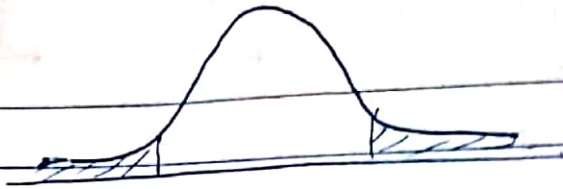


Οκτώβριος 2021



Ζήτημα 1

$y$  = ποσότητα που πωλείται,  $x_1$ : ηλικία του πελάτη

$x_2$ : δίδαγμα που πωλείται  
 $x_2 = 0$ : χαμηλή  
 $x_2 = 1$ : υψηλή

$$\hat{y} = 11.4 + 43.5 x_1 - 25.9 x_2 - 20.6 x_3$$

$\downarrow \hat{\beta}_1$        $\downarrow \hat{\beta}_2$        $\downarrow \hat{\beta}_3$

$$t_i^* = \frac{\hat{\beta}_i}{se(\hat{\beta}_i)} \sim t_{n-p} \quad \leftarrow \quad t_{n-p} = t_{14-4} = t_{10}$$

SST έχει 13 df το οποίο είναι  $n-1$  με βάση τη θεωρία, άρα  $n-1=13 \Rightarrow n=14$  παρατηρήσεις.  
 $k=3$  μεταβλητές,  $p=k+1=4$  παράμετροι

$$R^2 = 1 - \frac{n-1}{n-p} (1 - r^2) \approx 95.84\%$$

$$SSE = SST - SSR = 11,206$$

|       |       |                |
|-------|-------|----------------|
| $x_1$ |       |                |
| $x_2$ | -0.98 | 0.544          |
| $x_3$ | -1.92 | 0.0006 < 0.001 |

Μοντέλο  $y \sim x_1, x_2$  άρα  $k=2$ ,  $p=k+1=3$ ,  $df = n-p = 11$ .  
 $t$ -test με  $df=11$

|       |       |                              |
|-------|-------|------------------------------|
|       | $T$   | $P$                          |
| $x_1$ | 8.62  | 0.000                        |
| $x_2$ | -6.66 | $3.58 \cdot 10^{-5} < 0.001$ |

Ανομοιότητα για το μοντέλο  $y \sim x_1, x_2$

$$H_0: \beta_0 = \beta_1 = \beta_2 = 0$$

$$H_1: \text{έστω ένας συντελεστής} \neq 0$$

$$df_1 = 2, df_2 = 11$$

$$F^* = \frac{MSR}{MSE} = 44.68$$

$$p\text{-value} = 5.238 \cdot 10^{-6} < 0.001$$

$(F_{test, df_1, df_2}) = F_{2, 11}$

Μοντέλο  $y \sim x_1$ :  $\hat{y} = 186 + 22.8 x_1$

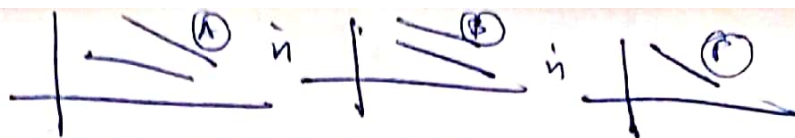
$$R^2_{pred} = 1 - \frac{PRESS}{SST} = 1 - \frac{255010}{349942} \approx 0.27$$

Έχουμε 1 μεταβλητή:  $df_1 = 1$ ,  $df_2 = n-p = 14 - (1+1) = 12$

$$F^* = 9.78 \rightarrow p\text{-value} = 0.0087$$

Τέλος επωλεστων συμπεριληψεις.

Επειραρχε αν



• Αν  $\beta_3 \neq 0$  τότε έχουμε (A).

Από την ανάλυση του μοντέλου  $y \sim x_1, x_2, x_3$  έχουμε για τη  $x_3$  ότι το  $p\text{-value} < 0.001$  άρα ~~ο~~ στατιστικά σημαντικός ο έλεγχος και  $\beta_3 \neq 0$ .

Άρα θα προσαρμόσουμε 2 διαφορετικές ευθείες.

