

# **CURSUL – III -**

**TESTE NEPARAMETRICE**

**DREAPTA DE REGRESIE**

## Teste neparametrice

Testul t pentru compararea mediilor depinde, în special pentru selecțiile de volum mic, de ipoteza că cele două populații sunt distribuite aproximativ normal și că dispersiile sunt practic egale.

De regulă, tehnicile statistice care se ocupă de variabilele continue se bazează pe ipoteza că variabila aleatoare are o distribuție normală de bază. Ipoteza nu este atât de restrictivă, deoarece de multe ori este posibil să o modificăm astfel încât să obținem alta, aproximativ normal distribuită. Suplimentar, dacă vom considera mediile, în concordanță cu teorema limită centrală, distribuția mediei probelor se apropie cu atât mai mult de distribuția normală, cu cât crește volumul probelor.

Și astfel, ca o concluzie practică, erorile se datorează mai curând lipsei de constanță a dispersiei sau lipsei de independență a variabilelor decât deviațiilor de la normalitate.

Pentru cazurile când nu stim distribuția variabilei, o cale alternativă este să aplicăm teste care nu necesită ipoteze despre tipul de distribuție.

Testele independente de distribuție, numite și teste de rang, înlocuiesc valorile variabilei cantitative observate cu rangurile lor. Testele neparametrice sunt valabile și pentru variabile normal distribuite, dar sunt mai puțin eficiente, pentru același prag de semnificație fiind necesare eșantioane mai mari decât pentru testele parametrice.

### Media și dispersia unui eșantion dintr-o populație finită.

Să considerăm o populație finită de  $N$  elemente, la care asociem numerele  $x_1, x_2, \dots, x_N$ . Dacă presupunem că toate elementele au aceeași probabilitate  $\frac{1}{N}$ , putem calcula media și dispersia populației:

$$(1) \quad \mu = E(X) = \sum_1^N x_i p_i = \frac{1}{N} \sum_1^N x_i$$

și

$$(2) \quad \sigma^2 = D(X) = E(X^2) - (E(X))^2 = \sum_1^N x_i^2 p_i - \left(\sum_1^N x_i p_i\right)^2 = \frac{1}{N} \sum_1^N x_i^2 - \frac{1}{N^2} \left(\sum_1^N x_i\right)^2 = \left(\frac{1}{N} - \frac{1}{N^2}\right) \sum_1^N x_i^2 - \frac{2}{N^2} \sum_{i \neq j} x_i x_j = \frac{N-1}{N^2} \sum_1^N x_i^2 - \frac{2}{N^2} \sum_{i \neq j} x_i x_j$$

Multimea tuturor selecțiilor posibile de mărimea  $n$  din populație va include:

$$(x_1, x_2, \dots, x_{n-1}, x_n)$$

$$(x_1, x_2, \dots, x_{n-1}, x_{n+1})$$

.

.

.

$$(x_{N-n+1}, x_{N-n+2}, \dots, x_N)$$

Aceste probe sunt formate prin alegerea a n elemente din N. Există  $C_N^n$  căi de a alege o astfel de probă. Încă o dată, presupunem că fiecare probă are aceeași probabilitate de a fi selectată,  $\frac{1}{C_N^n}$ .

Să considerăm media selecției j:  $\overline{X}_j = \frac{1}{n} \sum_{i=1}^n x_{ji}$  și să considerăm variabila aleatoare  $\overline{X} = (\overline{X}_j)_{j=1, C_N^n}$

Valoarea medie a variabilei  $\overline{X}$  este

$$E(\overline{X}) = \sum_{j=1}^{C_N^n} \overline{X}_j p_j = \frac{1}{C_N^n} \sum_{j=1}^{C_N^n} \overline{X}_j = \frac{1}{C_N^n} \left[ \frac{1}{n} (x_1 + x_2 + \dots + x_{n-1} + x_n) + \frac{1}{n} (x_1 + x_2 + \dots + x_{n-1} + x_{n+1}) + \dots + \frac{1}{n} (x_{N-n+1} + x_{N-n+2} + \dots + x_N) \right]$$

Acum să considerăm de câte ori intră în sumă orice  $x_i$  particular, să spunem  $x_1$ . Probele care conțin  $x_1$  se obțin prin selectarea a (n-1) alte elemente din populația disponibilă de (N-1) elemente și, aceasta se poate face în  $C_{N-1}^{n-1}$  moduri. Vor fi deci  $C_{N-1}^{n-1}$  probe conținând  $x_1$  și la fel se aplică pentru fiecare  $x_i$ .

$$C_N^n = \frac{N!}{n!(N-n)!} = \frac{N}{n} \frac{(N-1)!}{(n-1)!(N-n)!} = \frac{N}{n} C_{N-1}^{n-1}$$

În consecință

$$(3) \quad E(\overline{X}) = \frac{1}{C_N^n} \left( \frac{1}{n} C_{N-1}^{n-1} \sum_1^N x_i \right) = \frac{1}{N} \sum_1^N x_i = \mu$$

ceea ce înseamnă că media mediei probei este egală cu media populației.

Pentru calcularea dispersiei folosim identitatea

$$(4) D(\bar{X}) = E(\bar{X}^2) - (E(X))^2$$

$$\text{Să considerăm } E(\bar{X}^2) = \sum_{j=1}^{C_N^n} \bar{X}_j^2 p_j = \frac{1}{C_N^n} \sum_{j=1}^{C_N^n} \bar{X}_j^2$$

Mai departe

$$\sum_{j=1}^{C_N^n} \bar{X}_j^2 = \left[ \frac{1}{n} (x_1 + x_2 + \dots + x_{n-1} + x_n) \right]^2 + \dots + \left[ \frac{1}{n} (x_{N-n+1} + x_{N-n+2} + \dots + x_N) \right]^2$$

Când ridicăm la pătrat fiecare termen, fiecare  $x_i$  va deveni  $x_i^2$  și, după cum vedem, fiecare  $x_i$  apare de  $C_{N-1}^{n-1}$  ori. Astfel

$$(5) \sum_{j=1}^{C_N^n} \bar{X}_j^2 = \frac{1}{n^2} C_{N-1}^{n-1} (x_1^2 + x_2^2 + \dots + x_N^2)$$

Ridicarea la pătrat a sumei dă deasemenea termeni de forma  $x_i x_j$  și fiecare termen va apare de  $C_{N-2}^{n-2}$ .

