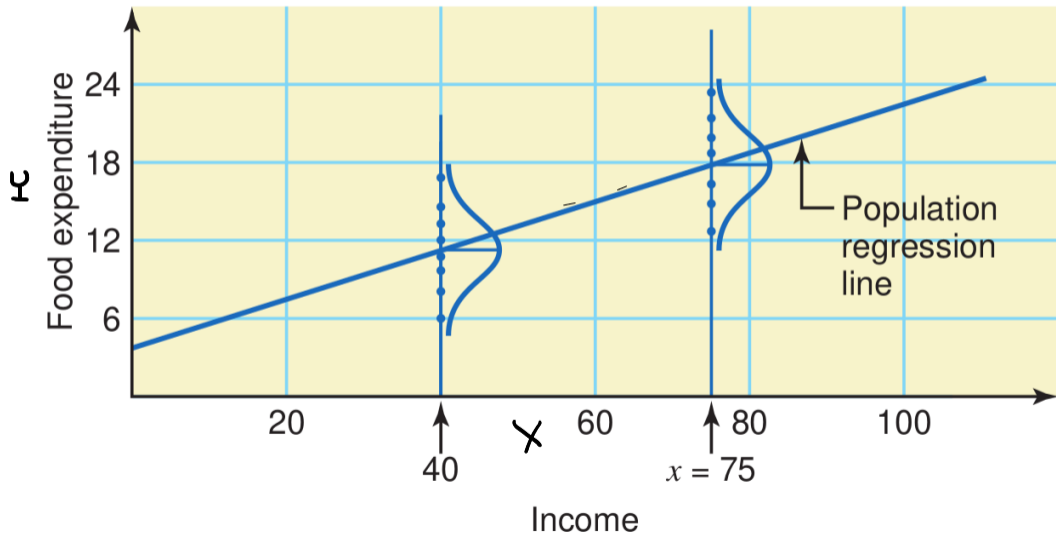


**MEM-205 Περιγραφική Στατιστική**  
Τμήμα Μαθηματικών και Εφ. Μαθηματικών, Πανεπιστήμιο Κρήτης

Κώστας Σμαραγδάκης (kesmarag@pm.me)

15-03-2023

# Απλή Γραμμική Παλινδρόμηση



## Δειγματικό μοντέλο απλής γραμμικής παλινδρόμησης

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

$$y = A + Bx$$

- ▶  $a$  είναι δειγματική προσέγγιση του  $A$
- ▶  $b$  είναι δειγματική προσέγγιση του  $B$
- ▶  $\hat{y}$  είναι η εκτιμώμενη τιμή του  $y$  για δοσμένο  $x$

## Τυχάιο σφάλμα του δειγματικού μοντέλου απλής γραμμικής παλινδρόμησης

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

# Απλή Γραμμική Παλινδρόμηση

Έστω το τυχαίο δείγμα

ανεξαρτητή

$\alpha, \beta$

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

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

εξαρτημένη

$(x_j, y_j) \rightarrow \hat{y}_j$

Για το τυχαίο σφάλμα του δειγματικού μοντέλου/απλής γραμμικής παλινδρόμησης έχουμε:

Προετοιμασία της τιμ  $e$

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

όπου η προσέγγιση του κάθε  $y_n$  δίνεται ως

$$\hat{y}_n = a + bx_n, \quad n = 1, \dots, N$$

## Άθροισμα τετραγωνικών σφαλμάτων

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

Άθροισμα τετραγωνικών σφαλμάτων συναρτήσει των παραμέτρων του δειγματικού μοντέλου

