



ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ
ΣΧΟΛΗ ΤΕΧΝΟΛΟΓΙΩΝ ΠΛΗΡΟΦΟΡΙΚΗΣ & ΕΠΙΚΟΙΝΩΝΙΩΝ
ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ



«ΑΝΑΛΥΤΙΚΗ ΔΕΔΟΜΕΝΩΝ, 6ο Εξάμηνο»

Ακαδημαϊκό Έτος 2018 – 2019

Διορία	28/06/2019
Ονοματεπώνυμο Φοιτητή (Αριθμός Μητρώου)	Σκάρος Γεώργιος (Π15128) Γκιολής Δημήτρης (Π16023) Κωστάκη Ειρήνη (Π16064)

ΣΚΑΡΟΣ ΓΕΩΡΓΙΟΣ (Π15128)

ΓΙΟΛΗΣ ΔΗΜΗΤΡΗΣ (Π16023)

ΚΩΣΤΑΚΗ ΕΙΡΗΝΗ (Π16064)



Πρόλογος

Για την υλοποίηση της εργασίας χρησιμοποιήθηκε η γλώσσα προγραμματισμού **R** και το πρόγραμμα **Rstudio**. Επίσης, η δημιουργία της βάσης δεδομένων του πρώτου ερωτήματος έγινε σε **postgresql**.

Η εκδόσεις που χρησιμοποιήθηκαν είναι οι εξής:

- R version 3.6.0
- PostgreSQL 9.6.6

Οι βιβλιοθήκες στην R που χρησιμοποιήθηκαν είναι οι εξής

- library(Metrics)
- library(class)
- library(caret)
- library(pROC)
- library(mlbench)
- library(factoextra)
- library(cluster)
- library(NbClust)
- library(arules)
- library(arulesViz)

Αν δεν είναι εγκατεστημένη οποιαδήποτε από τις παραπάνω βιβλιοθήκες αυτό μπορεί να γίνει με την εντολή `install.packages("")`

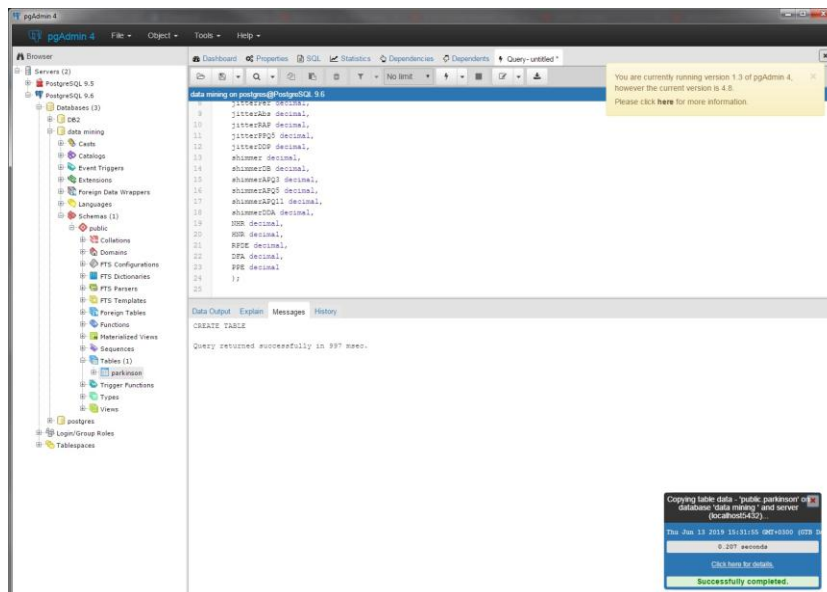
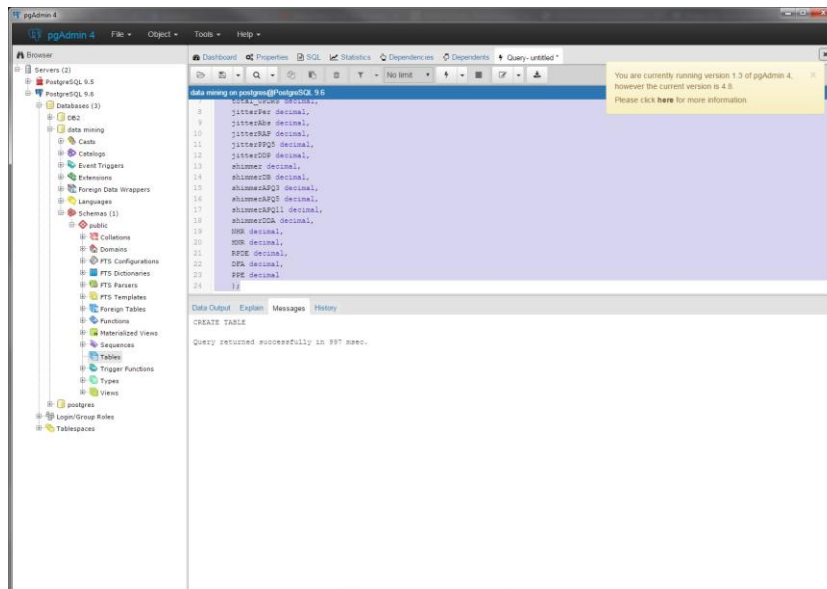
Η εργασία αποτελείται από 7 R scripts.



Εισαγωγή δεδομένων

Εισαγωγή στην PostgreSQL

Μετά την δημιουργία της κατάλληλης βάσης δεδομένων στο περιβάλλον μας (το CreateTable υπάρχει στα αρχεία) έγινε η εισαγωγή των δεδομένων στο περιβάλλον μας. Όλες οι στήλες της βάσης επιλέχθηκαν να είναι τύπου decimal μιας και όλα είναι δεκαδικοί αριθμοί με μόνες εξαιρέσεις την ηλικία του ασθενή, το φύλο του και το identifier του.



ΣΚΑΡΟΣ ΓΕΩΡΓΙΟΣ (Π15128)

ΓΙΟΛΗΣ ΔΗΜΗΤΡΗΣ (Π16023)

ΚΩΣΤΑΚΗ ΕΙΡΗΝΗ (Π16064)



Στην συνέχεια έγινε η εισαγωγή στο περιβάλλον της R.

The screenshot displays the RStudio interface. The main window shows a data frame with 10 rows and 14 columns. The columns are: subject, age, sex, test_time, motor_UPDRS, total_UPDRS, Jitter..., JitterAbs., JitterRAP, JitterPPQS, JitterDDP, Shimmer, ShimmerdB., ShimmerAPQ3, and ShimmerAPQ5. The Environment pane on the right shows the 'parkinsons_updrs' object with 5875 observations and 22 variables. The Console pane at the bottom shows the R version 3.5.0 (2018-04-23) and the command to read the CSV file: `parkinsons_updrs <- read.csv("C:/Users/George/desktop/unipi/6o/dataining/parkinson/parkinsons_updrs.data", quote="")`.

