

CURSUL – II –

STATISTICA

INTERVALE DE INCREDERE

VERIFICAREA IPOTEZELOR

STATISTICĂ MATEMATICĂ ȘI BIostatistică

Statistica matematică este principala aplicație a teoriei probabilităților. Metodele statistice constau, în esență, în elaborarea unor concluzii plauzibile privitoare la colectivități mari de fenomene, pe baza cunoașterii unui număr restrâns dintre acestea și extrapolării rezultatelor.

Legile care stau la baza statisticii și care permit aceste generalizări sunt teorema limită centrală și legea numerelor mari.

Într-o exprimare intuitivă, avem rezultatul că, dacă o variabilă aleatoare ξ este suma unui număr mare de variabile aleatoare independente, fiecare variabilă aleatoare având o pondere mică în sumă, atunci funcția de repartiție a variabilei aleatoare ξ este foarte apropiată de o funcție de repartiție normală.

Exprimat mai riguros și mai general, avem următoarea teoremă:

Teorema limită centrală (A.M. Leapunov)

Fie $\xi_1, \xi_2, \dots, \xi_n$ variabile aleatoare independente.

Fie $M(\xi_k) = a_k, D(\xi_k) = \sigma_k^2$ și $\rho_k^3 = M(|\xi_k - m_k|)^3$ când $k = \overline{1, n}$

Notăm $\sigma_{(n)}^2 = \sum_1^n \sigma_i^2$, $\rho_{(n)}^3 = \sum_1^n \rho_i^3$

Dacă $\lim_{n \rightarrow \infty} \frac{\rho_{(n)}}{\sigma_{(n)}} = 0$, atunci funcția de repartiție a variabilei

$$\frac{\xi_1 + \xi_2 + \dots + \xi_n - (a_1 + a_2 + \dots + a_n)}{\sigma_{(n)}}$$

tinde, când $n \rightarrow \infty$, către funcția $\Phi(x)$ a lui Laplace.

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$$

Teorema limită centrală este teorema fundamentală a teoriei erorilor. Laplace, Gauss și alți matematicieni, studiind repartiția erorilor, au ajuns la concluzia că funcția de repartiție normală poate fi luată drept model teoretic pentru cercetarea probabilistică a aproape tuturor fenomenelor naturii.

Teorema lui Cebâșev

Dacă $\xi_1, \xi_2, \dots, \xi_n$ sunt variabile aleatoare (discrete sau continue) independente ale căror dispersii sunt mai mici decât o constantă C, atunci oricare ar fi numărul pozitiv ε , probabilitatea inegalității

$$\left| \frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n} \right| \langle \mathcal{E} \rangle$$

tinde către 1, atunci când numărul variabilelor aleatoare tinde către infinit.

Demonstrație:

Să considerăm variabila aleatoare $\bar{\zeta} = \frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n}$. Având în vedere liniaritatea

operatorului de calcul a mediei avem $M(\bar{\zeta}) = \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}$.

Aplicând inegalitatea lui Cebâșev variabilei aleatoare $\bar{\zeta}$ se obține:

$$P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| \langle \mathcal{E} \rangle\right) \geq 1 - \frac{D\left(\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n}\right)}{\mathcal{E}^2}$$

Mai departe, din proprietățile operatorului D

$$D\left(\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n}\right) = \frac{D(\zeta_1) + D(\zeta_2) + \dots + D(\zeta_n)}{n^2} \leq \frac{C + C + \dots + C}{n^2} = \frac{nC}{n^2} = \frac{C}{n}$$

Deci

$$P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| \langle \mathcal{E} \rangle\right) \geq 1 - \frac{C}{n\mathcal{E}^2}$$

Trecând la limita pentru $n \rightarrow \infty$ obținem

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| \langle \mathcal{E} \rangle\right) \geq 1$$

și cum probabilitatea nu poate depăși 1,

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| \langle \mathcal{E} \rangle\right) = 1$$

Cel mai frecvent, în practică, variabilele aleatoare ζ_i au aceeași medie μ și concluzia teoremei devine

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \mu\right| \langle \mathcal{E} \rangle\right) = 1$$

În esență, teorema lui Cebâșev stabilește că, deși variabilele aleatoare independente pot lua valori îndepărtate față de mediile lor, media aritmetică a unui număr suficient de mare de astfel de

variabile aleatoare ia cel mai probabil valori apropiate de un număr constant $\frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}$ (sau μ atunci când mediile variabilelor sunt egale între ele).

Ca urmare, între comportarea fiecărei variabile aleatoare și comportarea mediilor lor există diferență esențială. Putem spune foarte precis ce valoare va lua media aritmetică a acestor variabile aleatoare. Explicația acestui fapt constă în aceea că abaterile diverselor variabile aleatoare sunt de semne diferite și, ca urmare, se compensează între ele.

TEORIA SELECȚIEI

Populații și selecții. Inferența statistică

În practică avem adesea nevoie să facem judecăți asupra unor mari colecții de rezultate posibile experimental ori a altor cantități, dar nu putem sau este extrem de scump, să examinăm toate aceste date. În astfel de cazuri, în loc să examinăm întregul set de date pe care îl numim în cele ce urmează *populație*, tragem concluziile după examinarea a o parte din ele, alese la întâmplare, parte pe care o numim *selecție*.

Procedul de obținere a probelor este numit tot selecție, iar procedul de extrapolare a concluziilor la întreaga populație este cunoscut ca *inferența statistică*.

Vom considera că o caracteristică dată a populației este o variabilă aleatoare pe un câmp de probabilitate (Ω, K, P) în care elementele lui Ω sunt chiar elementele populației, iar P este o probabilitate cunoscută sau nu.

Enumerarea valorilor observate ale caracteristicii urmărite și a frecvențelor lor relative definește *repartiția statistică a selecției*.

Teorema lui Leapunov, numită și teorema fundamentală a statisticii matematice, care justifică utilizarea metodei selecției stabilește că funcția de repartiție statistică a caracteristicilor selecțiilor tinde la funcția teoretică de repartiție a caracteristicii studiate când volmul selecției tinde la ∞ .

Exemplu 1

Putem dori să tragem concluzii despre evoluția rezistenței unei tulpini de germeni patogeni la un medicament dat și, în acest scop, examinăm rezultatele antibiogramelor făcute într-un eșantion de spitale într-o perioadă recentă (luniile de iarnă), comparată cu aceeași perioadă a anului precedent. Deși rezultatele obținute se referă la spitale și mai precis numai la o parte din ele, concluziile le extindem la scara întregii populații.

Exemplu 2

Rezultatele privind absorbția unui medicament după administrarea orală prin determinarea nivelurilor din plasma ale medicamentului la un lot de voluntari sănătoși le considerăm ca rezultate probabile pentru întreaga populație ce include și potențiali pacienți.

Populația poate fi infinită sau finită, în ultimul caz, numărul indivizilor populației – N - se mai numește și *volumul populației*. În mod similar, numărul de indivizi sau valori din cadrul unei probe este denumit *volumul probei* sau *volumul eșantionului*.

Valabilitatea concluziilor despre populație depinde de “reprezentativitatea” probei. Pentru populații finite aceasta înseamnă că fiecare membru al populației are aceeași șansă să fie selectat, când spunem că selecția este o selecție la întâmplare sau “selecție aleatoare”. Desigur că selecția unor voluntari sănătoși pentru determinarea parametrilor farmacocinetici ai unui medicament nu este din acest punct de vedere o selecție reprezentativă. În cazurile în care avem motive să credem că patologia căreia se adresează medicamentul nu afectează funcțiile metabolice și de excreție, această aproximare este acceptată pentru motivul că o selecție corectă ar implica loturi mult mai mari cu cheltuieli și timp de lucru mult crescute.

În practică, în studiile de bioechivalență, pentru reducerea volumului loturilor pe care se fac testările, se administrează amândouă medicamentele la toți membri lotului, în două perioade diferite. Fiecare component al lotului primește unul din medicamente în prima perioadă și celălalt în a doua perioadă.

Deoarece perioada de administrare poate influența și ea rezultatul experimentului, alegerea indivizilor care vor primi în prima perioadă primul medicament se face în mod aleator. În cazul când sunt mai multe perioade, de exemplu I-IV, și mai multe medicamente A, B, C, D se alcătuiește un tabel de felul

I	II	III	IV
A	B	D	C
B	C	A	D
C	D	B	A
D	A	C	B

așa zisul pătrat “latin”, unde observăm că fiecare literă apare o singură dată în fiecare linie și în fiecare coloană. Se numește pătrat latin deoarece, cum se va arata mai departe, în cazul în care mai intervine și o altă variabilă, de exemplu doza din fiecare medicament, se folosesc și litere grecești, alcătuiindu-se pătrate “greco-latine”.

Deasemenea, studiile de bioechivalență se fac tot pe voluntari sănătoși, pornind de la ipoteza că modificările de biodisponibilitate asociate stărilor patologice sunt aceleași pentru cele două medicamente testate, ceea ce, evident, este numai în parte adevărat.