$$SSE = \sum e_n^2 = \sum (y_n - \hat{y}_n)^2 = \sum (y_n - a - bx_n)^2$$

$$Q(a, b) = SSE = \sum_{n=1}^N (y_n - a - bx_n)^2$$

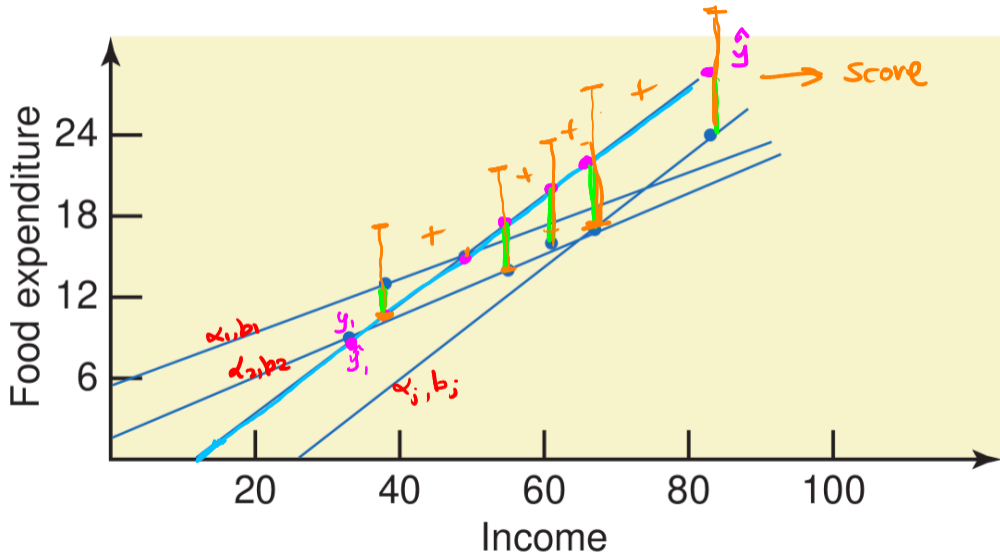
Εκτίμηση ελαχίστων τετραγώνων

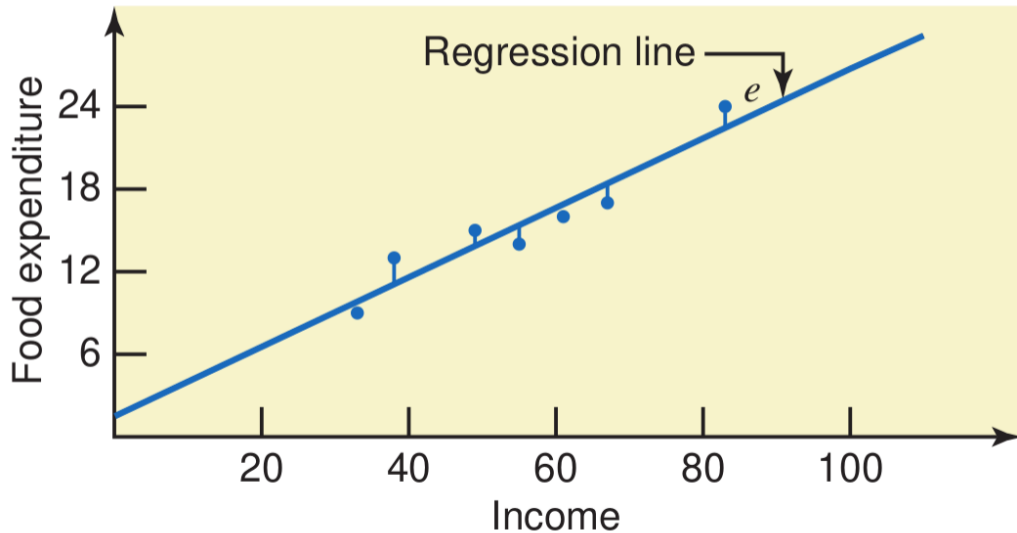
Ως εκτίμησεις των  $a, b$  λαμβάνουμε τις τιμές  $a^*, b^*$  που ελαχιστοποιούν το άθροισμα των τετραγωνικών σφαλμάτων.

$$a, b = \arg \min_{a', b'} Q(a', b')$$

$$a, b \in \mathbb{R} \text{ τ.ω. } Q(a', b') \geq Q(a, b) \quad \forall a', b' \in \mathbb{R}$$

