

# Tipuri de distribuții și compararea statistică

# Forma distribuțiilor de frecvențe

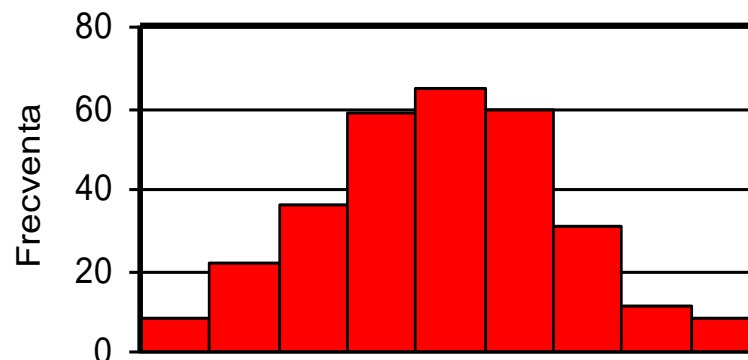
⊕ Simetrice (normale)

⊕ Deviate

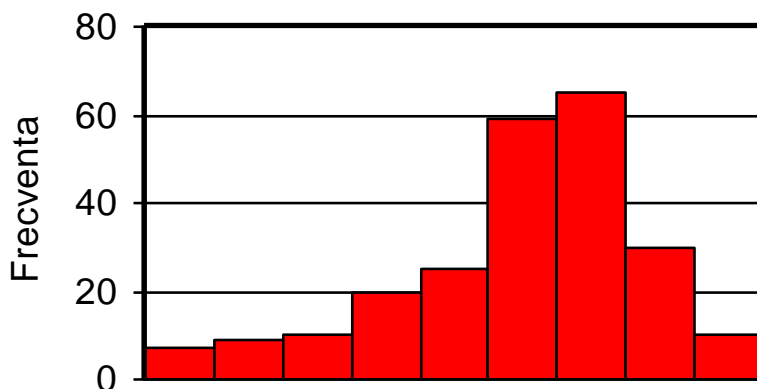
- pozitiv
- negativ

Majoritatea distribuțiilor întâlnite în practica medicală sunt simetrice sau deviate pozitiv

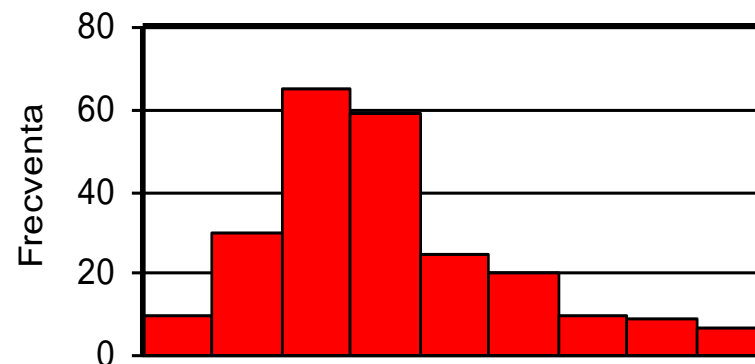
# Distribuții de frecvențe



**Normală**

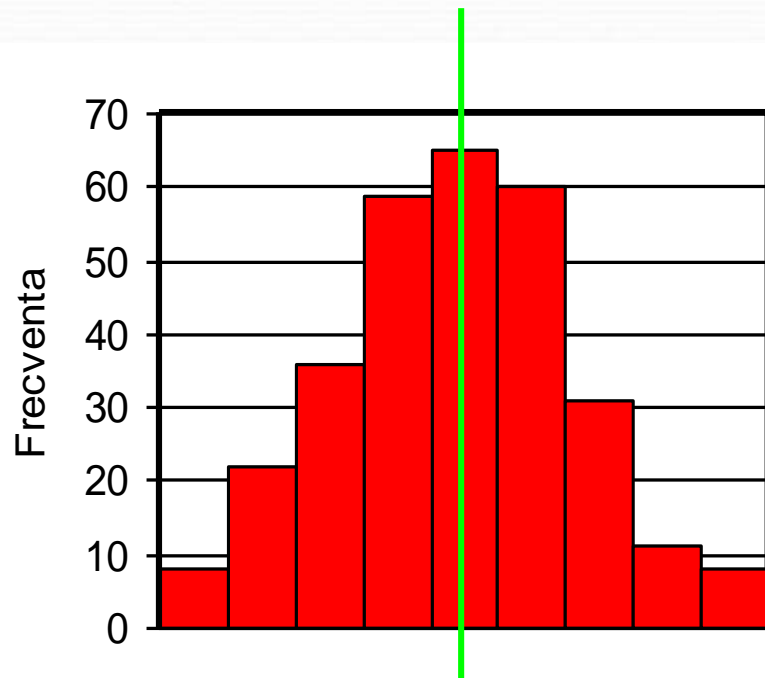


**Deviată negativ**

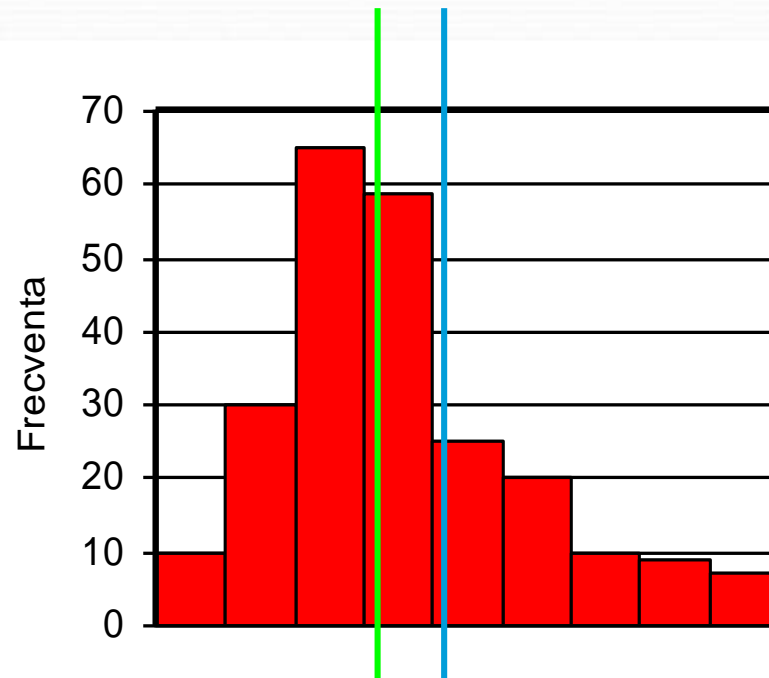


**Deviată pozitiv**

# Localizarea Mediei și Mediane



**media = mediana**



**mediana      media**

**mediana  $\neq$  media**

# Medie sau Mediană?

Când distribuția frecvențelor setului de date este:

- ✚ *simetrică* - *media* este în mod obișnuit cea mai potrivită unitate de măsură.
- ✚ *deviată* - *mediana* este în mod obișnuit cea mai potrivită unitate de măsură.

