

MEM-205 Περιγραφική Στατιστική
Τμήμα Μαθηματικών και Εφ. Μαθηματικών, Πανεπιστήμιο Κρήτης

Κώστας Σμαραγδάκης (kesmarag@gmail.com)

Θεωρία 8ης εβδομάδας

x
y

Αιτιοκρατικό μοντέλο

$$y = A + Bx$$

Πιθανοθεωρητικό μοντέλο - Μοντέλο απλής γραμμικής παλινδρόμησης

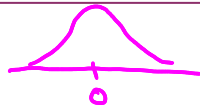
$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$y = A + Bx + \epsilon, \quad \epsilon : \text{όρος τυχαίου σφάλματος}$$

$\uparrow \quad \uparrow \quad \uparrow$
 $\alpha \quad b \quad 0$
 $A : \text{σταθερός όρος (constant term)}, \quad B : \text{κλίση (slope)}$

$$\hat{y} = \alpha + bx$$

$$y = A + Bx + \varepsilon, \quad \varepsilon \sim N(0, \sigma_\varepsilon^2)$$
$$\hat{y} = \alpha + bx$$

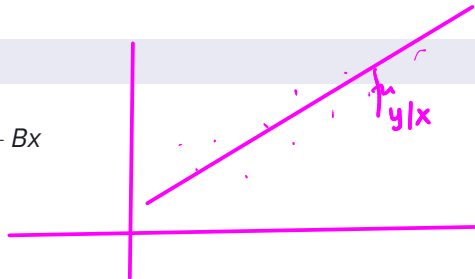


Παραδοχές

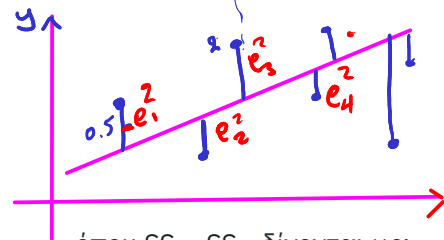
- ▶ Για δοσμένο x το ε ακολουθεί κανονική κατανομή με μηδενική μέση τιμή.
- ▶ Τα τυχαία σφάλματα διαφορετικών παρατηρήσεων είναι ανεξάρτητα.
- ▶ Για κάθε x οι κατανομές των τυχαίων σφαλμάτων παρουσιάζουν την ίδια τυπική απόκλιση.

Ευθεία παλινδρόμησης για τον πληθυσμό

$$\mu_{y|x} = A + Bx$$



Απλή Γραμμική Παλινδρόμηση - Εκτίμηση Ελαχίστων Τετραγώνων



$$\hat{y} = a + bx$$

$$b = \frac{SS_{xy}}{SS_{xx}}, \quad a = \bar{Y} - b\bar{X}$$

όπου SS_{xy} , SS_{xx} δίνονται ως:

$$SS_{yy} = \sum y_n^2 - \frac{(\sum y_n)^2}{N}$$

$$SS_{xy} = \sum_{n=1}^N x_n y_n - \frac{(\sum_{n=1}^N x_n)(\sum_{n=1}^N y_n)}{N}, \quad SS_{xx} = \sum_{n=1}^N x_n^2 - \frac{(\sum_{n=1}^N x_n)^2}{N}$$

Επιπλέον τα SS_{xy} και SS_{xx} μπορούν ισοδύναμα να υπολογισθούν ως:

$$SS_{xy} = \sum_{n=1}^N (x_n - \bar{X})(y_n - \bar{Y}), \quad SS_{xx} = \sum_{n=1}^N (x_n - \bar{X})^2$$

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

$$\sum (y_n - \hat{y}(x_n))^2$$

$$y = A + Bx + \epsilon, \quad \epsilon : \text{όρος τυχαίου σφάλματος}$$

- Για κάθε x έχουμε υποθέσει ότι το σφάλμα ϵ ακολουθεί την κανονική κατανομή $\mathcal{N}(0, \sigma_\epsilon)$.
- Η τυπική απόκλιση σ_ϵ του τυχαίου σφάλματος αναφέρεται στο πληθυσμό και κατά επέκταση η τιμή της δεν είναι γνωστή στις περισσότερες περιπτώσεις.

