

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

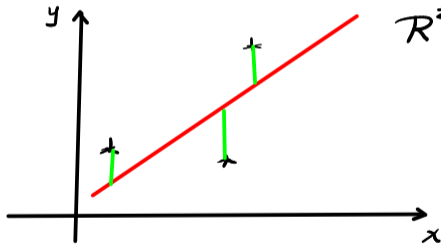
Κώστας Σμαραγδάκης (kesmarag@gmail.com)

20-03-2023

Συντελεστής Προσδιορισμού (Coefficient of Determination)

$$SST = SSR + SSE \Rightarrow SSR = SST - SSE$$

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



$$R^2: (x_i, y_i) \rightarrow [0, 1]$$

Συντελεστής Προσδιορισμού (Coefficient of Determination)

$$R^2 = \frac{SSR}{SST}$$

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

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

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

$$R^2 = \frac{\overbrace{\left[\frac{1}{N-1} SS_{xy} \right]^2}^{\text{Cov}[x,y]}}{\underbrace{\frac{1}{N-1} SS_{xx}}_{\text{Std}(x)} \underbrace{\frac{1}{N-1} SS_{yy}}_{\text{Std}(y)}}$$

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

$$SS_{xy} = \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$

$$SS_{xx} = \sum_{n=1}^N (x_n - \bar{x})^2$$

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

Συντελεστής Προσδιορισμού (Coefficient of Determination)

Παράδειγμα

Βρείτε τον συντελεστή προσδιορισμού του συνόλου δεδομένων:

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

$$N = 4$$

x	y	x ²	y ²	xy
0	1	0	1	0
1	3	1	9	3
2	4	4	16	8
5	4	25	16	20
8	12	30	42	31

$$\begin{aligned}SS_{xx} &= \sum x^2 - \frac{1}{N} (\sum x)^2 \\ &= 30 - \frac{1}{4} 8^2 = \\ &= 30 - 16 = 14\end{aligned}$$

$$SS_{yy} = 42 - \frac{1}{4} 12^2 = 6$$

$$\begin{aligned}SS_{xy} &= \sum xy - \frac{\sum x \sum y}{N} = \\ &= 31 - \frac{8 \cdot 12}{4} = 7\end{aligned}$$

$$R^2 = \frac{7^2}{6 \cdot 14} = \frac{7}{12}$$

Δειγματική Κατανομή της Κλίσης b

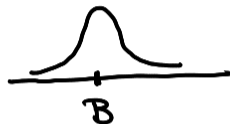
$$y = A + Bx + \varepsilon, \quad \mu_{y|x} = A + Bx, \quad A, B \text{ 'άγνωστα}$$

μέση τιμή, τυπική απόκλιση και κατανομή του b

$$\alpha, b, \quad b = \frac{SS_{xy}}{SS_{xx}}, \quad \alpha = \bar{y} - b\bar{x}$$

$$\hat{\mu}_{y|x} = \alpha + bx$$

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



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

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

$$b \sim \mathcal{N}(\mu_b, \sigma_b^2)$$

- ▶ Όταν το σ είναι άγνωστο δεν μπορούμε να υπολογίσουμε το σ_b

Εκτιμητήρια της τυπικής απόκλισης του b

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

αυτό μπαίνει
να υπολογιστεί.

$$b = \frac{SS_{xy}}{SS_{xx}} \quad \text{Τυχαία μεταβλητή}$$

$$SS_{xy} = \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y}) = \sum_{n=1}^N (x_n - \bar{x}) y_n - \bar{y} \sum_{n=1}^N (x_n - \bar{x}) \rightarrow 0$$

$$\sum_{n=1}^N (x_n - \bar{x}) = \sum x_n - N\bar{x} = N \underbrace{\frac{1}{N}}_{\bar{x}} \sum x_n - N\bar{x} = 0$$