În toate experimentele biologice, planificarea experimentului trebuie făcută în așa fel încât diferențele în tratament să nu coincidă cu diferențe în vârstă, sex, sau alți parametri. Dacă, de exemplu, femeile din lot primesc primul medicament și bărbații al doilea, se spune că diferențele de sex sunt “confundate” cu diferențele de tratament. În acest caz nu se poate spune dacă diferențele obținute se datorează tratamentului sau diferenței de sex.

Parametrii de selecție ai unei variabile aleatoare :

Dacă printr-un procedeu oarecare cuantificăm răspunsul culturilor microbiene la antibioticele din exemplul 1, sau dacă luăm în considerație concentrațiile de medicament în sânge, din al doilea exemplu, și probabilitățile ca valorile să aparțină unor intervale diferite, obținem o variabilă aleatoare X asociată cu rezultatul experimentului corespunzător.

Parametrii acestei variabile aleatoare sunt denumiți, prin abuz de limbaj, “parametri ai populației”.

Dacă în exemplul al doilea X_i este concentrația de medicament în sângele bolnavului i , la o oră de la administrare, la primul voluntar putem obține o valoare x_1 , pentru al doilea voluntar o valoare x_2 , etc. În acest fel găsim valorile x_1, x_2, \dots, x_n ale variabilelor aleatoare independente X_1, X_2, \dots, X_n .

Media de selecție este o variabilă aleatoare: $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$

Dacă distribuția lui X este normală $-N(\mu, \sigma)$, aceiași pentru fiecare i , datorită linearității operatorului E care definește media, obținem $M(\bar{X}) = \mu_{\bar{X}} = \mu$ adică valoarea pentru *media mediei de selecție este media populației*.

Dacă la datele experimentale se adaugă o constantă, $x'_i = x_i + a$, media de selecție crește cu aceeași constantă: $\bar{W} = \frac{\sum_1^n (X_i + a)}{n} = \bar{X} + a$

Similar, dacă fiecare valoare se înmulțește cu o constanta $Z_i = kX_i$, media de selecție \bar{Z} se

înmulțește cu aceeași constantă: $\bar{Z} = \frac{\sum_1^n kX_i}{n} = k\bar{X}$

Dispersia de selecție

Ca o măsură a abaterii datelor față de media de selecție, se introduce noțiunea de dispersie de

selecție $s_x^2 = \frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2$.

În aplicațiile practice, pentru reducerea numărului de calcule, formula se aduce la o altă formă și anume:

$$s_x^2 = \frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_1^n x_i^2 - 2\bar{X} \sum_1^n x_i + n\bar{X}^2 \right) = \frac{1}{n-1} \left(\sum_1^n x_i^2 - 2n\bar{X}^2 + n\bar{X}^2 \right) =$$

$$\frac{1}{n-1} \left(\sum_1^n x_i^2 - n\bar{X}^2 \right) = \frac{1}{n-1} \left(\sum_1^n x_i^2 - \frac{(\sum_1^n x_i)^2}{n} \right) \quad \text{Dacă}$$

$z_i = kx_i + a \Rightarrow s_z^2 = k^2 s_x^2$. Într-adevăr

$$s_z^2 = \frac{1}{n-1} \sum_1^n (z_i - \bar{Z})^2 = \frac{1}{n-1} \sum_1^n (kx_i + a - k\bar{X} - a)^2 = k^2 s_x^2$$

s_x se numește abaterea standard de selecție sau deviație standard, când nu este pericol de confuzie privind variabila aleatoare la care se referă folosindu-se și notația s_d .

Se mai utilizează în practică și noțiunile:

- Dispersia "populației" = $\frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2$ și respectiv
- Deviația standard a "populației", precum și "abaterea standard a mediei" (prescurtarea SEM – standard error of mean) definită prin raportul $SEM = \frac{s_x}{\sqrt{n}}$
- precum și coeficientul de variație $v = \frac{s_x}{\bar{X}} * 100$.

Covarianța de selecție

Covarianța de selecție se definește prin formula $s_{xy} = \frac{1}{n-1} \sum_1^n (x_i - \bar{X})(y_i - \bar{Y})$

Se observă că aceasta se mai poate scrie și sub altă formă, mai utilă în sensul simplificărilor de calcul în anumite aplicații.

$$s_{XY} = \frac{1}{n-1} \left(\sum_1^n x_i y_i - \bar{X} \sum_1^n y_i - \bar{Y} \sum_1^n x_i + n \bar{X} \bar{Y} \right) = \frac{1}{n-1} \left(\sum_1^n x_i y_i - n \bar{X} \bar{Y} - n \bar{X} \bar{Y} + n \bar{X} \bar{Y} \right) =$$

$$\frac{1}{n-1} \left(\sum_1^n x_i y_i - n \bar{X} \bar{Y} \right) = \frac{1}{n-1} \left(\sum_1^n x_i y_i - \frac{\sum_1^n x_i \sum_1^n y_i}{n} \right)$$

Coeficientul de corelație de selecție

Coeficientul de corelație de selecție se definește prin formula

$$\rho(x, y) = \frac{s_{XY}}{s_X s_Y} = \frac{\frac{1}{n-1} \sum_1^n (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2} \sqrt{\frac{1}{n-1} \sum_1^n (y_i - \bar{Y})^2}} = \frac{\sum_1^n (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\sum_1^n (x_i - \bar{X})^2} \sqrt{\sum_1^n (y_i - \bar{Y})^2}}$$

Proprietăți ale caracteristicilor de selecție

Considerăm în continuare o selecție de volum n dintr-o populație cu media μ și dispersia σ^2

Propoziție

Media mediei de selecție este egală cu media populației. $M(\bar{X}) = \mu$

Demonstrație:

$$M(\bar{X}) = \frac{M(\sum x_i)}{n} = \frac{\sum M(x_i)}{n} = \frac{n\mu}{n} = \mu$$

Propoziție

Media dispersiei de selecție este egală cu dispersia populației $M(s_X^2) = \sigma^2$

Demonstrație:

$$M(s_X^2) = M\left(\frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2\right) = \frac{1}{n-1} M\left(\sum_1^n x_i^2 - 2\bar{X} \sum_1^n x_i + \sum_1^n \bar{X}^2\right) =$$

$$\frac{1}{n-1} M\left(\sum_1^n x_i^2 - 2n\bar{X}^2 + n\bar{X}^2\right) = -\frac{1}{n-1} M\left(\sum_1^n x_i^2 - n\bar{X}^2\right) = \frac{1}{n-1} M\left(\sum_1^n x_i^2 - \frac{(\sum_1^n x_i)^2}{n}\right)$$

Dar, mai departe

$$M\left(\sum_1^n x_i^2\right) = n(\sigma^2 + \mu^2)$$

$$M\left(\sum_1^n x_i\right)^2 = M\left(\sum_1^n x_i^2 + 2\sum_{i \neq j}^n x_i x_j\right) = \sum_1^n M(x_i^2) + 2\frac{n(n-1)}{2}M(x_i)M(x_j) = \\ = n(\sigma^2 + \mu^2) + n(n-1)\mu^2 = n\sigma^2 + n^2\mu^2$$

și înlocuind în expresia lui $M(s_x^2)$ obținem $M(s_x^2) = \frac{n(\sigma^2 + \mu^2) - \sigma^2 - n\mu^2}{n-1} = \sigma^2$

Propoziție

Variabila aleatoare $Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$ tinde, când $n \rightarrow \infty$ către o variabilă $N(0,1)$

Aceasta este o consecință a teoremei limită centrală și este aplicabilă atât variabilelor continue cât și celor discrete.

Într-adevăr aplicând teorema lui Leapunov pentru variabilele aleatoare x_1, x_2, \dots, x_n obținem că:

$$\frac{x_1 + x_2 + \dots + x_n - (\mu_1 + \mu_2 + \dots + \mu_n)}{\sqrt{\sigma^2 + \sigma^2 + \dots + \sigma^2}} = \frac{n\bar{X} - n\mu}{\sqrt{n\sigma^2}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

tinde către o variabilă aleatoare $N(0,1)$.

Propoziție

Dacă variabila aleatoare X este normal distribuită, atunci variabila aleatoare $V = (n-1)\frac{s_x^2}{\sigma^2}$

este repartizată $\chi^2(n-1)$

Demonstrație:

$$V = (n-1)\frac{s_x^2}{\sigma^2} = \frac{\sum_1^n (x_i - \bar{X})^2}{\sigma^2} = \frac{\sum_1^n [(x_i - \mu) - (\bar{X} - \mu)]^2}{\sigma^2} = \\ = \frac{\sum_1^n (x_i - \mu)^2 - 2\sum_1^n (x_i - \mu)(\bar{X} - \mu) + \sum_1^n (\bar{X} - \mu)^2}{\sigma^2} = \frac{\sum_1^n (x_i - \mu)^2 - 2(n\bar{X} - n\mu)(\bar{X} - \mu) + n(\bar{X} - \mu)^2}{\sigma^2} = \\ = \frac{\sum_1^n (x_i - \mu) - n(\bar{X} - \mu)}{\sigma^2} = \sum_1^n \left(\frac{x_i - \mu}{\sigma} \right)^2 - \left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \right)^2$$

Dar variabila aleatoare $\frac{x_i - \mu}{\sigma}$ este repartizată $N(0,1)$ deoarece $M\left(\frac{x_i - \mu}{\sigma}\right) = \frac{M(x_i) - \mu}{\sigma}$ și

$D\left(\frac{x_i - \mu}{\sigma}\right) = \frac{D(x_i)}{\sigma^2} = 1$, iar $\frac{(\bar{X} - \mu)}{\frac{\sigma}{\sqrt{n}}}$ este repartizată tot $N(0,1)$ în conformitate cu teorema limită

centrală.