# Απλή Γραμμική Παλινδρόμηση - Εκτίμηση Ελαχίστων Τετραγώνων





$$Q(a, b) = \sum_{n=1}^N (y_n - a - bx_n)^2$$

$$\frac{\partial Q}{\partial a} = -\sum_{n=1}^N 2(y_n - a - bx_n) \quad , \quad \text{διαφορ} \quad \frac{\partial Q}{\partial a} = 0 \Rightarrow$$

$$\Rightarrow \sum_{n=1}^N (y_n - a - bx_n) = 0 \Rightarrow \sum_{n=1}^N y_n - Na - b \sum_{n=1}^N x_n = 0 \Rightarrow \frac{1}{N}$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N y_n - b \frac{1}{N} \sum_{n=1}^N x_n = \alpha \Rightarrow \boxed{\alpha = \bar{y} - b\bar{x}}$$

$$\frac{\partial Q}{\partial b} = -\sum_{n=1}^N 2(y_n - a - bx_n)x_n = 0$$

$$\sum_{n=1}^N (x_n y_n - \alpha x_n - b x_n^2) = 0 \Rightarrow \sum_{n=1}^N (x_n y_n - \bar{y} x_n + b \bar{x} x_n - b x_n^2) = 0$$



$$\Rightarrow \sum_{n=1}^N x_n y_n - \bar{y} \sum_{n=1}^N x_n + b \bar{x} \sum_{n=1}^N x_n - b \sum_{n=1}^N x_n^2 = 0$$

$$\sum_{n=1}^N x_n y_n - \frac{1}{N} \sum_{n=1}^N x_n \sum_{n=1}^N y_n + b \left[ \frac{1}{N} \left( \sum_{n=1}^N x_n \right)^2 - \sum_{n=1}^N x_n^2 \right] = 0$$

$$\Rightarrow b = \frac{\sum_{n=1}^N x_n y_n - \frac{1}{N} \sum_{n=1}^N x_n \sum_{n=1}^N y_n}{\sum_{n=1}^N x_n^2 - \frac{1}{N} \left( \sum_{n=1}^N x_n \right)^2}$$

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

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

όπου  $SS_{xy}$ ,  $SS_{xx}$  δίνονται ως:

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

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

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

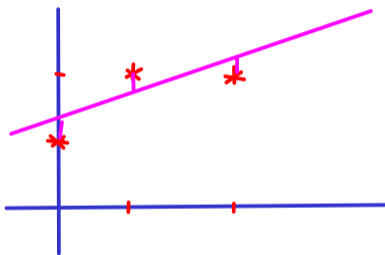
$$\sum_{n=1}^N (x_n - \bar{X})(y_n - \bar{Y}) \rightarrow \sum_{n=1}^N x_n y_n - \frac{(\sum_{n=1}^N x_n)(\sum_{n=1}^N y_n)}{N}$$

$$\sum x_n y_n - \bar{Y} \sum x_n - \bar{X} \sum y_n + N \bar{X} \bar{Y} =$$

$$\sum x_n y_n - \frac{1}{N} \sum y_n \sum x_n - \frac{1}{N} \sum x_n \sum y_n + N \frac{1}{N} \sum x_n \frac{1}{N} \sum y_n$$

## Παράδειγμα

Βρείτε τη εκτίμηση ελαχίστων τετραγώνων του μοντέλου γραμμικής παλινδρόμησης υποθέτοντας τα παρακάτω δεδομένα.



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

	x	y	xy	x <sup>2</sup>
	0	1	0	0
	1	2	2	1
	2	2	4	4
Sum	3	5	6	5

$$S_{xx} = \sum x_n^2 - \frac{1}{N} \left( \sum x_n \right)^2 = 5 - \frac{1}{3} 3^2 = 2$$

$$S_{xy} = \sum x_n y_n - \frac{1}{N} \left( \sum x_n \right) \left( \sum y_n \right) = 6 - \frac{1}{3} 3 \cdot 5 = 1$$

$$\left. \begin{array}{l} S_{xx} = 2 \\ S_{xy} = 1 \end{array} \right\} b = \frac{1}{2}$$

$$a = \bar{Y} - b\bar{X} = \frac{5}{3} - \frac{1}{2} \frac{3}{3} = \frac{7}{6}$$

$$\hat{y} = \frac{7}{6} + \frac{1}{2}x$$

$$\hat{y}(0) = \frac{7}{6} \quad y(0) = 1$$

$$\hat{y}_1 = \frac{7}{6} \quad y_1 = 1$$

$$e_1 = 1 - \frac{7}{6} = -\frac{1}{6}$$

$$\hat{y}_2 = \hat{y}(1) = \frac{7}{6} + \frac{1}{2} = \frac{10}{6}$$

$$y_2 = 2 \quad e_2 = 2 - \frac{10}{6} = \frac{1}{3}$$

$$\hat{y}_3 = \hat{y}(2) = \frac{7}{6} + 1 = \frac{13}{6}$$

$$y_3 = 2 \quad e_3 = 2 - \frac{13}{6} = -\frac{1}{6}$$

### Άσκηση

Βρείτε τη εκτίμηση ελαχίστων τετραγώνων του μοντέλου γραμμικής παλινδρόμησης υποθέτοντας τα παρακάτω δεδομένα.

