

Elemente de matematică aplicate în biologie

Motto

**Matematica se bucură de o poziție specială în raport cu celelalte științe pentru că legile ei sunt
absolut certe și indiscutabile**

(A. Einstein, Geometry and experience, Sidelight on Relativity, Dover Publication, New York, 1983)

Conf. Univ.Dr. Dana Constantinescu

1. Introducere

Argument

Matematica a câștigat și și-a menținut o poziție excepțională între științe pentru că rezultatele sale sunt obținute dintr-un număr mic de axiome (mai mult sau mai puțin evidente) printr-un lanț de raționamente. Deoarece e bazată pe o logică impecabilă, matematica furnizează științelor naturale un grad înalt de securitate (și claritate) care altfel nu poate fi atins. Din acest motiv, tratarea riguros matematică a acestora este de dorit și se realizează ori de câte ori e posibil. Mai mult decât atât, matematica este un mijloc de comunicare între oameni de știință și ingineri de diverse specialități. Ca rezultat, dacă o anumită ramură a științei este prezentată în formă riguros matematică, accesibilitatea și audiența ei sporește.

(I. D. Mayergoyz, Mathematical Models of hysteresis and their applications, Elsevier Science Inc. New York, 2003)

Deși dezvoltarea biologiei nu a fost influențată în mod esențial de dezvoltarea matematicii, în ultimele decenii este recunoscută importanța completării studiului descriptiv al unor fenomene sau mecanisme biologice cu aspecte legate de prelucrearea și interpretarea datelor obținute. Cea mai avansată formă a folosirii matematicii în biologie este biologia matematică. Ea își propune modelarea matematică a proceselor biologice și studiul modelelor folosind metode specifice matematicii. Pentru construirea și validarea modelelor matematice se pot folosi cercetări statistice.

Statistica dezvoltă tehnici și proceduri de înregistrare, descriere, analiză și interpretare a datelor experimentale sau a rezultatelor obținute din observarea unui proces social, economic, biologic etc., precum și vizualizarea datelor folosind softuri dedicate acestui scop. Cunoașterea unor elemente și principii de bază ale statisticii este importantă în momentul actual, permițând realizarea unor analize corecte a datelor și evitarea erorilor de interpretare a acestora. Strâns legată de statistica inferențială este teoria probabilităților, care furnizează metode și tehnici pentru stabilirea unor previziuni (inferențe statistice) referitoare la caracteristicile unei populații pornind de la rezultatele obținute din observarea unui eșantion al acesteia.

Biostatistica (combinație de cuvinte între biologie și statistică) este aplicarea statisticii într-un număr mare de domenii ale biologiei.

Biostatistica are drept obiectiv și fundamentarea teoretică a proiectării și controlului experimentelor biologice, mai ales în medicină și agricultură, deoarece ea analizează și interpretează date concrete și realizează inferențe asupra acestora.

Se consideră că principalii beneficiari ai biostatisticii sunt

- Sănătatea publică (studiul aspectelor epidemiologice, legate de nutriție, corelarea stării de sănătate și proprietățile mediului înconjurător, organizarea serviciilor de studiu al sănătății populației)
- Ecologia și previziunile ecologice (studiul influenței diverșilor factori asupra dinamicii populațiilor)
- Statistica genetică (studiază legătura între variațiile genotipului și ale fenotipului). Studiul genetic al populațiilor este folosit în agricultură pentru îmbunătățirea soiurilor de plante și animale, iar în genetica umană studiul statistic ajută la identificarea cauzelor care influențează predispoziția la anumite afecțiuni)
- Analiza secvențelor biologice (secvențe AND, secvențe de peptide...)

În cele ce urmează prezentăm unele aplicații directe ale statisticii matematice și ale teoriei probabilităților în descrierea unor fenomene simple ce apar în biologie și agricultură. Asocierea celor două domenii beneficiare ale matematicii nu este întâmplătoare, agricultura fiind în bună măsură biologie aplicată.

