output)
85     plot_cf(confusion_matrix(y_test, y_pred), "Neural Network")
86     print(confusion_matrix(y_test, y_pred))
87     print(classification_report(y_test, y_pred))
88     print(accuracy_score(y_test, y_pred))

```

## 2.6 Συμπεράσματα - Εφαρμογή σε απλό διακοσμητή chat room

Classifier	tag	precision	recall	f1-score	support
Multinomial Naive Bayes	ham	0.97	1.0	0.98	685
Multinomial Naive Bayes	spam	0.98	0.86	0.92	140
Decision Trees	ham	0.97	0.97	0.97	685
Decision Trees	spam	0.87	0.87	0.87	140
Extra Trees	ham	0.98	0.99	0.99	685
Extra Trees	spam	0.95	0.91	0.93	140
Logistic Regression	ham	0.99	0.98	0.98	685
Logistic Regression	spam	0.91	0.93	0.92	140
Neural Network	ham	0.98	0.99	0.98	685
Neural Network	spam	0.95	0.89	0.92	140

Από τα παραπάνω καταλήξαμε στην επιλογή των τριών πιο αποδοτικών μοντέλων μέσα από τα οποία κάθε μήνυμα του chat room, αφού περάσει την προεπεξεργασία που αναφέραμε, σε αυτό το σημείο πρέπει να τονίσουμε ότι έχουμε αποθηκεύσει ένα διάνυσμα που περιέχει το idf του train dataset, έτσι ώστε να μπορεί να υπολογιστεί το χαρακτηριστικό διάνυσμα για κάθε νέο μήνυμα και να μην περιοριστούμε στην χρήση μόνο των προτάσεων που δίνονται από το dataset. Σε περίπτωση που ένα μήνυμα χαρακτηριστεί ως spam από τουλάχιστον δύο μοντέλα τότε στέλνεται ένα μήνυμα προς όλους τους χρήστες το οποίο λέει ότι ένας συγκεκριμένος χρήστης προσπάθησε να στείλει ένα spam μήνυμα, σε αντίθετη περίπτωση το μήνυμα αποστέλεται κανονικά. Παρακάτω παραθέτονται εικόνες με παραδείγματα χρήσης του chat room για την κατασκευή του οποίου συμβουλευτήκαμε το παρακάτω άρθρο <https://www.geeksforgeeks.org/gui-chat-application-using-tkinter-in-python/>

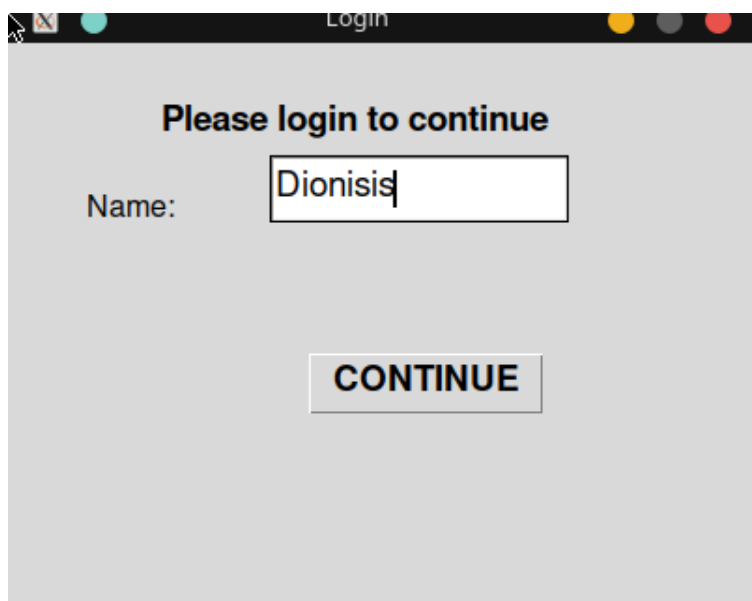


Figure 9: Σύνδεση Χρήστη στο chat room

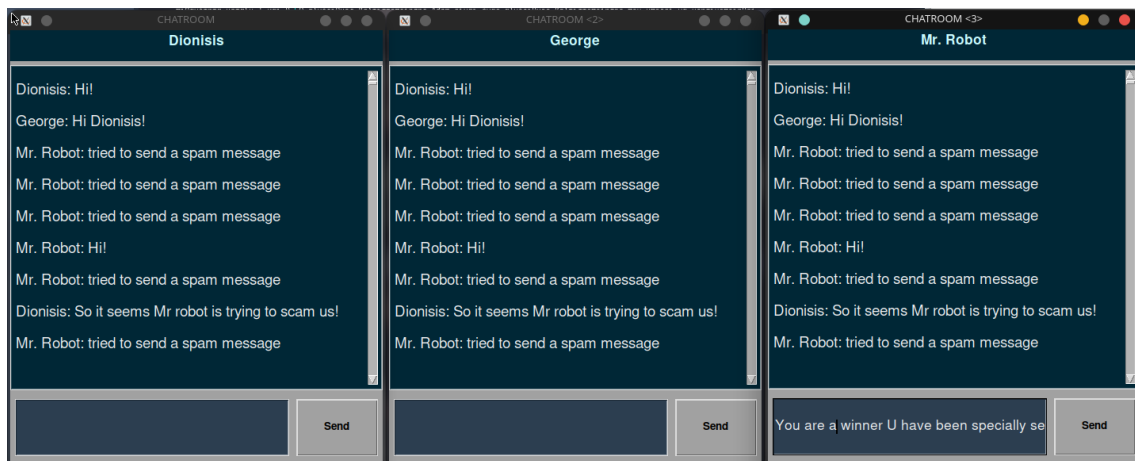


Figure 10: Παράδειγμα αποτροπής αποστολής ενός spam μηνύματος στο chat

### 2.6.1 Κώδικας python3 Server του chat room