$F_{1,10}$

Εναλλακτικός τρόπος:

$H_0: \beta_3 = 0$  ;  $M_0: y \sim x_1, x_2$

$H_1: \beta_3 \neq 0$  ;  $M_1: y \sim x_1, x_2, x_3$

$$F^* = \frac{(SSE_{M_0} - SSE_{M_1}) / 1}{SSE_{M_1} / (14 - 3 - 1)} = \frac{38358 - 11,206}{11,206 / 10} = 24.22988$$

$$F\text{-test}_{1,10} = 6.02 \cdot 10^{-4} < 0.001 \text{ άρα στατιστικά σημαντικός ο έλεγχος} \\ \text{συν. απορρίπτω } H_0 \text{ και δεχόμαι } H_1: \beta_3 \neq 0.$$

Άρα πρέπει να προσαρμόσων δύο ευθείες:

$$E(y_x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Προσαρμοσμένες ευθείες:

$$(1) \quad E(y_x) = \beta_0 + \beta_1 x_1$$

$$(2) \quad E(y_x) = \beta_0 + \beta_1 x_1 + \beta_2 + \beta_3 x_1 = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) x_1, \quad x_2 = 0$$

Ερμηνεία:

Αν η διαθεσιμότητα νερού είναι χαμηλή, τότε ο ρυθμός με τον οποίο αυξάνει ο ρυθμός φωτισμού ως προς την ηλεκτρική απεικόνιση είναι μεγαλύτερος σε σχέση με το αν η διαθεσιμότητα νερού ήταν χαμηλή.



## Ζήτημα 2

$$y \sim x_1, x_2, x_3, x_4, x_5$$

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5$$

## A ανάλυση

Πάνω στο πρόβλημα Συστήματος Κερών.

$$n = 51 \quad (df \text{ of SST} + 1 = 50 + 1 = 51) \quad \text{Παρατηρήσεις}$$

$$k = 5 \text{ μεταβλητές}, \quad p = k + 1 = 6 \text{ παράμετροι}$$

$$df = n - p = 50 - 45$$

$$x_3: t^* = 1.18, \quad p\text{-value} = 0.145$$

$$x_5: t^* = -1.98, \quad p\text{-value} = 0.054$$

$$R^2 = \frac{SSR}{SST} = 0.9724 \approx 97.24\%$$

$$F^* = 317.16 \rightarrow F_{5, df} = F_{5, 45} (F^*), \quad p\text{-value} = 6.928 \cdot 10^{-39}$$

μεταβλητές  $n-p=45$

Επιτυγχάνεται μεγάλο  $R^2$  που σημαίνει ότι το πρόβλημα ερμηνεύει κανονιστικά τα δεδομένα μας. Ωστόσο, από τα  $t$ -tests και τα σχετικά  $p$ -values φαίνεται κάποιες μεταβλητές ίσως να μην είναι σημαντικές για το πρόβλημα και άρα θα πρέπει να εξετάσει η αρίθμηση μοντέλων με λιγότερες μεταβλητές.

Επίσης, σε φαίνεται να υπάρχει πολλαπλότητα στις μεταβλητές  $x_1, x_2, x_3, x_4, x_5$  με βάση το κριτήριο VIF (από VIF values  $< 5$ )

## B ανάλυση

(i) Επιλέγουμε 5 και 7 μοντέλα λόγω των απλών τους.

(ii) Συγκρίνουμε  $M_5$  και  $M_7$   $\rightarrow k=3, p=4$

$$H_0: M_5: x_2, x_4, x_5 \quad (\text{Reduced}) \quad H_1: M_7: x_2, x_3, x_4, x_5 \quad (\text{Full model})$$

$$df_R - df_F = 1, \quad df_F = n - p = 51 - 5 = 46$$

$$F^* = \frac{(SSE_R - SSE_F) / 1}{SSE_F / df_F} \sim F_{(df_R - df_F), df_F} = F_{1, 46}$$

$$S = \frac{\sqrt{SSE}}{\sqrt{n-k-1}} \Rightarrow SSE = (n-k-1)S^2, \quad SSE_R = (51-3-1) \cdot 473969 = 24421859$$

$$SSE_F = (51-4-1) \cdot 467719 = 21215500, \quad F^* = 2.26 \rightarrow p\text{-value} = 0.139$$

Από έλεγχο  $F$  στατιστική ανεπαρκής που σημαίνει ότι δεν μπορούμε να απορρίψουμε  $H_0$  οπότε κρατάμε το  $M_5: y \sim x_2, x_4, x_5$

(iii) residuals vs norm quantiles

Με βάση το QQ plot, τα υποδομένα φάνε να είναι κοντά σε ευθεία γραμμή να ακολουθούν ικανοποιητικά κανονική κατανομή με εξαίρεση το σημείο 33 που μοιάζει με outlier.

