

$df_{SST} + 1 = 26 + 1$ ή από πίνακα σημειώνω στο διάγραμμα.

ΖΗΤΗΜΑ 1

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

$n = 27$, $x_2 = \begin{cases} 1, & \text{αν διατίκωση} = 1 \\ 0, & \text{αν διατίκωση} = 2 \end{cases}$

y : ποσότητα μεταλλού, x_1 : ταχύτητα παραγωγής, $x_3 = x_1 \cdot x_2$

Στα αποτελέσματα του regression analysis για το Μοντέλο 1: y vs x_1, x_2, x_3 φαίνονται οι έλεγχοι t που πραγματοποιούνται για τις 3 ανεξαρτητές:

$$x_j \text{ μεταβλητή: } t^* = \frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)}$$

$$k = 3 \text{ μεταβλητές} \Rightarrow p = k + 1 = 4 \text{ και } n - p = 27 - 4 = 23$$

Έχουμε $H_0: \beta_j = 0$, $H_1: \beta_j \neq 0$ με στατιστικό έλεγχο υποθέτουμε.

$$t^* = \frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} = \frac{\hat{\beta}_j}{se(\hat{\beta}_j)} \sim t_{n-p} = t_{23}$$

$$\text{αρα για } x_2: t^* = \frac{\hat{\beta}_2}{se(\hat{\beta}_2)} \xrightarrow{\text{SE coef}} \text{δίνεται } t^* = 3.19$$

$$\text{αρα } p\text{-value} \approx 0.004$$

$$\text{για } x_3: t^* = \frac{\hat{\beta}_3}{se(\hat{\beta}_3)} = \frac{-0.1767}{0.1288} \approx -1.37 \sim t_{23}$$

$$p\text{-value} \approx 0.183$$

$$\bar{R}^2 = R\text{-Sq}(\text{adj}) = 1 - \frac{n-1}{n-p} (1 - R^2) = 1 - \frac{26}{23} (1 - 0.945) = 0.9378 \approx 93.78\%$$

Για την ANOVA του μοντέλου 1.

$$\text{Έχουμε } SST = SSR + SSE \Rightarrow SSR = SST - SSE = 179,069 - 9,904 \Rightarrow$$

$$SSR = 169,165$$

Για το Μοντέλο 2: y vs x_1, x_2 , $n = 27$, $k = 2$ μεταβλητές, $p = k + 1 = 3$, $n - p = 27 - 3 = 24$

$$\bar{R}^2 = 1 - \frac{n-1}{n-p} (1 - R^2) \Rightarrow 1 - R^2 = \left(1 - \bar{R}^2\right) \frac{n-p}{n-1} \Rightarrow R^2 = 1 - \frac{n-p}{n-1} (1 - \bar{R}^2) =$$

$$= 1 - \frac{24}{26} (1 - 0.935) = 0.94 \pm 94\%$$

Στην ανάλυση ANOVA έχουμε ότι :

$$F^* = \frac{MSR}{MSE} = 188.57 \text{ και είναι έλεγχος } F \text{ που αφορά τον συνολικό έλεγχο υποτεθείται}$$

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

$$H_1 : \text{τουλάχιστον ένας συντελεστής } \beta_i \neq 0, i=0,1,2$$

$$F^* = 188.57 \sim F_{d, n-p} = F_{2, 24}$$

αριθμός μεταβλητών
μοντέλου

$$\text{όρα } p\text{-value} = 2 \cdot 10 \cdot 10^{-15}$$

Αποφά για την μεταβλητή x_2 : t-test με $t^* = \frac{\hat{\beta}_2}{se(\hat{\beta}_2)} = \frac{53.129}{8.210} = 6.47$
 και $t^* = 6.47 \sim t_{n-p} = t_{24} \rightarrow p\text{-value} = 1.08 \cdot 10^{-6}$

