

Αποστόλου Ιωάννης 3190013
Καλαντζής Ηλίας 3190068
Κωνσταντίνος Κατσάμης 31900237

Αφελής ταξινομητής Bayes: Bernoulli

Ο Naive Bayes Bernoulli εκπαιδεύεται υπολογίζοντας την πιθανοτήτων κάθε λέξης του λεξιλογίου που το έχει δοθεί στο στάδιο της εκπαίδευσης σε συνδυασμό με τα δεδομένα εκπαίδευσης. Έτσι όταν του δοθούν νέα δεδομένα υπολογίζει το γινόμενο των πιθανοτήτων κάθε λέξεις σύμφωνα με τις κριτικές που είχε συναντήσει στο παρελθόν και σύμφωνα αυτών των πιθανοτήτων καταλήγει στο αν το νέο σχόλιο που του δόθηκε είναι θετικό ή αρνητικό.

Σταθερές και πίνακες που χρησιμοποιήθηκαν:

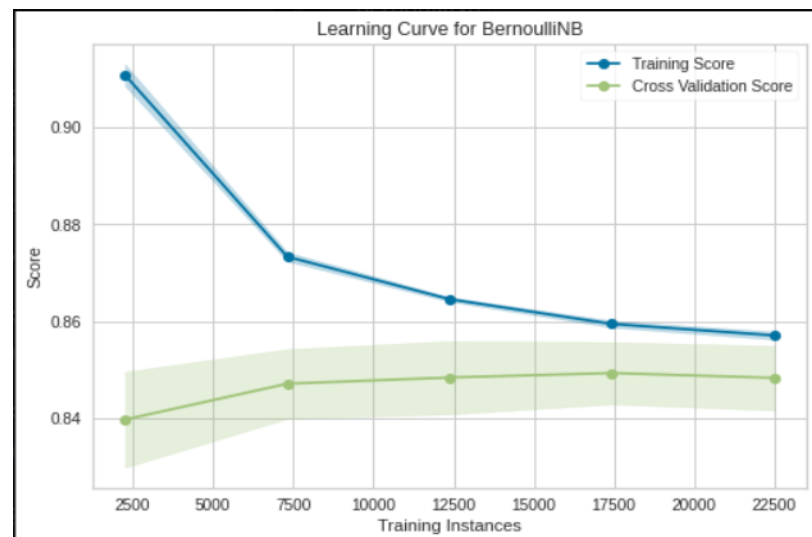
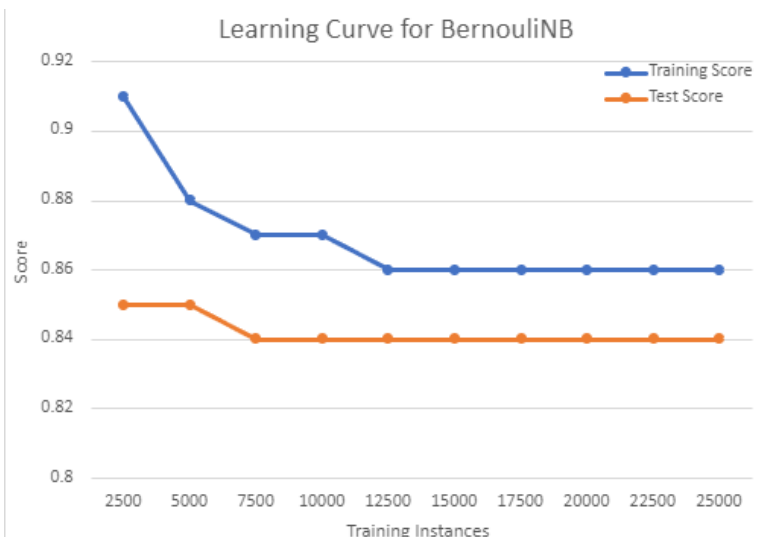
Κατα την υλοποίηση μας έπειτα από δοκιμές χρησιμοποιούμε λεξικό $M = 4000$ πιο συχνών λέξεων και δεδομένα ελέγχου, εκπαίδευσης μεγέθους 25000. Αυτό έχει σαν αποτέλεσμα να έχουμε 4 πίνακες τους:

X_{train_binary} , X_{test_binary} (μεγεθος 25000×4000) οι οποίοι περιέχουν τις κριτικές σε μορφή διανύσματος ιδιοτήτων ενώ Y_{train} , Y_{test} (μέγεθος 1×25000) είναι οι πίνακες που περιέχουν τις αποτιμήσεις κάθε κριτικής που απεικονίζουν με 1 τις θετικές και 0 της αρνητικές.

