

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

Κώστας Σμαραγδάκης (kesmarag@gmail.com)

05-03-2020

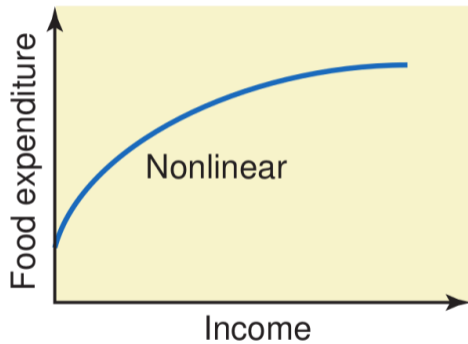
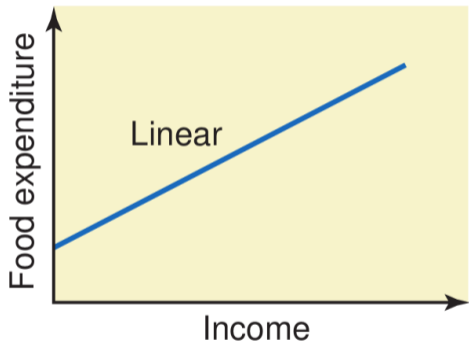
Παλινδρόμηση

Ένα μοντέλο παλινδρόμησης είναι μια μαθηματική εξίσωση που περιγράφει την σχέση μεταξύ δύο ή περισσότερων μεταβλητών. Το μοντέλο παλινδρόμησης με δύο μεταβλητές, μια ανεξάρτητη και μια εξαρτημένη ονομάζεται **μοντέλο απλής παλινδρόμησης**.

Απλή Γραμμική Παλινδρόμηση (Simple Linear Regression)

Ένα μοντέλο παλινδρόμησης το οποίο συνδέει με γραμμικό τρόπο την ανεξάρτητη με την εξαρτημένη μεταβλητή ονομάζεται **μοντέλο απλής γραμμικής παλινδρόμησης**.

Παλινδρόμηση (Regression)



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)

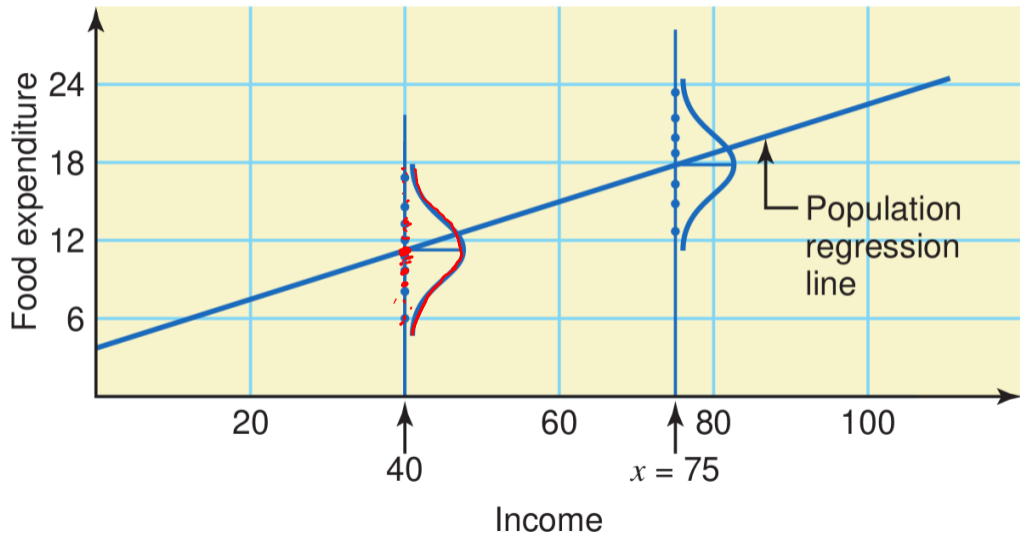
Παραδοχές

- ▶ Για δοσμένο x το ϵ ακολουθεί τυπική κανονική κατανομή. $\mathcal{N}(0, \sigma_\epsilon)$
- ▶ 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



Δειγματικό μοντέλο απλής γραμμικής παλινδρόμησης

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

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

Τυχαιο σφάλμα του δειγματικού μοντέλου απλής γραμμικής παλινδρόμησης

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

Έστω το τυχαίο δείγμα

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

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

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

όπου η προσέγγιση του κάθε \hat{y}_n δίνεται ως

$$\hat{y}_n = a + bx_n$$

Άθροισμα τετραγωνικών σφαλμάτων

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

Άθροισμα τετραγωνικών σφαλμάτων συναρτήσει των παραμέτρων του δειγματικού μοντέλου

$$Q(a, b) = \text{SSE} = \sum_{n=1}^N (y_n - \underbrace{a - bx_n}_{-\hat{y}_n})^2$$