Εκτιμήτρια της τυπικής απόκλισης των σφαλμάτων

$$e_n = y_n - \hat{y}_n$$

$$s_e = \sqrt{\frac{SSE}{N-2}}, \quad \text{SSE} = \sum_{n=1}^N (y_n - \hat{y}_n)^2 = \underline{SS_{yy}} - \underline{b} \underline{SS_{xy}}$$

$\sum (e_n - \bar{e})^2 = SSE$

$\bar{x} \quad \mu_x$
 $\bar{y} \quad \mu_y$

Διαστήματα εμπιστοσύνης.

Συνολικό άθροισμα τετραγώνων

$$SST = \sum_{n=1}^N (y_n - \bar{Y})^2 = SS_{yy}$$

Άθροισμα τετραγώνων παλινδρόμησης

$$SSR = \sum_{n=1}^N (\hat{y}_n - \bar{Y})^2$$

Συντελεστής Προσδιορισμού

$$R^2 = \frac{SSR}{\underline{SST}}, \quad 0 \leq R^2 \leq 1$$

- Ποσοτικοποιεί την αποτελεσματικότητα του μοντέλου.

$$\sum (y_n - \bar{y})^2 =$$

$$= \sum (\underbrace{y_n - \hat{y}_n}_{\text{SSE}} + \underbrace{\hat{y}_n - \bar{y}}_{\text{SSR}})^2$$

$$\overset{SS_{yy}}{\downarrow} \geq 0$$

$$SST = SSR + SSE$$

$$b = \frac{SS_{xy}}{SS_{xx}}$$

$$b = \frac{SS_{xy}}{SS_{yy}}$$

$$R^2 = \frac{SST - SSE}{SST} = \frac{b SS_{xy}}{SS_{yy}}, \quad 0 \leq R^2 \leq 1$$

$$= \sum (y_n - y_n)^2 + \sum (y_n - \bar{y})^2 + 2 \sum (y_n - \hat{y}_n)(\hat{y}_n - \bar{y})$$

Αντικαθιστώντας τη τιμή του b έχουμε το R^2 στη μορφή:

$$R^2 = \frac{SS_{xy}^2}{SS_{xx}SS_{yy}}$$

$$\sum (x_n - \bar{x})^2 \quad \sum (y_n - \bar{y})^2$$

Συντελεστής Γραμμικής Συσχέτισης - Pearson

- Συμβολίζεται με ρ όταν αφορά τον πληθυσμό.

$$\rho \in [-1, 1]$$

$$\sum (x_i - \bar{x})(y_i - \bar{y})$$

- Συμβολίζεται με r όταν αφορά ένα δείγμα.

$$r \in [-1, 1]$$

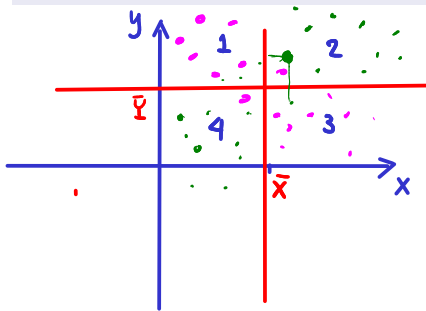
$$r = \frac{SS_{xy}}{\sqrt{SS_{xx}SS_{yy}}}$$

Γραμμική Συσχέτιση (Linear Correlation)

$$R^2 = \frac{SS_{xy}^2}{SS_{xx}SS_{yy}} \quad (\text{Συντελεστής Προσδιορισμού})$$

$$r = \frac{SS_{xy}}{\sqrt{SS_{xx}SS_{yy}}} \quad (\text{Συντελεστής Γραμμικής Συσχέτισης})$$

Σχέση μεταξύ συντελεστών γραμμικής συσχέτισης και προσδιορισμού



$$r = \text{sign}(SS_{xy})\sqrt{R^2}$$

$$x = \text{sign}(x) \cdot |x|$$

$$\sum (x_n - \bar{x})(y_n - \bar{y})$$

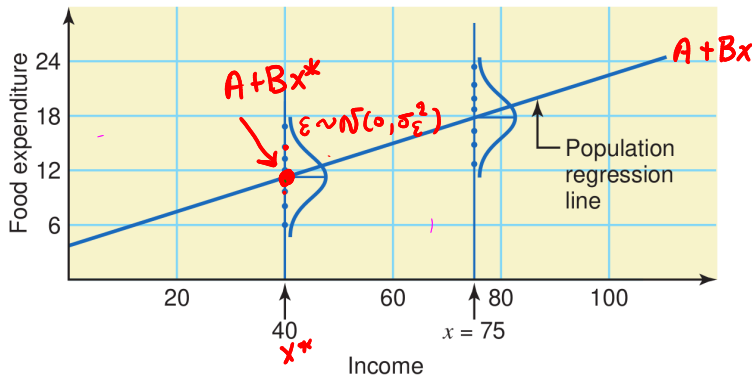
$$\text{sign}(x) = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0 \end{cases}$$