Deci, V este o sumă de $n-1$ pătrate de variabile de tip $N(0,1)$.

Propoziție

Dacă x_1, x_2, \dots, x_n este o selecție dintr-o populație normal distribuită, atunci variabila aleatoare

$T = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$ este repartizată Student cu n grade de libertate.

Demonstrație:

$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} = \frac{\frac{\bar{X} - \mu}{\sigma}}{\frac{s}{\sigma}} = \frac{\frac{\bar{X} - \mu}{\sigma}}{\sqrt{\frac{\sum_1^n (x_i - \bar{X})^2}{(n-1)\sigma^2}}} = \frac{Z}{\sqrt{\frac{V}{n-1}}}$$

unde $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ este repartizată $N(0,1)$, iar $V = \frac{\sum_1^n (x_i - \bar{X})^2}{\sigma^2}$ este repartizată $\chi^2(n-1)$.

Deci, T este repartizată Student cu $n-1$ grade de libertate.

Propoziție

Date fiind două selecții aleatoare independente $x_{11}, x_{12}, \dots, x_{1n_1}$ și $x_{21}, x_{22}, \dots, x_{2n_2}$ din populații

normal distribuite $N(\mu_1, \sigma_1)$ și $N(\mu_2, \sigma_2)$, variabila aleatoare $F = \frac{\frac{S_1^2}{\sigma_1^2}}{\frac{S_2^2}{\sigma_2^2}}$ este repartizată Fisher –

Snedecor $F(n_1 - 1, n_2 - 1)$

Demonstrație:

Avem într-adevar $F = \frac{S_1^2}{S_2^2} = \frac{\frac{\sum_1^{n_1} (x_{1i} - \bar{X}_1)^2}{(n_1 - 1)\sigma_1^2}}{\frac{\sum_1^{n_2} (x_{2i} - \bar{X}_2)^2}{(n_2 - 1)\sigma_2^2}}$ iar numărătorul și numitorul sunt repartizate, conform

propozitiei 2.3.5.4., respectiv $\frac{\chi^2(n_1 - 1)}{n_1 - 1}$ și $\frac{\chi^2(n_2 - 1)}{n_2 - 1}$.

Estimații

Teoria estimației urmărește evaluarea parametrilor unei repartiții în general cunoscute. Valorile numerice obținute se numesc *estimații* sau *estimatori*. Se obțin estimații punctuale în cazul în care se folosesc datele selecției pentru a obține valorile parametrilor și estimații ale intervalelor de încredere în cazul în care se determină un interval în care se află, cu o anumită probabilitate valoarea estimată.

Un estimator al parametrului θ se va nota cu $\hat{\theta}$. O estimație este nedeplasată dacă $M(\hat{\theta}) = \theta$, adică media estimației este egală chiar cu valoarea teoretică a parametrului estimat.

Conform proprietății 2.3.5.1, $M(\bar{X}) = \mu$ adică media de selecție este un estimator nedeplasat al mediei, iar conform proprietății 2.3.5.2., $M(s^2) = \sigma^2$ adică dispersia de selecție este un estimator nedeplasat al dispersiei.

Problema estimării intervalelor se reduce la găsirea unui interval de încredere (θ_L, θ_U) cu un coeficient de încredere $1 - \alpha$ astfel încât $P(\theta_L < \theta < \theta_U) = 1 - \alpha$.

Este de dorit ca $1 - \alpha$ să fie cât mai mare (de obicei este cuprins între 0,9 și 0,99) iar intervalul (θ_L, θ_U) să fie cât mai mic. În stabilirea intervalelor se utilizează caracteristicile numerice cuantile. Se numesc *cuantile de ordin β* valoarea x_β a variabilei aleatoare x pentru care $F(x_\beta) = P(x < x_\beta) = \beta$ adică valoarea variabilei aleatoare care are la stânga ei aria β sub curba densității de probabilitate.

Evident:

$$P\left(x < x_{\frac{\alpha}{2}}\right) = \frac{\alpha}{2} \qquad P\left(x < x_{1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2}$$

$$P\left(x_{\frac{\alpha}{2}} < x < x_{1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Pentru a estima un interval se alege $1 - \alpha$, se citesc din tabelele cuantilele, de exemplu $z_{1-\frac{\alpha}{2}}$ și $z_{\frac{\alpha}{2}}$ și se precizează intervalul. În prealabil, în funcție de mărimea pentru care se caută intervalul se precizează cu care din repartițiile cunoscute trebuie lucrat.

Estimarea intervalelor de încredere pentru medii

Cazul când se cunoaște dispersia.

Se consideră o populație repartizată normal $N(\mu, \sigma^2)$. Dacă se cunoaște dispersia se poate folosi faptul că $z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ este repartizată $N(0,1)$. Se notează cu z_{α} cuantila de ordinul α pentru repartiția

$N(0,1)$. Evident

$$P\left(z_{\frac{\alpha}{2}} < z < z_{1-\frac{\alpha}{2}}\right) = F\left(z_{1-\frac{\alpha}{2}}\right) - F\left(z_{\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Așadar intervalul $\left(z_{\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}}\right)$ este un interval de estimare cu coeficientul de încredere $1 - \alpha$. Din

anumite puncte de vedere este recomandabil să se utilizeze acele intervale care lasă atât la dreapta cât și la stânga lor aceeași arie, egală cu $\frac{\alpha}{2}$.

Deoarece repartiția $N(0,1)$ este simetrică față de axa Oy avem relația $z_{\frac{\alpha}{2}} = -z_{1-\frac{\alpha}{2}}$

Din relațiile

$$\begin{aligned} -z_{1-\frac{\alpha}{2}} < z < z_{1-\frac{\alpha}{2}} &\Rightarrow -z_{1-\frac{\alpha}{2}} < \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} < z_{1-\frac{\alpha}{2}} \Rightarrow -z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} < \bar{x} - \mu < z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} \Rightarrow \\ -\bar{x} - z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} < -\mu < -\bar{x} + z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} \end{aligned}$$

rezultă

$$\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

Așadar intervalul căutat este

$$(\theta_L, \theta_U) = \left(\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right)$$

Mărima $E = z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ poartă numele de *eroare* și servește la calculul numărului de experiențe

$$n = \left(\frac{z_{1-\frac{\alpha}{2}}}{E} \right)^2 \text{ atunci când este impusă eroarea și se alege un coeficient } 1 - \alpha$$

Metoda descrisă mai poate fi aplicată și în cazul în care x nu este repartizată normal deoarece z este repartizată $N(0,1)$ indiferent de repartiția variabilelor x_1, x_2, \dots, x_n (teorema limită centrală).

Cazul când dispersia este necunoscută

Dacă nu se cunoaște dispersia în estimarea intervalelor se utilizează dispersia de selecție care este un estimator nedeplasat al dispersiei deoarece $E(s^2) = \sigma^2$

Se consideră x_1, x_2, \dots, x_n o selecție dintr-o populație de tipul $N(\mu, \sigma^2)$.

Conform celor arătate anterior mărimea $T = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$ este repartizată $T(n-1)$ și, ca urmare

$$P\left(t_{n-1, \frac{\alpha}{2}} < T < t_{n-1, 1-\frac{\alpha}{2}}\right) = F\left(t_{n-1, 1-\frac{\alpha}{2}}\right) - F\left(t_{n-1, \frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Deoarece repartiția Student este simetrică față de origine $t_{n-1, 1-\frac{\alpha}{2}} = -t_{n-1, \frac{\alpha}{2}}$ și înlocuindu-l pe T în relația anterioară, se obține

$$P\left(t_{n-1, \frac{\alpha}{2}} < T < t_{n-1, 1-\frac{\alpha}{2}}\right) = P\left(t_{n-1, \frac{\alpha}{2}} < \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} < t_{n-1, 1-\frac{\alpha}{2}}\right) = 1 - \alpha$$

$$\text{și } \bar{X} - t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$$

Ca urmare intervalul căutat este

$$(\theta_L, \theta_U) = \left(\bar{X} - t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}, \bar{X} + t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \right)$$

În acest caz eroarea este

$$E = t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$$

Dacă numărul de experiențe este $n > 30$, se poate folosi aproximația

$$t_{n-1, 1-\frac{\alpha}{2}} = z_{1-\frac{\alpha}{2}}$$

Estimarea intervalului de încredere $1 - \alpha$ pentru diferenței a două medii

Se consideră două selecții din populații normal repartizate $N(\mu_1, \sigma_1^2)$ și $N(\mu_2, \sigma_2^2)$.

Cazul dispersiilor σ_1^2, σ_2^2 cunoscute.

