

Απλό γραμμικό μοντέλο

Για x ποσοτική και y ποσοτική εξαρτημένη μεταβλητή

$$E(y|x) = E(y_x) = b_0 + b_1 x = \mu_x$$

$$\Rightarrow \text{Εκτίμηση} : \hat{y}_i = \hat{b}_0 + \hat{b}_1 x_i$$

↙ τυχαίο σφάλμα

$$\text{Προσμενόμενη} : y_i = E(y_{x_i}) + \varepsilon_i = b_0 + b_1 x_i + \varepsilon_i$$

Για τα \hat{b}_0 και \hat{b}_1 : ελαχιστοποίηση του συνολικού τυχαίου σφάλματος

$$S(b_0, b_1) = \sum_{i=1}^n \underbrace{(y_i - E(y_{x_i}))}_{\varepsilon_i}^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

$$\Rightarrow \text{Προκύπτει για } (\hat{b}_0, \hat{b}_1) \text{ ότι} : \frac{\partial S(b_0, b_1)}{\partial b_0} \bigg|_{(\hat{b}_0, \hat{b}_1)} = 0$$

και

$$\frac{\partial S(b_0, b_1)}{\partial b_1} \bigg|_{(\hat{b}_0, \hat{b}_1)} = 0$$

όπου

$$\hat{b}_0 = -\hat{b}_1 \bar{x} + \bar{y}$$

$$\text{και } \hat{b}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x})$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$$

Υποθέτουμε για τυχαία σφάλματα ε_i :

$$1) E(\varepsilon_i) = 0$$

$$2) V(\varepsilon_i) = \sigma^2 \text{ (ομοσκεδαστικότητα)}$$

$$3) \text{cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i, j \text{ αλληλοεξαρτησά - ανεξάρτητα.}$$

$$4) \varepsilon_i \sim N(0, \sigma^2)$$

Κοινωνικές τ.μ.

$$1) Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \Rightarrow E(Y_i) = \mu_i$$

$$\Rightarrow V(Y_i) = V(\varepsilon_i) = \sigma^2$$

$$\Rightarrow \text{cov}(Y_i, Y_j) = \text{cov}(\varepsilon_i, \varepsilon_j) = 0$$

$$Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$$

$$2) \bar{Y} \sim N(\beta_0 + \beta_1 \bar{X}, \frac{\sigma^2}{n})$$

$$3) \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \Rightarrow E(\hat{\beta}_1) \stackrel{\text{anal.}}{=} \frac{1}{S_{xx}} E(S_{xy}) = \beta_1$$

$$\Rightarrow V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}}) \Rightarrow \text{πυκνότητα } \frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{S_{xx}}} \sim N(0, 1)$$

⊗ όπως σ^2 αγνώστου

από τα χρησιμοποιούμενα S^2 , η οποία δεν είναι απερίσπαστη

$$\text{επιμέτρηση} : S^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \stackrel{\text{e.i.}^2}{=} \frac{1}{n-2} \text{SSE} = \frac{1}{n-2} (S_{yy} - \frac{S_{xy}^2}{S_{xx}})$$

⊗ Ισχύει ότι :

$$S_{yy} = SST = SSE + SSR$$

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\bar{Y} - \hat{Y}_i)^2$$

$$** \sum_{i=1}^n (Y_i - \hat{Y}_i)(\hat{Y}_i - \bar{Y}) = 0$$

Αρα έχουμε ότι :

$$\frac{\sigma^2}{S_{xx}} \rightarrow \frac{S^2}{S_{xx}} = \hat{V}(\hat{\beta}_1)$$

εστimation

⊗ ωνικό σφάλμα
 $se(\hat{\beta}_1) = \sqrt{\hat{V}(\hat{\beta}_1)}$

$$\text{και ισχύει ότι } t = \frac{\hat{\beta}_1 - \beta_1}{S/\sqrt{S_{xx}}} \sim t_{n-2}$$

$$\text{⊗ για } \frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-2}$$

(t-test με $H_0: \beta_1 = \beta_1(0)$

και $H_1: \beta_1 \neq \beta_1(0)$

• p-value : $P(|t_{n-2}| > t)$

• για α δόσιμο
εξετάζουμε εάν t
αυτή συν περιοχή
όπου $P(|t| > t_{\alpha/2}) = \alpha$

$$[\hat{\beta}_1 - t_{n-2, \alpha/2} \frac{S}{\sqrt{S_{xx}}}, \hat{\beta}_1 + t_{n-2, \alpha/2} \frac{S}{\sqrt{S_{xx}}}]$$

4) Για των $\hat{\beta}_0 \sim N\left(\beta_0, \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}\right)$ μη ομοσκεδαστική

