

ΠΑΝΕΠΙΣΤΗΜΙΟ ΜΑΚΕΔΟΝΙΑΣ

ΤΜΗΜΑ ΕΦΑΡΜΟΣΜΕΝΗΣ ΠΛΗΡΟΦΟΡΙΚΗΣ

Εργασία πάνω στην επιβλεπόμενη μάθηση και την ταξινόμηση

Κωνσταντίνος Πασβάντης

10 Νοεμβρίου 2022



Περιεχόμενα

1	Εισαγωγή	2
2	Μερος Α'	3
2.1	Logistic Regression	3
2.2	Decision Tree Classifier	4
2.3	k-Nearest Neighbors	4
2.4	Linear Discriminant Analysis	5
2.5	Gaussian Naive Bayes Classifier	6
2.6	Support Vector Machines	7
2.7	Neural Networks	8
2.8	Συμπεράσματα	9
3	Μερος Β'	10
3.1	Logistic Regression	10
3.2	Decision Tree Classifier	10
3.3	k-Nearest Neighbors	11
3.4	Linear Discriminant Analysis	12
3.5	Gaussian Naive Bayes	12
3.6	Support Vector Machines	13
3.7	Neural Networks	14
4	Τελικά Συμπεράσματα	14

Κατάλογος Σχημάτων

1	Μετρικές για το Logistic Regression	3
2	Μετρικές για τα δέντρα αποφάσεων	4
3	Μετρικές για τους κ πλησιέστερους γείτονες	5
4	Μετρικές για Linear Discriminant Analysis	6
5	Μετρικές για Gaussian Naive Bayes	7
6	Μετρικές για Support Vector Machines	7
7	Μετρικές για Neural Networks	9
8	Μετρικές για Logistic Regression pt.B	10
9	Μετρικές για τα δέντρα αποφάσεων pt.B	11
10	Μετρικές για τους κ πλησιέστερους γείτονες pt.B	11
11	Μετρικές για Linear Discriminant Analysis pt.B	12
12	Μετρικές για Gaussian Naive Bayes pt.B	13
13	Μετρικές για Support Vector Machines pt.B	13
14	Μετρικές για Neural Networks pt.B	14

Κατάλογος Πινάκων

1	Προβλέψεις για Linear Regression	3
2	Προβλέψεις για Decision Trees	4
3	Προβλέψεις για κ πλησιέστερους γείτονες	5
4	Προβλέψεις για Linear Discriminant Analysis	6
5	Προβλέψεις για Gaussian Naive Bayes	6
6	Προβλέψεις για Support Vector Machines	8
7	Προβλέψεις για Neural Networks	8
8	Προβλέψεις για Logistic Regression pt.B	10
9	Προβλέψεις για Decision Trees pt.B	11
10	Προβλέψεις για κ πλησιέστερους γείτονες pt.B	12
11	Προβλέψεις για Linear Discriminant Analysis pt.B	12
12	Προβλέψεις για Gaussian Naive Bayes pt.B	12
13	Προβλέψεις για Support Vector Machines pt.B	13
14	Προβλέψεις για Neural Networks pt.B	14

1 Εισαγωγή

Δώθεντός ένα αρχείο με δεδομένα από διάφορες εταιρείες ανά τα έτη, καλούμαι να αντιμετωπίσω ένα απλό, αλλά ταυτόχρονα απαιτητικό πρόβλημα: την δημιουργία μοντέλων ταξινόμησης με την βοήθεια των διαφόρων τεχνικών της μηχανικής μάθησης.

Τα δεδομένα περιγράφουν τις διαφορετικές καταστάσεις στις οποίες βρισκόταν η κάθε εταιρεία και στο τέλος κάθε παρατήρησης, δίνεται και η τελική κατάσταση της, δηλαδή, αν τελικά είναι υγιής ή αν έχει κυρήξει χρεωκοπία.

Με αυτά τα δεδομένα λοιπόν, πρέπει να συνταχθεί κώδικας ο οποίος θα διαβάζει τα δεδομένα από το αρχείο, με τα οποία θα εκπαιδεύονται τα διαφορετικά μοντέλα μηχανικής μάθησης, και θα προβλέπει τις εταιρείες εκείνες οι οποίες θα κυρήξουν πτώχευση.

Οι απαιτήσεις του προβλήματος είναι αυτές που το καθιστούν απαιτητικό: τα μοντέλα ταξινόμησης πρέπει να προβλέπουν με ποσοστό επιτυχίας 62% τις εταιρείες οι οποίες θα πτωχεύσουν και με 70% τις εταιρείες οι οποίες είναι υγιείς.

Στην εργασία αυτή θα παρουσιαστούν τα μοντέλα επιβλεπόμενης μάθησης που χρησιμοποιήθηκαν με σκοπό την λύση του παραπάνω προβλήματος και θα συγκριθούν μεταξύ τους ως προς το ποιο μοντέλο έδωσε τα καλύτερα αποτελέσματα.

Η εργασία είναι χωρισμένη σε δύο μέρη. Στο πρώτο μέρος τα μοντέλα ταξινόμησης εκπαιδεύονται χωρίς να γίνει κάποια τροποποίηση ως προς τον αριθμό των δειγμάτων, ενώ στο δεύτερο μέρος, το *train set* δημιουργείται με τέτοιο τρόπο ώστε η αναλογία υγιών προς μη υγιών εταιρειών να είναι 3:1.

Πριν να ξεκινήσει το κύριο μέρος της εργασίας στο οποίο παραθέτονται οι αλγόριθμοι, οι μόνες τροποποιήσεις που κάνουμε στα δεδομένα για τώρα θα είναι να αφαιρεθεί η στήλη με την χρονία της κάθε παρατήρησης, αφού δεν μας παρέχει κάποια συγκεκριμένη πληροφορία με την οποία θα μπορούσαν να εκπαιδευτούν τα μοντέλα, και να γίνει *scale* ώστε οι τιμές των μεταβλητών να κυμαίνονται από 0 έως 1.

Επίσης θεωρείται δεδομένο ότι ο σκοπός του κάθε μοντέλου είναι να κατατάξει τις εξαρτημένες μεταβλητές του προβλήματος, στην συγκεκριμένη περίπτωση το αν μία εταιρεία χρεωκοπήσει ή όχι, με την βοήθεια των ανεξάρτητων μεταβλητών, οι οποίες θα είναι οι μεταβλητές των *Xtrain*. Με τον όρο μεταβλητές θα εννοούμε κάθε στήλη που υπάρχουν στα δεδομένα μας, ενώ κάθε εταιρεία θα αποτελεί και μια εγγραφή.