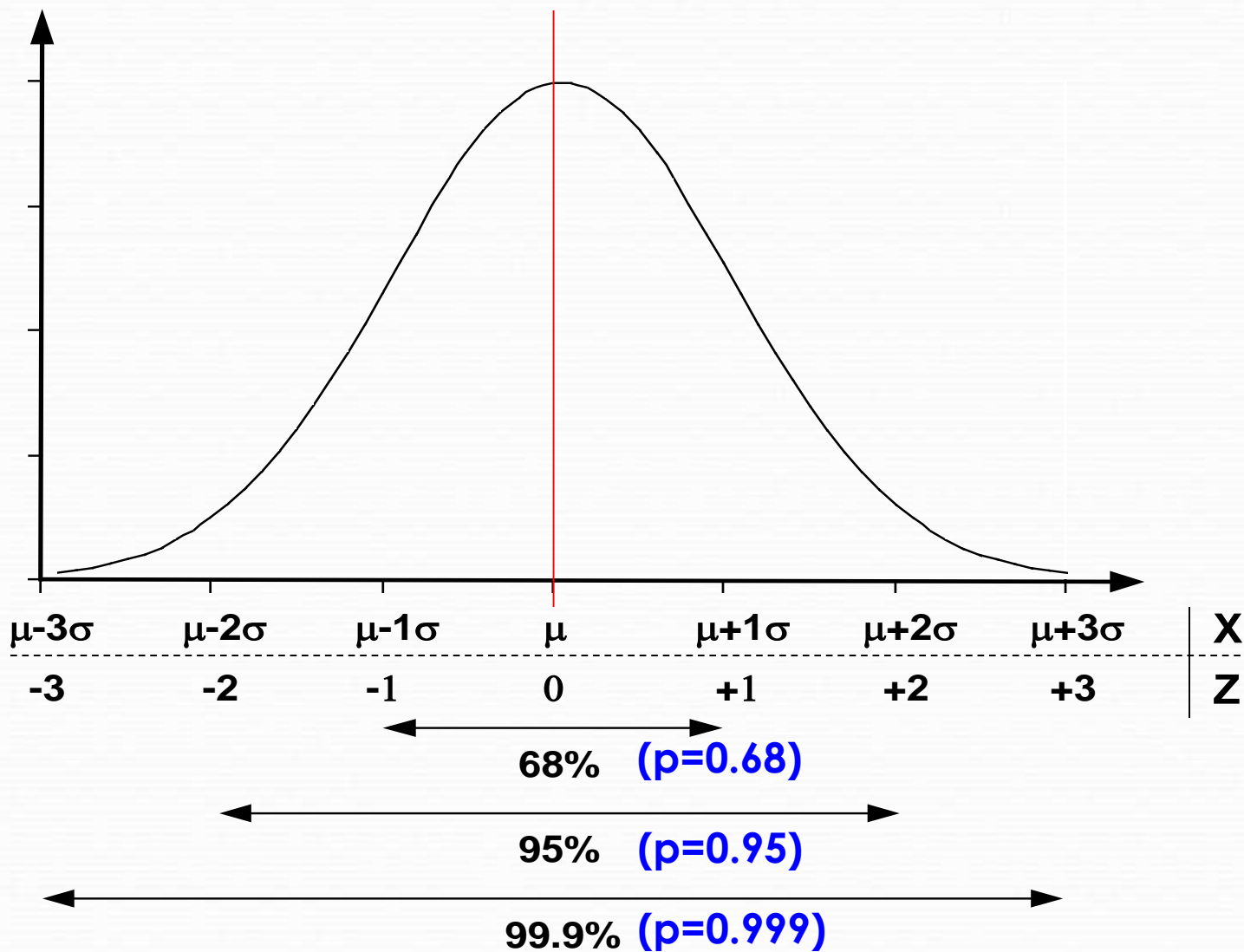
# Distribuția normală(I)

- ❖ Numită și distribuție Gaussiană este una din **cele mai importante** distribuții în statistică.
- ❖ Distribuțiile **multor** măsurători medicale aproximează distribuția normală **ex.** Tensiunea arterială, Colesterolul sanguin.
- ❖ Se definește printr-o medie și o variantă și se caracterizează ca fiind în formă de clopot și simetrică în jurul mediei.
- ❖ Va fi **mai plată dacă varianța este mai mare** și mai ascuțită dacă varianța este mai mică.

# Distribuția normală(II)

- ❖ Probabilitatea ca o variabilă aleatoare Normală să ia o valoare între:
  - medie și 1 deviație standard de ambele părți este 0.68,
  - medie și 1.96 x deviația standard de ambele părți este 0.95
  - medie și 2.58 x deviația standard de ambele părți este 0.99

# Distribuția normală



# Distribuția normală

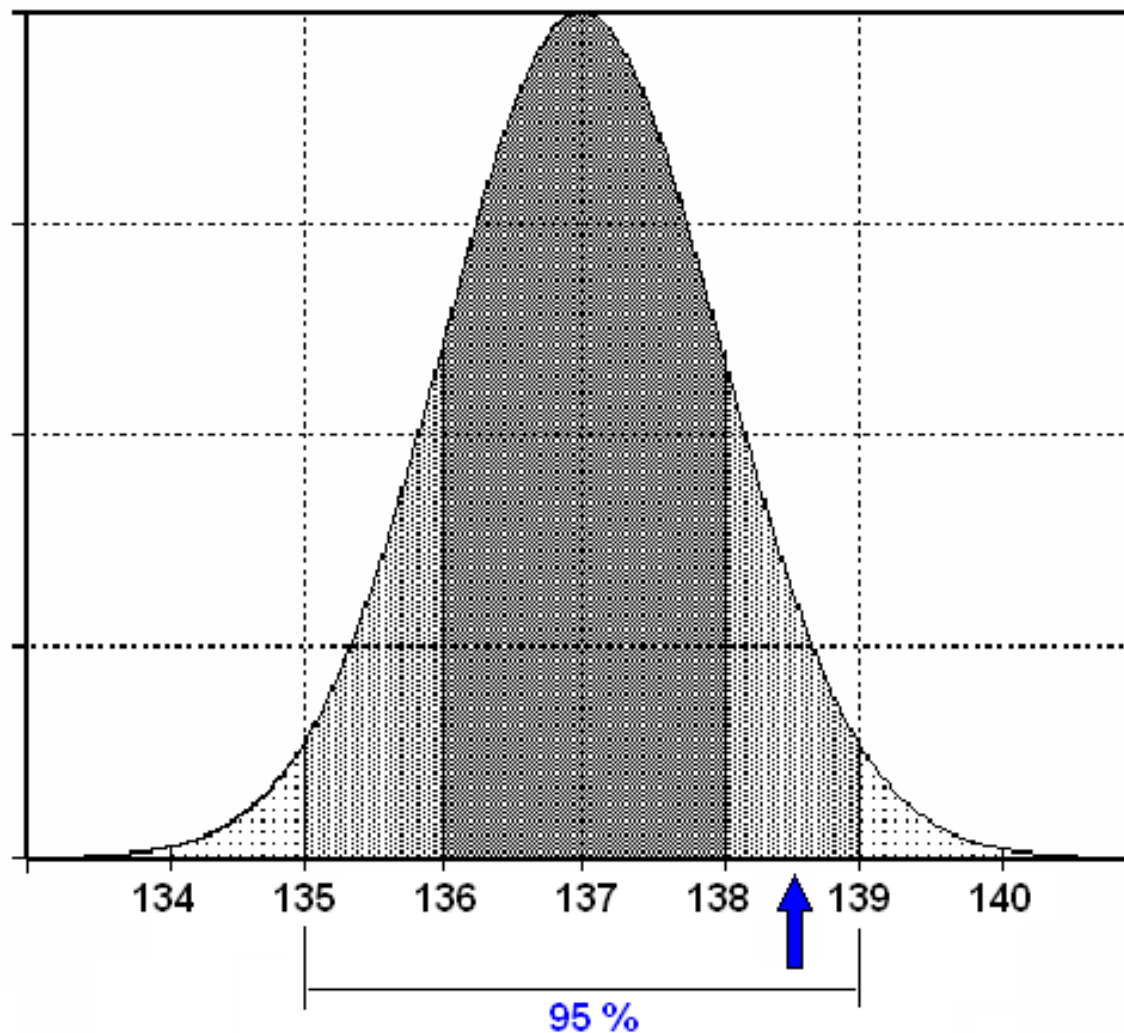
- ❖ Distribuția normală apare pentru date continue când **media și varianța populației sunt cunoscute**.
- ❖ Când acestea **nu sunt cunoscute** și când avem doar informații obținute pe eșantion despre medie și varianță atunci folosim **distribuția testului t Student**. Pe măsură ce mărimea eșantionului crește, el tinde spre o distribuție normală.
- ❖ **De fapt**, multe distribuții tind spre Normalitate, mai ales dacă se strâng numeroase eșantioane. De aceea, distribuția Normală a devenit o parte importantă a teoriei statistice.