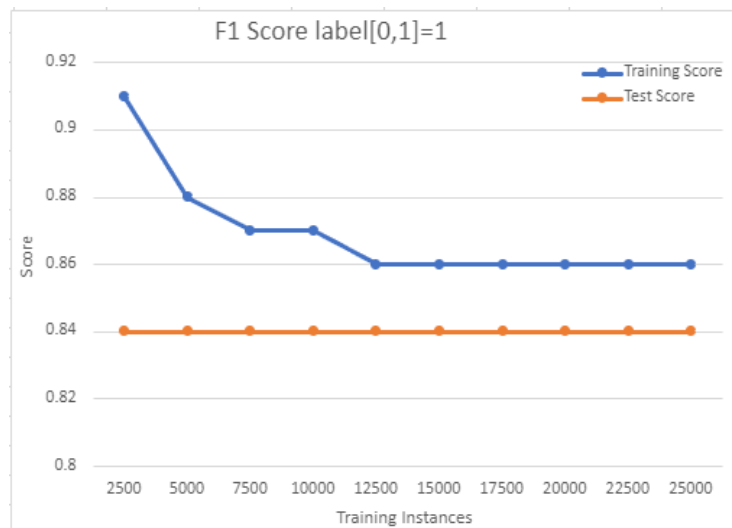
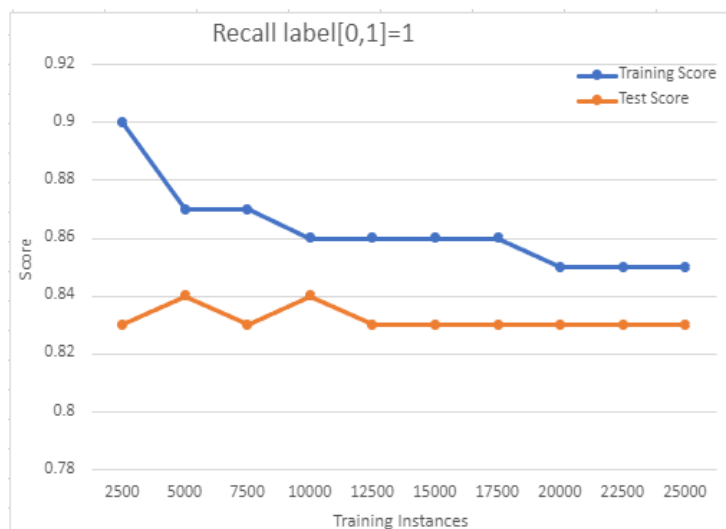
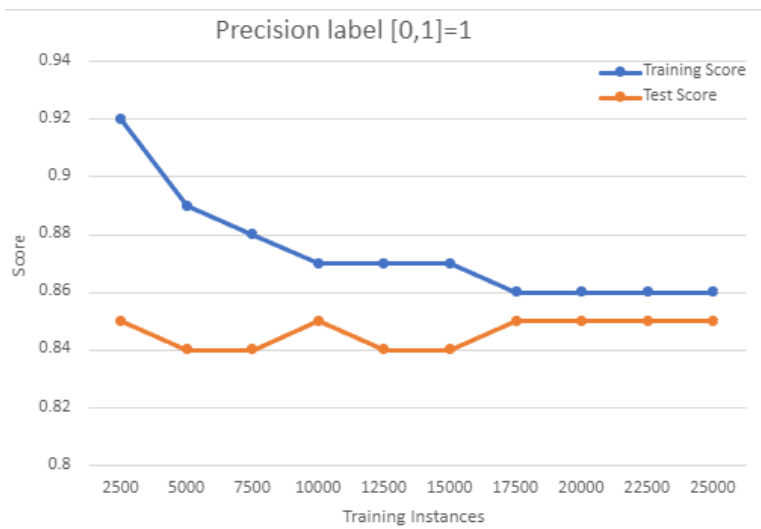
Τα δεδομένα φορτώνονται και μετατρέπονται στην σωστή μορφή από της `load_words`. Η εκπαίδευση γίνεται από την `fit` και η αποτίμηση από την `predict`.

Αποτελέσματα:

Η ακρίβεια μας “σταθεροποιήτε” στα 12500 στοιχεία όπου από εκεί και πέρα παραμένει σταθερή στο 86% ακρίβεια σε νέα δεδομένα και 84% σε δεδομένα εκπαίδευσης. Στα διαγράμματα από κάτω δεξιά έχουμε την δική μας υλοποίηση ενώ αριστερά την υλοποίηση της βιβλιοθήκης `sklearn`. Τα learning curve και των 2 περιπτώσεων είναι πολύ κοντά με την έτοιμη βιβλιοθήκη κοντά στα 12500 στοιχεία να έχει 86.5% ακρίβεια σε νέα δεδομένα και 85% σε στα δεδομένα όπου εκπαιδεύτηκε