```

1 import socket
2 import threading
3 from collections import Counter
4 import pandas as pd
5 from utils import *
6 import pickle
7 import torch.nn as nn
8 import torch
9 class BinaryClassification(nn.Module):
10     def __init__(self, input_dim):
11         super(BinaryClassification, self).__init__()
12         self.layer_1 = nn.Linear(input_dim, 512)
13         self.layer_2 = nn.Linear(512, 128)
14         self.layer_3 = nn.Linear(128, 32)
15         self.layer_out = nn.Linear(32, 1)
16         self.relu = nn.ReLU()
17
18
19     def forward(self, inputs):
20         x = self.relu(self.layer_1(inputs))
21         x = self.relu(self.layer_2(x))
22         x = self.relu(self.layer_3(x))
23         x = torch.sigmoid(self.layer_out(x))
24
25
26         return x
27 import numpy as np
28
29 # read columns template pandas series representing a tf idf vector for every message to be
30 # classified
31 template = pd.read_csv('data/test.csv',nrows=1)
32
33 # read idf needed for the tfidf vector representation of every new message for every term of
34 # the message
35 # that is in the corpus vocabulary tf * idf we used this method since we wanted to calculate
36 # the tfidf
37 # vector even for messages that have not been already represented
38 idf = pd.read_csv("data/idf.csv")
39
40 # load extra trees trained model
41 extra_trees = pickle.load(open('models/etc.pkl', 'rb'))
42 # load extra trees trained model
43 random_forest = pickle.load(open('models/rf.pkl', 'rb'))
44
45 # load trained neural net
46 model = BinaryClassification(len(template.iloc[0])-2)
47 model.load_state_dict(torch.load('models/nn.pt'))
48 model.eval()
49
50 PORT = 5000

```



```

50
51 # An IPv4 address is obtained
52 # for the server.
53 SERVER = socket.gethostbyname(socket.gethostbyname())
54
55 # Address is stored as a tuple
56 ADDRESS = (SERVER, PORT)
57
58 FORMAT = "utf-8"
59
60 # Lists that will contains
61 # all the clients connected to
62 # the server and their names.
63 clients, names = [], []
64
65 # Create a new socket for
66 # the server
67 server = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
68
69 # bind the address of the
70 # server to the socket
71 server.bind(ADDRESS)
72
73 # function to start the connection
74 def startChat():
75
76     print("server is working on " + SERVER)
77
78     # listening for connections
79     server.listen()
80
81     while True:
82
83         # accept connections and returns
84         # a new connection to the client
85         # and the address bound to it
86         conn, addr = server.accept()
87         conn.send("NAME".encode(FORMAT))
88
89         # 1024 represents the max amount
90         # of data that can be received (bytes)
91         name = conn.recv(1024).decode(FORMAT)
92
93         # append the name and client
94         # to the respective list
95         names.append(name)
96         clients.append(conn)
97         print(f"Name is :{name}")
98         # broadcast message
99         broadcastMessage(f"{name} has joined the chat!".encode(FORMAT))
100        conn.send('Connection successful!'.encode(FORMAT))
101
102        # Start the handling thread
103        thread = threading.Thread(target=handle, args=(conn, addr))
104        thread.start()
105
106        # no. of clients connected
107        # to the server
108        print(f"active connections {threading.activeCount()-1}")
109
110    # method to handle the
111    # incoming messages
112    def handle(conn, addr):
113
114        print(f"new connection {addr}")
115        connected = True
116
117        while connected:
118            # receive message
119            message = conn.recv(1024)
120            print(message)
121            # broadcast message
122            broadcastMessage(str.encode(is_spam(message)))
123
124            # close the connection
125            conn.close()

```

