

ΤΕΧΝΗΤΗ ΝΟΗΜΟΣΥΝΗ-Εργασία 2_η

Γενικά:

Το πρόγραμμα αποτελείται από δύο java αρχεία τα Test.java,Train.java

Main: Test.java

Arguments:

args[0]:path/to/folder/ru1

args[1]:πλήθος φακέλων που θα χρησιμοποιηθούν για training απ'τά 10.

Πχ C:\Users\User\Desktop\ru1 9

Train.java

ReadFile:

Έχει ορίσματα τον φάκελο ru1 και το πλήθος των testing φακέλων και Διαβάζει τα μηνύματα που περιέχονται στους training φακέλους και μετράει τον αριθμό εμφάνισης κάθε λέξης στα spams και στα hams αντιστοίχα.

Συγκεκριμένα:

Διαβάζουμε ένα ένα τα αρχεία μέσα στους φακέλους και για κάθε λέξη εάν είναι σε Spam `if(bool)` και δεν έχει εμφανιστεί ξανά στο ίδιο txt `!multOccuranceS.containsKey(words[j])` προσθέτουμε +1 στον αριθμό εμφανίσεων της στην spam. Αλλιώς αν είναι σε ham και εμφανίζεται πρώτη φορά στο txt προσθέτουμε +1 στην ham.

Props:

Έχει ορίσματα τον φάκελο ru1 και το πλήθος των testing φακέλων. Υπολογίζει την πιθανότητα κάθε λέξης με δεδομένο ότι το μήνυμα είναι ham και spam: $P(word_i | ham)$ και $P(word_i | spam)$

Εφαρμόζει εξομάλυνση laplace για αποφυγή μηδένισης γινομένου και log για αποφυγή underflow.

Αφού καλέσουμε την ReadFile για να αρχικοποιηθούν τα hashmaps **ham** και **spam** υπολογίζουμε για κάθε λέξη στην spam τον log της πιθανότητας $P(\text{wordi} | \text{spam})$

$:\log((\text{double})(\text{spam.get(i)}+1)/(\text{spam.size()}+\text{countWords}))$

Όπου

$\text{spam.get(i)}/\text{spam.size}():$

οι εμφανίσεις της λέξης στα spams / το σύνολο των διαφορετικών λέξεων στα spams

Ταυτόχρονα +1 στον αριθμητή ενώ +countWords το σύνολο των διαφορετικών λέξεων στον παρανομαστή για την εξομάλυνση Laplace.

Με τον ίδιο τρόπο υπολογίζουμε και τα $P(\text{wordi} | \text{ham})$.

Εάν η λέξη δεν υπάρχει καθόλου στα hams ορίζουμε την συχνότητα της λέξης 0 στον υπολογισμό της $\log(P(\text{wordi} | \text{ham}))$

$\text{if} (!\text{ham.containsKey(i)})\{$

$\text{WordHam} = \log((\text{double})(0+1)/(\text{ham.size()}+\text{countWords}));$

ενώ αντίστοιχα αν υπάρχει στην ham και όχι στην spam ορίζουμε την συχνότητα της λέξης 0 στον υπολογισμό της $\log(P(\text{wordi} | \text{spam}))$.

Έτσι τα WordHamProps και WordSamProps περιέχουν ως κλειδιά όλες τις λέξεις που εμφανίζονται στα μηνύματα και ως τιμή τις πιθανότητες που περιγράφηκαν παραπάνω για τα hams και τα props αντίστοιχα.

Test:

Η test είναι η main μας η οποία διαβάσει κάθε ένα testing φάκελο(10-training files) και για κάθε λέξη που βρίσκει αθροίζει τα log των πιθανοτήτων τους όπως υπολογίστηκαν στην Train ξεχωριστά για Spam και Ham :**totalS,totalH**.

Έπειτα στο άθροισμα totalS προσθέτει τον log του πλήθους των λέξεων των Spams και στο totalH των Hams ώστε να προκύψει ο αριθμητής για την δεσμευμένη πιθανότητα

$\log(P(\text{spam} | x_i)) = \log(P(\text{spam}) * P(x_i | \text{spam}) / P(x_i))$ και

$\log(P(\text{ham} | x_i)) = \log(P(\text{ham}) * P(x_i | \text{ham}) / P(x_i))$: **probS,probH**

Αφού συγκρίνει τα δύο αυτά \log **probS>probH** αυξάνει το πλήθος των **spams(spamsFound++)** ή **hams** που βρήκε το πρόγραμμα καθώς και το πλήθος των **spams** που ήταν πράγματι **spams (tp++)** ,και το πλήθος **ham** που ήταν πράγματι **ham (tn++)**.

Έπειτα υπολογίζει την τιμή της ανάκλησης ,της ακρίβειας,του F1 και του **accuracy:Recall,Precision,f1,accuracy**.