* προσοχή: οι εκτιμήσεις

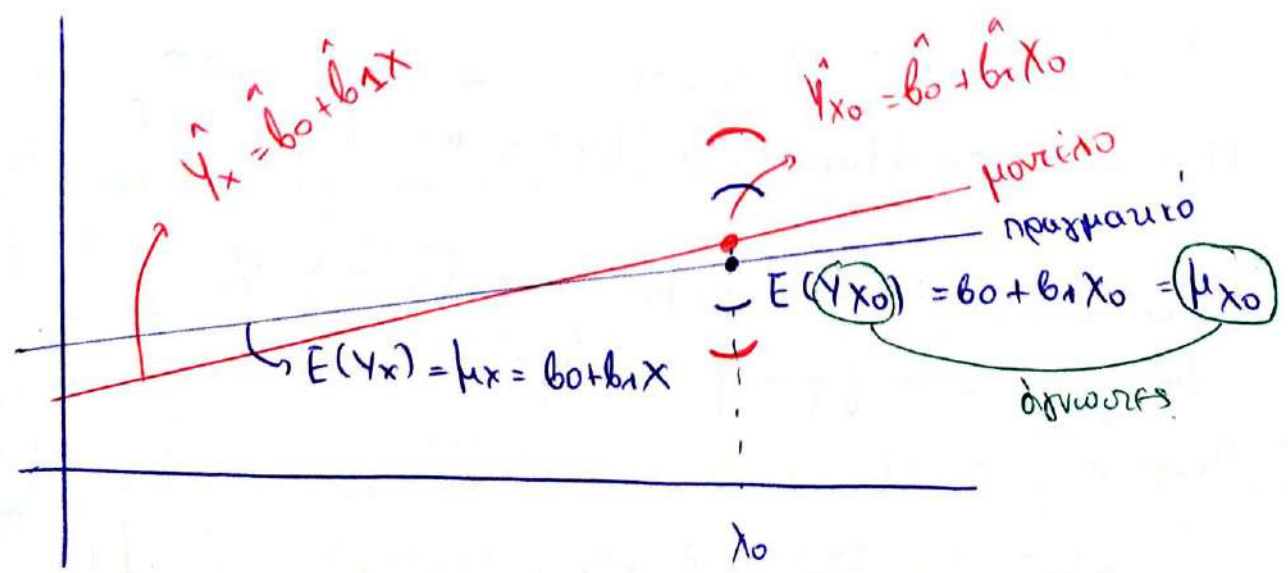
$\hat{\beta}_0$ και $\hat{\beta}_1$ δεν είναι αλληλοεξάρτητες: $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\bar{x}\sigma^2}{S_{xx}}$

5) Για των $\hat{y}_i \sim N\left(\mu_i, \sigma^2 \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}}\right)\right)$ μη ομοσκεδαστική

Συντελεστής προσδιορισμού R^2 : $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$
 ποσοστό μεταβολής τιμών y που εξηγείται από την x . Όταν $R^2 \rightarrow 1$, τότε η εξάρτηση είναι ισχυρή.

Συντελεστής συσχέτισης Pearson: $r_{xy} = \frac{S_{xy}}{\sqrt{S_{xx}}\sqrt{S_{yy}}} \in [-1, 1]$
 Για το από γραμμικό μοντέλο: $r_{xy}^2 = R^2$

Πρόβλεψη σημειακής τιμής \hat{y}_i και μέση τιμής



Διαστήματα εμπιστοσύνης προβλέψεων και μέσης τιμής προβλέψεων

$$E(\hat{Y}_{x_0}) = \mu_{x_0}$$

$$V(\hat{Y}_{x_0}) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

με γνωστούς
αποτελέσματα $\sim N(0,1)$
και αν αντί στα $\sigma \rightarrow S$

με Δ.Ε. : $\left[\hat{Y}_{x_0} \pm t_{\alpha/2} S \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}} \right]$

τότε $t = \frac{\hat{Y}_{x_0} - \mu_{x_0}}{S \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \sim t_{n-2}$

Αντίστοιχα για το $Y_{x_0} : \hat{Y}_{x_0} - Y_{x_0} \sim N\left(0, \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}\right)\right)$

και αντικαθιστώντας το $\sigma \rightarrow S$

$t = \frac{\hat{Y}_{x_0} - Y_{x_0}}{S \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \sim t_{n-2}$

με Δ.Ε. : $\left(\hat{Y}_{x_0} \pm t_{\alpha/2} S \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}} \right)$

Παραγωγή : το διάστημα εμπιστοσύνης της προβλέψεως Y_{x_0}
μεγαλύτερο από της μέσης τιμής μ_{x_0}

Παρατήρηση :

1) Μετατροπή μη γραμμικών σε γραμμικά :

$n \times Y_i = a e^{b x_i} + \varepsilon_i$
 $\Rightarrow \ln Y_i = b_0 + b_1 x_i + \varepsilon_i$

οπότε κάθε τι που υπολογίζουμε

θα πρέπει να το δούμε στο αρχικό

(γ) $1/Y \sim 1/x$

* όπως $E(Y^*) = E(1/Y) \neq 1/E(Y)$

και $E(Y^*) = E(\ln Y) \neq \ln E(Y)$

μόνο γραμμικά
μπορούμε να δώσουμε
τα Δ.Ε.

θα πρέπει
να το
φέρουμε
σε κατάλλη-
λη μορφή
ώστε να έχουμε
ανακάλυψη.

Ερμηνεία Συντελεστών

$$\hat{Y}_0 = \hat{b}_0 + \hat{b}_1 x_0$$

1) \hat{b}_1 : αν αυξηθεί η x κατά μια μονάδα, αναμένεται η y να αυξηθεί κατά b_1 (δεν υπάρχει ελαστικότητα στο x)

2) Αν υπάρχει ανεξέλεγκτος από μεταβλητή x^2 :

$$E(Y_x) = b_0 + b_1 x + b_2 x^2 + \text{άλλα μεταβλητές}$$

αν x αυξηθεί κατά μια μονάδα, τότε η διαφορά συν αναμενόμεν τιμή θα είναι :

$$E(Y_{x+1}) - E(Y_x) = b_1 + b_2 + 2xb_2$$

εξάρτηση από το x

Ελέγχος υποθέσεων άλλων γραμμικών μοντέλων

1) Residuals e_i vs fitted values : θέλουμε να παρατηρήσουμε τάσεις

2) Normal Q-Q plot μεταξύ e_i και \hat{e}_i

α) βρίσκουμε τα ποσοτικά σημεία του δείγματος (για τα e_i)

β) συγκρίνουμε τα ποσοτικά σημεία του \hat{e}_i

θέλουμε να ακολουθούν μια ευθεία.

