

Φροντιστήριο #3

Συντελεστής Προσδιορισμού

Έστω  $\hat{y}_i$  οι τιμές που προκύπτουν για τη μεταβλητή απόκρισης  $Y$  χρησιμοποιώντας την ευθεία παραγωγής

$$y = \beta_0 + \beta_1 x, \text{ δηλαδή } \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, i=1, \dots, n$$

Αυτές οι ποσότητες είναι οι ευτιμήσεις μας για την τιμή της μεταβλητής απόκρισης στις θέσεις  $x_i, i=1, \dots, n$  και καλούνται ευτιμημένες τιμές (ή ευτιμώμενες τιμές) της  $Y$  στις θέσεις  $x_i, i=1, \dots, n$ .

$$\text{Οι διαφορές } \hat{\epsilon}_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i), i=1, \dots, n$$

καλούνται ευτιμημένα ή ευτιμώμενα σφάλματα και εκφράζουν τις αποκλίσεις των παρατηρηθείσων τιμών  $y_i$  της  $Y$  (για την τιμή  $x_i$  της  $X$ ) από τις ευτιμημένες τιμές  $\hat{y}_i$ .

$$\text{Η ποσότητα } g(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \hat{\epsilon}_i^2 = SSE$$

είναι η ελάχιστη τιμή του αθροίσματος  $g(\beta_0, \beta_1)$  καλείται άθροισμα τετραγώνων των (ευτιμημένων) σφαλμάτων.  
SSE: sum of squares of Errors.

Είναι φανερό ότι:

- αν SSE μικρό η ευθεία θα περνάει "κοντά" από τα σημεία  $(x_i, y_i), i=1, \dots, n$
- αν SSE μεγάλο, η ευθεία δεν θα βρίσκεται "κοντά"

σε όλα τα σημεία  $(x_i, y_i), i=1, \dots, n$ .

-2-

Ένα μέτρο μεταβλητότητας των τιμών της  $Y$  γύρω από τη μέση τιμή τους  $\bar{y}$  είναι το συνολικό άθροισμα των τετραγώνων  $SST$ , όπου  $SST = \sum_{i=1}^n (y_i - \bar{y})^2$

Μπορούμε να γράψουμε:

$$SST = SSR + SSE$$

όπου  $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$  είναι το άθροισμα τετραγώνων της παλινδρόμησης (Regression Sum of Squares).

Έχουμε:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i = (\bar{y} - \hat{\beta}_1 \bar{x}) + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1 (x_i - \bar{x})$$

Άρα

$$SSR = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y})^2 = \sum_{i=1}^n [\bar{y} + \hat{\beta}_1 (x_i - \bar{x}) - \bar{y}]^2$$

$$= \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

Συνεπώς, για το  $SSR$ , έχουμε:

$$SSR = \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

Επίσης:

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

$$= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)$$

$$= \sum_{i=1}^n (y_i^2 - \hat{\beta}_0 y_i - \hat{\beta}_1 x_i y_i - \hat{\beta}_0 y_i + \hat{\beta}_0^2 + \hat{\beta}_0 \hat{\beta}_1 x_i - \hat{\beta}_1 x_i y_i + \hat{\beta}_1 \hat{\beta}_0 x_i + \hat{\beta}_1^2 x_i^2)$$

$$\begin{aligned}
&= \sum_{i=1}^n y_i^2 - \hat{\beta}_0 \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i \\
&\quad - \hat{\beta}_0 \sum_{i=1}^n y_i + n \hat{\beta}_0^2 + \hat{\beta}_0 \hat{\beta}_1 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i + \hat{\beta}_1 \hat{\beta}_0 \sum_{i=1}^n x_i \\
&\quad + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2
\end{aligned}$$