În consecință, putem scrie

$$(6) \frac{1}{C_N^n} \sum_{j=1}^{C_N^n} \bar{X}_j^2 = \frac{1}{C_N^n} \left[ \frac{1}{n^2} C_{N-1}^{n-1} (x_1^2 + x_2^2 + \dots + x_N^2) + \frac{2}{n^2} C_{N-2}^{n-2} (x_1 x_2 + \dots + x_{N-1} x_N) \right]$$

Pentru a înlocui în (4) punem  $(E(\bar{X}))^2$  în forma:

$$(7) (E(\bar{X}))^2 = \left[ \frac{1}{N} (x_1 + x_2 + \dots + x_{N-1} + x_N) \right]^2 = \frac{x_1^2 + x_2^2 + \dots + x_N^2}{N^2} + \frac{2(x_1 x_2 + \dots + x_{N-1} x_N)}{N^2}$$

Substituind (6) și (7) în (4), obținem:

$$(8) D(\bar{X}) = \left( \frac{1}{C_N^n} \frac{1}{n^2} C_{N-1}^{n-1} - \frac{1}{N^2} \right) (x_1^2 + x_2^2 + \dots + x_N^2) + \left( \frac{1}{C_N^n} \frac{2}{n^2} C_{N-2}^{n-2} - \frac{2}{N^2} \right) (x_1 x_2 + \dots + x_{N-1} x_N)$$

Coeficientul lui  $(x_1^2 + x_2^2 + \dots + x_N^2)$  se poate scrie ca

$$\frac{1}{C_N^n} \frac{1}{n^2} C_{N-1}^{n-1} - \frac{1}{N^2} = \frac{1}{C_{N-1}^{n-1}} \frac{1}{\frac{N}{n}} \frac{1}{n^2} C_{N-1}^{n-1} - \frac{1}{N^2} = \frac{1}{nN} - \frac{1}{N^2} = \frac{N-n}{nN^2} = \frac{N-n}{n(N-1)} \frac{N-1}{N^2}$$

și coeficientul lui  $(x_1 x_2 + \dots + x_{N-1} x_N)$  este

$$\frac{1}{C_{N-2}^{n-2}} \frac{2}{n(n-1)} \frac{2}{n^2} C_{N-2}^{n-2} - \frac{2}{N^2} = \frac{2(n-1)}{nN(N-1)} - \frac{2}{N^2} = -\frac{2}{N^2} \frac{N-n}{n(N-1)}$$

Apoi substituind aceste rezultate în (8), obținem:

$$(9) D(\bar{X}) = \frac{(N-n)}{n(N-1)} \left\{ \frac{N-1}{N^2} (x_1^2 + x_2^2 + \dots + x_N^2) - \frac{2}{N^2} (x_1x_2 + \dots + x_{N-1}x_N) \right\}$$

Partea din { } este exact  $\sigma^2$ , astfel încât

$$(10) D^2(\bar{X}) = \frac{\sigma^2}{n} \frac{N-n}{N-1} = \frac{\sigma^2}{n} \left( 1 - \frac{n-1}{N-1} \right)$$

Este de notat că dacă  $N \rightarrow \infty$ , atunci dispersia lui  $\bar{X} \rightarrow \frac{\sigma^2}{n}$ , forma ei obișnuită pentru o populație infinită, sau pentru experimentul de tip extracție din urnă cu întoarcerea bilelor extrase în urnă.

### Testul de rang Wilcoxon

Testul de rang Wilcoxon<sup>1</sup> este un test cu ipoteza nulă că două populații sunt identice, față de ipoteza alternativă că ele diferă printr-o translație lineară. Testul înlocuiește observațiile prin rangurile lor. Rangurile sunt repartizate la valorile din selecții în ordinea creșterii mărimii fără să țină cont de probele cărora le aparțin.

Să presupunem că o probă este de mărime  $n$  și alta de mărime  $N-n$ . Testul presupune că orice combinație de ranguri în aceste două grupuri este egal probabilă. Numărul total de moduri de grupare a rangurilor este  $C_N^n$ .

Consideram urmatorul exemplu

Nivelele plasmatice maxime ale ionului  $EDTA^{4-}$  după administrare i.m.

Tabelul 1.

Voluntar	CE	IA	BL	PM	MC	DP	SL
Prima zi	33,3	25,1	22,8	32,4	23,7	48,33	33,04
rangurile	9	3	1	7	2	11	8
a-3-a zi	25,4	31,2	28,4	39,2			
rangurile	4	6	5	10			

Privind rezultatele în a treia zi de tratament la proba de mărime  $n$ , suma rangurilor este  $4+6+5+10=25$ . Combinațiile de ranguri pentru care putem obține o sumă a rangurilor mai mică decât aceasta, pentru un  $n = 4$  dat sunt

$$1+2+3+4=10, 1+2+3+5=11, 1+2+3+6=13, 1+2+3+7=14, 1+2+3+8=15 \text{ etc.}$$

După cum se poate vedea nu este ușor să calculăm toate posibilitățile, astfel încât vom folosi faptul că media rangurilor unei probe este distribuită aproximativ normal cu parametri care sunt calculați în continuare.

Sunt disponibile tabelele care dau limitele de acceptare a ipotezei  $H_0$  pentru suma obținută, ca o funcție de  $n$ ,  $N$  și riscul asumat. Pentru exemplul nostru găsim în tabele, pentru  $\alpha = 0,05$ ,  $n_1 = 4$  și  $n_2 = 7$  intervalul  $11 - 25$ .