Considerăm o selecție aleatoare $x_{11}, x_{12}, \dots, x_{1n_1}$ din populația $N(\mu_1, \sigma_1^2)$ și o selecție $x_{21}, x_{22}, \dots, x_{2n_2}$ dintr-o populație $N(\mu_2, \sigma_2^2)$. Estimatorii nedeplasați ai mediilor μ_1 și μ_2 sunt:

$$\bar{X}_1 = \frac{\sum_{i=1}^{n_1} x_{1i}}{n_1} \text{ și } \bar{X}_2 = \frac{\sum_{i=1}^{n_2} x_{2i}}{n_2}$$

Considerând variabila aleatoare $\bar{X}_1 - \bar{X}_2$, ea este normal repartizată iar estimația și dispersia ei vor fi

$$M(\bar{X}_1 - \bar{X}_2) = M(\bar{X}_1) - M(\bar{X}_2) = \mu_1 - \mu_2 \text{ și } D(\bar{X}_1 - \bar{X}_2) = D(\bar{X}_1) + D(\bar{X}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \text{ unde am ținut}$$

cont că x_{1i} și x_{2i} sunt independente.

Mai departe, variabila aleatoare $z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{D(\bar{X}_1 - \bar{X}_2)}} = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$ este repartizată

$N(0,1)$.

Deoarece, $P\left(z_{\frac{\alpha}{2}} < z < z_{1-\frac{\alpha}{2}}\right) = 1 - \alpha$ și $z_{\frac{\alpha}{2}} = -z_{1-\frac{\alpha}{2}}$ rezultă

$$\left(\bar{X}_1 - \bar{X}_2\right) - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < \left(\bar{X}_1 - \bar{X}_2\right) + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Așadar, intervalul de estimație pentru diferența mediilor este

$$(\Theta_1, \Theta_2) = \left(\left(\bar{X}_1 - \bar{X}_2\right) - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \left(\bar{X}_1 - \bar{X}_2\right) + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

În acest caz, eroarea este $E = z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$.

Dispersii necunoscute dar presupuse egale

În cazul în care nu cunoaștem dispersiile dar știm că sunt egale $\sigma_1^2 = \sigma_2^2 = \sigma^2$ utilizăm dispersia ponderată de selecție

$$s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} = \frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (x_{2i} - \bar{X}_2)^2}{n_1+n_2-2}$$

ca un estimator nedeplasat pentru σ^2 .

Avem într-adevăr,

$$M(s_p^2) = \frac{(n_1-1)M(s_1^2) + (n_2-1)M(s_2^2)}{n_1+n_2-2} = \frac{(n_1-1)\sigma_1^2 + (n_2-1)\sigma_2^2}{n_1+n_2-2} = \sigma^2$$

În continuare vom arăta că mărimea $T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$ este repartizată $T(n_1 + n_2 - 2)$

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sigma_{\bar{X}_1 - \bar{X}_2}}$$

Se observă că $T = \frac{\sigma_{\bar{X}_1 - \bar{X}_2}}{\frac{s_p}{\sigma_{\bar{X}_1 - \bar{X}_2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$ este raportul între o variabilă aleatoare repartizată $N(0,1)$ și

deoarece

$$\frac{s_p}{\sigma_{\bar{X}_1 - \bar{X}_2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = \frac{s_p}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = \frac{s_p}{\sigma} = \sqrt{\frac{s_p^2}{\sigma^2}} =$$

$$\sqrt{\frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (x_{2i} - \bar{X}_2)^2}{(n_1 + n_2 - 2)\sigma^2}} = \sqrt{\frac{\sum_{i=1}^{n_1} \left(\frac{x_{1i} - \bar{X}_1}{\sigma}\right)^2 + \sum_{i=1}^{n_2} \left(\frac{x_{2i} - \bar{X}_2}{\sigma}\right)^2}{n_1 + n_2 - 2}}$$

variabila $\frac{s_p}{\sigma_{\bar{X}_1 - \bar{X}_2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$ este de tipul $\sqrt{\frac{\chi^2(n_1 + n_2 - 2)}{n_1 + n_2 - 2}}$

Dar $\sum_1^{n_1} \left(\frac{x_{1i} - \bar{X}_1}{\sigma} \right)^2$ este repartizat $\chi^2(n_1 - 1)$ iar $\sum_1^{n_2} \left(\frac{x_{2i} - \bar{X}_2}{\sigma} \right)^2$ este repartizat $\chi^2(n_2 - 1)$, deci T este repartizat $T(n_1 + n_2 - 2)$ și

$$P\left(t_{n_1+n_2-2, \frac{\alpha}{2}} < T < t_{n_1+n_2-2, 1-\frac{\alpha}{2}} \right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Deoarece repartiția Student este simetrică $t_{n_1+n_2-2, \frac{\alpha}{2}} = -t_{n_1+n_2-2, 1-\frac{\alpha}{2}}$ rezultă că

$$\bar{X}_1 - \bar{X}_2 - t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} < \mu_1 - \mu_2 < \bar{X}_1 - \bar{X}_2 + t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Deci, $(\Theta_1, \Theta_2) = \left(\bar{X}_1 - \bar{X}_2 - t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \bar{X}_1 - \bar{X}_2 + t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right)$ cu eroarea

$$E = t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}.$$

Estimarea intervalelor de încredere pentru dispersie

Considerăm o selecție de volum n dintr-o populație normală $N(\mu, \sigma^2)$. Conform celor arătate anterior variabila aleatoare $v = \frac{(n-1)s^2}{\sigma^2}$ este repartizată $\chi^2(n-1)$ și ca urmare

$$P\left(\chi_{n-1, \frac{\alpha}{2}}^2 < v < \chi_{n-1, 1-\frac{\alpha}{2}}^2 \right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Deci, $\chi_{n-1, \frac{\alpha}{2}}^2 < (n-1) \frac{s^2}{\sigma^2} < \chi_{n-1, 1-\frac{\alpha}{2}}^2$ și $\frac{(n-1)s^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{n-1, \frac{\alpha}{2}}^2}$.

Estimarea intervalului de încredere raportul a două dispersii

Se consideră selecția aleatoare $x_{11}, x_{12}, \dots, x_{1n_1}$ dintr-o populație $N(\mu_1, \sigma_1^2)$ și o selecție $x_{21}, x_{22}, \dots, x_{2n_2}$ dintr-o populație $N(\mu_2, \sigma_2^2)$.

Conform cu cele arătate anterior, raportul $F = \frac{s_1^2}{s_2^2}$ este repartizat $F(n_1 - 1, n_2 - 1)$ și deci

$$P\left(f_{n_1-1, n_2-1, \frac{\alpha}{2}} < F < f_{n_1-1, n_2-1, 1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Rezultă că $\frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, \frac{\alpha}{2}} < \frac{\sigma_2^2}{\sigma_1^2} < \frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, 1-\frac{\alpha}{2}}$, iar intervalul de estimare pentru raportul dispersiilor este:

$$(\Theta_L, \Theta_U) = \left(\frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, \frac{\alpha}{2}}, \frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, 1-\frac{\alpha}{2}} \right)$$

Aplicație: Utilizarea intervalelor de încredere în studiile de comparare a biodisponibilității medicamentelor¹

La introducerea în terapie de către un producător a unui medicament ce reprezintă o reproducere a altui medicament deja în uz, se pune problema comparării biodisponibilității acestora. În practică se cere ca raportul ariilor de sub curbele concentrațiilor plasmatice ale celor două medicamente să se afle în intervalul 0,8 - 1,25.

$$0,8 < \frac{\mu_{AUC}^T}{\mu_{AUC}^R} < 1,25$$

unde indicele T se referă la medicamentul testat și R desemnează medicamentul referință.

Atunci însă când ariile de sub curbă prezintă variabilități intra și interindividuale considerabile (determinările de biodisponibilitate se fac pe loturi de circa 10 – 20 de voluntari sănătoși) este de preferat a se determina un interval de încredere pentru media ariei realizată de medicamentul nou.

Pornind de la faptul ca $T = \frac{(\bar{X}_R - \bar{X}_T) - (\mu_R - \mu_T)}{s_p \sqrt{\frac{1}{n_R} + \frac{1}{n_T}}}$ este repartizată $T(n_R + n_T - 2)$ se deduce un

interval de încredere cu probabilitatea $1 - \alpha$ pentru $\mu_T - \mu_R$

$$\bar{X}_T - \bar{X}_R - t_{1-\frac{\alpha}{2}} < \mu_T - \mu_R < \bar{X}_T - \bar{X}_R + t_{1-\frac{\alpha}{2}}$$

unde am notat $s = s_p \sqrt{\frac{1}{n_R} + \frac{1}{n_T}}$.

După cum se va arăta mai departe, această estimare este puțin utilă în caz că s_p reprezintă practic intervariabilitatea, iar interschimbabilitatea care necesită bioechivalență trebuie să se bazeze pe intravariabilitatea.

STATISTICĂ MATEMATICĂ ȘI BIOSTATISTICĂ

Statistica matematică este principala aplicație a teoriei probabilităților. Metodele statistice constau, în esență, în elaborarea unor concluzii plauzibile privitoare la colectivități mari de fenomene, pe baza cunoașterii unui număr restrâns dintre acestea și extrapolării rezultatelor.

