

ΘΕΜΑΤΑ 2021

ΖΗΤΗΜΑ 1

Ισχύει $x_2 = 0 \rightarrow \text{χαμηλή}$, $x_2 = 1 \rightarrow \text{υψηλή}$, οπότε
 $x_3 = 0 \rightarrow \text{χαμηλή}$, $x_3 = x_1 \rightarrow \text{υψηλή}$

Το μοντέλο $E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$ ισοδυναμεί με

$$E(y_x) = \beta_0 + \beta_1 x_1 \rightarrow \text{χαμηλή} \quad \eta \quad \mu \epsilon$$

$$E(y_u) = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) x_1 \rightarrow \text{υψηλή}.$$

Συγκρίνοντας τις σχέσεις αυτές, γίνεται αντιληπτό πως (I) εάν $\beta_3 \neq 0$, τότε οι συντελεστές του x_1 διαφέρουν στις δύο περιπτώσεις και ως εκ τούτου απαιτούνται δύο ειδικές διαφορετικών κλίσεων για την προσάρμογή στα δεδομένα. (II) Εάν $\beta_3 = 0$, αλλιώς $\beta_2 \neq 0$, τότε ξανά δύο ειδικές θα πρέπει να προσαρμόζονται στα δεδομένα, μία για κάθε είδος (υψηλή/χαμηλή), μόνο που στην περίπτωση αυτή η κλίση τους θα είναι κοινή και ίση με β_1 . (Παρά-Ημένες ειδικές). Τέλος, (III) εάν $\beta_3 = 0$ και $\beta_2 = 0$, τότε οι δύο σχέσεις ταυτίζονται, οπότε μία μόνο ειδικά αρκεί για την προσάρμογή στα δεδομένα.

Οι έλεγχοι που προτείνονται στα πλαίσια αυτά είναι οι ακόλουθοι: αρχικά, ο έλεγχος της υπόθεσης $H_0: \beta_3 = 0$ έναντι της $H_1: \beta_3 \neq 0$, η οποία ισοδυναμεί με την μηδενική απουσία αλληλεπίδρασης μεταξύ των x_1 και x_2 . Εάν η υπόθεση αυτή απορριφθεί, τότε είναι βέβαιο πως μελετάται η περίπτωση (I). Διαφορετικά, ακολουθεί ο έλεγχος της υπόθεσης $H_0: \beta_2 = 0$ έναντι της $H_1: \beta_2 \neq 0$. Απόρριψη της υπόθεσης αυτής ισοδυναμεί με

των περιπτώσεων (II), ενώ αποδοχή της με των περιπτώσεων (III).

$$\bar{R}^2 = 1 - \left[\frac{(1-R^2)(n-1)}{n-k-1} \right] = 1 - \left[\frac{0.032 \cdot 13}{10} \right] = 0.958,$$

$$\text{ήτοι } \bar{R}^2 = 95.8\%$$

$$SSE = SST - SSR = 349942 - 338736 \Rightarrow SSE = 11206$$

$$\text{Για } y \text{ versus } x_1: R_{\text{pred}}^2 = 1 - \frac{\text{PRESS}}{SST} \Rightarrow$$

$$R_{\text{pred}}^2 = 1 - \frac{255010}{349942} \Rightarrow R_{\text{pred}}^2 = 27.13\%$$

Όπως υποδεικνύει και το διάγραμμα, για την περιγραφή των επιφάνων απαιτούνται δύο διαφορετικές επιδείξεις. Αυτό μπορεί να φανεί από τα ήδη υπάρχοντα δεδομένα, αλλά προκύπτει και από τους ελέγχους που προτάθηκαν προηγουμένως:

$$H_0: \beta_3 = 0, \quad F = \frac{SSE_2 - SSE_3}{SSE_3/(n-p)} = \frac{38358 - 11206}{11206/10} \Rightarrow$$

$\Rightarrow F = 24.229$ με $p\text{-value} = 0.0006$, οπότε η H_0 απορρίπτεται και για 2 επιδείξεις διαφορετικές επιδούς απαιτούνται.

ΖΗΤΗΜΑ 2

$$A) R^2 = \frac{SSR}{SST} = \frac{3.53751 \cdot 10^{14}}{3.6379 \cdot 10^{14}} = 97,24\%$$

Φαίνεται οι X_5 και X_1 να μην έχουν καλή προσάρτηση
βόθρ του F-test. Το κατά πόσο είναι απαραίτητες για
το μοντέλο σπαστεί περισσότερη διερεύνηση.

B) Βόθρ του ελέγχου G δέχουμε όσο το δυνατότερο μικρό
G και αντίστοιχα S και AIC. Αντίθετα, οι τιμές του
 R^2 και R^2_{pred} επιθυμώμε να είναι όσο το δυνατό
υψηλότερες. Βόθρ αυτών, 2 μοντέλα που σίγουρα ξεχω-
ρίζουν είναι τα M5 και M7.

Ισχύει $SSE = S^2(n-k-1)$, άρα

$$SSE_5 = 1.0558 \cdot 10^{13} \quad \text{και} \quad SSE_7 = 1.0063 \cdot 10^{13}$$

Τα μοντέλα διαφέρουν κατά 1 μεταβλητή, άρα ισχύει

$$F = \frac{SSE_5 - SSE_7}{SSE_7 / (n-p)} = \frac{(1.0558 - 1.0063) \cdot 10^{13}}{1.0063 \cdot 10^{13} / 46} = 2.26,$$

με p-value = 0.139, το οποίο είναι αρκετά υψηλό, άρα
επιθυμώμε το μοντέλο M5, βάσει του κριτηρίου F.