subject	age	sex	test_time	motor_UPDRS	total_UPDRS	Jitter...	JitterAbs.	JitterRAP	JitterPPQS	JitterDDP	Shimmer	ShimmerdB.	ShimmerAPQ3	ShimmerAPQ5
1	72	0	5.6431	28.199	34.398	0.00662	3.380e-05	0.00401	0.00317	0.01204	0.02565	0.230	0.01438	0.0
2	72	0	12.6660	28.447	34.894	0.00300	1.680e-05	0.00132	0.00130	0.00395	0.02024	0.179	0.00994	0.0
3	72	0	19.6810	28.695	35.389	0.00481	2.462e-05	0.00205	0.00308	0.00616	0.01675	0.181	0.00734	0.0
4	72	0	25.6470	28.905	35.810	0.00528	2.657e-05	0.00191	0.00364	0.00573	0.02309	0.327	0.01106	0.0
5	72	0	33.6420	29.187	36.375	0.00335	2.014e-05	0.00093	0.00130	0.00278	0.01703	0.176	0.00679	0.0
6	72	0	40.6520	29.435	36.870	0.00333	2.290e-05	0.00119	0.00139	0.00357	0.02227	0.214	0.01006	0.0
7	72	0	47.6490	29.682	37.363	0.00422	2.404e-05	0.00212	0.00221	0.00637	0.04352	0.445	0.02376	0.0
8	72	0	54.6400	29.928	37.857	0.00476	2.471e-05	0.00226	0.00239	0.00678	0.02191	0.212	0.00979	0.0
9	72	0	61.6690	30.177	38.353	0.00432	2.854e-05	0.00156	0.00207	0.00468	0.04296	0.371	0.01774	0.0
10	72	0	68.6880	30.424	38.849	0.00496	2.702e-05	0.00258	0.00253	0.00773	0.03610	0.300	0.02030	0.0

ΣΚΑΡΟΣ ΓΕΩΡΓΙΟΣ (Π15128)

ΓΙΟΛΗΣ ΔΗΜΗΤΡΗΣ (Π16023)

ΚΩΣΤΑΚΗ ΕΙΡΗΝΗ (Π16064)

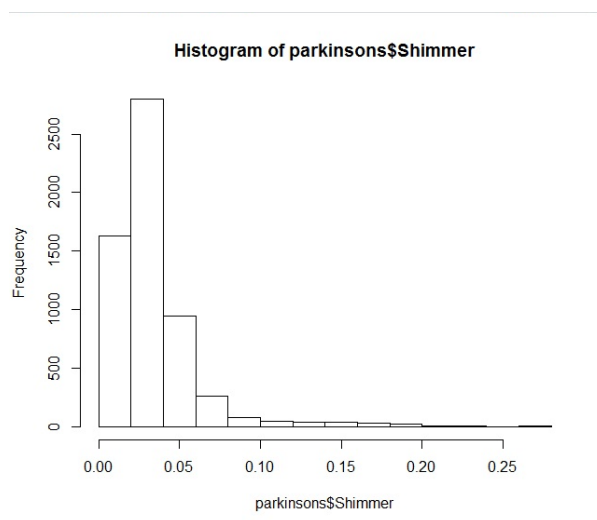
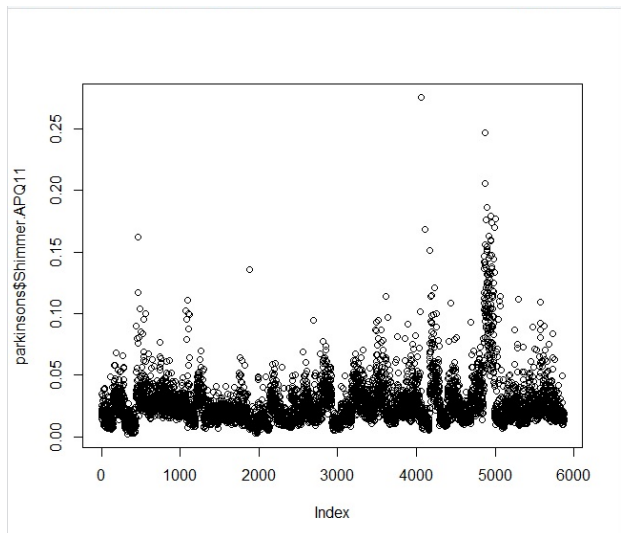
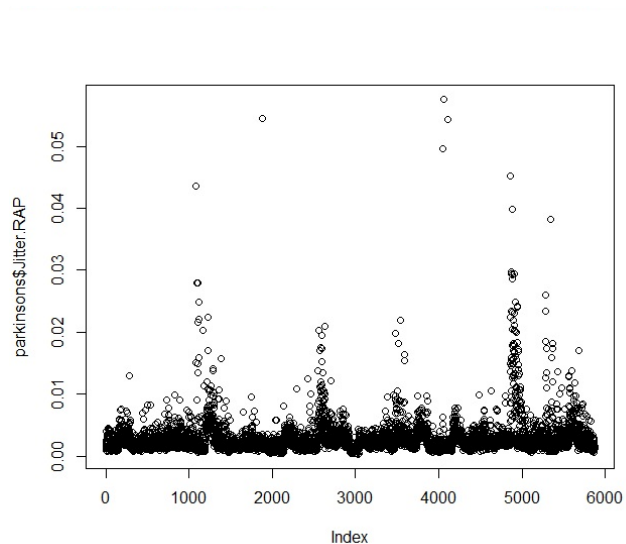
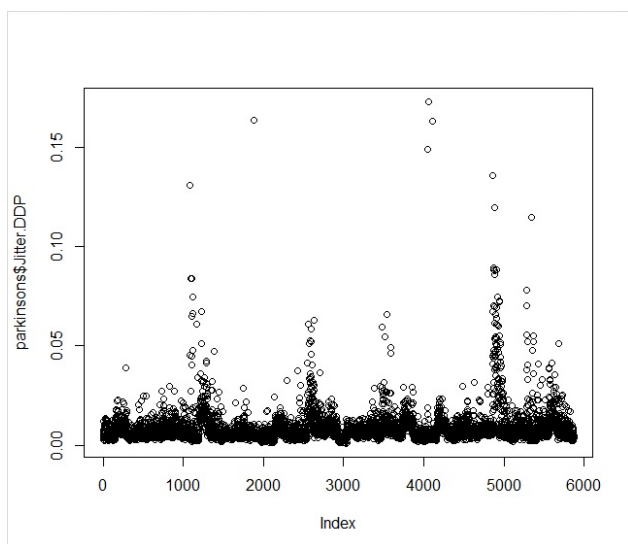


Προεπεξεργασία δεδομένων

Εύρεση και διαχείριση μη λογικών τιμών Outliers

Εύρεση τιμών

Σε αυτό το βήμα έγιναν scatter plots για όλες τις στήλες των δεδομένων με αποτέλεσμα να βρεθούν μη λογικές και ακραίες τιμές. Παρακάτω θα παραθέσω μερικά από τα plots που έγιναν ενδεικτικά. Τα υπόλοιπα screenshots υπάρχουν στον φάκελο screens του παραδοτέου.



ΣΚΑΡΟΣ ΓΕΩΡΓΙΟΣ (Π15128)

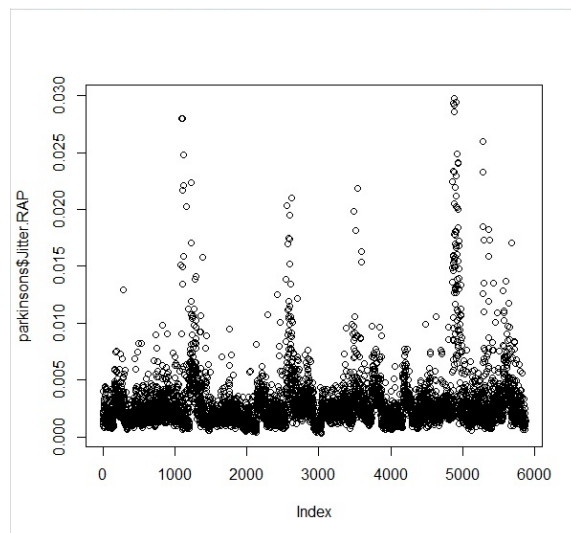
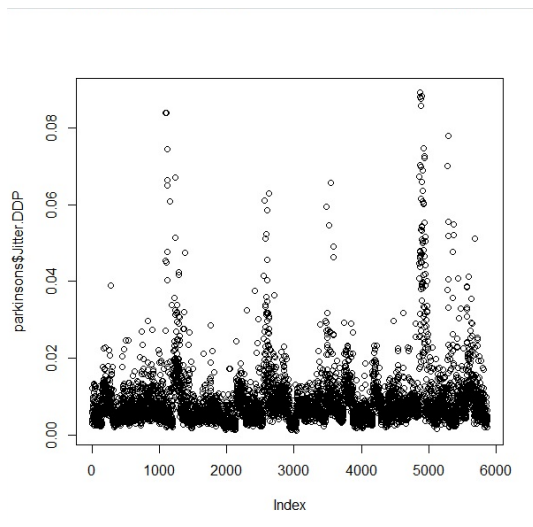
ΓΙΟΛΗΣ ΔΗΜΗΤΡΗΣ (Π16023)

