

O **variabilă discretă** este o variabilă care are un număr finit sau infinit dar numărabil de valori; o astfel de variabilă poate avea valori finite sau infinite.

Un **eșantion** este o submulțime a populației. Dintr-un punct de vedere teoretic fiecare individ are aceleași șanse de a aparține eșan-
tionului, și orice grup particular de indivizi este ales în mod independent pentru a face parte din eșantion. Dacă aceste condiții sunt orice valoare dintr-un interval real, incluzând orice valoare posibilă îndeplinite atunci avem un **eșantion aleator simplu**.

Datele sunt valori variabilei colectate de la fiecare individ din eșan-
tion.

O **populație** este o mulțime de obiecte (numite și indivizi) ale căror proprietăți vor fi analizate.

O **variabilă nominală** este o variabilă care numește sau descrie un individ dintr-o populație fără a putea asigura o ordine naturală acestor valori..

O **variabilă ordinală** este o variabilă ale cărei valori pot fi ordonate în mod natural.

Deviația standard a eșantionului, s , este rădăcina pătrată a dispersiei eșantionului.

Dispersia eșantionului, s^2 , n fiind dimensiunea eșantionului, este:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)^2}{n-1},$$

Median (Me) este valoarea din mijloc când datele din eșantion sunt sortate.

Momentul central de ordin k al populației este

$$\mu_k = M[(X - \mu)^k].$$

Momentul central de ordin k al eșantionului X este

$$m_k = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)^k}{n}.$$

Valoarea P este probabilitatea de a obține un rezultat cel puțin încredere (de la fel de neobișnuit (extrem) ca rezultatul obținut din eșantion. Când valoarea P este "mică" putem respinge H_0 . (Legea tare numerelor mari) Fie $(X_n)_{n \geq 1}$ un șir de variabile aleatoare independente și identic distribuite cu media μ și dispersia σ^2 . Atunci

$$P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu\right) = 1.$$

Un interval de încredere cu nivelul de încredere $(1 - \alpha)$ pentru media unei populații cu dispersia cunoscută este

$$\left(\bar{x}_n - z^* \frac{\sigma}{\sqrt{n}}, \bar{x}_n + z^* \frac{\sigma}{\sqrt{n}}\right),$$

unde z^* este valoarea critică asociată cu $\alpha/2$. Mai mult, acest interval este exact pentru o populație distribuită normal și aprox-
imativ altfel, când eșantionul este suficient de mare ($n \geq 30$). Fie $(X_n)_{n \geq 1}$ un șir de variabile aleatoare independente având dispersii finite, uniform mărginite, i. e. $D^2[X_n] \leq c$, pentru orice $n \geq 1$. Atunci

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n M[X_i]\right| < \epsilon\right) = 1.$$

Testarea ipotezelor statistice este un proces prin care se ia o decizie între două ipoteze opuse.

Ipotezele statistice sunt formulate în așa fel încât întotdeauna una din ele este falsă iar cealaltă adevărată.

Una dintre ipoteze este testată sperând că se poate arăta că este puțin probabil ca ea să fie adevărată, ceea ce implică faptul că cealaltă ipoteză este probabil adevărată.

Ecuția acestei drepte (linia SD) este

$$y - \bar{y}_n = m(x - \bar{x}_n),$$

unde \bar{x}_n și \bar{y}_n sunt mediile.

Scorul testului este statistica ce corespunde valorii P ; **valoarea critică**, pe de altă parte, corespunde nivelului de semnificație.

Concluzia testului depinde de rezultatul comparării celor două valori.

Media și dispersia mediei de selecție, \bar{x}_n , sunt μ și σ^2/n :

$$M[\bar{x}_n] = \mu, D^2[\bar{x}_n] = \frac{\sigma^2}{n}.$$

În plus, pentru valori mari ale lui n (≥ 30), distribuția mediei de selecție este normală, i. e.

$$\bar{x}_n \sim N(\mu, \sigma^2/n).$$

variabila aleatoare care are drept valori toate mediiile de selecție posibile, $\bar{x}_n^{(k)}$, pentru eșantioane de dimensiune n , se numește **Un test parametric** presupune că populația urmează o anumită distribuție și inferează asupra parametrilor acelei distribuții.

Un parametru este o valoare numerică care privește întreaga pop- **O statistică** este un parametru calculat pentru un eșantion în locul întregii populații.