Επίσης το σύνολο των δεδομένων που έχουμε χωρίσει το αρχικό *dataset* για να κάνουμε τις προβλέψεις αποτελείται από 56 μη υγιείς και 2088 υγιείς επιχειρήσεις. Οπότε αφού θέλουμε ποσοστό επιτυχίας 62% για τις μη υγιείς, πρέπει να προβλεφθούν σωστά τουλάχιστον 35 από τις 56, και με παρόμοιο σκεπτικό, τουλάχιστον 1462 από τις 2088 υγιείς. Οπότε για να βρεθεί το ποσοστό πρόβλεψης υγιών και μη υγιών επιχειρήσεων σε κάθε μοντέλο, αρκεί να διαιρέσουμε τον αριθμό *TP* με το 56 και τον αριθμό *TN* με το 2088 αντίστοιχα.

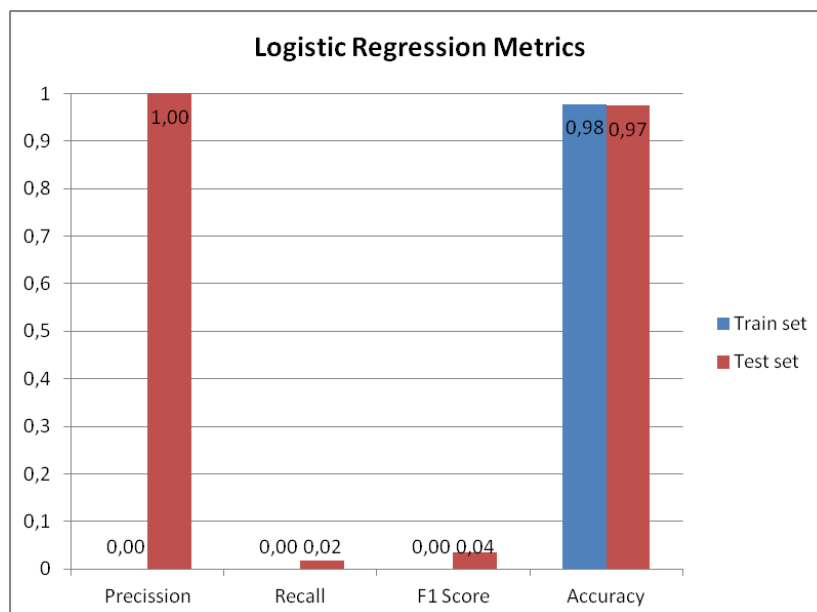
2 Μέρος Α΄

2.1 Logistic Regression

Το μοντέλο αυτό ψάχνει τις καλύτερες δυνατές παραμέτρους ώστε να δημιουργήσει μια σιγμοειδή συνάρτηση της οποίας το σύνολο τιμών ανήκει στο $(0, 1)$. Εφόσον τα πιθανά αποτελέσματα που ζητάμε είναι δυο, είτε να πτωχεύσει είτε να μην πτωχεύσει μια εταιρεία, αν συμβολίσουμε με 1 την κατάσταση στην οποία μια εταιρεία πτωχεύει και με 0 την κατάσταση στην οποία η εταιρεία είναι υγιής, η σιγμοειδής συνάρτηση που αναφέρθηκε νωρίτερα μπορεί να φανεί ιδιαίτερα αποτελεσματική.

Στο training phase, το μοντέλο, έχοντας σαν δεδομένα όλες τις ζητούμενες μεταβλητές μαζί με τα αποτελέσματά τους, αρχικά προσαρμόζει μία τυχαία συνάρτηση της προηγούμενης μορφής σε αυτά, και στην συνέχεια ελέγχει αν υπάρχει κάποια άλλη συνάρτηση η οποία εξυπηρετεί τον σκοπό αυτό καλύτερα. Τελικά, αφού δημιουργήσει την συνάρτηση η οποία ανταποκρίνεται πιο ικανοποιητικά στα δεδομένα, εξετάζει ποια είναι η πιθανότητα μια παρατήρηση με κάποιες τιμές μεταβλητών να είναι κλάσης 1 ή 0. Αν η πιθανότητα αυτή είναι πάνω από 50%, τότε η προβλεπόμενη κλάση για αυτή την παρατήρηση είναι η 1, ενώ αν είναι κάτω από 50%, είναι η 0.

Τα αποτελέσματα από αυτό το μοντέλο δεν κατάφεραν να ικανοποιήσουν κανέναν από τους δύο περιορισμούς που μας ζητήθηκαν, κυρίως λόγω της αδυναμίας του μοντέλου να προσδιορίσει ποιές εταιρείες θα χρεοκωπήσουν.



Σχήμα 1: Μετρικές για το Logistic Regression

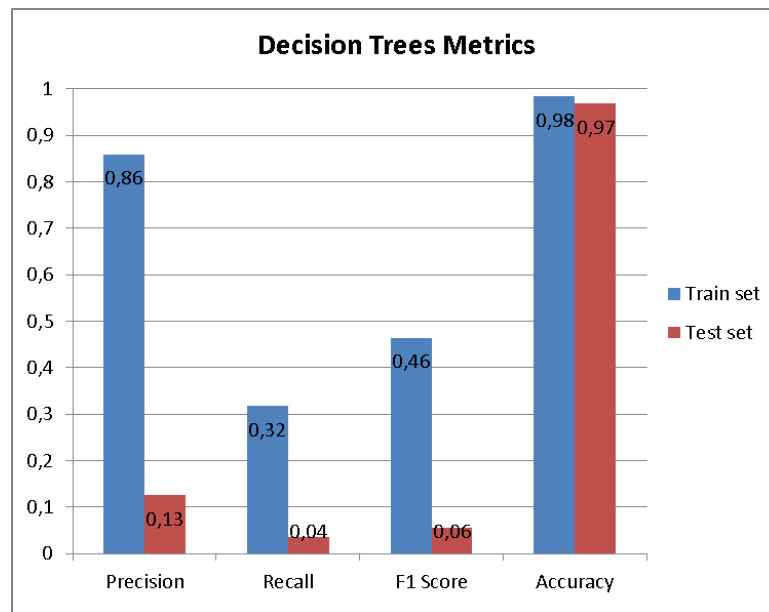
Πίνακας 1: Προβλέψεις για Linear Regression

	TP	TN	FP	FN
Train	0	8380	0	192
Test	1	2088	0	55

Όπως μπορούμε να δούμε η τιμή για το accuracy είναι πολύ ψηλά, κυρίως λόγω του μεγάλου αριθμού υγιών επειρηρήσεων που προβλέπει το μοντέλο. Εφόσον επίσης βγήκε σαν αποτέλεσμα πρόβλεψης να χρεοκοπήσει μόνο μία εταιρεία, η οποία τελικά ήταν όντως χρεοκωπημένη, εξηγεί και το γεγονός ότι το precision στο test set είναι ίσο με την μονάδα.