```

126
127 # method for broadcasting
128 # messages to the each clients
129 def broadcastMessage(message):
130     for client in clients:
131         client.send(message)
132
133
134 def is_spam(message):
135     # read 3 of the best models and pass the text through them
136     # if at least one method considers that the message is spam instead of the message
137     # an error is broadcasted
138
139     # byte to str
140     message = message.decode()
141
142     # extract user from message
143     user = ""
144     for i in range(0, len(message)):
145         user += str(message[i])
146         if message[i] == ":":
147             break
148
149     # isolate message without the user:
150     new_mes = message[i+1:]
151
152     # preprocess message
153     processed = message_processing(new_mes)
154
155     # calculate parameters
156     length = len(new_mes)
157     has_num = has_numbers_text(new_mes)
158     has_cur = has_currency_text(new_mes)
159
160     # count term frequency for every word of the message
161     tf = Counter(processed)
162
163     # calculate tf-idf weight using corpus idf
164     for col in template.columns:
165         template[col].values[:] = 0.0
166
167
168     # update template vector (zeros initialized)
169     template.has_num.replace(0, has_num)
170     template.has_money.replace(0, has_cur)
171     template.length.replace(0, length)
172
173     # convert message to tf-idf vector
174     # load random forest
175     for key in tf:
176         try:
177             # operation to raise exception if word is not in corpus vocabulary
178             a = (template[key])
179             # update corresponding tf-idf value
180             template[key] = tf.get(key) * idf[key].iloc[[0]].values
181         except:
182             print("word not in corpus vocabulary")
183
184     # sum of the models output (range 0 to 3 )
185     sum = 0
186     # extract vector to be inserted in models
187     input = template.iloc[0][2:].copy().to_numpy().reshape(-1, len(template.iloc[0][2:]))
188     # extra trees output
189     sum += extra_trees.predict(input)
190     # random forest output
191     sum += random_forest.predict(input)
192     # neural net output
193     with torch.no_grad():
194         sum += torch.round(model(torch.Tensor(input.astype(float)))).item()
195     # if no more than 1 model classified the message as spam return it as is
196     if (sum < 2):
197         return(message)
198     # if at least two model marked the message as spam
199     else:
200         # Print in console that a spam message was found by the user
201         print("SPAM Message found from : ", user)

```

```

202         # return warning to be displayed so that all users know that a user tried to scam them
203         return(str(user)+" tried to send a spam message")
204
205     # call the method to
206     # begin the communication
207     startChat()

```

### 2.6.2 Κώδικας python3 Client του chat room

