

6. Variancia analízis

Több minta szórásnégyzetének (variációjának) összehasonlításán alapul a statisztika egyik nagy fejezete, a variancia analízis. A vizsgálatok célja ennek alkalmazásakor ugyanaz, mint a két mintára kiterjedő statisztikai próbáké volt: sokaságok egyezésének vagy eltérésének valószínűsítése. Meg kell jegyezni, hogy amikor a mintaszórások eltérés-valószínűségét F próbával határozzuk meg, a mintaelemek legyenek függetlenek és normális eloszlásúak.

A módszer lényegét számos variancia analitikus eljárás közül a legegyszerűbbön, az "egytényezős" variancia elemzésen mutatjuk be.

Több sokasággal foglalkozunk, amelyekről feltesszük, hogy $N(\mu_i, \sigma^2)$ eloszlásúak, ahol μ_i az i -edik sokaság várható értéke, σ^2 pedig a sokaságok megegyező variációjára. Kérdésünk az, a minták elkerülhetetlen eltérése véletlen-e, avagy érvényesült valami olyan hatás, aminek alapján a sokaságok nem tekinthetők megegyezőnek.

Szabatosan a

$$H_0 : \mu_i = \mu \quad i = 1, 2, \dots, m \quad (6.1)$$

nullhipotézis elfogadásáról vagy elvetéséről van szó. Vegyük észre, hogy különös módon most szórások összehasonlításával középértékek eltéréséről ítélkezünk.

Az eljárást példán mutatjuk be. Tegyük fel, abban kívánunk dönteni, hogy három, L1, L2 és L3 laboratórium egyenlő megbízhatóan dolgozik-e, avagy a laboratóriumokból érkező eredményeket fenntartással kell fogadni. A vizsgálathoz a három laboratórium 1,5 tömeg% ként tartalmazó gázolajat kap, amelyet ugyanazzal a (megegyező szórású) szabványos módszerrel kell megvizsgálnia. L1 labor $n_1 = 3$ párhuzamos mérést végez, L2 labor $n_2 = 5$ -öt, L3 $n_3 = 4$ -et. A beküldött eredményeket a 6.1 táblázatban bemutatott elrendezésű táblázatba foglaljuk. Itt x_{ij} jelenti a j -edik laboratórium i -edik mérési értékét. (A számszerű értékek a 6.3 táblázatban találhatóak).

6.1 táblázat. A variancia analízis alapadatai

	j = 1	j = 2	j = 3	Sorösszegek /átlagok
i = 1	x_{11}	x_{12}	x_{13}	
i = 2	x_{21}	x_{22}	x_{23}	
i = 3	x_{31}	x_{32}	x_{33}	
i = 4		x_{42}	x_{43}	
i = 5		x_{52}		
Oszlopösszeg	$\sum_{i=1}^{n_1} x_{i1}$	$\sum_{i=1}^{n_2} x_{i2}$	$\sum_{i=1}^{n_3} x_{i3}$	$\sum_{j=1}^m \sum_{i=1}^{n_j} x_{ij}$
Elemzészám	n_1	n_2	n_3	$n = \sum_{j=1}^3 n_j$
Szab.fok	$n_1 - 1$	$n_2 - 1$	$n_3 - 1$	$n - 3$
Átlag	\bar{x}_1	\bar{x}_2	\bar{x}_3	$\bar{x} = \frac{1}{n} \sum_{j=1}^3 \sum_{i=1}^{n_j} x_{ij}$
Eltérésnégyzet- összeg	$\sum_{i=1}^3 (x_{i1} - \bar{x}_1)^2$	$\sum_{i=1}^5 (x_{i2} - \bar{x}_2)^2$	$\sum_{i=1}^4 (x_{i3} - \bar{x}_3)^2$	$\sum_{j=1}^m \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$ $= SS_{\text{intra}}$
Variancia	$\frac{\sum_{i=1}^3 (x_{i1} - \bar{x}_1)^2}{n_1 - 1}$	$\frac{\sum_{i=1}^5 (x_{i2} - \bar{x}_2)^2}{n_2 - 1}$	$\frac{\sum_{i=1}^4 (x_{i3} - \bar{x}_3)^2}{n_3 - 1}$	$\frac{1}{n - 3} \sum_{j=1}^m \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$ $= MS_{\text{intra}}$
Átlag csoportok között	$\frac{1}{n} \sum_{j=1}^m \sum_{i=1}^{n_j} x_{ij}$			
Eltérésnégyzet- összeg csoportok között	$SS_{\text{inter}} = \sum_{j=1}^m n_j (\bar{x}_j - \bar{x})^2$			
Variancia csoportok között	$MS_{\text{inter}} = \frac{1}{m - 1} \sum_{j=1}^m n_j (\bar{x}_j - \bar{x})^2$			