# DIFERENTE SEMNIFICATIVE SI NESEMNIFICATIVE

a) Exemplu:

- |                           |             |
|---------------------------|-------------|
| • BAIETI                  | FETE        |
| • $n = 25$                | $n = 25$    |
| • $X = 137 \text{ cm}$    | $X = 138.5$ |
| • $s = 5 \text{ cm}$      | $s = 5$     |
| • $s_x = 1 \text{ cm}$    | $s_x = 1$   |
| • $(135, 139) \dots 95\%$ |             |

## g) Distribuția mediilor esantioanelor



# DIFERENTE SEMNIFICATIVE SI NESEMNIFICATIVE

a) Exemplu:

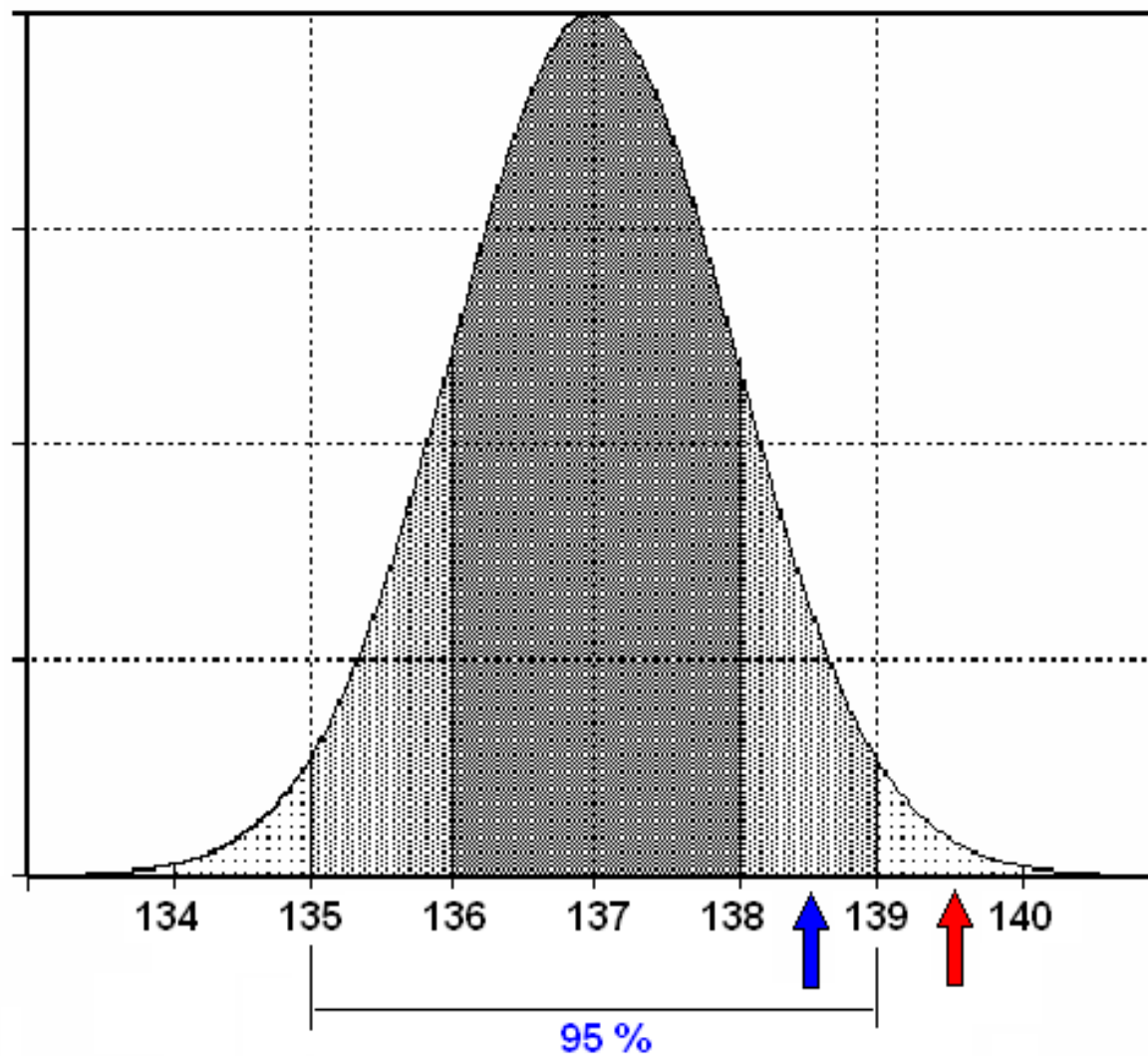
- |                       |                             |
|-----------------------|-----------------------------|
| • BAIETI              | FETE                        |
| • $n = 25$            | $n = 25$                    |
| • $X = 137$ cm        | $X = 138.5$                 |
| • $s = 5$ cm          | $s = 5$                     |
| • $s_x = 1$ cm        | $s_x = 1$                   |
| • $(135, 139)$ ...95% | <b>dif. nesemnificative</b> |

# DIFERENTE SEMNIFICATIVE SI NESEMNIFICATIVE

a) Exemplu:

- |                       |              |             |
|-----------------------|--------------|-------------|
| • BAIETI              | FETE         |             |
| • $n = 25$            | $n = 25$     |             |
| • $X = 137$ cm        | $X = 138.5$  | $X = 139.5$ |
| • $s = 5$ cm          | $s = 5$      |             |
| • $s_x = 1$ cm        | $s_x = 1$    |             |
| • $(135, 139)$ ...95% | dif. nesemn. |             |

## g) Distributia mediilor esantioanelor



# DIFERENTE SEMNIFICATIVE SI NESEMNIFICATIVE

a) Exemplu:

- |                       |             |                    |
|-----------------------|-------------|--------------------|
| • BAIETI              | FETE        |                    |
| • $n = 25$            | $n = 25$    |                    |
| • $X = 137$ cm        | $X = 138.5$ | $X = 139.5$        |
| • $s = 5$ cm          | $s = 5$     |                    |
| • $s_x = 1$ cm        | $s_x = 1$   |                    |
| • $(135, 139)$ ...95% | d. nesemn.  | dif. semnificative |