```

1  # import all the required modules
2  import socket
3  import threading
4  from tkinter import *
5  from tkinter import font
6  from tkinter import ttk
7
8  # import all functions /
9  # everything from chat.py file
10 from chat import *
11
12 PORT = 5000
13 SERVER = "127.0.1.1"
14 ADDRESS = (SERVER, PORT)
15 FORMAT = "utf-8"
16
17 # Create a new client socket
18 # and connect to the server
19 client = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
20 client.connect(ADDRESS)
21
22
23 # GUI class for the chat
24 class GUI:
25     # constructor method
26     def __init__(self):
27
28         # chat window which is currently hidden
29         self.Window = Tk()
30         self.Window.withdraw()
31
32         # login window
33         self.login = Toplevel()
34         # set the title
35         self.login.title("Login")
36         self.login.resizable(width=False, height=False)
37         self.login.configure(width=400, height=300)
38         # create a Label
39         self.pls = Label(self.login,
40                          text="Please login to continue",
41                          justify=CENTER,
42                          font="Helvetica 14 bold")
43
44         self.pls.place(relheight=0.15, relx=0.2, rely=0.07)
45         # create a Label
46         self.labelName = Label(self.login, text="Name: ", font="Helvetica 12")
47
48         self.labelName.place(relheight=0.2, relx=0.1, rely=0.2)
49
50         # create a entry box for
51         # tyoing the message
52         self.entryName = Entry(self.login, font="Helvetica 14")
53
54         self.entryName.place(relwidth=0.4, relheight=0.12, relx=0.35, rely=0.2)
55
56         # set the focus of the cursor
57         self.entryName.focus()
58
59         # create a Continue Button
60         # along with action
61         self.go = Button(self.login,
62                          text="CONTINUE",
63                          font="Helvetica 14 bold",
64                          command=lambda: self.goAhead(self.entryName.get()))
65
66         self.go.place(relx=0.4, rely=0.55)
67         self.Window.mainloop()
68

```

```

69     def goAhead(self, name):
70         self.login.destroy()
71         self.layout(name)
72
73         # the thread to receive messages
74         rcv = threading.Thread(target=self.receive)
75         rcv.start()
76
77     # The main layout of the chat
78     def layout(self, name):
79
80         self.name = name
81         # to show chat window
82         self.Window.deiconify()
83         self.Window.title("CHATROOM")
84         self.Window.resizable(width=False,
85                                height=False)
86         self.Window.configure(width=470,
87                                height=550,
88                                bg="#002736")
89         self.labelHead = Label(self.Window, bg="#002736",
90                                fg="#C9FBFF",
91                                text=self.name,
92                                font="Helvetica 13 bold",
93                                pady=5)
94
95         self.labelHead.place(relwidth=1)
96         self.line = Label(self.Window,
97                            width=450,
98                            bg="#A1A1A1")
99
100        self.line.place(relwidth=1,
101                        rely=0.07,
102                        relheight=0.012)
103
104        self.textCons = Text(self.Window,
105                              width=20,
106                              height=2,
107                              bg="#002736",
108                              fg="#EAECEE",
109                              font="Helvetica 14",
110                              padx=5,
111                              pady=5)
112
113        self.textCons.place(relheight=0.745,
114                            relwidth=1,
115                            rely=0.08)
116
117        self.labelBottom = Label(self.Window,
118                                  bg="#A1A1A1",
119                                  height=80)
120
121        self.labelBottom.place(relwidth=1,
122                                rely=0.825)
123
124        self.entryMsg = Entry(self.labelBottom,
125                               bg="#2C3E50",
126                               fg="#EAECEE",
127                               font="Helvetica 13")
128
129        # place the given widget
130        # into the gui window
131        self.entryMsg.place(relwidth=0.74,
132                            relheight=0.06,
133                            rely=0.008,
134                            relx=0.011)
135
136        self.entryMsg.focus()
137
138        # create a Send Button
139        self.buttonMsg = Button(self.labelBottom,
140                                text="Send",
141                                font="Helvetica 10 bold",
142                                width=20,
143                                bg="#A1A1A1",
144                                command=lambda: self.sendButton(self.entryMsg.get()))

```