A táblázatban látható két variancia, az MS_{intra} és MS_{inter} érték közül az első a kénmeghatározó módszer szórását, véletlen hibáját becsli. A második a laborok középértékeinek eltérését tükrözi azok közös középértékétől. Belátható, hogy ha a középértékek egymástól jobban eltérnek, mint amennyit a módszer szórása megenged, akkor a laboratóriumok között szignifikáns eltérés van. A döntés az MS_{intra} és MS_{inter} varianciák F próbáján alapul. Ha a kapott F nagyobb, mint a kritikus $F(\alpha, v_1, v_2)$ érték, akkor a (6.1) nullahipotézist elvetjük.

A variancia analízisnek ezeket a lépéseit a 6.2 táblázat mutatja.

6.2 táblázat. A variancia analízis eredményei

	SS	Szab.fok	MS	\hat{F}	p
Csoportokon belül	SS_{intra}	v_{intra}	MS_{intra}		
Csoportok között	SS_{inter}	v_{inter}	MS_{inter}	--	--
Összesen	SS_{total}	v_{total}	--	--	--

A táblázat legelső sorában az

$$SS_{\text{intra}} + SS_{\text{inter}} = SS_{\text{total}} = \sum_{j=1}^m \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2 \quad (6.2)$$

egyenlőségnek kell (matematikai okokból) teljesülnie. Ez hasznos ellenőrzési lehetőség. Ugyanez áll a szabadsági fokokra is. A táblázatban szereplő p érték azt adja meg, hogy a kapott \hat{F} hányadosnál *nagyobb* értékek mely valószínűséggel fordulnak elő. A statisztikus döntést nyilván ennek alapján is meg lehet hozni. A variancia analízis algoritmusai azonban legtöbbször kéri az α tévedési valószínűséget és megadják F kritikus értékét.

A bevezetésben bemutatott példa számszerű eredményeit a 6.3 táblázat mutatja be.

6.3 táblázat. Egytényezős variancia analízis

	L1	L2	L3	sorösszeg	Sorátlag
x1.	1,5	1,6	1,3		
x2.	1,55	1,72	1,3		
x3.	1,47	1,4	1,4		
x4.		1,48	1,45		
x5.		1,55			
x_{ij} összegek	4,52	7,75	5,45	17,72	
Mérésszámok	3	5	4	12	
Szab.fokok	2	4	3	9	
Átlagok	1,506666667	1,55	1,3625		1,476666
Eltérésnégyzet-összeg	0,0032666	0,0588	0,016875	0,0789416	
Variációk	0,0016333	0,0147	0,005625		0,008771
VARIANCIA ANALÍZIS					
Tényezők	SS	df	MS	\hat{F}	p-érték
Csoportok között	0,081725	2	0,0408625	4,658661459	0,040851902
Csoporton belül	0,078941667	9	0,008771296		
Összesen	0,160666667	11			

A kritikus F érték 5% tévedést megengedve, egyoldalas kérdésfeltevésnél 4.256 lenne. Ennél a kapott F érték nagyobb, így a nullahipotézist, miszerint a laboratóriumok egyformán dolgoznak elvetjük. p értékből látjuk, hogy a döntés nem módfelett biztos, hiszen, ha "igazságosabbak" akarunk lenni, és csak 3% tévedést vállalnánk, a laboratóriumokat már nem tartanók különbözőnek.