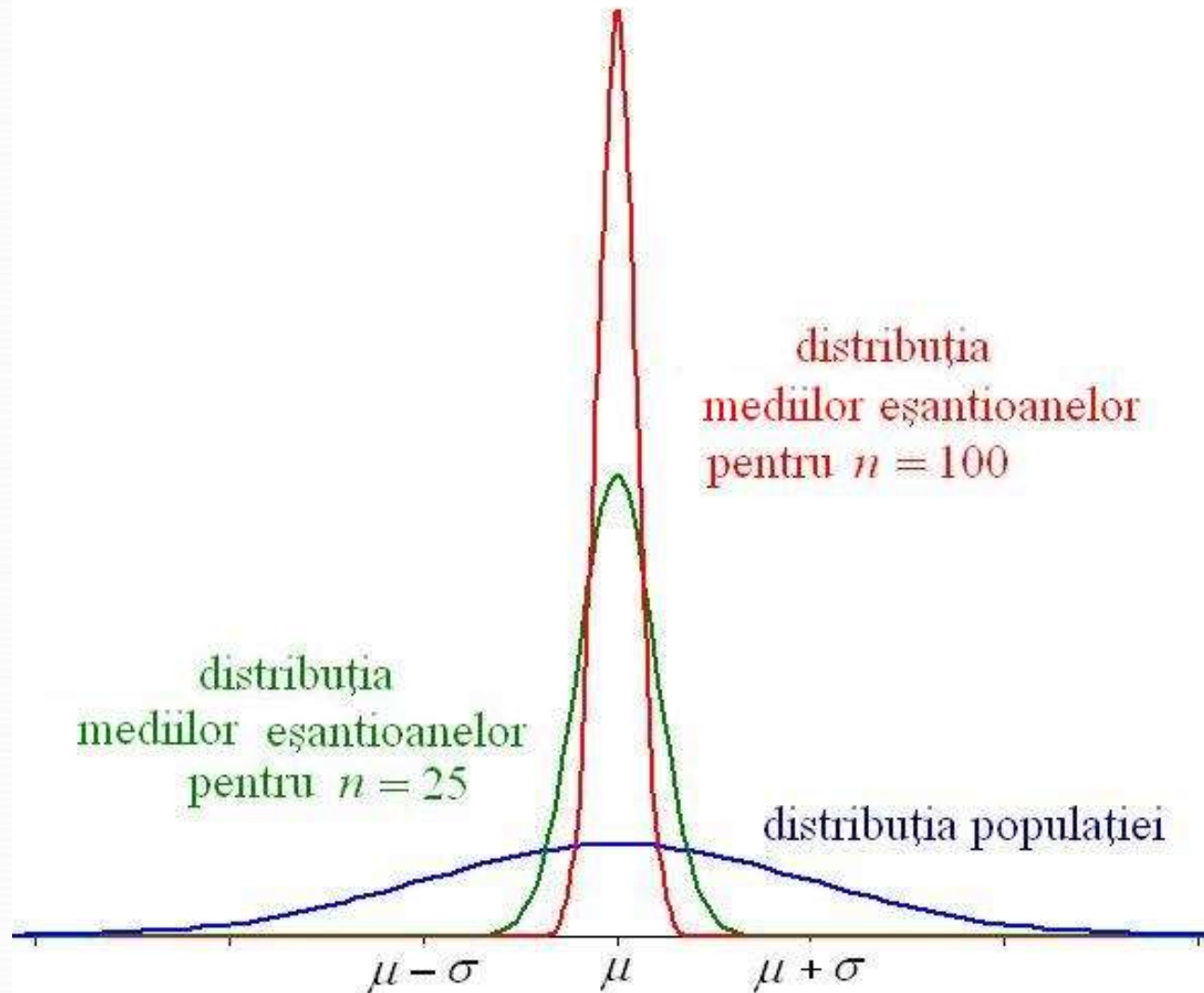
## **b) DEFINITII**

- **DIFERENTE NESEMNIFICATIVE**
- **Au probabilitate mare sa apara din intamplare**
- **Cauza: Variabilitatea de esantionare**
- **Cele doua esantioane apartin aceleiasi populatii**
  
- **DIFERENTE SEMNIFICATIVE**
- **Au probabilitate mica sa apara din intamplare**
- **Trebuie atribuite altei cauze**

# Teorema limită centrală (1)

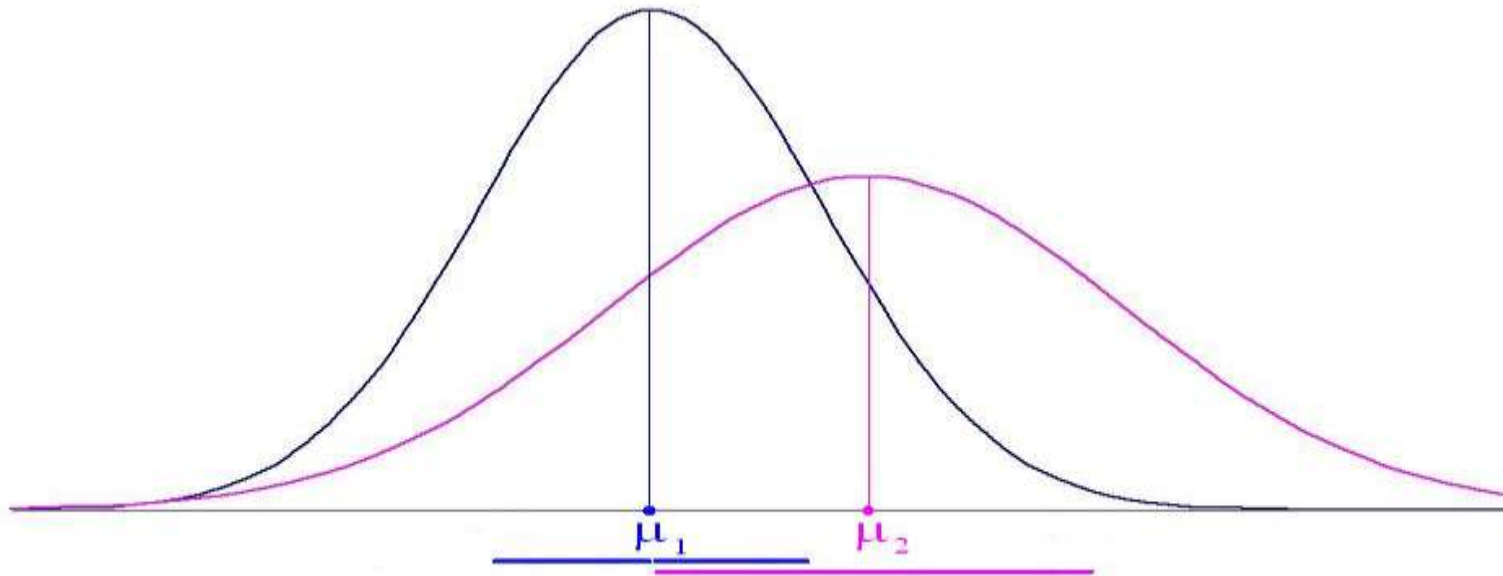
- Dacă se extrag eșantioane de volum  $n$  dintr-o populație, atunci pentru valori “mari” ale lui  $n$  mediile de eșantion sunt distribuite (aproximativ) normal.
- În caz că  $X$  are o distribuție normală  $N(\mu, \sigma^2)$ , atunci
$$M = (X_1 + X_2 + \dots + X_n) / n$$
are o distribuție normală  $N(\mu, \sigma^2/n)$ .
- Iar dacă variabila aleatoare  $X$  este distribuita aproximativ normal, atunci  $M$  va fi distribuita normal chiar și pentru valori “mici” ale lui  $n$ .

# Teorema limită centrală (2)



# Compararea populațiilor (1)

- Atunci când comparăm între ele două populații distribuite normal, comparația se poate face la nivelul mediilor  $\mu_1$  și  $\mu_2$  și/sau la nivelul varianțelor  $\sigma_1^2$  și  $\sigma_2^2$ .



**Diferă populațiile între ele?**

# Testarea statistică de semnificație (1)

- Este o metodă de stabilire a gradului de plauzibilitate (“adevărului”) al unei afirmații
- Explicații necesare:

testare = procedură standard de examinare, aplicată nediscriminatoriu tuturor indivizilor din eșantion

statistică = formulă de calcul cu datele provenite din eșantion

semnificație = înțeles “precis”

1 este “semnificativ” diferit de 0 ?

Da, la nivelul km. Nu, la nivelul  $\mu\text{m}$ . Răspuns subiectiv.

# Testarea statistică de semnificație (2)

- Robert A. Fisher (1925):  
validitatea unei ipoteze științifice este stabilită pe baza unui singur test, cu opțiunea de a nu emite o judecată definitivă atunci când rezultatul testului nu este “suficient de limpede”.

Sunt posibile doar două opțiuni:

fie vom “respinge ipoteza nulă”,

fie vom amâna decizia (nu sunt suficiente date pentru a trage vreo concluzie).

# Testarea statistică de semnificație (3)

- Ipotezele științifice se refera la populații teoretice, care au de obicei un număr infinit de indivizi și sunt reprezentate de distribuții continue.
- Ipoteza științifică este înlocuită printr-o ipoteză statistică, exprimată prin intermediul parametrului acelei populații (cum ar fi proporția, media ...)
- Valoarea parametrului este estimată prin exploatarea datelor obținute dintr-un eșantion extras din populație, apoi este comparată cu o valoare “așteptată”.
- Discrepanța dintre cele două va influența “credința” noastră în validitatea ipotezei științifice.

# Testarea statistică de semnificație (4)

- Ideea testării statistice de semnificație este simplă: ipoteza statistică va servi ca alternativa la o alta ipoteza – așa-numita “ipoteză nula” – care este luata în considerare doar pentru a fi respinsă.

Prin acceptarea inițială a adevărului ipotezei nule vor rezulta anumite consecințe logice/statistice, iar acestea vor fi confruntate cu datele observate.

Orice dovadă aflată în contradicție cu ipoteza nulă va servi ca justificare a alternativei.