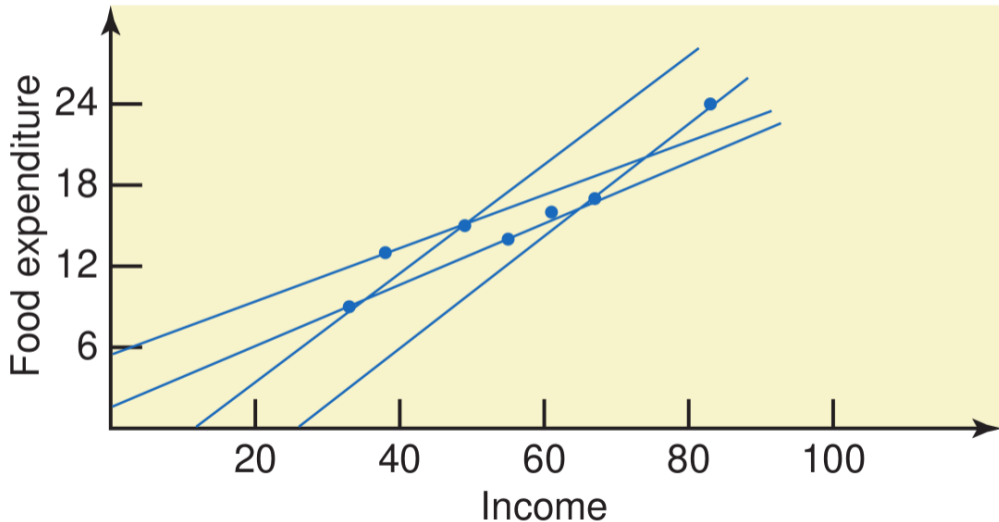
Εκτίμηση ελαχίστων τετραγώνων

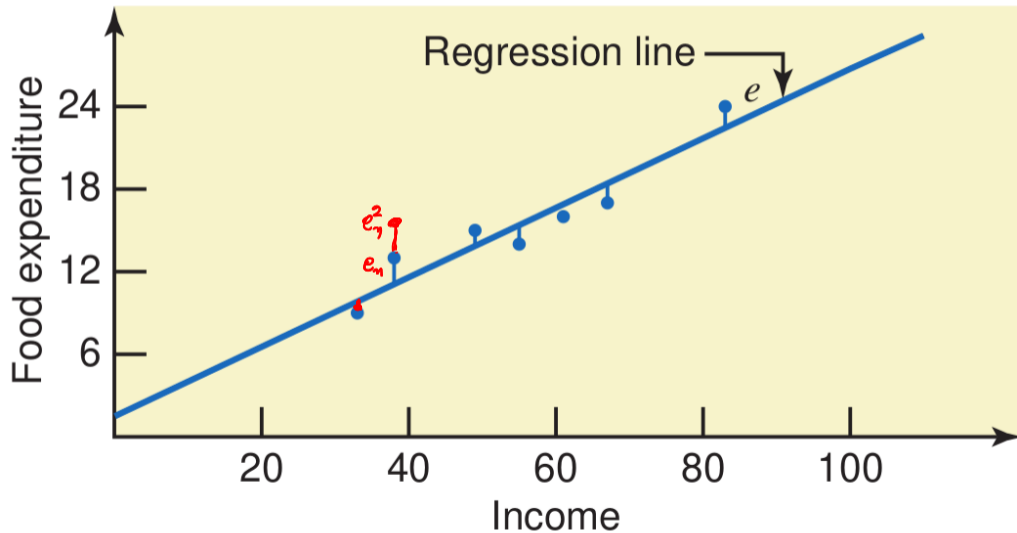
Ως εκτίμησεις των a, b λαμβάνουμε τις τιμές a^*, b^* που ελαχιστοποιούν το άθροισμα των τετραγωνικών σφαλμάτων.

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



Απλή Γραμμική Παλινδρόμηση - Εκτίμηση Ελαχίστων Τετραγώνων





Απλή Γραμμική Παλινδρόμηση - Εκτίμηση Ελαχίστων Τετραγώνων

$$\begin{aligned}
 & \sum_{n=1}^N e_n = 0 \\
 Q(a, b) &= \sum_{n=1}^N (y_n - a - bx_n)^2 \\
 \frac{\partial Q}{\partial a} &= -2 \sum_{n=1}^N (y_n - a - bx_n) = 0 \Leftrightarrow \sum_{n=1}^N y_n - a \sum_{n=1}^N 1 - b \sum_{n=1}^N x_n = 0 \Leftrightarrow \sum_{n=1}^N y_n - aN - b \sum_{n=1}^N x_n = 0 \\
 a &= -b \sum_{n=1}^N x_n + \sum_{n=1}^N y_n \Leftrightarrow \boxed{a = \bar{Y} - b\bar{X}}
 \end{aligned}$$

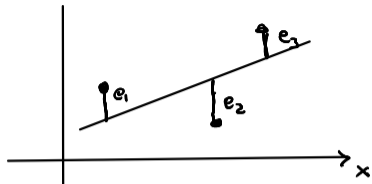
$$\frac{\partial Q}{\partial b} = -2 \sum_{n=1}^N x_n (y_n - a - bx_n) = 0 \Leftrightarrow \sum_{n=1}^N x_n y_n - a \sum_{n=1}^N x_n - b \sum_{n=1}^N x_n^2 = 0 \Leftrightarrow$$

$$\Leftrightarrow \sum_{n=1}^N x_n y_n - \bar{Y} \sum_{n=1}^N x_n + b\bar{X} \sum_{n=1}^N x_n - b \sum_{n=1}^N x_n^2 = 0$$

$$b \left(\bar{X} \sum_{n=1}^N x_n - \sum_{n=1}^N x_n^2 \right) = \bar{Y} \sum_{n=1}^N x_n - \sum_{n=1}^N x_n y_n$$