2.2 Decision Tree Classifier

Σκοπός των δέντρων αποφάσεων είναι να κατηγοριοποιήσουν τα δεδομένα κάνοντας προβλέψεις, οι οποίες είναι βασισμένες πάνω σε απαντήσεις ερωτήσεων που αφορούσαν τα δεδομένα τα οποία εκπαιδεύτηκαν. Πιο συγκεκριμένα, τα δέντρα αποφάσεων δημιουργούν μία ρίζα από την οποία προσπαθούν να χωρίσουν τα δεδομένα με τον καλύτερο δυνατό τρόπο. Ο τρόπος που το κάνουν αυτό είναι να δημιουργήσουν λογικές ερωτήσεις ώστε να καταλήξουν σε μία απόφαση. Από την ρίζα δημιουργούνται κόμβοι, οι οποίοι με την σειρά τους δημιουργούν επιπλέον κόμβους που απεικονίζουν το τί απόφαση πρέπει να παρθεί για να συνεχίσουν σε κατώτερο επίπεδο. Στο κατώτατο επίπεδο βρίσκονται τα φύλλα κάθε δέντρου, τα οποία αντιπροσωπεύουν τις πιθανές κατηγορίες αποτελεσμάτων (στην συγκεκριμένη περίπτωση χρεωκοπημένη και μη χρεωκοπημένη).



Σχήμα 2: Μετρικές για τα δέντρα αποφάσεων

Όπως μπορούμε να δούμε, ο ταξινομητής αυτός δίνει καλύτερα αποτελέσματα σε σχέση με το Logistic Regression, αλλά πάλι δεν ικανοποιεί τις απαιτήσεις του προβλήματός μας. Επίσης ξανασχολιάζουμε, το γεγονός ότι το accuracy στο test set είναι πάλι ψηλά δεν σημαίνει ότι το μοντέλο βγάζει σωστά αποτελέσματα. Για να το καταλάβουμε καλύτερα αυτό, ο τύπος που υπολογίζει το accuracy είναι :

$$acc = \frac{TP + TN}{TP + TN + FP + FN}$$

. Όμως εδώ εφόσον το μοντέλο δεν προβλέπει αρκετά TP , TN και FP σε σχέση με τον αριθμό των TP (όπως φαίνεται και από τον παρακάτω πίνακα), λογικό είναι η τιμή του κλάσματος να είναι υψηλή.

Πίνακας 2: Προβλέψεις για Decision Trees

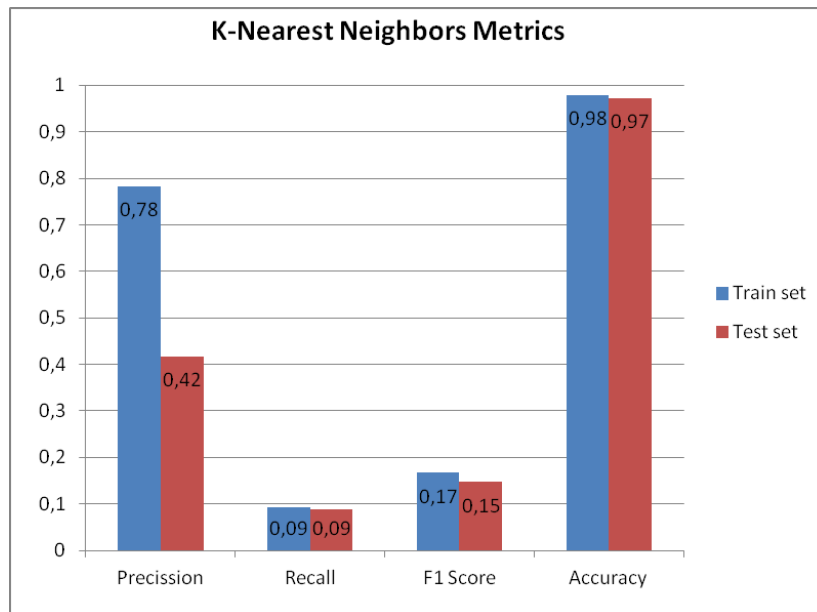
	TP	TN	FP	FN
Train	61	8370	10	131
Test	2	2074	14	54

2.3 k-Nearest Neighbors

Η βασική αρχή και το σκεπτικό αυτού του μοντέλου είναι ότι μία εγγραφή η οποία βρίσκεται κοντά σε κάποια άλλη με συγκεκριμένη κλάση, θα έχει την ίδια κλάση μαζί της.

Ο αλγόριθμος των k πλησιέστερων γειτόνων (στην περίπτωση μας 5 πλησιέστερων γειτόνων εφόσον δεν αλλάξαμε την παράμετρο k), για να προβλέψει την κλάση μίας εγγραφής λειτουργεί ως εξής: βρίσκει όλες τις αποστάσεις αυτής της εγγραφής από τις υπόλοιπες με τις οποίες έχουμε εκπαιδεύσει το μοντέλο. Συνήθως επιλέγεται η ευκλείδεια. Αφού υπολογίσει όλες τις αποστάσεις, βρίσκει τις 5 εγγραφές οι οποίες έχουν την μικρότερη απόσταση από την ζητούμενη. Η κλάση που εμφανίζεται τις περισσότερες φορές

από αυτές τις 5 εγγραφές, είναι ουσιαστικά και η κλάση την οποία προβλέπει το μοντέλο. Τέλος, τοποθετεί την εγγραφή την οποία προέβλεψε σε αυτή την κλάση και με παρόμοιο σχεπτικό συνεχίζει ο αλγόριθμος.