# Testarea statistică de semnificație (5)

- Ipoteza alternativă trebuie să exprime o inegalitate sau o discrepanță (în niciun caz o egalitate/coincidență)
- Ipoteza nulă exprimă, în aceiași termeni ca și alternativa, o egalitate.
- Ipoteza alternativă este cea pe care o dorim confirmată ca adevărată
- Notății tradiționale:  $(H_1)$  sau  $(H_a)$  alternativa  
 $(H_o)$  ipoteza nulă

# IPOTEZE STATISTICE

- **a) IPOTEZA NULA (DE ZERO)**

- $H_0 : X_1 = X_2$
- Nu sunt diferite semnificative între cele două valori (esantioane)
- Semnul “=” nu este în sens matematic, ci statistic

- **b) IPOTEZE ALTERNATIVE**

- $H_1 : X_1 \neq X_2$  (bilaterală)
- $X_1 > X_2$  ,  $X_1 < X_2$  (unilaterală)

## ● PRAG DE SEMNIFICATIE

- a) DEFINITIE:

- valoarea probabilitatii sub care incepem sa consideram diferentele ca fiind semnificative

- b) VALOARE UZUALA:

- $\alpha = 0.05 = 5 \%$

- c) NIVEL DE INCREDERE (confidenta)

- $1 - \alpha = 0.95 = 95 \%$

- 1.4. COEFICIENTUL  $p$

$p$  = probabilitatea ca diferentele observate sa fi aparut din intamplare

## 1.5. DECIZIA

- *Daca  $p > 0.05 \Rightarrow$  Diferente Nesemnificative, (N),  $H_0$  este acceptata*

-----

- *Daca  $p < 0.05 \Rightarrow$  Diferente Semnificative, (S),  $H_0$  este respinsa*
- *Daca  $p < 0.01 \Rightarrow$  Diferente Foarte Semnificative, (F),  $H_0$  respinsa*
- *Daca  $p < 0.001 \Rightarrow$  Diferente Extrem de Semnificative, (E),  $H_0$  respinsa*

# Testarea statistică de semnificație (6)

- Exemple de ipoteze științifice:

Vârsta “foarte înaintată” este un predictor semnificativ al maladiei Alzheimer

Medicamentul A ajută pacienții să se însănătoșească mai bine decât medicamentul B

Pacienții își revin în urma unui tratament standard

- Exemplu “negativ”:

Efectele medicamentului E asupra pacienților bărbați și femei sunt similare.

# Testarea statistică de semnificație (7)

- Ipotezele statistice corespunzătoare:

Incidența maladiei Alzheimer este mai mare la persoanele de vârstă foarte înaintată (prin comparație cu persoanele de vârstă înaintată)

Proporția pacienților însănătoșiți dintre cei tratați cu medicamentul A este mai mare decât proporția corespunzătoare pentru medicamentul B

Starea medie de sanatate a pacienților, în urma unui tratament standard, este mai bună decât înaintea începerii tratamentului

# Testarea statistică de semnificație (8)

- Ipotezele nule corespunzătoare:

Incidența maladiei Alzheimer la persoanele de vârstă foarte înaintată este aceeași cu cea la persoanele de vârstă înaintată

Proporția pacienților însănătoșiți dintre cei tratați cu medicamentul A este egală cu cea corespunzătoare pentru medicamentul B

Starea medie de sănătate a pacienților, în urma unui tratament standard, nu suferă nicio schimbare

Conform lui R. A. Fisher, ipoteza nulă este “ridicăată” – ca un complement al ipotezei alternative – doar pentru a fi respinsă, iar prin respingerea ei vom accepta ca „adevărată” ipoteza științifică inițială

# T. S. S. – abordarea clasica (1)

O testare statistica a semnificatiei se efectuează în cinci pași consecutivi:

- Pasul 1: Specificăm ipoteza alternativă, apoi ipoteza nulă.
- Pasul 2: Alegem statistica adaptată situației concrete.
- Pasul 3: Alegem nivelul de semnificație, și pe baza sa calculăm pragul de separare (între valorile “acceptabile” și cele considerate ca “inacceptabile”).
- Pasul 4: Calculăm valoarea statisticii, folosind efectiv datele din eșantion (ales aleator).
- Pasul 5: Decidem, prin compararea valorii calculate cu pragul dat de nivelul de semnificație, dacă să respingem sau nu ipoteza nulă.

## T. S. S. – abordarea clasica (2)

- Discuția în jurul testării statistice de semnificație începe cu ultimul pas.
- În acesta un decidentul (d-voastră) va trebui fie să respingă ipoteza nulă  $H_0$   
(și prin urmare să accepte ipoteza alternativă  $H_a$ ),  
fie să nu respingă pe  $H_0$ .
- În realitate  $H_0$  este fie adevărată, fie falsă – dar decidentul nu cunoaște situația reală.
- Patru posibilități:

# T. S. S. – abordarea clasica (3)

**Realitatea** (necunoscută)

	Realitatea (necunoscută)	
	$H_0$ este falsa	$H_0$ este adevarata
Decizia	Respingem $H_0$	Eronată (eroare de tipul I)
	Nu respingem $H_0$	Corectă!

# T. S. S. – abordarea clasica (4)

- În testarea statistică de semnificație importanța maximă o are eroarea de tipul I. Probabilitatea ei, numărul

$$\alpha = \text{Prob}(\text{decizie eronată} \mid H_0 \text{ este adevărată})$$

este **nivelul de semnificație** a cărui valoare a fost aleasă anterior (la Pasul 3).

Fiecare decident dorește să păstreze nivelul de semnificație cât mai mic posibil – întrucât este de fapt probabilitatea de a face o eroare!

Astfel că valori cum este 0,05 sunt destul de des întâlnite, iar în științele medicale se recomandă alegerea unor valori mai mici, de exemplu 0,001.

# Exemplu (1)

( $H_a$ ) Ritmul cardiac mediu al pacienților hipertensivi scade în urma administrării medicamentului,

( $H_o$ ) Ritmul cardiac mediu al pacienților hipertensivi, în urma administrării medicamentului, nu suferă nici o schimbare.

Exprimăm formal ipotezele de mai sus astfel:

$$(\mathbf{H}_a) \quad \mu_a < \mu_b$$

$$(\mathbf{H}_o) \quad \mu_a = \mu_b$$

unde  $\mu_b$ , respectiv  $\mu_a$  reprezintă ritmul cardiac mediu înainte, respectiv după administrarea medicamentului.

## Exemplu (2)

Datele pe care le obținem apar în mod natural împerechiate; mai precis, pentru fiecare pacient măsurăm ritmul cardiac înainte (  $x_b$  ) și după (  $x_a$  ) administrarea medicamentului.

Evident, am putea calcula diferența  $d = x_b - x_a$  și am putea considera că medicamentul este:

eficace pentru pacientul nostru dacă  $d > 0$  ,

ineficace dacă  $d = 0$  (nu se constată schimbare)

și dăunător dacă  $d < 0$  .

## Exemplu (3)

Să notăm cu  $\delta$  diferența medie; atunci testarea statistică de semnificație de mai sus este înlocuită prin

$$(H_a) \quad \delta > 0$$

$$(H_o) \quad \delta = 0$$

Dacă presupunem că ritmul cardiac al pacienților hipertensivi, și înainte, și după administrarea medicamentului, este distribuit normal, atunci rezultă că diferențele  $d$  sunt și ele distribuite normal.

Statistica adaptată situației este 
$$t = \frac{m}{s / \sqrt{n}}$$

în care  $m$  este media eșantionului diferențelor,  $s$  este abaterea standard,  $n$  este volumul eșantionului.