Καμπύλες precision, recall, F1 στην περίπτωση μας τα δεδομένα ελέγχου και εκπαίδευσης μας αποτελούνταν από 25000 στοιχεία στο κάθε ένα εκ των οποίων τα 12500 ήταν θετικά και τα υπόλοιπα 12500 ήταν αρνητικά. Συνεπώς τα δείγματα μας ήταν χωρισμένα 50-50 με αποτέλεσμα και οι 3 καμπύλες να έχουν σχεδόν τα ίδια αποτελέσματα με την precision στο 86% και 85% για δεδομένα εκπαίδευσης και ελέγχου αντίστοιχα. Την recall με 85% και 83% για δεδομένα εκπαίδευσης και ελέγχου και την F1 όπου είχε 86% και 85%

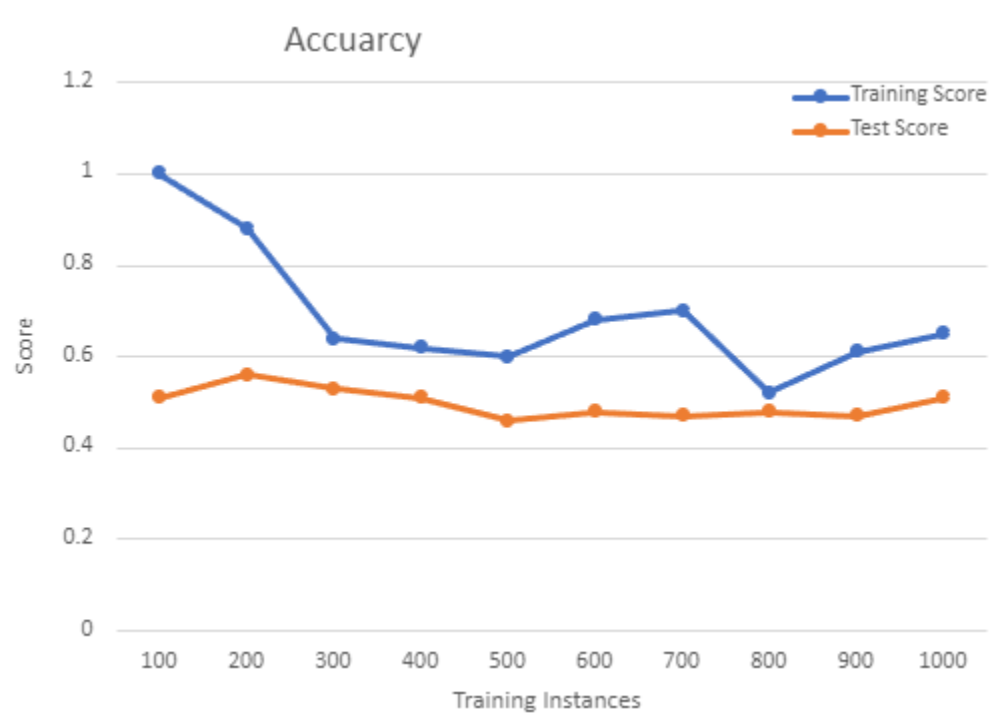


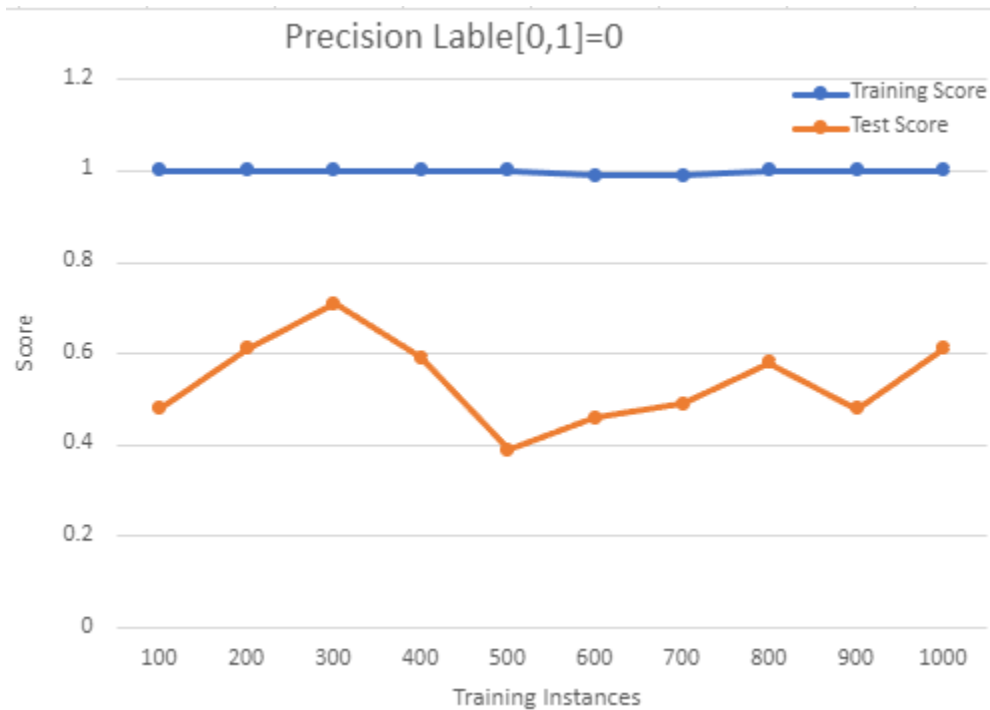
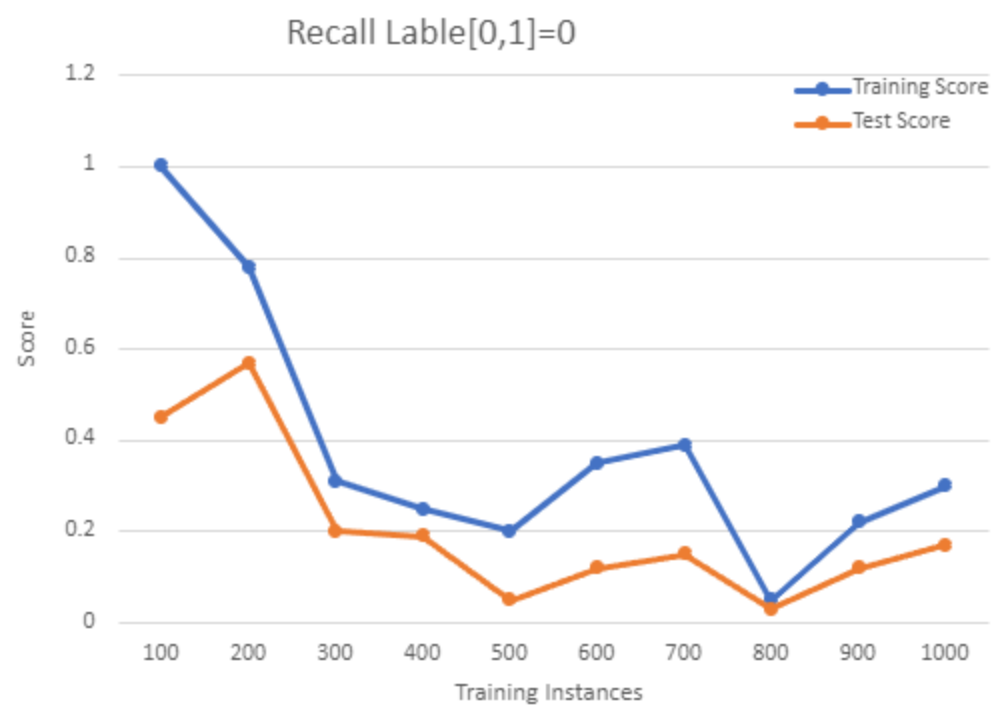
ID3

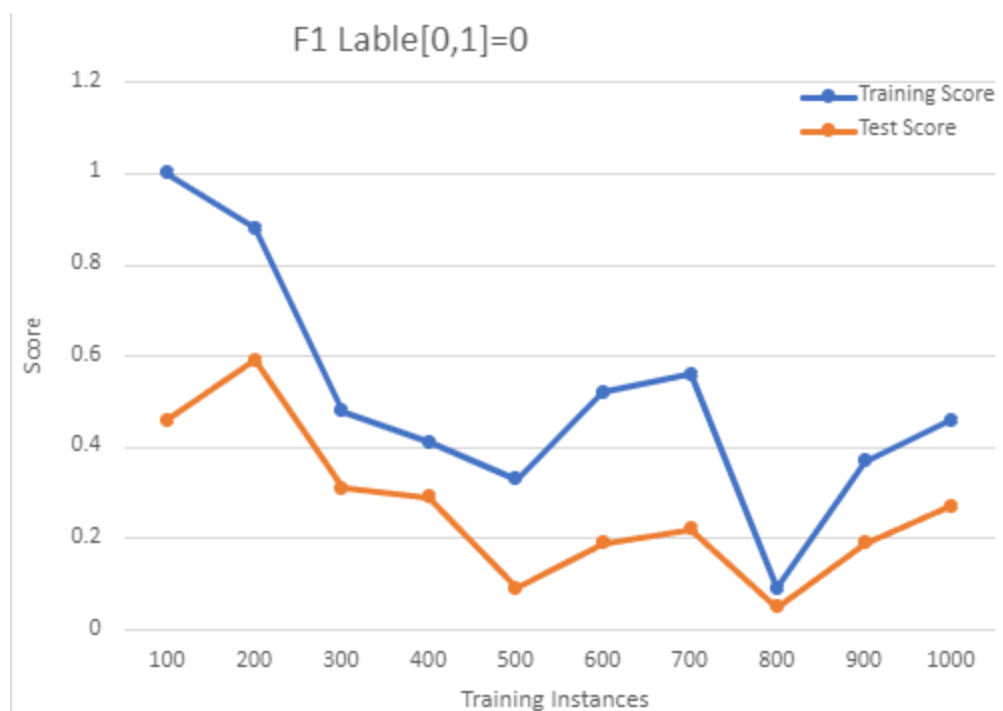
Ο ID3 είναι ένας αλγόριθμος στον οποίο δημιουργούμε ένα δέντρο αποφάσεων. Για να έχουμε σωστές αποφάσεις θα πρέπει να υπάρξει μια σωστή αναλογία μεταξύ βάθους του δέντρου, καθώς οι υπολογιστικοί πόροι δεν είναι απεριόριστοι.