Legile care stau la baza statisticii și care permit aceste generalizări sunt teorema limită centrală și legea numerelor mari.

Într-o exprimare intuitivă, avem rezultatul că, dacă o variabilă aleatoare ξ este suma unui număr mare de variabile aleatoare independente, fiecare variabilă aleatoare având o pondere mică în sumă, atunci funcția de repartiție a variabilei aleatoare ξ este foarte apropiată de o funcție de repartiție normală.

Exprimat mai riguros și mai general, avem următoarea teoremă:

Teorema limită centrală (A.M. Leapunov)

Fie $\xi_1, \xi_2, \dots, \xi_n$ variabile aleatoare independente.

Fie $M(\xi_k) = a_k, D(\xi_k) = \sigma_k^2$ și $\rho_k^3 = M(|\xi_k - m_k|)^3$ când $k = \overline{1, n}$

Notăm $\sigma_{(n)}^2 = \sum_1^n \sigma_i^2$, $\rho_{(n)}^3 = \sum_1^n \rho_i^3$

Dacă $\lim_{n \rightarrow \infty} \frac{\rho_{(n)}}{\sigma_{(n)}} = 0$, atunci funcția de repartiție a variabilei

$$\frac{\xi_1 + \xi_2 + \dots + \xi_n - (a_1 + a_2 + \dots + a_n)}{\sigma_{(n)}}$$

tinde, când $n \rightarrow \infty$, către funcția $\Phi(x)$ a lui Laplace.

$$\Phi(x) = \frac{1}{2\pi} \int_{-\infty}^z e^{-\frac{t^2}{2}} dt$$

Teorema limită centrală este teorema fundamentală a teoriei erorilor. Laplace, Gauss și alți matematicieni, studiind repartiția erorilor, au ajuns la concluzia că funcția de repartiție normală poate fi luată drept model teoretic pentru cercetarea probabilistică a aproape tuturor fenomenelor naturii.

Teorema lui Cebâșev

Dacă $\xi_1, \xi_2, \dots, \xi_n$ sunt variabile aleatoare (discrete sau continue) independente ale căror dispersii sunt mai mici decât o constantă C, atunci oricare ar fi numărul pozitiv ε , probabilitatea inegalității

$$\left| \frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n} \right| < \varepsilon$$

tinde către 1, atunci când numărul variabilelor aleatoare tinde către infinit.

Demonstrație:

Să considerăm variabila aleatoare $\bar{\zeta} = \frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n}$. Având în vedere liniaritatea

operatorului de calcul a mediei avem $M(\bar{\zeta}) = \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}$.

Aplicând inegalitatea lui Cebâșev variabilei aleatoare $\bar{\zeta}$ se obține:

$$P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| < \varepsilon\right) \geq 1 - \frac{D\left(\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n}\right)}{\varepsilon^2}$$

Mai departe, din proprietățile operatorului D

$$D\left(\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n}\right) = \frac{D(\zeta_1) + D(\zeta_2) + \dots + D(\zeta_n)}{n^2} \leq \frac{C + C + \dots + C}{n^2} = \frac{nC}{n^2} = \frac{C}{n}$$

Deci

$$P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| < \varepsilon\right) \geq 1 - \frac{C}{n\varepsilon^2}$$

Trecând la limita pentru $n \rightarrow \infty$ obținem

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| < \varepsilon\right) \geq 1$$

și cum probabilitatea nu poate depăși 1,

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}\right| < \varepsilon\right) = 1$$

Cel mai frecvent, în practică, variabilele aleatoare ζ_i au aceeași medie μ și concluzia teoremei devine

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{\zeta_1 + \zeta_2 + \dots + \zeta_n}{n} - \mu\right| < \varepsilon\right) = 1$$

În esență, teorema lui Cebâșev stabilește că, deși variabilele aleatoare independente pot lua valori îndepărtate față de mediile lor, media aritmetică a unui număr suficient de mare de astfel de

variabile aleatoare ia cel mai probabil valori apropiate de un număr constant $\frac{M(\zeta_1) + M(\zeta_2) + \dots + M(\zeta_n)}{n}$ (sau μ atunci când mediile variabilelor sunt egale între ele).

Ca urmare, între comportarea fiecărei variabile aleatoare și comportarea mediilor lor există diferență esențială. Putem spune foarte precis ce valoare va lua media aritmetică a acestor variabile aleatoare. Explicația acestui fapt constă în aceea că abaterile diverselor variabile aleatoare sunt de semne diferite și, ca urmare, se compensează între ele.

TEORIA SELECȚIEI

Populații și selecții. Inferența statistică

În practică avem adesea nevoie să facem judecăți asupra unor mari colecții de rezultate posibile experimental ori a altor cantități, dar nu putem sau este extrem de scump, să examinăm toate aceste date. În astfel de cazuri, în loc să examinăm întregul set de date pe care îl numim în cele ce urmează *populație*, tragem concluziile după examinarea a o parte din ele, alese la întâmplare, parte pe care o numim *selecție*.

Procedeele de obținere a probelor este numit tot selecție, iar procedeele de extrapolare a concluziilor la întreaga populație este cunoscut ca *inferența statistică*.

Vom considera că o caracteristică dată a populației este o variabilă aleatoare pe un câmp de probabilitate (Ω, K, P) în care elementele lui Ω sunt chiar elementele populației, iar P este o probabilitate cunoscută sau nu.

Enumerarea valorilor observate ale caracteristicii urmărite și a frecvențelor lor relative definește *repartiția statistică a selecției*.

Teorema lui Leapunov, numită și teorema fundamentală a statisticii matematice, care justifică utilizarea metodei selecției stabilește că funcția de repartiție statistică a caracteristicilor selecțiilor tinde la funcția teoretică de repartiție a caracteristicii studiate când volmul selecției tinde la ∞ .

Exemplu 1

Putem dori să tragem concluzii despre evoluția rezistenței unei tulpini de germeni patogeni la un medicament dat și, în acest scop, examinăm rezultatele antibiogramelor făcute într-un eșantion de spitale într-o perioadă recentă (luniile de iarnă), comparată cu aceeași perioadă a anului precedent. Deși rezultatele obținute se referă la spitale și mai precis numai la o parte din ele, concluziile le extindem la scara întregii populații.

Exemplu 2

Rezultatele privind absorbția unui medicament după administrarea orală prin determinarea nivelurilor din plasma ale medicamentului la un lot de voluntari sănătoși le considerăm ca rezultate probabile pentru întreaga populație ce include și potențiali pacienți.

Populația poate fi infinită sau finită, în ultimul caz, numărul indivizilor populației – N- se mai numește și *volumul populației*. În mod similar, numărul de indivizi sau valori din cadrul unei probe este denumit *volumul probei* sau *volumul eșantionului*.

Valabilitatea concluziilor despre populație depinde de “reprezentativitatea” probei. Pentru populații finite aceasta înseamnă că fiecare membru al populației are aceeași șansă să fie selectat, când spunem că selecția este o selecție la întâmplare sau “selecție aleatoare”. Desigur că selecția unor voluntari sănătoși pentru determinarea parametrilor farmacocinetici ai unui medicament nu este din acest punct de vedere o selecție reprezentativă. În cazurile în care avem motive să credem că patologia căreia se adresează medicamentul nu afectează funcțiile metabolice și de excreție, această aproximare este acceptată pentru motivul că o selecție corectă ar implica loturi mult mai mari cu cheltuieli și timp de lucru mult crescute.

În practică, în studiile de bioechivalență, pentru reducerea volumului loturilor pe care se fac testările, se administrează amândouă medicamentele la toți membri lotului, în două perioade diferite. Fiecare component al lotului primește unul din medicamente în prima perioadă și celălalt în a doua perioadă.

Deoarece perioada de administrare poate influența și ea rezultatul experimentului, alegerea indivizilor care vor primi în prima perioadă primul medicament se face în mod aleator. În cazul când sunt mai multe perioade, de exemplu I-IV, și mai multe medicamente A, B, C, D se alcătuiește un tabel de felul

I	II	III	IV
A	B	D	C
B	C	A	D
C	D	B	A
D	A	C	B

așa zisul pătrat “latin”, unde observăm că fiecare literă apare o singură dată în fiecare linie și în fiecare coloană. Se numește pătrat latin deoarece, cum se va arata mai departe, în cazul în care mai intervine și o altă variabilă, de exemplu doza din fiecare medicament, se folosesc și litere grecești, alcătuiindu-se pătrate “greco-latine”.

Deasemenea, studiile de bioechivalență se fac tot pe voluntari sănătoși, pornind de la ipoteza că modificările de biodisponibilitate asociate stărilor patologice sunt aceleași pentru cele două medicamente testate, ceea ce, evident, este numai în parte adevărat.

În toate experimentele biologice, planificarea experimentului trebuie făcută în așa fel încât diferențele în tratament să nu coincidă cu diferențe în vârstă, sex, sau alți parametri. Dacă, de exemplu, femeile din lot primesc primul medicament și bărbații al doilea, se spune ca diferențele de sex sunt “confundate” cu diferențele de tratament. În acest caz nu se poate spune dacă diferențele obținute se datorează tratamentului sau diferenței de sex.