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

$$\text{Άρα } b = \frac{\sum_{n=1}^N (x_n - \bar{x}) y_n}{\sum_{n=1}^N (x_n - \bar{x})^2} = \sum_{n=1}^N \frac{(x_n - \bar{x})}{\sum_{k=1}^N (x_k - \bar{x})^2} y_n = \sum_{n=1}^N C_n y_n$$

\swarrow αρ. βθ. \nwarrow Τυχαίες μεταβλητές
 C_n y_n

$$y_n = \overbrace{\alpha + bx_n} + \varepsilon$$

X_1, X_2 . ανεξάρτητες τότε $\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2)$

$$X = \alpha Y \Rightarrow \text{Var}(X) = \alpha^2 \text{Var}(Y)$$

$$\text{Var}(b) = \sum_{n=1}^N C_n^2 \text{Var}(y_n) = \sigma_e^2 \sum_{n=1}^N C_n^2$$

$$\sum_{n=1}^N C_n^2 = \frac{\sum_{n=1}^N (x_n - \bar{x})^2}{\left(\sum_{k=1}^N (x_k - \bar{x})^2 \right)^2} = \frac{\sum_{n=1}^N (x_n - \bar{x})^2}{\left[\sum_{n=1}^N (x_n - \bar{x})^2 \right]^2} = \frac{1}{SS_{xx}}$$

$$\text{Var}(b) = \frac{\sigma_e^2}{SS_{xx}}$$

$$\sigma_b^2 = \frac{\sigma_e^2}{SS_{xx}}$$

$$\sigma_b = \frac{\sigma_e}{\sqrt{SS_{xx}}}$$

$$\underline{b \sim t(B, s_b^2)} \quad \eta \quad b \sim N(B, \sigma_b^2)$$

Το $(1 - \alpha) * 100\%$ διάστημα εμπιστοσύνης για το Β είναι:

$$Be \quad [b - ts_b, b + ts_b]$$

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

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

- Περιθώριο σφάλματος: $E = ts_b$

Παράδειγμα

Για επτά νοικοκυριά μιας πόλης έχουμε τα ακόλουθα ζεύγη ετήσιου εισοδήματος και εξόδων σίτισης

$$\begin{array}{c} X \quad Y \\ \{(55, 14), (83, 24), (38, 13), (61, 16), (33, 9), (49, 15), (67, 17)\} \end{array}$$

1. Βρείτε την προσεγγιστική ευθεία γραμμικής παλινδρόμησης ($\hat{y} = a + bx$) χρησιμοποιώντας τη μέθοδο ελαχίστων τετραγώνων.
2. Υπολογίστε το 95% διάστημα εμπιστοσύνης για την παραμετρο B του πληθυσμού ($y = A + Bx$).

$$\textcircled{1} \quad b = \frac{SS_{xy}}{SS_{xx}}, \quad a = \bar{y} - b\bar{x}$$

$$N = 7$$

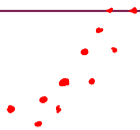
$$df = N - 2 = 5$$

$$\textcircled{2} \quad B \in \left[b - \overset{2.571}{\underset{1}{t}} s_b, b + \overset{2.571}{\underset{1}{t}} s_b \right]$$