Σχήμα 3: Μετρικές για τους k πλησιέστερους γείτονες

Πίνακας 3: Προβλέψεις για k πλησιέστερους γείτονες

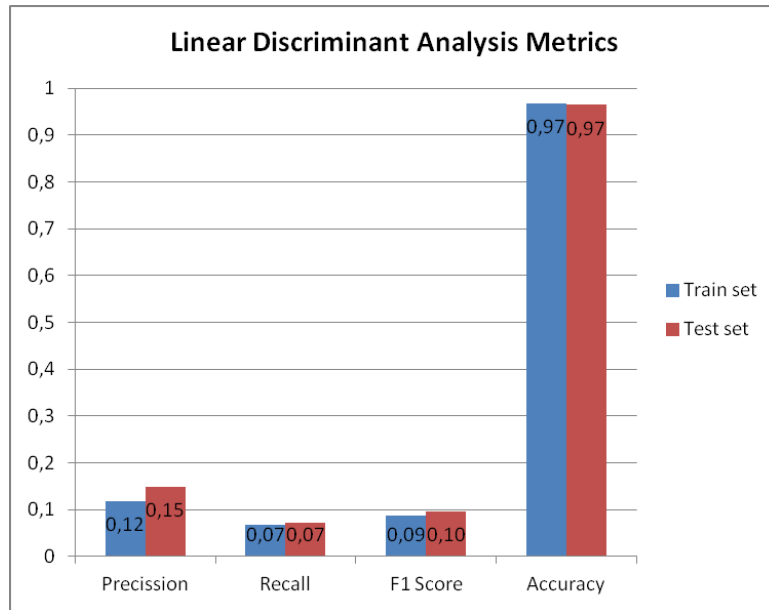
	TP	TN	FP	FN
Train	18	8375	5	174
Test	5	2081	7	51

Εδώ το μοντέλο παρόλο που μπόρεσε να βρει 5 από τις 56 μη υγιείς επιχειρήσεις στο *testset*, ενώ ο προηγούμενος αλγόριθμος μπόρεσε να βρει 7, βλέπουμε ότι τα στατιστικά των μετρικών είναι κατα πλειοψηφία καλύτερα όσον αφορά τον αλγόριθμο των Decision Trees. Σε αυτό βοήθησε ότι ο αριθμός των TN αυξήθηκε, ενώ αυτός των FP μειώθηκε σημαντικά.

2.4 Linear Discriminant Analysis

Η συνοπτική περιγραφή αυτής της μεθόδου μπορεί να αποδοθεί στον αγγλικό όρο "dimensionality reduction".

Ψάχνει μια διαχωριστική γραμμή που χωρίζει τα δεδομένα με τον καλύτερο δυνατό τρόπο, αφού τα προβάλλει πάνω σε αυτή τη γραμμή. Αυτό γίνεται με την ταυτόχρονη ικανοποίηση δύο συνθηκών: την μεγιστοποίηση απόστασης μεταξύ των μέσων κάθε κλάσης και την ελαχιστοποίηση της διασποράς που έχει η κάθε κλάση. Σε πολυδιάστατα προβλήματα όπως το δικό μας, ο όρος μέσος και διασπορά προκύπτουν από τις παραμέτρους της πολυδιάστατης κανονικής κατανομής. Εφόσον βρεθεί η διαχωριστική γραμμή, οι προβλέψεις γίνονται εκτιμώντας την πιθανότητα ένα καινούργιο σύνολο παρατηρήσεων να ανήκει σε κάποια κλάση. Αφού η διαχωριστική γραμμή δημιουργήθηκε με τέτοιο τρόπο ώστε να ικανοποιήσει τις δύο συνθήκες, το μοντέλο όταν βλέπει μια καινούργια εγγραφή, μπορεί να την τοποθετήσει στο σημείο της γραμμής στο οποίο έχει την μεγαλύτερη πιθανότητα να ανήκει, λόγω των τιμών των μεταβλητών της.



Σχήμα 4: Μετρικές για Linear Discriminant Analysis

Μπορούμε να δούμε ότι σε σύγκριση με τα προηγούμενα μοντέλα, δεν αποδίδει και τόσο καλά, εκτός βέβαια από το μοντέλο Logistic Regression. Επίσης παρατηρούμε και τον μεγάλο αριθμό *FP* και *FN* στο train set, πράγμα που δικαιολογεί γιατί το precision του train set είναι τόσο χαμηλό.

Πίνακας 4: Προβλέψεις για Linear Discriminant Analysis

	TP	TN	FP	FN
Train	13	8283	97	179
Test	4	2065	23	52

2.5 Gaussian Naive Bayes Classifier

Το μοντέλο στηρίζεται στο Θεώρημα του Bayes, έχοντας σαν βασικές αρχές ότι κάθε μεταβλητή του συνόλου δεδομένων είναι ανεξάρτητη και ισόνομη με κάθε άλλη μεταβλητή, και ότι όλες ακολουθούν κανονική κατανομή.

Για κάθε εγγραφή που θέλει να προβλέψει, υπολογίζει ποιά είναι η πιθανότητα να ανήκει στην κλάση 1 (μη υγιής) και ποιά να ανήκει στην κλάση 0 (υγιής) χρησιμοποιώντας τον τύπο του Bayes:

$$p(y | \mathbf{X}) = \frac{p(y)p(\mathbf{X} | y)}{p(\mathbf{X})}$$

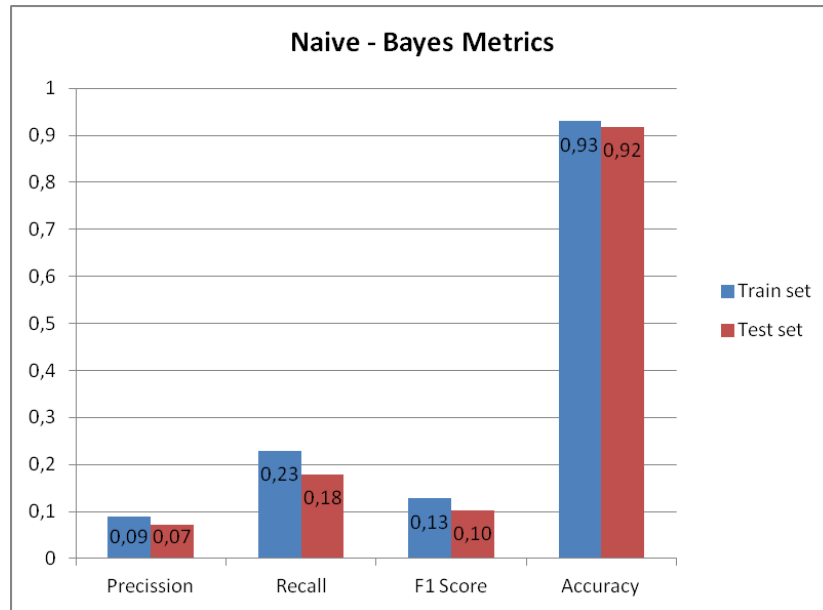
όπου εδώ $y = 0$ ή 1 , οι κλάσεις του προβλήματος και $\mathbf{X} = (X_1, X_2, \dots, X_{11})$, οι μεταλητές της κάθε παρατήρησης.