(Teorema limită centrală, Lindeberg-Lévy) Fie $(X_n)_{n \geq 1}$ un șir de variabile aleatoare independente și identic distribuite cu media μ și dispersia σ^2 . Atunci

$$\frac{\frac{1}{n} \sum_{i=1}^n X_i - \mu}{\frac{\sigma}{\sqrt{n}}} \rightarrow N(0, 1) \text{ sau}$$

$$\lim_{n \rightarrow \infty} P \left(\left| \frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \right| \leq a \right) = \frac{1}{\sqrt{2\pi}} \int_{-a}^a \exp -t^2/2 dt.$$

Un interval de încredere pentru un parametru θ cu $(1 - \alpha)$ nivel de încredere este definit cu ajutorul a două statistici L și U astfel ca

$$P(L \leq \theta \leq U) = (1 - \alpha).$$

Testul F - Inferență asupra raportului dintre dispersii

• Testul se desfășoară astfel:

1. Formulăm **ipoteza nulă**:

$$H_0: \frac{\sigma_1^2}{\sigma_2^2} = 1$$

2. Formulăm **ipoteza alternativă** conform informațiilor obținute din eșantion. Putem avea două tipuri de ipoteză alternativă

$$H_a: \frac{\sigma_1^2}{\sigma_2^2} > 1 \quad (\text{asimetrică la dreapta}) \text{ pentru un un test one-tailed}$$

$$H_a: \frac{\sigma_1^2}{\sigma_2^2} \neq 1 \quad (\text{simetrică}) \text{ pentru un un test two-tailed}$$

Testul T - inferență asupra mediilor a două populații (σ_1^2, σ_2^2 necunoscute)

Testul decurge astfel:

1. Formulăm **ipoteza nulă**, care susține că diferența mediilor celor două populații ia o anumită valoare:

$$H_0: \mu_1 - \mu_2 = m_0$$

2. Formulăm **ipoteza alternativă** conform datelor din eșantioane. Putem avea trei tipuri de ipoteză alternativă

$$H_a: \mu_1 - \mu_2 < m_0 \quad (\text{asimetrică la stânga}) \text{ sau}$$

$$H_a: \mu_1 - \mu_2 > m_0 \quad (\text{asimetrică la dreapta}) \text{ sau}$$

$$H_a: \mu_1 - \mu_2 \neq m_0 \quad (\text{simetrică}).$$

Ipotezele asimetrice se mai numesc **one-tailed**, iar cea simetrică **two-tailed**.

3. Alegem nivelul de semnificație $\alpha \in \{1\%, 5\%\}$.

• Putem utiliza testul Z chiar dacă cele două populații sunt doar aproximativ normal distribuite, dacă cele două eșantioane sunt suficient de mari ($n_1, n_2 \geq 30$).

• Testul se desfășoară astfel:

1. Formulăm **ipoteza nulă**, care susține că diferența mediilor celor două populații ia o valoare fixată:

$$H_0: \mu_1 - \mu_2 = m_0$$

2. Formulăm **ipoteza alternativă** conform informațiilor obținute din eșantioane.

6. Comparăm valoarea critică cu scorul z ; dacă scorul z aparține **zonei de respingere**, atunci acceptăm H_a și respingem H_0 . Zonele de respingere sunt:

$$(-\infty, z^*] \text{ pentru } H_a \text{ asimetrică la stânga,}$$

$$[z^*, +\infty) \text{ pentru } H_a \text{ asimetrică la dreapta,}$$

$$(-\infty, -|z^*|] \cup [|z^*|, +\infty) \text{ pentru } H_a \text{ simetrică.}$$

Dacă scorul z nu aparține zonei de respingere spunem că **nu există suficiente dovezi cu nivelul de semnificație α pentru a respinge ipoteza nulă (încercarea de a respinge H_0 eșuează).**

6. Comparăm valoarea critică cu scorul t ; dacă scorul t aparține **zonei de respingere**, atunci acceptăm H_a și respingem H_0 . Zonele de respingere sunt:

$$(-\infty, t^*] \text{ pentru } H_a \text{ asimetrică la stânga,}$$

$$[t^*, +\infty) \text{ pentru } H_a \text{ asimetrică la dreapta,}$$

$$(-\infty, -|t^*|] \cup [|t^*|, +\infty) \text{ pentru } H_a \text{ simetrică.}$$