Διαστήματα εμπιστοσύνης για τις τιμές της εξαρτημένης μεταβλητής

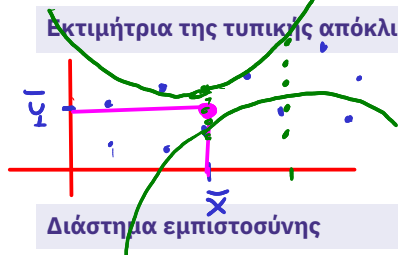
1. Για δοσμένο x^* ποιο είναι το διάστημα εμπιστοσύνης $(1-\alpha)*100\%$ για τη μέση τιμή $\mu_{y|x^*}$;
2. Για δοσμένο x^* ποιο είναι το διάστημα εμπιστοσύνης $(1-\alpha)*100\%$ για την τιμή μιας συγκεκριμένης παρατήρησης y^* ;

$$\alpha = 0.05$$

$$y = A + Bx + \varepsilon$$



Διάστημα Εμπιστοσύνης για την εκτίμηση της $\mu_{y|x^*}$



$$s_{\hat{\mu}_{y|x^*}} = s_e \sqrt{\frac{1}{N} + \frac{(x^* - \bar{X})^2}{SS_{xx}}}$$

$$\begin{aligned} \hat{y}_{1x^*} &= A + Bx^*, \quad \sigma_{\mu_{y|x^*}} \\ \hat{\mu}_{y|x^*} &= a + bx^* \\ \underline{\underline{S_{\hat{\mu}_{y|x^*}}}}} \end{aligned}$$

Το $(1 - \alpha) * 100\%$ διάστημα εμπιστοσύνης για την $\mu_{y|x^*}$ είναι:

95%

$$[\hat{\mu}_{y|x^*} - ts_{\hat{\mu}_{y|x^*}}, \hat{\mu}_{y|x^*} + ts_{\hat{\mu}_{y|x^*}}]$$

όπου το t λαμβάνεται από την t_{df} , $df = \underline{N - 2}$ έτσι ώστε

$$P(T < t) = \underline{1 - \alpha/2}$$

► Περιθώριο σφάλματος: $E = ts_{\hat{\mu}_{y|x^*}}$

Εκτιμήτρια της τυπικής απόκλιση του \hat{y}^*

$$s_{\hat{y}^*} = s_e \sqrt{1 + \frac{1}{N} + \frac{(x^* - \bar{X})^2}{SS_{xx}}}$$

Διάστημα εμπιστοσύνης

Το $(1 - \alpha) * 100\%$ διάστημα εμπιστοσύνης για την y^* είναι:

$$[\hat{y}^* - ts_{\hat{y}^*}, \hat{y}^* + ts_{\hat{y}^*}]$$

όπου το t λαμβάνεται από την t_{df} , $df = N - 2$ έτσι ώστε

$$P(T < t) = 1 - \alpha/2$$

► Περιθώριο σφάλματος: $E = ts_{\hat{y}^*}$

$$y = A + Bx^* + \varepsilon = \mu_{y|x^*} + \varepsilon$$

X, Y ανεξάρτητες
τ.μ.

$$Z = X + Y$$

$$E\{Z\} = E\{X\} + E\{Y\}$$

$$\sigma_Z^2 = \sigma_X^2 + \sigma_Y^2$$

$$\sigma_y^2 = \sigma_\varepsilon^2 \left(\frac{1}{N} + \frac{(x^* - \bar{X})^2}{SS_{xx}} \right) + \sigma_\varepsilon^2$$