Εφόσον παίρνουμε σαν δεδομένο ότι οι μεταλητές είναι ανεξάρτητες και ισόνομες, μπορούμε να πούμε ότι $p(\mathbf{X}) = p(X_1)p(X_2)\dots p(X_{11})$ και $p(\mathbf{X} | y) = p(X_1 | y)p(X_2 | y)\dots p(X_{11} | y)$. Επίσης αφού κάθε μεταβλητή ακολουθεί κανονική κατανομή, οι πιθανότητες $p(X_1)\dots$ μπορούν να υπολογιστούν από την συνάρτηση πυκνότητας πιθανότητας της κατανομής.

Πίνακας 5: Προβλέψεις για Gaussian Naive Bayes

	TP	TN	FP	FN
Train	44	7928	452	148
Test	10	1958	130	46

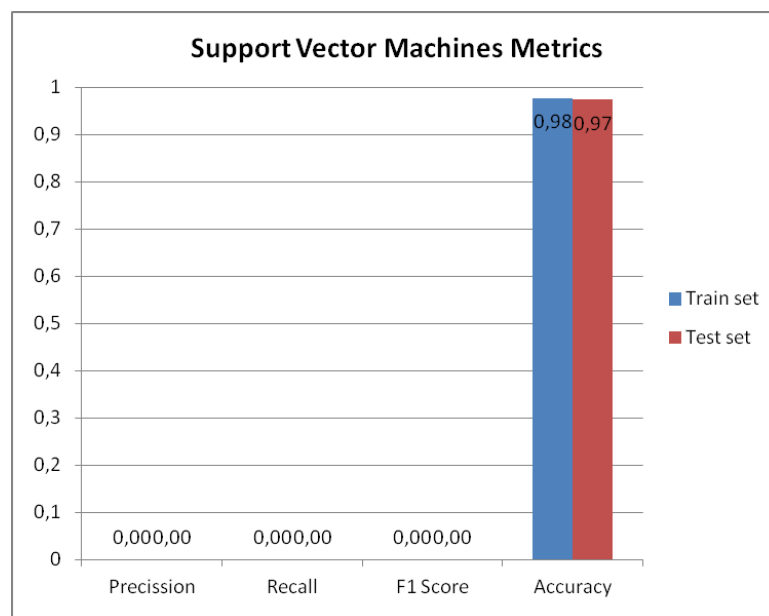
Το μοντέλο βλέπουμε ότι μπόρεσε να προβλέψει 10 από τις 56 μη υγιείς επιχειρήσεις, δηλαδή περισσότερες από κάθε άλλο μοντέλο, αλλά για ακόμα μία φορά το ποσοστό είναι πολύ πιο κάτω από το 62%.



Σχήμα 5: Μετρικές για Gaussian Naive Bayes

2.6 Support Vector Machines

Ο αλγόριθμος που χρησιμοποιούν τα support vector machines είναι παρόμοιος με αυτόν που χρησιμοποιεί το linear discriminant analysis . Δηλαδή, ψάχνει την καλύτερη γραμμή που χωρίζει τα δεδομένα, αυτή την φορά όμως δεν προβάλλει τις παρατηρήσεις πάνω στην ευθεία, αλλά δημιουργεί ένα υπερεπίπεδο. Για να βρει την καλύτερη δυνατή γραμμή, αναζητάει ποια από όλες τις πιθανές γραμμές που χωρίζουν τα δεδομένα δημιουργεί το μέγιστο υπερεπίπεδο. Δηλαδή, βρίσκει από κάθε κλάση το σημείο το οποίο είναι πιο κοντά στην γραμμή, το οποίο ονομάζεται διάνυσμα υποστήριξης. Σκοπός είναι να μεγιστοποιηθεί η απόσταση μεταξύ των διανυσμάτων υποστήριξης και της γραμμής που δημιουργείται, ώστε να δημιουργηθεί και το μέγιστο υπερεπίπεδο.



Σχήμα 6: Μετρικές για Support Vector Machines

Το μοντέλο αυτό βλέπουμε ότι δεν κατάφερε να προβλέψει καμία επειχόμενη που θα χρεοκοπήσει, είτε αυτό είναι αλήθεια είτε ψέμα.

Πίνακας 6: Προβλέψεις για Support Vector Machines

	TP	TN	FP	FN
Train	0	8380	0	192
Test	0	2088	0	56

2.7 Neural Networks

Ένα νευρωνικό δίκτυο αποτελείται από το 1ο επίπεδο (input layer), το οποίο αποτελείται από νευρώνες οι οποίοι συμβολίζουν την κάθε μεταβλητή που υπάρχει στα δεδομένα του προβλήματος, τα κρυφά επίπεδα (hidden layers) τα οποία αποτελούνται από βοηθητικούς νευρώνες ώστε να εκπαιδευτεί το μοντέλο, και το τελευταίο επίπεδο (output layer) το οποίο αποτελείται από τόσους νευρώνες όσες είναι και οι κλάσεις των δεδομένων.

Οι νευρώνες του κάθε επιπέδου συνδέονται με όλους τους νευρώνες του επόμενου με κάποιο βάρος. Στο 1ο επίπεδο υπάρχουν οι τιμές της κάθε μεταβλητής μιας εγγραφής, οι οποίες τροφοδοτούν το επόμενο επίπεδο. Η διαδικασία ξεκινάει χρησιμοποιώντας την 1η εγγραφή. Συμβολίζω με A_i την ποσότητα που παράγεται για το 2ο επίπεδο :

$$A_i = \sum_{j=1}^N b_j A_j$$

, όπου

- N ο αριθμός των μεταβλητών,
- b_j η τιμή βάρους που συνδέεται ο νευρώνας j του προηγούμενου επιπέδου με τον i νευρώνα,
- A_j η τιμή που υπάρχει στον j νευρώνα του προηγούμενου επιπέδου.

Εδώ να σημειωθεί ότι παίρνουμε την απλή περίπτωση όπου δεν υπάρχει *bias*.

Εφόσον υπολογισθεί για κάθε νευρώνα η τιμή A_i , τότε η τιμή κάθε νευρώνα του 2ου επιπέδου θα είναι ίση με $f(A_i)$, όπου f είναι η συνάρτηση ενεργοποίησης του συγκεκριμένου επιπέδου.

Με αυτό το σκεπτικό το νευρωνικό δίκτυο υπολογίζει τις τιμές των νευρώνων κάθε επιπέδου συναρτήσει των τιμών που έχουν οι νευρώνες στο προηγούμενο επίπεδο.