residuals vs fitted

Φάνε να υπάρχει ένα trend στα residuals, για τιμές μεγαλύτερες των fitted values. Πιθανό για μικρές τιμές φάνε να μην είναι ενός band όπως και θέλουμε.

hat values vs index

$$\frac{2p}{n} = \frac{2 \cdot 3}{51} = 0.1176$$

Για  $h_{ii} > \frac{2p}{n}$  σημαίνει επιρροή  $\rightarrow$  έχω κάποιους

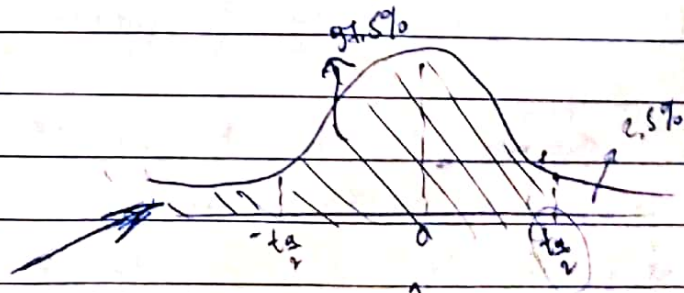
dffits vs index

$$2\sqrt{\frac{p}{n}} = 2\sqrt{\frac{3}{51}} \approx 0.485 \quad \text{αν } dffits_i > 2\sqrt{\frac{p}{n}} \text{ τότε : παρατηρητής ίσως σημείο επιρροής}$$

παρατηρητές 5 και 41 σημεία επιρροής.

partial residuals plot (vs  $x_2$  and  $x_4$ )

Βλέπουμε ότι και για τις 2 μεταβλητές  $x_2$  και  $x_4$  τα σημεία συρμύζιουν κατά μια ευθεία. Ειδικά για την  $x_4$ , ωστόσο και για  $x_2$  υπάρχει trend. Καταλήγουμε ότι χρειάζεται και τις δύο μεταβλητές στο μοντέλο.



(iv) Μοντέλο 3:  $y \sim x_2, x_4$

$\alpha: 95\%: (-920,677, \dots)$  ως πρόβλεψη μιας νέας παρατηρήσεως  $\hat{y}_{x_0}$ , εάν  $\hat{y}_{x_0} = 110,976$  και  $x_0'(X'X)^{-1}x_0 = 0,0377$

Από εμπειρίες έχουμε General Linear Model:  $\hat{y}_{x_0} - y_{x_0} \sim N(0, \sigma^2(x_0'(X'X)^{-1}x_0 + 1))$  και  $y_{x_0} \in \hat{y}_{x_0} \pm t_{\alpha/2} \cdot S \cdot (x_0'(X'X)^{-1}x_0 + 1)^{1/2}$ ,  $S = 499,000$  από πίνακα.  $t_{\alpha/2} \approx 2,02$  Το ίδιο βρίσκω για  $qt(0,975,48) \approx 2,01$  και με υπολογίζω.



**Ζήτηση 3**

Poisson:  $f(y) = \frac{e^{-\mu} \cdot \mu^y}{y!}$ ,  $y=0,1,2,\dots$ ;  $\ln f(y) = -\mu + y \ln \mu - \ln y!$   
 Γινόμενο ολικών:  $g(y|x) = \ln y = b \cdot x$   
 ελαττωμένη deviance:  $D_n(\hat{\theta}) = -2(\hat{l}_n - \hat{l}_{rop}) = 2 \sum_{i=1}^n [y_i \ln(\frac{y_i}{\hat{\mu}_i})]$ ,  $AIC = -2 \hat{l}_n + 2$   
 $Y \rightarrow$  αριθμός περιστατικών ανά μονάδα,  $X_1$  επιρροή,  $X_2$  αλληλεπίδραση,  $X_3$  επηρεάζει τον κωδικό

**(i) Συστάσεις:**

MONTEO 3:  $X_3$ :  $\frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} \rightarrow z\text{-test}$   $p\text{-value} = 2 \cdot \text{pnorm}(\text{abs}(-2.887), \text{mean}=0, \text{sd}=1, \text{lower.tail}=\text{FALSE}) = 0.0039$