Fie  $R$  suma rangurilor și  $\bar{R}$  media rangurilor probei de mărime  $n$ . Conform (1), valoarea medie a lui  $R$  este  $E(\bar{R}) = \frac{1}{N} \sum_{i=1}^N x_i$ . În cazul nostru  $x_i$  sunt rangurile de  $N$  valori însemnând numerele  $1, 2, \dots, N$ . În consecință

$$E(\bar{R}) = \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{N} (1 + 2 + \dots + N) = \frac{1}{N} \frac{N(N+1)}{2} \Rightarrow E(\bar{R}) = \frac{N+1}{2}$$

Calculul lui  $\sigma^2$  dă:

$$\begin{aligned} \sigma^2 = D(X) &= E(X^2) - (E(X))^2 = \frac{1}{N} \sum_{i=1}^N x_i^2 - \frac{1}{N^2} \left( \sum_{i=1}^N x_i \right)^2 = \frac{1}{N} \sum_{i=1}^N i^2 - \frac{1}{N^2} \left( \sum_{i=1}^N i \right)^2 = \\ &= \frac{1}{N} \frac{N(N+1)(2N+1)}{6} - \frac{1}{N^2} \left( \frac{N(N+1)}{2} \right)^2 = \frac{N^2-1}{12} \end{aligned}$$

Dispersia lui  $\bar{R}$  se obține prin înlocuirea lui  $\sigma$  în (10)

$$(11) \quad D(\bar{R}) = \frac{\sigma^2}{n} \left( 1 - \frac{n-1}{N-1} \right) = \frac{N^2-1}{12n} \frac{N-n}{N-1} = \frac{(N+1)(N-n)}{12n}$$

În concluzie, variabila aleatoare  $\frac{\bar{R} - E(\bar{R})}{\sqrt{D(\bar{R})}} = \frac{\bar{R} - \frac{N+1}{2}}{\sqrt{\frac{(N+1)(N-n)}{12n}}}$  va fi repartizată

aproximativ  $N(0,1)$ .

Kruskal și Wallis<sup>2</sup> au observat că aproximația este îmbunătățită când valoarea  $\alpha$  este mai mare de 0,02 prin aducerea lui  $\bar{R}$  mai aproape de media lui cu  $\frac{1}{2n}$ .

În literatura medicală și biologică testul se mai numește Mann – Whitney și se utilizează notațiile  $n = n_1$  și  $N - n = n_2$  ( $n_1 \leq n_2$ ).

Când cel puțin unul din numerele  $n_1$  și  $n_2$  sunt mai mici decât 10, distribuția de probabilitate a sumei rangurilor pozitive R se poate calcula direct. Intervalele de încredere cu diverse probabilități (0,95; 0,99; etc.) pentru R se găsesc în tabele.

În exemplul nostru  $n = 4$ ,  $N = 11$ ,  $R = 25$ ,  $\bar{R} = \frac{25}{4} = 6,25$  și

$$z = \frac{\bar{R} - \frac{N+1}{2}}{\sqrt{\frac{(N+1)(N-n)}{12n}}} = \frac{6,25 - \frac{11+1}{2}}{\sqrt{\frac{(11+1)(11-4)}{12*4}}} = \frac{0,25}{\sqrt{\frac{7}{4}}} = 0,19$$

Valoarea obținută ne asigură că nu apare o acumulare a EDTA la orice nivel de risc  $\alpha$  din cele uzual utilizate.

Dacă facem corecția pentru continuitate

$$z = \frac{\bar{R} - \frac{N+1}{2} + \frac{1}{2n}}{\sqrt{\frac{(N+1)(N-n)}{12n}}} = \frac{6,25 - \frac{11+1}{2} + \frac{1}{8}}{\sqrt{\frac{(11+1)(11-4)}{12*4}}} = \frac{0,375}{\sqrt{\frac{7}{4}}} = 0,285$$

concluzia nu se schimbă.

### **Ajustarea pentru valori egale în testul Wilcoxon**

Dacă apar egalități, o alternativă pentru neglijarea lor este de a repartiza la aceste observații media rangurilor pe care le-ar fi primit dacă nu erau egale.

Să considerăm un grup de k egalități. Numerele întregi  $m+1, m+2, \dots, m+k$  sunt înlocuite cu media lor.

$$\frac{(m+1) + (m+2) + \dots + (m+k)}{k} = \frac{km + \frac{k(k+1)}{2}}{k} = m + \frac{k+1}{2}$$

Suma pătratelor ( $x_1^2 + x_2^2 + \dots + x_N^2$ ) este astfel redusă prin

$$(m+1)^2 + (m+2)^2 + \dots + (m+k)^2 - k \left[ m + \frac{(k+1)}{2} \right]^2 =$$

$$km^2 + 2(1+2+\dots+k)m + (1^2 + 2^2 + \dots + k^2) - km^2 - km(k+1) - \frac{k(k+1)^2}{4} =$$

$$2 \frac{k(k+1)}{2} m + \frac{k(k+1)(2k+1)}{6} - k(k+1)m - \frac{k(k+1)^2}{4} = \frac{k(k+1)}{12} (4k+2-3k-3) = \frac{(k-1)k(k+1)}{12} = \frac{T}{12}$$

Suma rangurilor rămâne neschimbată.