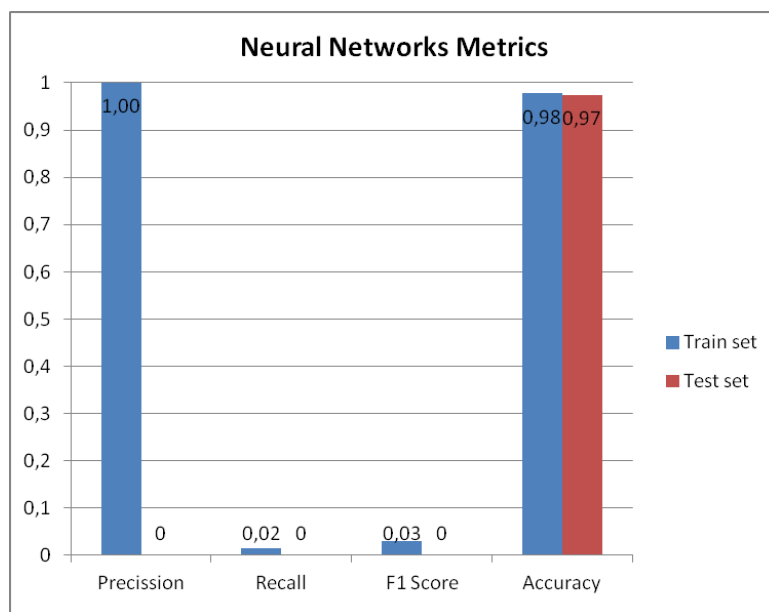
Στο τέλος η τιμή που θα έχουν οι νευρώνες του τελικού επιπέδου, θα καθορίσει και την κλάση που προβλέπει το νευρωνικό για την συγκεκριμένη εγγραφή. Αν η κλάση που προέβλεψε το μοντέλο για αυτή την εγγραφή είναι διαφορετική από την πραγματική, τότε με μια διαδικασία που ονομάζεται back propagation το νευρωνικό δίκτυο, αφού υπολογίσει το σφάλμα υπολογισμού, αλλάζει μερικά από τα βάρη με τέτοιο τρόπο ώστε να προκύψει καλύτερο αποτέλεσμα για τις επόμενες εγγραφές.

Αυτή η διαδικασία επαναλαμβάνεται για όλες τις εγγραφές στο train set ώστε να καταλήξει με το καλύτερο δυνατό μοντέλο πρόβλεψης νέων εγγραφών για κάποιο δείγμα που δεν έχει ξαναδεί.

Αξίζει να σημειωθεί ότι το μοντέλο δεν κατάφερε να προβλέψει καμία από τις 56 χρεοκοπημένες εταιρείες και ότι τα αποτελέσματα πρόβλεψης για το test set είναι ακριβώς τα ίδια με το support vector machines.

Πίνακας 7: Προβλέψεις για Neural Networks

	TP	TN	FP	FN
Train	3	8380	0	189
Test	0	2088	0	56



Σχήμα 7: Μετρικές για Neural Networks

2.8 Συμπεράσματα

Από ότι καταλάβαμε μέσα από τα διαγράμματα των μετρικών και τους πίνακες με τις προβλέψεις που έκανε το κάθε μοντέλο, κανένα δεν ήταν σε θέση ούτε να προβλέψει τις εταιρείες που θα χρεοκωπήσουν με ποσοστό πάνω από 18%.

Το καλύτερο μοντέλο όσον αφορά τουλάχιστον αυτά τα δεδομένα που είχαμε, αποδείχθηκε να είναι το Gaussian Naïve Bayes αφού προέβλεψε με ποσοστό 17% τις χρεοκωπημένες εταιρείες.

3 Μερos Β'

Η αδυναμία των μοντέλων στο 1ο μέρος της εργασίας οφειλόταν στο δείγμα, και συγκεκριμένα στην μεγάλη διαφορά υγιών και μη υγιών εταιρειών. Ο μεγάλος αριθμός υγιών εταιρειών έπαιξε καθοριστικό ρόλο στα αποτελέσματα κάθε ταξινόμητης.

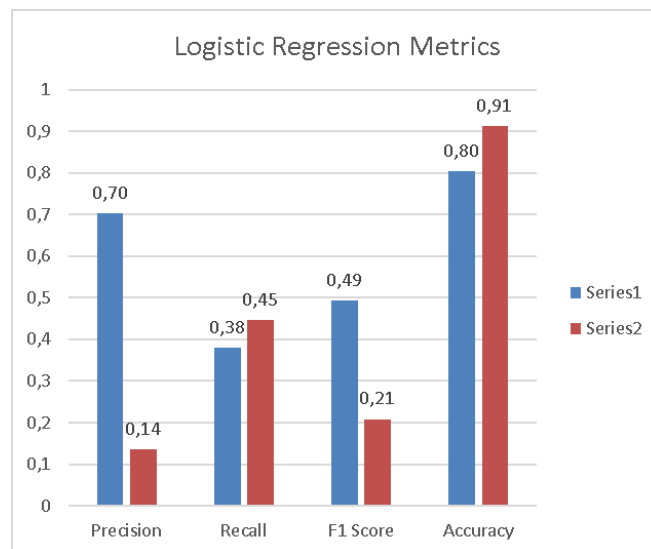
Όπως προαναφέρθηκε και στην εισαγωγή, στο Β' μέρος θα παρουσιαστούν ακριβώς τα ίδια μοντέλα, με την μοναδική διαφορά ότι τα δεδομένα που δίνουμε για εκπαίδευση αυτή τη φορά θα έχουμε αναλογία 3:1 μεταξύ υγιών και μη υγιών εταιρειών. Σε αυτό το κεφάλαιο δεν θα περιγραφούν ξανά το τι κάνει το κάθε μοντέλο, αλλά θα σχολιαστούν τα αποτελέσματα κάθε μοντέλου μετά από αυτή την μικρή τροποποίηση.

3.1 Logistic Regression

Βλέπουμε από τον πρώτο κιόλας αλγόριθμο την διαφορά στα αποτελέσματα. Από εκεί που το ποσοστό πρόβλεψης των μη υγιών εταιρειών ήταν κάτω από 2%, τώρα ανέβηκε στο 44%. Μπορεί και πάλι να μην ικανοποιήθηκε η απαίτηση του προβλήματος, αλλά τουλάχιστον είμαστε θετικά προδιαθετιμένοι για τα αποτελέσματα των υπολοίπων αλγόριθμων.