Για το μοντέλο 3 : y vs x_1 : $n = 27$, $k = 1$, $p = k+1 = 2$, $n-p = 25$

$$R^2_{\text{pred}} = 1 - \frac{\text{PRESS}}{\text{SST}} \quad (1)$$

$$\text{όρα } R^2 = \frac{\text{SSR}}{\text{SST}} \Rightarrow \text{SST} = \frac{\text{SSR}}{R^2} = \frac{149,661}{0.829} = 180,532$$

$$\text{όρα } \text{SSE} = \text{SST} - \text{SSR} = 180,532 - 149,661 = 30,871$$

$$\text{και από (1)} : R^2_{\text{pred}} = 1 - \frac{\text{PRESS}}{\text{SST}} = 1 - \frac{34,546.9}{180,532} \approx 0.8086 = 80.86\%$$

$$MSE = \frac{\text{SSE}}{n-p} = \frac{30,871}{25} = 1,234.84$$

$$\text{και } F^* = \frac{MSR}{MSE} = \frac{149,661}{1,234.84} = 121.20 \sim F_{d, n-p} = F_{1, 25}$$

αριθμός μεταβλητών
μοντέλου

$$\text{όρα } p\text{-value} = 4.45 \cdot 10^{-11}$$

Μέσω του μοντέλου $E(y_x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$ μπορούμε να ελέγξουμε αν χρειάζεται να προσεφθεσθούν δύο διαφ. ενδείξεις, δύο παράλληλες ή μια κοινή γραμμή δύο παραμρφικές διαδικασίες ($x_2=0$, $x_2=1$).

Αντί γράφω

• Αν $\beta_3 \neq 0$:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 x_3$$

• Όταν $x_2 = 0$: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1$

• Όταν $x_2 = 1$: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 + \hat{\beta}_3 x_1 = (\hat{\beta}_0 + \hat{\beta}_2) + (\hat{\beta}_1 + \hat{\beta}_3) x_1$
Δηλαδή δύο ξεχωριστές ευθείες.

• Αν $\beta_3 = 0$

• Αν $\beta_2 = 0$: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_2 x_2$

• Αν $\beta_1 \neq 0$: 1 κοινή ευθεία για τις 2 διαδικασίες

• Όταν $x_2 = 0$: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1$

• Όταν $x_2 = 1$: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 = (\hat{\beta}_0 + \hat{\beta}_2) + \hat{\beta}_1 x_1$
Δηλαδή δύο παράλληλες ευθείες.

Το μοντέλο 3 βλέπουμε ότι έχει αρκετά χαμηλότερους δείκτες προσδιορισμού σε σχέση με τα 2 πρώτα μοντέλα που φαίνονται να εξηγούν καλύτερα τα δεδομένα μας.

Θα συγκρίνουμε τα 2 μοντέλα (1 & 2) μέσω ενός στατιστικού ~~ελέγχου~~ ^{υποθέτου} :

$H_0: \beta_3 = 0$: $M_2: y \sim x_1, x_2$

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

$$F^* = \frac{(SSE_{M_2} - SSE_{M_1})/1}{SSE_{M_1}/(27-3-1)} = \frac{10714 - 9904}{9904/23} = 1.88 \approx F_{d,23} = F_{1,23}$$

αρχ $p\text{-value} = 0.184$

↓
Διαφορά παραμέτρων στα 2 μοντέλα

φαίνεται, λοιπόν, ότι ο έλεγχος είναι στατιστικά μη σημαντικός που σημαίνει ότι δεν μπορούμε να απορρίψουμε την H_0 , δηλαδή το μοντέλο με x_1, x_2 (M_2) το οποίο είναι και το μοντέλο που θα επιλεγούμε ως βέλτιστο. Δηλαδή $E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
δηλ. 2 παράλληλες ευθείες θα ήταν η προσρμογή με βάση τα παραπάνω.

(πχ για επίπεδο σημαντικότητας 5%)

Όσο, δίνεται στον εκφώνημα ένα χάρη που φαίνεται να εξετάζει την

προαρκούν του μοντέλου 1 (δω μη παράλληλες ευθείες),
 Παρ'όλαυτα βλέπουμε ότι και αυτό δεν απέχει πολύ από το
 προηγούμενο γενάριο καθώς οι ευθείες δείχνουν να έχουν παρόμοια
 κλίση.