7. Összefüggések vizsgálata

A matematikai statisztika eddig tárgyalt fejezetei többnyire egy valószínűségi változóval foglalkoztak, ha pedig többel, akkor is feltételezték azok egymástól való függetlenségét. Nyilvánvaló ugyanakkor, hogy vizsgált rendszereket leíró, azokra ható változók közötti összefüggések elsődrendű fontosságúak.

Az összefüggések többféle szempontból tárgyalhatók, pl. abból, hogy okságiak-e, avagy nem azok, hogy van-e róluk előzetes ismeretünk avagy csak tapasztalati leírásunk stb. A továbbiakban aszerint tegyünk különbséget, hogy két (vagy több) *valószínűségi változó* összefüggésével kell foglalkozni, avagy nem valószínűségi (*determinisztikus*) *változók* hatnak egy valószínűségi változóra. Ez utóbbi változó legtöbbször azért tekintendő valószínűségi változónak, mert pontos, valódi értékét rátelepedett hiba terheli:

$$Y = Y_{\text{valódi}} + \varepsilon \quad (7.1)$$

Ebben az esetben Y eloszlása megegyezik ε eloszlásával, annyi különbséggel, hogy ha ε várható értéke 0, akkor Y -é $Y_{\text{valódi}}$.

7.1 Valószínűségi változó függése determinisztikus változó(k)tól.

Bizonyos, hogy a gyakorlatban változót hibátlanul nem lehet mérni vagy beállítani, elvben tehát determinisztikus változó nincs. Az azonban mégis elfogadható, hogy egyes független változók két-három nagyságrenddel pontosabbak, mint a függő, így nem okoz nagy hibát, ha azokat determinisztikusnak tekintjük.

7.1.1 A legkisebb négyzetek elve

Akár valószínűségi változó függ determinisztikus változó(k)tól, akár determinisztikus, szükségünk van egy matematikai összefüggésre (modellre), amelyik a függést leírja. Jelöljük a modellt F -fel. Ilyen modell lehet egy origón áthaladó vagy általános helyzetű egyenes, egy exponenciális függvény stb. A modellnek állandói (konstansai, paraméterei) vannak (meredekség, tengelymetszet) és független változói. Legyen az előbbiek jele a_1, a_2, \dots utóbbiaké x_1, x_2, \dots , de egyszerűség kedvéért tekintsünk most egyetlen x -t. A független változókat "prediktoroknak" vagy "regresszoroknak" is nevezik. A számított (jósolt, predikált) érték jel legyen y .

A modell tehát teljesen általánosan így fest:

$$Y = F(x_1, x_2, \dots, a_1, a_2, \dots)$$

Ha a paraméterek ismertek, beállított független változónál y kiszámítható. A feladatot azonban általában meg kell előznie a paraméterek meghatározása (*becslése*), egyfajta "kalibráció". Ismert x értékeknél párhuzamos kísérletekben meghatározzunk y mért értékeket, és a paramétereket tekintjük ismeretleneknek. Ha a mérések pontosak lennének, bármelyikből

ki lehetne számítani az ismeretlen a_1, a_2, \dots paramétereket. A kapott y értékeket azonban ismeretlen hibával mérjük:

$$\begin{aligned} F(x_1, a_1, a_2, \dots) &= y_1 + \varepsilon_1 \\ F(x_2, a_1, a_2, \dots) &= y_2 + \varepsilon_2 \\ &\dots \\ F(x_n, a_1, a_2, \dots) &= y_n + \varepsilon_n \end{aligned} \quad (7.2)$$

ezért a (7.2) egyenletekből mérésről mérésre más a paraméterek adódnának. Megállapodás szerint azokat az a_1, a_2, \dots értékeket fogadjuk el optimálisnak, amelyeknél a mért és számított értékek különbségnégyzeteinek összege minimális:

$$Q = \sum_{i=1}^n \left(y_i - \hat{y}_i \right)^2 = \min \quad (7.3)$$

\hat{y} az $F(x, a, b)$ modellel számított értéket jelöli.