Πίνακας 8: Προβλέψεις για Logistic Regression pt.B

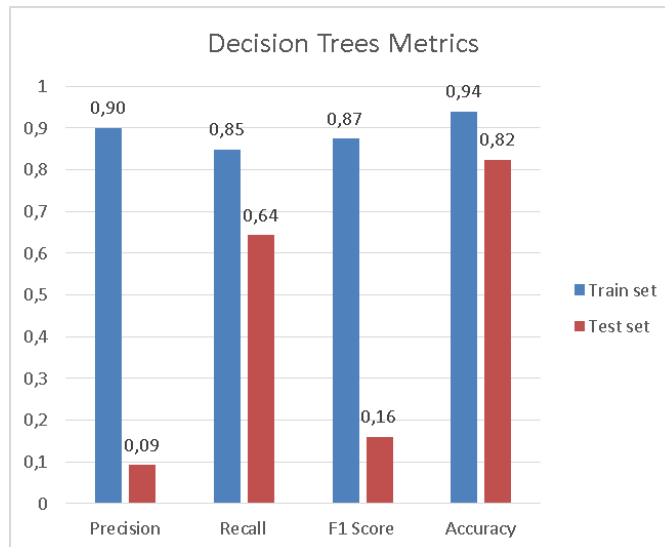
	TP	TN	FP	FN
Train	73	545	31	119
Test	25	1929	159	31



Σχήμα 8: Μετρικές για Logistic Regression pt.B

3.2 Decision Tree Classifier

Εδώ έχουμε την πρώτη φορά που ικανοποιούνται οι απαιτήσεις του προβλήματος. Βρέθηκαν 36 από τις 56 μη υγιείς επιχειρήσεις και 1731 από τις 2088 υγιείς. Βέβαια το *precision* εδώ πέρα είναι αρκετά μικρό, διότι το μοντέλο προέβλεψε συνολικά ότι 357 εταιρείες θα χρεοκοπήσουν, αυξάνοντας έτσι προφανώς τον αριθμό *FP*.



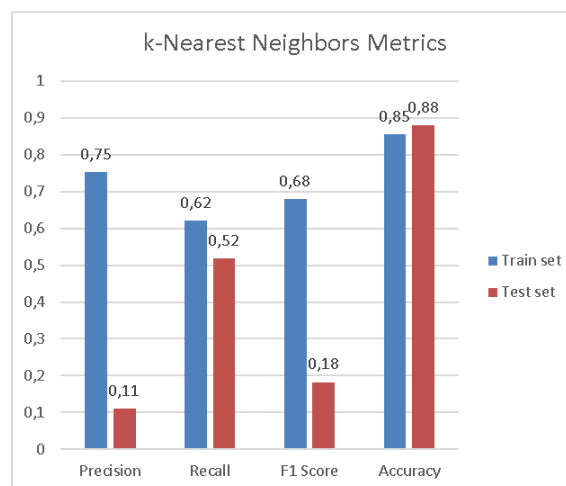
Σχήμα 9: Μετρικές για τα δέντρα αποφάσεων pt.B

Πίνακας 9: Προβλέψεις για Decision Trees pt.B

	TP	TN	FP	FN
Train	163	558	18	29
Test	36	1731	357	20

3.3 k-Nearest Neighbors

Με ποσοστό επιτυχίας 52% πάλι δεν ικανοποιούνται οι απαιτήσεις του προβλήματος. Ίσως με διαφορετικό αριθμό k , δηλαδή αν αλλάζαμε τον αριθμό που δέχεται το μοντέλο ώστε να υπολογίσει τις k πλησιέστερες αποστάσεις να παίρναμε καλύτερα αποτελέσματα, αλλά με τις default επιλογές ($k=5$), ο αλγόριθμος φαίνεται να αποτυγχάνει.

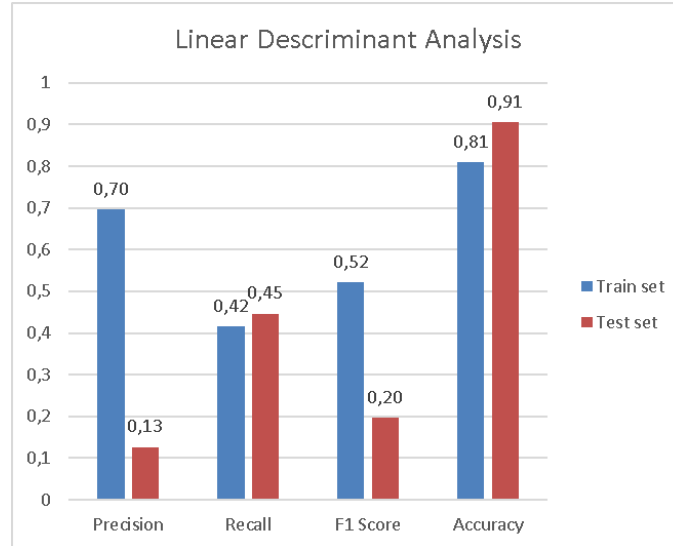


Σχήμα 10: Μετρικές για τους k πλησιέστερους γείτονες pt.B

Πίνακας 10: Προβλέψεις για x πλησιέστερους γείτονες pt.B

	TP	TN	FP	FN
Train	119	537	39	73
Test	29	1855	233	27

3.4 Linear Discriminant Analysis



Σχήμα 11: Μετρικές για Linear Discriminant Analysis pt.B

Το συγκεκριμένο μοντέλο βλέπω ότι προβλέπει λιγότερα TP σε σχέση με τα 3 προηγούμενα, αλλά βλέπω επίσης ότι προβλέπει πολύ περισσότερες εταιρείες οι οποίες είναι υγιείς. Δεν παίζει όμως μεγάλο ρόλο καθώς η κύρια απαίτηση του προβλήματος δεν ικανοποιείται.

Πίνακας 11: Προβλέψεις για Linear Discriminant Analysis pt.B

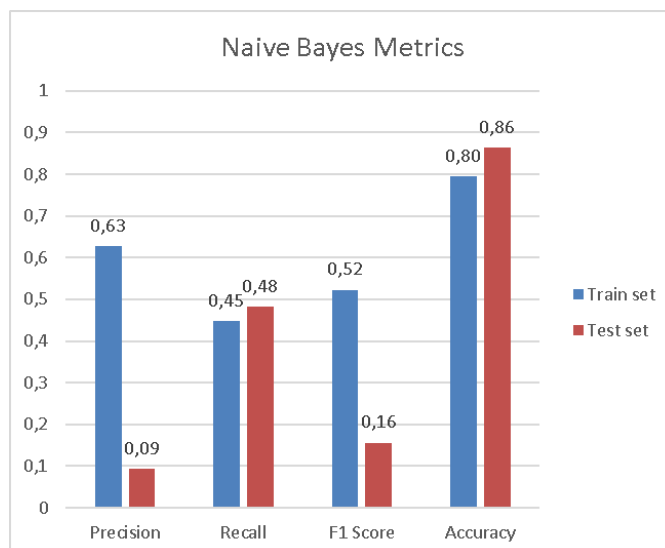
	TP	TN	FP	FN
Train	80	541	35	112
Test	25	1916	172	31

3.5 Gaussian Naive Bayes

Μπορεί η μετρική *recall* να βγάζει καλύτερο αποτέλεσμα συγκριτικά με το Linear Discriminant Analysis αλλά όλες οι υπόλοιπες μετρίες βγάζουν χειρότερα αποτελέσματα.