Διάστημα Εμπιστοσύνης του Β

cum. prob	t _{.50}	t _{.75}	t _{.90}	t _{.95}	t _{.99}	t _{.995}	t _{.9975}	t _{.999}	t _{.9995}	t _{.9999}	
one-tail	0.50	0.25	0.20	0.15	0.10	0.05	0.025	0.01	0.005	0.001	
two-tails	1.00	0.50	0.40	0.30	0.20	0.10	0.05	0.02	0.01	0.002	
df											
1	0.000	1.000	1.376	1.963	3.078	6.314	12.71	31.82	63.66	318.31	636.62
2	0.000	0.816	1.061	1.386	1.886	2.920	4.303	6.965	9.925	22.327	31.599
3	0.000	0.765	0.978	1.250	1.638	2.353	3.182	4.541	5.841	10.215	12.924
4	0.000	0.741	0.941	1.190	1.533	2.132	2.776	3.747	4.604	7.173	8.610
5	0.000	0.727	0.920	1.156	1.476	2.015	2.571	3.365	4.032	5.893	6.869
6	0.000	0.718	0.906	1.134	1.440	1.943	2.447	3.143	3.707	5.208	5.959
7	0.000	0.711	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	0.000	0.706	0.889	1.108	1.397	1.860	2.306	2.896	3.355	4.501	5.041
9	0.000	0.703	0.883	1.100	1.383	1.833	2.262	2.821	3.250	4.297	4.781
10	0.000	0.700	0.879	1.093	1.372	1.812	2.228	2.764	3.169	4.144	4.587
11	0.000	0.697	0.876	1.088	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	0.000	0.695	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.930	4.318
13	0.000	0.694	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.852	4.221
14	0.000	0.692	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.787	4.140
15	0.000	0.691	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733	4.073
16	0.000	0.690	0.865	1.071	1.337	1.746	2.120	2.583	2.921	3.686	4.015
17	0.000	0.689	0.863	1.069	1.333	1.740	2.110	2.567	2.898	3.646	3.965
18	0.000	0.688	0.862	1.067	1.330	1.734	2.101	2.552	2.878	3.610	3.922
19	0.000	0.688	0.861	1.066	1.328	1.729	2.093	2.539	2.861	3.579	3.883
20	0.000	0.687	0.860	1.064	1.325	1.725	2.086	2.528	2.845	3.552	3.850
21	0.000	0.686	0.859	1.063	1.323	1.721	2.080	2.518	2.831	3.527	3.819
22	0.000	0.686	0.858	1.061	1.321	1.717	2.074	2.508	2.819	3.505	3.792
23	0.000	0.685	0.858	1.060	1.319	1.714	2.069	2.500	2.807	3.485	3.768
24	0.000	0.685	0.857	1.059	1.318	1.711	2.064	2.492	2.797	3.467	3.745
25	0.000	0.684	0.856	1.058	1.316	1.708	2.060	2.485	2.787	3.450	3.725
26	0.000	0.684	0.856	1.058	1.315	1.706	2.056	2.479	2.779	3.435	3.707
27	0.000	0.684	0.855	1.057	1.314	1.703	2.052	2.473	2.771	3.421	3.690
28	0.000	0.683	0.855	1.056	1.313	1.701	2.048	2.467	2.763	3.408	3.674
29	0.000	0.683	0.854	1.055	1.311	1.699	2.045	2.462	2.756	3.396	3.659
30	0.000	0.683	0.854	1.055	1.310	1.697	2.042	2.457	2.750	3.385	3.646
40	0.000	0.681	0.851	1.050	1.303	1.684	2.021	2.423	2.704	3.307	3.551
60	0.000	0.679	0.848	1.045	1.296	1.671	2.000	2.390	2.660	3.232	3.460
80	0.000	0.678	0.846	1.043	1.292	1.664	1.990	2.374	2.639	3.195	3.416
100	0.000	0.677	0.845	1.042	1.290	1.660	1.984	2.364	2.626	3.174	3.390
1000	0.000	0.675	0.842	1.037	1.282	1.646	1.962	2.330	2.581	3.098	3.300
Z	0.000	0.674	0.842	1.036	1.282	1.645	1.960	2.326	2.576	3.090	3.291
	0%	50%	60%	70%	80%	90%	95%	98%	99%	99.8%	99.9%
	Confidence Level										

$$\mathbb{R}^2 \subset \mathbb{R}^1$$



$$y = e^x \Rightarrow \ln y = x \overset{1}{\ln e} \Rightarrow \ln y = x$$

$$y_n = A + Bx_n + \varepsilon_n$$

$$\ln y_n = A + Bx_n + \varepsilon_n$$

► Όταν τα x και y δεν συνδέονται με γραμμικό τρόπο.

