

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

Κώστας Σμαραγδάκης (kesmarag@gmail.com)

23-03-2020

Αιτιοκρατικό μοντέλο

$$y = A + Bx$$

Πιθανοθεωρητικό μοντέλο - Μοντέλο απλής γραμμικής παλινδρόμησης

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

A : σταθερός όρος (constant term), B : κλίση (slope)

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

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

Απλή Γραμμική Παλινδρόμηση - Εκτίμηση Ελαχίστων Τετραγώνων

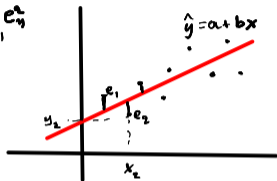
$$y = \mathbf{A} + \mathbf{B}x + \epsilon$$

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

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

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

$$\min_{a, b} \sum_{n=1}^N e_n^2$$



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

όπου SS_{xy} , SS_{xx} δίνονται ως:

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

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

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

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

- ▶ 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.

Εκτιμητήρια της τυπικής απόκλισης των σφαλμάτων

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

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

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

- Ποσοτικοποιεί την αποτελεσματικότητα του μοντέλου.

$$SST = SSR + SSE$$

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

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

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

Συντελεστής Γραμμικής Συσχέτισης - Pearson

- ▶ Συμβολίζεται με ρ όταν αφορά τον πληθυσμό.

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

- ▶ Συμβολίζεται με r όταν αφορά ένα δείγμα.

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

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

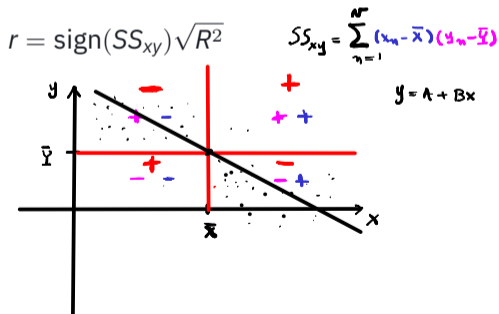
Γραμμική Συσχέτιση (Linear Correlation)

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

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

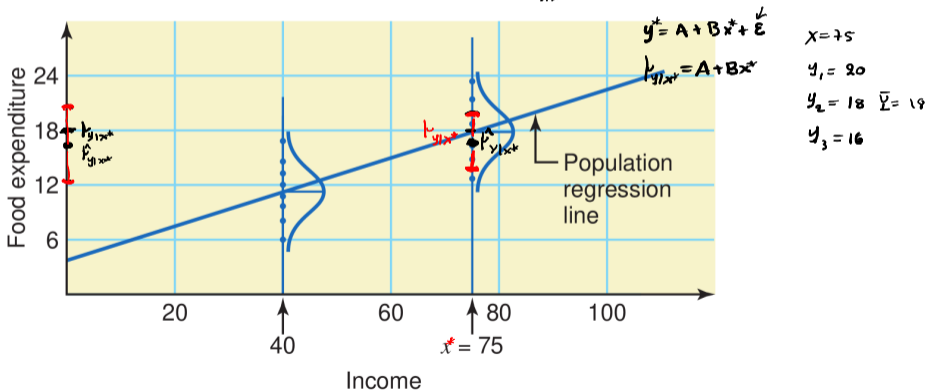
Σχέση μεταξύ συντελεστών γραμμικής συσχέτισης και προσδιορισμού

$$-1 \leq r \leq 1$$



Διαστήματα εμπιστοσύνης για τις τιμές της εξαρτημένης μεταβλητής

1. Για δοσμένο x^* ποιο είναι το διάστημα εμπιστοσύνης $(1-\alpha)^*100\%$ για τη μέση τιμή $\mu_{y|x^*}$; $x^* = 75$ $\mu_{y|x^*} = 19$ $\hat{y}_{y|x^*} = a + bx^*$ $S_{\hat{y}_{y|x^*}}$
2. Για δοσμένο x^* ποιο είναι το διάστημα εμπιστοσύνης $(1-\alpha)^*100\%$ για την τιμή μιας συγκεκριμένης παρατήρησης y^* ; $\hat{y}^* = a + bx^*$ $S_{y^*} > S_{\hat{y}^*}$



Διαστήματα εμπιστοσύνης για τις τιμές της εξαρτημένης μεταβλητής

$$y_i = A + Bx_i + \varepsilon_i, \quad i=1, \dots, N \quad \frac{1}{N} \sum_{i=1}^N y_i = \frac{1}{N} \sum_{i=1}^N A + \frac{1}{N} \sum_{i=1}^N Bx_i + \frac{1}{N} \sum_{i=1}^N \varepsilon_i \Leftrightarrow \bar{Y} = A + B\bar{X} + \bar{\varepsilon} \quad \mathbb{E}\{\bar{Y}\} = A + B\mathbb{E}\{\bar{X}\} + \mathbb{E}\{\bar{\varepsilon}\} = A + B\mu_x$$