Πολλαπλό Σταθμικό

Μοντέλο

ΒΑΣΙΚΑ ΣΤΟΙΧΗΑ

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon} \quad \text{με} \quad \underline{\varepsilon} \sim N_n(0, \sigma^2 I_n)$$

⊙ $E(\varepsilon_i) = 0$

⊙ $\text{Var}(\varepsilon_i) = \sigma^2$ ομοσκεδαστικότητα

⊙ $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$

↳ ελαχιστοποίηση $\underline{\varepsilon}$

$$S(\underline{\beta}) = \dots = \sum_{i=1}^n \varepsilon_i^2$$

$$\frac{\partial S(\underline{\beta})}{\partial \underline{\beta}} = 0 \Rightarrow \underline{\hat{\beta}} = (\underline{X}'\underline{X})^{-1} \underline{X}'\underline{Y} \Rightarrow \underline{\hat{Y}} = \underline{X}\underline{\hat{\beta}} = \underline{X}(\underline{X}'\underline{X})^{-1} \underline{X}'\underline{Y}$$

$$\Rightarrow \underline{\hat{Y}} = \underline{H}\underline{Y}$$

* υπολοιπα: $\underline{e} = \underline{Y} - \underline{\hat{Y}} = (\underline{I} - \underline{H})\underline{Y}$
 $= (\underline{I} - \underline{H})(\underline{X}\underline{\beta} + \underline{\varepsilon}) = (\underline{I} - \underline{H})\underline{\varepsilon}$

→ τιτρες προβορής

Τι κατανομές ακολουθούν οι βασικές μεταβλητές;

1) Υπόθεση: $\underline{\varepsilon} \sim N_n(0, \sigma^2 I)$, $\varepsilon_i \sim N(0, \sigma^2)$ κ.λ.π.

2) $\underline{Y} \sim N_n(\underline{X}\underline{\beta}, \sigma^2 I_n)$ ⊙ $V(\underline{Y}) = E((\underline{Y} - \underline{X}\underline{\beta})(\underline{Y} - \underline{X}\underline{\beta})') = \dots = E(\underline{\varepsilon}\underline{\varepsilon}')$
 $= E((\underline{\varepsilon} - E(\underline{\varepsilon}))(\underline{\varepsilon} - E(\underline{\varepsilon}))') = V(\underline{\varepsilon})$

3) για την $\underline{\hat{\beta}} = (\underline{X}'\underline{X})^{-1} \underline{X}'\underline{Y} \rightarrow E(\underline{\hat{\beta}}) = A E(\underline{Y}) = A \underline{X}\underline{\beta} = \underline{\beta}$

αρα $\underline{\hat{\beta}} \sim N_p(\underline{\beta}, \sigma^2 (\underline{X}'\underline{X})^{-1}) \rightarrow V(\underline{\hat{\beta}}) = V(A\underline{Y}) \overset{\text{τιτρες}}{=} AV(\underline{Y})A' = \sigma^2 AA' = \sigma^2 (\underline{X}'\underline{X})^{-1}$
 $\hat{\beta}_j \sim N(\beta_j, \sigma^2 C_{jj})$ με $\text{cov}(\hat{\beta}_j, \hat{\beta}_k) = \sigma^2 C_{jk}$

→ $p = k+1$ → τίτρες = αριθμός παραμέτρων
 ↳ μεταβλητές

όμως σ^2 αγνωστο

↳ ανεξικατορία με $S^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - (p)} = \frac{SSE}{n - p} = \frac{\sum \varepsilon_i^2}{n - p}$ ($E(S^2) = \sigma^2$)
 αο αττό σταθμικό ήταν $p=2$

αρα όταν θα κάνουμε έλεγχο, αυτός θα γίνεται με την

μεταβλητή:
$$t_j = \frac{\hat{\beta}_j - \beta_j}{s\sqrt{C_{jj}}} \sim t_{n-p}$$

και $se(\hat{\beta}_j) = s\sqrt{C_{jj}}$

4)
$$\underline{\hat{Y}} = H\underline{Y} = X\underline{\hat{\beta}} \sim N_n(X\underline{\beta}, \sigma^2 H)$$

* $V(\hat{Y}) = V(HY) = E\{(HY - E(HY))(HY - E(HY))'\} = HV(Y)H' = \sigma^2 HH' = \sigma^2 H$

γιατί H
συμμετρικοί
και αναστρέψιμοι

Προσοχή: η \hat{Y}_i με την \hat{Y}_j συνεισώσασα συσχετισμένες
 $\text{cov}(\hat{Y}_i, \hat{Y}_j) = \text{cov}(\hat{Y}_i, Y_i) = \sigma^2 h_{ij}$

↓

$$\hat{Y}_i = \underline{\beta}' x_i = \sum_{j=1}^n h_{ij} Y_j$$
 και άρα $\hat{Y}_i \sim N(\underline{\beta}' x_i, \sigma^2 h_{ii})$

5) Για τα υπόλοιπα $\underline{e} \sim N_n(\underline{0}, \sigma^2(I-H))$

$$e_i \sim N(0, \sigma^2(1-h_{ii}))$$

δεν υπάρχει πια
η ομοσχευσιμότητα

και πλέον δεν υπάρχει ανεξαρτησία

$$\text{cov}(e_i, e_j) = -\sigma^2 h_{ij}$$

* SOS: δεν χρησιμοποιούμε ως απρόβλεπτη εκτιμήτρια του σ^2 των $\hat{\sigma}^2 = \frac{e'e}{n}$
 αφού $\int \frac{E(e_i^2)}{n} = \int \frac{V(e_i) + E(e_i)^2}{n} = \int \frac{\sigma^2(1-h_{ii})}{n} = \frac{\sigma^2}{n} \sum (1-h_{ii}) \stackrel{*}{=} \frac{\sigma^2}{n} (n-p) = \frac{\sum e_i^2}{n} \neq \sigma^2$

* $\text{tr}(H) = \text{tr}(X(X'X)^{-1}X') \stackrel{\text{tr}(AB)=\text{tr}(BA)}{=} \text{tr}((X'X)(X'X)^{-1}) = \text{tr}(I) = p \neq \sigma^2$

άρα
$$\sum_{i=1}^n h_{ii} = p$$

* τα αποτελέσματα αυτά προκύπτουν από την ανάλυση
ελαχίστων τετραγώνων.