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \quad \text{από την (1D)}$$

$$(x_1^{(1)}, \dots, x_1^{(K)}, y_1), (x_2^{(1)}, \dots, x_2^{(K)}, y_2)$$

$$y = A + \mathbf{x}^T \mathbf{B} + \epsilon$$

$$\mathbf{x} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(K)} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} B^{(1)} \\ B^{(2)} \\ \vdots \\ B^{(K)} \end{bmatrix}$$

$$y = A + \sum_{j=1}^K x^{(j)} B^{(j)}$$

$$y = A + B^{(1)}x^{(1)} + \dots + B^{(K)}x^{(K)} + \epsilon$$

Ευθεία παλινδρόμησης για τον πληθυσμό

$$\mu_{y|\mathbf{x}} = A + \mathbf{x}^T \mathbf{B}$$

Δειγματικό μοντέλο απλής γραμμικής παλινδρόμησης

$$\hat{y} = a + \mathbf{x}^T \mathbf{b}$$

- ▶ a είναι δειγματική προσέγγιση του A
- ▶ $\mathbf{b} = [b^{(1)}, b^{(2)}, \dots, b^{(K)}]^T$ είναι δειγματική προσέγγιση του \mathbf{B}
- ▶ \hat{y} είναι η εκτιμώμενη τιμή του y για δοσμένο $\mathbf{x} = [x^{(1)}, x^{(2)}, \dots, x^{(K)}]^T$

Τυχαίο σφάλμα του δειγματικού μοντέλου απλής γραμμικής παλινδρόμησης

$$e = y - \hat{y}$$

Έστω το τυχαίο δείγμα

$$\{(x_1^{(1)}, \dots, x_1^{(K)}, y_1), (x_2^{(1)}, \dots, x_2^{(K)}, y_2), \dots, (x_N^{(1)}, \dots, x_N^{(K)}, y_N)\}$$

$\hat{y}_1 \quad \hat{y}_2 \quad \alpha, \hat{b} \rightarrow \hat{y}_N$

Για το τυχαίο σφάλμα του δειγματικού μοντέλου πολλαπλής γραμμικής παλινδρόμησης έχουμε:

$$e_n = y_n - \hat{y}_n, \quad n = 1, \dots, N$$

όπου η προσέγγιση του κάθε y_n δίνεται ως

$$\hat{y}_n = a + \mathbf{x}_n^T \mathbf{b}$$

Άθροισμα τετραγωνικών σφαλμάτων

$$\text{SSE} = \sum_{n=1}^N e_n^2$$

$$\mathbf{p} = \begin{bmatrix} a \\ b^{(1)} \\ b^{(2)} \\ \vdots \\ b^{(K)} \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(K)} \\ 1 & x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(K)} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_N^{(1)} & x_N^{(2)} & \dots & x_N^{(K)} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Προσέγγιση ελαχίστων τετραγώνων

$$\mathcal{Q}(\mathbf{p}) = \mathbf{e}^T \mathbf{e} = \mathbf{y}^T \mathbf{y} - 2\mathbf{p}^T \mathbf{X}^T \mathbf{y} + \mathbf{p}^T \mathbf{X}^T \mathbf{X} \mathbf{p}$$

$$\mathbf{p} = \arg \min_{\mathbf{p}'} \mathcal{Q}(\mathbf{p}')$$

$$\mathbf{p} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$e_n = y_n - \hat{y}_n = y_n - \alpha - \tilde{x}_n^T \tilde{b}, \quad n = 1, \dots, N$$

$$\begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} \alpha \\ \alpha \\ \vdots \\ \alpha \end{bmatrix} - \underbrace{\begin{bmatrix} x_1^{(1)} & \dots & x_1^{(k)} \\ x_2^{(1)} & \dots & x_2^{(k)} \\ \vdots & \vdots & \vdots \\ x_N^{(1)} & \dots & x_N^{(k)} \end{bmatrix}}_{X} \begin{bmatrix} b^{(1)} \\ \vdots \\ b^{(k)} \end{bmatrix}$$

$\tilde{e} \quad \quad \tilde{y}$

$$\boxed{e = y - Xp}$$

$$Q = \sum e_n^2 = e^T e = (y - Xp)^T (y - Xp) =$$

$\tilde{x} \in \mathbb{R}^{N, k+1}$

$\begin{bmatrix} 1 & x_1^{(1)} & \dots & x_1^{(k)} \\ 1 & x_2^{(1)} & \dots & x_2^{(k)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N^{(1)} & \dots & x_N^{(k)} \end{bmatrix} \in \mathbb{R}^{N, k+1}$