$$\hat{y}_{y|x^*} = a + bx^*, \quad a = \bar{Y} - b\bar{X} \rightarrow \hat{\beta}_{y|x^*} = \bar{Y} + b(x^* - \bar{X})$$

$$G_x^2 = \text{var}\{X\}, \quad \text{cov}\{X, Y\} = \mathbb{E}\{(X - \mathbb{E}\{X\})(Y - \mathbb{E}\{Y\})\} \quad \text{var}\{X+Y\} = \text{var}\{X\} + \text{var}\{Y\} + 2\text{cov}\{X, Y\} \quad \text{var}\{aX\} = a^2 \text{var}\{X\}$$

$$\text{var}\{\hat{\beta}_{y|x^*}\} = \text{var}\{\bar{Y}\} + (x^* - \bar{X})^2 \text{var}\{b\} + 2(x^* - \bar{X})\text{cov}\{\bar{Y}, b\}$$

$$\text{var}\{\bar{Y}\} = \text{var}\left\{\frac{1}{N} \sum_{i=1}^N y_i\right\} = \frac{1}{N^2} \sum_{i=1}^N \text{var}\{y_i\} = \frac{1}{N^2} G_y^2 N = \frac{G_y^2}{N}, \quad \text{var}\{b\} = \frac{G_y^2}{SS_{xx}}$$

$$b = \frac{SS_{xy}}{SS_{xx}}, \quad SS_{xy} = \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^N (x_i - \bar{x})y_i - \bar{y} \sum_{i=1}^N (x_i - \bar{x}) = \sum_{i=1}^N (x_i - \bar{x})y_i, \quad b = \sum_{i=1}^N c_i y_i, \quad c_i = \frac{x_i - \bar{x}}{SS_{xx}}$$

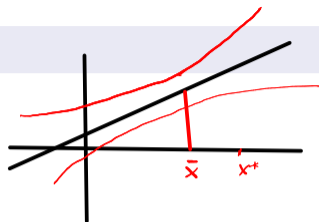
$$\sum_{i=1}^N c_i = 0, \quad \sum_{i=1}^N c_i x_i = 1 \quad (\text{άσκηση})$$

$$\begin{aligned} \text{cov}\{\bar{Y}, b\} &= \mathbb{E}\{(\bar{Y} - \mathbb{E}\{\bar{Y}\})(b - \mathbb{E}\{b\})\} = \mathbb{E}\{\bar{\varepsilon} \cdot (b - B)\} = \mathbb{E}\{\bar{\varepsilon} \cdot (\sum_{i=1}^N c_i y_i - B)\} = \mathbb{E}\{\bar{\varepsilon} \cdot (\sum_{i=1}^N c_i (A + Bx_i + \varepsilon_i) - B)\} = \\ &= \mathbb{E}\{\bar{\varepsilon} (A \sum_{i=1}^N c_i + B \sum_{i=1}^N c_i x_i + \sum_{i=1}^N c_i \varepsilon_i - B)\} = \mathbb{E}\{\bar{\varepsilon} (\sum_{i=1}^N c_i \varepsilon_i)\} = \mathbb{E}\{\bar{\varepsilon} \sum_{i=1}^N c_i \varepsilon_i\} = \sum_{i=1}^N c_i \mathbb{E}\{\bar{\varepsilon} \varepsilon_i\} = \\ &= \sum_{i=1}^N c_i \mathbb{E}\{\varepsilon_i\} \frac{1}{N} \sum_{m=1}^N \mathbb{E}\{\varepsilon_m\} = \sum_{i=1}^N c_i \mathbb{E}\left\{\frac{1}{N} \sum_{m=1}^N \varepsilon_m\right\} = G_y^2 \sum_{i=1}^N c_i = 0. \end{aligned}$$

$$\text{var}\{\hat{\beta}_{y|x^*}\} = G_{\hat{\beta}_{y|x^*}}^2 = G_y^2 \left(\frac{1}{N} + \frac{(x^* - \bar{x})^2}{SS_{xx}} \right)$$

Εκτιμήτρια της τυπικής απόκλιση του $\hat{\mu}_{y|x^*}$

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



Διάστημα εμπιστοσύνης

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

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

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

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

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

Διάστημα Εμπιστοσύνης για την εκτίμηση συγκεκριμένης τιμής της y

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

$$\hat{y}^* = a + b x^*$$
$$y^k = \underbrace{a + b x^k}_{\hat{y}^k} + \varepsilon^k$$

$$G_{\hat{y}^k | x^k} = G_{\varepsilon}^2 \left(\frac{1}{N} + \frac{(x^k - \bar{X})^2}{SS_{xx}} \right)$$

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

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

Διάστημα εμπιστοσύνης

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

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

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

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

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

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

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

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

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

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

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

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

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

Παράδειγμα

Να βρεθεί το δειγματικό μοντέλο γραμμικής παλινδρόμησης για το σύνολο δεδομένων

$$\{(1, -1, 1), (0, -1, -1), (2, 0, 2), (1, 1, 2)\}$$

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:

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