Εκτίμηση Μέγιστης Πιθανοφάνειας

$$l = \ln L = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} (\underline{y} - \underline{X}\underline{b})' (\underline{y} - \underline{X}\underline{b})$$

$$\Rightarrow \frac{\partial l}{\partial \underline{b}} = 0 \Rightarrow \hat{\underline{b}} = (\underline{X}'\underline{X})^{-1} \underline{X}'\underline{y}$$

$$\frac{\partial l}{\partial \sigma^2} = 0 \Rightarrow \hat{\sigma}^2 = \frac{\underline{e}'\underline{e}}{n} = \frac{SSE}{n} \quad \left(\begin{array}{l} \text{δεν είναι μερόσημο} \\ \text{ενώ η } S^2 = \frac{SSE}{n-p} \text{ είναι} \end{array} \right)$$

Μέγιστη πιθανοφάνεια μοντέλων

$$\hat{l} = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln\left(\frac{SSE}{n}\right) - \frac{n}{2} \quad \circ \text{ εξαρτάται από το σφάλμα SSE}$$

Έλεγχος

① Έλεγχος του λόγου των πιθανοφάνων

$$H_0: M_0, H_1: M_1 \text{ με } M_0 < M_1$$

$$-2(\hat{l}_0 - \hat{l}_1) = n \ln\left(\frac{SSE_0}{SSE_1}\right) \sim \chi_d^2 \quad \hookrightarrow \text{διαφορά παραμέτρων.}$$

② Στατιστικός έλεγχος F για $H_0: M_0, H_1: M_1$ με $M_0 < M_1$

$$F^* = \frac{(SSE_0 - SSE_1) / d}{SSE_1 / (n-p)} \sim F_{d, (n-p)} \quad \begin{array}{l} \rightarrow \text{διαφορά αριθμ. παραμέτρων} \\ \rightarrow \text{παραμετροί } H_1 \end{array}$$

απορρίπτουμε την H_0 για μικρά p-values: $p = P(F > F^*)$

③ Στατιστικός έλεγχος $t \equiv F$ έλεγχος για μοντέλο με διαφορά μιας μεταβλητής

$$H_0: M_0 (b_j = 0), H_1: M_1 (b_j \neq 0) \Rightarrow F = \frac{(SSE_0 - SSE_1) / 1}{SSE_1 / (n-p)} = t_{n-p}^2 \sim F_{1, n-p}$$

$$p\text{-value} = P(F > F^*) = P(|t| > t^*) \quad \hookrightarrow \text{όσο πιο μεγάλη } F, \text{ τόσο μικρότερη η p-value και τόσο μεγαλύτερη το σφάλμα SSE}_0$$

Diagnostics with residuals

1) Κανονικά υπόλοιπα $\underline{e} \sim N_n(0, \sigma^2(I-H))$, $e_i \sim N(0, \sigma^2(1-h_{ii}))$

Παρατήρηση 1: Δεν υπάρχει ομοσχεδαικότητα καθώς h_{ii} διαφέρει $\forall i$. (σε αντίθεση με τα ϵ_i που έχουν ομοσχεδαικότητα με σ^2)

$$\begin{aligned} \textcircled{*} e_i &= y_i - \sum_{j=1}^n h_{ij} y_j \\ &= \epsilon_i - \sum_{j=1}^n h_{ij} \epsilon_j \end{aligned}$$

$$\textcircled{*} \text{ισχύει ότι } \frac{1}{n} \leq h_{ii} \leq \frac{1}{c_i}$$

2) Standardized residuals

$$r_i = \frac{e_i}{\sqrt{1-h_{ii}}} \quad \begin{array}{l} \text{έχει} \\ \text{variance 1} \\ \text{αλλά ότι είναι normally distributed} \end{array}$$

$$S^2 = \frac{SSE}{n-p}$$

3) Deleted residuals \equiv standardized residuals όπως κάναμε fit το μοντέλο σε όλα ως παρατηρήσεις εκτός από την i .

$$r_i^! = \frac{e_i}{S_{(i)} \sqrt{1-h_{ii}}} \sim t_{n-p-1}$$

$$S_{(i)}^2 = \frac{SSE(i)}{n-p-1}$$

4) Press residuals: Εάν αφαιρέσουμε την i -οστή παρατήρηση προκύπτει μια επιπλέον εκτίμηση $\hat{y}_{(i)}$

$$\text{press_residual} = y_i - \hat{y}_{(i)} = \frac{e_i}{1-h_{ii}}$$

$\textcircled{*}$ Έλεγχος: Δεν πρέπει να παρατηρήσουμε patterns

α) $\rightarrow e_i \text{ vs } e_{i-1}$

$\rightarrow e_i \text{ vs } i$

$\rightarrow e_i \text{ vs } x_i$

$\rightarrow e_i \text{ vs } \hat{y}_i \rightarrow$ δεν

ισχύει ότι ανεξαρτητως ομοσχεδαικότητας
 $\text{COV}(\underline{e}, \hat{\underline{y}}) = 0$

Όσο πιο μικρό το press residual, τόσο πιο μεγάλη είναι η ικανότητα του μοντέλου να δίνει σωστές προβλέψεις για άγνωστες παρατηρήσεις.