ΚΩΣΤΑΚΗ ΕΙΡΗΝΗ (Π16064)

Διαχείριση τιμών

Στην περίπτωση μας βλέπουμε ότι υπάρχουν αρκετές τιμές που είναι μακριά από το σύνολο των δεδομένων. Ο τρόπος διαχείρισης που επιλέχθηκε για την περίπτωση μας είναι να ανατεθεί σε αυτές τις τιμές η μέση τιμή της στήλης αντί να τις διαγράψουμε διότι το σύνολο των δεδομένων μας είναι αρκετά μικρό. Παρακάτω παραθέτω δύο από τα «καθαρισμένα» παραδείγματα.

-παρατήρηση: Σε αυτό το σημείο μπορούμε να παρατηρήσουμε ότι οι τιμές των δεδομένων ανάμεσα 4900 και 5000 έχουν πολύ ακραίες τιμές σε σύγκριση με τις υπόλοιπες. Αν το δείγμα μας ήταν μεγαλύτερο αυτές οι τιμές θα έπρεπε να αφαιρεθούν τελείως ή να κανονικοποιηθούν αλλά στην περίπτωση μας αυτό θα άλλαζε πολύ το δείγμα μας οπότε αφέθηκαν ως έχει



Κοινωνικοποίηση του πίνακα

Σε αυτό το βήμα χρησιμοποιήθηκε η συνάρτηση scale για να κανονικοποιηθούν όλες οι στήλες πέραν του φύλου. Σε αυτό το σημείο έκρινα σωστό να παραλείψω την στήλη με το identifier του κάθε ασθενή καθώς μας είναι άχρηστη την υπόλοιπη διαδικασία.

Μείωση του όγκου δεδομένων

Πέρα από την παραπάνω παρατήρηση που έγινε για την μείωση του πλήθους εγγραφών, ένας ακόμα τρόπος για να μειωθούν τα δεδομένα είναι να παραλείψουμε κάποια από τις στήλες με κριτήριο της συσχέτισης της (correlation) που έχει με τις άλλες. Αν το correlation είναι κοντά στο 1 αυτό σημαίνει ότι οι 2 στήλες είναι ανάλογες με αποτέλεσμα να μην χρειάζονται και οι δύο στο data set μειώνοντας έτσι την πολυπλοκότητα.



Classification/regression

Linear regression

Αρχικά με την συνάρτηση `Sample` χώρισα το data set σε test και train(20% test 80% train). Σε αυτό το σημείο πρέπει να σημειωθεί ότι το data set έχει δύο στήλες στόχους πράγμα που σημαίνει ότι δοκίμασα όλους τους πιθανούς συνδυασμούς για να φτιάξω το γραμμικό μοντέλο. Εξηγώντας, οι 4 επιλογές είναι οι εξής :

- στήλη στόχος το `motor_UPDRS`
- στήλη στόχος το `total_UPDRS`
- στήλη στόχος το `motor_UPDRS`, συμπεριλαμβάνοντας την στήλη `total_UPDRS`
- στήλη στόχος το `total_UPDRS`, συμπεριλαμβάνοντας την στήλη `motor_UPDRS`

Αφού έφτιαξα το μοντέλο με την συνάρτηση `lm()` και για τις 4 περιπτώσεις εκπαίδευσα με το ίδιο train (πειράζοντας βέβαια τις στήλες όπου χρειαζόταν) τα 4 μοντέλα μου.

```
> summary(m_motor)

Call:
lm(formula = motor_UPDRS ~ age + sex + test_time + jitter... +
    jitter.Abs. + jitter.RAP + jitter.PPQ5 + jitter.DDP + Shimmer +
    Shimmer.db. + Shimmer.APQ3 + Shimmer.APQ5 + Shimmer.APQ11 +
    Shimmer.DDA + NHR + HNR + RPDE + DFA + PPE, data = train)

Residuals:
    Min       1Q   Median       3Q      Max
-3.5345 -0.7311 -0.1056  0.7535  2.3457

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.05938    0.01721   3.451 0.000564 ***
age          0.19631    0.01421  13.816 < 2e-16 ***
sex         -0.19062    0.03422  -5.571 2.68e-08 ***
test_time    0.08245    0.01341   6.149 8.46e-10 ***
jitter...    0.17131    0.03554   4.820 1.48e-06 ***
jitter.Abs.  -0.40331    0.04004  -10.072 < 2e-16 ***
jitter.RAP   -12.83119   13.03694  -0.984 0.325059
jitter.PPQ5  -0.01025    0.03006  -0.341 0.733180
jitter.DDP   13.06298   13.03680   1.002 0.316391
Shimmer      0.07877    0.06216   1.267 0.205121
Shimmer.db.  0.05758    0.08368   0.688 0.491451
Shimmer.APQ3 -0.01588    0.12010  -0.132 0.894847
Shimmer.APQ5 -0.20976    0.08258  -2.540 0.011116 *
Shimmer.APQ11 0.26144    0.04945   5.287 1.30e-07 ***
Shimmer.DDA  -0.24525    0.12243  -2.003 0.045214 *
NHR          -0.15106    0.03494  -4.323 1.57e-05 ***
HNR          -0.22925    0.03124  -7.339 2.52e-13 ***
RPDE         0.02681    0.01932   1.388 0.165275
DFA          -0.21210    0.01711  -12.393 < 2e-16 ***
PPE          0.17615    0.02819   6.249 4.50e-10 ***

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9177 on 4674 degrees of freedom
Multiple R-squared:  0.1622,    Adjusted R-squared:  0.1588
F-statistic: 47.61 on 19 and 4674 DF,  p-value: < 2.2e-16
> |

> summary(m_total)

Call:
lm(formula = total_UPDRS ~ age + sex + test_time + jitter... +
    jitter.Abs. + jitter.RAP + jitter.PPQ5 + jitter.DDP + Shimmer +
    Shimmer.db. + Shimmer.APQ3 + Shimmer.APQ5 + Shimmer.APQ11 +
    Shimmer.DDA + NHR + HNR + RPDE + DFA + PPE, data = train)

Residuals:
    Min       1Q   Median       3Q      Max
-2.6473 -0.6977 -0.1735  0.6927  2.4758

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.08887    0.01712   5.191 2.18e-07 ***
age          0.24367    0.01413  17.239 < 2e-16 ***
sex         -0.27870    0.03404  -8.188 3.40e-16 ***
test_time    0.08617    0.01334   6.460 1.16e-10 ***
jitter...    0.15050    0.03536   4.257 2.11e-05 ***
jitter.Abs.  -0.27724    0.03984  -6.959 3.89e-12 ***
jitter.RAP   -10.24644   12.96939  -0.790 0.4295
jitter.PPQ5  -0.03733    0.02990  -1.248 0.2120
jitter.DDP   10.42207   12.96925   0.804 0.4217
Shimmer      0.06079    0.06184   0.983 0.3256
Shimmer.db.  -0.01185    0.08325  -0.142 0.8868
Shimmer.APQ3 -0.02995    0.11948  -0.251 0.8021
Shimmer.APQ5 -0.08461    0.08216  -1.030 0.3031
Shimmer.APQ11 0.15365    0.04919   3.124 0.0018 **
Shimmer.DDA  -0.19089    0.12180  -1.567 0.1171
NHR          -0.14891    0.03476  -4.284 1.87e-05 ***
HNR          -0.26390    0.03108  -8.492 < 2e-16 ***
RPDE         0.04232    0.01922   2.202 0.0277 *
DFA          -0.21899    0.01703  -12.862 < 2e-16 ***
PPE          0.12555    0.02804   4.477 7.74e-06 ***

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.913 on 4674 degrees of freedom
Multiple R-squared:  0.1769,    Adjusted R-squared:  0.1736
F-statistic: 52.88 on 19 and 4674 DF,  p-value: < 2.2e-16
> |
```