Μεγάλο Βάθος & λίγα δεδομένα:

Όταν έχουμε μεγάλο βάθος και λίγα δεδομένα, καταφέρνει το δέντρο να βγάλει πολύ καλά αποτελέσματα στα δεδομένα που εκπαιδεύτηκε επειδή έχει μεγάλο βάθος, όμως σε καινούργια τα ποσοστά δεν είναι εξίσου καλά. Παράλληλα παρατηρούμε πως υπάρχει overfitting.

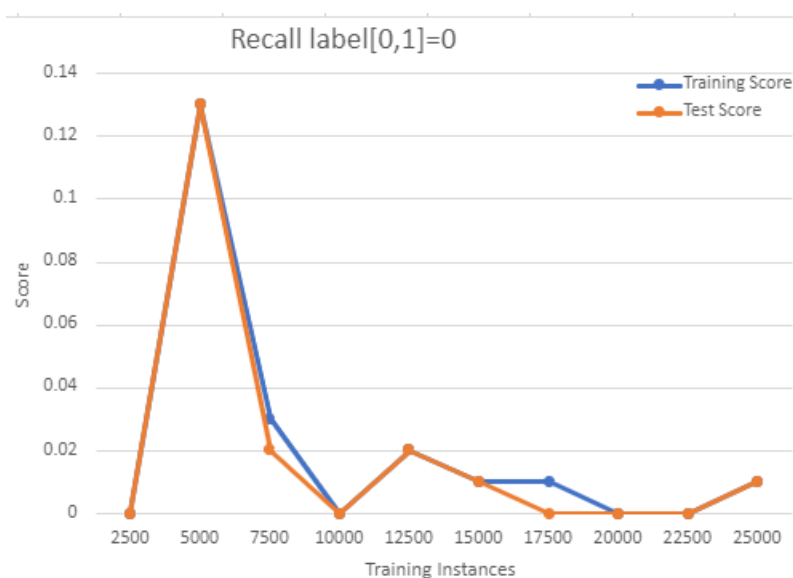


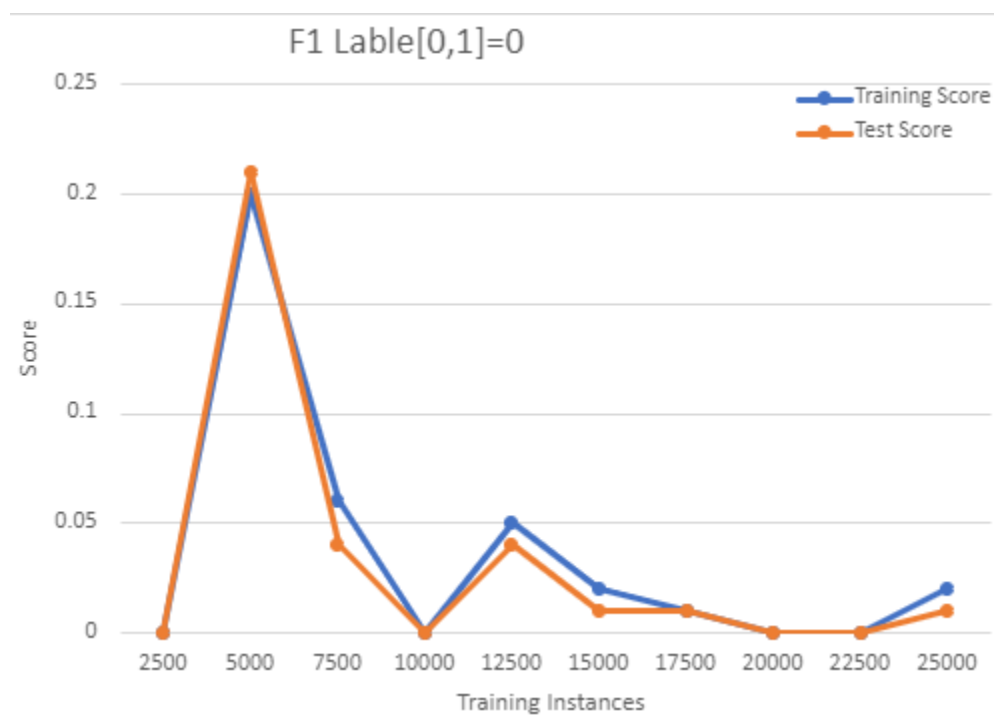
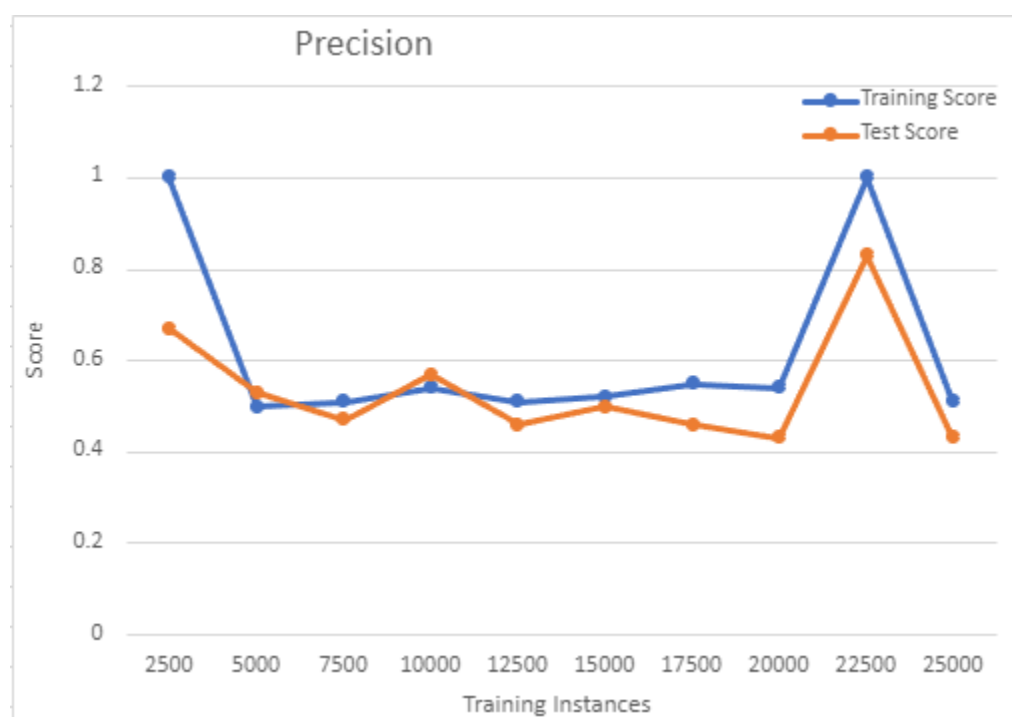




Μικρό Βάθος & πολλά δεδομένα:

Στα πολλά δεδομένα με μικρό βάθος παρατηρούμε πως το δέντρο βαραινει πάρα πολλή από την μια πλευρά του και τα περισσότερα δεδομένα καταλήγουν σε ένα γενικό κόμβο, όπου είναι τα ποσοστά ενός μεγάλου μέρους των δεδομένων και δεν μπορούμε να βγάλουμε σαφή συμπεράσματα. Για παράδειγμα στα δικά μας δεδομένα που είναι 50-50 οι πιθανότητες, καταλήγει ο μεγαλύτερος όγκος των δεδομένων στο τελευταίο φύλλο όπου έχει την πιθανότητα 0.5-0.5 και άρα δεν βγάζουμε σχεδόν καθόλου σωστά αποτελέσματα. Εδώ παρακάτω βλέπουμε τα δεδομένα για 25.000 στοιχεία (δλδ όλα τα στοιχεία) και το βάθος 5.





Random Forest

Ο Random Forest, είναι ένας αλγόριθμος που συνδυάζει ταξινομητές που κάνουν διαφορετικά λάθη. Συγκεκριμένα, δημιουργούμε διάφορες παραλλαγές του συνόλου εκπαίδευσης, που έχουν τον ίδιο αριθμό παραδειγμάτων με το αρχικό. Οι επιλογές των παραδειγμάτων γίνονται με επανατοποθέτηση. Χωρίζουμε επίσης το σύνολο ιδιοτήτων σε μικρότερα υποσύνολα, αλλά αυτή τη φορά, κρατάμε ένα μικρότερο αριθμό ιδιοτήτων (ίδιο για όλα τα υποσύνολα) για κάθε υποσύνολο.