2. Aplicații ale statisticii descriptive în biologie și în agricultură

Statistica matematică se ocupă cu descrierea și analiza numerică a fenomenelor (sociale, economice, științifice etc). Statistica operează cu date care se pot colecta din surse existente sau se pot obține prin observații și studii experimentale.

Datele statistice sunt în fapt observații codificate realizate asupra unei mulțimi de elemente de aceeași natură, mulțime care se numește **populație statistică**. O populație poate fi finită sau infinită. Numărul de elemente al unei populații finite se numește *volumul populației*.

Elementele populației (indivizii) sunt purtătoare de informații. Indivizii pot fi persoane (de exemplu formând populația unei localități), agenți economici, obiecte (de exemplu mijloacele fixe ale unui agent economic, piese produse sau comercializate), evenimente (de exemplu operațiuni bancare), opinii (relative la servicii, calitatea unui produs), etc.

Caracteristica populației este trăsătura comună a elementelor sale care este supusă studiului statistic. În statistica matematică ea este cuantificată prin valori numerice. Deoarece o caracteristică variază de la individ la individ, ea poate fi considerată ca o funcție $X: P \rightarrow R$, unde P este populația statistică.

O caracteristică poate fi discretă (dacă valorile sale formează o mulțime finită) sau continuă (în cazul când caracteristica poate lua orice valoare reală).

De exemplu, caracteristica ce indică numărul de piese defecte din fiecare lot este o discretă, în timp ce profitul unei firme sau volumul încasărilor pot fi interpretate ca și caracteristici continue.

Un fenomen deosebit de important este cuantificarea fenomenelor sociale, adică transpunerea în limbaj numeric a caracteristicilor acestor fenomene pentru a înlesni compararea, analiza și sinteza lor, precum și pentru a face prognoze asupra lor.

Problema cuantificării fenomenelor sociale este o problemă de bază a științelor sociale, în condițiile creșterii exigențelor față de determinările științifice ale acestora.

Există *fenomene sociale măsurabile prin natura lor*, de exemplu fenomenele demografice, fenomenele economice, diverse fenomene politice sau culturale

Fenomenele sociale măsurabile cu aproximație se referă în special la opiniile și comportamentele colectivităților umane. În acest caz măsurarea nu poate fi efectuată decât prin compararea intensităților cu care se manifestă acestea la diverse persoane, adică prin realizarea unei scări de mărimi numită *scalogramă*.

Un exemplu de scalogramă care reprezintă intensitatea opiniilor este cea care conține trei niveluri: cu totul de acord, de acord, nu sunt de acord.

Statistica matematică operează cu fenomene cuantificabile numeric, deci fiecărui element al unei scalograme i se asociază un număr.

Demersul statistic are două niveluri: descrierea statistică (statistica descriptivă) și inferența statistică (statistica inferențială).

Statistica descriptivă se ocupă cu înregistrarea, gruparea, prelucrarea și prezentarea datelor obținute prin investigație și pe această bază descrie fenomenul studiat. În studiul statistic descriptiv toate elementele populației sunt luate în considerație. Scopul statisticii descriptive este îndepărtarea detaliilor neimportante și focalizarea atenției asupra unor aspecte de interes și anume:

- precizarea valorii în jurul căreia sunt centrate datele
- descrierea împrăștierea acestora în jurul valorii centrale
- vizualizarea datelor cu ajutorul histogramelor
- analiza corelației între fenomene

Statistica inferențială are ca obiect de studiu investigarea prin sondaj: din întreaga populație se selectează un eșantion reprezentativ asupra căruia se fac măsurători sau observații legate de o anumită caracteristică a populației. Pe baza rezultatelor obținute se fac *inferențe statistice* (adică se formulează concluzii) asupra parametrilor populației. Statistica inferențială folosește deci informația rezultată din studierea unui eșantion pentru a obține concluzii referitoare la întreaga populație din care a fost selectat eșantionul. Aceste concluzii nu sunt de tip determinist ci se obțin folosind metode și tehnici ale teoriei probabilităților, teorie ce conține mecanisme de măsurare și analiză a incertitudinii

legate de evenimentele viitoare. Această incertitudine este exprimată cu ajutorul nivelelor de încredere.