Όπως

$$\begin{aligned}
& - \hat{\beta}_0 \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i + 2 \hat{\beta}_1 \hat{\beta}_0 \sum_{i=1}^n x_i + n \hat{\beta}_0^2 + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2 \\
&= -\hat{\beta}_0 (n \hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i) - \hat{\beta}_1 \sum_{i=1}^n x_i y_i + 2 \hat{\beta}_1 \hat{\beta}_0 \sum_{i=1}^n x_i + n \hat{\beta}_0^2 \\
&\quad + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2 \\
&= -n \hat{\beta}_0^2 - \hat{\beta}_0 \hat{\beta}_1 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i + 2 \hat{\beta}_1 \hat{\beta}_0 \sum_{i=1}^n x_i \\
&\quad + n \hat{\beta}_0^2 + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2 \\
&= \hat{\beta}_0 \hat{\beta}_1 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2 \\
&= \hat{\beta}_0 \hat{\beta}_1 \sum_{i=1}^n x_i - \hat{\beta}_1 (\hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2) + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2 \\
&= \hat{\beta}_0 \hat{\beta}_1 \sum_{i=1}^n x_i - \hat{\beta}_1 \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1^2 \sum_{i=1}^n x_i^2 + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2 = 0.
\end{aligned}$$

Άρα,

$$SSE = \sum_{i=1}^n y_i^2 - \hat{\beta}_0 \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i y_i$$

Η ποσότητα:  $s^2 = MSE = \frac{SSE}{n-2} = \frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2$



$$= \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

είναι αμερόληπτη εκτιμήτρια της παραμέτρου  $\sigma^2$  δηλ.  $E(s^2) = \sigma^2 = \text{var}(Y_i) = \text{var}(\epsilon_i)$  λήγεται μέσω άθροισης τετραγώνων των υπολοίπων ή μέσω τετραγωνικό υπόλοιπο ή μέσω τετραγωνικό σφάλμα (error mean square, residual mean square).

Χρησιμοποιούμε τον λόγο  $\frac{SSR}{SST}$  ως ένα δείκτη ποιότητας του μοντέλου γραμμικής παλινδρόμησης

$$y = \hat{\beta}_0 + \hat{\beta}_1 x.$$

Η ποσότητα  $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$  καλείται

συντελεστής προσδιορισμού (coefficient of determination) του γραμμικού μοντέλου. Συχνά εκφράζεται ως ποσοστό, παίρνει τιμές μεταξύ 0 και 1.

Τιμές  $R^2 \rightarrow 1$ : η ευθεία παλινδρόμησης περνάει πολύ κοντά από τα περισσότερα σημεία, ενώ

$R^2 \rightarrow 0$  πολλά σημεία βρίσκονται μακριά από την ευθεία και θα πρέπει να αναζητηθεί μία άλλη σχέση της ανεξαρτητής και της εξαρτημένης μεταβλητής (ενδεχομένως, μη-γραμμική).

Ο συντελεστής  $R^2$  εκφράζει το ποσοστό της σωματικής -5- διασποράς των τιμών της εξαρτημένης μεταβλητής  $Y$  η οποία ερμηνεύεται (εξηγείται) από την ανεξάρτητη μεταβλητή  $X$  (μέσω της ευθείας παλινδρόμησης).

Μόνο για το απλό γραμμικό μοντέλο ισχύει ότι:

$R^2 = r^2$ , όπου  $r$  είναι ο συντελεστής γραμμικής συσχέτισης των  $X$  και  $Y$ .

Έλεγχος στατιστικής σημαντικότητας του συντελεστή

$R^2$ . Ελέγχουμε  $H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$  (F-test)

Στην απλή παλινδρόμηση το F-test είναι ακριβώς ισοδύναμο με το t-test για τη στατιστική σημαντικότητα του συντελεστή  $\beta_1$ .