Η ερμηνεία που δίνουμε στο διάγραμμα είναι ότι όταν για δεδομένη τιμή
 ταχύτητας παραγωγής, ~~μεταβαίνουμε~~ μεταβαίνουμε από τη διαδικασία 2
 ($x_2=0$) στη διαδικασία 1 (~~$x_2=1$~~) τότε έχουμε αύξηση
 στην ποσότητα μεγάλου. ~~Επαρκώς, κάποιος ίδιος, να προτιμάμε να γίνει η~~
 Τελικά με βάση το διάγραμμα ~~και~~ ο ρυθμός αύξησης της
 ποσότητας μεγάλου για τη διαδικασία 1 ($x_2=1$) είναι μικρότερος (κλίση
 ευθείας) από τον αντίστοιχο για τη διαδικασία 2 ($x_2=0$).

Για το Μοντέλο 2 που βρήκαμε ως βέλτιστο θα ίσχυαν τα ίδια, όχι
 όμως ~~είναι~~ οι ρυθμοί αύξησης της ποσότητας μεγάλου είναι
 διαφορετικοί.

ΖΗΤΗΜΑ 4

$$f(y) = \frac{e^{-\mu_x} \cdot \mu_x^y}{y!}, \quad y=0,1,2. \quad n(x) = g(\mu_x) = \ln \mu_x = \theta'x \Rightarrow \mu_x = e^{\theta'x}$$

$$D_n(\hat{\theta}) = -2(\hat{\ell}_n - \hat{\ell}_{kup}) = 2 \sum_{i=1}^n \left[y_i \ln \frac{y_i}{\hat{\mu}_i} \right] \quad AIC = -2\hat{\ell}_n + 2d$$

$n=12$, X_1 και X_2 \rightarrow φύλο, $x_2 = \begin{cases} 0, & \text{άντρας} \\ 1, & \text{γυναίκα} \end{cases}$
 \swarrow αριθμός \swarrow διαδικασία φαρμάκου
 \swarrow \swarrow

Για τις p-values όταν αφορά τους ελέγχους Wald έχουμε: ~~άλλα και στο Z-test~~

$$H_0: \beta_j = 0 \quad H_1: \beta_j \neq 0$$

$$z_j = \frac{\hat{\beta}_j}{se(\hat{\beta}_j)} \sim N(0,1)$$

από το Μοντέλο 2: $z_1 = \frac{-14.835}{0.0629} \rightarrow p\text{-value} = 2 \cdot p\text{norm}(\text{abs}(\frac{-14.835}{0.0629}), \text{mean}=0, \text{sd}=1, \text{lower.tail}=\text{FALSE})$

$$z_2 = \frac{\hat{\beta}_2}{se(\hat{\beta}_2)} = \frac{-0.013193}{0.0629} = -0.210 \rightarrow p\text{-value} = 0.834$$

Ορίσας για Μοντέλο 1:

$$z_1 = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{-0.12962}{0.00873} = -14.848 \rightarrow p\text{-value} = 7.17 \cdot 10^{-50}$$

■

για ~~Μοντέλο~~ Μοντέλο 0 ισχύει:

$$D_0(\hat{\beta}) = -2(\hat{\ell}_0 - \hat{\ell}_{top}) \Rightarrow \hat{\ell}_{top} = \frac{D_0(\hat{\beta})}{2} + \hat{\ell}_0 = \frac{272.305}{2} + (-222.219)$$

$$\Rightarrow \hat{\ell}_{top} = -85.567$$

$$\text{Για Μοντέλο 1: } D_1(\hat{\beta}) = -2(\hat{\ell}_1 - \hat{\ell}_{top}) \Rightarrow \hat{\ell}_1 = \hat{\ell}_{top} - \frac{D_1(\hat{\beta})}{2} =$$

$$= -85.567 - \frac{67.695}{2} \Rightarrow \hat{\ell}_1 = \cancel{50.110} - 120.415$$

$$AIC_1 = -2\hat{\ell}_1 + 2 \cdot d \xrightarrow{\text{όλοι οι παράμετροι}} = -2 \cdot (-120.415) + 2 \cdot 2 = 244.83$$

$$AIC_0 = -2\hat{\ell}_0 + 2 \cdot d = -2 \cdot (-222.219) + 2 \cdot 1 = 446.44$$

Το M_1 σε σχέση με M_0 έχουν 1 μεταβλητή extra και παράμετρο διαφορά, δηλ. η διαφορά σε β.ε ~~είναι~~ είναι 1.

$$D_0(\hat{\beta}) - D_1(\hat{\beta}) = 272.305 - 67.695 = 204.61$$

$H_0: M_0$

$H_1: M_1$

$$D_0(\hat{\beta}) - D_1(\hat{\beta}) \sim \chi^2_1 = \chi^2_1$$

\rightarrow διαφορά σε παράμετρος 2 μοντέλων.