Σημαντικό ρόλο στην αξιοπιστία του μοντέλου παίζουν οι τιμές **p-value** και **R-squared** που βλέπουμε κάτω δεξιά στις εικόνες. Το **p-value** θέλουμε να είναι μικρότερο του 0,05 και το **R-squared** όσο το δυνατόν πιο κοντά γίνεται στο 1.

Με τα μοντέλα που έφτιαξα κάνω πρόβλεψη και για τις 4 περιπτώσεις με την συνάρτηση `predict()` και δημιουργώ 4 προβλεπτές και προσπαθώ να προβλέψω την αντίστοιχη τιμή της κάθε στήλης στόχου κάθε φορά. Οι δύο παραπάνω προβλεπτές έχουν πολύ κακό ποσοστό

ΣΚΑΡΟΣ ΓΕΩΡΓΙΟΣ (Π15128)

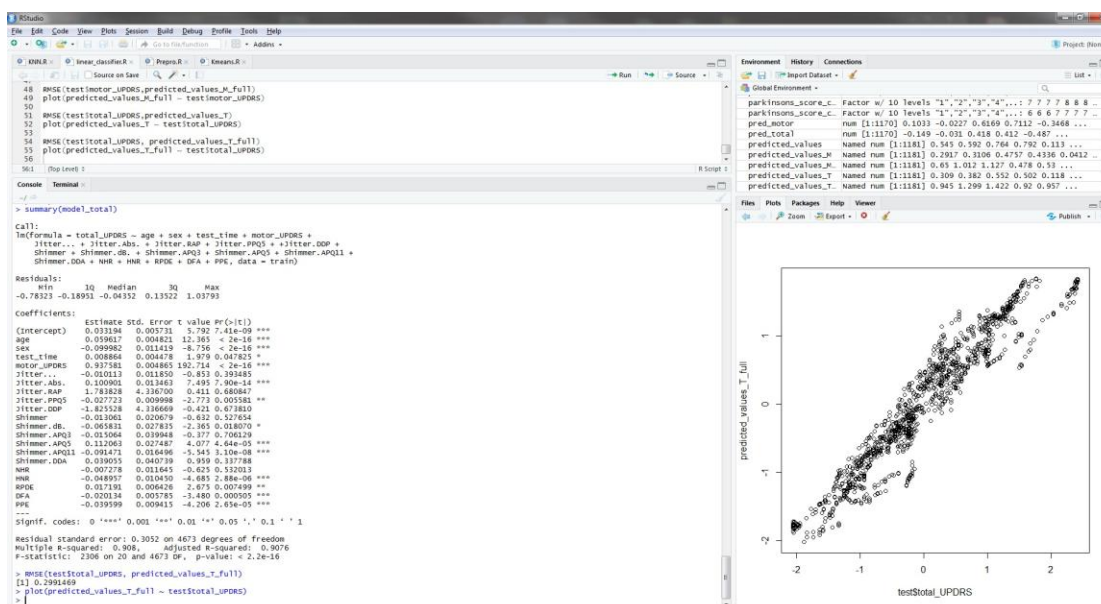
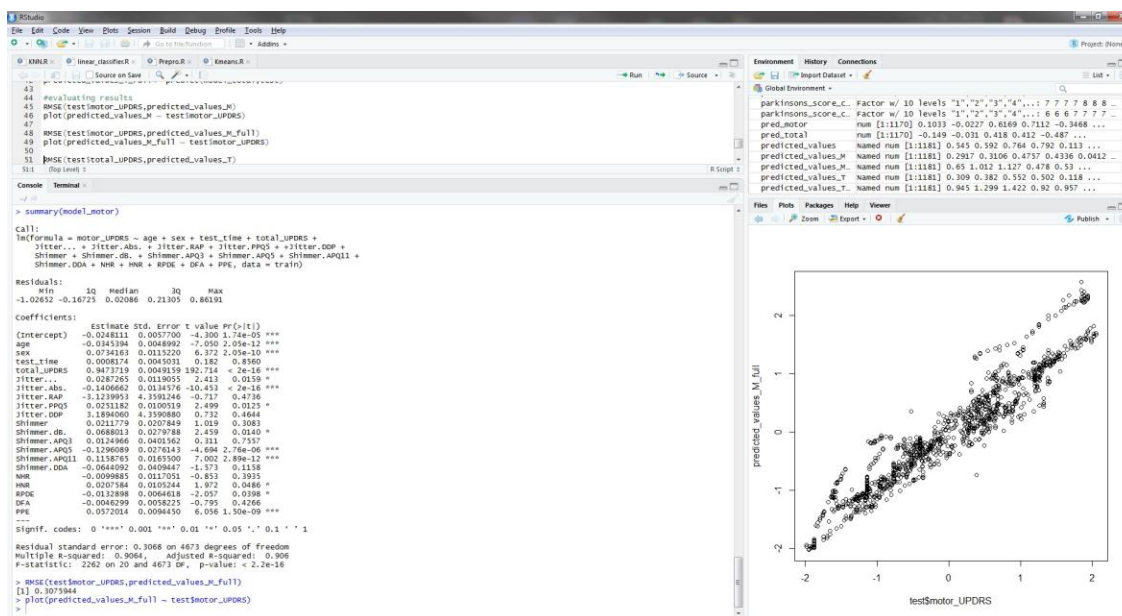
ΓΙΟΛΗΣ ΔΗΜΗΤΡΗΣ (Π16023)

ΚΩΣΤΑΚΗ ΕΙΡΗΝΗ (Π16064)



πρόβλεψης αφού το R-squared είναι πολύ κοντά στο 0 οπότε δεν θα δείξω παραδείγματα από scatter plot από αυτούς τους δύο προβλεπτές.

Οι άλλες δύο περιπτώσεις έχουν μια ικανοποιητική πρόβλεψη αλλά και πάλι όχι αρκετά καλή έτσι ώστε να είναι ο αλγόριθμος που θα προτιμηθεί από αυτούς που υλοποιήθηκαν. Παρακάτω παραθέτω τα Summary ,RMSE values ,plots για τις δύο πιο καλά προσαρμοσμένες περιπτώσεις του προβλήματος μας.



