
Ανάλυση Κατηγορικών Δεδομένων

Εργασία 2

Ονοματεπώνυμο : Ιωάννα Τσέτσι

A.M.: s6160095

Διδάσκων : Αντωνία Κορρέ

June 30, 2022

Το Ινδικό Ινστιτούτο του Διαβήτη, πραγματοποίησε μία ερεύνα με σκοπό τη διερεύνηση της σημασίας κάποιων προγνωστικών παραγόντων στην πρόβλεψη της πιθανότητας για διαβήτη σε γυναίκες ηλικίας από 21 ετών και πάνω που προέρχονται από την φυλή Pima. Οι προγνωστικοί παράγοντες που μετρήθηκαν παρουσιάζονται στον Πίνακα 1 ¹.

Για την υλοποίηση της ανάλυσης θα χρησιμοποιήσουμε το στατιστικό πακέτο R. Τα δεδομένα βρίσκονται στην βιβλιοθήκη `mlbench` και το σετ δεδομένων ονομάζεται `Pima Indians Diabetes`.

Εφόσον έχουμε εγκαταστήσει τις απαραίτητες βιβλιοθήκες μπορούμε να προχωρήσουμε στην στατιστική ανάλυση. Αρχικά βλέπουμε ότι στα δεδομένα έχουμε αρκετές ελλείψεις τιμές (`missing values`, `NA`). Για την ανάλυσή των δεδομένων θα τις αφαιρέσουμε από το αρχείο ².

Αριθμός Μεταβλήτης	Όνομα	Τύπος	Σήμασια
1	Pregnant	Αριθμητική	Πόσες φορές έμεινε έγκυος
2	Glucose	Αριθμητική	Συγκέντρωση γλυκόζης στο πλάσμα
3	Pressure	Αριθμητική	Αρτηριακή πίεση (mm Hg)
4	Triceps	Αριθμητική	Πάχος δερματικής πτυχής (mm)
5	Insulin	Αριθμητική	Ινσουλίνη ορού (μ U/ml)
6	Mass	Αριθμητική	Δείκτης μάζας σώματος (βάρος σε kg/(ύψος σε m) ²)
7	Pedigree	Αριθμητική	Score για την πιθανότητα διαβήτη από κληρονομικά αίτια
8	Age	Αριθμητική	Ηλικία (σε έτη)
9	Diabetes	Κατηγορική	Νόσησης από διαβήτη

Table 1: Μεταβλήτες Αρχείου Pima Indians Diabetes

¹Για τις εντολές του έλεγχου του τυπού κάθε μεταβλήτης βλέπε το Παράρτημα [Listing 3](#)

²Για τις εντολές για την αφαίρεση NA βλέπε στο Παράρτημα [Listing 2](#)

Στην συνέχεια θα προχωρήσουμε στη ομαδοποίηση δυο μεταβλητών³. Η πρώτη είναι η ομαδοποίηση της μεταβλητής ηλικίας (Age) στις ομάδες 20-30, 31-40, 41-50 και 50+.

Ενώ η άλλη μεταβλητή είναι ο αριθμός των κυήσεων (Pregnant) ομαδοποιημένες στις ομάδες 0-5, 6-10 και 10+. Τα δεδομένα παρουσιάζονται αναλυτικά στον Πίνακα 2⁴.

Αριθμός Μεταβλητής	Όνομα	Τύπος	Σήμασια	Τιμές
9	Diabetes	Κατηγορική	Νόσησης από διαβήτη	negative positive
8	age	Κατηγορική	Ομαδοποίηση της μεταβλητής Age	"20-30", "31-40" "41-50", "50+"
1	pregnant	Κατηγορική	Ομαδοποίηση της μεταβλητής Pregnant	"0-5", "6-10" "10+"

Table 2: Μεταβλητές Κατηγορικών Δεδομένων .

Σκοπός της μελέτης είναι να φτιάξουμε ένα μοντέλο ώστε να ερμηνεύσουμε την πιθανότητα για διαβήτη σε γυναίκες ηλικίας 21 ετών και πάνω σε σχέση με τους προγνωστικούς παραγόντες που έχουν συλλεχθεί. Για την ανάλυση αυτή θα χρησιμοποιήσουμε το λογιστικό μοντέλο παλινδρομησης (logistic regression) για την διτιμή μεταβλητή diabetes που εκφράζει το ενδεχόμενο νοσησης από διαβήτη. Έστω Y η μεταβλητή diabetes έχουμε ότι

$$Y_i = \begin{cases} 1 & \text{να έχει καποία διαβήτη με πιθανοτητα } p \\ 0 & \text{να μην έχει καποία διαβήτη με πιθανοτητα } 1-p \end{cases} \quad (1)$$

Η λογιστική σύναρτηση είναι η συνάρτηση συνδέσης συνδεσή σε αυτό το γενικεύμενο γραμμικό μοντέλο και για τα δεδομένα μας έχουμε

$$\text{logit}(p) = \text{logit}\left(\frac{p}{1-p}\right) = \beta' \mathbf{X} = \beta_0 + \beta_1 X_1 + \dots + \beta_{11} X_{11}$$

³Για τις εντολές για την ομαδοποίηση των μεταβλητών βλέπε Παράρτημα [Listing 4](#)

⁴Για τις εντολές του τύπου κάθε μεταβλητής βλέπε στο Παράρτημα [Listing 5](#)

Τα μοντέλο που προκύπτει με τα δεδομένα μας είναι το : ⁵

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-8.990364	1.170461	-7.681	1.58e-14 ***
glucose	0.039232	0.005865	6.689	2.24e-11 ***
pressure	-0.002747	0.011892	-0.231	0.8173
triceps	0.015717	0.017320	0.907	0.3642
insulin	-0.000692	0.001351	-0.512	0.6085
mass	0.063485	0.027425	2.315	0.0206 *
pedigree	1.051124	0.434989	2.416	0.0157 *
age31-40	0.802485	0.395535	2.029	0.0425 *
age41-50	1.377709	0.536629	2.567	0.0102 *
age50+	1.236659	0.633326	1.953	0.0509 .
pregnant6-10	-0.150064	0.424600	-0.353	0.7238
pregnant10+	0.930171	0.779488	1.193	0.2327

Table 3: Μοντέλο λογιστικής παλινδρόμησης με όλες τις μεταβλητές

$$\beta' = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \\ \hat{\beta}_4 \\ \hat{\beta}_5 \\ \hat{\beta}_6 \\ \hat{\beta}_7 \\ \hat{\beta}_8 \\ \hat{\beta}_9 \\ \hat{\beta}_{10} \\ \hat{\beta}_{11} \end{bmatrix} \approx \begin{bmatrix} -8.990364 \\ 0.039232 \\ -0.002747 \\ 0.015717 \\ -0.000692 \\ 0.063485 \\ 1.051124 \\ 0.802485 \\ 1.377709 \\ 1.236659 \\ -0.150064 \\ 0.930171 \end{bmatrix}$$

⁵Για τις εντολές για την λογιστική παλινδρόμηση όλων των μεταβλητών βλέπε στο Παράρτημα [Listing 6](#)

Απο τον Πίνακα 3 βλέπουμε ότι η σταθερά $\hat{\beta}_0$ έχει την τιμή -8.90364 . Αυτό σημαίνει ότι το log odds για να νοσήσει κάποια από διαβήτη η οποία ανήκει στην ηλικιακή ομάδα $20 - 30$ και ο αριθμός των κύησεων της είναι από $0 - 5$ είναι -8.90364 .

Για την παράμετρο glucose $\hat{\beta}_1$ βλέπουμε ότι η 1 μοναδιαία αύξηση της συγκέντρωσης γλυκόζης στο πλάσμα αυξάνει τα log odd κατά 0.039232 , όταν ολές οι υπολοιπές μεταβλητές παραμένουν σταθερές.

Ενώ για την παράμετρο pressure $\hat{\beta}_2$ έχουμε ότι 1 μοναδιαία αύξηση της αρτηριακής πίεσης (mm Hg) μειώνει τα log odds κατά -0.002747 , όταν ολές οι υπολοιπές μεταβλητές είναι παραμενουν σταθερές.

Ακόμα για την παράμετρο triceps $\hat{\beta}_3$ έχουμε ότι 1 μοναδιαία αύξηση του πάχους δερματικής πτύχης (mm) αυξάνει τα log odds κατά 0.015717 , όταν ολές οι υπολοιπές μεταβλητές είναι παραμενουν σταθερές.

Επιπλέον για την παράμετρο insulin $\hat{\beta}_4$ παρατηρούμε ότι 1 μοναδιαία αύξηση του ινσουλινη ορού (mu U/ml) μειώνει τα log odds κατά -0.000692 , όταν ολές οι υπολοιπές μεταβλητές είναι παραμενουν σταθερές.

Αντιστοίχα για την παράμετρο mass $\hat{\beta}_5$ παρατηρούμε ότι 1 μοναδιαία αύξηση του δείκτη μάζας σώματος (βάρος σε kg/(ύψος σε m)²) αυξάνει τα log odds κατά 0.063485 , όταν ολές οι υπολοιπές μεταβλητές είναι παραμενουν σταθερές.

Ενώ για την παράμετρο pedigree $\hat{\beta}_6$ παρατηρούμε ότι 1 μοναδιαία αύξηση του Score για την πιθανότητα διαβήτη από κληρονομικά αίτια αυξάνει τα log odds κατά 1.051124 όταν ολές οι υπολοιπές μεταβλητές είναι παραμενουν σταθερές.

Ακόμα για τις γυναίκες που ανήκουν στην ηλικιακή ομάδα $31-40$ από την παράμετρο age31-40, $\hat{\beta}_7$, έχουμε

$$\hat{\beta}_0 + \hat{\beta}_7 = -8.90364 + 0.802485 = -8.187879$$

Δηλαδή τα log odds για τις γυναίκες που νοσούν από διαβήτη και ανήκουν στην ηλικιακή κατηγορία $31-40$ μειώνονται κατά -8.187879 από τις γυναίκες που ανήκουν στην ηλικιακή κατηγορία $20-30$ και νοσούν από διαβήτη.

Ενώ για τις γυναίκες που ανήκουν στην ηλικιακή ομάδα 41-50 απο την παραμετρο age41-50 , $\hat{\beta}_8$, έχουμε

$$\hat{\beta}_0 + \hat{\beta}_8 = -8.90364 + 1.377709 = -7.525931$$

Δηλαδή τα log odds για τις γυναίκες που νοσούν απο διαβήτη και ανήκουν στην ηλικιακή κατηγορία 41-50 μείωνονται κατα -7.525931 απο τις γυναίκες που ανήκουν στην ηλικιακή κατηγορία 20-30 και νοσούν απο διαβήτη.

Για τις γυναίκες που ανήκουν στην ηλικιακή ομάδα 50 και πάνω απο την παραμετρο age50+ , $\hat{\beta}_9$, έχουμε

$$\hat{\beta}_0 + \hat{\beta}_9 = -8.90364 + 1.236659 = -7.666981$$

Δηλαδή τα log odds για τις γυναίκες που νοσούν απο διαβήτη και ανήκουν στην ηλικιακή κατηγορία 50 και πάνω μείωνονται κατα -7.666981 απο τις γυναίκες που ανήκουν στην ηλικιακή κατηγορία 20-30 και νοσούν απο διαβήτη.