Εναλλακτικά, δράφουμε:

$H_0$ : η εξίσωση παλινδρόμησης δεν εξηγεί καθόλου τις μεταβολές της  $Y$  ή το ποσοστό της διασποράς που εξηγείται για την  $Y$  είναι μηδέν

vs

$H_1$ : η εξίσωση παλινδρόμησης εξηγεί ένα μέρος των μεταβολών της  $Y$  ή το ποσοστό της ερμηνεύμενης διασποράς της  $Y$  είναι μεγαλύτερο του μηδέν.

Παρατήρηση: Ο συντελεστής γραμμικής συσχέτισης  $r$  έχει το ίδιο πρόσημο με την ευκρινή τιμή της κλίσης του απλού μοντέλου παλινδρόμησης  $b_1 = \hat{\beta}_1$ .

Ο έλεγχος μαθαίνει έλεγχος ανάλησης διακύμα- -6-  
 νου - analysis of variance test (ANOVA test)  
 και για τη διεύρεση του μπορούμε αρχικά να  
 κατασκευάσουμε τον ακόλουθο πίνακα που μαθαίνει  
 πίνακας ανάλησης διακύματος (ANOVA table).

<u>Πηγή</u>	<u>Άθροισμα</u>	<u>β.ε.</u>	<u>Μέσα</u>	<u>Λόγος</u>	<u><math>F_{1, n-2}</math></u>
<u>Μεταβλητότητα</u>	<u>Τετραγώνων</u>		<u>Τετραγώνων</u>		
Παλινδρόμηση	SSR	1	$SSR/1 = MSR$	$\frac{MSR}{MSE} = F_{1, n-2}$	
Κατάλοιπα	SSE	$n-2$	$SSE/n-2 = MSE$		
Σύνολο	SST	$n-1$			

Υπολογίζουμε τον λόγο  $F_{1, n-2}$ .

Απορρίπτουμε την  $H_0$  σε εσο  $\alpha$  αν

$$F_{1, n-2} > \tilde{F}_{1, n-2, \alpha} \quad \text{όπου } \tilde{F}_{1, n-2, \alpha} \text{ είναι το}$$

$\alpha$ -ποσοστιαίο σημείο της κατανομής  $\tilde{F}_{1, n-2}$  με

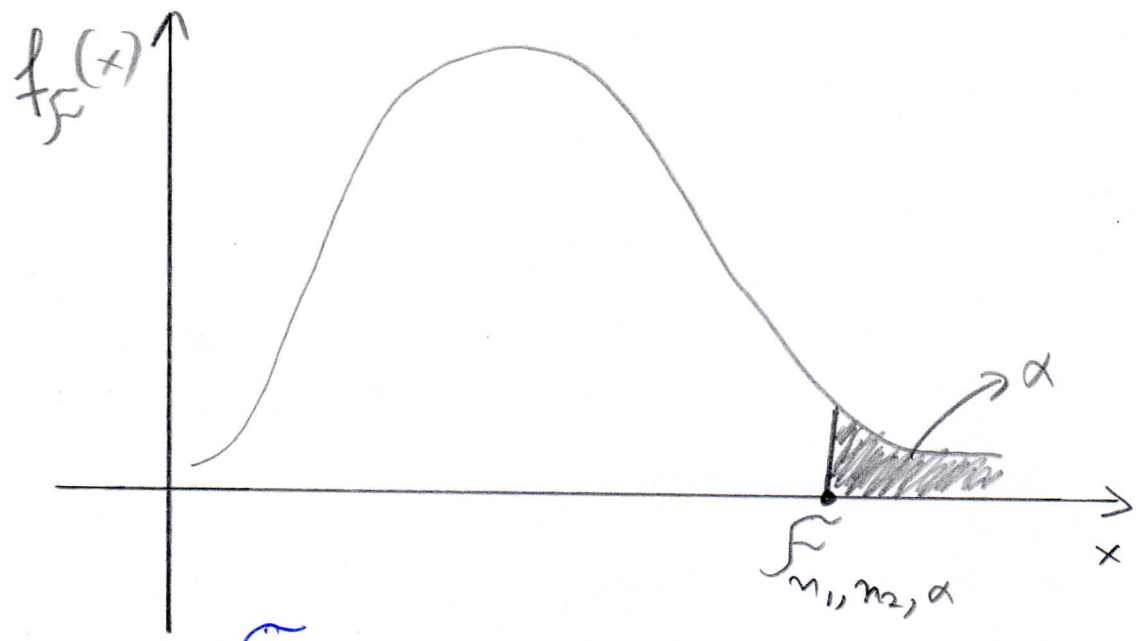
$n_1 = 1$  και  $n_2 = n - 2$  β.ε.