```

145         self.buttonMsg.place(relx=0.77,
146                               rely=0.008,
147                               relheight=0.06,
148                               relwidth=0.22)
149
150
151         self.textCons.config(cursor="arrow")
152
153         # create a scroll bar
154         scrollbar = Scrollbar(self.textCons)
155
156         # place the scroll bar
157         # into the gui window
158         scrollbar.place(relheight=1,
159                        relx=0.974)
160
161         scrollbar.config(command=self.textCons.yview)
162
163         self.textCons.config(state=DISABLED)
164
165         # function to basically start the thread for sending messages
166         def sendButton(self, msg):
167             self.textCons.config(state=DISABLED)
168             self.msg = msg
169             self.entryMsg.delete(0, END)
170             snd = threading.Thread(target=self.sendMessage)
171             snd.start()
172
173         # function to receive messages
174         def receive(self):
175             while True:
176                 try:
177                     message = client.recv(1024).decode(FORMAT)
178
179                     # if the messages from the server is NAME send the client's name
180                     if message == 'NAME':
181                         client.send(self.name.encode(FORMAT))
182                     else:
183                         # insert messages to text box
184                         self.textCons.config(state=NORMAL)
185                         self.textCons.insert(END,
186                                             message+"\n\n")
187
188                         self.textCons.config(state=DISABLED)
189                         self.textCons.see(END)
190                 except:
191                     # an error will be printed on the command line or console if there's an error
192                     print("An error occurred!")
193                     client.close()
194                     break
195
196         # function to send messages
197         def sendMessage(self):
198             self.textCons.config(state=DISABLED)
199             while True:
200                 message = (f"{self.name}: {self.msg}")
201                 client.send(message.encode(FORMAT))
202                 break
203
204
205         # create a GUI class object
206         g = GUI()

```

### 3 Παράρτημα

Παρακάτω ακολουθεί μια βιβλιοθήκη που περιέχει συναρτήσεις που χρησιμοποιούνται συχνά από διάφορα scripts.

```

1 import torch
2 import pandas as pd
3 from torch.utils.data import Dataset
4 from nltk import word_tokenize
5 from nltk.stem import PorterStemmer
6 from nltk.corpus import stopwords as sw
7 import re
8
9
10 # text processing
11 ps = PorterStemmer()
12 stop_words = set(sw.words("english"))
13
14 # check if pandas series text field has a number
15 def has_numbers(row):
16     text = row["text"]
17     string = word_tokenize(text)
18     for word in string:
19         if bool(re.search(r'\d', word)):
20             return 1
21     return 0
22
23 # check if pandas series text field has a currency symbol
24 def has_currency(row):
25     text = row["text"]
26     string = word_tokenize(text)
27     for word in string:
28         if('$' in word or '€' in word or '£' in word):
29             return 1
30     return 0
31
32 # check if a string has numbers
33 def has_numbers_text(text):
34     string = word_tokenize(text)
35     for word in string:
36         if bool(re.search(r'\d', word)):
37             return 1
38     return 0
39
40 # check if a string has currency symbols
41 def has_currency_text(text):
42     string = word_tokenize(text)
43     for word in string:
44         if('$' in word or '€' in word or '£' in word):
45             return 1
46     return 0
47
48 # calculate text field length from pandas series
49 def sms_len(row):
50     return len(row["text"])
51
52 # stem words that are not stop words and are alphabetic strings from text field of pandas
53 # series
54 def processing(row):
55     text = row["text"]
56     tokens = word_tokenize(text)
57     stemmed_tokens = []
58     stemmed_tokens = [ps.stem(word.lower()) for word in tokens if (word not in stop_words and
59 word.isalpha())]
60     joined = (" ".join(stemmed_tokens))
61
62     return joined
63
64 # stem words that are not stop words and are alphabetic strings from string
65 def message_processing(text):
66     tokens = word_tokenize(text)
67     stemmed_tokens = []
68     stemmed_tokens = [ps.stem(word.lower()) for word in tokens if (word not in stop_words and
69 word.isalpha())]
70     joined = (" ".join(stemmed_tokens))
71     return word_tokenize(joined)
72
73 # list or numpy dataset to torch tensor

```