ΣΚΑΡΟΣ ΓΕΩΡΓΙΟΣ (Π15128)

ΓΙΟΛΗΣ ΔΗΜΗΤΡΗΣ (Π16023)

ΚΩΣΤΑΚΗ ΕΙΡΗΝΗ (Π16064)

KNN

Πριν το μοντέλο του προβλεπτή

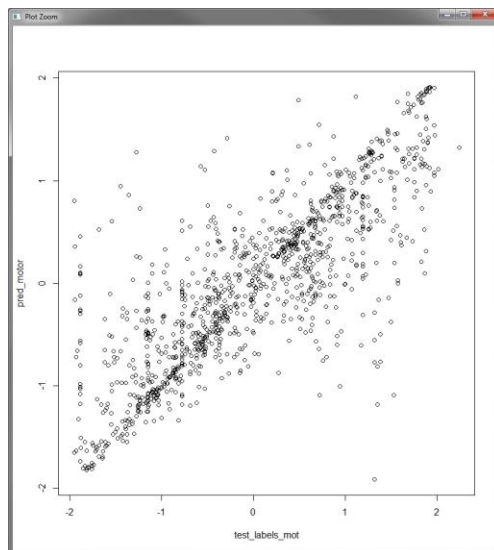
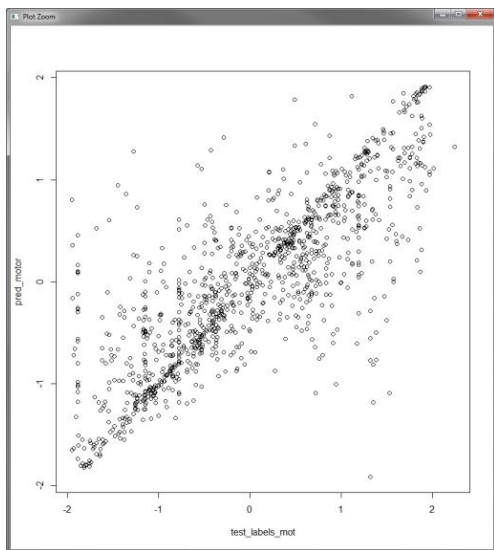
Σε αυτό το ερώτημα εργάστηκα παρόμοια με το προηγούμενο ως αναφορά την στήλη στόχο.

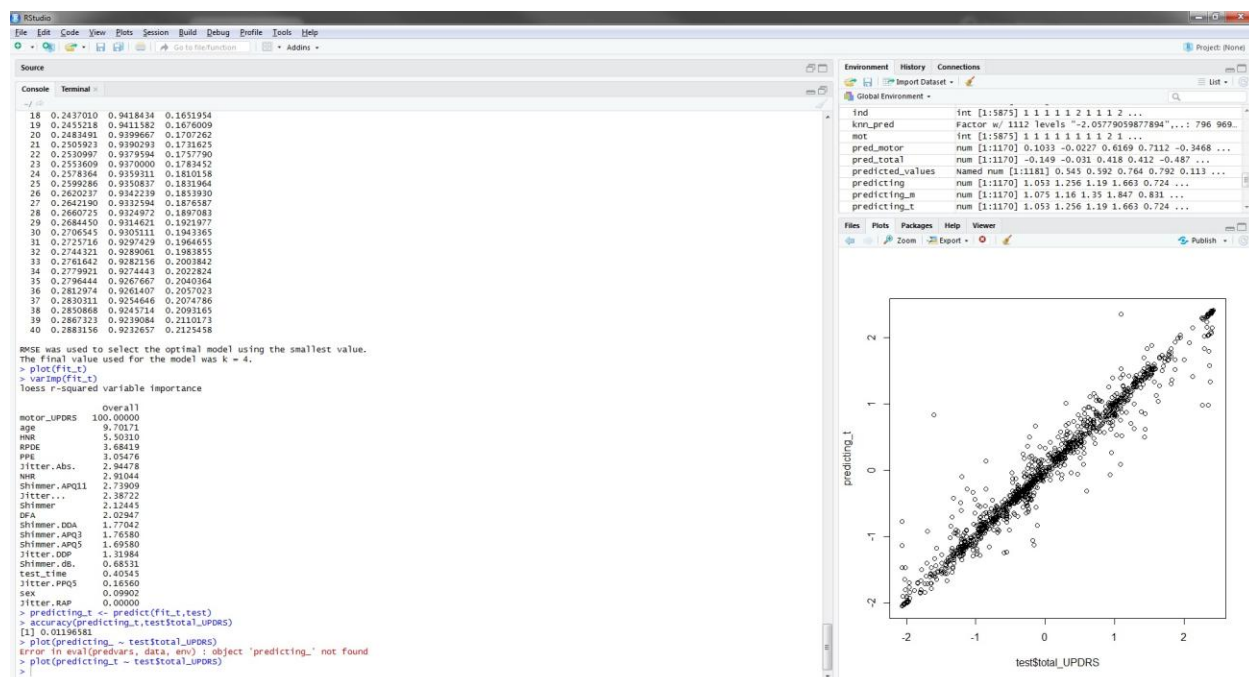
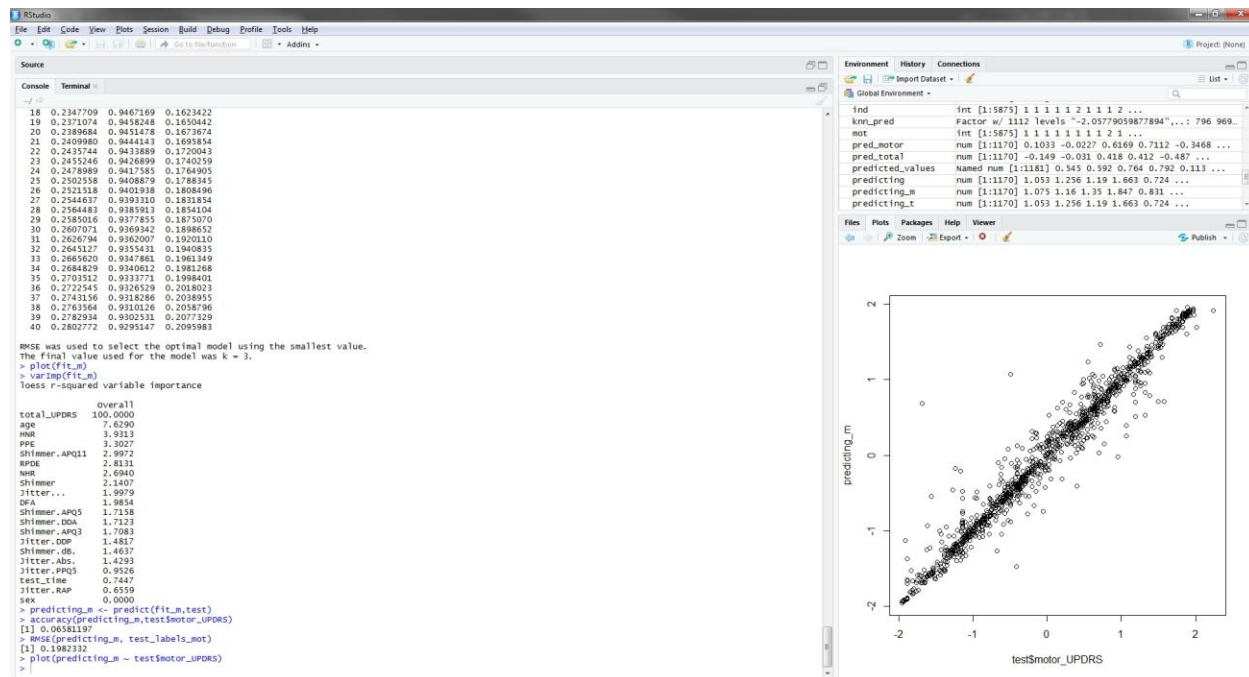
Αρχικά δημιουργήθηκε μια συνάρτηση για την διαδικασία του cross validation η οποία παίρνει το train set το κόβει σε 10 κομμάτια και ελέγχει να βρει το καλύτερο ποσοστό λάθους αφού έχει βρει το καθένα ξεχωριστά. Αυτή η διαδικασία επαναλαμβάνεται με 3 διαφορετικά διπλώματα (folds) του σετ μας ψάχνοντας να βρει όσο το δυνατόν καλύτερο αποτέλεσμα.