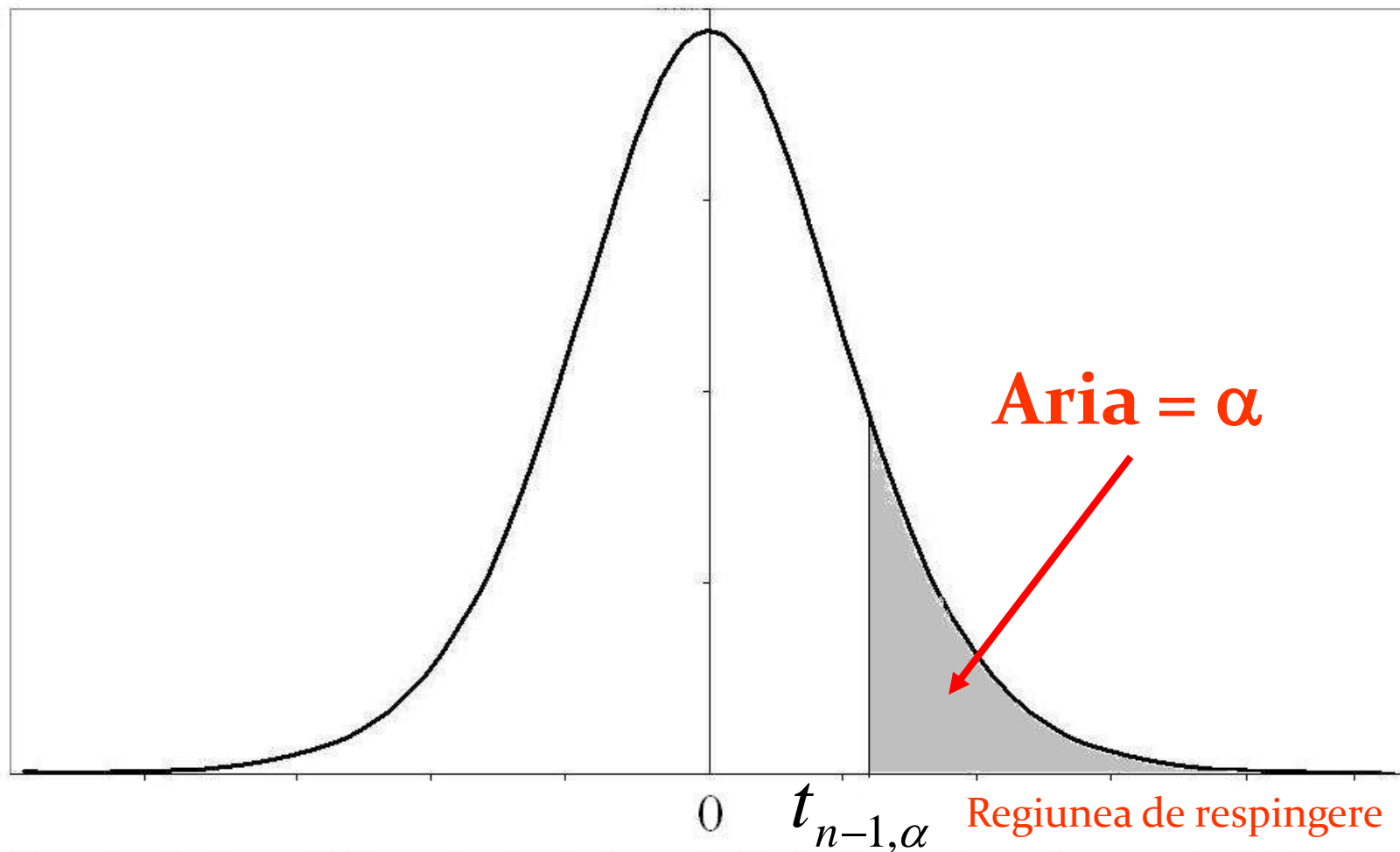
## Exemplu (4)

Știm că valorile  $t$  sunt distribuite Student, cu  $n-1$  grade de libertate.

Marea majoritate a lor se plasează în jurul lui 0. Medicamentul fiind hipertensiv, ne așteptăm ca diferențele  $d$  să fie pozitive, deci valorile lui  $m$  (și implicit ale lui  $t$ ) să fie pozitive.

Alegem nivelul de semnificație  $\alpha = 0.10$ . Ca rezultat vom obține o valoare prag  $t_{n-1,\alpha}$

## Exemplu (5)



## Exemplu (6)

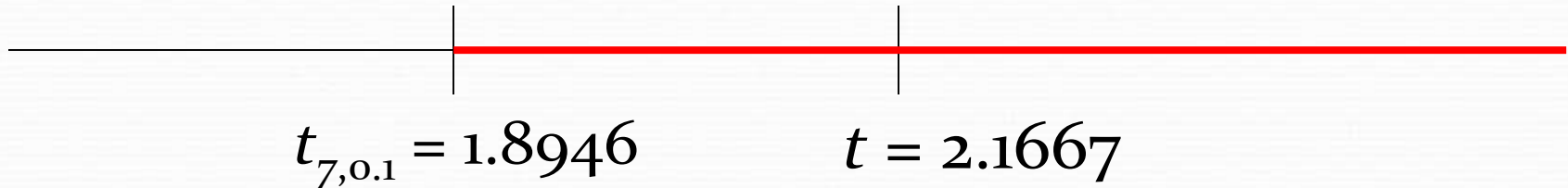
Pentru  $n = 8$  avem  $t_{n-1,\alpha} = t_{7,0.1} = 1.8946$   
Folosim acum datele obținute dintr-un eșantion (8 indivizi):

Înainte (b/m)	După (b/m)	Diferența
66	58	8
69	65	4
...	...	...
73	66	7

Calculăm  $m = 3.375$ ,  $s = 4.406$ ,  $t = 2.1667$ .

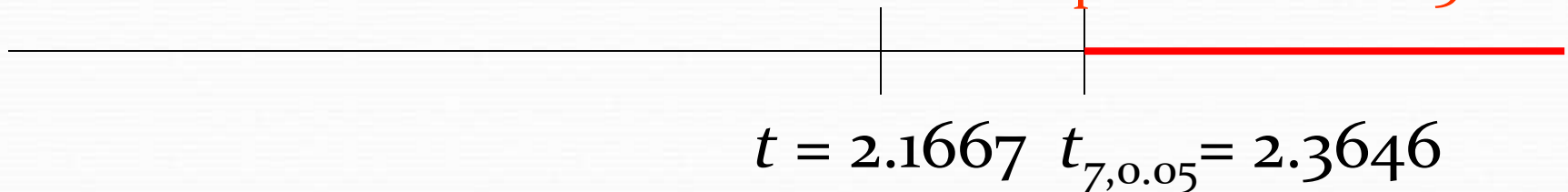
## Exemplu (7)

Regiunea de respingere pentru  $\alpha = 0.1$



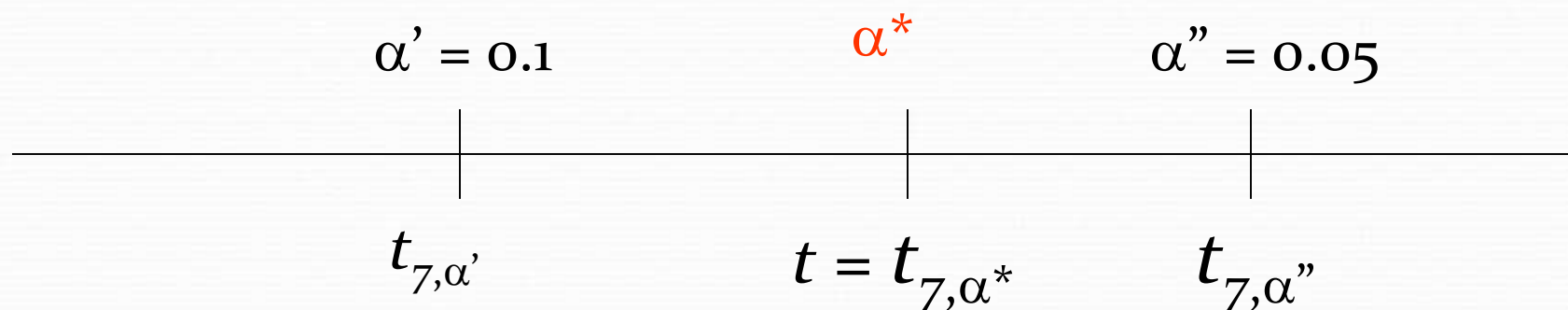
Concluzie: respingem ipoteza nulă!

Regiunea de respingere  
pentru  $\alpha = 0.05$



Concluzie: nu putem respinge ipoteza nulă!

## Exemplu (8)



Nivelul de semnificație particular  $\alpha^*$  este cunoscut ca valoarea p a ipotezei alternative.

Este cel mai mic nivel de semnificație care ne permite să acceptăm ca adevărată ipoteza alternativă – prin respingerea ipotezei nule – bazându-ne doar pe datele din eșantionul ales.

Este riscul de a accepta ca adevărată ipoteza alternativă – ceea ce dorim de fapt – atunci când ipoteza nula este cea adevărată.

# T. S. S. – abordarea moderna

O testare statistica de semnificatie se efectuează astăzi în doi pași:

- Pasul 1: Specificăm ipoteza alternativă, apoi ipoteza nulă.
- Pasul 2: Calculăm valoarea p a ipotezei alternative, apoi interpretăm aceasta ca risc.

Pentru calculul valorii p se utilizează în Excel funcția

`TTEST(domeniu1, domeniu2, lateralitate, tip)`

tip = 1 înseamnă test pereche

# $\chi^2$ ca test statistic de semnificație (1)

- Compară două variabile aleatoare (cu câte 2 valori fiecare) și le evaluează independența statistică.
- De exemplu, variabila  $X$  este “Medicamentul”  
iar valorile sale sunt “MedA”, “MedB”;  
pe de altă parte, variabila  $Y$  este “Starea  
pacientului”,  
cu valorile “vindecat”, “nevindecat”.
- Independența statistică a lui  $X$  și  $Y$  corespunde dependenței liniare a liniilor (sau a coloanelor) din tabelul extins de contingență.