```

71 class CustomDataset(Dataset):
72     """ Custom Spam or Ham Dataset. """
73
74     def __init__(self, x, y):
75         """
76         Args:
77             x (numpy array): # tf-idf vector representation of message
78             y (numpy array): # class spam or ham
79         """
80         self.x = torch.Tensor(x)
81         self.y = torch.Tensor(y)
82
83     def __len__(self):
84         return len(self.x)
85
86     def __getitem__(self, idx):
87         return self.x[idx], self.y[idx]
88
89 # read train or test dataset
90 def read_data(id):
91     '''
92     id = test or train
93     '''
94     # Read Train Data
95     df = pd.read_csv('data/' + id + '.csv', index_col=0)
96     df.category = pd.factorize(df.category)[0]
97     X, y = df.iloc[:, 1:].to_numpy(), df["category"].to_numpy()
98     X = torch.Tensor(X)
99     y = torch.Tensor(y)
100     return X, y
101
102 # plot confusion matrix
103 def plot_cf(conf_matrix, name, show=False):
104     import matplotlib.pyplot as plt
105     fig, ax = plt.subplots(figsize=(7.5, 7.5))
106     ax.matshow(conf_matrix, cmap=plt.cm.Blues, alpha=0.3)
107     for i in range(conf_matrix.shape[0]):
108         for j in range(conf_matrix.shape[1]):
109             ax.text(x=j, y=i, s=conf_matrix[i, j], va='center', ha='center', size='xx-large')
110
111     plt.xlabel('Predictions', fontsize=18)
112     plt.ylabel('Actuals', fontsize=18)
113     plt.title('Confusion Matrix ' + name, fontsize=18)
114     if show:
115         plt.show()
116     else:
117         plt.savefig("plots/" + name + '_cf.png')
118
119 # decorate axis from plot
120 def decorate_axis(ax, remove_left=True):
121     ax.spines['top'].set_visible(False)
122     ax.spines['right'].set_visible(False)
123     ax.spines['left'].set_visible(False)
124     ax.get_xaxis().tick_bottom()
125     ax.get_yaxis().tick_left()
126     ax.tick_params(axis='x', direction='out')
127     ax.tick_params(axis='y', length=0)
128
129     ax.grid(axis='y', color="0.9", linestyle='--', linewidth=1)
130     ax.grid(axis='x', color="0.9", linestyle='--', linewidth=1)
131     ax.set_axisbelow(True)

```

## 4 Βιβλιογραφία

- [1] Tiago A. Almeida, José María Gómez Hidalgo, and Akebo Yamakami. “Contributions to the study of SMS spam filtering: new collection and results”. In: *DocEng '11*. 2011.
- [2] Wes McKinney. “Data Structures for Statistical Computing in Python”. In: *Proceedings of the 9th Python in Science Conference*. Ed. by Stéfan van der Walt and Jarrod Millman. 2010, pp. 56–61. DOI: [10.25080/Majora-92bf1922-00a](https://doi.org/10.25080/Majora-92bf1922-00a).
- [3] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. ” O’Reilly Media, Inc.”, 2009.
- [4] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [5] Adam Paszke et al. “PyTorch: An Imperative Style, High-Performance Deep Learning Library”. In: *Advances in Neural Information Processing Systems 32*. Ed. by H. Wallach et al. Curran Associates, Inc., 2019, pp. 8024–8035. URL: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [6] Geoffrey I. Webb. “Naïve Bayes”. In: *Encyclopedia of Machine Learning*. Ed. by Claude Sammut and Geoffrey I. Webb. Boston, MA: Springer US, 2010, pp. 713–714. ISBN: 978-0-387-30164-8. DOI: [10.1007/978-0-387-30164-8\\_576](https://doi.org/10.1007/978-0-387-30164-8_576). URL: [https://doi.org/10.1007/978-0-387-30164-8\\_576](https://doi.org/10.1007/978-0-387-30164-8_576).
- [7] Ashraf M. Kibriya et al. “Multinomial Naive Bayes for Text Categorization Revisited”. In: *AI 2004: Advances in Artificial Intelligence*. Ed. by Geoffrey I. Webb and Xinghuo Yu. Berlin, Heidelberg: Springer Berlin Heidelberg, 2005, pp. 488–499. ISBN: 978-3-540-30549-1.
- [8] Xindong Wu et al. “Top 10 algorithms in data mining”. In: *Knowledge and information systems* 14.1 (2008), pp. 1–37.
- [9] Jaak Simm, Ildefons Magrans de Abril, and Masashi Sugiyama. *Tree-Based Ensemble Multi-Task Learning Method for Classification and Regression*. 6. The Institute of Electronics, Information and Communication Engineers, 2014, pp. 1677–1681. URL: <http://CRAN.R-project.org/package=extraTrees>.
- [10] David R Cox. “The regression analysis of binary sequences”. In: *Journal of the Royal Statistical Society: Series B (Methodological)* 20.2 (1958), pp. 215–232.
- [11] Warren S McCulloch and Walter Pitts. “A logical calculus of the ideas immanent in nervous activity”. In: *The bulletin of mathematical biophysics* 5.4 (1943), pp. 115–133.
- [12] Diederik Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *International Conference on Learning Representations* (Dec. 2014).
- [13] Agnes Lydia and Sagayaraj Francis. “Adagrad - An Optimizer for Stochastic Gradient Descent”. In: Volume 6 (May 2019), pp. 566–568.