Κριτήρια Επιλογής Μοντέλου

1) Δείχνει $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \leq 1$

Όσο περισσότερες μεταβλητές προσθέτουμε, τόσο θα μικραίνει το SSE. Μας ενδιαφέρουν οι υψηλές αυξομειώσεις για να καταλάβουμε ότι η μεταβλητή που προστίθεται ή αφαιρείται έχει μεγάλη στατιστική εξάρτηση με τον Y

2) Adjusted R^2 : $\bar{R}^2 = 1 - \frac{(SSE)/(n-p)}{SST/(n-1)}$

M_1 : μοντέλο με p μεταβλητές

M_2 : μοντέλο με p+q μεταβλητές

$SSE_2 < SSE_1 \Rightarrow R_2^2 > R_1^2$

αυξάνεται μόνο αν προσέχει μια μεταβλητή που βελτιώνει πολύ το μοντέλο

αν ισχύει $\bar{R}_2 > \bar{R}_1 \Rightarrow \frac{SSE_2}{n-p-q} < \frac{SSE_1}{n-p} \dots \Rightarrow \boxed{\frac{(SSE_1 - SSE_2)/q}{SSE_2/(n-p-q)} > 1}$

Εάν η τιμή ως εξερχομάρζωνος F είναι > 1

τότε η προσθήκη q μεταβλητών είναι σημαντική για το μοντέλο.

3) $R^2_{\text{prediction}}$: βασίζεται στα υπόλοιπα PRESS

$R^2_{\text{pred}} = 1 - \frac{\text{PRESS}}{SST}$ όπου $\text{PRESS} = \sum_{i=1}^n \left(\frac{e_i}{1-h_{ii}} \right)^2$

Όσο πιο μικρή είναι η τιμή PRESS, τόσο λιγότερο επιρροάζει το μοντέλο η αφαίρεση μιας παρατηρήσης $\hat{e} \Rightarrow$ επομένως το μοντέλο έχει υψηλή δυνατότητα πρόβλεψης και άρα τόσο πιο μεγάλο το R^2_{pred}

4) Cp-Mallows

$$C_p = \frac{SSE(p)}{SSE(p')} (n-p') + 2p - n$$

→ υπο εξέταση μοντέλο
→ πλήρη μοντέλο: p' παράμετροι

Θέλουμε $C_p \leq p$ $\otimes C_p = p$ θα το πλήρη μοντέλο.

5) $AIC = -2\hat{\ell} + 2q$

→ αριθμός παραμέτρων μοντέλου

$$AIC = n \ln \left(\frac{SSE}{n} \right) + 2p + n \ln(2\pi) + n + 2$$

(normal $\hat{\ell}$)

\otimes ο έλεγχος F συμφωνεί με το AIC λόγω SSE

→ Θέλουμε να είναι μικρό και παράληλα

Κριτήριο Πολλαπλότητας

οι μεταβλητές να είναι στατιστικά σημαντικές

$$VIF_j = \frac{1}{1 - R_j^2} \quad \text{αν } VIF > 5 \Rightarrow \text{πολλαπλότητα}$$

Κριτήρια για Σημεία Ενρίπσης

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_j - \bar{x})^2}$$

Όταν x_i μακριά από \bar{x} τότε h_{ii} υψηλό

1) $h_{ii} = x_i' (x'x)^{-1} x_i$ αν $h_{ii} > 2p/n$ η παρατήρηση i θεωρείται σημείο ενρίπσης.
(hat values)

2) Απόσταση Cook $D_i = \frac{e_i^2 h_{ii}}{p S^2 (1 - h_{ii})^2} = \frac{r_i^2 h_{ii}}{p(1 - h_{ii})} \gg 1$ τότε θεωρείται σημείο ενρίπσης
απόσταση $\hat{\beta}$ από $\hat{\beta}_{(i)}$
→ έχει αμείωμένη ισχύ παρατήρηση.

3) $DFBETAS_{ji} = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\sqrt{S_{ii}^2 C_{jj}}}$

4

αν $|DFBETAS_{ji}| > 2/\sqrt{n}$ τότε η i-οστή παρατήρηση έχει επιρροή στον συντελεστή β_j

4) $DFFITS_i = \frac{\hat{y}_i - \hat{y}_{(i)}}{\sqrt{S_{ii}^2 h_{ii}}}$, $|DFFITS_i| > 2\sqrt{p/n}$
 τότε η i-οστή παρατήρηση ίσως είναι influential point.

Διαγνωστικά Γραφήματα

για την ανακάλυψη των μεταβλητών

1) Added variable plots

$$y = X\beta + \varepsilon$$

$$y = X\beta + \delta W + \varepsilon^*$$

extra μεταβλητή

$$\text{Αν } e \sim (I-H)W$$

σημαίνει ελθία
 γραφή, τότε

η W είναι απαραίτητη
 στο μοντέλο

$$(y - \hat{y})$$

$$e = \hat{\delta}(I-H)W + (I-H)\varepsilon^*$$

υπόλοιπα
 όταν κάνουμε regression
 της y ως προς X

υπόλοιπα όταν
 κάνουμε
 regression της W
 ως προς X

2) Partial residual plots

$$E(y) = X^+ \beta^+ = X\beta + \delta W$$

$$\text{residual : } e^+ = y - X^+ \hat{\beta}^+$$

$$= (y - X\hat{\beta}) - \hat{\delta}W$$

$$= \tilde{e} - \hat{\delta}W$$

$$\tilde{e} = e^+ + \hat{\delta}W$$

plot : \tilde{e} και W

⇒ αναμένουμε γραμμική σχέση

⇒ αν δεν είναι, τότε η μεταβλητή
 χρειάζεται μετασχηματισμό.