Πειράματα:

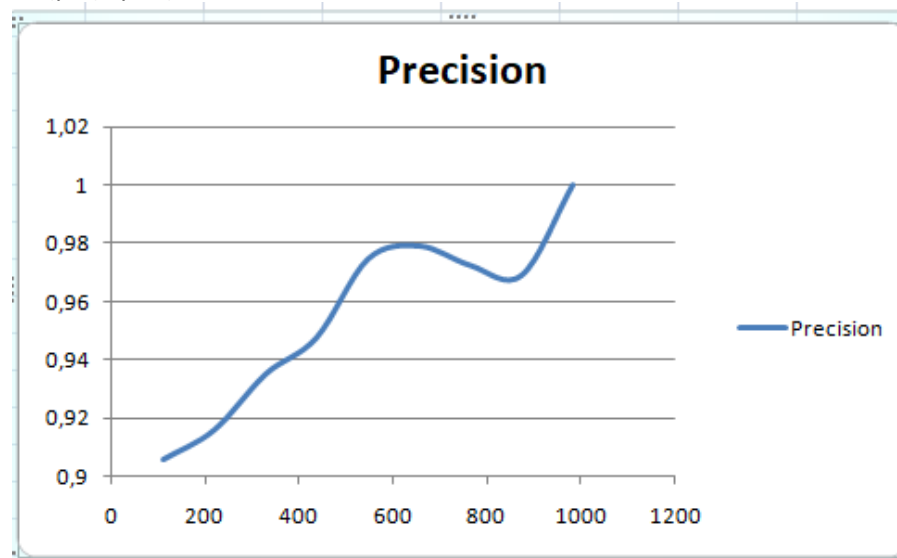
Τρέχοντας το πρόγραμμα για τον φάκελο PU1 και για διαφορετικό πλήθος δεδομένων εκπαίδευσης κάθε φορά έχουμε τον παρακάτω πίνακα .

Η μεταβλητή X αναφέρεται στο πλήθος των διαφορετικών λέξεων που χρησιμοποιήθηκαν για training ανάλογα με το πλήθος των parts τ ορίσαμε ως είσοδο.

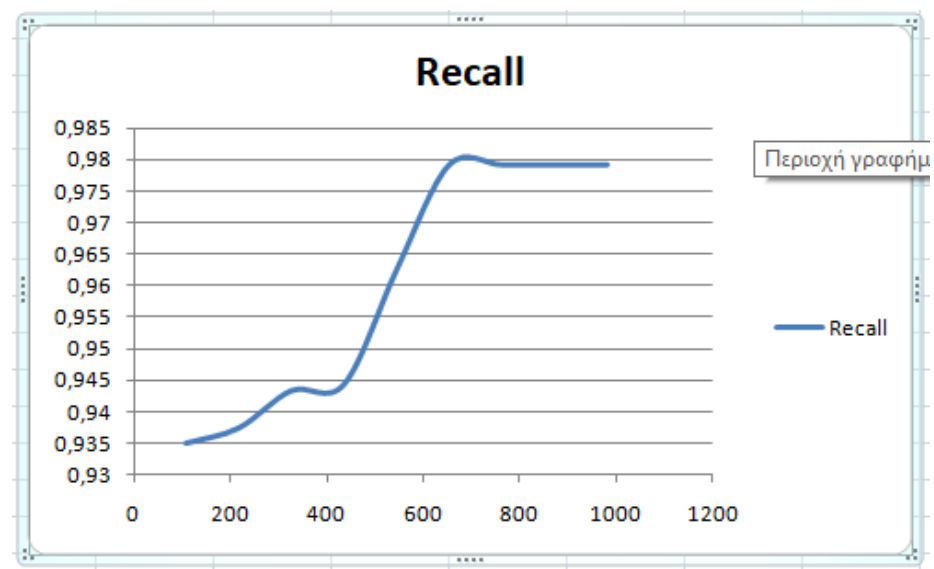
X	Precision	Recall	f1	Accuracy
109	0,9058	0,9351	0,9202	0,9286
218	0,916	0,9375	0,9266	0,9346
327	0,9351	0,9434	0,9392	0,9462
436	0,9477	0,9444	0,946	0,9525
545	0,9746	0,9625	0,9685	0,9724
654	0,9791	0,9791	0,9791	0,9816
763	0,9724	0,9791	0,9757	0,9785
872	0,969	0,9791	0,974	0,977
981	1	0,9791	0,9894	0,9908

Οι καμπύλες που προκύπτουν απ τις παραπάνω τιμές είναι οι εξείς:

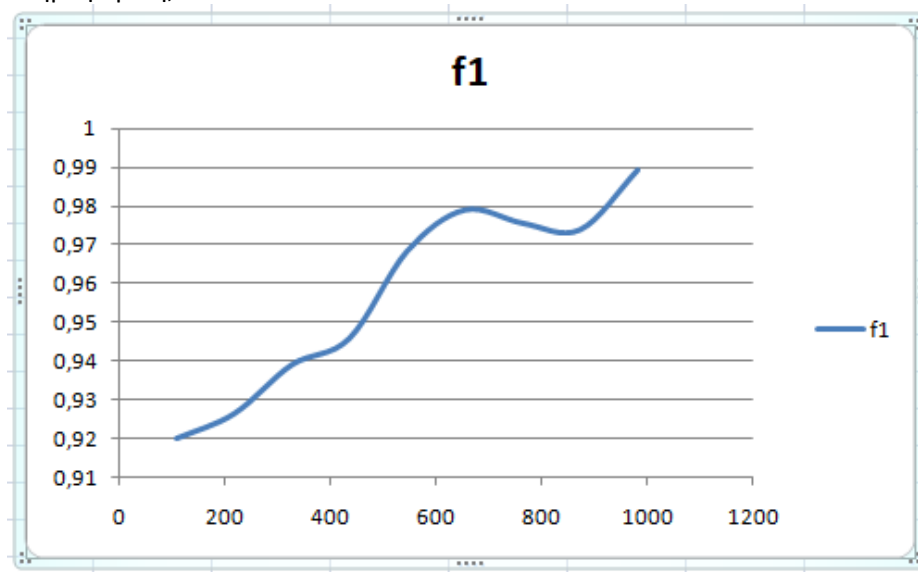
Precision:



Recall:



F1:



Accuracy:

