

CURSUL AL II-LEA

2. Indicatori statistici

2.1. Serii de valori. Așa cum s-a văzut în cursul anterior, uneori este necesar să urmărim mai întâi o singură variabilă numerică din multitudinea de variabile înregistrate într-un tabel de date. În acest caz, datele numerice pe care le avem la dispoziție sunt un simplu **șir de numere** asociate, fiecare din ele, unui individ.

Aceste șiruri de numere rezultate din datele culese le vom numi **serii statistice** sau **serii de date** sau **serii de valori**.

Ceea ce trebuie urmărit în primul rând la o **serie de valori** este modul în care valorile din serie sunt distribuite în plaja de valori între un minim și un maxim, cum se distribuie în jurul mediei, care este tendința centrală a seriei, care sunt valorile cel mai des întâlnite, etc.

Caracterizarea sintetică a unei serii de valori este dată de așa numiții **indicatori statistici**, între care media, deviația standard, mediana, etc, indicatori pe care îi vom descrie în continuare.

Definiție: Indicatorii statistici sunt **numere reale**, care sintetizează o parte din informația conținută de o serie de valori, dând posibilitatea aprecierii globale a întregii serii, în loc să ținem cont de fiecare valoare din șir.

Așa cum se va vedea în acest curs, fiecare indicator urmărește să scoată în evidență proprietăți diferite ale șirului de valori.

Astfel, prin combinarea mai multor indicatori, obținem informații relevante și sintetice despre valorile șirului. Dacă în locul șirului propriu-zis, folosim o serie de indicatori statistici, o parte din informație se pierde. Totuși, de obicei se pierde ceea ce este nesemnificativ, accidental, indicatorii statistici reținând doar esențialul. De aici și utilitatea și importanța lor în statistică.

În cele ce urmează, valorile din șirul de numere ce constituie o serie de valori le vom nota cu

$$X: x_1, x_2, \dots, x_n \text{ sau } Y: y_1, y_2, \dots, y_n$$

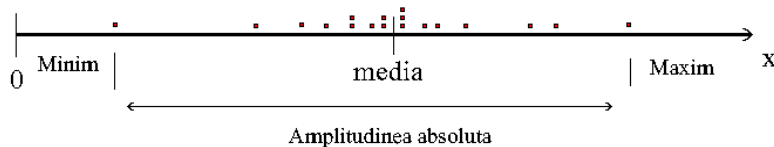
sau notații asemănătoare folosind alte litere ale alfabetului.

De exemplu, în loc să spunem că cele 10 valori ale glicemiei la cei zece pacienți dintr-un lot sunt: 88, 97, 103, 89, 93, 105, 98, 105, 88, 103, vom scrie în loc de *Glicemie* litera *X*, și în locul fiecărui număr din cele zece, simbolurile x_1, x_2, \dots, x_{10} . Deci, x_1 ține locul lui 88, x_2 pe cel al lui 97, etc. Aceste notații le folosim pentru a ușura înțelegerea formulelor de calcul pentru unii indicatori.

Valori extreme, amplitudine

Cel mai ușor de căutat și de înțeles ca semnificație sunt indicatorii **Minim** și **Maxim** care sunt cei ce ne indică *plaja de valori* pe care se întinde seria de valori. Minim este cea mai mică valoare din serie, iar Maxim este cea mai mare.

Amplitudinea absolută, este diferența dintre maximum și minimumul unei serii de valori și ne dă informații despre *lărgimea* plajei de valori pe care se întind datele din serie (vezi **figura 1.1**). O serie de valori cu o amplitudine mare indică o plajă de valori întinsă datorată fie unei dispersii sau împrăstieri mari a datelor, fie simplului fapt că sunt multe valori. Dacă două serii de valori au același număr de valori, dar una are o amplitudine mai mare, atunci valorile ei sunt mai împrăstiate.



$$\text{Amplitudinea relativa} = (\text{Maxim} - \text{Minim})/\text{media}$$

Figura 1.1. Indicatorii medie, minim, maxim, amplitudine absolută și amplitudine relativă.

De cele mai multe ori, valorile minimă și maximă dintr-o serie nu se înscriu în limitele de normalitate, ceea ce nu înseamnă neapărat că seria conține valori anormale. Totuși, de obicei, cele mai îndepărtate câteva valori, atât cele mai mici cât și cele mai mari trebuie verificate pentru a ne asigura că nu este vorba de date eronate.

De exemplu, deși se consideră că valorile normale pentru latența semnalului nervos pe nervul optic între stimularea retinei și răspunsul cortical sunt situate aproximativ între 90 ms și 115 ms, un eșantion de indivizi sănătoși poate să producă o serie de valori care are și una sau câteva excepții. De aceea, din 20 sau 30 de valori, una poate fi 88 ms iar alta 117 ms, majoritatea fiind însă între 90 și 115 ms.

2.2. Valori medii. Media aritmetică a unei serii de valori. Este un indicator simplu și în același timp foarte sintetic, fiind un foarte bun indiciu al valorii în jurul căreia se grupează datele. Se notează cu litera m sau, dacă seria de valori este notată cu o majusculă ca X sau Y , media se notează cu \bar{X} sau \bar{Y} . Formula este cea cunoscută:

$$\bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n} = m \quad (1.1)$$

Definiție:

Media aritmetică unei serii de valori este raportul dintre suma valorilor seriei și numărul lor.

Media este indicatorul care arată tendința centrală a seriei de valori, și de obicei arată unde tind datele să se aglomereze. De cele mai multe ori, valorile din serie sunt situate în majoritate în apropierea mediei, iar o mai mică parte din ele sunt situate mult în stânga sau în dreapta mediei. O situație a valorilor din serie față de medie se poate observa din așa-numitul grafic punctual de dispersie, din care este dat un exemplu în **figura 1.2**



Figura 1.2. Cele mai multe valori sunt de obicei mai apropiate de medie.

Dar nu totdeauna datele din seria de valori se situează preponderent în apropierea mediei. Mai rar, și oarecum mai forțat, ne putem întâlni și cu situații în care datele din serie se situează preponderent în stânga și dreapta, departe de medie și doar o mică parte dintre ele se situează aproape de medie, așa cum se observă în **figura 1.3**.

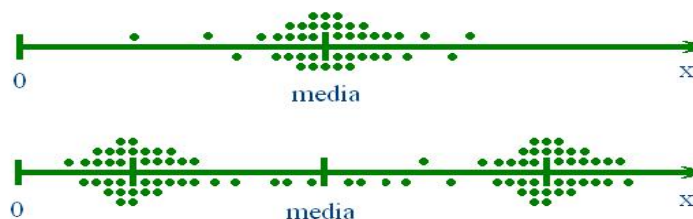


Figura 1.3. Uneori, cele mai multe valori sunt sub medie și peste medie, destul de departe de aceasta. În seriile de mai sus, avem aceeași medie, dar este evident că nu avem aceeași situație. Valorile din seria de jos sunt mai împrăștiate.

Astfel, dacă în același lot sunt cuprinși indivizi hipertiroidieni și hipotiroidieni, și se măsoară la fiecare concentrația hormonului tiroidian T_4 , vom observa că hipotiroidienii au preponderent valori în stânga mediei, cei mai mulți destul de departe de medie, iar hipertiroidienii au preponderent valori în dreapta, tot departe de medie.

De fapt într-un asemenea caz, în zona centrală lipsesc exact ceea ce am spune că sunt normalii, adică indivizi care au valori pentru T_4 ușor peste medie și ușor sub medie, și care nu au fost incluși într-un astfel de lot.

Evident că un eșantion așa de eterogen nu este folosit prea des în statistică pentru că, așa cum vom vedea, în acest caz este foarte indicat să se constituie două eșantioane distincte pentru cele două categorii de pacienți. Totuși, asemenea situații, chiar dacă de obicei nu sunt indicate și sunt puțin artificiale, există. Situația de mai sus este ilustrată în **figura 1.3**.

O formulă simplificată pentru media aritmetică este dată de:

$$\bar{X} = \frac{x_1 \cdot F_1 + x_2 \cdot F_2 + \dots + x_n \cdot F_n}{F_1 + F_2 + \dots + F_n}$$

unde cu n am notat numărul de valori diferite din seria de valori, iar F_1, F_2, \dots, F_n sunt frecvențele de apariție în serie ale valorilor x_1, x_2, \dots, x_n .

Această formulă se spune că este formula pentru **media ponderată**. Nu trebuie să credem că media ponderată calculată cu formula de mai sus și media aritmetică calculată cu formula (1.1), sunt indicatori diferiți. Ambele medii sunt în realitate identice. Media ponderată se calculează de obicei mai simplu și deci nu reprezintă decât o formă mai simplă de calcul al mediei aritmetice.

Prin faptul că este un indicator extrem de fidel al tendinței centrale al unei serii statistice, media este un indicator extrem de mult utilizat în statistică. Media aritmetică are dezavantajul că este sensibilă la valori extreme fie foarte mici, fie foarte mari. Adăugarea unei singure valori (sau a câtorva) mult mai mari decât celelalte, modifică sensibil media aritmetică.

De asemenea, dacă datele sunt distribuite în jurul mediei puternic asimetric, media își pierde din puterea de a evoca tendința centrală, în aceste cazuri fiind mult mai utilă mediana (vezi mai jos).

2.3. Împrăștiere. Valorile dintr-o serie de valori pot fi mai aglomerate în jurul mediei sau mai dispersate, adică la distanțe mari de medie. Un mod de a măsura aceste abateri de la medie este să se facă diferența între toate aceste valori și media lor. Unele abateri vor fi pozitive, altele negative. Ele nu pot fi adunate, deoarece, prin adunare dau suma 0.

Dispersia. Un mod de a ocoli faptul că suma abaterilor absolute este 0, este ridicarea la pătrat a acestora înainte de a fi adunate, pentru a face să dispară semnele negative la unele și pozitive la altele.

Suma obținută, ar trebui împărțită la numărul de abateri pentru a se obține o medie. În realitate, din motive teoretice foarte bine întemeiate, dar mai greu de explicat în cuvinte simple, împărțirea se face la $n-1$ și nu la n . Motivul pentru care se face acest lucru va fi înțeles mai bine în contextul unor noțiuni enunțate la cursul despre teoria estimației. Valoarea care se obține astfel se numește **dispersie** și este un indicator al gradului de împrăștiere al seriei. Dispersia se notează cu D și are formula:

$$D = \frac{(x_1 - \bar{X})^2 + (x_2 - \bar{X})^2 + \dots + (x_n - \bar{X})^2}{n - 1}$$

După cum se observă, numărătorul fracției din definiția dispersiei este cu atât mai mare cu cât abaterile individuale de la medie sunt mai mari și deci este natural să considerăm că o valoare mare a dispersiei arată o împrăștiere mare a valorilor din serie.