αρκ ~~p-value~~ για $204.61 \sim \chi^2_1 \rightarrow p\text{-value} = 2.06 \cdot 10^{-46}$.

αρκ βελτιστικά επιλεγχοί και απορρίπτω $H_0 \Rightarrow$ επιλέγω M_1 over M_0

$$\text{Ορίσας } D_1(\hat{\beta}) - D_2(\hat{\beta}) = 0.043 \sim \chi^2_1 \rightarrow p\text{-value} = 0.836$$

μη διατ. επιλεγχοί και δεν απορρίπτω $H_0 \Rightarrow$ επιλέγω M_1 over M_2

Μοντέλο 1:

$$R^2_D = \left(1 - \frac{D_1(\hat{\beta})}{D_0}\right) \cdot 100\% = \left(1 - \frac{67.695}{272.305}\right) \cdot 100\% = 75.14\%$$

(ii)

στο M_2 με βάση έλεγχο Wald η χ^2 ήταν μη στατιστ. σημαντική
Επίσης M_2 έχει το χαμηλότερο AIC και deviance ελάχιστα
διαφορετικά από M_2 ενώ ο χ^2 έλεγχος της διαφοράς των
deviance μεταξύ M_1 και M_2 έδειξε ότι επιλέγω το M_1
ως το βέλτιστο μοντέλο. Επίσης και ο δείκτης ψευδο R^2 -deviance
έδειξε μερική αύξηση πληρότητας από M_0 στο M_1 και ελάχιστη αύξηση από το M_1 στο M_2 .

Άρα: καλύτερο μοντέλο:

$$M_1: E(y) = \mu_x = e^{\beta'x} = e^{\beta_0 + \beta_1 \cdot x_1}$$

και το προσαρμοσμένο μοντέλο: $\hat{\mu}_x = e^{\hat{\beta}_0 + \hat{\beta}_1 \cdot x_1} = e^{0.37677 - 0.12962 \cdot x_1}$

σηλαδή προέκυψε ανεξαρτησία από το φύλο (x_2 μεταβλητή).

(iii) Θέλουμε ένα 0.95-ΔΕ για e^{β_1} του M_1 .

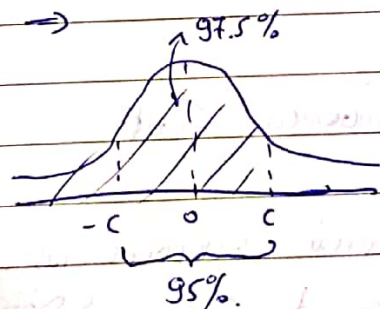
$$z_1 = \frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} \sim N(0,1) \text{ ασυμπτωτικά}$$

Θέλουμε βαθιά c τέτοια ώστε: $P[-c \leq z_1 \leq c] = 0.95 \rightarrow c = Z^{-1}(0.975)$ και ΔΕ:

$$-c \leq \frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} \leq c \Rightarrow -c \cdot se(\hat{\beta}_1) \leq \hat{\beta}_1 - \beta_1 \leq c \cdot se(\hat{\beta}_1) \Rightarrow$$

$$\Rightarrow \hat{\beta}_1 - c \cdot se(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + c \cdot se(\hat{\beta}_1) \text{ και}$$

$$e^{\hat{\beta}_1 - c \cdot se(\hat{\beta}_1)} \leq e^{\beta_1} \leq e^{\hat{\beta}_1 + c \cdot se(\hat{\beta}_1)}$$



$c \approx 1.96$ από $\Phi_{norm}(0.975)$ (R εντολή)

$\hat{\beta}_1 = -0.12962$, $se(\hat{\beta}_1) = 0.00873$ άρα

$\hat{\beta}_1 - c \cdot se(\hat{\beta}_1) = -0.1467$, $\hat{\beta}_1 + c \cdot se(\hat{\beta}_1) = -0.1125$

και $e^{\hat{\beta}_1 - c \cdot se(\hat{\beta}_1)} = e^{-0.1467} = 0.864$, $e^{\hat{\beta}_1 + c \cdot se(\hat{\beta}_1)} = e^{-0.1125} = 0.894$.