Προβλεπτής

Όμοια με τον γραμμικό προβλεπτή φτιάχνω 4 διαφορετικά μοντέλα για τις ανάγκες μου και τα υλοποιώ και συγκρίνω τα αποτελέσματα είτε οπτικά (με κάποιο plot) είτε από την τιμή της συνάρτησης accuracy. Το k επιλέχθηκε από τον αλγόριθμο αφού ήταν έτσι φτιαγμένος ώστε να δοκιμάζει διαφορετικά k και να επιλέγει αυτό που δίνει το καλύτερο αποτέλεσμα.

Παρακάτω δείχνω plots από τα μοντέλα και τα accuracy scores όπως και τα Variable importance της κάθε περίπτωσης. Από το τελευταίο μπορούμε να δούμε ποιες από τις μεταβλητές δεν παίζουν πολύ ρόλο στην πρόβλεψη που γίνεται. Τα δύο πρώτα plots είναι από τις περιπτώσεις που δεν είχαμε την δεύτερη μεταβλητή στόχο στις μεταβλητές πρόβλεψης και τα επόμενα είναι με τις μεταβλητές. Είναι φανερό το γεγονός ότι όταν βάζουμε τις μεταβλητές στόχους μέσα τα αποτελέσματα είναι πολύ καλύτερα.





ΣΚΑΡΟΣ ΓΕΩΡΓΙΟΣ (Π15128)

ΓΙΟΛΗΣ ΔΗΜΗΤΡΗΣ (Π16023)

ΚΩΣΤΑΚΗ ΕΙΡΗΝΗ (Π16064)



Clustering

Kmeans

Αρχικά παραλήφθηκε η στήλη στόχος με όλους τους πιθανούς συνδυασμούς.

Ο kmeans είναι ένας αλγόριθμος ο οποίος λειτουργεί με κεντροειδή, αυτό σημαίνει ότι η πρώτη διαδικασία είναι να τρέψουμε έναν αλγόριθμο για να μας υπολογίσει αυτά τα κεντροειδή. Μετά από δύο αλγορίθμους που χρησιμοποιήθηκαν το ιδανικό k υπολογίστηκε ότι είναι 3 , δηλαδή ότι το σύνολο των δεδομένων μας χωρίζεται σε 3 clusters.

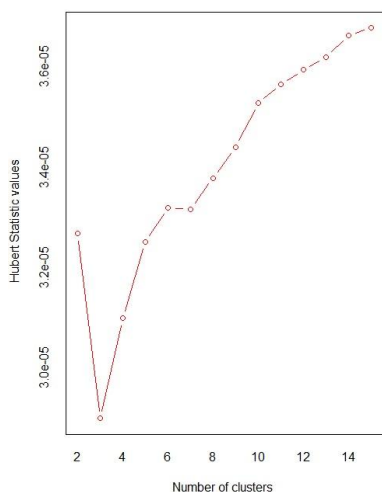
```
> NbClust(data = parkinsons_clust, diss = NULL, distance = "euclidean",
+         min.nc = 2, max.nc = 15, method = "kmeans")
*** : The Hubert index is a graphical method of determining the number of clusters.
      In the plot of Hubert index, we seek a significant knee that corresponds to a
      significant increase of the value of the measure i.e the significant peak in Hubert
      index second differences plot.

*** : The D index is a graphical method of determining the number of clusters.
      In the plot of D index, we seek a significant knee (the significant peak in Dindex
      second differences plot) that corresponds to a significant increase of the value of
      the measure.

*****
* Among all indices:
* 6 proposed 2 as the best number of clusters
* 12 proposed 3 as the best number of clusters
* 1 proposed 4 as the best number of clusters
* 1 proposed 5 as the best number of clusters
* 1 proposed 7 as the best number of clusters
* 1 proposed 14 as the best number of clusters
* 2 proposed 15 as the best number of clusters

***** conclusion *****

* According to the majority rule, the best number of clusters is 3
```



Αφού έτρεξα τον ολόθερμο για κάθε ένα από τα σύνολα πηρά την πρωτοβουλία να χωρίσω κάθε ένα από τα χαρακτηριστικά στόχους σε bins (ίσα σαν μέγεθος διαστήματα) για να μπορέσω να παρατηρήσω ποσοστά ευστοχίας του αλγορίθμου ,το οποίο θα παραθέσω παρακάτω μαζί με τα clusters που δημιουργήθηκαν.

[illegible][illegible]

```
> table(parkinsons_score_col_M, clusters_M$cluster)/nrow(parkinsons_clust)
```

```
parkinsons_score_col_M      1      2      3
1 0.222978723 0.001531915 0.096000000
2 0.257531915 0.024680851 0.161191489
3 0.128340426 0.004765957 0.102978723
```

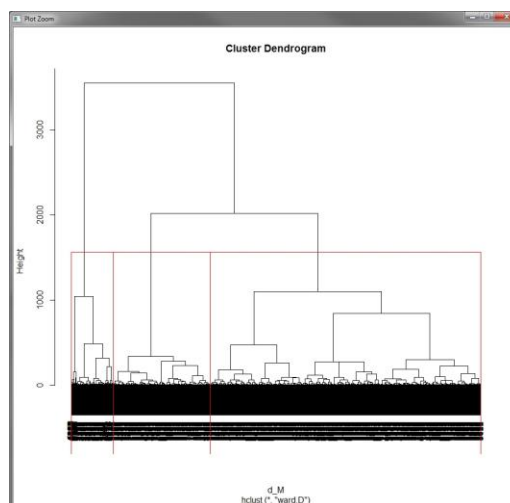
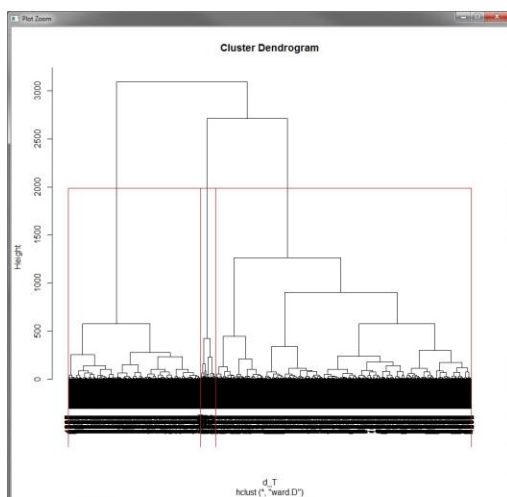
```
> table(parkinsons_score_col_T, clusters_T$cluster)/nrow(parkinsons_clust)
```

```
parkinsons_score_col_T      1      2      3
1 0.086127660 0.001531915 0.199829787
2 0.175829787 0.025191489 0.307574468
3 0.095489362 0.004255319 0.104170213
```