Παρατηρήσεις

Αν έχουμε μοντέλο που περιέχει και indicator μεταβλητές - κατηγορικές

$$\text{π.χ. } Y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 \quad (1)$$

όπου $X_2 = \begin{cases} 0 \\ 1 \end{cases}$ με βάση κάποια συνθήκη που διαχωρίζει τα δεδομένα σε δύο ομάδες.
και $X_3 = X_1 \cdot X_2$ μεταβλητή αλληλεπίδρασης.

→ Μπορούμε να προσεγγίσουμε το μοντέλο απευθείας
→ ή μπορούμε να κάνουμε διαχωρισμό, μίμωσι ή κλίση ως ειδικά παλινδρόμηση, ή η ροή με τους άξονες αλλάζοντας για μια από τις δύο ομάδες.

$$\text{Για } X_2 = 0 : E(Y_I) = b_0 + b_1 X_1 \quad (2)$$

$$\text{Για } X_2 = 1 : E(Y_{II}) = (b_0 + b_3) + (b_1 + b_3) X_1 \quad (3)$$

Με ελέγχους F ελέγχουμε : (ανάλυση ανασυσυν R).

- 1) $H_0: b_3 = 0$, $H_1: b_3 \neq 0$ στο μοντέλο (1)
- 2) Αν $b_3 \neq 0$ τότε υπάρχει αλληλεπίδραση και έτσι δίνουμε δύο διαφορετικές ερωτήσεις
- 3) Αν $b_3 = 0$ (συνήθως σημαίνει η H_0), ελέγχουμε την $H_0: b_2 = 0$ και $H_1: b_2 \neq 0$ στο $Y = b_0 + b_1 X_1 + b_2 X_2$
- 4) Αν $b_2 \neq 0$ τότε έχουμε δύο παράλληλες ευθείες
Αν $b_2 = 0$ τότε έχουμε μια κοινή ευθεία.

Poisson Παλινδρόμηση

- Η τ.μ Y ακολουθεί κατανομή Poisson με σάρωση πυκνότητας πιθανότητας: $f(y) = \frac{e^{-\mu} \mu^y}{y!}$, $y=0,1,2,\dots$

$$E(Y) = \mu$$

$$V(Y) = \mu$$

- Για την παράμετρο μ :
(αναμενόμενη τιμή)

$$\mu \rightarrow \mu_x = e^{x' \underline{\beta}} > 0$$

$$\Rightarrow \ln \mu_x = x' \underline{\beta} = g(\mu_x)$$

$$\text{με } \boxed{g(x) = \ln x}$$

log link function

⊛ στο γραμμικό μοντέλο είχαμε: $y \sim N(\mu_x, \sigma^2)$
 $g(\mu_x) = \mu_x = x' \underline{\beta}$

- Οι παράμετροι $\hat{\underline{\beta}}$ προκύπτουν από την μεγιστοποίηση της πιθανοφάνειας: $L = \prod_{i=1}^n f(y_i, \mu_i)$ όπου $\mu_i = \mu_{x_i} = e^{x_i' \underline{\beta}}$

$$\ell = \ln L = \sum_{i=1}^n [-e^{x_i' \underline{\beta}} + y_i x_i' \underline{\beta} - \ln(y_i!)]$$

$$\frac{\partial}{\partial \underline{\beta}} \Rightarrow$$

$$x' (y - \hat{\mu}) = 0$$

μια γραμμική ως προς x

επαναληπτική μέθοδος

όχι όπως γραμμικό μοντέλο $x' (y - x \hat{\underline{\beta}}) = 0$

- Ερμηνεία Συντελεστών: $e^{\hat{\beta}_j}$ πολλαπλασιαστική αλλαγή των αναμενόμενων τιμών y (στο μ δηλαδή), όταν το x_j αυξάνεται κατά 1 μονάδα. (πρέπει να πάρω το εκθετικό της τιμής που δίνει η R να ως ανεξάρτητες)

Τεταυροί Έλεγχοι Συμπυκνωμένης Μεταβλητών

1) Τεταυροί Έλεγχος Wald

$$\hat{\beta}_j \sim N(\beta_j, \hat{V}(\hat{\beta}_j))$$

$$\frac{\hat{\beta}_j - \beta_j}{(\hat{V}(\hat{\beta}_j))^{1/2}} \sim N(0,1)$$

με $H_0: \beta_j = 0, H_1: \beta_j \neq 0$

→ προκύπτει από την Hessian μήτρα των δεύτερων παραγώγων $\frac{\partial^2 \ell}{\partial \beta_j \partial \beta_j}$

2) Likelihood ratio test : $\hat{L} = \max_{\beta} \ell(\beta_0, \dots, \beta_k)$
 $-2(\hat{\ell}_0 - \hat{\ell}_1) \sim \chi^2_d$ $\hat{\ell}_0 = \max_{\beta} \ell(\beta_0, \dots, \beta_k)$

↳ nested στο M_0 με d διαστάσεις μεταβλητές.

3) Deviance log-likelihood ratio statistic

H_0 : our model ~ restriction $g(\mu_i) = X_i^T \beta$

H_1 : saturated model ~ δα επιβάλλονται περιορισμοί = 0 αριθμός παραμετήρων = αριθμός Παραμετήρων.