Λίγα στοιχεία για την κατανομή  $\tilde{F}_{n_1, n_2}$

Η γραφική της παράσταση έχει την ακόλουθη  
 μορφή (του σχήματος) για διάφορες τιμές των  
 $n_1, n_2$ . Στο σχήμα καταγράφουμε και το σημείο



$F_{n_1, n_2, \alpha}$



Το σημείο  $F_{n_1, n_2, \alpha}$  είναι τ.ω.  $P(F > F_{n_1, n_2, \alpha}) = \alpha$ .

Παραδείγματα

① Μία τράπεζα ενδιαφέρεται να μελετήσει την αποτελεσματική συμπεριφορά των πελατών της. Θεωρούμε ότι η ετήσια αποταμίευση των πελατών της εξαρτάται από το ετήσιο εισόδημά τους. Σε δείγμα 10 πελατών έχουμε τα ακόλουθα στοιχεία (τιμές σε χιλιάδες €).

Ετήσια αποταμίευση (Y)	5	3	5	4	6	1	2	8	2	3
Ετήσια εισοδήματα (X)	50	31	28	45	50	32	36	55	26	47

(α) Να επισημάνετε και να ερμηνεύσετε τους συντελεστές του γραμμικού υποδείγματος.

(β) Ποια η τιμή του συντελεστή προσδιορισμού και ποια η ερμηνεία του.

(δ) Να παραταποποιηθεί έλεγχος για εσσ  $\alpha = 5\%$  -8-  
 για τη στατιστική σημαντικότητα του συντελεστή  
 προσδιορισμού (ή ολόκληρης της εξίσωσης παλινδρό-  
 μησης).

$$\text{Δίνονται: } \sum_{i=1}^{10} x_i = 400, \quad \sum_{i=1}^{10} y_i = 39, \quad \sum_{i=1}^{10} x_i^2 = 17000$$

$$\sum_{i=1}^{10} y_i^2 = 193, \quad \sum_{i=1}^{10} x_i y_i = 1700.$$

Λύση

(α) Οι συντελεστές του γραμμικού υποδείγματος

$$Y = \beta_0 + \beta_1 X + \varepsilon \text{ δίνονται:}$$

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^{10} x_i y_i - \sum_{i=1}^{10} x_i \sum_{i=1}^{10} y_i}{n \sum_{i=1}^{10} x_i^2 - \left( \sum_{i=1}^{10} x_i \right)^2} = \frac{10 \cdot 1700 - 400 \cdot 39}{10 \cdot 17000 - (400)^2} = 0,14$$

$$\text{και } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \frac{39}{10} - 0,14 \cdot \frac{400}{10} = -1,7.$$

Έ συνεπώς, η ευτερώρηνη ευθεία παλινδρόμησης της

Y πάνω στη X θα είναι:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X = -1,7 + 0,14 X.$$

Ερμηνεία των συντελεστών:

$\hat{\beta}_1 = 0,14$  Για κάθε αύξηση κατά 1000 € των  
 ετήσιων εισοδημάτων του πελάτη της τράπεζας



-9-  
Οι ετήσιες αποταμιεύσεις που αναμένεται να έχουν μία αύξηση περίπου 14% με 140 ευρώ (0,14 · 1000 €).