Astfel

$$\sigma^2 = \frac{1}{N} \sum_1^N x_i^2 - \frac{1}{N^2} \left( \sum_1^N x_i \right)^2 = \frac{1}{N} \left( \frac{N(N+1)(2N+1)}{6} - \frac{T}{12} \right) - \frac{1}{N^2} \left[ \frac{N(N+1)}{2} \right]^2 =$$

$$\frac{2N(N+1)(2N+1) - T - 3N(N+1)^2}{12N} = \frac{N(N+1)(4N+2-3N-3) - T}{12N} = \frac{N(N^2-1) - T}{12N}$$

și

$$D(\bar{R}) = \frac{N(N^2-1) - T}{12nN} \frac{N-n}{n-1}$$

## Teste referitoare la perechi de observații

### a) Testul semnelor

Să considerăm nivelele plasmatice maxime  $x_i$  ale unui medicament după o primă administrare la un număr de  $n$  voluntari sănătoși și  $y_i$  nivelele plasmatice maxime după trei zile de tratament. Fie  $\rho(x, y)$  probabilitatea de apariție a valorilor  $x$  și  $y$ . Dacă medicamentul nu se acumulează în organism, cele două seturi de concentrații sunt selecții ale aceleiași populații și  $\rho(x_i, y_i) = \rho(y_i, x_i)$  pentru toate perechile.

Aceasta implică simetria lui  $\rho(x, y)$  față de linia  $y - x = 0$ .

Să definim variabila aleatoare  $z = y - x$ .

Avem că  $P(y < x) = P(y > x) = \frac{1}{2}$  sau  $P(y - x < 0) = P(y - x > 0) = \frac{1}{2}$  care este mai

departe echivalent cu  $P(z < 0) = P(z > 0) = \frac{1}{2}$ . Astfel  $z$  va avea o mediană zero.

Mai departe definim variabilele  $z_i$  după cum urmează  $z_i = 1$  pentru  $z_i > 0$



și  $z_i = 0$  pentru  $z_i < 0$ .

Presupunem continuitatea distribuției de grup originală  $\rho(x, y)$ ,  $z$  va fi deasemenea continuă, și “intersecțiile” (cazurile  $x_i = y_i$ ) vor avea probabilitatea zero.

$z_i$  sunt independente, astfel încât suntem în situația binomială de a face  $n$  încercări independente, probabilitatea de success  $z_i = 1$  fiind  $\frac{1}{2}$  la fiecare încercare. Astfel,  $\sum_1^n z_i$  are o distribuție binomială cu parametrii  $p = \frac{1}{2}$  și  $n$ .

Distribuția de grup  $\rho(x, y)$  poate fi diferită în fiecare încercare, însă de fiecare dată  $P(z_i = 1) = \frac{1}{2}$  și astfel distribuția lui  $\sum_1^n z_i$  va fi neschimbată.

Alternativa ipotezei nule este ca în locul lui  $x_i$  să avem  $x_i' = x_i - d_i$ , ceea ce înseamnă că fiecare  $x_i$  descrește cu o cantitate  $d_i$ , unde  $d_i > 0$ . În acest caz  $\rho(x', y)$  nu va mai fi simetric, ci deplasat spre stânga și  $P(z_i > 0) = P(y_i - x_i' > 0) = P(y_i > x_i) > \frac{1}{2}$ .

Astfel,  $P(z_i > 0)$  nu va mai fi în mod necesar constantă și distribuția lui  $\sum_1^n z_i$  nu va mai fi o distribuție binomială.

Testul semnelor, dă pentru probabilitatea a  $k$  diferențe pozitive

$$P\left(\sum_1^n z_i \geq \frac{k}{n}, p = \frac{1}{2}\right) = \sum_{i=k}^n z_i C_n^i \left(\frac{1}{2}\right)^i \left(1 - \frac{1}{2}\right)^{n-i} = \frac{1}{2^n} \sum_{i=k}^n C_n^i = \frac{1}{2^n} \sum_{i=k}^n C_n^{n-i} = \frac{1}{2^n} \sum_{j=0}^{n-k} C_n^j$$

În cazurile simple, pentru  $k$  și  $n$  mici, această probabilitate se poate calcula direct.

Pentru valori mai mari, se poate folosi aproximația normală.

Sa luăm în considerare valorile nivelelor plasmatice ale ionului  $EDTA^{4-}$  (Tabelul 1) după administrarea i.m. la patru voluntari sănătoși.

Voluntar	CE	IA	BL	PM
Prima zi	33,3	25,1	22,8	32,4
a-3-a zi	25,4	31,2	28,4	39,2
$z_i$	-7,9	+6,1	+5,6	+6,8
$z_i$	0	1	1	1

Avem

$$P\left(\sum_1^4 z_i > \frac{3}{4}, p = \frac{1}{2}\right) = \frac{1}{2^4} \sum_{j=0}^{4-4} C_4^j = \frac{1}{2^4} C_4^0 = \frac{1}{2^4} = 0,06$$

ceea ce înseamnă că putem accepta ipoteza nulă privind egalitatea constantei de eliminare în prima zi cu cea din ziua a treia.

### b) Testul Wilcoxon pentru observații perechi

Wilcoxon a propus deasemenea un test pentru determinări pare în care rangurile sunt atribuite mărimii absolute a diferențelor și apoi se dă rangurilor semnul diferențelor.

Ipoteza nulă este că distribuția diferențelor este simetrică față de zero, astfel orice rang este pozitiv sau negativ cu aceeași probabilitate. Valorile egale primesc ca rang media rangurilor grupului.

Numărul total de moduri de sume de ranguri ce se pot obține este  $2^N$ .

Să atașăm rangurilor  $i$  variabilele aleatoare  $d_i$  ce iau valorile  $d_i=1$  când  $i$  este pozitiv și  $d_i=0$  când  $i$  este negativ.