Υπολογίζεται η πιθανοφάνεια

για H_0 : $\hat{\mu}_i = \hat{y}_i = e^{X_i^T \hat{\beta}}$

για H_1 : $\tilde{\mu}_i = y_i$

$\tilde{\mu}_i = y_i$ $\frac{\partial \ell}{\partial \mu_i} = 0$

και υπολογίζουμε $D(\hat{\beta}) = -2 \sum \ell(\hat{\beta}) - \tilde{\ell}$

$$D(\hat{\beta}) = 2 \sum_{i=1}^n y_i \ln y_i / \hat{\mu}_i \sim \chi^2_{n-p}$$

$$2 \left\{ \sum_{i=1}^n y_i \ln(y_i / \hat{\mu}_i) - \sum_{i=1}^n (y_i - \hat{\mu}_i) \right\}$$

"0"

Γενικά για $H_0: M^*$ και $H_1: M$ ($M^* \subset M$) : $-2(\ell(\hat{\beta}^*) - \ell(\hat{\beta})) \sim \chi^2_q$
 $P(\hat{\beta}^*) - D(\hat{\beta}) \sim \chi^2_q$
 (Εάν επιβληθεί περιορισμός τύπου $\beta_j = 0$)

Επιλογές Ενδείκτης Ποιότητας

- 1) AIC = $-2 \sum_{i=1}^n [-e^{x_i' \hat{\beta}} + y_i x_i' \hat{\beta} - \ln(y_i!)] + 2q$
- 2) Υπαρξάν random splits R^2 αλλά δεν είναι αξιολογικοί
- 3) Έλεγχος Deviance πορείας σε σχέση με null deviance \Rightarrow δέλωμε να ελευθερώσω το άνω μινιμοζέση.

Residuals

1) Pearson $r_i^P = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{\mu}_i}}$

2) Standardised $r_i^{PS} = \frac{r_i^P}{\sqrt{1 - h_{ii}}}$

όπου h_{ii} diagonal
 $\hat{H} = \hat{W}^{1/2} X (X' \hat{W} X)^{-1} X' \hat{W}^{1/2}$

3) Deviance residual

$$r_i^D = \text{sgn}(y_i - \hat{\mu}_i) [2 \{ y_i \ln(y_i / \hat{\mu}_i) - (y_i - \hat{\mu}_i) \}]^{1/2}$$

$$\Rightarrow D(\hat{\beta}) = \sum_{i=1}^n (r_i^D)^2$$

\Rightarrow Standardize Deviance

4) Likelihood residuals (studentized R) residual
 εξαρτώνται (r_i^{PS}) και (r_i^D)

$$r_i^{PS} = \frac{r_i^D}{\sqrt{1 - h_{ii}}}$$

\Rightarrow Partial residual plots : $r_{ji}^T = \frac{y_i - \hat{\mu}_i}{\hat{\mu}_i} + \hat{\beta}_j x_{ij}$

- αν είναι εύκολο να σπυρξώ
 τότε η x_j είναι ανεξάρτητο
- αν δώ υπάρχει pattern, χρησιμοποιούμε μετασχηματισμούς.

\Rightarrow Added variable plots (ίδιο με partial plots)

Offset συν Παλινδρόμηση Poisson

Πρέπει να δάσουμε υπόψη το μέγεθος του πηλυσμού στον οποίο συμβαίνουν τα γεγονότα Y_i

$$\mu = E(Y) = \text{αναμενόμενος αριθμός γεγονότων} / \text{ανά μονάδα (κατ.)}$$

π.χ για 7 αεροπορικές εταιρείες με κοινά χαρακτηριστικά, ο αναμενόμενος αριθμός ατυχημάτων θα είναι $\mu_i^* = 7\mu_1$

- Αν η αναμενόμενη τιμή μ_i αλλάξει με τα χαρακτηριστικά του ατυχήμενου, τότε το n_i δηλώνει τον αριθμό των ατυχημάτων με τα κοινά χαρακτηριστικά.

$$\text{άρα } \mu_i^* = n_i \mu_i = n_i \exp(x_i' \underline{\theta})$$

$$\ln(\mu_i^*) = \ln(n_i) + \underline{\theta}' x_i$$

$$\Rightarrow \ln(\mu_i^*/n_i) = \underline{\theta}' x_i = \eta_i$$

Logistic regression

* Link function : μας δίνει την γραμμική προβλεπόμενη

- $Y \sim \text{Normal}$: $\eta_x = g(\mu_x) = X'b$ (identity link function)
- $Y \sim \text{Poisson}$: $\eta_x = g(\mu_x) = \ln \mu_x = X'b$ (log link function)
- $Y \sim \text{Bernoulli} (n_x=1)$ ή $Y \sim \text{Binomial} (n_x > 1)$

$$\eta_x = g(\mu_x) = \ln \left(\frac{\mu_x}{n_x - \mu_x} \right) = X'b$$

Στις περιπτώσεις της λογιστικής παλινδρόμησης, η εξαρτημένη μεταβλητή Y είναι διακριτή και ακολουθεί διωνυμική κατανομή. Δηλαδή η συνάρτηση πιθανότητας για Y ενισχίες σε n ανεξάρτητες δοκιμές bernoulli είναι :

$$P(Y) = \binom{n}{y} p^y (1-p)^{n-y} = \exp \left\{ y \ln \left[\frac{p}{1-p} \right] + n \ln(1-p) + \ln \binom{n}{y} \right\}$$

⊛ p : πιθανότητα ενισχίας

μορφή exponential family
όπου $\theta = \ln \left[\frac{p}{1-p} \right]$ δείχνει την link function

Logistic Regression Model :

$$\ln \left(\frac{p_i}{1-p_i} \right) = b_0 + b_1 x_{i1} + \dots + b_k x_{ik}, \quad i=1, \dots, n$$

όπου $p_i = P_{xi} = \frac{\exp(b_0 + \dots)}{1 + \exp(b_0 + \dots)}$

και $\mu_i = E(y_i) = n_i p_i = n_i \frac{e^{x_i' \underline{b}}}{1 + e^{x_i' \underline{b}}}$

Η συνάρτηση πιθανότητας είναι :

$$l = \ln L = \sum_{i=1}^n \left\{ y_i x_i' \underline{b} - n_i \ln(1 + e^{x_i' \underline{b}}) + \ln \binom{n_i}{y_i} \right\}$$

όπου $\frac{\partial l}{\partial b_j} = 0$ μας δίνει αναγκαίως εξισώσεις για τους ανεξάρτητους.