$\hat{\beta}_0 = -1,7$  Αν υποθέσουμε ότι ένας πελάτης δεν έχει ετήσια εισοδήματα, οι ετήσιες αποταμιεύσεις που αναμένεται να μειωθούν κατά περίπου 1700 ευρώ (-1,7 · 1000 €).

(β) Έχουμε  $R^2 = \frac{SSR}{SST}$

Όπως,  $SSR = \hat{\beta}_1^2 \left( \sum_{i=1}^{n=10} x_i^2 - n \bar{x}^2 \right)$

$$= 0,14^2 (17000 - 10 \cdot 40^2) = 19,6$$

$$SST = \sum_{i=1}^{n=10} y_i^2 - n \bar{y}^2 = 193 - 10 \left( \frac{39}{10} \right)^2$$
$$= 193 - 152,1 = 40,9$$

$$SSE = SST - SSR = 40,9 - 19,6 = 21,3$$

Άρα,  $R^2 = \frac{SSR}{SST} = \frac{19,6}{40,9} = 0,479 = 47,9\%$

Ερμηνεία  $R^2$ : Το 48% περίπου της σωματικής μετα-

βλητότητας των ετήσιων αποταμιεύσεων των πελατών της τράπεζας οφείλονται στα ετήσια εισοδήματα τους και το υπόλοιπο 52% οφείλεται σε άλλους παράγοντες (ή/και σε μεταβλητές που

δεν έχω προβλεφθεί από το μοντέλο).

-10-

(γ) Μπορούμε να φτιάξουμε τον πίνακα ανάγνωσης διαυδράτους.

<u>Πηγή</u> <u>Μεταβλητότητας</u>	<u>Αθροίσματα</u> <u>Τετραγώνων</u>	<u>Β.ε</u>	<u>Μέσα</u> <u>Τετράγωνα</u>	<u>F</u> <u>1,8</u>
Παχιδρόμηση	19,6	1	19,6	$F_{1,8} = \frac{19,6}{2,66}$
Κατάλοιπα	21,3	8	2,66	
Σύνολο	40,9	9		= 7,36

Ελέγχουμε  $H_0: \beta_1 = 0$

$H_1: \beta_1 \neq 0$

Απορρίπτουμε την  $H_0$  σε ε.σ.σ.  $\alpha = 5\%$  αν

$$F_{1,8} = 7,36 > F_{1,8,0.05} = 5,32$$

ισχύει η ανισότητα άρα τα δεδομένα παρέχουν ισχυρές ενδείξεις απόρριψης της  $H_0$  συνεπώς φαίνεται ότι τα ετήσια εισοδήματα

επιφέρουν στατιστικά σημαντική επίδραση (τη μεταβλητότητα) των ετήσιων αποταμιεύσεων των πελατών της τράπεζας.

(δ) (Εχτρα ερώτημα)

Να υπολογιστεί το μέσο τετραγωνικό σφάλμα της εξίσωσης παλινδρόμησης.

Έχουμε:  $s^2 = MSE = \frac{SSE}{n-2} = \frac{21,3}{8} = 2,6625$

Το ωπικό σφάλμα των καταλοίπων είναι:  $s = \sqrt{s^2} = \sqrt{2,6625} = 1,631$ . (Τα σφάλματα εμπίπτουν περίπου  $\pm 1,64\mu$  γύρω από την ευθεία παλινδρόμησης). (κατά μέσο όρο)

Παρατήρηση (για τον έλεγχο F)

Το πηλίκο  $F_{1, n-2} = \frac{MSR}{MSE}$  μπορεί να γραφεί

εναλλακτικά:

$$F_{1, n-2} = \frac{\frac{R^2}{1}}{\frac{(1-R^2)}{(n-2)}} = \frac{SSR/1}{SSE/(n-2)} = \frac{MSR}{MSE}$$

$$= \frac{R^2(n-2)}{(1-R^2)}, \text{ αν γνωρίζουμε την τιμή του } R^2$$

συντελεστή προσδιορισμού  $R^2$ .