Ez a követelmény a legkisebb négyzetek elve. Gyakorlati alkalmazására a következő fejezet ad példát.

Az ε_i hibákról nemcsak azt szokták feltételezni, hogy várható értékük 0, hanem azt is, szórásuk megegyezik. Ez az ún. *homoszkedasztikus* eset. Ha ugyanis a mérési hibák x változó mentén változnak (*heteroszkedasztikus* eset), a fellépő nagy eltérések (azok négyzetei) aránytalanul eltorzítják a minimum helyét, ezzel a paraméterek értékét. Ilyen esetben az eltéréseket *súlyozni* szokás, amivel a minimum követelmény így alakul:

$$\sum_{i=1}^n w_i (y_i - F(x_i, a_1, a_2, \dots))^2 = \min \quad (7.4)$$

A súly általában az adott x változóértéknél érvényes variancia reciproka:

$$w_i = \frac{1}{s_i^2} \quad (7.5)$$

A súlyozott legkisebb négyzetek módszerére más esetekben, pl. az y változóra alkalmazott transzformáció miatt is szükség lehet.

7.1.2 Egyenes paramétereinek becslése (lineáris regresszió)

a) Az egyenes állandói

A lineáris regressziónál az egyenes ismert egyenletének érvényességét tételezzük fel:

$$y = F(x, a, b) = a + bx \quad (7.6)$$

A paraméterek becslésére n darab x_i, y_i értékpárt használunk fel. A becslés gondolatmenetének megfelelően minimálni kell a mért és számított y értékek eltérése négyzetének összegét :

$$Q = \sum_{i=1}^n \left(y_i - \hat{y}_i \right)^2 = \sum_{i=1}^n (y_i - a - bx_i)^2 = \min \quad (7.7)$$

(\hat{y} az $F(x,a,b)$ modellel számított értéket jelöli. A további összefüggésekben az egyszerűség kedvéért a szummázás jelénél az i indexet elhagyjuk, sőt, ahol nem zavaró, az index a változók mellől is hiányzik).

Az a és b paraméterek függvényében Q négyzetösszeg nyilvánvalóan ott lesz minimális, ahol Q -nak a és b szerinti parciális deriváltjai 0 értékűek lesznek. Fenn kell tehát állnia, hogy

$$\frac{\partial Q}{\partial a} = -2 \sum (y_i - a - bx_i) = 0 \quad (7.8)$$

$$\frac{\partial Q}{\partial b} = -2 \sum (y_i - a - bx_i)x_i = 0. \quad (7.9)$$

A kapott egyenleteket egyszerűsítve, az összegezéseket tagonként végrehajtva és azokat rendezve az

$$a n + b \sum x = \sum y \quad (7.10)$$

$$a \sum x + b \sum x^2 = \sum xy \quad (7.11)$$

lineáris egyenletrendszer adódik, amelyből megoldás után a meredekségre a

$$\hat{b} = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \quad (7.12)$$

összefüggés, a tengelymetszetre pedig (7.9) egyenlet n -nel való osztása után az

$$\hat{a} = \bar{y} - b\bar{x} \quad (7.13)$$

képlet adódik.