Să considerăm suma rangurilor positive  $s = \sum d_i i$ . Media ei va fi

$$E(s) = E\left(\sum_1^N d_i i\right) = \sum_1^N i E(d_i)$$

$$\text{Dar } E(d_i) = 1 * \frac{1}{2} + 0 * \frac{1}{2} = \frac{1}{2} \text{ și } E(s) = \sum_1^N i \frac{1}{2} = \frac{N(N+1)}{4}$$

$$E(s^2) = E\left(\sum_1^N i d_i\right)^2 = E\left(\sum_1^N i^2 d_i^2 + 2 \sum_{i \neq j} i j d_i d_j\right) = \sum_1^N i^2 E(d_i^2) + 2 \sum_{i \neq j} i j E(d_i d_j)$$

$$\text{Însă } E(d_i^2) = 1^2 * \frac{1}{2} + 0^2 * \frac{1}{2} = \frac{1}{2} \text{ și } E(d_i d_j) = 0 * 0 * \frac{1}{4} + 0 * 1 * \frac{1}{4} + 1 * 0 * \frac{1}{4} + 1 * 1 * \frac{1}{4} = \frac{1}{4}$$

În consecință

$$E(s^2) = \frac{1}{2} \sum_1^N i^2 + \frac{1}{4} \sum_{i \neq j} 2ij = \frac{1}{2} \sum_1^N i^2 + \frac{1}{4} \left[ \left(\sum_1^N i\right)^2 - \sum_1^N i^2 \right]$$

Acum putem calcula dispersia lui  $s$

$$D(s) = E(s^2) - (E(s))^2 = \frac{1}{2} \sum_1^N i^2 + \frac{1}{4} \left[ \left( \sum_1^N i \right)^2 - \sum_1^N i^2 \right] - \frac{1}{4} \left( \sum_1^N i \right)^2 = \frac{1}{4} \sum_1^N i^2 = \frac{N(N+1)(2N+1)}{24}$$

În cazul în care apar egalități,  $\frac{(k-1)k(k+1)}{48}$  trebuie să fie scăzut pentru fiecare grup

de egalitati. O alternativă este de a scoate toate valorile egale din probă.

Să considerăm acum observațiile pare din experimentul ce a dus la datele din tabelul 3.

Tabelul 3 Nivelele plasmatice maxime ale  $EDTA^{4-}$  după administrarea i.m.

Voluntar	CE	IA	BL	PM	
Prima zi	33,3	25,1	22,8	32,4	
a-3-a zi	25,4	31,2	28,4	39,2	
Diferența	-7,9	+6,1	+5,6	+6,8	
$d_i$	0	1	1	1	
Rangul	-4	2	1	3	S=3+2+1=6

În acest caz  $N=4$  și  $z = \frac{s - E(s)}{\sqrt{D(s)}} = \frac{s - \frac{N(N+1)}{4}}{\sqrt{\frac{N(N+1)(2N+1)}{24}}} = \frac{6-5}{\sqrt{\frac{4*5*9}{24}}} = 0,27$  care este foarte

apropiat de valorile obținute anterior.

### c) Testul H

Testul H, sau testul Kruskal – Wallis<sup>3</sup> este o generalizare a testului Wilcoxon în cazul a k probe,  $k > 2$ . La fel ca și în testul Wilcoxon, observațiile primesc ranguri, și media rangurilor  $R_i$  se calculează pentru fiecare grup.

$$E(\bar{R}_i) = \frac{N+1}{2} \text{ și } D^2(\bar{R}_i) = \frac{(N+1)(N-n_i)}{12n_i}$$

Raportul  $\frac{\bar{R}_i - E(\bar{R}_i)}{\sqrt{D^2(\bar{R}_i)}}$  va fi repartizat  $N(0,1)$ , conform teoremei limita centrala.

Kruskal și Wallis au arătat că suma pătratelor lor, cu un factor de ponderare

$$\left(1 - \frac{n_i}{N}\right) \text{ are aproximativ distribuția } \chi^2(k-1) :$$

$$H = \sum \left[ \frac{\bar{R}_i - \frac{N+1}{2}}{\sqrt{\frac{(N+1)(N-n_i)}{12n}}} \right]^2 \left(1 - \frac{n_i}{N}\right) \cong \chi^2(k-1)$$

Dacă apar valori egale, H trebuie să fie împărțit la factorul  $1 - \frac{\sum T}{N^3 - N}$  unde  $T = (k-1)k(k+1)$  este calculat pentru fiecare grup de legături.

Pentru probe mici aproximația nu este prea bună și Kruskal și Wallis au dat tabele pentru  $k=3$  și  $n_i \leq 5$ .

Să aplicăm testul pentru același experiment, considerând două grupuri de observații după prima administrare și un grup de observații după a – 5- a administrare:

Nivelele plasmatice maxime ale ionului  $EDTA^{4-}$  după administrarea i.m. sunt în tabelul1.

$$\bar{R}_1 = \frac{9+3+1+7}{4} = 5, \bar{R}_2 = \frac{2+11+8}{3} = 7 \text{ și } \bar{R}_3 = \frac{4+6+5+10}{4} = 6,25$$

$$H = \sum \left[ \frac{\bar{R}_i - \frac{N+1}{2}}{\sqrt{\frac{(N+1)(N-n_i)}{12n}}} \right]^2 \left(1 - \frac{n_i}{N}\right) = \left[ \frac{5 - \frac{11+1}{2}}{\sqrt{\frac{(11+1)(11-4)}{12*4}}} \right]^2 \left(1 - \frac{4}{11}\right) + \left[ \frac{7 - \frac{11+1}{2}}{\sqrt{\frac{(11+1)(11-3)}{12*3}}} \right]^2 \left(1 - \frac{3}{11}\right) + \left[ \frac{6,25 - \frac{11+1}{2}}{\sqrt{\frac{(11+1)(11-4)}{12*4}}} \right]^2 \left(1 - \frac{4}{11}\right) = \frac{4}{7} \frac{7}{11} + \frac{3}{8} \frac{8}{11} + \frac{6,25 * 4}{7} \frac{7}{11} = \frac{9,5}{11} = 0,86$$

Dat fiindcă  $\chi_{2;0,05}^2 = 0,103$  valoarea obținută pentru test aparține zonei de acceptare, ipoteza ca grupurile sunt selectate din aceeași populație este acceptată.

### Alegerea între testele laplaciene și testele neparametrice

Testele nonparametrice au o putere mai mică decât cele clasice, deoarece înlocuirea valorilor cu rangurile lor semnifică pierderea a o parte din informație. De exemplu am spune ca doi boxeri sunt de aceiasi valoare deoarece fiecare a câștigat câte 5 meciuri din 10 întâlniri dintre ei. În condiția în care în ultima întâlnire A l-a omorât pe B, concluzia trebuie schimbată, deoarece diferența de valoare între ei la ultimul meci a fost cu mult mai mare decât celelate diferențe.