De fapt, este bine de reținut că:

- La medii aproximativ egale, este mai împrăștiată seria cu dispersia mai mare.
- La dispersii aproximativ egale, este mai împrăștiată seria cu media mai mică.

Dispersia are dezavantajul că se exprimă cu unitățile de măsură ale valorilor din serie, ridicate la pătrat, și are în general valori foarte mari comparativ cu abaterea medie. De exemplu, dacă valorile din serie se măsoară în mg/l, atunci dispersia se măsoară în mg^2/l^2 , ceea ce este în mod evident extrem de nenatural.

În plus, dacă abaterile absolute au o medie, de exemplu în jurul lui 10, dispersia va avea o valoare în jurul lui 100, adică exagerat de mare în comparație cu abaterile absolute. De aceea se mai folosește un alt indicator, numit abatere standard care este radicalul dispersiei.

Abaterea standard. Se notează cu σ și are formula:

$$\sigma = \sqrt{D} \text{ sau } \sigma = \sqrt{\frac{(x_1 - \bar{X})^2 + (x_2 - \bar{X})^2 + \dots + (x_n - \bar{X})^2}{n - 1}}$$

Acest indicator se exprimă cu aceeași unitate de măsură ca și valorile din seria considerată și este un indicator foarte fidel al împrăștierii seriei. Abaterile standard, nu are dezavantajele dispersiei, adică unitatea de măsură este aceeași cu a valorilor din serie, și, are o valoare comparabilă cu abaterile individuale de la medie.

Exemplu de calcul:

Să presupunem că am măsurat zilnic tensiunea arterială sistolică la doi pacienți timp de 10 zile, obținând pentru fiecare următoarele valori:

- 170, 180, 160, 180, 190, 190, 180, 190, 170, 190, pentru primul pacient și
- 160, 170, 190, 160, 190, 190, 200, 180, 180, 180, pentru al doilea.

Lăsând la o parte studiul modului cum evoluează de la zi la zi tensiunea pacienților, care este bineînțeles importantă, să ne propunem să determinăm care are tensiunea cu valori mai împrăștiate, indiferent de evoluția în timp.

Notând prima serie cu X iar pe a doua cu Y se constată ușor că ambele au media 180 (datele nu sunt reale, au fost deliberat alese ca să simplifice calculele). Atunci, vom avea pentru abaterile de la medie și pentru pătratele lor următoarele valori:

- $x_i - \bar{X}$: -10, 0, -20, 0, 10, 10, 0, 10, -10, 10. $\bar{X} = 180$.
- $y_i - \bar{Y}$: -20, -10, 10, -20, 10, 10, 20, 0, 0, 0. $\bar{Y} = 180$.
- $(x_i - \bar{X})^2$: 100, 0, 400, 0, 100, 100, 0, 100, 100, 100.
- $(y_i - \bar{Y})^2$: 400, 100, 100, 400, 100, 100, 400, 0, 0, 0.

Deci vom avea pentru D_x :

$$D_x = \frac{(x_1 - \bar{X})^2 + (x_2 - \bar{X})^2 + \dots + (x_{10} - \bar{X})^2}{10 - 1} = \frac{400 + 6 \cdot 100}{9} = \frac{1000}{9} = 111,1$$

și cu un calcul absolut analog, $D_y = 1600 / 9 = 177,7$. Se observă că, în timp ce abaterile de la medie sunt de ordinul zecilor, dispersiile sunt de ordinul sutelor, ceea ce este destul de nenatural, și în plus, după cum am mai spus, unitatea de măsură este cu totul alta.

Pentru abaterile standard, vom avea:

$$\sigma_x = \sqrt{D_x} = \sqrt{111,1} = 10,5$$

$$\sigma_y = \sqrt{D_y} = \sqrt{177,7} = 13,3$$

calculele fiind făcute cu aproximație. Deci, este ceva mai împrăștiată seria Y .

De fapt, este bine de reținut că:

- La medii aproximativ egale, este mai împrăștiată seria cu deviația standard mai mare.
- La deviații standard aproximativ egale, este mai împrăștiată seria cu media mai mică.

Ce se întâmplă însă dacă mediile și deviațiile sunt foarte diferite? Atunci o bună apreciere se obține dacă se folosește raportul deviației standard față de medie, exprimat în procente, acest raport fiind un alt indicator al împrăștierii valorilor dintr-o serie. Acest indicator se numește **coeficient de variație**.

Coeficientul de variație. Este raportul dintre deviația standard și medie, atunci când media este diferită de 0 și se exprimă în procente:

$$C.V. = \frac{\sigma}{\bar{X}}$$

Pentru seriile de mai sus, coeficientul de variație este mai mare pentru cea mai împrăștiată, adică pentru cea cu deviația standard mai mare:

- $C.V._x = 10,5 / 180 = 0,058 = 5,8 \%$.
- $C.V._y = 13,3 / 180 = 0,073 = 7,3\%$.

Totuși, seriile de mai sus sunt comparabile cu ajutorul abaterilor standard, deoarece au aceeași medie, și, așa cum s-a văzut, la medii egale sau aproximativ egale, are valorile mai împrăștiate seria cu abaterea standard mai mare.

Aprecierea cu ajutorul coeficientului de variație se face mai ales atunci când două serii de valori au medii mult diferite și deviațiile standard pot să nu ne dea o indicație suficient de utilă. De exemplu, măsurând **latența și amplitudinea semnalului electric pe nervul optic** la 120 de pacienți cu *scleroză multiplă*, s-au obținut următoarele rezultate:

- *Latența* medie: 113,6
- Abaterea standard a *latenței*: 14,7
- *Amplitudinea* medie: 2,68
- Abaterea standard a *amplitudinii*: 2,03

Dacă dorim să apreciem împrăștierea valorilor din cele două serii, abaterile standard nu ne sunt de ajutor. Într-adevăr, latența are o abatere standard mult mai mare decât amplitudinea, dar și media latenței este cu mult mai mare decât aceea a amplitudinii. De aceea, în acest caz, doar coeficientul de variație ne permite o apreciere corectă a împrăștierilor, în vederea comparării lor:

- Pentru latență: $C.V._{latența} = \frac{14,7}{113,6} = 0,129 = 12,9\%$
- Pentru amplitudine: $C.V._{amplitudine} = \frac{2,03}{2,68} = 0,757 = 75,7\%$

Se observă că valorile amplitudinii sunt cu mult mai împrăștiate decât cele ale latenței. Acest fapt se datorează atât unei variabilități biologice mai mari la amplitudine decât la latență, cât și unei variabilități datorate aparatelor de măsură, care măsoară latența cu mai multă precizie, în timp ce la măsurarea amplitudinii, erorile de măsurare sunt mai mari.

Coeficientul de variație este cel mai fidel indicator al împrăștierii unei serii statistice, dar are și el un inconvenient, este cu atât mai fidel cu cât mediile sunt mai depărtate de 0.

La medii foarte apropiate de 0 își pierde din fidelitate și nu este indicat să fie folosit. Acest lucru se întâmplă mai ales atunci când valorile din serie sunt și negative și pozitive, și când, din acest motiv, media poate fi aproape de 0.

2.4. Indicatori de asimetrie. Atunci când valorile unei serii sunt distribuite nesimetric în jurul mediei, acest fapt este imposibil de surprins cu ajutorul indicatorilor de dispersie. De aceea, s-au introdus indicatori care să pună în evidență și acest aspect al seriilor de valori: excentricitatea, sau asimetria. Va trebui să ținem cont atât de numărul de valori care sunt în stânga și în dreapta mediei, cât și depărtarea lor față de medie.

Mediana. Este un indicator al tendinței centrale, și anume este valoarea de mijloc, într-o serie de valori.

Definiție:

Mediana este acea valoare dintr-o serie de valori, pentru care exact jumătate din ele sunt mai mici decât ea, iar jumătate mai mari.

Altfel spus, este valoarea măsurată pentru individul din mijloc, dacă indivizii pe care s-au făcut măsurătorile ar fi ordonați crescător. Pentru o înțelegere mai ușoară, să luăm un exemplu cu numai 10 înregistrări: tensiunea arterială maximă la un bolnav în 10 zile:

150, 160, 160, 170, 160, 170, 150, 160, 170, 160.

Dacă se așază aceste valori într-un șir crescător, obținem:

150, 150, 160, 160, 160, 160, 160, 170, 170, 170.

În acest caz, mediana se ia între a cincia și a șasea valoare din acest șir ordonat, adică 160. Dacă aceste două valori de mijloc diferă, se ia media lor aritmetică. Dacă numărul de măsurători este impar atunci mediana este chiar valoarea de mijloc, care în acest caz este unică.

De fapt, mediana este importantă în primul rând la serii de valori cu foarte multe înregistrări, caz în care se poate lucra direct pe tabelul de frecvență, sau chiar pe tabelul pe clase.

Pentru a exemplifica modul cum se caută mediana pe tabelul de frecvență, vom lua **tabelul 1.3**, în care sunt centralizate vârstele a 234 de pacienți, fiecare valoare a vârstei având o anumită frecvență absolută F_i , o frecvență relativă f_i și o frecvență relativă cumulată crescător, f_{icc} (vezi mai sus, pentru amănunte).

Valoarea medianei se culege din coloana întâi, a vârstelor, dar pentru a ști care valoare trebuie aleasă, trebuie să privim pe ultima coloană, a frecvențelor cumulate, f_{icc} , în dreptul frecvenței cumulate de 50%.

Se observă că, pe coloana frecvențelor cumulate, nu există frecvența de 50%, dar, există frecvența de 47,9%, care este prea mică, și frecvența de 53,8%, care este prea mare. În acest caz, mediana se citește din dreptul primei frecvențe cumulate crescător care depășește 50%, în cazul nostru, în dreptul frecvenței de 53,8%, și pe coloana **Vârsta** citim 55 ani. Deci, vârsta mediană este 55 ani.

Tabelul 1.3. Vârstele a 234 de pacienți centralizate într-un tabel de frecvență

Vârsta	F_i	f_i	f_{icc}	Vârsta	F_i	f_i	f_{icc}
26	1	0.4%	0.4%	53	10	4.3%	43.2%
28	1	0.4%	0.9%	54	11	4.7%	47.9%
29	1	0.4%	1.3%	55	14	6.0%	53.8%
30	2	0.9%	2.1%	56	11	4.7%	58.5%
31	2	0.9%	3.0%	57	9	3.8%	62.4%
32	1	0.4%	3.4%	58	19	8.1%	70.5%
35	3	1.3%	4.7%	59	5	2.1%	72.6%
36	2	0.9%	5.6%	60	9	3.8%	76.5%
37	3	1.3%	6.8%	61	13	5.6%	82.1%
38	1	0.4%	7.3%	62	5	2.1%	84.2%
40	3	1.3%	8.5%	63	4	1.7%	85.9%
41	5	2.1%	10.7%	64	4	1.7%	87.6%
42	1	0.4%	11.1%	65	6	2.6%	90.2%
43	4	1.7%	12.8%	66	2	0.9%	91.0%
44	10	4.3%	17.1%	67	4	1.7%	92.7%
45	6	2.6%	19.7%	68	3	1.3%	94.0%
46	6	2.6%	22.2%	69	4	1.7%	95.7%
47	5	2.1%	24.4%	70	1	0.4%	96.2%
48	13	5.6%	29.9%	71	2	0.9%	97.0%
49	2	0.9%	30.8%	72	2	0.9%	97.9%
50	4	1.7%	32.5%	74	1	0.4%	98.3%
51	6	2.6%	35.0%	77	2	0.9%	99.1%
52	9	3.8%	38.9%	78	2	0.9%	100.0%
				Total	234	100.0%	

Deci, vom spune că jumătate dintre pacienți au vârstele cuprinse între 26 și 55 ani și jumătate au vârstele mai mari decât 55 ani. Această alegere este permisă în cazul acesta al vârstelor care se înregistrează cu valori întregi.

Mediana este un indicator al tendinței centrale, ca și media, dar oferă mai puțină informație decât aceasta din urmă. La distribuțiile echilibrate, la care valorile din serie se dispun aproximativ simetric în stânga și în dreapta mediei, media și mediana sunt foarte apropiate, deci folosirea medianei este superfluă. Dacă însă mediana este mult în stânga sau în dreapta mediei, distribuția se zice că este excentrică.

De exemplu, venitul median este mai informativ decât venitul mediu deoarece distribuția veniturilor într-o populație este foarte excentrică, fiind foarte mulți indivizi cu salarii foarte mici și foarte puțini indivizi cu salarii foarte mari.

Cuartilele. În mod asemănător cu căutarea medianei, se poate pune problema căutării unor valori pentru care să avem un sfert din valorile seriei mai mici și respectiv, mai mari.

Definiție:

Cuartila Q_1 este acea valoare dintr-o serie de valori, pentru care 25% din valorile seriei sunt sub Q_1 și 75%, peste.

Pentru tabelul de frecvențe 1.3, cuartila Q_1 se caută în dreptul frecvenței relative cumulate crescător de 25%. În tabel găsim procentul de 24,4% și în dreptul lui vârsta de 47 de ani, precum și frecvența de 29,9 și în dreptul ei vârsta de 48 de ani. Vom lua tot vârsta care corespunde primului procent peste 25%, adică 48 de ani.

Definiție:

Cuartila Q_3 este acea valoare dintr-o serie de valori, pentru care 75% din valorile seriei sunt sub Q_3 și 25%, peste.

Pentru tabelul 1.3, cuartila Q_3 se ia din dreptul frecvenței relative cumulate crescător de 75%. Poate fi luată cu aproximație, 60 ani.

Care este utilitatea medianei și cuartilelor în aprecierea simetriei distribuției? Pentru a sublinia utilitatea indicatorilor Q_1 și Q_3 , să considerăm șirul vârstelor:

- cel mai tânăr pacient,
- Q_1 ,
- Vârsta mediană,
- Q_3 ,
- cel mai în vârstă pacient.

Pentru **tabelul 1.3**, obținem șirul: 26 ani, 48 ani, 55 ani, 60 ani, 69 ani.

- Se observă că sfertul (25%) pacienților cei mai tineri este situat în zona 26 - 48 de ani adică într-o plajă de 22 de ani.
- Sfertul următor, este între 48 și 55 de ani, adică pe un interval de doar 7 ani.
- Al treilea sfert este situat între 55 și 60 de ani, adică pe 5 ani,
- Cei mai în vârstă 25 % din pacienți sunt între 60 și 69 de ani, pe un interval de 9 ani.

Putem să spunem că vârstele pacienților se distribuie **ușor asimetric**, deoarece:

1. Sfertul cel mai tânăr se distribuie pe o plajă de 22 de ani, iar cel mai în vârstă pe o plajă de doar 9 ani.
2. Sfertul al doilea se distribuie pe 7 ani, iar al treilea doar pe 5 ani.

În cadrul laboratorului, alte exemple vor arăta utilitatea acestor indicatori.

Să mai observăm că mediana este într-un fel “cuartila de 50%”, adică Q_2 . Se spune că există trei cuartile: Q_1 , mediana, Q_3 .

Decile Uneori, loturi mai mari de multe sute de indivizi trebuie urmărite foarte atent în ceea ce privește modul cum sunt distribuite valorile și de aceea s-au introdus indicatorii **decile**, care sunt de o acuratețe mai bună decât cuartilele. Sunt 9 decile, fiecare corespunzând unui procent de 10%, 20%, ... 90% din lot, asemănător cu cuartilele. Decila 5, sau de 50%, este de fapt mediana.

Centilele (percentilele) sunt mai rar folosite, în studii pe mii de cazuri, de obicei de un interes mai larg, național, internațional, în studii epidemiologice, și sunt corespunzătoare procentelor de 1%, 2%,...99% din lot. Centila de 25% este cuartila Q_1 , cea de 50% este mediana, iar cea de 75% este cuartila Q_3 . Centilele de 10%, 20%,...90%, sunt cele nouă decile. Centilele dau o imagine destul de exactă a distribuției valorilor dintr-o serie de valori foarte mare. Nu are rost să calculăm centile pentru serii cu câteva sute de valori, pentru că erorile sunt prea mari și imaginea obținută este deformată.

Modul. Dintre frecvențele absolute apărute într-un tabel de frecvențe, una este maximă. Clasa sau valoarea corespunzătoare acestei frecvențe maxime se numește **mod**. Modul este de obicei un indicator al tendinței centrale. În tabelul 1.2. modul este clasa de la 55 la 60 de ani, cu frecvența absolută 53. De obicei, frecvențele absolute au tendința de a crește către mod, după care urmează o descreștere continuă. Modul este deci o indicație relativă la maximumul frecvențelor absolute. Sunt însă distribuții la care se înregistrează creșteri și descreșteri astfel încât pot apare două moduri sau chiar mai multe. Aceste distribuții sunt mai rare și au un caracter cu totul special. Ele se numesc distribuții **bimodale** sau **multimodale**, după caz.

Este un indicator care poartă în el puțină informație despre datele seriei. Modul este mult influențat de fluctuații aleatoare și nu este prea recomandat pentru a aprecia tendința centrală a valorilor dintr-o serie. Mai mult, unele distribuții pot fi multimodale, caz în care modul nu mai indică prea mult despre tendința centrală.

Excentricitate. (Engl. Skew, Skweness). Este un indicator al asimetriei și este luat de diverși autori cu diverse formule. Distribuțiile cu excentricitate pozitivă sunt mai des întâlnite decât cele cu excentricitate negativă. În medicină, parametrii fiziologici sunt în majoritate modificați în diverse afecțiuni în sensul că au valori peste normal. Astfel, tensiunea arterială o vom întâlni la valori normale, crescute sau scăzute. Cum indivizi cu valori foarte mari, vom întâlni cu atât mai rar cu cât valoarea este mai mare, distribuția va avea o coadă spre dreapta. La fel la mulți alți parametri cum ar fi bilirubina, transaminazele, colesterolul, lipemia, etc.

Totuși, vom întâlni și parametri care se distribuie cu asimetrie stânga în patologii: hemoglobina, calcemia, sodiul ionic, etc. Hemoglobina, de exemplu, se poate distribui cu frecvență mai mare la valori relativ normale și cu frecvențe din ce în ce mai mici pe măsură ce coborâm la valori mai mici. Chiar dacă avem o patologie de tip anemie, ne așteptăm ca frecvența în jurul a 9-10 să fie mai mare decât frecvența în jurul a 7-8, frecvență care ne așteptăm să fie foarte mică.

Excentricitatea unei serii de valori x_1, x_2, \dots, x_n , se calculează cu formula:

$$sk = \frac{\sum_{i=1}^n (x_i - \bar{X})^3}{n\sigma^3}$$

Cu cât o distribuție este mai simetrică cu atât sk tinde la 0. Ca o regulă generală, **la distribuțiile cu excentricitate pozitivă, media este mai mare decât mediana**. Evident, media este mai mică decât mediana la distribuțiile cu excentricitate negativă. Există cazuri rare în care regula de mai sus nu este valabilă.

Sunt multe alte formule pentru alți coeficienți de excentricitate și când vorbim despre excentricitate, trebuie să menționăm la ce coeficient de excentricitate ne referim. Uneori se folosește un coeficient de asimetrie care măsoară diferența dintre medie și mediană, eventual raportată la abaterea standard sau la intervale intercuartilice ($Q_3 - Q_1$). Indiferent ce formulă se folosește, o excentricitate egală cu zero, sau foarte apropiată de zero, este un indiciu al simetriei repartiției valorilor din serie. Din contră, excentricități mult diferite de 0, peste 0,15 -0,20, sau mai jos de -0,15 -0,20 sunt indicii ale asimetriei. Dăm mai jos, cu titlu facultativ, câteva formule pentru coeficienți de excentricitate.

$$sk_1 = \frac{\bar{X} - Mo}{\sigma} \quad sk_2 = \frac{3(\bar{X} - Me)}{\sigma} \quad sk_3 = \frac{2(Q_3 + Q_1 - 2Me)}{Q_3 - Q_1} \quad sk_4 = \frac{(Q_3 - Me) - (Me - Q_1)}{(Q_3 - Me) + (Me - Q_1)}$$

Boltirea (facultativ). Boltirea este un indicator care se bazează pe lungimea cozilor unei distribuții. Cele cu cozi relativ mari se numesc **leptocurtice** iar cele cu cozi relativ mici se numesc **platicurtice** (vezi **figura 1.4**). Formula de calcul a boltirii este:

$$k = \frac{\sum_{i=1}^n (x_i - \bar{X})^4}{n\sigma^4} - 3$$

Așa cum se va vedea în capitolul despre repartiții, boltirea este un indicator util în aprecierea apropierii repartiției de repartiția normală. Distribuțiile din **figura 1.4** au aceeași medie, aceeași dispersie, aproximativ aceeași excentricitate dar diferă mult ca boltire.

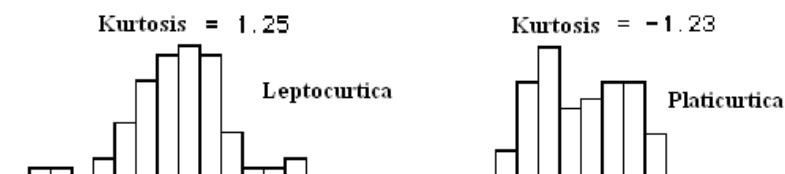


Figura 1.4. Distribuție leptocurtică și distribuție platicurtică.

2.5. Clasificarea indicatorilor

Indicatorii statistici poartă în ei, fiecare, o anumită cantitate de informație, din seria de valori pentru care au fost calculați. Așa cum s-a văzut în paragrafele precedente, unii indicatori ne dau informații despre tendința centrală a valorilor din serie, alții ne dau informații despre împrăștierea valorilor, alții ne dau indicații despre simetria valorilor din serie, boltirea ne dă indicații despre lungimea cozilor distribuției, etc.

Informația oferită de indicatorii statistici este redundantă, în sensul că, de exemplu, împrăștierea valorilor din serie este indicată și de dispersie și de abaterea standard și de amplitudinea absolută și de coeficientul de variație, etc. Totuși, fiecare din ei aduce o mică informație specifică, deci, nu ne putem lipsi de unul sau altul

dintre indicatorii statistici. Uneori trebuie folosiți unii dintre indicatori, fiind cei mai eficienți, alții trebuie folosiți alții.

Pentru a avea o idee despre modul cum trebuie folosiți indicatorii statistici, ei sunt clasificați în câteva categorii mai importante, categorii care vor fi exemplificate mai jos, insistând pe aceia care sunt cei mai importanți, restul fiind indicatori mai rar folosiți, numai în cazuri speciale.

Indicatori ai tendinței centrale. Cei mai importanți indicatori ai tendinței centrale sunt **media, mediana și modul**. Media indică tendința centrală atunci când seria de valori este repartizată simetric în jurul ei și când valorile nu au o dispersie exagerat de mare. În cazul seriilor de valori distribuite foarte asimetric, tendința centrală nu mai este indicată de către medie, ci de către mediană.

Modul, este un indicator al tendinței centrale, la seriile unimodale, adică atunci când în tabelul de frecvențe există un singur maxim. Dacă avem o serie multimodală, modul își pierde calitatea de indicator al tendinței centrale.

Indicatori ai împrăștierii. Folosiți mai des în practică, și deci mai importanți, sunt **dispersia, abaterea standard și coeficientul de variație**.

Abaterea standard este indicatorul folosit cel mai des pentru aprecierea împrăștierii, dar atunci când mediile diferă mult, este mai util coeficientul de variație. Dispersia este folosită ca măsură a împrăștierii în testele statistice (vezi capitolul dedicat testelor statistice).

Indicatori ai asimetriei. **Mediana și cuartilele** sunt cel mai mult folosite pentru aprecierea asimetriei valorilor dintr-o serie. De fapt, mediana se folosește în combinație cu media pentru aprecierea asimetriei. O mediană mult diferită de medie indică asimetrie puternică, iar o mediană foarte apropiată de medie indică o tendință spre simetrie.

Cuartilele, se folosesc în combinație cu mediana și indicatorii minim și maxim, pentru aprecierea simetriei.

Indicatorii statistici fundamentali. Sunt indicatorii care poartă în ei cea mai mare cantitate de informație din informația conținută de seria de valori.

La seriile de valori distribuite relativ simetric, indicatorii statistici fundamentali sunt **media și deviația standard**. În capitolul dedicat repartițiilor, se va vedea că, dacă o serie de valori are o repartiție normală și are suficient de multe valori, cei doi indicatori, poartă în ei aproape toată informația. Astfel, dacă o serie de valori de acest tip are media \bar{X} și deviația standard σ , scrierea încetățenită este $\bar{X} \pm \sigma$

La seriile distribuite asimetric, deși se consideră ca indicatori fundamentali tot media și deviația standard, sunt mai utile **mediana și cuartilele**. În acest caz, este încetățenită scrierea medianei M și a cuartilelor Q_1 și Q_3 în forma $M [Q_1; Q_3]$. De exemplu, dacă o serie puternic asimetrică are mediana 2,45, iar cuartilele sunt $Q_1=1,54$ și $Q_3=6,23$, acest fapt se precizează astfel: 2,45 [1,54; 6,23].

3. Chestiuni de examen:

1. Definiția și formula mediei

2. Formula deviației standard și a coeficientului de variație

3. Definiția medianei și a cuartilelor Q_1, Q_3

4. Media unei serii de valori numerice este:

- A. Suma valorilor împărțită la numărul lor
- B. Mai mare decât valoarea minimă din serie
- C. Mai mică decât valoarea maximă din serie
- D. Un indicator al tendinței centrale a valorilor seriei

5. Media unei serii de valori numerice are următoarele proprietăți:

- A. Este egală cu cea mai mică valoare din serie
- B. Dacă schimbăm o valoare din serie, mărind-o, media se schimbă, mărindu-se
- C. Dacă schimbăm o valoare din serie, mărind-o, media se schimbă, micșorându-se
- D. Dacă ștergem o valoare din serie, media rămâne nemodificată

6. Media unei serii de valori numerice este un indicator al:
- A. **Tendinței centrale a valorilor seriei**
 - B. Împrăștierei valorilor seriei
 - C. Plaja de valori între care sunt cuprinse valorile seriei
 - D. Media nu este indicator statistic
7. Dispersia unei serii de valori numerice este un indicator al:
- A. Tendinței centrale a valorilor seriei
 - B. **Împrăștierei valorilor seriei**
 - C. Plaja de valori între care sunt cuprinse valorile seriei
 - D. Simetriei distribuției valorilor seriei în jurul mediei
8. Dispersia unei serii de valori numerice are printre dezavantaje:
- A. **Se măsoară cu unitatea de măsură a valorilor seriei, ridicată la pătrat**
 - B. **Are valori prea mari, comparativ cu abaterile individuale de la medie**
 - C. Indică și tendința centrală a valorilor seriei
 - D. Nu se poate calcula cu exactitate
9. Abaterea standard unei serii de valori numerice are printre avantajele:
- A. **Se măsoară cu unitatea de măsură a valorilor seriei**
 - B. **Are valori comparabile cu abaterile individuale de la medie**
 - C. Indică și tendința centrală a valorilor seriei
 - D. Nu se poate calcula dacă dispersia este negativă
10. Dacă două serii de valori au aproximativ aceeași medie, atunci:
- A. **Este mai împrăștiată cea cu dispersia mai mare**
 - B. Este mai împrăștiată cea cu abaterea standard mai mică
 - C. Sunt la fel de împrăștiate
 - D. Nu se pot compara împrăștierea cu ajutorul dispersiei în acest caz
11. Dacă două serii de valori au medii foarte diferite, atunci:
- A. Este mai împrăștiată cea cu dispersia mai mare
 - B. Este mai împrăștiată cea cu abaterea standard mai mare
 - C. **Nu se pot compara nici cu ajutorul dispersiei și nici cu ajutorul abaterii standard**
 - D. Au aceeași împrăștiere
12. Dacă media unei serii de valori este 10 și dispersia 4, atunci coeficientul de variație este:
- A. 40%
 - B. **20%**
 - C. 80%
 - D. 10%
13. Dacă mediile a două serii de valori sunt foarte diferite, iar abaterile standard sunt tot foarte diferite, atunci este mai împrăștiată :
- A. **Cea cu coeficientul de variație mai mare**
 - B. **Cea cu raportul dintre abaterea standard și medie mai mare**
 - C. Cea cu coeficientul de variație mai mic
 - D. Împrăștierea celor două serii de valori nu se pot compara
14. Mediana unei serii de valori numerice este:
- A. Egală cu media
 - B. Un grafic
 - C. **Un număr**
 - D. Un tabel de frecvență
15. Mediana unei serii de valori numerice este:
- A. **Valoarea pentru care jumătate din valorile seriei sunt mai mari și jumătate mai mici**
 - B. Valoarea situată la mijloc, între minimul seriei și maximumul seriei
-

- C. Valoarea cea mai frecvent întâlnită printre valorile seriei
D. Un indicator al excentricității valorilor seriei
16. Dacă o serie de valori are în componență 21 de numere, atunci, pentru aflarea medianei, se ordonează valorile crescător și se ia:
A. Valoarea a 11-a din șirul ordonat
B. Media între valorile a 10 și a 11-a
C. Media între valorile a 11 și a 12-a
D. Valoarea a 10-a din șirul ordonat
17. Dacă o serie de valori are în componență 24 de numere, atunci, pentru aflarea medianei, se ordonează valorile crescător și se ia:
A. Valoarea a 12-a din șirul ordonat
B. Media între valorile a 11-a și a 12-a
C. Media între valorile a 12-a și a 13-a
D. Valoarea a 13-a din șirul ordonat
18. Cuartila întâi a unei serii de valori este:
A. Valoarea din seria ordonată situată la 25% din numărul de valori al seriei
B. Valoarea din seria ordonată situată la 75% din numărul de valori al seriei
C. Valoarea numerică pentru care un sfert din valorile seriei ordonate sunt mai mici
D. Valoarea numerică pentru care un sfert din valorile seriei sunt mai mici
19. Cuartila a treia a unei serii de valori este:
A. Valoarea din seria ordonată situată la 25% din numărul de valori al seriei
B. Valoarea din seria ordonată situată la 75% din numărul de valori al seriei
C. Valoarea numerică pentru care un sfert din valorile seriei ordonate sunt mai mici
D. Valoarea numerică pentru care trei sferturi din valorile seriei ordonate sunt mai mari
20. Referitor la indicatorii decile, este adevărat:
A. Avem exact nouă decile
B. Avem exact 99 de decile
C. Decila 50 este mediana
D. Decila a treia este mediana
21. Indicatorii statistici fundamentali sunt:
A. Dispersia și media
B. Media și abaterea standard
C. Abaterea standard și mediana
D. Mediana și cuartilele
22. Indicatorii de dispersie (sau de împrăștiere) sunt:
A. Amplitudinea, media, dispersia și mediana
B. Abaterea standard, media, dispersia și mediana
C. Amplitudinea, media, dispersia și abaterea standard
D. Abaterea standard, dispersia și coeficientul de variație
23. Care din următorii indicatori statistici ajută la aprecierea asimetriei:
A. Mediana, media și excentricitatea
B. Mediana, cuartilele și excentricitatea
C. Mediana, cuartilele și media
D. Mediana, dispersia și excentricitatea
24. Indicatorii statistici pentru tendința centrală a valorilor unei serii de valori sunt:
A. Media, dispersia și mediana
B. Media, abaterea standard și modul
C. Media, dispersia și excentricitatea
D. Media, mediana și modul
-