Dacă scorul t nu aparține zonei de respingere vom spune că **nu există suficiente dovezi cu nivelul de semnificație α pentru a respinge ipoteza nulă (încercarea de a respinge H_0 eșuează).**

3. Comparăm valoarea critică cu scorul z ; dacă scorul z aparține **zonei de respingere**, atunci acceptăm H_a și respingem H_0 . Zonele de respingere sunt:

$$(-\infty, z^*] \text{ pentru } H_a \text{ asimetrică la stânga,}$$

$$[z^*, +\infty) \text{ pentru } H_a \text{ asimetrică la dreapta,}$$

Dimensiunea efectului este mărimea diferenței descoperite în eșantionul aleator (dacă există).

Nivelul de semnificație, α , este probabilitatea (condiționată) maximă pe care ne-o asumăm drept risc de face o eroare de tipul I.

Puterea testului, este 1 minus probabilitatea de a face o eroare de tipul II.

Un test neparametric numit și **distribution-free** sau **parametru-free** se bazează pe puține fapte - de obicei distribuția și parametrii săi (medie, dispersie) nu sunt cunoscuți.

Decizia asupra H_0	se respinge nu se respinge	validitatea ipotezei H_0	
		adevărată	falsă
Decizia asupra H_0	se respinge	Eroare de tip I (fals pozitiv)	Corect (adevărat pozitiv)
	nu se respinge	Corect (adevărat negativ)	Eroare de tip II (fals negativ)

3. Alegem nivelul de semnificație $\alpha \in \{1\%, 5\%\}$.

4. Calculăm **scorul F (statistica testului)**

$$F = \frac{s_1^2}{s_2^2}$$

5. Determinăm valoarea critică corespunzătoare pentru α

$$F^* = qf(1 - \alpha, n_1 - 1, n_2 - 1) \text{ pentru } H_a \text{ asimetrică la dreapta}$$

$$F_2^* = qf(\alpha/2, n_1 - 1, n_2 - 1), F_1^* = qf(1 - \alpha/2, n_1 - 1, n_2 - 1)$$

pentru H_a simetrică.

4. Calculăm **scorul t (statistica testului)**

a) dacă dispersiile sunt egale:

$$t = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$$

unde $s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$, iar numărul de grade de libertate este $df = n_1 + n_2 - 2$

b) dacă dispersiile sunt diferite:

$$t = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

numărul de grade de libertate fiind $df = \min(n_1 - 1, n_2 - 1)$.

Putem avea trei tipuri de ipoteză alternativă

$$H_a: \mu_1 - \mu_2 < m_0 \quad (\text{asimetrică la stânga}) \text{ sau}$$

$$H_a: \mu_1 - \mu_2 > m_0 \quad (\text{asimetrică la dreapta}) \text{ sau}$$

$$H_a: \mu_1 - \mu_2 \neq m_0 \quad (\text{simetrică}).$$

Ipotezele asimetrice se mai numesc **one-tailed**, iar cea simetrică **two-tailed**.

3. Alegem un nivel de semnificație $\alpha \in \{1\%, 5\%\}$.

1. Formulăm mai întâi **ipoteza nulă**, care susține că media populației ia o anumită valoare:

$$H_0: \mu = \mu_0$$

2. Formulăm **ipoteza alternativă** conform informațiilor obținute din eșantion. Putem avea trei tipuri de ipoteză alternativă

$$H_a: \mu < \mu_0 \quad (\text{asimetrică la stânga}) \text{ sau}$$

$$H_a: \mu > \mu_0 \quad (\text{asimetrică la dreapta}) \text{ sau}$$

$$H_a: \mu \neq \mu_0 \quad (\text{simetrică}).$$

Ipotezele asimetrice se mai numesc **one-tailed**, iar cea simetrică **two-tailed**.