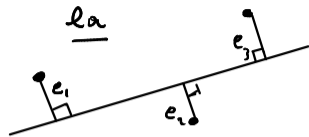
$$\begin{aligned}
 & \frac{1}{N} \sum x_n \cdot \sum x_n \\
 & \frac{1}{N} \left(\sum x_n \right)^2 \\
 & \frac{1}{N} \sum x_n \sum y_n
 \end{aligned}$$

Απλή Γραμμική Παλινδρόμηση - Εκτίμηση Ελαχίστων Τετραγώνων



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

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

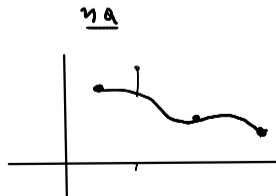


όπου SS_{xy} , SS_{xx} δίνονται ως:

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

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

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



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

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

$$= \sum_{n=1}^N x_n y_n - \frac{1}{N} \sum_{n=1}^N x_n \sum_{n=1}^N y_n - \frac{1}{N} \sum_{n=1}^N x_n \sum_{n=1}^N y_n + N \frac{1}{N^2} \sum_{n=1}^N x_n \sum_{n=1}^N y_n$$

$$= \sum_{n=1}^N x_n y_n - \frac{\sum_{n=1}^N x_n \sum_{n=1}^N y_n}{N}$$

Παράδειγμα

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

	x	y	xy	x ²
	0	1	0	0
	1	2	2	1
	2	2	4	4
Σ	3	5	6	5

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

$$\sum_{n=1}^3 x_n = 3 \quad \sum_{n=1}^3 y_n = 5 \quad \sum_{n=1}^3 x_n y_n = 6 \quad \sum_{n=1}^3 x_n^2 = 5$$

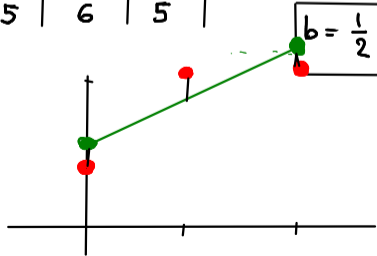
$$\bar{X} = \frac{3}{3} = 1 \quad \bar{Y} = \frac{5}{3}$$

$$\bar{Y} = \frac{5}{3}$$

$$SS_{xy} = 6 - \frac{3 \cdot 5}{3} = 6 - 5 = 1 \quad S_{xx} = 5 - \frac{9}{3} = 5 - 3 = 2$$

$$a = \frac{5}{3} - \frac{1}{2} = \frac{7}{6}$$

$$b = \frac{1}{2}$$



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

$$x=0 \quad \hat{y} = \frac{7}{6}$$

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

Άσκηση

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

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

Διανυσματική μορφή

Έστω διανύσματα $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ στήλες με στοιχεία της παρατηρήσεις της ανεξάρτητης και της εξαρτημένης μεταβλητής αντίστοιχα. Το μοντέλο απλής γραμμικής παλινδρόμησης δίνει εκτιμήσεις για τις τιμές της εξαρτημένης μεταβλητής που αντιστοιχούν στις παρατηρήσεις της ανεξάρτητης μεταβλητής που περιέχονται στο \mathbf{x} :

$$\hat{\mathbf{y}} = a\mathbf{u} + b\mathbf{x} \quad \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = a \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + b \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

όπου $\mathbf{u} \in \mathbb{R}^N$ διάνυσμα στήλη με στοιχεία άσους.
Κατά επέκταση έχουμε:

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

καθώς και

$$\text{SSE} = \mathbf{e}^T \mathbf{e} = [e_1, e_2, \dots, e_n] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = \sum_{n=1}^n e_n^2$$

$$Q(a, b) = (\mathbf{y} - a\mathbf{u} - b\mathbf{x})^T (\mathbf{y} - a\mathbf{u} - b\mathbf{x})$$

$$\bar{X} = \frac{1}{N} \mathbf{u}^T \mathbf{x}, \quad \bar{Y} = \frac{1}{N} \mathbf{u}^T \mathbf{y}$$

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

$$SS_{xy} = (\mathbf{x} - \bar{X}\mathbf{u})^T (\mathbf{y} - \bar{Y}\mathbf{u}), \quad SS_{xx} = (\mathbf{x} - \bar{X}\mathbf{u})^T (\mathbf{x} - \bar{X}\mathbf{u})$$

Παράδειγμα

Βρείτε τη εκτίμηση ελαχίστων τετραγώνων του μοντέλου γραμμικής παλινδρόμησης υποθέτοντας τα παρακάτω δεδομένα κάνοντας χρήση των διανυσματικών εκφράσεων.

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