Parametrii de selecție ai unei variabile aleatoare :

Dacă printr-un procedeu oarecare cuantificăm răspunsul culturilor microbiene la antibioticele din exemplul 1, sau dacă luăm în considerație concentrațiile de medicament în sânge, din al doilea exemplu, și probabilitățile ca valorile să aparțină unor intervale diferite, obținem o variabilă aleatoare X asociată cu rezultatul experimentului corespunzător.

Parametrii acestei variabile aleatoare sunt denumiți, prin abuz de limbaj, “parametri ai populației”.

Dacă în exemplul al doilea X_i este concentrația de medicament în sângele bolnavului i , la o oră de la administrare, la primul voluntar putem obține o valoare x_1 , pentru al doilea voluntar o valoare x_2 , etc. În acest fel găsim valorile x_1, x_2, \dots, x_n ale variabilelor aleatoare independente X_1, X_2, \dots, X_n .

Media de selecție este o variabilă aleatoare: $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$

Dacă distribuția lui X este normală $-N(\mu, \sigma)$, aceiași pentru fiecare i , datorită linearității operatorului E care definește media, obținem $M(\bar{X}) = \mu_{\bar{X}} = \mu$ adică valoarea pentru *media mediei de selecție este media populației*.

Dacă la datele experimentale se adaugă o constantă, $x_i' = x_i + a$, media de selecție crește cu aceeași constantă: $\bar{W} = \frac{\sum_1^n (X_i + a)}{n} = \bar{X} + a$

Similar, dacă fiecare valoare se înmulțește cu o constanta $Z_i = kX_i$, media de selecție \bar{Z} se

înmulțește cu aceeași constantă:
$$\bar{Z} = \frac{\sum_1^n kX_i}{n} = k\bar{X}$$

Dispersia de selecție

Ca o măsură a abaterii datelor față de media de selecție, se introduce noțiunea de dispersie de

selecție
$$s_x^2 = \frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2$$
.

În aplicațiile practice, pentru reducerea numărului de calcule, formula se aduce la o altă formă și anume:

$$s_x^2 = \frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_1^n x_i^2 - 2\bar{X} \sum_1^n x_i + n\bar{X}^2 \right) = \frac{1}{n-1} \left(\sum_1^n x_i^2 - 2n\bar{X}^2 + n\bar{X}^2 \right) =$$

$$\frac{1}{n-1} \left(\sum_1^n x_i^2 - n\bar{X}^2 \right) = \frac{1}{n-1} \left(\sum_1^n x_i^2 - \frac{(\sum_1^n x_i)^2}{n} \right) \quad \text{Dacă}$$

$z_i = kx_i + a \Rightarrow s_z^2 = k^2 s_x^2$. Într-adevăr

$$s_z^2 = \frac{1}{n-1} \sum_1^n (z_i - \bar{Z})^2 = \frac{1}{n-1} \sum_1^n (kx_i + a - k\bar{X} - a)^2 = k^2 s_x^2$$

s_x se numește abaterea standard de selecție sau deviație standard, când nu este pericol de confuzie privind variabila aleatoare la care se referă folosindu-se și notația s_d .

Se mai utilizează în practică și noțiunile:

- Dispersia "populației" = $\frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2$ și respectiv
- Deviația standard a "populației", precum și "abaterea standard a mediei" (prescurtarea SEM – standard error of mean) definită prin raportul $SEM = \frac{s_x}{\sqrt{n}}$
- precum și coeficientul de variație $v = \frac{s_x}{\bar{X}} * 100$.

Covarianța de selecție

Covarianța de selecție se definește prin formula
$$s_{xy} = \frac{1}{n-1} \sum_1^n (x_i - \bar{X})(y_i - \bar{Y})$$

Se observă că aceasta se mai poate scrie și sub altă formă, mai utilă în sensul simplificărilor de calcul în anumite aplicații.

$$s_{XY} = \frac{1}{n-1} \left(\sum_1^n x_i y_i - \bar{X} \sum_1^n y_i - \bar{Y} \sum_1^n x_i + n \bar{X} \bar{Y} \right) = \frac{1}{n-1} \left(\sum_1^n x_i y_i - n \bar{X} \bar{Y} - n \bar{X} \bar{Y} + n \bar{X} \bar{Y} \right) =$$

$$\frac{1}{n-1} \left(\sum_1^n x_i y_i - n \bar{X} \bar{Y} \right) = \frac{1}{n-1} \left(\sum_1^n x_i y_i - \frac{\sum_1^n x_i \sum_1^n y_i}{n} \right)$$

Coeficientul de corelație de selecție

Coeficientul de corelație de selecție se definește prin formula

$$\rho(x, y) = \frac{s_{XY}}{s_X s_Y} = \frac{\frac{1}{n-1} \sum_1^n (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2} \sqrt{\frac{1}{n-1} \sum_1^n (y_i - \bar{Y})^2}} = \frac{\sum_1^n (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\sum_1^n (x_i - \bar{X})^2} \sqrt{\sum_1^n (y_i - \bar{Y})^2}}$$

Proprietăți ale caracteristicilor de selecție

Considerăm în continuare o selecție de volum n dintr-o populație cu media μ și dispersia σ^2

Propoziție

Media mediei de selecție este egală cu media populației. $M(\bar{X}) = \mu$

Demonstrație:

$$M(\bar{X}) = \frac{M(\sum x_i)}{n} = \frac{\sum M(x_i)}{n} = \frac{n\mu}{n} = \mu$$

Propoziție

Media dispersiei de selecție este egală cu dispersia populației $M(s_X^2) = \sigma^2$

Demonstrație:

$$M(s_X^2) = M\left(\frac{1}{n-1} \sum_1^n (x_i - \bar{X})^2\right) = \frac{1}{n-1} M\left(\sum_1^n x_i^2 - 2\bar{X} \sum_1^n x_i + \sum_1^n \bar{X}^2\right) =$$

$$\frac{1}{n-1} M\left(\sum_1^n x_i^2 - 2n\bar{X}^2 + n\bar{X}^2\right) = -\frac{1}{n-1} M\left(\sum_1^n x_i^2 - n\bar{X}^2\right) = \frac{1}{n-1} M\left(\sum_1^n x_i^2 - \frac{\left(\sum_1^n x_i\right)^2}{n}\right)$$

Dar, mai departe

$$M\left(\sum_1^n x_i^2\right) = n(\sigma^2 + \mu^2)$$

$$M\left(\sum_1^n x_i\right)^2 = M\left(\sum_1^n x_i^2 + 2\sum_{i \neq j}^n x_i x_j\right) = \sum_1^n M(x_i^2) + 2\frac{n(n-1)}{2}M(x_i)M(x_j) = \\ = n(\sigma^2 + \mu^2) + n(n-1)\mu^2 = n\sigma^2 + n^2\mu^2$$

și înlocuind în expresia lui $M(s_x^2)$ obținem $M(s_x^2) = \frac{n(\sigma^2 + \mu^2) - \sigma^2 - n\mu^2}{n-1} = \sigma^2$

Propoziție

Variabila aleatoare $Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$ tinde, când $n \rightarrow \infty$ către o variabilă $N(0,1)$

Aceasta este o consecință a teoremei limită centrală și este aplicabilă atât variabilelor continue cât și celor discrete.

Într-adevăr aplicând teorema lui Leapunov pentru variabilele aleatoare x_1, x_2, \dots, x_n obținem că:

$$\frac{x_1 + x_2 + \dots + x_n - (\mu_1 + \mu_2 + \dots + \mu_n)}{\sqrt{\sigma^2 + \sigma^2 + \dots + \sigma^2}} = \frac{n\bar{X} - n\mu}{\sqrt{n\sigma^2}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

tinde către o variabilă aleatoare $N(0,1)$.

Propoziție

Dacă variabila aleatoare X este normal distribuită, atunci variabila aleatoare $V = (n-1)\frac{s_x^2}{\sigma^2}$

este repartizată $\chi^2(n-1)$

Demonstrație:

$$V = (n-1)\frac{s_x^2}{\sigma^2} = \frac{\sum_1^n (x_i - \bar{X})}{\sigma^2} = \frac{\sum_1^n [(x_i - \mu) - (\bar{X} - \mu)]^2}{\sigma^2} = \\ = \frac{\sum_1^n (x_i - \mu)^2 - 2\sum_1^n (x_i - \mu)(\bar{X} - \mu) + \sum_1^n (\bar{X} - \mu)^2}{\sigma^2} = \frac{\sum_1^n (x_i - \mu)^2 - 2(n\bar{X} - n\mu)(\bar{X} - \mu) + n(\bar{X} - \mu)^2}{\sigma^2} = \\ = \frac{\sum_1^n (x_i - \mu) - n(\bar{X} - \mu)}{\sigma^2} = \sum_1^n \left(\frac{x_i - \mu}{\sigma} \right)^2 - \left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \right)^2$$

Dar variabila aleatoare $\frac{x_i - \mu}{\sigma}$ este repartizată $N(0,1)$ deoarece $M\left(\frac{x_i - \mu}{\sigma}\right) = \frac{M(x_i) - \mu}{\sigma}$ și

$D\left(\frac{x_i - \mu}{\sigma}\right) = \frac{D(x_i)}{\sigma^2} = 1$, iar $\frac{(\bar{X} - \mu)}{\frac{\sigma}{\sqrt{n}}}$ este repartizată tot $N(0,1)$ în conformitate cu teorema limită

centrală.