În realizarea unei cercetări statistice se parcurg de obicei următoarele etape:

- *colectarea datelor* care se realizează prin metode specifice obiectivului și condițiilor cercetării. În funcție de tipul de analiză folosit (descriptivă sau inferențială) se folosește întreaga populație sau doar un eșantion.

- *procesarea datelor* înseamnă cuantificarea lor numerică și obținerea seriilor de date.

- *analiza datelor* se realizează prin metode și tehnici specifice statisticii matematice. Această etapă necesită o cunoaștere profundă a filosofiei ce stă în spatele fiecărei metode deoarece este posibil să se obțină rezultate ne semnificative statistic atunci când ipotezele de lucru sau condițiile de aplicare a metodelor nu sunt îndeplinite.

- *interpretarea rezultatelor* este diferită în statistica descriptivă și în cea inferențială. În primul caz se obțin informații concrete și clare despre populația studiată, în al doilea caz validarea rezultatelor obținute este realizată prin compararea cu ce se știa sau se bănuia în domeniul respective. În unele situații analiza statistică dezvăluie corelații între fenomene, legături care ar fi fost greu sau chiar imposibil de observat fără eficientul mecanism statistico-matematic.

În momentul de față există o vastă informație statistică la nivel global, datorată în principal dezvoltării continue a tehnologiei calculatoarelor. Realizarea și folosirea corectă a bazelor de date reprezintă o preocupare importantă în mediul economic și nu numai. Soft-urile statistice joacă un rol important în analiza datelor. Ele îmbină proceduri statistice clasice și moderne cu tehnici de grafică interactivă. Multe soft-uri au două versiuni: una profesională și una academică. Literatura de specialitate califică drept foarte performante, printre altele, următoarele pachete de programe:

- S-PLUS (<http://www.insightful.com/products/splus/>)
- XploRe (<http://www.xplorettech.com/index.pl>)
- Statistica (<http://www.statsoft.com/>)
- SPSS (<http://www.spss.com/>)

2.1. Serii de date și distribuții de frecvențe

Considerăm o populație statistică P finită de volum N pentru care o caracteristică C este codificată de valorile numerice x_1, x_2, \dots, x_N , nu neapărat diferite.

Sirul finit de numere se notează

$$X : x_1, x_2, \dots, x_N$$

și se numește *serie de date*.

Exemplu: $X : 0, 1, 0, 0, 2$ este o serie de date care poate fi interpretată o funcție $X : \{a, b, c, d, e\} \rightarrow \{0, 1, 2\}$, unde $X(a) = 0$, $X(b) = 1$, $X(c) = 0$, $X(d) = 0$, $X(e) = 2$.

În acest caz populația este $P = \{a, b, c, d, e\}$. Deoarece identitatea indivizilor din populație nu este interesantă din punct de vedere statistic, aceasta este neglijată în etapele următoare.

Definiție: Distribuția de frecvențe (sau variabila statistică) asociată caracteristicii C a populației P de volum N este

$$X = \begin{pmatrix} x_1 & x_2 & x_3 & \dots & x_k \\ n_1 & n_2 & n_3 & \dots & n_k \end{pmatrix}$$

unde x_j , $j \in \{1, 2, \dots, k\}$ sunt valorile diferite înregistrate pentru caracteristica C iar n_j , $j \in \{1, 2, \dots, k\}$ reprezintă numărul indivizilor populației caracterizați de valoarea x_j .

Numărul n_j se numește frecvența absolută de apariție a valorii x_j .