► Αλλά υπάρχει μετασχηματισμός $g : y \rightarrow y'$ τέτοιος ώστε x και y' να μπορούν να περιγραφούν με ένα γραμμικό μοντέλο.

$$y_n^* = A + Bx_n + \varepsilon_n$$

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

Παραδείγματα

- $y = e^x$
- $y = x^2$
- $y = \frac{1}{x}$
- $y = \log(x)$

$$\exp(\ln \hat{y}) = e^{\alpha + bx}$$

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

Παράδειγμα

Βρείτε κατάλληλο μετασχηματισμό για το παρακάτω σύνολο δεδομένων ώστε να είναι εφαρμόσιμο το μοντέλο γραμμικής παλινδρόμησης.

$$\{(0, 1), (1, 2), (4, 14), (5, 25), (6, 35)\} \rightarrow (s, s^2), (6, 6^2 - 1)$$

$$(0, 0^2 + 1) \downarrow (1, 1^2 + 1) \rightarrow (4, 4^2 - 2)$$

$$y = x^2 \quad \text{ή} \quad \sqrt{y} = x$$

$$\{(0, 1), (1, \sqrt{2}), (4, \sqrt{14}), (5, 5), (6, \sqrt{35})\} \rightarrow \mathbb{R}^2 \approx \mathbb{1}$$

↳ α, b τ.ω. SSE να γίνει ελάχιστο.

$$\sqrt{y} = \alpha x + b \Rightarrow y = (\alpha x + b)^2$$

Γραμμική Συσχέτιση (Linear Correlation)

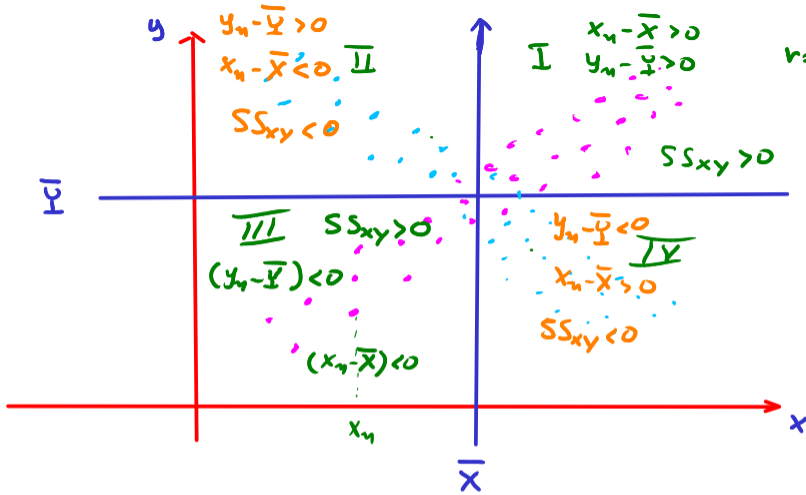
$$R^2 = \frac{SS_{xy}^2}{SS_{xx}SS_{yy}} \quad (\text{Συντελεστής Προσδιορισμού})$$

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

Σχέση μεταξύ συντελεστών γραμμικής συσχέτισης και προσδιορισμού

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

$$\text{sign}(SS_{xy}) = \begin{cases} 1, & SS_{xy} \geq 0 \\ -1, & SS_{xy} < 0 \end{cases}$$