Testul Z - Inferență asupra mediei unei populații (σ cunoscută)

1. Formulăm **ipoteza nulă**, care susține că media populației ia o anumită valoare:

$$H_0: \mu = \mu_0$$

2. Formulăm **ipoteza alternativă** conform informațiilor obținute din eșantion. Putem avea trei tipuri de ipoteză alternativă

$$H_a: \mu < \mu_0 \quad (\text{asimetrică la stânga}) \text{ sau}$$

$$H_a: \mu > \mu_0 \quad (\text{asimetrică la dreapta}) \text{ sau}$$

$$H_a: \mu \neq \mu_0 \quad (\text{simetrică}).$$

Ipotezele asimetrice se mai numesc **one-tailed**, iar cea simetrică **two-tailed**.

$$(-\infty, -|z^*|] \cup [|z^*|, +\infty) \text{ pentru } H_a \text{ simetrică.}$$

Dacă scorul z nu aparține zonei de respingere spunem că **nu există suficiente dovezi cu nivelul de semnificație α pentru a respinge ipoteza nulă (încercarea de a respinge H_0 eșuează).**

6. Comparăm valoarea critică cu scorul F ; dacă scorul F aparține **zonei de respingere**, atunci acceptăm H_a și respingem H_0 . Zonele de respingere sunt:

$$[F^*, +\infty) \text{ pentru } H_a \text{ asimetrică la dreapta,}$$

$$(0, F_2^*] \cup [F_1^*, +\infty) \text{ pentru } H_a \text{ simetrică}$$

Dacă scorul F nu aparține zonei de respingere spunem că **nu există suficiente dovezi cu nivelul de semnificație α pentru a respinge ipoteza nulă (încercarea de a respinge H_0 eșuează).**

5. Determinăm valoarea critică corespunzătoare

$$t^* = qt(\alpha, df) \text{ pentru } H_a \text{ asimetrică la stânga } (t^* < 0),$$

$$t^* = qt(1 - \alpha, df) \text{ pentru } H_a \text{ asimetrică la dreapta } (t^* > 0),$$

$$t^* = -qt(\alpha/2, df) = qt(1 - \alpha/2, df) \text{ pt. } H_a \text{ simetrică } (t^* > 0).$$

6. Comparăm valoarea critică cu scorul z ; dacă scorul z aparține **zonei de respingere**, atunci acceptăm H_a și respingem H_0 . Zonele de respingere sunt:

$$(-\infty, t^*] \text{ pentru } H_a \text{ asimetrică la stânga,}$$

$$[t^*, +\infty) \text{ pentru } H_a \text{ asimetrică la dreapta,}$$

$$(-\infty, -|t^*|] \cup [|t^*|, +\infty) \text{ pentru } H_a \text{ simetrică.}$$

4. Calculăm **scorul z (statistica testului)**

$$z = \frac{(\bar{x}_{n_1} - \bar{x}_{n_2}) - m_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

5. Determinăm valoarea critică corespunzătoare lui α

$$z^* = qnorm(\alpha) \text{ pentru } H_a \text{ asimetrică la stânga } (z^* < 0),$$

$$z^* = qnorm(1 - \alpha) \text{ pentru } H_a \text{ asimetrică la dreapta } (z^* > 0),$$

$$z^* = -qnorm(\alpha/2) = qnorm(1 - \alpha/2) \text{ pt. } H_a \text{ simetrică } (z^* > 0).$$

Testul T - Inferență asupra mediei unei populații (σ necunoscută)

3. Alegem un nivel de semnificație $\alpha \in \{1\%, 5\%\}$.

4. Calculăm **scorul t (statistica testului)**

$$t = \frac{\bar{x}_n - \mu_0}{s/\sqrt{n}}$$

5. Determinăm valoarea critică corespunzătoare lui α

$$t^* = qt(\alpha, n - 1) \text{ pentru } H_a \text{ asimetrică la stânga } (t^* < 0),$$

$$t^* = qt(1 - \alpha, n - 1) \text{ pentru } H_a \text{ asimetrică la dreapta } (t^* > 0),$$

$$t^* = -qt(\alpha/2, n - 1) = qt(1 - \alpha/2, n - 1) \text{ pt. } H_a \text{ simetrică } (t^* > 0)$$

3. Alegem un nivel de semnificație $\alpha \in \{1\%, 5\%\}$.

4. Calculăm **scorul z (statistica testului)**

$$z = \frac{\bar{x}_n - \mu_0}{\sigma/\sqrt{n}}$$

5. Determinăm valoarea critică corespunzătoare lui α

$$z^* = qnorm(\alpha) \text{ pentru } H_a \text{ asimetrică la stânga } (z^* < 0),$$

$$z^* = qnorm(1 - \alpha) \text{ pentru } H_a \text{ asimetrică la dreapta } (z^* > 0),$$

$$z^* = -qnorm(\alpha/2) = qnorm(1 - \alpha/2) \text{ pentru } H_a \text{ simetrică } (z^* > 0)$$