Σε ό,τι αφορά το Μοντέλο 3, το A.E για τη βήμερακή
πρόβλεψη είναι το $[\hat{\mu} - \underbrace{t_{0.025, 46} \cdot S \cdot (\hat{x}'(X'X)^{-1}\hat{x} + 1)^{1/2}}_{\delta}, \dots]$

$$\text{Άρα } 110976 - \delta = -920677 \Rightarrow \delta = 1031653$$

$$\text{Άρα το άλλο άκρο είναι το } 110976 + \delta = 1142629$$

2ΗΡΗΜΑ 3

$$AIC_0 = -2\hat{\ell}_0 + 2 \cdot 1 \Rightarrow AIC_0 = 3673.56$$

$$D_0 - D_1 = 0.93 \Rightarrow -2(\hat{\ell}_0 - \hat{\ell}_1) = 0.93 \Rightarrow \hat{\ell}_1 = \hat{\ell}_0 - \frac{0.93}{2} \Rightarrow$$

$$\Rightarrow \hat{\ell}_1 = -1836.245$$

$$AIC_1 = -2\hat{\ell}_1 + 2 \cdot 2 \Rightarrow AIC_1 = 3676.49$$

$$D_2 - D_3 = 920.1, p\text{-val} < 0.001 \rightarrow \text{απόρριψη } H_2$$

$$R_D^2 = 1 - D(\hat{\beta})/D_0$$

$$\text{Για το μοντέλο 2: } R_D^2 = 1 - 3507.3/3510.73 \Rightarrow$$

$$R_D^2 = 0.09\%$$

$$\text{Confint: } Z_{\alpha/2} = 1.96, \text{ όρα είναι το}$$

$$[5.889 \cdot 10^{-3} - 1.96 \cdot 1.434 \cdot 10^{-3}, 5.889 \cdot 10^{-3} + 1.96 \cdot 1.434 \cdot 10^{-3}] =$$

$$[3.07836 \cdot 10^{-3}, 8.69964 \cdot 10^{-3}] \text{ για τον } \hat{\beta}_2, \text{ όρα}$$

$$1.00308 < e^{\hat{\beta}_2} < 1.0087, \text{ με μέση τιμή}$$

$$e^{\hat{\beta}_2} = 1.0059$$

Ερμηνεία: Εάν η απόσταση από το πιο κοντινό μερί αυξάνει κατά 1 μονάδα, τότε το γρήγορο αέριο φούρνι αυξάνεται κατά 0.59% (αφού παύεται επί 1.0059)

2ΗΤΗΜΑ 4

λη $\frac{P_i}{1-P_i} = \beta_0 + \beta_1 X_1^i + \dots + \beta_k X_k^i$, για αν i -οστή παρατήρηση

X_1 : μισθός ανδρών X_2 : επίπεδο εκπαίδευσης $X_3 = 1$: φοιτητής

i) Για β_2 $AIC_2 = -2\hat{\ell}_2 + 2 \cdot 3 \Rightarrow 1178.8 - 6 = -2\hat{\ell}_2 \Rightarrow \hat{\ell}_2 = -586.4$

$AIC_3 = -2\hat{\ell}_3 + 2 \cdot 2 \Rightarrow AIC_3 = 2258.084$

$D_2 - D_1 = 1.6$, p -value = 0.206 μικρή διαφορά deviance, σχέση με p -value \rightarrow το M_1 απορρίπτεται

$D_3 - D_2 = -2(\hat{\ell}_3 - \hat{\ell}_2) = 1081.284$, p -value < 0.001 , όρα λόγω της πολύ υψηλής διαφοράς των Deviance και του πολύ μικρού p -value, το M_3 απορρίπτεται

Δεχόμαστε, λοιπόν, το M_2 , γεγονός που έρχεται σε σύμφωνια και με το κριτήριο AIC αλλά και με τους ελέγχους Wald (για όλα τα β ισχύει p -value < 0.001).

ii) $z_{0.02} = 1.96$, ισχύει $\beta_3 \in [\hat{\beta}_3 - z_{0.02} se(\hat{\beta}_3), \hat{\beta}_3 + z_{0.02} se(\hat{\beta}_3)]$
 $= [-0.6814 - 1.96 \cdot 0.17, -0.6814 + 1.96 \cdot 0.17] = [-1.0146, -0.3482]$

Άρα $e^{\beta_3} \in [0.3625, 0.706]$

iii) Δεδομένου πως $e^{\hat{\beta}_3} = 0.506$, φαίνεται πως ένας φοιτητής έχει 49.4% περισσότερη πιθανότητα να αποθηρεύσει της πιθανότητας σε σχέση με κάποιον που δεν είναι φοιτητής.

(iv) Οι καμπύλες ROC αντιστοιχούν ένα δείγμα καρτίς προσομοιώντας τις τυχαίες λογιστικές παλινδρόμους και παρίστανται ότι δείχνει την ευαισθησία από την ειδικότητα του μοντέλου. Συγκεκριμένα, όσο μεγαλύτερο είναι το AUC (area under the curve) της ROC, τόσο υψηλότερη είναι η προβλεπτική ικανότητα του μοντέλου.

Από τις δοσμένες καμπύλες φαίνεται ξεκάθαρα πως η προβλεπτική ικανότητα του Μοντέλου 3 δεν είναι καλή (το AUC του είναι σχεδόν 0.5, δηλαδή αυτό του random classifier), ενώ αντίθετα είναι πολύ καλή για τα Μοντέλα 1 και 2, αφού το AUC είναι συντριβάνη.

⊕ αναφορά στο CI του προγράμματος εξωτελισμού
δηλ. κυμαίνεται από 30% έως 64% περίπου