$r \in \{-1, 1\}$
 $r = \text{sign}(SS_{xy}) \sqrt{R^2} \in [-1, 1]$

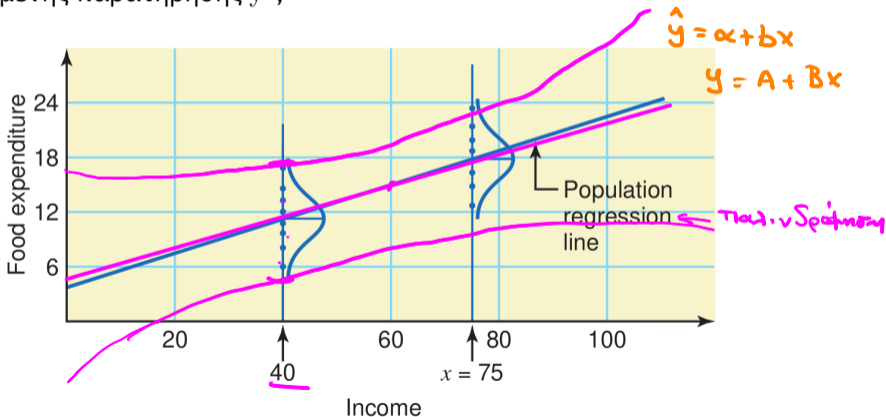
$$SS_{xy} = \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$



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^* ;

$x_n \rightarrow \hat{y}_n$
 (x_n, y_n)



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

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

το πινάκιο απόκλιση των σεφαιρών.

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

$$a + bx$$

Διάστημα εμπιστοσύνης

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

$$\mu_{y|x^*} \in [\hat{\mu}_{y|x^*} - t_{s_{\hat{\mu}_{y|x^*}}}, \hat{\mu}_{y|x^*} + t_{s_{\hat{\mu}_{y|x^*}}}]$$

$$B : \sigma_B = \frac{\sigma_e}{\sqrt{SS_{xx}}}$$

$$S_B = \frac{Se}{\sqrt{SS_{xx}}}$$

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

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

t - κατανομή σ_e άγνωστο.

- Περιθώριο σφάλματος: $E = t_{s_{\hat{\mu}_{y|x^*}}}$

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

$$\hat{y}_* = \alpha + bx_* + e_{x_*}$$

$$y_* = A + Bx_* + E_{x_*}$$

↑
δεν θα έχει πάντα την ίδια τιμή.

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

$$\hat{\mu}_{y|x_*} = \alpha + bx_*$$

Διάστημα εμπιστοσύνης

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

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

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

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

$$\sigma_{\hat{y}_*}^2 = \sigma_{\hat{\mu}_{y|x_*}}^2 + \sigma_{e_{x_*}}^2 \Rightarrow S_{\hat{y}_*}^2 = S^2 \hat{\mu}_{y|x_*}^2 + S^2 e$$

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

$$S_{\hat{y}_*}^2 = S_e^2 \left(\frac{1}{N} + \frac{(x_* - \bar{X})^2}{SS_{xx}} \right) + S_e^2 = S_e^2 \left[1 + \frac{1}{N} + \frac{(x_* - \bar{X})^2}{SS_{xx}} \right]$$

$$y^* \quad \hat{y}^*$$

$$y^* = \underbrace{A + Bx^*}_{\hat{y}^*} + \varepsilon_{x^*}$$

$$h_{y|x^*}$$

$$\hat{h}_{y|x^*}$$

$$h_{y|x^*} = E[y^*]$$

$$\boxed{y^*} \quad S_e \sqrt{1 + \frac{1}{N} + \frac{(x^* - \bar{X})^2}{SS_{xx}}}$$

$$\hat{y}^* = \underbrace{a + bx^*}_{\hat{h}_{y|x^*}} + e_{x^*}$$

$$\boxed{h_{y|x^*}} \quad S_e \sqrt{\frac{1}{N} + \frac{(x^* - \bar{X})^2}{SS_{xx}}}$$

Παράδειγμα

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

$N = 4$

$x^* = 1.5$

$$S_{\hat{y}^*} = S_e \sqrt{1 + \frac{1}{4} + \frac{(x^* - \bar{x})^2}{SS_{xx}}}$$

$$df = N - 2 = 2 \quad \pm$$

$$y^* \in [\hat{y}^* - \pm S_{\hat{y}^*}, \hat{y}^* + \pm S_{\hat{y}^*}]$$