άρα 0.95 ΔΕ για e^{β_1} :

$$0.864 \leq e^{\beta_1} \leq 0.894$$

$e^{\hat{\beta}_1} = e^{-0.12962} = 0.878$

$e^{\hat{\beta}_1}$ είναι εντός του 95% ΔΕ του e^{β_1}

Παρατηρούμε ότι η τιμή του $e^{\hat{\beta}_1}$ είναι εντός του 95% ΔΕ του e^{β_1}
όπως αναμέναμε.

Με: $\hat{\mu}_x = e^{\hat{\beta}_0 + \hat{\beta}_1 x}$

Για αύξηση κατά 1 μονάδα της x_1 , δηλ. $x_1' = x_1 + 1$:

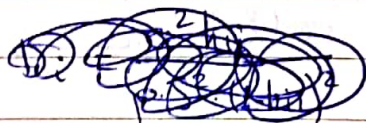
$$\hat{\mu}_{x'} = e^{\hat{\beta}_0 + \hat{\beta}_1 x'} = e^{\hat{\beta}_0 + \hat{\beta}_1 (x_1 + 1)} = e^{\hat{\beta}_1} \cdot e^{\hat{\beta}_0 + \hat{\beta}_1 x_1} = e^{\hat{\beta}_1} \cdot \hat{\mu}_x$$

Δηλ. ποσ/στηκε η τιμή του $\hat{\mu}_x$ με $e^{\hat{\beta}_1} = 0.878$ ή
 $1 - 0.878 = 0.122 = 12.2\%$
 ισοδύναμα μειώθηκε κατά

άρθ. ανήκοντος τη δόσολογία φαρμάκου κατά 1 μονάδα προκύπτει
 μείωση 12.2% στον αριθμό ασθενών με θετική απάνκριση.

(iv) Για τα υποδογία deviance παρατηρούμε ότι γενικά κινούνται όπως θα
 περιμέναμε, δηλ. πάνω στην ευθεία γραμμή με ελαφρώς τα βήματα 30 και 35
 που ίσως είναι σημεία επιρροής και κάποια βήματα κατά την μέση του
 άξονα x των normal quantiles. ~~Ενδεχομένως είναι περισσότερα από 2 σημεία
 επιρροής στο 30 το έρθαν τα βήματα από πιο κάτω στην ευθεία
 όποτε ο βόρος ή υποπείρθηκε κέντρο ή κανονισμός~~ Να σημειωθεί εδώ πως
 τα υποδογία deviance δεν κατανομισαν σύμφωνα με την κανονική κατανομή
 και χρησιμοποιούν στον εντοπισμό άτυπων ή απόμακρων παρατηρήσεων.

Απόσταση Cook:



Γενικά θεωρούμε ως σημεία επιρροής εκείνα με Cook Distance μεγαλύτερο
 του 1. Στο δόθεν διάγραμμα, φαίνεται μόνο η παρατήρηση 30
 να έχει σχετικά υψηλή Cook Distance αλλά ~~και 35~~ δεν μπορεί
 να χαρακτηριστεί σημείο επιρροής αφού $D_{30} < 1$ και κατά το 0.5.

Τελικά, δεν μπορούμε να χαρακτηρίσουμε κάποιο σημείο ως σημείο
 επιρροής ^{με βάση Cook Distance} και χρειαζόμαστε ίσως κάποιο άλλο κριτήριο εύρεσης σημείων
 επιρροής για να εφευρούμε περαιτέρω το σημείο 30 και ενδεχομένως και
 το σημείο 35.