Αφού ετοιμάσουμε τα σύνολα που θα χρησιμοποιήσουμε, καλούμε ανάλογα με τα δέντρα που θα φτιάξουμε τόσες φορές τον ταξινομητή, και κάθε αποτέλεσμα που προκύπτει το τοποθετούμε σε έναν πίνακα. Έπειτα, αφού τελειώσει η διαδικασία με τον ταξινομητή, παίρνουμε τις επικρατέστερες τιμές από κάθε πίνακα.

Σημαντικές παράμετροι που παίζουν ρόλο στον Random Forest, είναι το πλήθος δέντρων που θα υλοποιήσει, και ο βαθμός στον οποίο θα γίνουν οι τυχαίες επιλογές στα υποσύνολα έτσι ώστε να μην δημιουργούνται τα ίδια μονοπάτια στον ταξινομητή.

Στα παραδείγματα μας, λόγω χρόνου, κρατούσαμε το 30% των παραδειγμάτων και δημιουργούσαμε λίγα δέντρα.

Αποτελέσματα:

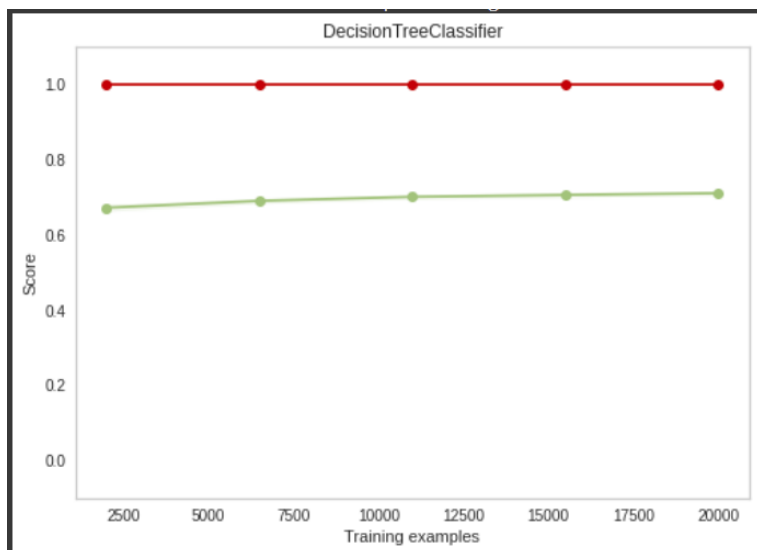
Χρησιμοποιώντας τον ID3 που υλοποιήσαμε, η καλύτερη ακρίβεια που πετύχαμε είναι 54% για 7500 στοιχεία στα δεδομένα εκπαίδευσης. Στα νέα δεδομένα, η καλύτερη ακρίβεια είναι 71% για 12000 στοιχεία.

Παρακάτω, φαίνονται οι καμπύλες ορθότητας, ακρίβειας και ανάκλησης.