Observații: 1. Din definiția frecvențelor relative rezultă că

$$\sum_{j=1}^k n_j = n_1 + n_2 + \dots + n_k = N.$$

2. Unei caracteristici i se poate asocia și **distribuția frecvențelor relative**

$$X_r = \begin{pmatrix} x_1 & x_2 & x_3 & x_k \\ f_1 & f_2 & f_3 & f_k \end{pmatrix}, \quad f_j = \frac{n_j}{N}.$$

În acest caz $\sum_{j=1}^k f_j = 1$. Frecvența relativă f_j poate fi interpretată ca fiind probabilitatea ca valoarea x_j să fie luată de caracteristica C , iar distribuția frecvențelor relative este în fapt o variabilă aleatoare.

Exemplu: Pentru seria de date $X: 0, 1, 2, 5, 2, 3, 3, 2$

distribuția de frecvențe este $X = \begin{pmatrix} 0 & 1 & 2 & 3 & 5 \\ 1 & 1 & 3 & 2 & 1 \end{pmatrix}$ iar cea a frecvențelor relative este

$$X_r = \begin{pmatrix} 0 & 1 & 2 & 3 & 5 \\ 1/8 & 1/8 & 3/8 & 2/8 & 1/8 \end{pmatrix}$$

2.2. Reprezentări grafice

Graficul corespunzător unei serii statistice se numește diagramă. Cazul seriilor pentru care caracteristica este măsurată cantitativ (și exprimată prin numere reale) se întâlnesc în mod current următoarele reprezentări grafice:

- reprezentarea cu segmente verticale;
- histograma cu bare
- poligonul frecvențelor
- reprezentarea cu sectoare circulare

a) **Reprezentarea cu segmente verticale (histograma cu segmente)** se folosește pentru serii cu un număr redus de date, de obicei numere întregi.

Pentru distribuția de frecvențe $X_r = \begin{pmatrix} x_1 & x_2 & x_3 & x_k \\ n_1 & n_2 & n_3 & n_k \end{pmatrix}$, histograma cu segmente, sau reprezentarea cu segmente, este familia de segmente verticale ce unesc punctele de coordonate $(x_i, 0)$ și (x_i, n_i) unde $i \in \{1, 2, \dots, k\}$

Exemplu: Pentru $X = \begin{pmatrix} 1 & 3 & 2 & 4 & 5 \\ 3 & 2 & 4 & 3 & 1 \end{pmatrix}$ reprezentarea cu segmente verticale este prezentată în figura 2.1.

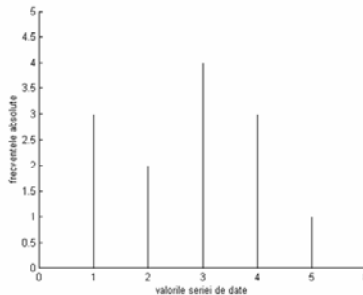


Figura 2.1. Histograma cu segmente

b) **Histograma cu bare** se folosește pentru seriile cu un număr mare de date ce nu sunt neapărat numere întregi. Ea se realizează astfel:

- se determina valoarea minimă, x_{\min} și valoarea maximă x_{\max} a seriei de date

- se divide segmentul $[x_{\min}, x_{\max}]$ prin puncte echidistante cu pasul $h = \frac{x_{\max} - x_{\min}}{n}$, unde n este numărul de intervale ales de analistul seriei. Punctele de diviziune sunt $x_j = x_{\min} + j \cdot h$, unde $j \in \{0, 1, 2, \dots, n\}$
- se calculează câte valori ale seriei aparțin fiecărui interval $I_j = [x_j, x_{j+1})$. Acest număr, notat n_j , se numește frecvența clasei I_j .
- Deasupra fiecărui interval I_j se trasează un dreptunghi cu baza I_j și înălțimea proporțională cu n_j . Pentru determinarea înălțimii dreptunghiului se poate folosi formula $H_j = \frac{n_j}{h \cdot N}$.

Obiecul grafic rezultat din alăturarea acestor dreptunghiuri se numește *histograma cu bare a seriei de date* sau *histograma distribuției de frecvențe*, pentru că ilustrează modul în care sunt distribuite datele. Un exemplu de histogramă cu bare este dat în Figura 2.2.

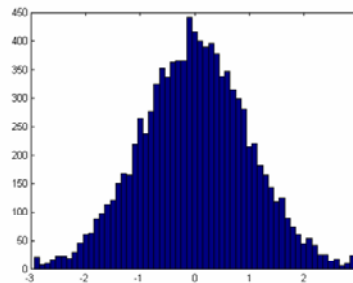


Figura 1.2. Histograma cu bare

O problemă legată de generarea histogramelor este legată de precizarea numărului de intervale de diviziune. În perioada de început a statisticii computaționale numărul de intervale era proporțional cu \sqrt{N} . În unele programe statistice el este ales proporțional cu $\log_2 N$. Cea mai bună idee este să generăm histograma corespunzătoare mai multor numere de intervale și să le comparăm.

c) Poligonul frecvențelor se obține unind vârfurile segmentelor verticale în cazul reprezentării cu segmente. În cazul reprezentării din Figura 2.1, poligonul de frecvențe, A, B, C, D, E este dat în figura 2.3.

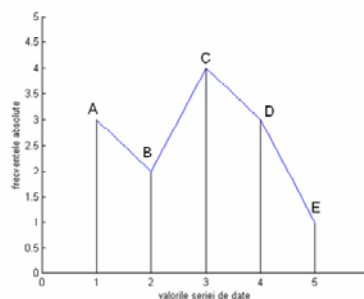


Figura 2.3. Poligon de frecvențe

d) Reprezentarea cu sectoare circulare este folosită pentru obținerea rapidă a unei viziuni globale asupra importanței relative a diverselor clase ale statisticii, interpretarea lor fiind ușurată de colorarea diferită a diverselor clase. În general această reprezentare este folosită pentru seriile cu un număr mic de clase.

Reprezentarea se realizează astfel:

- se determină clasele seriei și numărul de valori ale seriei din fiecare clasă (frecvențele absolute ale claselor)

- pe un cerc se consideră sectoare circulare proporționale cu frecvențele fiecărei clase. Unghiul la centru corespunzător clasei cu frecvența absolută n_j este $\theta_j = \frac{n_j}{360 \cdot N}$.

e) **Reprezentarea polară** se folosește atunci când caracteristica statistică prezintă o anumită periodicitate. De exemplu date înregistrate calendaristic (numarul de nasteri înregistrate în fiecare lună) sau date referitoare la aspecte geografice (intensitatea vântului ce bate din anumite direcții). Ea se construiește astfel: pe semidrepte cu aceeași origine și care împart planul într-un număr de sectoare egale (acest număr se stabilește în funcție de caracterul seriei statistice) se consideră segmente ce pornesc din origine, proporționale cu frecvențele absolute ale claselor și se unesc extremitățile acestor segmente. Se obține un poligon închis în care clasele cu frecvență mai mare sunt reprezentate prin vârfuri aflate la distanță mai mare față de origine.

2.3. Indicatori statistici

2.3.1. Indicatori de poziție (de nivel, de localizare)

a) **media aritmetică** $\bar{x} = \frac{x_1 n_1 + x_2 n_2 + \dots + x_k n_k}{N}$

Media aritmetică este sensibilă față de valorile extreme ale seriei, ea devenind nereprezentativă dacă termenii seriei sunt foarte împrăstiați. Omogenitatea colectivității este o condiție a reprezentativității, pentru orice tip de mărime medie.

b) **media armonică** $x_{arm} = \frac{N}{\frac{n_1}{x_1} + \frac{n_2}{x_2} + \dots + \frac{n_k}{x_k}}$

Media armonică este influențată de prezența valorilor individuale mici și de frecvența acestora.

Media armonică se utilizează pentru exprimarea tendinței centrale în funcție de scopul cercetării și mai ales în funcție de natura obiectivă dintre valorile variabilei numerice observate.

În economie este folosită la calculul productivității, pentru calculul indicelui (sintetic) al prețurilor mărfurilor și tarifelor serviciilor (care sintetizează indicii individuali ai acestor prețuri și tarife).

c) **media geometrică** $x_g = \sqrt[n]{x_1^{n_1} \cdot x_2^{n_2} \cdot \dots \cdot x_k^{n_k}}$

Media geometrică este folosită mai rar ca indicator statistic, îndeosebi când termenii prezintă o evidentă concentrare către valorile cele mai mici sau când se urmărește să se acorde o importanță deosebită valorilor individuale reduse.

Dacă cel puțin o valoare individuală este nulă sau negativă, calculul mediei geometrice este lipsit de sens. Ea nu poate fi folosită dacă în cadrul seriei există cel puțin un termen negativ, deoarece expresia devine imaginară.

Media geometrică mai este denumită și *medie de ritm*, fiind folosită pentru calculul ritmului mediu de creștere. Un exemplu de folosire a mediei geometrice ca indicator statistic este dat în exemplul următor:

Exemplu O colonie de microorganisme a fost studiată pe parcursul a două zile. S-a constatat că masa sa inițială era 10 g, după o zi era 20 g iar a treia zi era 160 g. Să se calculeze ritmul mediu de creștere al coloniei.

Masa coloniei s-a dublat în prima zi și s-a multiplicat de 8 ori în a doua zi. Dacă se calculează rapid media aritmetică se constată că, în medie, ritmul de creștere este $\frac{2+8}{2} = 5$.

Acest rezultat în mod evident este incorect deoarece, în acest caz după o zi colonia ar avea , $10 \cdot 5 = 50g$, iar după două zile ar avea $50 \cdot 5 = 250g$, ceea ce nu este adevărat

Dimpotrivă, dacă indicele mediu de dinamică se determină ca media geometrică a dinamicilor individuale se obține următoarea valoare: $\bar{x}_g = \sqrt{2 \cdot 8} = 4$. Acesta este un rezultat mult mai corect decât cel anterior deoarece pornind de la 10 g colonia ar avea $10 \cdot 4 = 40g$ după prima zi (ceea ce nu e adevărat) și $40 \cdot 4 = 160g$ după a doua zi. Acest rezultat verifică datele problemei, deci ritmul mediu de creștere este egal cu media geometrică a ritmurilor intermediare de creștere, adică este 4.

d) **mediانا** seriei de date $X : x_1, x_2, \dots, x_N$ cu termenii ordonați crescător este numărul

$$me = \begin{cases} \frac{x_{N+1}}{2} & \text{daca } N \text{ este impar} \\ \frac{x_{N/2} + x_{1+N/2}}{2} & \text{daca } N \text{ este par} \end{cases}$$

Mediana este o valoare ce caracterizează “centrul” seriei de date. În cazul când N este par mediana nu este obligatoriu valoare a seriei de date.

Are proprietatea că suma frecvențelor valorilor mai mici ca me este egală cu suma frecvențelor mai mari ca me .

Este utilizată în studiul fertilității, mortalității, determinarea duratei de viață.

e) **modul** (moda su dominanta) este valoarea cu cea mai mare frecvență de apariție (care este la modă). Există repartiții unimodale (cu un singur mod), bimodale (cu două moduri) etc.

Valoare modală este influențată de mărimea valorilor din centrul seriei (la distribuțiile unimodale) sau din centrul îngrămădirii de observații (la distribuțiile plurimodale). Celelalte valori nu au nici o influență asupra ei.

Distribuțiile bimodale (cu două frecvențe maxime) reprezintă o situație rar întâlnită, care impune separarea unităților colectivității în două distribuții de frecvențe.

2.3.2. Indicatorii variației (împrăștierii)

Indicatorii tendinței centrale nu dau nici o explicație asupra împrăștierii, respectiv a modului în care termenii seriei se abat între ei sau de la medie. Astfel, apare necesitatea calculării unor noi indicatori care rezolvă:

- verificarea reprezentativității mediei ca valoare tipică a seriei de distribuție;
- verificarea gradului de omogenitate al seriei;
- verificarea sistematizării informațiilor prin gruparea statistică;
- caracterizarea gradului și formei de variație a unei variabile statistice.

Acești indicatori care dau o caracterizare precisă a unei serii statistice prin care se poate cunoaște variația valorilor individuale (cum se grupează aceste valori în jurul valorii medii, dacă sunt apropiate sau îndepărtate de această valoare), se numesc *indicatorii variației*. Ei sunt:

a) **amplitudinea** este diferența dintre cea mai mare și cea mai mică valoare a seriei de date (sau a distribuției de frecvențe)

b) **abaterea medie absolută** $e_x = \frac{1}{N} \sum_{i=1}^k n_i |x_i - \bar{x}|$

c) **varianța (dispersia)** $s^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$

d) **abaterea medie pătratică (standard)** $s = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$

Propoziție Dispersia și abaterea medie pătratică ale unei distribuții de frecvențe

$X = \begin{pmatrix} x_1 & x_2 & x_3 & \dots & x_k \\ n_1 & n_2 & n_3 & \dots & n_k \end{pmatrix}$, unde $\sum_{i=1}^k n_i = N$ se calculează folosind formulele

$$s^2 = \frac{\sum_{i=1}^k x_i^2 \cdot n_i}{N} - \left(\frac{\sum_{i=1}^k x_i \cdot n_i}{N} \right)^2, \text{ respectiv } s = \sqrt{\frac{\sum_{i=1}^k x_i^2 \cdot n_i}{N} - \left(\frac{\sum_{i=1}^k x_i \cdot n_i}{N} \right)^2}.$$

Dispersia este un indice de variație ce dă indicații privind împrăștierea valorilor seriei în jurul valorii medii. Cu cât este mai mică dispersia, cu atât valorile seriei statistice se grupează mai mult în jurul valorii medii. În acest caz media este un indicator statistic relevant pentru studiul seriei. O dispersie mare arată că elementele eșantionului au o împrăștiere mare și valoarea medie nu dă informații relevante despre serie.

Dispersia este influențată de mărimea valorilor din seria de date. Dacă valorile sunt mari, dispersia poate fi și ea mare, dar cazul seriilor de date cu valori mici dispersia poate avea valori mici chiar dacă datele nu sunt grupate în jurul mediei. De aceea, pentru studiul împrăștierii se folosește coeficientul de variație care nu este influențat în mod esențial de mărimea termenilor seriei de date.

e) **coeficientul de variație** $CV = \frac{s}{\bar{x}}$

Coeficientul de variație are valori cuprinse în intervalul [0,1]. El este cel mai sintetic indicator al împrăștierii.

Cu cât coeficientul de variație e mai aproape de 0, cu atât seria este mai omogenă și media este mai reprezentativă. Dacă este mai apropiat de 1, împrăștieria valorilor este mare și media nu este un indicator reprezentativ.

Practica utilizării coeficientului de variație a stabilit pragul de trecere de la starea de omogenitate la cea de eterogenitate: În literatura de specialitate se avansează nivelul de 35 - 40 % ca limită maximă admisibilă pentru coeficientul de variație.

- Dacă $CV \leq 0.35$, populația este *omogenă* și media este un indicator relevant.
- Dacă $CV > 0.35$, populația este *eterogenă* și media nu este un indicator relevant

În analizele financiare coeficientul de variație este o măsură a riscului relativ.