Această pierdere de informație este reală în cazul testelor neparametrice atunci când efectiv variabilele aleatoare sunt repartizate normal și au dispersiile egale. În caz contrar se poate întâmpla ca un test neparametric să fie chiar mai eficient decât cele parametrice.

În altă ordine de idei, aplicarea testelor neparametrice în cazul selecțiilor de volume mari, este foarte laborioasă. Ca urmare, conduita de urmat în alegerea unui tip sau altul de test ar fi după cum urmează:

1. În cazul eșantioanelor mici sunt de preferat testele neparametrice deoarece calculele sunt mai rapide și eficiența este comparabilă cu cea a testelor clasice.
2. Când se știe că selecțiile aparțin la populații repartizate normal și cu dispersii egale, testele clasice sunt mai eficiente.
3. Când nu se cunosc repartițiile variabilelor, alegerea și concluziile se vor face în funcție de alte informații privitoare la experiment.
4. Când se știe că variabilele aleatoare testate nu sunt repartizate normal sau este vorba de variabile care se bazează pe o scală arbitrară (“scoruri”) sau clasificări pe criterii preponderant calitative (de exemplu “ameliorarea” stării subiecților tratați) se apelează la testele neparametrice.

## Regresia liniară

Dacă reprezentarea grafică a două mărimi ce sunt observate simultan sugerează o dependență liniară, ajungem la problema determinării dreptei ce descrie “cel mai bine” această dependență. După cum s-a discutat la capitolul privind extremele funcțiilor de mai multe variabile, o soluție a acestei probleme o constituie “dreapta prin cele mai mici pătrate”, dreapta pentru care suma pătratelor distantelor de la ea la punctele experimentale este minimă. Această soluție consideră punctele ca fiind “exacte”. Problema capătă cu totul altă înfățișare atunci când punctele experimentale sunt considerate valori ale unor variabile aleatoare, devenind o problemă de statistică matematică și analiză numerică în același timp.

Examinăm în continuare cazul cel mai simplu când valorile variabilei  $x$  (care în cele mai multe cazuri corespunde timpului) nu sunt afectate de erori și, pentru fiecare valoare a lui  $x$  corespund un număr de valori  $y$ , determinate într-un singur experiment printr-o metodă afectată de erori întâmplătoare:

$$y_{11}, y_{12}, \dots, y_{1n_1}, \text{ pentru } x_1$$

.

.

$$y_{i1}, y_{i2}, \dots, y_{in_i}, \text{ pentru } x_i, i=1, 2, \dots, k$$

Cazul când pentru orice  $i$  avem  $n_i = 1$  este relativ mai simplu, dar este de subliniat că și în cazul când aceștia sunt diferiți de 1 poate fi tratat în aceeași manieră admitând că între perechile  $(x_i, y_i)$  să fie și perechi cu același  $x_i$ .

Să admitem că pentru un  $x$  fixat, valoarea măsurată  $y$  este o variabilă aleatoare cu următoarea structură:

$$(1) \quad y = \eta + \varepsilon = \alpha x + \beta + \varepsilon$$

distribuită normal cu dispersia  $\sigma^2$  și media  $\eta = \alpha + \beta x$

Problema care ne-o punem este aceea ca, din datele experimentale  $y_i$ , să obținem niște estimări  $a, b$  și  $s^2$  pentru  $\alpha, \beta$  și  $\sigma^2$ , și să determinăm distribuțiile acestor estimări.

Ecuția de regresie estimată este deci:

$$(2) \quad y = a + bx$$

Metoda celor mai mici pătrate dă valorile a și b care minimizează suma pătratelor deviațiilor (erorilor) între valorile observate  $y_i$  și cele prezise de ecuația de regresie (2):

$$(3) SS_E = \sum (Y_i - y_i)^2 = \sum (Y_i - a - bx_i)^2$$

Metoda este în principal datorată lui Gauss. Pentru aflarea parametrilor a și b, nu este necesară ipoteza privind distribuția normală a erorilor, dar aceasta este necesară pentru construirea unor intervale de încredere și pentru testarea unor ipoteze privind aceiași estimatori. Metoda celor mai mici pătrate oferă avantajul că estimatorii pe care îi dă sunt deplasați și au o dispersie minimă în clasa estimatorilor nedepasați.

Valorile lui a și b care minimizează suma pătratelor erorilor sunt soluțiile sistemului

$$\frac{\partial SS}{\partial a} = 0 \text{ și } \frac{\partial SS}{\partial b} = 0$$

$$\frac{\partial SS}{\partial a} = -2 \sum (Y_i - a - bx_i) = 0 \quad \text{și} \quad \frac{\partial SS}{\partial b} = -2 \sum (Y_i - a - bx_i)x_i = 0$$

ceea ce este echivalent cu

$$(4) na + b \sum x_i = \sum Y_i \text{ și}$$

$$a \sum x_i + b \sum x_i^2 = \sum x_i Y_i$$

Rezolvând sistemul prin regula lui Cramer se obțin ca estimatori pentru  $\alpha$  și  $\beta$ :

$$a = \frac{\sum Y_i \sum x_i^2 - \sum x_i \sum x_i Y_i}{n \sum x_i^2 - (\sum x_i)^2} \text{ și } b = \frac{n \sum x_i Y_i - \sum x_i \sum Y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

Numărătorul expresiei lui b poate fi scris și în forma

$$n \sum x_i y_i - \sum x_i \sum y_i = n \left( \sum x_i y_i - \sum \frac{x_i}{n} \sum y_i \right) = n \sum (x_i - \bar{X}) y_i$$

Deoarece  $\sum (x_i - \bar{X}) = 0$  și  $\bar{Y} \sum (x_i - \bar{X}) = 0$ , mai putem scrie

$$\sum (x_i - \bar{X}) y_i = \sum (x_i - \bar{X}) y_i - \bar{Y} \sum (x_i - \bar{X}) = \sum (x_i - \bar{X}) (y_i - \bar{Y})$$

Similar, după cum se poate ușor verifica, avem:

$$n \sum x_i^2 - (\sum x_i)^2 = n \sum (x_i - \bar{X})^2$$

În consecință, o formă alternativă pentru b este  $b = \frac{\sum (x_i - \bar{X})(y_i - \bar{Y})}{\sum (x_i - \bar{X})^2}$ .