> |

Hierarchical

Αυτός ο αλγόριθμος ξεκινά και ενώνει γειτονικά στοιχεία μέχρι να καταλήξει σε ένα σύνολο. Ανάλογα με το ποσό ψηλά ενώνονται τα υποσύνολα καταλαβαίνουμε τόσα είναι τα clusters του αλγορίθμου. Στο παράδειγμα μας έχουν χωριστεί τα σύνολα στο δέντρο με κόκκινο χρώμα και αφού έχει γίνει αυτό κόβουμε το δέντρο και έχουμε τα 3 cluster μας και μπορούμε με τον ίδιο τρόπο που έγινε και στον kmeans να υπολογίσουμε ποσοστά επιτυχίας.



```
> table(parkinsons_hier_score_M, groups_M)
      groups_M
parkinsons_hier_score_M  1    2    3
1    656 1158   69
2    469 1720  416
3    271 1000  116

> table(parkinsons_hier_score_T, groups_T)
      groups_T
parkinsons_hier_score_T  1    2    3
1    811  867   11
2   2016  796  176
3    902  266   30

> |

> table(parkinsons_hier_score_M, groups_M)/nrow(parkinsons)
      groups_M
parkinsons_hier_score_M  1    2    3
1 0.11165957 0.19710638 0.01174468
2 0.07982979 0.29276596 0.07080851
3 0.04612766 0.17021277 0.01974468

> table(parkinsons_hier_score_T, groups_T)/nrow(parkinsons)
      groups_T
parkinsons_hier_score_T  1    2    3
1 0.138042553 0.147574468 0.001872340
2 0.343148936 0.135489362 0.029957447
3 0.153531915 0.045276596 0.005106383

> |
```



Association rule mining

Σε αυτό το βήμα επειδή τα δεδομένα μας είναι αριθμητικά επρεπε να τα κάνουμε κατηγορικά χρησιμοποιώντας την συνάρτηση cut(). Αυτή η συνάρτηση μας επιτρέπει να χωρίσουμε ένα σύνολο από numeric δεδομένα σε ίσα διαστήματα έτσι ώστε να μπορεί ο apriori() αλγόριθμος που θα χρησιμοποιήσουμε στην συνέχεια να υπολογίσει πιθανότητες εμφάνισης. Το πλήθος της κάθε κατηγορίας που σχηματίζεται υπολογίστηκε εμπειρικά από τα histograms της κάθε μεταβλητής. Παρακάτω παραθέτω ένα στιγμιότυπο από την εκτέλεση του αλγορίθμου οπού υπάρχουν οι καλύτερες 50 συσχετίσεις ταξινομημένες ως αναφορά την εμπιστοσύνη που δείχνουν στον σύνολο των δεδομένων άλλα και την τιμή lift.

```
> inspect(head(sort(parkinson_rules, by = "lift"), 150))
```

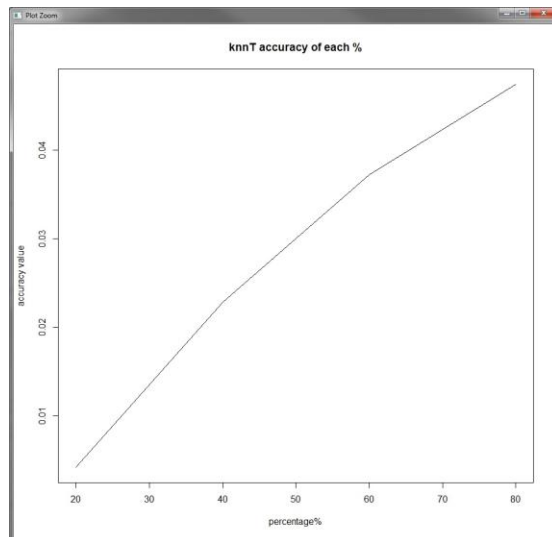
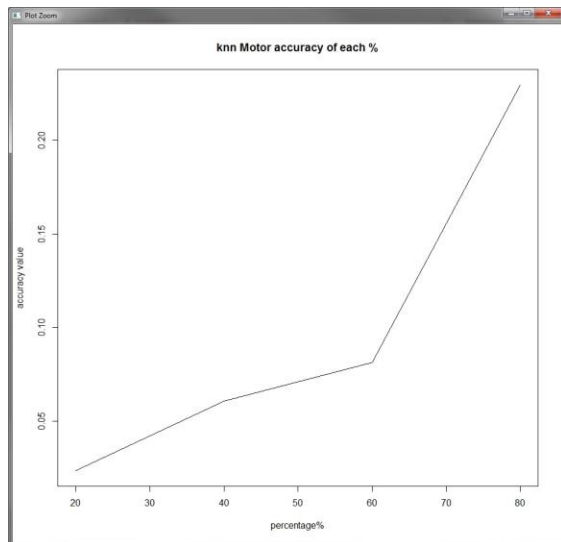
	lhs	rhs	support	confidence	lift	count
[1]	{jitter.RAP=3}	=> {jitter.DDP=3}	0.01038298	1	96.31148	61
[2]	{jitter.DDP=3}	=> {jitter.RAP=3}	0.01038298	1	96.31148	61
[3]	{age=6,total_UPDRS=1}	=> {motor_UPDRS=1}	0.02144681	1	36.26543	126
[4]	{total_UPDRS=1,HNR=7}	=> {motor_UPDRS=1}	0.01140426	1	36.26543	67
[5]	{age=6,total_UPDRS=1,PPE=2}	=> {motor_UPDRS=1}	0.01157447	1	36.26543	68
[6]	{age=6,total_UPDRS=1,Shimmer=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[7]	{age=6,total_UPDRS=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[8]	{age=6,total_UPDRS=1,Shimmer.db.=1}	=> {motor_UPDRS=1}	0.01276596	1	36.26543	75
[9]	{age=6,total_UPDRS=1,Shimmer.DDA=1}	=> {motor_UPDRS=1}	0.01242553	1	36.26543	73
[10]	{age=6,total_UPDRS=1,jitter.Abs.=1}	=> {motor_UPDRS=1}	0.01514894	1	36.26543	89
[11]	{age=6,total_UPDRS=1,Shimmer.APQ3=1}	=> {motor_UPDRS=1}	0.01412766	1	36.26543	83
[12]	{age=6,total_UPDRS=1,Shimmer.APQ5=1}	=> {motor_UPDRS=1}	0.01582979	1	36.26543	93
[13]	{age=6,total_UPDRS=1,jitter.RAP=1}	=> {motor_UPDRS=1}	0.02076596	1	36.26543	122
[14]	{age=6,total_UPDRS=1,jitter.DDP=1}	=> {motor_UPDRS=1}	0.02076596	1	36.26543	122
[15]	{age=6,total_UPDRS=1,jitter.PPQ5=1}	=> {motor_UPDRS=1}	0.02093617	1	36.26543	123
[16]	{age=6,total_UPDRS=1,HNR=1}	=> {motor_UPDRS=1}	0.02144681	1	36.26543	126
[17]	{age=6,total_UPDRS=1,jitter...=1}	=> {motor_UPDRS=1}	0.02144681	1	36.26543	126
[18]	{total_UPDRS=1,jitter.Abs.=1,HNR=7}	=> {motor_UPDRS=1}	0.01021277	1	36.26543	60
[19]	{total_UPDRS=1,Shimmer.APQ5=1,HNR=7}	=> {motor_UPDRS=1}	0.01021277	1	36.26543	60
[20]	{total_UPDRS=1,jitter.RAP=1,HNR=7}	=> {motor_UPDRS=1}	0.01140426	1	36.26543	67
[21]	{total_UPDRS=1,jitter.DDP=1,HNR=7}	=> {motor_UPDRS=1}	0.01140426	1	36.26543	67
[22]	{total_UPDRS=1,jitter.PPQ5=1,HNR=7}	=> {motor_UPDRS=1}	0.01140426	1	36.26543	67
[23]	{total_UPDRS=1,HNR=1,HNR=7}	=> {motor_UPDRS=1}	0.01140426	1	36.26543	67
[24]	{total_UPDRS=1,jitter...=1,HNR=7}	=> {motor_UPDRS=1}	0.01140426	1	36.26543	67
[25]	{age=6,total_UPDRS=1,jitter.Abs.=1,PPE=2}	=> {motor_UPDRS=1}	0.01123404	1	36.26543	66
[26]	{age=6,total_UPDRS=1,Shimmer.APQ3=1,PPE=2}	=> {motor_UPDRS=1}	0.01004255	1	36.26543	59
[27]	{age=6,total_UPDRS=1,Shimmer.APQ5=1,PPE=2}	=> {motor_UPDRS=1}	0.01089362	1	36.26543	64
[28]	{age=6,total_UPDRS=1,jitter.RAP=1,PPE=2}	=> {motor_UPDRS=1}	0.01157447	1	36.26543	68
[29]	{age=6,total_UPDRS=1,jitter.DDP=1,PPE=2}	=> {motor_UPDRS=1}	0.01157447	1	36.26543	68
[30]	{age=6,total_UPDRS=1,jitter.PPQ5=1,PPE=2}	=> {motor_UPDRS=1}	0.01157447	1	36.26543	68
[31]	{age=6,total_UPDRS=1,HNR=1,PPE=2}	=> {motor_UPDRS=1}	0.01157447	1	36.26543	68
[32]	{age=6,total_UPDRS=1,jitter...=1,PPE=2}	=> {motor_UPDRS=1}	0.01157447	1	36.26543	68
[33]	{age=6,total_UPDRS=1,Shimmer=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[34]	{age=6,total_UPDRS=1,Shimmer=1,Shimmer.db.=1}	=> {motor_UPDRS=1}	0.01021277	1	36.26543	60
[35]	{age=6,total_UPDRS=1,Shimmer=1,Shimmer.DDA=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[36]	{age=6,total_UPDRS=1,Shimmer=1,Shimmer.APQ3=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[37]	{age=6,total_UPDRS=1,Shimmer=1,Shimmer.APQ5=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[38]	{age=6,total_UPDRS=1,jitter.RAP=1,Shimmer=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[39]	{age=6,total_UPDRS=1,jitter.DDP=1,Shimmer=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[40]	{age=6,total_UPDRS=1,jitter.PPQ5=1,Shimmer=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[41]	{age=6,total_UPDRS=1,Shimmer=1,HNR=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[42]	{age=6,total_UPDRS=1,jitter...=1,Shimmer=1}	=> {motor_UPDRS=1}	0.01038298	1	36.26543	61
[43]	{age=6,total_UPDRS=1,Shimmer.db.=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01174468	1	36.26543	69
[44]	{age=6,total_UPDRS=1,Shimmer.APQ11=1,Shimmer.DDA=1}	=> {motor_UPDRS=1}	0.01140426	1	36.26543	67
[45]	{age=6,total_UPDRS=1,jitter.Abs.=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01157447	1	36.26543	68
[46]	{age=6,total_UPDRS=1,Shimmer.APQ3=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01191489	1	36.26543	70
[47]	{age=6,total_UPDRS=1,Shimmer.APQ5=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[48]	{age=6,total_UPDRS=1,jitter.RAP=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[49]	{age=6,total_UPDRS=1,jitter.DDP=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[50]	{age=6,total_UPDRS=1,jitter.PPQ5=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[51]	{age=6,total_UPDRS=1,Shimmer.APQ11=1,HNR=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[52]	{age=6,total_UPDRS=1,jitter...=1,Shimmer.APQ11=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[53]	{age=6,total_UPDRS=1,Shimmer.db.=1,Shimmer.DDA=1}	=> {motor_UPDRS=1}	0.01191489	1	36.26543	70
[54]	{age=6,total_UPDRS=1,jitter.Abs.=1,Shimmer.db.=1}	=> {motor_UPDRS=1}	0.01174468	1	36.26543	69
[55]	{age=6,total_UPDRS=1,Shimmer.db.=1,Shimmer.APQ3=1}	=> {motor_UPDRS=1}	0.01259574	1	36.26543	74
[56]	{age=6,total_UPDRS=1,Shimmer.db.=1,Shimmer.APQ5=1}	=> {motor_UPDRS=1}	0.01276596	1	36.26543	75
[57]	{age=6,total_UPDRS=1,jitter.RAP=1,Shimmer.db.=1}	=> {motor_UPDRS=1}	0.01276596	1	36.26543	75