## $\chi^2$ ca test statistic de semnificație (2)

- Ipoteza alternativă este că există o asocierie între variabile.
- Ipoteza nulă este că nu.
- În Excel, calculul valorii p se face cu funcția  
 $\text{CHITEST}(\text{observate}, \text{așteptate})$
- Folosirea funcției este precedată de calculul valorilor “așteptate” într-un tabel de contingență “paralel” cu cel ce conține datele inițiale.

# Testul hi-pătrat: exemplu (1)

- Datele în tabelul de contingență

	Medicament A	Medicament B (Placebo)
Vindecați	55	40
Nevindecați	25	35

- Ipoteza alternativă: există o asocierie între variabilele “Starea de sănătate” și “Tip medicament”
- Ipoteza nulă: nu există asocierie (variabilele sunt statistic independente)

# Testul hi-pătrat: exemplu (2)

■ Tabelul extins cu totaluri:

	Medic. A	Placebo	Total rânduri
Vindecați	55	40	95
Nevindecați	25	35	60
Total coloane	80	75	155

# Testul hi-pătrat: exemplu (3)

- Tabelul “paralel” (calculat în situația ipotezei nule adevărate) cu datele obținute din totaluri:

	Medicament A	Medicament B (Placebo)
Vindecați	49.03	45.97
Nevindecați	30.98	29.03

$$49.03 = \frac{total\_r\acute{a}nd \times total\_coloana}{total\_general} = \frac{95 \times 80}{155}$$

# Testul hi-pătrat: exemplu (4)

## ■ Statistica:

$$X^2 = \sum_{\text{celule}} \frac{(\text{observ} - \text{astept})^2}{\text{astept}}$$

are o distribuție aproximativ  $\chi^2$  ((nr.rânduri-1)(nr.coloane-1))

## ■ Valoarea p calculată cu funcția CHITEST:

0.0489

Concluzie: asocierea este “semnificativă” (clasic)

sau

Există o asociere semnificativă între “Tipul medicamentului” și “Starea de sănătate”

# Testul hi-pătrat: corecția Yates

Testul dă rezultate “corecte” în caz că toate valorile așteptate sunt  $> 5$ . Altfel, valoarea  $p$  este incorectă.

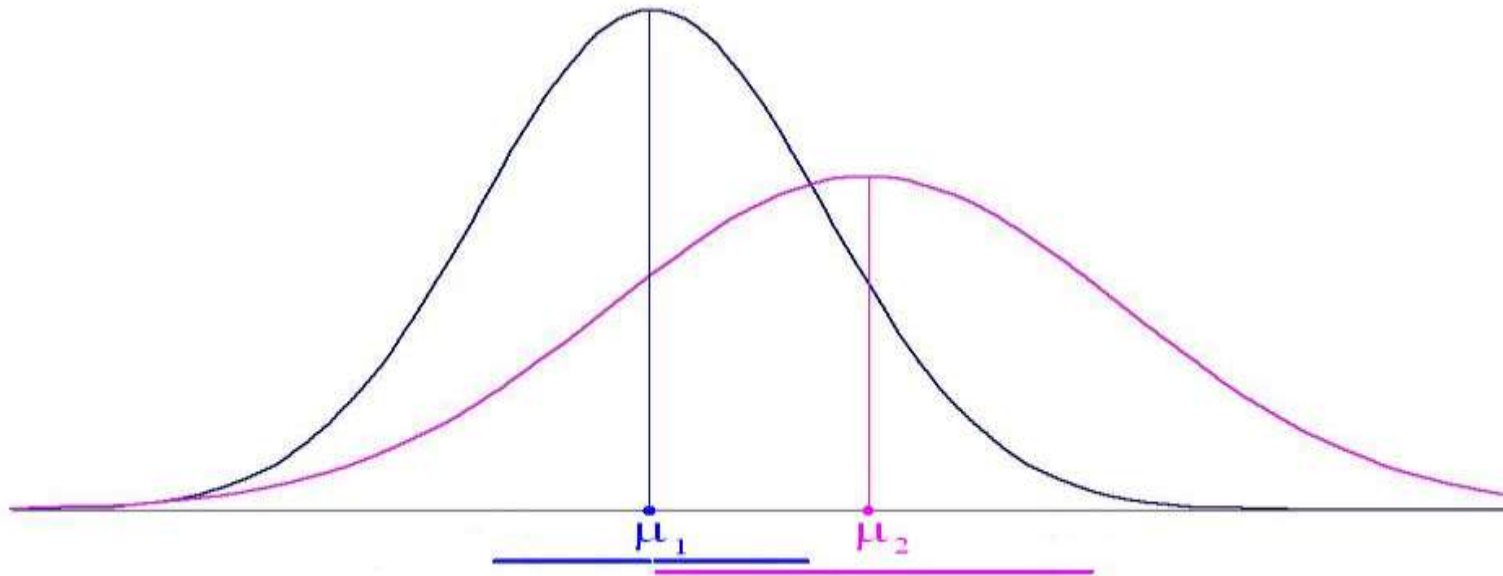
Varianta cu corecția Yates:

$$X^2 = \sum_{\text{celule}} \frac{(\text{observ} - \text{astept} - 0.5)^2}{\text{astept}}$$

produce valori  $p$  ușor mai mari, dar mai corecte.

# Compararea populațiilor (1)

- Atunci când comparăm între ele două populații distribuite normal, comparația se poate face la nivelul mediilor  $\mu_1$  și  $\mu_2$  și/sau la nivelul varianțelor  $\sigma_1^2$  și  $\sigma_2^2$ .



Diferă populațiile între ele?

# Compararea populațiilor (2)

- Nu cunoaștem nici mediile  $\mu_1$  și  $\mu_2$  și nici varianțele  $\sigma_1^2$  și  $\sigma_2^2$ . Ele sunt estimate, pe baza a două eșantioane de volume  $n_1$  și respectiv  $n_2$ , prin mediile de eșantion  $m_1$  și  $m_2$
- Să presupunem că  $m_2 > m_1$ . Are sens să considerăm ipoteza alternativă

$$(H_a) \mu_2 > \mu_1$$

și să-i calculăm valoarea p.

Reamintim că valoarea p este riscul de a accepta ipoteza alternativă atunci când de fapt ipoteza nulă este adevărată.

Valoarea p o putem calcula imediat, în Excel, cu TTEST.

# Compararea populațiilor (3)

- Ne-am pus problema comparării a două populații luând în considerare mediile (sau proporțiile, eventual varianțele lor). Am luat în considerare parametrii care determină populațiile: mediile ( $\mu$ ), proporțiile ( $\pi$ ), varianțele ( $\sigma^2$ ).
  - Multe dintre metodele de comparație se bazează pe ipoteza “fundamentală” că anumite variabile sunt distribuite (aproximativ) normal. Acestea sunt cunoscute sub numele de teste parametrice.
  - Există însă situații în care:
    - fie nu cunoaștem deloc felul în care sunt distribuite variabilele,
    - fie normalitatea distribuției este încălcată flagrant.
- În asemenea situații, pentru compararea populațiilor este posibil să folosim teste care nu presupun nimic despre tipul de distribuție, cu alte cuvinte **teste neparametrice**.

# Compararea populațiilor (4)

- În **testele de rang** valorile numerice ale variabilelor – obținute din eșantion – sunt înlocuite prin rangurile lor.
- Să prezentăm, ca exemplu, testul Wilcoxon.
- Ipoteza alternativă de la care plecăm, într-o exprimare generală, este următoarea:

( $H_a$ ) distribuția valorilor variabilei numerice (care ne interesează) este asimetrică în raport cu o.

Îi vom opune ipoteza nulă:

( $H_o$ ) distribuția valorilor variabilei numerice este simetrică în raport cu o.