Deci, V este o sumă de $n-1$ pătrate de variabile de tip $N(0,1)$.

Propoziție

Dacă x_1, x_2, \dots, x_n este o selecție dintr-o populație normal distribuită, atunci variabila aleatoare

$T = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$ este repartizată Student cu n grade de libertate.

Demonstrație:

$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} = \frac{\frac{\bar{X} - \mu}{\sigma}}{\frac{s}{\sigma}} = \frac{\frac{\bar{X} - \mu}{\sigma}}{\sqrt{\frac{\sum_1^n (x_i - \bar{X})^2}{(n-1)\sigma^2}}} = \frac{Z}{\sqrt{\frac{V}{n-1}}}$$

unde $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ este repartizată $N(0,1)$, iar $V = \frac{\sum_1^n (x_i - \bar{X})^2}{\sigma^2}$ este repartizată $\chi^2(n-1)$.

Deci, T este repartizată Student cu $n-1$ grade de libertate.

Propoziție

Date fiind două selecții aleatoare independente $x_{11}, x_{12}, \dots, x_{1n_1}$ și $x_{21}, x_{22}, \dots, x_{2n_2}$ din populații

normal distribuite $N(\mu_1, \sigma_1)$ și $N(\mu_2, \sigma_2)$, variabila aleatoare $F = \frac{\frac{S_1^2}{\sigma_1^2}}{\frac{S_2^2}{\sigma_2^2}}$ este repartizată Fisher –

Snedecor $F(n_1 - 1, n_2 - 1)$

Demonstrație:

Avem într-adevar $F = \frac{S_1^2}{S_2^2} = \frac{\frac{\sum_1^{n_1} (x_{1i} - \bar{X}_1)^2}{(n_1 - 1)\sigma_1^2}}{\frac{\sum_1^{n_2} (x_{2i} - \bar{X}_2)^2}{(n_2 - 1)\sigma_2^2}}$ iar numărătorul și numitorul sunt repartizate, conform

propozitiei 2.3.5.4., respectiv $\frac{\chi^2(n_1 - 1)}{n_1 - 1}$ și $\frac{\chi^2(n_2 - 1)}{n_2 - 1}$.

Estimații

Teoria estimației urmărește evaluarea parametrilor unei repartiții în general cunoscute. Valorile numerice obținute se numesc *estimații* sau *estimatori*. Se obțin estimații punctuale în cazul în care se folosesc datele selecției pentru a obține valorile parametrilor și estimații ale intervalelor de încredere în cazul în care se determină un interval în care se află, cu o anumită probabilitate valoarea estimată.

Un estimator al parametrului θ se va nota cu $\hat{\theta}$. O estimație este nedeplasată dacă $M(\hat{\theta}) = \theta$, adică media estimației este egală chiar cu valoarea teoretică a parametrului estimat.

Conform proprietății 2.3.5.1, $M(\bar{X}) = \mu$ adică media de selecție este un estimator nedeplasat al mediei, iar conform proprietății 2.3.5.2., $M(s^2) = \sigma^2$ adică dispersia de selecție este un estimator nedeplasat al dispersiei.

Problema estimării intervalelor se reduce la găsirea unui interval de încredere (θ_L, θ_U) cu un coeficient de încredere $1 - \alpha$ astfel încât $P(\theta_L < \theta < \theta_U) = 1 - \alpha$.

Este de dorit ca $1 - \alpha$ să fie cât mai mare (de obicei este cuprins între 0,9 și 0,99) iar intervalul (θ_L, θ_U) să fie cât mai mic. În stabilirea intervalelor se utilizează caracteristicile numerice cuantile. Se numesc *cuantile de ordin β* valoarea x_β a variabilei aleatoare x pentru care $F(x_\beta) = P(x < x_\beta) = \beta$ adică valoarea variabilei aleatoare care are la stânga ei aria β sub curba densității de probabilitate.

Evident:

$$P\left(x < x_{\frac{\alpha}{2}}\right) = \frac{\alpha}{2}$$

$$P\left(x < x_{1 - \frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2}$$

$$P\left(x_{\frac{\alpha}{2}} < \bar{x} < x_{1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Pentru a estima un interval se alege $1 - \alpha$, se citesc din tabelele cuantilele, de exemplu $x_{1-\frac{\alpha}{2}}$ și $x_{\frac{\alpha}{2}}$ și se precizează intervalul. În prealabil, în funcție de mărimea pentru care se caută intervalul se precizează cu care din repartițiile cunoscute trebuie lucrat.

Estimarea intervalelor de încredere pentru medii

2.4.1.1. Cazul când se cunoaște dispersia.

Se consideră o populație repartizată normal $N(\mu, \sigma^2)$. Dacă se cunoaște dispersia se poate folosi faptul că $z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ este repartizată $N(0,1)$. Se notează cu z_{α} cuantila de ordinul α pentru repartiția

$N(0,1)$. Evident

$$P\left(z_{\frac{\alpha}{2}} < z < z_{1-\frac{\alpha}{2}}\right) = F\left(z_{1-\frac{\alpha}{2}}\right) - F\left(z_{\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Așadar intervalul $\left(z_{\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}}\right)$ este un interval de estimare cu coeficientul de încredere $1 - \alpha$. Din

anumite puncte de vedere este recomandabil să se utilizeze acele intervale care lasă atât la dreapta cât și la stânga lor aceeași arie, egală cu $\frac{\alpha}{2}$.

Deoarece repartiția $N(0,1)$ este simetrică față de axa Oy avem relația $z_{\frac{\alpha}{2}} = -z_{1-\frac{\alpha}{2}}$

Din relațiile

$$\begin{aligned} -z_{1-\frac{\alpha}{2}} < z < z_{1-\frac{\alpha}{2}} &\Rightarrow -z_{1-\frac{\alpha}{2}} < \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} < z_{1-\frac{\alpha}{2}} \Rightarrow -z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} < \bar{x} - \mu < z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} \Rightarrow \\ &-\bar{x} - z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} < -\mu < -\bar{x} + z_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} \end{aligned}$$

rezultă

$$\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

Așadar intervalul căutat este

$$(\theta_L, \theta_U) = \left(\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right)$$

Mărimea $E = z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ poartă numele de *eroare* și servește la calculul numărului de experiențe

$$n = \left(\frac{z_{1-\frac{\alpha}{2}}}{E} \right)^2 \text{ atunci când este impusă eroarea și se alege un coeficient } 1 - \alpha$$

Metoda descrisă mai poate fi aplicată și în cazul în care x nu este repartizată normal deoarece z este repartizată $N(0,1)$ indiferent de repartiția variabilelor x_1, x_2, \dots, x_n (teorema limită centrală).

Cazul când dispersia este necunoscută

Dacă nu se cunoaște dispersia în estimarea intervalelor se utilizează dispersia de selecție care este un estimator nedeplasat al dispersiei deoarece $E(s^2) = \sigma^2$

Se consideră x_1, x_2, \dots, x_n o selecție dintr-o populație de tipul $N(\mu, \sigma^2)$.

Conform celor arătate anterior mărimea $T = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$ este repartizată $T(n-1)$ și, ca urmare

$$P\left(t_{n-1, \frac{\alpha}{2}} < T < t_{n-1, 1-\frac{\alpha}{2}}\right) = F\left(t_{n-1, 1-\frac{\alpha}{2}}\right) - F\left(t_{n-1, \frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Deoarece repartiția Student este simetrică față de origine $t_{n-1, 1-\frac{\alpha}{2}} = -t_{n-1, \frac{\alpha}{2}}$ și înlocuindu-l pe T în relația anterioară, se obține

$$P\left(t_{n-1, \frac{\alpha}{2}} < T < t_{n-1, 1-\frac{\alpha}{2}}\right) = P\left(t_{n-1, \frac{\alpha}{2}} < \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} < t_{n-1, 1-\frac{\alpha}{2}}\right) = 1 - \alpha$$

$$\text{și } \bar{X} - t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$$

Ca urmare intervalul căutat este

$$(\theta_L, \theta_U) = \left(\bar{X} - t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}, \bar{X} + t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \right)$$

În acest caz eroarea este

$$E = t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$$

Dacă numărul de experiențe este $n > 30$, se poate folosi aproximația

$$t_{n-1, 1-\frac{\alpha}{2}} = z_{1-\frac{\alpha}{2}}$$

Estimarea intervalului de încredere $1 - \alpha$ pentru diferenței a două medii

Se consideră două selecții din populații normal repartizate $N(\mu_1, \sigma_1^2)$ și $N(\mu_2, \sigma_2^2)$.

Cazul dispersiilor σ_1^2, σ_2^2 cunoscute.