(7.12) képlet könnyen számítható tényezőket tartalmaz. Mind számlálóját, mind nevezőjét aritmetikai műveletekkel átalakítható úgy is, hogy a képlet jobban megjegyezhető és többet mondó alakú legyen:

$$\hat{b} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad (7.14)$$

A paraméterbecslés ismertett elve és a (7.10), (7.11).egyenletek *többszörös* lineáris összefüggések paramétereinek becslésére is általánosíthatók. Az

$$y = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_m x_m \quad (7.15)$$

lineáris modell paramétereinek becslései n darab mérésből az $m + 1$ ismeretlenes

$$\begin{aligned}
 a_0 n + a_1 \sum x_{1i} + a_2 \sum x_{2i} + \dots + a_m \sum x_{mi} &= \sum y_i \\
 a_0 \sum x_{1i} + a_1 \sum x_{1i}^2 + a_2 \sum x_{1i} x_{2i} + \dots + a_m \sum x_{1i} x_{mi} &= \sum x_{1i} y_i \\
 a_0 \sum x_{2i} + a_1 \sum x_{1i} x_{2i} + a_2 \sum x_{2i}^2 + \dots + a_m \sum x_{2i} x_{mi} &= \sum x_{2i} y_i \quad (7.16) \\
 \dots \\
 a_0 \sum x_{mi} + a_1 \sum x_{1i} x_{mi} + a_2 \sum x_{2i} x_{mi} + \dots + a_m \sum x_{mi}^2 &= \sum x_{mi} y_i
 \end{aligned}$$

lineáris egyenletrendszer megoldásával lehet megkapni.

b) Szórásbecslések

Tételezzük fel, hogy az y mennyiség valóban lineáris függvénye x független változónak. Ebben az esetben csak a mérési hiba az oka annak, hogy a mért pontok nem esnek pontosan a becsült egyenesre. Ebből következik, hogy ebben az esetben az

$$s = \sqrt{\frac{\sum (y - \hat{y})^2}{n - 2}} \quad (7.17)$$

mennyiség a *mérési hiba becslése*. (A szabadsági fok azért $n - 2$, mert az egyenes két állandója két megkötést jelent az y mért értékek között). Ez a becslés egyébként

$$\begin{aligned}
 s &= \sqrt{\frac{1}{n - 2} \left(\sum (y - \bar{y})^2 - \frac{[\sum (x - \bar{x})(y - \bar{y})]^2}{\sum (x - \bar{x})^2} \right)} = \\
 &= \sqrt{\frac{1}{n - 2} \left(\sum (y - \bar{y})^2 - b \sum (x - \bar{x})(y - \bar{y}) \right)} \quad (7.18)
 \end{aligned}$$

módon is számítható.

Felvetődik ezután a *paraméterek* és a kapott paraméterekkel *számított* \hat{y} értékek szórásának becslése. Ezek rendre a következők:

$$s_a = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum (x - \bar{x})^2}} \quad (7.19)$$

$$s_b = s \sqrt{\frac{1}{\sum (x - \bar{x})^2}} \quad (7.20)$$

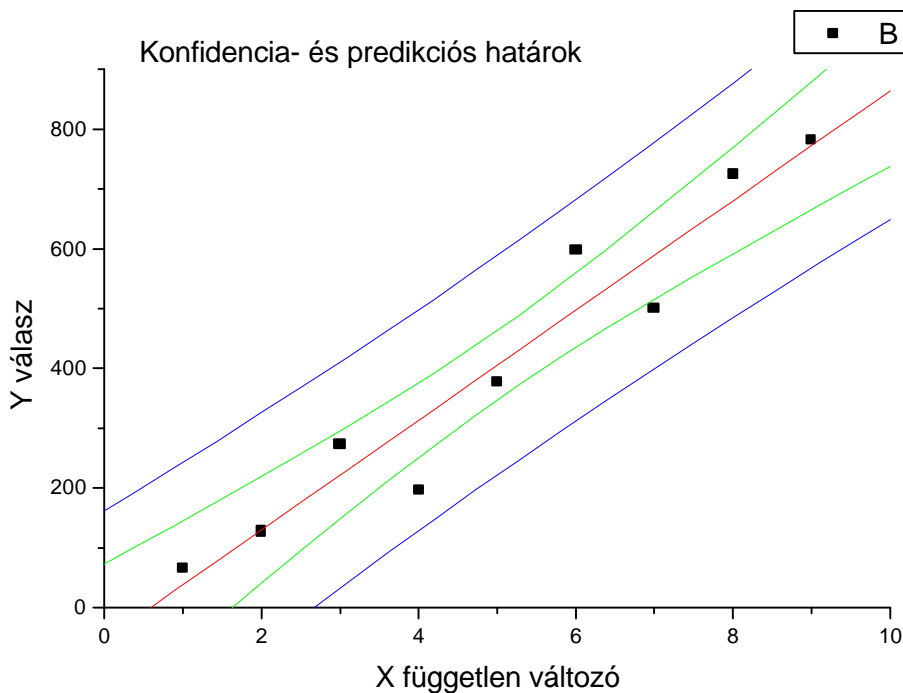
Normális eloszlású y értékek esetén adott valószínűségű konfidencia tartományt is megadhatunk a paraméterekhez.