MONTEO 2:  $X_2$ :  $z\text{-test}$   $p\text{-value} = 0.4526$

MONTEO 1:  $X_2$ :  $p\text{-value} = 0.3462$

MONTEO 1:  $\hat{l}_1 = j$ ,  $AIC = j$

για  $M_0$ :  $D_0(\hat{\theta}) = -2(\hat{l}_0 - \hat{l}_{rop}) \Rightarrow D_0(\hat{\theta}) = -2\hat{l}_0 + 2\hat{l}_{rop} \Rightarrow \hat{l}_{rop} = \frac{D_0(\hat{\theta})}{2} + \hat{l}_0$   
 $\Rightarrow \hat{l}_{rop} = -80.415$

για  $M_1$ :  $D_1(\hat{\theta}) = -2(\hat{l}_1 - \hat{l}_{rop}) \Rightarrow \hat{l}_1 = \hat{l}_{rop} - \frac{D_1(\hat{\theta})}{2} = -1835.315$

$AIC_1 = -2\hat{l}_1 + 2 \cdot d = 2 \cdot (-1835.315) + 2 \cdot 2 = 3674.63$

για  $M_0$ :  $D_0(\hat{\theta}) = -2(\hat{l}_0 - \hat{l}_{rop})$   $AIC_0 = -2\hat{l}_0 + 2 \cdot d = -2 \cdot (-1835.78) + 2 = 3673.56$

$D(\hat{\theta}) \sim \chi^2_{n-p}$   
 $M_1$ :  $D_1(\hat{\theta}) - D_0(\hat{\theta}) \sim \chi^2_{q=1}$   $\rightarrow p\text{-value} = 0.343$   $\Rightarrow$   $\text{FALSE}$   
 (66 given for  $M_0$ )  $H_0: M_0, H_1: M_1$   
 Αφαιρούμε τον αριθμό ελαττωμένης  $\rightarrow$  δεσφύρι το  $M_0$  (μικρότερο)

$M_2$  (66 given for  $M_1$ ):  $H_0: M_1, H_1: M_2$  ( $M_1 \leq M_2$ )  
 $D_2(\hat{\theta}) - D_1(\hat{\theta}) \sim \chi^2_{q=1} = \chi^2$   $\rightarrow p\text{-value} = 0.1138$   
 Αφαιρούμε τον αριθμό ελαττωμένης  $\Rightarrow$  δεσφύρι να απορριψω  $H_0$  και δεσφύρι το  $M_1$

$$R^2_D = \left(1 - \frac{D_2(\hat{\theta})}{D_0}\right) \cdot 100\% = \left(1 - \frac{3504.3}{3510.73}\right) \cdot 100\% = 0.0977\%$$

**(ii) Βρίσκω  $M_0$  και  $M_3$  πάλι ελαττωμένης με BIC Deviance.**

$H_0: M_0 \rightarrow$   $H_1: M_3 \rightarrow$   
 $D_3(\hat{\theta}) - D_0(\hat{\theta}) \sim \chi^2_{q=3} \rightarrow p\text{-value} = 6.96 \cdot 10^{-200} \ll 0.001 \rightarrow$   $\text{FALSE}$   
 ελαττωμένης και απορριψω  $H_0 \Rightarrow$  ελαττωμένης  $M_3$

Επειδή  $AIC_3 < AIC_0 \rightarrow M_3$  καλύτερο  
 αλλά: για τα  $z\text{-tests}$   $p\text{-values}$  του  $M_3$  όλοι ήταν μεγαλύτεροι επιρροής

(iii)  $e^{\hat{\beta}_2}$  του  $M_3$ :  $e^{\hat{\beta}_2} = e^{5.8892 \cdot 10^{-3}} = 1.0059$   $\rightarrow$   $\text{FALSE}$   
 Αφαιρούμε τον αριθμό ελαττωμένης  $\rightarrow$  δεσφύρι να απορριψω  $H_0$  και δεσφύρι το  $M_3$  (μικρότερο)



Zirupa 4

(A)  $Y \sim \text{Stochastisch}$ ,  $f(y) = \binom{n}{y} p^y (1-p)^{n-y}$ ,  $y = 0, 1, 2, \dots, n$  Parameter:  $p, n$   
 $Y_X \sim b(n_X, p_X)$  ( $n_X > 1$  Stochastisch & deterministisch)

$$\eta(x) = g(\mu_x) = \ln \frac{p_x}{n_x - p_x} = b_0 + b_1 x_1 + \dots + b_k x_k$$

$$= E(y) = \mu_x, \quad V(y) = \mu_x(1 - \mu_x)$$

link function.