# Testul Wilcoxon (1)

- Conform teoriei generale, vom încerca să “deducem” consecințe logice ale acceptării adevărului ipotezei nule, apoi să evaluăm dacă datele provenite din eșantion sunt “compatibile” cu consecințele.
- Să începem prin a analiza datele numerice  $x_1, x_2, \dots, x_n$  provenite dintr-un eșantion de volum  $n$ . Unele valori vor fi pozitive, altele vor fi negative, câteva chiar egale cu 0. Să presupunem că  $m$  dintre ele sunt nenule.
- Vom ordona crescător valorile nenule, luate în modul, apoi le vom înlocui cu rangurile lor:

$$|x_{(1)}| \leq |x_{(2)}| \leq \dots \leq |x_{(m)}|$$

# Testul Wilcoxon (2)

- Să notăm cu  $T_+$  suma rangurilor valorilor pozitive, și cu  $T_-$  suma rangurilor valorilor negative. Dacă acceptăm ideea că ipoteza nulă este adevărată, atunci  $T_+$  și  $T_-$  n-ar trebui să difere prea mult între ele. Pe de altă parte, suma lor ar trebui să fie egală cu suma tuturor rangurilor. Ar trebui să ne așteptăm ca atât  $T_+$  cât și  $T_-$  să fie apropiate de  $\frac{m(m+1)}{4}$ .
- Cu cât  $T_+$  diferă mai mult de  $\frac{m(m+1)}{4}$ , cu atât ipoteza nulă devine mai implauzibilă și drept urmare vom fi înclinați să acordăm credit alternativei ( $H_a$ ).

# Testul Wilcoxon (3)

- Calculul valorii  $p$  a ipotezei alternative se bazează pe faptul că statistica

$$\frac{T_+ - m(m+1)/4}{\sqrt{m(m+1)(2m+1)/24}}$$

este distribuită (cel puțin pentru valori “mari” ale lui  $n$ ) aproximativ normal standard.

Testul Mann-Whitney este în esență echivalent cu testul Wilcoxon.

# Testul Wilcoxon (4)

- Exemplu:

- Comanda

**MEANS valori semn**

din **EpiInfo** da

valoarea  $p = 0.0139$

	A	B	C
1	valori	semn	rang
2	1.9	+	4.5
3	1.9	+	4.5
4	3.5	+	7
5	0.5	+	2
6	-3.6	-	8
7	-5.2	-	9
8	0.6	+	3
9	-2.1	-	6
10	-0.4	-	1

# Analiza varianței - ANOVA (1)

- Deseori se pune problema comparării prin medii a mai mult de două populații, sau a unei populații stratificate în mai mult de două straturi.

În asemenea situații se poate aplica o generalizare a testului  $t$  pentru două populații, cunoscută sub numele de **analiza varianței** sau testul ANOVA.

- Prima aplicare a analizei varianței s-a făcut într-o situație în care se analizau recoltele obținute în urma tratării solului cu diferite feluri de îngrășăminte. Se păstrează, tradițional, unele dintre notațiile/noțiunile folosite atunci (cum este “media tratamentului”).

# Analiza varianței - ANOVA (2)

- Pentru a explica modul în care se efectuează analiza varianței, să luăm în considerare mai multe populații, fiecare populație având o medie și o varianță proprie (evident, necunoscute).
- Extragem, din fiecare populație, câte un eșantion, conform schemei următoare:

# Analiza varianței - ANOVA (3)

<u>Populația 1</u>		<u>Populația <math>k</math></u>		<u>Populația <math>K</math></u>
media $\mu_1$		media $\mu_k$		media $\mu_K$
varianța $\sigma_1^2$		varianța $\sigma_k^2$		varianța $\sigma_K^2$
Eșantion de volum $n_1$		Eșantion de volum $n_k$		Eșantion de volum $n_K$
media de eșantion $m_1$		media de eșantion $m_k$		media de eșantion $m_K$
varianța de eșantion $s_1^2$		varianța de eșantion $s_k^2$		varianța de eșantion $s_K^2$

# Analiza varianței - ANOVA (4)

Analiza varianței se efectuează pentru o ipoteză nulă ( $H_0$ ): nu există diferențe între mediile populațiilor care va trebui respinsă, pentru a se confirma ipoteza alternativă

( $H_a$ ): cel puțin două dintre mediile diferă între ele (adică cel puțin două dintre populații diferă prin medii).

Ca de obicei în problemele de testare de ipoteze, admitem pentru moment că ipoteza nulă ar fi adevărată, și deducem consecințe logice ale ei.

# Analiza varianței - ANOVA (5)

Exemplu: acțiunea unui medicament asupra indivizilor din patru categorii de vârstă, timp de 60 de zile, exprimată în scăderea procentuală a nivelului colesterolului:

Sub 20 ani	20-39 ani	40-59 ani	Peste 60 ani
15	22	17	13
17	25	22	8
31	20	28	19
7	36	15	16
19	22	10	22
20	12	2	<u>media = 15.60</u>
<u>media = 18.17</u>	9	8	
	41	<u>media = 14.57</u>	
	17		
	<u>media = 22.67</u>		

# Analiza varianței - ANOVA (6)

Avem  $N = 27$ ,  $K = 4$ .

Rezultatele oferite de *EpiInfo* sunt următoarele:

ANOVA, a Parametric Test for Inequality of Population Means  
(For normally distributed data only)

Variation	SS	df	MS	F Statistic
Between	305.4376	3	101.8125	1.3414
Within	1745.7576	25	75.9021	
Total	2051.1852	26		

Sum of Squares

P-Value = 0.2822

Mean of Squares

# Analiza varianței - ANOVA (7)

Valoarea  $p$  fiind 0.2822, respingerea ipotezei nule este improprie (chiar dacă discrepanța între medii ni s-ar părea suficient de mare).

Nu dispunem de suficiente date pentru a trage concluzia că scăderea procentuală a nivelului colesterolului depinde de categoria de vârstă.

(Dar nici nu putem trage concluzia că nu depinde de categoria de vârstă!)

# TESTE STATISTICHE

Teste uzuale

- **2.1.a. O MEDIE EXPERIMENTALA ( $n > 30$ ) CU O VALOARE TEORETICĂ (media populației)**
- **TESTUL Z**
- $H_o : X_E = \mu_o$
- condiții:  $\sigma$  cunoscut, distribuție normală
- **TEST DE SEMNIFICATIE, PARAMETRIC**

- **2.1.b. O MEDIE EXPERIMENTALA ( $n \leq 30$ ) CU O VALOARE TEORETICĂ (media populației)**
- **TESTUL t (nepereche, “pooled”)**
- $H_o : X_E = \mu_o$
- Condiții: distribuție normală
- **TEST DE SEMNIFICATIE, PARAMETRIC**

- **2.2.a. DOUĂ MEDII EXPERIMENTALE,  
SERII MARI ( $n > 30$ )**

- **TESTUL  $z$**

- $H_0 : X_1 = X_2$

- condiții: varianțele cunoscute, distribuții normale

- **TEST DE SEMNIFICATIE, PARAMETRIC**

- **2.2.b. DOUĂ MEDII EXPERIMENTALE, SERII INDEPENDENTE, MICI ( $n \leq 30$ ) (indivizi diferiți)**
- **TESTUL t nepereche (“pooled”)**
- $H_0 : X_1 = X_2$
- Condiții: distribuții normale; 2 variante de lucru:
  - Varianțe egale ( $s_1 = s_2$ )
  - Varianțe diferite
- **TEST DE SEMNIFICATIE, PARAMETRIC**

- **2.2.c. DOUĂ MEDII EXPERIMENTALE,  
SERII PERECHE , MICI ( $n \leq 30$ )  
(aceeași indivizi, două condiții diferite)**
- **TESTUL t pereche (“paired”, “matched”)**
- $H_0 : X_1 = X_2$
- Condiții: distribuții normale
- **TEST DE SEMNIFICATIE, PARAMETRIC**

## ***2.2.d. DOUĂ SERII CU DISTRIBUTIE NECUNOSCUTA SAU NEGAUSSIANA***

- TESTUL MANN – WHITNEY ('u')
- $H_0 : X_1 = X_2$
- TEST DE SEMNIFICATIE, NEPARAMETRIC

- **2.3. RANGURI - DOUA SERII**
- **a) SERII INDEPENDENTE (nepereche)**
- TESTUL WILCOXON - 'RANK SUM' (suma rangurilor)
- **b) SERII DEPENDENTE (serii pereche)**
- TEST WILCOXON - 'SIGN - RANK' (semn-rang)
- $H_o : Me_1 = Me_2$
- TEST DE SEMNIFICATIE, NEPARAMETRIC, PENTRU VARIABLE ORDINALE