$$\begin{aligned}
 \hat{a} \pm t_{n-2, \alpha} s_a \\
 \hat{b} \pm t_{n-2, \alpha} s_b
 \end{aligned} \quad (7.21)$$

A becsült paraméterekkel bármely x^* változóhoz kiszámítható egy \hat{y} érték várható értékének szórása:

$$s_{\hat{y}} = s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x - \bar{x})^2}} \quad (7.22)$$

Ez az x^* -től függő mennyiség a regressziós egyenes felett és alatt megadja a konfidencia "övet", tájékoztat arról, milyen határok között mozog a regressziós egyenes $1 - \alpha$ valószínűséggel. (7.1 ábra). Mint látható, a konfidencia tartomány az egyenes "súlypontjában" (az $x^* = \bar{x}$ pontban) legkeskenyebb és az egyenes két széle felé nő.



7.1 ábra. Regressziós egyenes és a megbízhatósági övek. Fekete négyzetek: mért értékek. Egyenes: regressziós egyenes. Belső öv: konfidencia határok, külső öv: predikciós határok.

Bármely x^* változónál *jövőben mért* y várható helyének bizonytalansága nagyobb. A *jóslási (predikciós) szórás*:

$$s_{y^*} = s \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x - \bar{x})^2}} \quad (7.23),$$

amit a 7.1 ábrán a külső öv mutat meg.

c) Az illesztés jósága

A 7.2 ábrát tanulmányozva látható, hogy az $y_i - \hat{y}_i$ távolság nem az egyetlen, amelyiket definiálni lehet. Beszélhetünk az $y_i - \bar{y}$ távolságról, és az $\hat{y}_i - \bar{y}$ távolságokról is. Belátható, hogy az $\hat{y}_i - \bar{y}$ távolságok

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (7.24)$$

négyzetösszege ill. az azokból számított "*modell okozta szórás*" arról vall, miért térnek el az y mért értékek átlaguktól amiatt, hogy azok x függvényei. Az is érthető, hogy ha ez a szórás összemérhető a kísérleti szórást becsülő, $y_i - \hat{y}_i$ különbségek