$$= x' \beta \quad (1)$$

$$\mu_x = E(y) = n_x p_x, \quad V(y) = n_x p_x (1 - p_x)$$

Αντικείμενα της (1) ως προς  $\rho_x$ :

$$P_i = \frac{e^{b_0 + b_1 x_{i1} + \dots + b_k x_{ik}}}{1 + e^{b_0 + b_1 x_{i1} + \dots + b_k x_{ik}}} = \frac{e^{x_i \beta}}{1 + e^{x_i \beta}}$$

for  $E(y_i) = n_i p_i = n_i \frac{e^{x_i \beta}}{1 + e^{x_i \beta}}$

(B) m: πωλητές προϊόντων, αποθηκευτική χωρητικότητα κερπών  $Y = \begin{cases} 1, & \text{ναι} \\ 0, & \text{όχι} \end{cases}$   
 $X_1$ : μισθός υπαλλήλου πωλητή κερπών,  $X_2$ : ετήσιο εισόδημα,  $X_3 = \begin{cases} 1, & \text{φοιτητής} \\ 0, & \text{όχι φοιτητής} \end{cases}$

$$(i) \quad Z = \frac{\hat{\beta}_j - \beta_j(0)}{se(\hat{\beta}_j)} \sim N(0, 1) \text{ asymptotically}$$

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

$z_2 = \frac{1.181 \cdot 10^{-5}}{9.397 \cdot 10^{-6}} \rightarrow p\text{-value} = 2 \cdot pnorm(|z_2|, lower.tail=FALSE) = 0.209$

απλ. έλεγχος με κριτ. τιμή:  $\beta_2 = 0$  → Δεν απορρίπτουμε ή απορρίπτουμε  $H_0$ . Απόδοξη  $H_0$   
 $z_3 \rightarrow p\text{-value} = 0.136$

### MONTEAO 2:

$$Z_3 = -4.009 \rightarrow p\text{-value} = 6.098 \cdot 10^{-5}$$

MONTELO 3:

$$z_3 = \frac{\hat{\beta}_3}{\text{se}(\hat{\beta}_3)} = 3.482 \rightarrow p\text{-value} = 0.000437$$

~~Handwritten scribbles and crossed-out text.~~

Monte Carlo 1:  $AIC_1 = 1179.2$

$$\hat{\ell}_1 = 4 - \frac{1179.2}{2} = -586.6$$

Montano 2:  $l_2 = d - \frac{AIC_2}{2} = 3 - \frac{1178.8}{2} = -586.4$

Modelo 3:  $AIC_3 = -2l_3 + 2d = -2 \cdot (-1127.042) + 2 \cdot 2 = 2258.084$

Τα μοντέλα 1 και 2 έχουν χαμηλότερο AIC σε σχέση με το μοντέλο 3.  
 Άρα θα ταιγαρίσω  $M_1$  vs  $M_2$ .

$$H_0: M_2 \sim x_1, x_3 \quad C \quad H_1: M_1 \sim x_1, x_2, x_3$$

$$D_2 - D_1 \sim \chi^2_{\text{αριθμός παραμέτρων}} = \chi^2_1$$

$$D_2 - D_1 = 1172.8 - 1171.2 = 1.6 \sim \chi^2_1$$

p-value = 0.2059  $\Rightarrow$  από μη στατ. σημαντικός έλεγχος, οπότε δεν πρέπει να απορρίψω την  $H_0 \Rightarrow$  επιλέγω  $M_2 \sim x_1, x_3 \rightarrow$  το καλύτερο μοντέλο

(ii) 95% ΔΕ για  $e^{\beta_3}$  του τριανού παραμέτρου  $M_2$

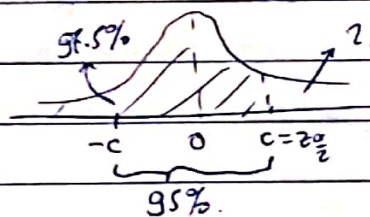
$$Z = \frac{\hat{\beta}_j - \beta_j(0)}{se(\hat{\beta}_j)} \rightarrow z_3 = \frac{\hat{\beta}_3 - \beta_3}{se(\hat{\beta}_3)} \sim N(0,1) \text{ ασυμμετρικά}$$

$$P[-c \leq z_3 \leq c] = 0.95 \rightarrow \cancel{1 - 0.05} c = Z^{-1}(\cancel{0.975}) 0.975 = z_{\frac{\alpha}{2}} \text{ και ΔΕ}$$