$\begin{bmatrix} \alpha \\ b^{(1)} \\ \vdots \\ b^{(k)} \end{bmatrix} \in \mathbb{R}^{k+1, 1}$

\tilde{p}

$$= y^T y - \underbrace{y^T X p}_{y \cdot (Xp)} - \underbrace{p^T X^T y}_{(Xp)^T y} + p^T X^T X p =$$

p_j
 $\begin{cases} \rightarrow j=1 \rightarrow a \\ \rightarrow j \neq 1 \rightarrow b^{(i)} \end{cases}$

$$= y^T y - 2 p^T X^T y + p^T X^T X p$$

$$\frac{\partial Q}{\partial p_j} = -2 \frac{\partial}{\partial p_j} (p^T X^T y) + \frac{\partial}{\partial p_j} (p^T X^T X p)$$

$$\frac{\partial}{\partial p_j} (p^T X^T y) = \frac{\partial}{\partial p_j} \left(\sum_i p_i (X^T y)_i \right) = (X^T y)_j \left(\frac{\partial p_j}{\partial p_j} \right) = \underline{(X^T y)_j}$$

$p^T X^T y = p \cdot (X^T y)$
 $= \sum_{i=1}^n p_i (X^T y)_i$

$$\frac{\partial}{\partial p_j} (p^T X^T X p) = \begin{bmatrix} \text{---} \end{bmatrix} \begin{bmatrix} \text{---} \end{bmatrix} \quad p^T X^T X p = (X p)^T X p =$$

$$= \frac{\partial}{\partial p_j} \left(\sum_{i=1}^N \left(\sum_{q=1}^K X_{iq} p_q \right)^2 \right) = \quad = (X p) \cdot (X p) =$$

$$2 \sum_{i=1}^N \left(\sum_{q=1}^K X_{iq} p_q \right) X_{ij} = 2 \left[(X^T X) p \right]_j$$

$$-2 (X^T y) + 2 (X^T X) p = 0$$

$$(X^T X) p = X^T y \Rightarrow p = (X^T X)^{-1} X^T y$$

Παράδειγμα

Να βρεθεί το δειγματικό μοντέλο γραμμικής παλινδρόμησης για το σύνολο δεδομένων

$$y = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 2 \end{bmatrix} \in \mathbb{R}^4 \quad \{(1, -1, 1), (0, -1, -1), (2, 0, 2), (1, 1, 2)\}$$

$$X = \begin{bmatrix} 1 & 1 & -1 \\ 1 & 0 & -1 \\ 1 & 2 & 0 \\ 1 & 1 & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 3}$$

$$P = \begin{bmatrix} \alpha \\ b^{(1)} \\ b^{(2)} \end{bmatrix} \in \mathbb{R}^3$$

$$P = \underbrace{(X^T X)^{-1}}_{\in \mathbb{R}^{3 \times 3}} X^T \cdot y \quad \begin{matrix} \in \mathbb{R}^{3 \times 4} \\ \in \mathbb{R}^{4 \times 1} \end{matrix} \in \mathbb{R}^3$$

$(0, 2, ?)$

Άσκηση

Δείξτε ότι η εκτίμηση ελαχίστων τετραγώνων

$$\mathbf{p} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

στη περίπτωση της απλής γραμμικής παλινδρόμησης οδηγεί, όπως περιμένουμε, στις εκτιμήσεις των παραμέτρων:

$$b = \frac{SS_{xy}}{SS_{xx}}, \quad a = \bar{Y} - b\bar{X}$$

Εκτιμήτρια της τυπικής απόκλιση του \hat{y}^*

$$s_{\hat{y}^*} = s_e \sqrt{1 + \frac{1}{N} + \frac{(x^* - \bar{X})^2}{SS_{xx}}}$$

Διάστημα εμπιστοσύνης

Το $(1 - \alpha) * 100\%$ διάστημα εμπιστοσύνης για την y^* είναι:

$$[\hat{y}^* - ts_{\hat{y}^*}, \hat{y}^* + ts_{\hat{y}^*}]$$

όπου το t λαμβάνεται από την t_{df} , $df = N - 2$ έτσι ώστε

$$P(T < t) = 1 - \alpha/2$$

- Περιθώριο σφάλματος: $E = ts_{\hat{y}^*}$