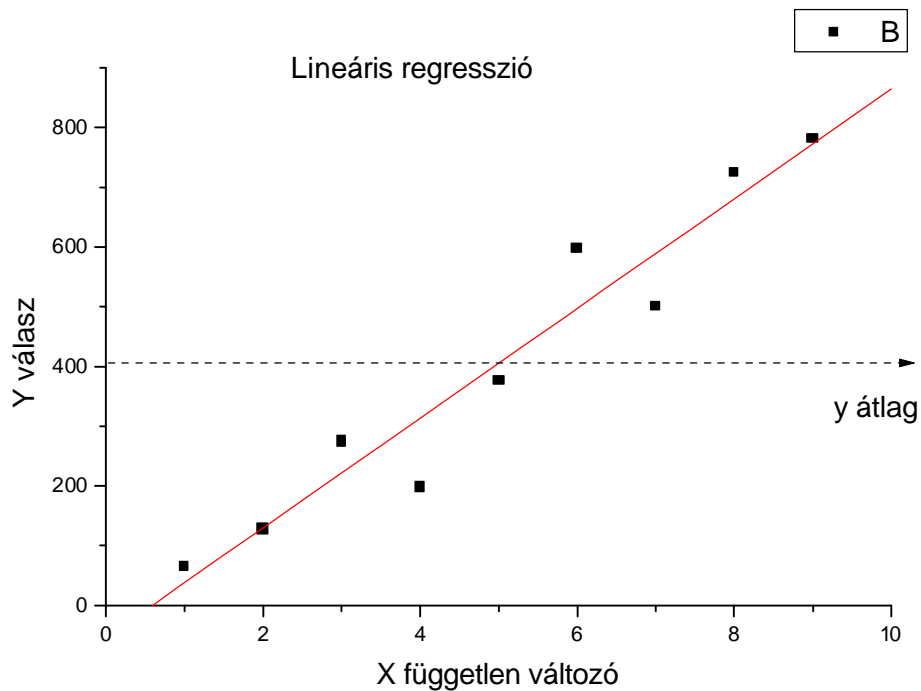
$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (7.25)$$

négyzetösszegéből számított "*reziduális szórással*", akkor kétséges a függés léte. Ezért számos esetben a regressziós számítást variancia analízis (l. 6. pont) követi, ahol ennek a két szórásnak négyzetét (varianciáját) F próbával hasonlítják össze. Minél nagyobb F, annál biztosabb a függés. Megjegyezhető, hogy a *modell okozta* eltérésnégyzetösszeg és a *reziduális* eltérésnégyzetösszeg kiadja az $y_i - \bar{y}$ távolságok

$$\sum_{i=1}^n (y_i - \bar{y})^2 \quad (7.26)$$

totális négyzetösszegét, mert

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}) \quad (7.27)$$



7.2 ábra. A regresszió megítéléséhez

A regresszió jóságát szokás a modell okozta eltérésnégyzetösszeg (7.24) és a teljes (totális) eltérésnégyzetösszeg (7.26) hányadosával is jellemezni.

$$r^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (7.28)$$

Ez a mennyiség azt adja meg, hogy az y értékek x -menti változásának hányadrésze tulajdonítható a lineárisnak tekintett függésnek. Az r^2 hányados az r korrelációs együttható (l. 7.2 pont) négyzete. Nem túl érzékeny mutató. Ha a mért pontok valamennyien pontosan az egyenesen vannak, értéke 1, de szemmeláthatóan szóró és nem is lineárisan függő mért értékeknél is viszonylag magas (0.9 feletti) lehet.

7.1.3 Nemlineáris paraméterbecslés

Nemlineáris (pl. hatványfüggvénnyel leírható, reciprokos, exponenciális) összefüggések paramétereinek becslésére a legkisebb négyzetek módszere ugyancsak alkalmazható. Az optimális paramétereket megadó

$$\sum_{i=1}^n (y_i - F(x_i, a_1, a_2, \dots))^2 = \min \quad (7.29)$$

kritériumban az F függvény *nemlineáris*, így a deriválás után (ha az lehetséges) kapott (7.8), (7.9) összefüggésekre emlékeztető egyenletek nem lineárisak és megoldásukhoz a numerikus matematika erre alkalmas módszerei használhatók.