Πίνακας 12: Προβλέψεις για Gaussian Naive Bayes pt.B

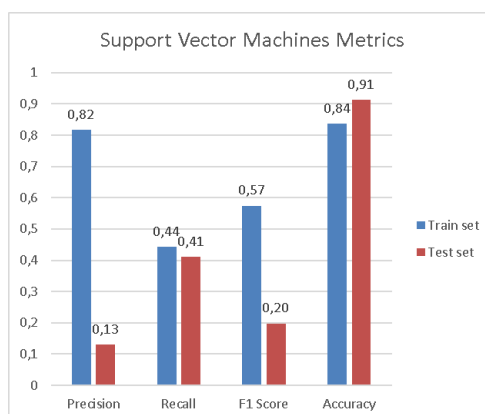
	TP	TN	FP	FN
Train	86	525	51	106
Test	27	1825	263	29



Σχήμα 12: Μετρικές για Gaussian Naive Bayes pt.B

3.6 Support Vector Machines

Εν,ω κανονικά θα περίμενε κανείς αυτό το μοντέλο να μας δώσει τα καλύτερα δυνατά αποτελέσματα, αφού δουλεύει καλύτερα με μικρότερο αριθμό δεδομένων, παρατηρούμε ότι κάτι τέτοιο δεν ισχύει. Μάλιστα, το ποσοστό επιτυχημένων προβλέψεων για τις μη υγιείς εταιρείες είναι ίσο με 41%, λιγότερο από κάθε μοντέλο στο Β' Μέρος της άσκησης έως τώρα.



Σχήμα 13: Μετρικές για Support Vector Machines pt.B

Πίνακας 13: Προβλέψεις για Support Vector Machines pt.B

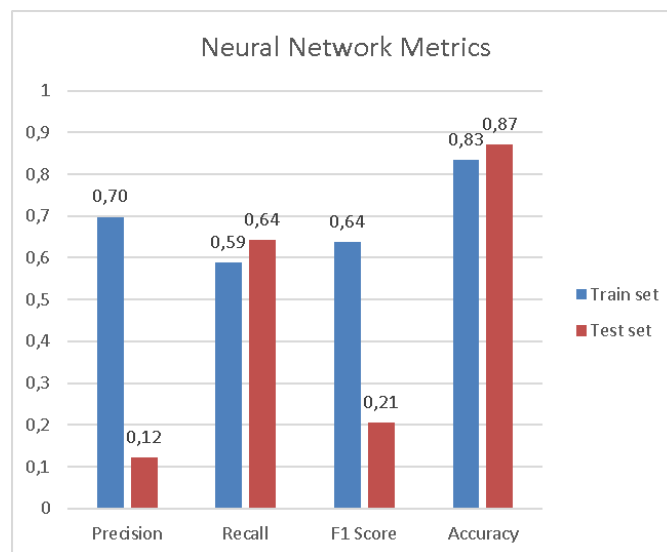
	TP	TN	FP	FN
Train	85	557	19	107
Test	23	1935	153	33

3.7 Neural Networks

Εδώ παρατηρούμε ότι έχουμε και δεύτερο ταξινομητή ο οποίος προβλέπει με ικανοποιητικό ποσοστό τις χρεωκοπημένες εταιρείες, και μάλιστα είναι καλύτερος από τον πρώτο αφού βλέπουμε ότι το ποσοστό ανέρχεται στο 64 %. Να σημειωθεί ότι δεν πειράξαμε σε μεγάλο βαθμό το νευρωνικό, δηλαδή δεν βάλουμε πολλά κρυφά επίπεδα, αντιθέτως με μόνο ένα κρυφό επίπεδο με 4 νευρώνες βλέπουμε ότι έχουμε επιθυμητό αποτέλεσμα. %.

Πίνακας 14: Προβλέψεις για Neural Networks pt.B

	TP	TN	FP	FN
Train	113	527	49	79
Test	36	1830	258	20



Σχήμα 14: Μετρικές για Neural Networks pt.B

4 Τελικά Συμπεράσματα

Παρά τις προσδοχίες που είχαμε για τα μοντέλα ταξινόμησης μετά την αποτυχία που είχαν όταν τροφοδοτήθηκαν με όλα τα δεδομένα εκπαίδευσης, ότι θα βρισκόταν σε θέση να προβλέψουν σε ικανοποιητικό ποσοστό την κλάση κάθε εταιρείας, βλέπουμε ότι κάτι τέτοιο δεν συμβαίνει. Μπορούμε να συνειδητοποιήσουμε λοιπόν ότι η αναλογία δύο κλάσεων στις οποίες ένας ταξινομητής εκπαιδεύεται και κάνει προβλέψεις παίζει σημαντικό ρόλο στα αποτελέσματά που παράγει.

Τα μόνα μοντέλα που κατάφεραν να βγάλουν ικανοποιητικά αποτελέσματα ήταν αυτά των δέντρων αποφάσεων και των νευρωνικών δικτύων. Όλα τα άλλα μοντέλα δεν ήταν σε θέση να προβλέψουν με ποσοστό πάνω από 62% τις χρεωκοπημένες εταιρείες, καθιστώντας τα άχρηστα για μελλονική χρήση. Το ποιο μοντέλο είναι καλύτερο από κάποιο άλλο μπορεί να βρεθεί κοιτώντας τα γραφήματα τα οποία παρουσιάστηκαν στο Β' Μερους της εργασίας.

Εδώ αξίζει να σημειωθεί ότι ενώ βρέθηκαν μόνο δύο μοντέλα που παρήγαγαν επαρκή αποτελέσματα, αυτό δεν σημαίνει ότι με διαφορετική επιλογή παραμέτρων θα εξακολουθούσαν τα υπόλοιπα να μην δίνουν σωστά αποτελέσματα. Για παράδειγμα τα support vector machines μπορεί να μας εξυπηρετούσαν αν διαλέγαμε διαφορετικό kernel . Γενικά εφόσον χρησιμοποιήσαμε τις default παραμέτρους κάθε μοντέλου δεν ήταν καθόλου σίγουρο ότι με αυτά τα δεδομένα θα έχουμε κάποιο σωστό αποτέλεσμα.

Συνοψίζοντας, τα καλύτερα μοντέλα για αυτό το πρόβλημα από ότι φαίνεται είναι αυτό των νευρωνικών δικτύων και των δέντρων αποφάσεων, τα οποία προβλέπουν με ποσοστό 64% τις εταιρείες.