Ερμηνεία Συντελεστών

$$\frac{\hat{p}}{1-\hat{p}} = e^{x' \underline{b}} \quad \text{και} \quad \hat{p} = \frac{e^{x' \underline{b}}}{1 + e^{x' \underline{b}}}$$

- αν $b_j < 0$ τότε ο λόγος odds θα υπολογιστεί κατά $e^{\hat{b}_j}$
- αν $b_j > 0$ τότε ο λόγος odds θα πολλαπλασιαστεί κατά $e^{\hat{b}_j}$ όταν η μεταβλητή x_j αυξηθεί κατά μια μονάδα.

⊗ Αν μια μεταβλητή είναι η lux και όχι η x, η ερμηνεία αλλάζει.

Αν η η x αυξηθεί κατά ένα ποσοστό π.χ. 10%. τότε :

$$\Delta(\log(\text{odds})) = b [\ln(x+0.1x) - \ln x] = b \ln(1.1x/x) = b \cdot \ln(1.1)$$

άρα η πολλαπλασιαστική μεταβολή των odds θα γίνει $\exp(b \ln(1.1))$

Odds ratio

$$\frac{\frac{\hat{p}_1}{1-\hat{p}_1}}{\frac{\hat{p}_2}{1-\hat{p}_2}} = \exp(\chi_1 - \chi_2)' \underline{\hat{\beta}}$$

και εαν υποθέσουμε ότι έχουμε μια ποσοτική μεταβλητή χ_2 και μια κατηγορική χ_1 τότε :

$$= \frac{\exp(\beta_0 + \hat{\beta}_2 \chi_2)}{\exp(\beta_0 + \hat{\beta}_1 + \hat{\beta}_2 \chi_2)} \rightsquigarrow \text{για οποιοδήποτε } \chi_2 \text{ το odds ratio είναι } e^{\hat{\beta}_1}$$

Έλεγχος

1) Wald statistic : $\underline{\hat{\beta}} \sim N_p(\underline{\beta}, \hat{V}(\hat{\beta}))$

κατάλληλη προσέγγιση

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

όπου $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$ με

$100(1-\alpha)\% \text{ ΔΕ.}$
 $\hat{\beta}_j \pm z_{\alpha/2} se(\hat{\beta}_j)$

→ και για το odds ratio : $\exp(\hat{\beta}_j \pm z_{\alpha/2} se(\hat{\beta}_j))$

2) Deviance

M_0 (μοντέλο με περιορισμούς) } $M_0 \subset M_1$
 M_1 πλήρες μοντέλο

Διαφορά Deviance : $D_0 - D_1 = -2(\hat{\ell}_0 - \hat{\ell}_s) + 2(\hat{\ell}_1 - \hat{\ell}_s)$
 $= -2(\hat{\ell}_0 - \hat{\ell}_1) \sim \chi^2_{p_1 - p_2}$

* πάντα δεικό που M_0 έχει λιγότερες παραμέτρους από M_1 και άρα (deviates) ανώτερη πληρότητα από το saturated

* όσο πιο μεγάλη η διαφορά, τόσο πιο πιθανό να απορριψουμε το M_0 .

→ διαφορά παραμέτρων

Deviance

Binomial Data

$$Y_i \sim b(n_i, p)$$

$$\mu_i = E(Y_i) = n_i p_i$$

$$D(\hat{\beta}) = -2 \sum_{i=1}^n Y_i \ln\left(\frac{\hat{p}_i}{p_i}\right) + (n_i - Y_i) \ln\left(\frac{1 - \hat{p}_i}{1 - p_i}\right)$$

$$\text{για } \hat{p}_i = \frac{Y_i}{n_i} \text{ και } \hat{\mu}_i = \frac{Y_i}{n_i}$$

Παρατήρηση: υπάρχει εξάρτηση ως D

από των προβλεπόμεν τιμή $\hat{\mu}_i$ και των

Παρατηρήσει Y_i . Όποτε η deviance

ένος μοντέλου για τα binomial

data έχει νόημα να των εξετάσουμε

και μόνον της. Όσο πιο χαμηλή, τόσο

πιο κοντά είναι η παρατηρήσει Y_i

στον προβλεπόμεν $\hat{\mu}_i$

Binary Data

$$Y \sim B(p_i)$$

$$\mu_i = E(Y_i) = p_i$$

$$D(\hat{\beta}) = -2 \sum_{i=1}^n \{ \hat{\mu}_i \logit(\hat{\mu}_i) + \ln(1 - \hat{\mu}_i) \}$$

Παρατήρηση: η D

εξαρτάται μόνο από των προβλεπόμεν τιμή $\hat{\mu}_i$.

Αρα δεν έχει νόημα να των εξετάσουμε

μόνον της. Μόνο με

σχεσίση των deviance δύο μοντέλων.

Επιλογή Μοντέλου

$$AIC = 2 \sum_{i=1}^n \left[n_i \ln(1 + e^{X_i' \hat{\beta}}) - Y_i X_i' \hat{\beta} - \ln\left(\frac{n_i}{Y_i}\right) \right] + 2p$$

Diagnostics with residuals: ισχύουν τα αντιστοιχα.
 * Χρήση half normal QQ plot για ανίχνευση outliers

Likelihood residuals: συνδυασμός των pearson και deviance

residuals. Προσέχεται των μεταβολή της deviance αν παρατηρείται η i-οση παρατήρηση. Όσο πιο μεγάλη μεταβολή, τόσο πιο πιθανό να είναι outlier.