Egyes esetekben nem kell ehhez a nehéz eljáráshoz folyamodni. Az

$$y = a_0 + a_1x + a_2x^2 + a_3x^3 \quad (7.30)$$

függvényben például helyettesíthetünk. Legyen $x_1 = x$, $x_2 = x^2$, $x_3 = x^3$. Ezzel a (7.30) egyváltozós összefüggés az

$$y = a_0 + a_1x_1 + a_2x_2 + a_3x_3 \quad (7.31)$$

háromváltozós, ám lineáris függvénné alakult, amelyikből a lineáris regresszió szabályai szerint a keresett paraméterek meghatározhatók. (l. a 7.16 egyenletrendszer). Hasonlóan lehet eljárni pl. az $y = a_0 + a_1 \ln x$ modell esetén is.

Függvények gyakran úgy linearizálhatók, hogy a transzformáció az y függő változót is érinti. Az

$$y = ae^{-bx} \quad (7.32)$$

összefüggés mindkét oldalát logaritmizálva az

$$\ln y = \ln a - bx \quad (7.33)$$

lineáris függvényhez jutunk, amelynek paraméterei lineáris regresszióval becsülhetők.

Az

$$\ln k = c - a/T \quad (7.34)$$

összefüggés $y = \ln k$ és $x = 1/T$ helyettesítéssel lineárisá alakítható.

A mért értékeket érintő átalakításoknál azonban figyelni kell arra a következményre, hogy az eredetileg (alkalmasint) egyenlő nagyságú hibák a transzformációk után eltérőkké, sőt esetleg aszimmetrikusakká válnak, így a becsült paraméterek torzítottak lehetnek és hibáikról gondos munka után lehet nyilatkozni. A *súlyozott legkisebb négyzetek módszerének* alkalmazása mindenképpen indokolt.

7.2 Valószínűségi változók összefüggése

Két valószínűségi változó ugyancsak összefügghet. Az összefüggés abban nyilvánul meg, hogy az egyik változó növekedése vagy csökkenése együttjár a másik változó csökkenésével vagy növekedésével. Ezt az összefüggést a kovariancia méri:

$$C(X,Y) = E[(X - E(X))(Y - E(Y))] \quad (7.35)$$

Ez a mennyiség pozitív, ha az X és Y valószínűségi változók együtt mozognak, negatív, ha ellentétesen. Szokás a kovarianciát a két változó szórásával osztva a -1 és $+1$ határok közé szorítani. A kapott mennyiség a korrelációs együttható:

$$\rho = \frac{C(X, Y)}{\sigma_X \sigma_Y} \quad (7.36)$$

Minél inkább megközelíti ρ a $+1$ vagy -1 értéket, annál szorosabb a két változó közötti összefüggés. Ha el is éri, a két változó egymás lineáris függvénye. Ha a korrelációs együttható 0 , akkor a két változó *korrelálatlan*. Ha két változó független, akkor korrelálatlan is. Fordítva a megállapítás nem érvényes. Attól, hogy ρ 0 -értékű, a két változó között még lehet függvénykapcsolat. Kivételt ez alól a normális eloszlású változók képviselnek.

A korrelációs együtthatót mintából az

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}} = \frac{\sum xy - \sum x \sum y / n}{\sqrt{(\sum x^2 - (\sum x)^2 / n)(\sum y^2 - (\sum y)^2 / n)}} \quad (7.37)$$

statisztika becsli.

Ha a korrelációs együttható szignifikáns 0 értékét akarjuk megvizsgálni, a

$$\hat{t} = |r| \sqrt{\frac{n-2}{1-r^2}} \quad (7.38)$$

értéket kell kiszámítani. Ha ez nagyobb, mint $t_{n-2, \alpha}$, akkor a $H_0: r = 0$ nullahipotézist **a** tévedési valószínűséggel elvetjük.

Korrelált valószínűségi változók lineáris összefüggését lehet regresszióval vizsgálni. A regressziós egyenes azonban más helyzetű, ha Y-t X függvényében, avagy X-et Y függvényében vizsgáljuk.

Általában fennáll, hogy a korreláció nem jelent szükségképpen oksági kapcsolatot. A függő és független változó fogalma ebben a környezetben gyakran értelmezhetetlen.