Επίδραση του δείγματος στην ποιότητα του αποτελέσματος

Σε αυτό το ερώτημα κριθήκαμε να ξανατρέξουμε τους αλγορίθμους με κομμάτια από το σύνολο των δεδομένων και να δούμε αν οι αποδώσεις του αλγόριθμου έπεσαν σε σχέση με το αρχικό σύνολο δεδομένων.

Είναι προφανές ότι η απόδοση των αλγορίθμων πέφτει όσο μικραίνουμε το data set, αυτό συμβαίνει διότι το μοντέλο δεν έχει αρκετά δεδομένα έτσι ώστε να μπορέσει να εκπαιδευτεί επαρκώς και να βγάλει ικανοποιητικά αποτελέσματα.

Αφού χώρισα το σύνολο σε κομμάτια της τάξεως των 20,40,60,80% ξαναέτρεξα τους αλγόριθμους και αυτά είναι κάποια από τα αποτελέσματα. Συγκεκριμένα από τον αλγόριθμο knn από τον ιεραρχικό αλγόριθμο.



Τέλος θεώρησα περιττό να κάνω αυτή την διαδικασία για το βήμα 5 μιας και δεν θα μπορούσαμε εύκολα να δούμε διαφορές στα αποτελέσματα αφού κάθε φορά θα είναι διαφορετικές οι συσχετισεις λόγω μικρου μεγέθους δεδομένων.



```
> scoreTablem20
      groups_T_20
parkinsons_hier_score_T_20 1 2 3
1 0.224211424 0.063086104 0.003410060
2 0.371696505 0.105711850 0.025575448
3 0.156862745 0.051150895 0.005967604

> scoreTablem20
      groups_T_20
parkinsons_hier_score_T_20 1 2 3
1 0.224211424 0.063086104 0.003410060
2 0.371696505 0.105711850 0.025575448
3 0.156862745 0.051150895 0.005967604

> scoreTablem40
      groups_T_40
parkinsons_hier_score_T_40 1 2 3
1 0.2313013403 0.0687418936 0.0030263727
2 0.3674881107 0.1175961954 0.0259403372
3 0.1552096844 0.0562040640 0.0008646779

> scoreTablem40
      groups_T_40
parkinsons_hier_score_T_40 1 2 3
1 0.2313013403 0.0687418936 0.0030263727
2 0.3674881107 0.1175961954 0.0259403372
3 0.1552096844 0.0562040640 0.0008646779

> scoreTablem60
      groups_T_60
parkinsons_hier_score_T_60 1 2 3
1 0.246490736 0.028635598 0.002245929
2 0.363559798 0.113700168 0.029758563
3 0.125210556 0.067939360 0.005334082

> scoreTablem60
      groups_T_60
parkinsons_hier_score_T_60 1 2 3
1 0.246490736 0.028635598 0.002245929
2 0.363559798 0.113700168 0.029758563
3 0.125210556 0.067939360 0.005334082

> scoreTablem80
      groups_T_80
parkinsons_hier_score_T_80 1 2 3
1 0.159293917 0.124202467 0.003190132
2 0.180986814 0.290940026 0.038068907
3 0.043173118 0.152275627 0.005954913

> scoreTablem80
      groups_T_80
parkinsons_hier_score_T_80 1 2 3
1 0.159293917 0.124202467 0.003190132
2 0.180986814 0.290940026 0.038068907
3 0.043173118 0.152275627 0.005954913

> |
```