$$-c \leq \frac{\hat{\beta}_3 - \beta_3}{se(\hat{\beta}_3)} \leq c \Rightarrow -c \cdot se(\hat{\beta}_3) \leq \hat{\beta}_3 - \beta_3 \leq c \cdot se(\hat{\beta}_3) \Rightarrow \hat{\beta}_3 + c \cdot se(\hat{\beta}_3) \geq \beta_3 \geq \hat{\beta}_3 - c \cdot se(\hat{\beta}_3)$$

$$\Rightarrow \beta_3 \in [\hat{\beta}_3 \pm c \cdot se(\hat{\beta}_3)] \rightarrow e^{\beta_3} \in e^{\hat{\beta}_3 \pm c \cdot se(\hat{\beta}_3)}$$

$$c = 1.96 \text{ από } q_{\text{norm}}(0.975)$$



$$\hat{\beta}_3 = -0.6814, \quad se(\hat{\beta}_3) = 0.17$$

$$\text{και } e^{\hat{\beta}_3 \pm c \cdot se(\hat{\beta}_3)} \rightarrow e^{-0.6814 - 1.96 \cdot 0.17} \approx 0.363$$

$$\rightarrow e^{-0.6814 + 1.96 \cdot 0.17} = 0.701$$

$$\text{αρα } 95\% \text{ ΔΕ για } e^{\beta_3}: \quad 0.363 \leq e^{\beta_3} \leq 0.701$$

(iii)  $e^{\hat{\beta}_i}$  για  $M_2$

$$M_2: \hat{\beta}_1 = 5.982 \cdot 10^{-3}, \quad \hat{\beta}_3 = -6.814 \cdot 10^{-1}$$

$$e^{\hat{\beta}_1} \approx 1.006, \quad e^{\hat{\beta}_3} \approx 0.506$$



$$p_x = E(Y_x) = \pi_x p_x = \pi_x \frac{e^{x_1' \beta}}{1 + e^{x_1' \beta}}$$

$$\frac{p_x}{1 - p_x} = \frac{e^{x_1' \beta}}{1 + e^{x_1' \beta}} \cdot \frac{1 + e^{x_1' \beta}}{1 + e^{x_1' \beta} - e^{x_1' \beta}} = \frac{e^{x_1' \beta}}{1 + e^{x_1' \beta} - e^{x_1' \beta}} = e^{x_1' \beta} = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}$$

Αν φοιτητής ( $x_3 = 1$ ) :  $\frac{p_x}{1 - p_x} = e^{\beta_0 + \beta_1 x_1 + \beta_2} = e^{\beta_2} \cdot e^{\beta_0 + \beta_1 x_1}$

Αν όχι φοιτητής ( $x_3 = 0$ ) :  $\frac{p_x}{1 - p_x} = e^{\beta_0 + \beta_1 x_1}$

$$\frac{\left(\frac{\hat{p}_x}{1 - \hat{p}_x}\right)_{x_3=1}}{\left(\frac{\hat{p}_x}{1 - \hat{p}_x}\right)_{x_3=0}} = e^{\hat{\beta}_2} \approx 0.506 \approx 50.6\% \rightarrow \text{δισκ. που φοιτητής ή μη δαύτην πηλαιοπλάσσει με παρόμοια 0.506}$$

Άρα αν κάποιος είναι φοιτητής (εν σχέση με έναν μη φοιτητή) συνεπάγεται μείωση κατά 49.4% στη σχετική πιθανότητα αποτυχής της πωλητικής κάρτας. Ανεκτιμώμενο αποτέλεσμα για έναν φοιτητή.

(iv) Το εμβαδόν κάτω από την ROC curve είναι το AUC. Καλή/επιτυχία πρόβλεψη ενός μοντέλου σημαίνει πως υπάρχουν τμήματα του ορίου που με υψηλή ευαισθησία και ~~καλή~~ <sup>αποδοτικότητα</sup> υψηλή ειδικότητα. Τότε η ROC πλησιάζει την άνω αριστερή γωνία του τετραγώνου του σχήματος. Ένας δείκτης που μετρά κατά πόσο πλησιάζει αυτή τη γωνία, είναι το εμβαδόν κάτω από την καμπύλη (area under curve AUC) με μέγιστη τιμή το 1. Εδώ βλέπουμε ότι  $M_1, M_2$  έχουν την καλύτερη ROC και ίδια βασικά. Ο βόρος  $M_2$  έχει λιγότερες μεταβλητές  $\Rightarrow$  επιλογή ως καλύτερο.