Οσό αναφορα τον αριθμό κυήσεων έχουμε που είχε η κάθε γυναίκα βλέπουμε απο την μεταβλητή pregnant6-10 , $\hat{\beta}_{10}$, οτί

$$\hat{\beta}_0 + \hat{\beta}_{10} = -8.90364 - 0.150064 = -9.053704$$

Δηλαδή τα log odds για τις γυναίκες που νοσούν απο διαβήτη και ο αριθμός των κυήσεων τους είναι απο 6-10 μείωνονται κατα -7.666981 απο τις γυναίκες που ο αριθμός των κυήσεων τους είναι απο 0-5 και νοσούν απο διαβήτη.

Τέλος βλέπουμε οτι για τις γυναίκες με αριθμό κυήσεων πάνω απο 10 στην παραμετρο pregnant10+ , $\hat{\beta}_{11}$, οτί

$$\hat{\beta}_0 + \hat{\beta}_{11} = -8.90364 + 0.930171 = -7.973469$$

Δηλαδή τα log odds για τις γυναίκες που νοσούν απο διαβήτη και ο αριθμός των κυήσεων τους είναι απο 10 και πάνω μείωνονται κατα -7.973469 απο τις γυναίκες που ο αριθμός των κυήσεων τους είναι απο 0-5 και νοσούν απο διαβήτη.

Στην συνέχεια θα προχωρήσουμε στον έλεγχο αν το μοντέλο μας είναι προσαρμοζεί καλύτερα απο ένα μοντέλο με μονο ορό την σταθερα. Ο ελέγχος που κάνουμε για την διαφορά των αποκλίσεων των δυο μοντέλων είναι

$$D_n(\hat{\beta}) - D_m(\hat{\beta}) \rightarrow x_{df_{null}-df_m}^2$$

⁶Απο οπού και συμπαίρνουμε οτί υπαρχει σημαντική διαφορά ανάμεσα στο μοντέλο μας απο το μοντέλο απο την σταθερά.

Στον Πίνακα 3 παρουσιάζονται και τα αποτελέσματα για τους συντελεστές του μοντέλου. Απο την δευτερη στηλή εχουμε το τυπικό σφάλμα για κάθε παραμέτρο (Std.Error), ενώ στην τρίτη στηλή εχουμε τον στατιστικό έλεγχου Z για τον έλεγχο

$$H_0 : \beta_i = 0 \text{ vs } H_1 : \beta_i \neq 0$$

και στην τελευταία στήλη παρουσιάζονται οι αντίστοιχες p-τιμές για κάθε έλεγχο αντίστοιχα. Αυτο που παρατηρούμε είναι οτι αρκετές απο τις παραμέτρους έχουμε ισχύρες στατιστικές ενδείξεις για να αποδεχουμε την μηδενική υποθεση ,δηλαδή οτί $\beta_i = 0$,οπως για παράδειγμα η παραμτρεος pressure.Ενώ για άλλές παραμέτρους βλέπουμε οτί ορίακα μπορούμε να τις θεωρησουμε στατιστικά σημαντικες για το μοντέλο μας($p - value \approx a (= 0.05)$).

Για την επιλογή των μεταβλητών που θα χρησιμοποιήσουμε στο μοντέλο είναι το κριτήριο Bayesian Information Criterion (BIC).Ακόμα η μέθοδος επιλογής που θα χρησιμοποιήσουμε είναι σε κάθε βήμα να γίνεται και η εφαρμογή και των δυο τεχνικών,δηλαδή με την μεθοδο επιλογης προς τα εμπρός και πρός τα πίσω.Συμφώνα με το BIC το επίλεγμενο μοντέλο περιέχει τις εξής παραμετρους 7:

- glucose : Συγκέντρωση γλυκόζης στο πλάσμα
- mass :Δείκτης μάζας σώματος
- pedigree :Score για την πιθανότητα διαβήτη από κληρονομικά αίτια
- age : Ηλικία

⁶Για τις εντολές για τον ελεγχος προσαρμογής για το μοντέλο βλέπε στο Παράρτημα [Listing 7](#)

⁷Για τις εντολές για την επιλογή του μοντέλου βλέπε στο Παράρτημα [Listing 8](#)

Τα μοντέλο που προκύπτει σύμφωνα με το κριτήριο BIC είναι το ⁸:

$$\begin{aligned} \text{logit}\left(\frac{\hat{P}(\text{νοσει απο διαβήτη})}{\hat{P}(\text{δεν νοσει απο διαβήτη})}\right) = & -9.023004 + 0.037482 \times (\text{glucose}) + 0.074965 \times (\text{mass}) \\ & + 1.088329 \times (\text{pedigree}) + 0.798086 \times (\text{group_age31} - 40) \\ & + 1.553980 \times (\text{group_age41} - 50) + 1.291923 \times (\text{group_age50} +) \end{aligned}$$

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-9.023004	1.009742	-8.936	< 2e-16 ***
glucose	0.037482	0.005036	7.443	9.82e-14 ***
mass	0.074965	0.020504	3.656	0.000256 ***
pedigree	1.088329	0.427625	2.545	0.010926 *
age31-40	0.798086	0.342608	2.329	0.019836 *
age41-50	1.553980	0.416508	3.731	0.000191 ***
age50+	1.291923	0.526703	2.453	0.014173 *

Table 4: Μοντέλο λογιστικής παλινδρόμησης σύμφωνα με το κριτήριο BIC.

Απο τον Πίνακα 4 και την τελευταία στήλη βλέπουμε ότι όλες οι παραμετροι είναι στατιστικά σημαντικές σε επιπέδο σημαντικότητας 5%.

Επιπλέον έχουμε ότι η σταθερά $\hat{\beta}_0$ έχει την τιμή -9.023004 . Αυτό σημαίνει ότι το log odds για να νοσήσει κάποια απο διαβήτη η οποία ανήκει στην ηλικιακή ομάδα 20 – 30 είναι -9.023004 .

Για την παράμετρο glucose $\hat{\beta}_1$ βλέπουμε ότι η 1 μοναδιαία αύξηση της συγκέντρωσης γλύκοζης στο πλάσμα αυξάνει τα log odd κατά 0.037482 , όταν όλες οι υπολοιπές μεταβλητές παραμένουν σταθερές.

Αντιστοίχα για την παράμετρο mass $\hat{\beta}_2$ παρατηρούμε ότι 1 μοναδιαία αύξηση του δείκτη μάζας σώματος (βάρους σε kg/(ύψος σε m)²) αυξάνει τα log odds κατά 0.074965, όταν όλες οι υπολοιπές μεταβλητές είναι παραμενουν σταθερές.

⁸Για τις εντολές για την λογιστική παλινδρόμηση με τις μεταβλητές απο το BIC βλέπε στο Παράρτημα [Listing 9](#)

Ενώ για την παράμετρο pedigree $\hat{\beta}_3$ παρατηρούμε ότι 1 μοναδιαία αύξηση του Score για την πιθανότητα διαβήτη από κληρονομικά αίτια αυξάνει τα log odds κατά 1.088329 όταν όλες οι υπολοιπές μεταβήτες είναι παραμενουν σταθερές.

Ακόμα για τις γυναίκες που ανήκουν στην ηλικιακή ομάδα 31-40 από την παράμετρο age31-40, $\hat{\beta}_4$, έχουμε

$$\hat{\beta}_0 + \hat{\beta}_4 = -9.023004 + 0.798086 = -8.224918$$

Δηλαδή τα log odds για τις γυναίκες νοσούν από διαβήτη που ανήκουν στην ηλικιακή κατηγορία 31-40 μειώνονται κατά -8.187879 από τις γυναίκες που ανήκουν στην ηλικιακή κατηγορία 20-30 και νοσούν από διαβήτη.

Ενώ για τις γυναίκες που ανήκουν στην ηλικιακή ομάδα 41-50 από την παράμετρο age41-50, $\hat{\beta}_5$, έχουμε

$$\hat{\beta}_0 + \hat{\beta}_5 = -9.023004 + 1.553980 = -7.469024$$

Δηλαδή τα log odds για τις γυναίκες νοσούν από διαβήτη και ανήκουν στην ηλικιακή κατηγορία 41-50 μειώνονται κατά -8.224918 από τις γυναίκες που ανήκουν στην ηλικιακή κατηγορία 20-30 και νοσούν από διαβήτη.

Για τις γυναίκες που ανήκουν στην ηλικιακή ομάδα 50 και πάνω από την παράμετρο age50+, $\hat{\beta}_6$, έχουμε

$$\hat{\beta}_0 + \hat{\beta}_6 = -9.023004 + 1.291923 = -7.731081$$

Δηλαδή τα log odds για τις γυναίκες νοσούν από διαβήτη και ανήκουν στην ηλικιακή κατηγορία 50 και πάνω μειώνονται κατά -7.731081 από τις γυναίκες που ανήκουν στην ηλικιακή κατηγορία 20-30 και νοσούν από διαβήτη.

Για μια γυναίκα ηλικίας 35 χρονών έχουμε ότι η πιθανότητα να νοσήσει μια γυναίκα από διαβήτη με όλες τις μεταβλητές να είναι σταθερές (στις μέσες τιμές τους) είναι ⁹

$$\text{logit}\left(\frac{\hat{P}(\text{νοσει απο διαβήτη}|\text{age}=35)}{\hat{P}(\text{δεν νοσει απο διαβήτη}|\text{age}=35)}\right) = 0.3591437$$

Ενώ τα odds ratio για μια γυναίκα να νοσήσει από διαβήτη στην ηλικία των 35 είναι ¹⁰

$$\exp(0.3591437) = 1.432103$$

Ενώ για μια γυναίκα ηλικίας 45 χρονών έχουμε ότι η πιθανότητα να νοσήσει μια γυναίκα από διαβήτη με όλες τις μεταβλητές να είναι σταθερές (στις μέσες τιμές τους) είναι ¹¹

$$\text{logit}\left(\frac{\hat{P}(\text{νοσει απο διαβήτη}|\text{age}=45)}{\hat{P}(\text{δεν νοσει απο διαβήτη}|\text{age}=45)}\right) = 0.544088$$

Ενώ τα odds ratio για μια γυναίκα να νοσήσει από διαβήτη στην ηλικία των 45 είναι ¹²

$$\exp(0.544088) = 1.723036$$

Ακόμα τα odds ratio για την εμφάνιση διαβήτη μεταξύ της γυναίκας 35 με την εμφάνιση του διαβήτη για την γυναίκα 45 ετών είναι ¹³

$$\frac{\exp(0.3591437)}{\exp(0.544088)} = \frac{1.432103}{1.723036} = 0.8311509$$

Δηλαδή οι γυναικές που ανήκουν στην ηλικιακή ομάδα 45 ετών αυξάνονται τα odds ratio κατά 8.3% για την εμφάνιση του διαβήτη σε σχέση τις γυναίκες που ανήκουν στην ηλικιακή ομάδα 35 ετών.

⁹Για τις εντολές για την πρόβλεψη για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 35, αποτέλεσμα σε log odd βλέπε στο Παράρτημα [Listing 10](#)

¹⁰Για τις εντολές για τα Odds Ratio για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 35 βλέπε στο Παράρτημα [Listing 11](#)

¹¹Για τις εντολές για την Πρόβλεψη για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 45, αποτέλεσμα σε log odds βλέπε στο Παράρτημα [Listing 12](#)

¹²Για τις εντολές για Odds Ratio για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 45 βλέπε στο Παράρτημα [Listing 13](#)

¹³Για τις εντολές για τα Odds Ratio εμφάνισης διαβήτη μεταξύ της γυναίκας 35 με την εμφάνιση του διαβήτη για την γυναίκα 45 ετών βλέπε στο Παράρτημα [Listing 14](#)

Στην συνέχεια θα προχωρήσουμε στην δημιουργία ενός Confusion matrix. Ο πίνακας αυτός στην ουσία μας επιτρέπει να οπτικοποιήσουμε την αποδόση ενός αλγορίθμου, στην περιγραφή μας του λογιστικού μοντέλου που χρησιμοποιήσαμε παραπάνω.

¹⁴Αυτό που προσπαθούμε να κάνουμε είναι να πρόβλεψουμε πότε ο αλγόριθμος για την δίτιμη παραμετρο του διαβήτη θα δώσει ως αποτέλεσμα $y=0$ ή $y=1$. Για την προβλέψη ότι η γυναίκα θα νοσήσει από διαβήτη ($y=1$) υποθέτουμε ότι η πιθανότητα θα είναι μεγαλύτερη από ένα όριο $\hat{\pi}_0$ που θα έχουμε ορίσει εμείς δηλαδή $\hat{\pi}_i > \hat{\pi}_0$ και για την πρόβλεψή ότι η γυναίκα δεν θα νοσήσει ($y=0$) υποθέτουμε ότι $\hat{\pi}_i < \hat{\pi}_0$. Εμείς για $\hat{\pi}_0$ θα το ορίσουμε ισό με $\hat{\pi}_0 = 0.5$

¹⁵Τα αποτελέσματα που παίρνουμε παρουσιάζονται στον παρακάτω πίνακα:

Prediction, $\pi_0 = 0.5$		
Actual	neg	pos
neg	111	30
pos	15	40

Table 5: Confusion matrix.

¹⁶Από όπου και συμπαίρνουμε ότι το ποσοστό ορθής ταξινόμησης εκείνων που νοσήσαν από διαβήτη είναι 57.15%.

$$\text{Sensitivity(Ευαισθησία)} = P(\hat{y} = 1|y = 1) = \frac{40}{40 + 30} = 57.14286\%$$

Ενώ το ποσοστό ορθής ταξινόμησης εκείνων που δεν νοσήσαν είναι 88.1%.

$$\text{Specificity(Ειδικότητα)} = P(\hat{y} = 0|y = 0) = 1 - \frac{15}{111 + 15} = 88.09524\%$$

Αυτό που βλέπουμε είναι ότι η ειδικότητα είναι αρκετά μεγάλη.

¹⁴Για τις εντολές Λογιστική παλινδρόμηση για την πρόβλεψη της μεταβλητής diabete βλέπε στο Παράρτημα [Listing 15](#)

¹⁵Για τις εντολές Δημιουργία Confusion Matrix βλέπε στο Παράρτημα [Listing 16](#)

¹⁶Υπολογισμός Ευαισθησία (Sensitivity) και Ειδικότητα (Specificity) βλέπε στο Παράρτημα [Listing 17](#)

Ακόμα υπάρχουν τα δεδομένα για τρεις γυναίκες με συγκεκριμένα χαρακτηριστικά όπως παρουσιάζονται στον Πίνακα 5.

Γυναίκες	1	2	3
Pregnant	3	2	0
Glocose	95	55	80
Pressure	50	62	70
Triceps	21	17	30
Insulin	70	100	110
Mass	30	33	35
Pedigree	0.5	0.7	0.9
Age	50	40	35

Table 6: Δεδομένα για τα 95% διαστήματα εμπιστοσύνης

Για τα παραπάνω δεδομένα θα χρειαστεί να δημιουργήσουμε σημειακά 95% διαστήματα πρόβλεψης για την πιθανότητα εμφάνισης διαβήτη στις τρεις αυτές γυναίκες ¹⁷.

Γυναίκες	Πιθανότητα Πρόβλεψης	Ανώ όριο δ.ε	Κάτω όριο δ.ε
1	0.2469113	0.4225588	0.1280815
2	0.1369591	0.2534081	0.06907139
3	0.06238096	0.1475301	0.02493914

Table 7: 95% σημειακά διαστήματα πρόβλεψης

Για την Γυναίκα 1 με τα χαρακτηριστικά όπως δίνονται στον πίνακα 6 έχουμε ότι η πιθανότητα πρόβλεψης να νοσήσει είναι 0.24 και με σημειακά διαστήματα (0.12,0.42). Ενώ για την δεύτερη γυναίκα η πιθανότητα πρόβλεψης να νοσήσει είναι 0.13 και με σημειακά διαστήματα (0.07,0.25). Τέλος για την τρίτη γυναίκα η πιθανότητα πρόβλεψης να νοσήσει είναι 0.06 και με σημειακά διαστήματα (0.02,0.14).

¹⁷Για τις εντολές για τα σημειακά δ.ε. βλέπε στο Παράρτημα για την πρώτη γυναίκα [Listing 18](#), για την 2 γυναίκα [Listing 19](#) και για την τρίτη γυναίκα [Listing 20](#)

ΠΑΡΑΡΤΗΜΑ

```
1
2 rm(list=ls(all=TRUE))
3
4 install.packages("mlbench")
5 library("mlbench")
6
7 data(PimaIndiansDiabetes2)
8 data<-PimaIndiansDiabetes2
```

Listing 1: Εγκατάσταση βιβλιοθηκών και δεδομένων

ΚΩΔΙΚΑΣ ΕΡΩΤΗΜΑΤΟΣ 1

```
1 sum(is.na(data))
2 # [1] 652
3
4 new.data<-na.omit(data)
5 sum(is.na(new.data))
6 # [1] 0
```

Listing 2: Ελέγχος ελλειπών τιμών και αφαίρεσή τους

```
1 str(new.data)
2 #'data.frame': 392 obs. of 9 variables:
3 # $ pregnant: num 1 0 3 2 1 5 0 1 1 3 ...
4 # $ glucose : num 89 137 78 197 189 166 118 103 115 126 ...
5 # $ pressure: num 66 40 50 70 60 72 84 30 70 88 ...
6 # $ triceps : num 23 35 32 45 23 19 47 38 30 41 ...
7 # $ insulin : num 94 168 88 543 846 175 230 83 96 235 ...
8 # $ mass : num 28.1 43.1 31 30.5 30.1 25.8 45.8 43.3 34.6
9 # $ pedigree: num 0.167 2.288 0.248 0.158 0.398 ...
10 # $ age : num 21 33 26 53 59 51 31 33 32 27 ...
11 # $ diabetes: Factor w/ 2 levels "neg","pos": 1 2 2 2 2 2 1 2
12 # - attr(*, "na.action")= 'omit' Named int [1:376] 1 2 3 6 8 10
13 # ... attr(*, "names")= chr [1:376] "1" "2" "3" "6" ...
```

Listing 3: Ελέγχος του τύπου των μεταβλητών

```

1
2 new.data$age<-cut(new.data$age, breaks=c(20,31,41,50,100),
3     labels=c("20-30","31-40","41-50","50+"))
4
5 levels(new.data$age)
6 # [1] "20-30" "31-40" "41-50" "50+"
7
8 new.data$pregnant<-cut(new.data$pregnant, breaks=c(-1,5,10,17),
9     labels=c("0-5","6-10","10+"))
10
11 levels(new.data$pregnant)
12 > levels(new.data$pregnant)
13 # [1] "0-5" "6-10" "10+"

```

Listing 4: Ομαδοποίηση των μεταβλητής Age και Pregnant

```

1 str(new.data)
2 # 'data.frame':   392 obs. of  9 variables:
3 #  $ pregnant: Factor w/ 3 levels "0-5","6-10","10+": 1 1 1 1 1
4 #    1 1 1 1 1 ...
5 #  $ glucose : num  89 137 78 197 189 166 118 103 115 126 ...
6 #  $ pressure: num  66 40 50 70 60 72 84 30 70 88 ...
7 #  $ triceps : num  23 35 32 45 23 19 47 38 30 41 ...
8 #  $ insulin : num  94 168 88 543 846 175 230 83 96 235 ...
9 #  $ mass     : num  28.1 43.1 31 30.5 30.1 25.8 45.8 43.3 34.6
10 #    39.3 ...
11 #  $ pedigree: num  0.167 2.288 0.248 0.158 0.398 ...
12 #  $ age      : Factor w/ 4 levels "20-30","31-40",...: 1 2 1 4 4
13 #    4 1 2 2 1 ...
14 #  $ diabetes: Factor w/ 2 levels "neg","pos": 1 2 2 2 2 2 1 2
15 #    1 ...
16 # - attr(*, "na.action")= 'omit' Named int [1:376] 1 2 3 6 8 10
17 #    11 12 13 16 ...
18 # ..- attr(*, "names")= chr [1:376] "1" "2" "3" "6" ...

```

Listing 5: Ελέγχος του τύπου των μεταβλητών για το τελικό αρχείο.

```

1
2 mylogit <- glm(formula = diabetes~glucose+ pressure+ triceps+
3   insulin+mass+pedigree+age+pregnant,
4   family = "binomial", data = new.data)
5   summary(mylogit)
6
7 #Call:
8 #glm(formula = diabetes ~ glucose + pressure + triceps +
9   insulin +
10   mass + pedigree + age + pregnant, family = "binomial",
11   data = new.data)
12 #
13 #Deviance Residuals:
14 #      Min       1Q   Median       3Q      Max
15 #-2.6284   -0.6525   -0.3584    0.5869    2.6191
16 #
17 #Coefficients:
18 #      Estimate Std. Error z value Pr(>|z|)
19 #(Intercept)  -8.990364   1.170461  -7.681 1.58e-14 ***
20 #glucose       0.039232   0.005865   6.689 2.24e-11 ***
21 #pressure     -0.002747   0.011892  -0.231  0.8173
22 #triceps       0.015717   0.017320   0.907  0.3642
23 #insulin      -0.000692   0.001351  -0.512  0.6085
24 #mass          0.063485   0.027425   2.315  0.0206 *
25 #pedigree      1.051124   0.434989   2.416  0.0157 *
26 #age31-40      0.802485   0.395535   2.029  0.0425 *
27 #age41-50      1.377709   0.536629   2.567  0.0102 *
28 #age50+        1.236659   0.633326   1.953  0.0509 .
29 #pregnant6-10  -0.150064   0.424600  -0.353  0.7238
30 #pregnant10+   0.930171   0.779488   1.193  0.2327
31 #---
32 #Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
33 #
34 #(Dispersion parameter for binomial family taken to be 1)
35 #
36 #      Null deviance: 498.10  on 391  degrees of freedom
37 #Residual deviance: 340.65  on 380  degrees of freedom
38 #AIC: 364.65
39 #
40 #Number of Fisher Scoring iterations: 5

```

Listing 6: Λογιστική Παλινδρόμηση με όλες τις μεταβλητές.

```

1 with(mylogit, pchisq(null.deviance - deviance,
2 df.null - df.residual, lower.tail = FALSE))
3
4 #[1] 4.462115e-28

```

Listing 7: Έλεγχος προσαρμογής για το μοντέλο.

ΚΩΔΙΚΑΣ ΕΡΩΤΗΜΑΤΟΣ 2

```

1
2 n <- dim(new.data)[1]
3 model_bic <- step(mylogit, trace=TRUE, direction = 'both', k =
  log(n))
4
5 #Start:  AIC=412.31
6 #diabetes ~ glucose + pressure + triceps + insulin + mass +
  pedigree +
7   age + pregnant
8 #
9 #           Df Deviance    AIC
10 #- pregnant  2   342.90 402.61
11 #- age       3   349.12 402.86
12 #- pressure  1   340.70 406.39
13 #- insulin  1   340.91 406.60
14 #- triceps  1   341.47 407.16
15 #- mass     1   346.14 411.83
16 #<none>      340.65 412.31
17 #- pedigree 1   346.80 412.49
18 #- glucose  1   395.09 460.78
19 #
20 #Step:  AIC=402.61
21 #diabetes ~ glucose + pressure + triceps + insulin + mass +
  pedigree +
22 #   age
23 #
24 #           Df Deviance    AIC
25 #- pressure  1   342.92 396.66
26 #- insulin  1   343.29 397.03
27 #- triceps  1   343.62 397.36
28 #- age      3   360.41 402.20
29 #- mass     1   348.62 402.36
30 #<none>      342.90 402.61
31 #- pedigree 1   349.52 403.26
32 #+ pregnant  2   340.65 412.31
33 #- glucose  1   397.69 451.43
34 #
35 #
36 #
37 #
38 #
39 #
40 #
41 #
42 #

```



```

43 #
44 #
45 #Step:   AIC=396.66
46 #diabetes ~ glucose + triceps + insulin + mass + pedigree + age
47 #
48 #           Df Deviance    AIC
49 #- insulin   1   343.30 391.07
50 #- triceps   1   343.64 391.41
51 #<none>      342.92 396.66
52 #- mass      1   348.91 396.68
53 #- age       3   361.34 397.17
54 #- pedigree  1   349.60 397.37
55 #+ pressure  1   342.90 402.61
56 #+ pregnant  2   340.70 406.39
57 #- glucose   1   398.23 446.00
58 #
59 #Step:   AIC=391.07
60 #diabetes ~ glucose + triceps + mass + pedigree + age
61 #
62 #           Df Deviance    AIC
63 #- triceps   1   344.05 385.85
64 #- mass      1   348.96 390.76
65 #<none>      343.30 391.07
66 #- pedigree  1   349.80 391.59
67 #- age       3   361.77 391.63
68 #+ insulin   1   342.92 396.66
69 #+ pressure  1   343.29 397.03
70 #+ pregnant  2   340.95 400.66
71 #- glucose   1   412.61 454.40
72 #
73 #Step:   AIC=385.85
74 #diabetes ~ glucose + mass + pedigree + age
75 #
76 #           Df Deviance    AIC
77 #<none>      344.05 385.85
78 #- pedigree  1   350.93 386.76
79 #- age       3   363.70 387.59
80 #+ triceps   1   343.30 391.07
81 #+ insulin   1   343.64 391.41
82 #+ pressure  1   344.04 391.81
83 #- mass      1   358.36 394.19
84 #+ pregnant  2   341.80 395.54
85 #- glucose   1   413.84 449.67

```

Listing 8: Επιλογή μοντέλου με το κριτήριο BIC

```

1
2 mylogit1 <- glm(formula =diabetes ~ glucose + mass + pedigree +
3               age,
4               family = "binomial", data = new.data)
5
6 #Call:
7 #glm(formula = diabetes ~ glucose + mass + pedigree + age,
8 #     family = "binomial",
9 #     data = new.data)
10 #
11 #Deviance Residuals:
12 #    Min       1Q   Median       3Q      Max
13 # -2.7102  -0.6399  -0.3706   0.6409   2.6590
14 #
15 #Coefficients:
16 #              Estimate Std. Error z value Pr(>|z|)
17 #(Intercept)  -9.023004   1.009742  -8.936  < 2e-16 ***
18 #glucose       0.037482   0.005036   7.443 9.82e-14 ***
19 #mass          0.074965   0.020504   3.656 0.000256 ***
20 #pedigree      1.088329   0.427625   2.545 0.010926 *
21 #age31-40      0.798086   0.342608   2.329 0.019836 *
22 #age41-50      1.553980   0.416508   3.731 0.000191 ***
23 #age50+        1.291923   0.526703   2.453 0.014173 *
24 #---
25 #Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
26 #
27 # (Dispersion parameter for binomial family taken to be 1)
28 #
29 #    Null deviance: 498.10  on 391  degrees of freedom
30 #Residual deviance: 344.05  on 385  degrees of freedom
31 #AIC: 358.05
32 #
33 #Number of Fisher Scoring iterations: 5

```

Listing 9: Λογιστική Παλινδρόμηση για τις μεταβλήτες σύμφωνα με το κριτήριο BIC .

ΚΩΔΙΚΑΣ ΕΡΩΤΗΜΑΤΟΣ 3

```

1 as.numeric(predict(mylogit1 ,
2 newdata=data.frame(
3 age=factor(factor('31-40', levels=c('20-30', '31-40', '41-50', '
    50+'))),
4 glucose =mean(new.data$glucose ),
5 pressure=mean(new.data$pressure),
6 triceps =mean(new.data$triceps ),
7 insulin =mean(new.data$insulin ),
8 mass=mean(new.data$mass),
9 pedigree =mean(new.data$pedigree )),type="response"))
10 # [1] 0.3591437

```

Listing 10: Πρόβλεψη για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 35, αποτέλεσμα σε log odds

```

1 as.numeric(predict(mylogit1 ,
2 newdata=data.frame(
3 age=factor(factor('41-50', levels=c('20-30', '31-40', '41-50', '
    50+'))),
4 glucose =mean(new.data$glucose ),
5 pressure=mean(new.data$pressure),
6 triceps =mean(new.data$triceps ),
7 insulin =mean(new.data$insulin ),
8 mass=mean(new.data$mass),
9 pedigree =mean(new.data$pedigree )),type="response"))
10 # [1] 0.544088