Considerăm o selecție aleatoare $x_{11}, x_{12}, \dots, x_{1n_1}$ din populația $N(\mu_1, \sigma_1^2)$ și o selecție $x_{21}, x_{22}, \dots, x_{2n_2}$ dintr-o populație $N(\mu_2, \sigma_2^2)$. Estimatorii nedeplasați ai mediilor μ_1 și μ_2 sunt:

$$\bar{X}_1 = \frac{\sum_{i=1}^{n_1} x_{1i}}{n_1} \text{ și } \bar{X}_2 = \frac{\sum_{i=1}^{n_2} x_{2i}}{n_2}$$

Considerând variabila aleatoare $\bar{X}_1 - \bar{X}_2$, ea este normal repartizată iar estimăția și dispersia ei vor fi

$$M(\bar{X}_1 - \bar{X}_2) = M(\bar{X}_1) - M(\bar{X}_2) = \mu_1 - \mu_2 \text{ și } D(\bar{X}_1 - \bar{X}_2) = D(\bar{X}_1) + D(\bar{X}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \text{ unde am ținut}$$

cont că x_{1i} și x_{2i} sunt independente.

Mai departe, variabila aleatoare $z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{D(\bar{X}_1 - \bar{X}_2)}} = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$ este repartizată

$N(0,1)$.

Deoarece, $P\left(z_{\frac{\alpha}{2}} < z < z_{1-\frac{\alpha}{2}}\right) = 1 - \alpha$ și $z_{\frac{\alpha}{2}} = -z_{1-\frac{\alpha}{2}}$ rezultă

$$\left(\bar{X}_1 - \bar{X}_2\right) - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < \left(\bar{X}_1 - \bar{X}_2\right) + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Așadar, intervalul de estimăție pentru diferența mediilor este

$$(\Theta_1, \Theta_2) = \left(\left(\bar{X}_1 - \bar{X}_2\right) - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \left(\bar{X}_1 - \bar{X}_2\right) + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

În acest caz, eroarea este $E = z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$.

Dispersii necunoscute dar presupuse egale

În cazul în care nu cunoaștem dispersiile dar știm că sunt egale $\sigma_1^2 = \sigma_2^2 = \sigma^2$ utilizăm dispersia ponderată de selecție

$$s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} = \frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (x_{2i} - \bar{X}_2)^2}{n_1+n_2-2}$$

ca un estimator nedeplasat pentru σ^2 .

Avem într-adevăr,

$$M(s_p^2) = \frac{(n_1-1)M(s_1^2) + (n_2-1)M(s_2^2)}{n_1+n_2-2} = \frac{(n_1-1)\sigma_1^2 + (n_2-1)\sigma_2^2}{n_1+n_2-2} = \sigma^2$$

În continuare vom arăta că mărimea $T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$ este repartizată $T(n_1 + n_2 - 2)$

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sigma_{\bar{X}_1 - \bar{X}_2}}$$

Se observă că $T = \frac{\sigma_{\bar{X}_1 - \bar{X}_2}}{\frac{s_p}{\sigma_{\bar{X}_1 - \bar{X}_2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$ este raportul între o variabilă aleatoare repartizată $N(0,1)$ și

deoarece

$$\frac{s_p}{\sigma_{\bar{X}_1 - \bar{X}_2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = \frac{s_p}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = \frac{s_p}{\sigma} = \sqrt{\frac{s_p^2}{\sigma^2}} =$$

$$\sqrt{\frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (x_{2i} - \bar{X}_2)^2}{(n_1 + n_2 - 2)\sigma^2}} = \sqrt{\frac{\sum_{i=1}^{n_1} \left(\frac{x_{1i} - \bar{X}_1}{\sigma}\right)^2 + \sum_{i=1}^{n_2} \left(\frac{x_{2i} - \bar{X}_2}{\sigma}\right)^2}{n_1 + n_2 - 2}}$$

variabila $\frac{s_p}{\sigma_{\bar{X}_1 - \bar{X}_2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$ este de tipul $\sqrt{\frac{\chi^2(n_1 + n_2 - 2)}{n_1 + n_2 - 2}}$

Dar $\sum_1^{n_1} \left(\frac{x_{1i} - \bar{X}_1}{\sigma} \right)^2$ este repartizat $\chi^2(n_1 - 1)$ iar $\sum_1^{n_2} \left(\frac{x_{2i} - \bar{X}_2}{\sigma} \right)^2$ este repartizat $\chi^2(n_2 - 1)$, deci T este repartizat $T(n_1 + n_2 - 2)$ și

$$P\left(t_{n_1+n_2-2, \frac{\alpha}{2}} < T < t_{n_1+n_2-2, 1-\frac{\alpha}{2}} \right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Deoarece repartiția Student este simetrică $t_{n_1+n_2-2, \frac{\alpha}{2}} = -t_{n_1+n_2-2, 1-\frac{\alpha}{2}}$ rezultă că

$$\bar{X}_1 - \bar{X}_2 - t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} < \mu_1 - \mu_2 < \bar{X}_1 - \bar{X}_2 + t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Deci, $(\Theta_1, \Theta_2) = \left(\bar{X}_1 - \bar{X}_2 - t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \bar{X}_1 - \bar{X}_2 + t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right)$ cu eroarea

$$E = t_{n_1+n_2-2, 1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}.$$

Estimarea intervalelor de încredere pentru dispersie

Considerăm o selecție de volum n dintr-o populație normală $N(\mu, \sigma^2)$. Conform celor arătate anterior variabila aleatoare $v = \frac{(n-1)s^2}{\sigma^2}$ este repartizată $\chi^2(n-1)$ și ca urmare

$$P\left(\chi_{n-1, \frac{\alpha}{2}}^2 < v < \chi_{n-1, 1-\frac{\alpha}{2}}^2 \right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Deci, $\chi_{n-1, \frac{\alpha}{2}}^2 < \frac{(n-1)s^2}{\sigma^2} < \chi_{n-1, 1-\frac{\alpha}{2}}^2$ și $\frac{(n-1)s^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{n-1, \frac{\alpha}{2}}^2}$.

Estimarea intervalului de încredere raportul a două dispersii

Se consideră selecția aleatoare $x_{11}, x_{12}, \dots, x_{1n_1}$ dintr-o populație $N(\mu_1, \sigma_1^2)$ și o selecție $x_{21}, x_{22}, \dots, x_{2n_2}$ dintr-o populație $N(\mu_2, \sigma_2^2)$.

Conform cu cele arătate anterior, raportul $F = \frac{s_1^2}{s_2^2}$ este repartizat $F(n_1 - 1, n_2 - 1)$ și deci

$$P\left(f_{n_1-1, n_2-1, \frac{\alpha}{2}} < F < f_{n_1-1, n_2-1, 1-\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha$$

Rezultă că $\frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, \frac{\alpha}{2}} < \frac{\sigma_2^2}{\sigma_1^2} < \frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, 1-\frac{\alpha}{2}}$, iar intervalul de estimare pentru raportul dispersiilor este:

$$(\Theta_L, \Theta_U) = \left(\frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, \frac{\alpha}{2}}, \frac{s_2^2}{s_1^2} f_{n_1-1, n_2-1, 1-\frac{\alpha}{2}} \right)$$

Aplicație: Utilizarea intervalelor de încredere în studiile de comparare a biodisponibilității medicamentelor¹

La introducerea în terapie de către un producător a unui medicament ce reprezintă o reproducere a altui medicament deja în uz, se pune problema comparării biodisponibilității acestora. În practică se cere ca raportul ariilor de sub curbele concentrațiilor plasmatice ale celor două medicamente să se afle în intervalul 0,8 - 1,25.

$$0,8 < \frac{\mu_{AUC}^T}{\mu_{AUC}^R} < 1,25$$

unde indicele T se referă la medicamentul testat și R desemnează medicamentul referință.

Atunci însă când ariile de sub curbă prezintă variabilități intra și interindividuale considerabile (determinările de biodisponibilitate se fac pe loturi de circa 10 – 20 de voluntari sănătoși) este de preferat a se determina un interval de încredere pentru media ariei realizată de medicamentul nou.

Pornind de la faptul ca $T = \frac{(\bar{X}_R - \bar{X}_T) - (\mu_R - \mu_T)}{s_p \sqrt{\frac{1}{n_R} + \frac{1}{n_T}}}$ este repartizată $T(n_R + n_T - 2)$ se deduce un

interval de încredere cu probabilitatea $1 - \alpha$ pentru $\mu_T - \mu_R$

$$\bar{X}_T - \bar{X}_R - t_{1-\frac{\alpha}{2}} < \mu_T - \mu_R < \bar{X}_T - \bar{X}_R + t_{1-\frac{\alpha}{2}}$$

unde am notat $s = s_p \sqrt{\frac{1}{n_R} + \frac{1}{n_T}}$.

După cum se va arăta mai departe, această estimare este puțin utilă în caz că s_p reprezintă practic intervariabilitatea, iar interschimbabilitatea care necesită bioechivalență trebuie să se bazeze pe intravariabilitatea.

¹W.J.Westlake: *Use of confidence intervals in analysis of comparative bioavailability trials*, *J. Pharm. Sci.*, 61 (8), 1340 – 1, 1972.