Putem verifica ușor că b este un estimator nedeplasat pentru  $\beta$ . Presupunem valoarea așteptată  $y_i$  dată de ecuația  $\alpha + \beta x_i$ , pentru un  $x = x_i$ . Atunci:

$$\begin{aligned} E(b) &= \frac{\sum (x_i - \bar{X})E(y_i)}{\sum (x_i - \bar{X})^2} = \frac{\sum (x_i - \bar{X})(\alpha + \beta x_i)}{\sum (x_i - \bar{X})^2} = \alpha \frac{\sum (x_i - \bar{X})}{\sum (x_i - \bar{X})^2} + \beta \frac{\sum (x_i - \bar{X})x_i}{\sum (x_i - \bar{X})^2} = \\ &= 0 + \beta \frac{\sum x_i^2 - \frac{1}{n}(\sum x_i)^2}{\sum (x_i - \bar{X})^2} = \beta \end{aligned}$$

Dispersiile lui a și b pot fi obținute direct, deoarece sunt funcții liniare de  $y_i$ , care valori sunt presupuse independente și distribuite normal, cu dispersia  $\sigma^2$ :

$$D(b) = D \left[ \frac{\sum (x_i - \bar{X})y_i}{\sum (x_i - \bar{X})^2} \right] = \frac{\sum (x_i - \bar{X})^2 D(y_i)}{\left( \sum (x_i - \bar{X})^2 \right)^2} = \frac{\sigma^2}{\sum (x_i - \bar{X})^2}$$

Din prima ecuație a sistemului (4) avem:  $a = \bar{Y} - b\bar{X}$ .

$$\begin{aligned} D(a) &= D \left( \frac{\sum y_i}{n} \right) + \bar{X}^2 D(b) = \frac{1}{n^2} \sum D(y_i) + \bar{X}^2 \frac{\sigma^2}{\sum (x_i - \bar{X})^2} = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum (x_i - \bar{X})^2} \right) = \\ &= \sigma^2 \frac{\sum x_i^2 - \frac{(\sum x_i)^2}{n} + \frac{(\sum x_i)^2}{n^2}}{n \sum (x_i - \bar{X})^2} = \frac{\sum x_i^2}{n} \frac{\sigma^2}{\sum (x_i - \bar{X})^2} = \frac{\sum x_i^2}{n} D(b) \end{aligned}$$

$$\text{Deci, } S_a = \sqrt{\frac{\sum x_i^2}{n} S_b^2}$$

### Estimații și ipoteze asupra coeficientului b

Coeficientul b are o importanță deosebită și prin aceea că el reprezintă o măsură a corelării între x și y.

1. Coeficientul b este, după cum s-a arătat, repartizat normal cu media  $\beta$  și dispersia

$$\frac{\sigma^2}{\sum (x_i - \bar{X})^2}$$



2. Dacă  $y_i$  sunt punctele experimentale, iar  $Y_i$  estimările lor teoretice,  $Y_i = a + bx_i$ , suma pătratelor erorilor va fi  $SS_E = \sum (y_i - Y_i)^2$ . Vom arăta că:

$$E\left(\frac{SS_E}{n-2}\right) = \sigma^2$$

Pentru a demonstra aceasta relație plecăm de la definiția sumei erorilor

$$SS_E = \sum [y_i - (a + bx_i)]^2 = \sum [(y_i - \bar{Y}) + (\bar{Y} - a - bx_i)]^2 = \sum [(y_i - Y) + (a + b\bar{X} - a - bx_i)]^2 =$$

$$= \sum [(y_i - \bar{Y}) - b(x_i - \bar{X})]^2 = \sum (y_i - \bar{Y})^2 - 2b \sum (x_i - \bar{X})(y_i - \bar{Y}) + b^2 \sum (x_i - \bar{X})^2 =$$

Dar  $b = \frac{\sum (x_i - \bar{X})(y_i - \bar{Y})}{\sum (x_i - \bar{X})^2}$  și putem înlocui  $\sum (x_i - \bar{X})(y_i - \bar{Y}) = b \sum (x_i - \bar{X})^2$ .

Deci avem  $SS_E = \sum (y_i - \bar{Y})^2 - b^2 \sum (x_i - \bar{X})^2 = A - B$

Calculăm separat  $E(A)$  și  $E(B)$ .

$$E(A) = E\left[\sum (y_i - \bar{Y})^2\right] = E\left(\sum y_i^2 - n\bar{Y}^2\right) = E\left(\sum y_i^2\right) - nE(\bar{Y}^2)$$

În continuare, folosind identitatea  $D(Y) = E(Y^2) + (E(Y))^2$  și faptul că  $E(\bar{Y}) = \alpha + \beta\bar{X}$  și

$$D(\bar{Y}) = \frac{\sigma^2}{n} \text{ obținem}$$

$$E(A) = \sum [(\alpha + \beta x_i)^2 + \sigma^2] - n\left[(\alpha + \beta\bar{X})^2 + \frac{\sigma^2}{n}\right] = n\sigma^2 - \frac{n\sigma^2}{n} + \beta^2 \sum (x_i - \bar{X})^2 =$$

$$= (n-1)\sigma^2 + \beta^2 \sum (x_i - \bar{X})^2$$

Mai departe,

$$E(B) = \sum (x_i - \bar{X})^2 E(b^2) = \sum (x_i - \bar{X})^2 [D(b) + (E(b))^2] = \sum (x_i - \bar{X})^2 \left(\frac{\sigma^2}{\sum (x_i - \bar{X})^2} + \beta^2\right)$$

și deci,

$$E(SS_E) = (n-1)\sigma^2 + \beta^2 \sum (x_i - \bar{X})^2 - \beta^2 \sum (x_i - \bar{X})^2 - \sigma^2 = (n-2)\sigma^2$$

3. Variabila aleatoare  $\frac{SS_E}{\sigma^2}$  este repartizată  $\chi^2(n-2)$ .

Pe baza acestor trei proprietăți putem estima intervalele de încredere pentru  $\beta$  și verifica ipoteze asupra valorilor sale.