```

Listing 11: Πρόβλεψη για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 45, αποτέλεσμα σε log odds

```

1 exp(0.3591437)
2 # [1] 1.432103

```

Listing 12: Odds Ratio για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 35

```

1 exp(0.544088)
2 # [1] 1.723036

```

Listing 13: Odds Ratio για την εμφάνιση διαβήτη για μια γυναίκα ηλικίας 45

```

1 exp(0.3591437)/exp(0.544088)
2 # [1] 0.8311506

```

Listing 14: Odds Ratio εμφάνισης διαβήτη μεταξύ της γυναίκας 35 με την εμφάνιση του διαβήτη για την γυναίκα 45 ετών

ΚΩΔΙΚΑΣ ΕΡΩΤΗΜΑΤΟΣ 4

```

1 n <- nrow(new.data)
2 train <- sample(n, n/2)
3 glm.fit <- glm(formula =diabetes ~ glucose + mass + pedigree +
4   age,
5               data = new.data,subset = train, family =
6   binomial)
7 glm.probs <- predict(glm.fit, new.data[-train, ], type = "
8   response")
9 glm.pred <- rep("neg", n/2)
10 glm.pred[glm.probs>0.5] <- "pos"

```

Listing 15: Λογιστική παλινδρόμηση για την πρόβλεψη της μεταβλητής diabetes

```

1 table(glm.pred)
2 # glm.pred
3 # neg pos
4 # 141 55
5
6 confTab <- table(glm.pred, diabetes [-train])
7 confTab
8
9
10 #glm.pred neg pos
11 # neg 111 30
12 # pos 15 40

```

Listing 16: Δημιουργία Confusion Matric

```

1
2 # Sensitivity
3 ((confTab[2,2])/(confTab[1,2]+confTab[2,2]))*100
4 #[1] 57.14286
5
6 #Specificity
7 (1-(confTab[2,1])/(confTab[1,1]+confTab[2,1]))*100
8 #[1] 88.09524

```

Listing 17: Υπολογισμός Ευαισθησία (Sensitivity) και Ειδικότητα (Specificity)

ΚΩΔΙΚΑΣ ΕΡΩΤΗΜΑΤΟΣ 5

```

1 newdata1<- with(new.data,
2   data.frame(pregnant=factor("0-5",levels=c( "0-5" ,"6-10", "
3     10+")),glucose=95,presure=50,
4   triceps=21,insulin=70,mass=30,pedigree=0.5,
5   age=factor('41-50',levels=c('20-30','31-40','41-50','50+'))))
6
7
8 data1<- cbind(newdata1,predict(mylogit1, newdata = newdata1,
9   type="link", se=TRUE))
10
11 newdata4<- within(data1,{
12   PredictedProb <- plogis(fit)
13   LL <- plogis(fit - (1.96 * se.fit))
14   UL <- plogis(fit + (1.96 * se.fit))
15 })
16 (newdata4)
17 #   pregnant glucose pressure triceps insulin mass pedigree
18 #1      0-5      95       50      21      70   30      0.5
19
20 #   age      fit      se.fit  residual.scale      UL
21 #   LL
22 # 1  41-50  -1.115154  0.4096304      1  0.4225588
23 #1      0.1280815
24
25 #   PredictedProb
26 #1      0.2469113

```

Listing 18: Προβλέψη για την Γυναίκα 1 και δ.ε 95%

```

1 newdata2<- with(new.data,
2   data.frame(pregnant=factor("0-5",levels=c( "0-5" ,"6-10", "
3     10+")),glucose= 80,presure=62,
4   triceps=30,insulin=100,mass=35,pedigree=0.7,
5   age=factor('31-40',levels=c('20-30','31-40','41-50','50+'))))
6
7 data2<- cbind(newdata2,predict(mylogit1, newdata = newdata2,
8   type="link", se=TRUE))
9
10 newdata5<- within(data2,{
11   PredictedProb <- plogis(fit)
12   LL <- plogis(fit - (1.96 * se.fit))
13   UL <- plogis(fit + (1.96 * se.fit))
14 })
15 (newdata5)
16 #   pregnant  glucose  pressure  triceps  insulin  mass
17 #   0-5        80        62        30        100    35    0.7
18
19 #   age      fit      se.fit  residual.scale
20 # 31-40 -1.84078  0.387889          1
21
22 #   UL      LL  PredictedProb
23 # 0.2534081 0.06907139    0.1369591

```

Listing 19: Προβλέψη για την Γυναίκα 2 και 95% δ.ε

```

1 newdata3<- with(new.data,
2   data.frame(pregnant=factor("0-5",levels=c( "0-5" ,"6-10", "
3     10+")),glucose= 55,presure=70,
4     triceps=17,insulin= 110,mass=33,pedigree=0.9,
5     age=factor('31-40',levels=c('20-30','31-40','41-50','50+'))))
6
7   data3<- cbind(newdata3,predict(mylogit1, newdata = newdata3,
8     type="link", se=TRUE))
9
10  newdata6<- within(data3,{
11    PredictedProb <- plogis(fit)
12    LL <- plogis(fit - (1.96 * se.fit))
13    UL <- plogis(fit + (1.96 * se.fit))
14  })
15  (newdata6)
16
17  #   pregnant  glucose  pressure  triceps  insulin  mass
18    pedigree
19  # 1      0-5      55      70      17      110      33      0.9
20
21  # age      fit      se.fit  residual.scale
22  # 31-40 -2.710084 0.4877437      1
23
24  # UL      LL      PredictedProb
25  #0.1475301 0.02493914 0.06238096

```

Listing 20: Προβλέψη για την Γυναίκα 3 και 95% δ.ε