$$\{(0, 2), (1, 1), (1, 2), (2, 4)\}$$

$X, Y$  τυχαίες μεταβλητές

$$Y = A + BX + \epsilon, \quad \epsilon: \text{όρος τυχαίου σφάλματος}$$

$$\hat{Y} = \alpha + bX + e \quad S_e^2 - \text{διασπορά του } e$$

- ▶ Για κάθε  $x$  έχουμε υποθέσει ότι το σφάλμα  $\epsilon$  ακολουθεί την κανονική κατανομή  $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$ .
- ▶ Η τυπική απόκλιση  $\sigma_\epsilon$  του τυχαίου σφάλματος αναφέρεται στο πληθυσμό και κατά επέκταση η τιμή της δεν είναι γνωστή στις περισσότερες περιπτώσεις.

## Εκτιμητήρια της τυπικής απόκλισης των σφαλμάτων

$$\begin{aligned} t_x &\rightarrow X \\ t_y &\rightarrow Y \end{aligned}$$

$e_n$  εξαρτάται από  $\bar{X}, \bar{Y}$

$$S_e = \sqrt{\frac{SSE}{N-2}}, \quad SSE = \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

χάσω

2 βαθμούς ελευθερ..

$$\alpha = \bar{Y} - b\bar{X}$$

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

Γιατί εμφανίζεται το  $N - 2$ ;

$$SSE = \sum (e_n - 0)^2$$



$$s_e = \sqrt{\frac{SSE}{N-2}}$$

$$s_e = \sqrt{\frac{SS_{yy} - b * SS_{xy}}{N-2}} = \sqrt{\frac{SSE}{N-2}}$$

όπου:

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

Υπενθυμίζεται

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

Εάν είχαμε γνώση των δεδομένων του πληθυσμού θα μπορούσαμε να υπολογίσουμε την τυπική απόκλιση των τυχαίων σφαλμάτων από τη σχέση:

$$\sigma_{\epsilon} = \sqrt{\frac{SS_{yy} - B * SS_{xy}}{N_p}}$$

όπου σε αυτή την περίπτωση θα είχαμε:

$$SS_{yy} = \sum_{n=1}^{N_p} (y_n - \mu_y)^2, \quad SS_{xy} = \sum_{n=1}^{N_p} (x_n - \mu_x)(y_n - \mu_y)$$

## Συντελεστής Προσδιορισμού (Coefficient of Determination)

Συνολικό Άθροισμα τετραγώνων

$$\sum_{j=1}^N (x_j, y_j) \rightarrow \alpha, b \rightarrow \hat{y}_j$$

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

Άθροισμα τετραγώνων παλινδρόμησης

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

Συντελεστής Προσδιορισμού

$$R^2 = \frac{SSR}{SST}, \quad 0 \leq R^2 \leq 1 \text{ (γιατί;)}$$

- Ποσοτικοποιεί την αποτελεσματικότητα του μοντέλου.

## Συντελεστής Προσδιορισμού (Coefficient of Determination)

$$SST = SSR + SSE$$

$$R^2 = \frac{SSR}{SSR + SSE}$$

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

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

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

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

### Παράδειγμα

Βρείτε τον συντελεστή προσδιορισμού του συνόλου δεδομένων:

$$\{(0, 1), (1, 3), (2, 4), (5, 4)\}$$

Μέση τιμή, τυπική απόκλιση και κατανομή του  $b$

$$\mu_b = B, \quad \sigma_b = \frac{\sigma_\epsilon}{\sqrt{SS_{xx}}}$$

$$b \sim \mathcal{N}(\mu_b, \sigma_b)$$

- ▶ Όταν το  $\sigma$  είναι άγνωστο δεν μπορούμε να υπολογίσουμε το  $\sigma_b$

Εκτιμητήρια της τυπικής απόκλιση του  $b$

$$s_b = \frac{s_e}{\sqrt{SS_{xx}}}$$