### a) Cazul dispersiilor cunoscute

În cazul în care se cunoaște dispersia erorilor de măsurare  $D(\varepsilon_i) = D(y_i) = \sigma^2$  se folosește faptul că variabila aleatoare  $z = \frac{b - \beta}{\sqrt{D(b)}} = \frac{b - \beta}{\left[ \frac{\sigma^2}{\sum (x_i - \bar{X})^2} \right]^{\frac{1}{2}}}$  este repartizată  $N(0,1)$ .

### b) Cazul dispersiilor necunoscute

În acest caz se înlocuiește dispersia lui  $b$ :  $\frac{\sigma^2}{\sum (x_i - \bar{X})^2}$  cu estimatorul numit “dispersia

de selecție”:  $S_b = \frac{\sum (y_i - Y_i)^2}{\sum (x_i - \bar{X})^2} = \frac{SS_E}{\sum (x_i - \bar{X})^2}$ .

Variabila aleatoare

$$T = \frac{b - \beta}{\left[ \frac{SS_E}{(n-2)\sum (x_i - \bar{X})^2} \right]^{\frac{1}{2}}} = \frac{\frac{b - \beta}{\sigma_b}}{\left[ \frac{SS_E}{(n-2)\sigma^2} \right]^{\frac{1}{2}}} = \frac{Z}{\sqrt{\frac{\chi^2_{n-2}}{n-2}}}$$

este repartizată Student cu  $n-2$  grade de libertate.

Ca urmare putem determina intervalele în care se află  $\beta$  cu diverse probabilități sau verifica ipoteze privind valoarea lui, exact cum este utilizat testul  $t$  pentru testarea ipotezei privind media necunoscută.

Intervalul de încredere pentru  $\beta$  este  $b - t_{n-2, 1-\frac{\alpha}{2}} S_b < \beta < b + t_{n-2, 1-\frac{\alpha}{2}} S_b$

### Estimarea dispersiei drepte de regresie

Considerăm un punct  $x_0$  fixat și punctul corespunzător lui:  $y_0$ , pe dreapta de regresie  $y$

$$y = \alpha + \beta x + \varepsilon = a + bx$$

$$\bar{Y} = a + b\bar{X}$$

$$y_0 = a + bx_0 = \bar{Y} - b\bar{X} + bx_0$$

$y_0 = \bar{Y} + b(x_0 - \bar{X})$  estimatia lui  $y_0$  este o variabilă aleatoare distribuită normal.

Avem  $E(Y_0) = \bar{Y}_0 = \alpha + \beta x_0$  și

$$D(y_0) = \sigma_{y_0}^2 = \sigma_y^2 + \sigma_b^2(x_0 - \bar{X})^2 = \frac{\sigma^2}{n} + \frac{\sigma^2}{\sum(x_i - \bar{X})^2}(x_0 - \bar{X})^2$$

Estimând valoarea lui  $\sigma^2$  prin  $s = \frac{SS_E}{n-2}$  avem  $s_{y_0}^2 = s^2 \left[ \frac{1}{n} + \frac{(x_0 - \bar{X})^2}{\sum(x_i - \bar{X})^2} \right]$ .

Variabila aleatoare  $T = \frac{y_0 - (\alpha + \beta x_0)}{s_{y_0}}$  este repartizată Student cu  $n-2$  grade de libertate și permite calculul intervalelor de încredere pentru  $\alpha + \beta x_0$ .

Dispersia  $s_{y_0}$  depinde de distanța între  $x_0$  și  $\bar{X}$ , valoarea sa fiind minimă atunci când  $x_0 = \bar{X}$ . În acest caz,  $y_0 = \bar{Y}$  și  $s_{y_0} = s_y$ .

Facem observația că dispersia determinată în punctul  $y_0$  este dispersia datorată regresiei. Valorile experimentale nu sunt însă valori ale regresiei  $y_0 = \bar{Y} + b(x_0 - \bar{x})$ , estimate drepte de regresie. În acest caz, valoarea individuală determinată diferă față de valoarea  $Y_0$  printr-o eroare  $\varepsilon$ , a cărei dispersie este egală cu  $\sigma^2$ , variabilitatea datelor individuale față de valorile corespunzătoare regresie  $Y$ .

Ca urmare, valorile individuale vor avea dispersia  $\sigma_{y_0}^2 = \sigma^2 + \frac{\sigma^2}{n} + \sigma^2 \frac{(x_0 - \bar{X})^2}{\sum(x_i - \bar{X})^2}$

ceea ce, pentru valorile de selecție devine  $s_{y_0}^2 = s^2 \left[ 1 + \frac{1}{n} + \frac{(x_0 - \bar{X})^2}{\sum(x_i - \bar{X})^2} \right]$ .

<sup>1</sup>F.Wilcoxon: *Individual comparisons by ranking methods, Biometrics Bul.*,180-83,1947

<sup>2</sup>W.H.Kruskal, W.Allen Wallis: *Use of ranks in one-criterion analysis of variance, J. Am. Stat. Assoc.*,47,583-621,1952

<sup>3</sup>W.H.Kruskal, W.A.Wallis; *use of ranks in the one - criterion analysis of variance, J.Am.Stat.Assoc.*,47,583-621,1952