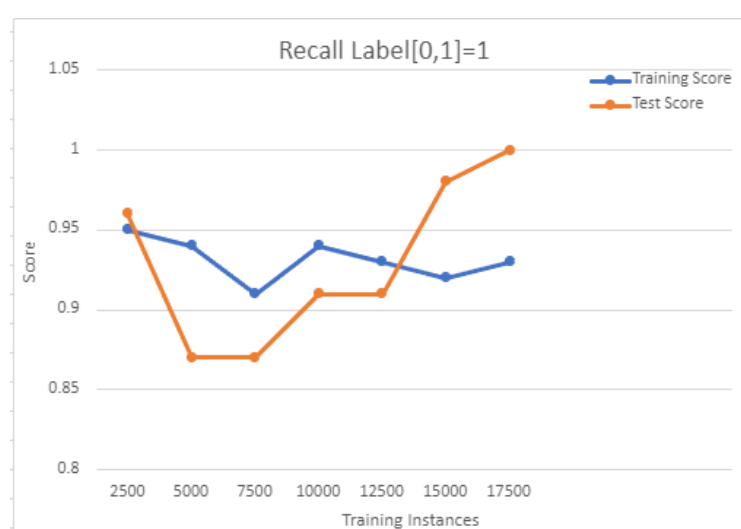
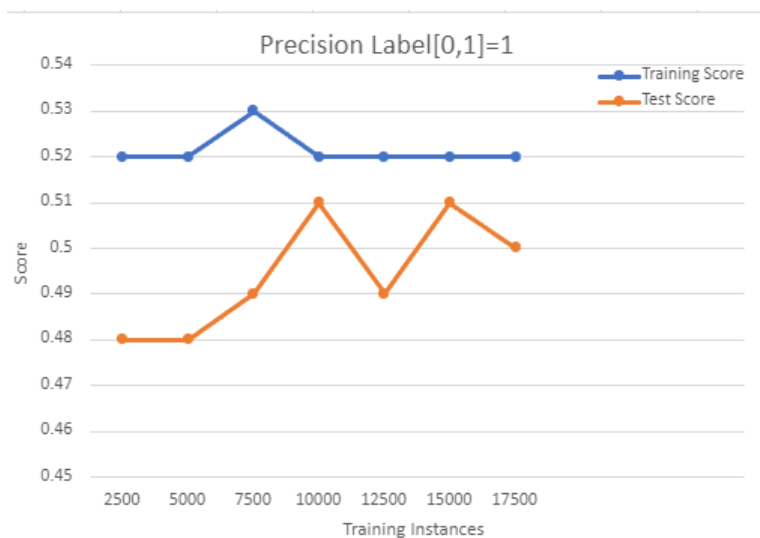
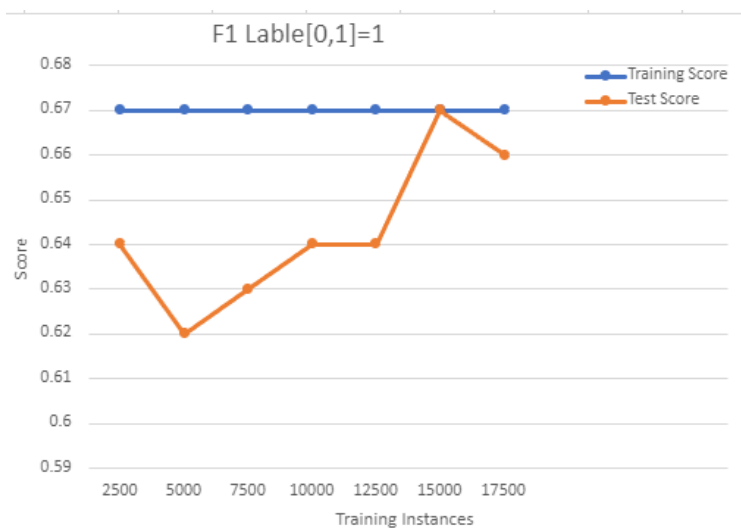
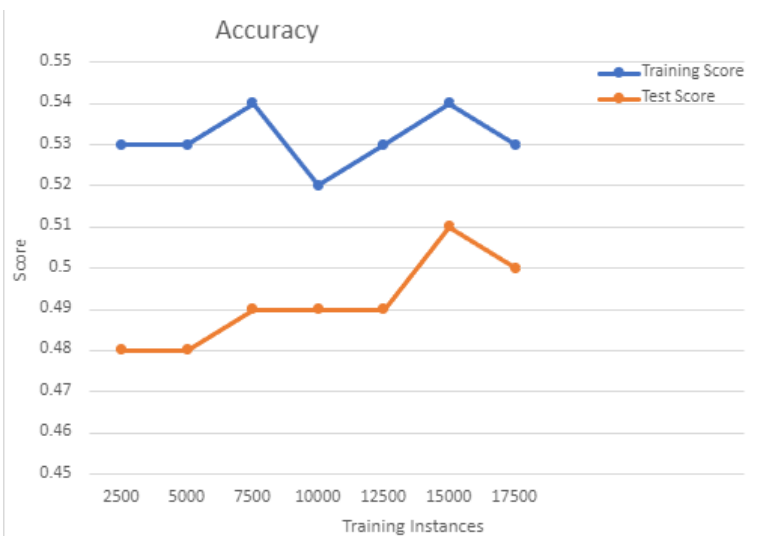
Στην πρώτη εικόνα, φαίνεται το score σε συναρτηση με τα παραδείγματα του Random Forest με τον έτοιμο ID3 της βιβλιοθήκης.

Χρησιμοποιώντας τον έτοιμο ID3, η καλύτερη ακρίβεια που πετύχαμε είναι 0.52% για 22500 στοιχεία.

Ακολουθεί η καμπύλη του έτοιμου ID3:



Ακολουθούν οι καμπύλες στα νέα δεδομένα, του Random Forest με τον δικό μας ID3:



Ακολουθούν οι καμπύλες στα νέα δεδομένα με τον έτοιμο ID3:



Παρατήρηση: Στην υλοποίηση του Random Forest με την έτοιμη ID3 (Decision Tree) παρατηρούμε αύξηση στην ακρίβεια 10%

Σημείωση: Για να γίνει import ο id3 μέσα στον κωδικά του Random Forest, θα πρέπει να γίνει στο colab κατά τη διάρκεια εκτέλεσης του προγράμματος. Δεν γίνεται αυτοματα. Γι' αυτόν τον λόγο αντιγράψαμε τον id3 μέσα στο αρχείο.