Exemple

1 Cantitatea de deșuri organice produse la o ferma în decursul a 100 zile consecutive a fost înregistrată în tabelul de mai jos

Cantitatea de deșuri produse zilnic "x _i "	Numarul de zile în care s-a produs cantitatea de deșuri "n _i "	Frecvența relativă n _i / 100
0	5	
1	15	
2	23	
3	22	
4	16	
5	9	
6	5	
7	5	

- Să se completeze coloana frecvențelor relative;
- Să se deseneze histograma cu segmente verticale asociată datelor din table.
- Să se calculeze indicatorii de poziție (media, mediana, modul) și indicatorii de împrăștiere (dispersia, abaterea standard și coeficientul de variație)
- Să se interpreteze datele obținute

Rezolvare: a) $X = \left(\begin{array}{cccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ \frac{5}{100} & \frac{15}{100} & \frac{23}{100} & \frac{22}{100} & \frac{16}{100} & \frac{9}{100} & \frac{5}{100} & \frac{5}{100} \end{array} \right)$

- b) Histograma cu segmente este

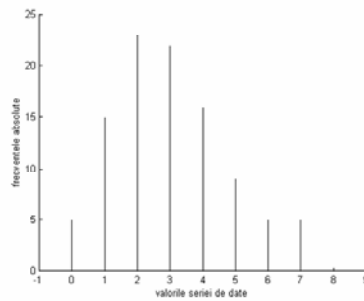


Figura 2.4 Histograma cu segmente a seriei de date din exercițiul 1

c) Indicatorii de poziție sunt:

- media $\bar{x} = \frac{0 \cdot 5 + 1 \cdot 15 + 2 \cdot 23 + 3 \cdot 22 + 4 \cdot 16 + 5 \cdot 9 + 6 \cdot 5 + 7 \cdot 5}{100} = 3.01$

- mediana se calculează ținând cont ca sunt 100 termeni în serie. Dacă scriem termenii seriei în ordine crescătoare, repetându-i de atâtea ori cât indică frecvența absolută obținem $x_{50} = x_{51} = 3$. Deci

$me(X) = \frac{x_{50} + x_{51}}{2} = \frac{3 + 3}{2} = 3.$

-modul este $mo(X) = 2$ pentru că această valoare are cel mai mare număr de apariții.

Indicatorii de poziție sunt

-dispersia: $s^2 = \frac{1}{100} \cdot [5 \cdot 0^2 + 15 \cdot 1^2 + 23 \cdot 2^2 + 22 \cdot 3^2 + 16 \cdot 4^2 + 9 \cdot 5^2 + 5 \cdot 6^2 + 5 \cdot 7^2 - \frac{2.85^2}{100}] = 3.0499$

- abaterea standard $s = \sqrt{s^2} = \sqrt{3.0499} = 1.7463$

-coeficientul de variație $CV = \frac{s}{\bar{x}} = 0.58.$

2. Vârsta persoanelor dintr-o comunitate a fost înregistrată și datele au fost grupate în tabelul de mai jos.

Vârsta (în ani)	Număr persoane
[0,5)	5
[5,10)	12
[10,15)	33
[15,20)	71
[20,25)	119
[25,30)	175
[30,35)	185
[35,40)	158
[40,45)	122
[45,50)	69
[50,55)	35
[55,60)	11
≥60	5
Total	1000

a) Să se deseneze histograma cu bare a acestei serii de date (vârstele mai mari de 60 ani se identifică cu intervalul [60,65).

b) Identificând fiecare interval cu mijlocul său, să se constituie seria statistică a vârstelor celor 1000 de persoane din comunitate. Să se determine media, mediana și dispersia acestei serii.

Rezolvare:

a) histograma este

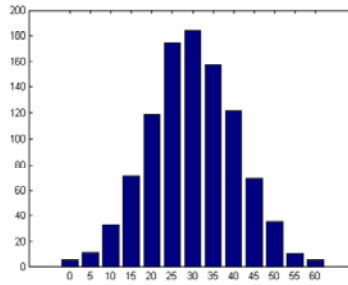


Figura 2.5 Histograma cu bare a seriei de date din exercițiul 2

b) Seria de date este

$$X = \begin{pmatrix} 2.25 & 7.25 & 12.25 & 17.25 & 22.25 & 27.25 & 32.25 & 37.25 & 42.25 & 47.25 & 52.25 & 57.25 & 62.25 \\ 5 & 12 & 33 & 71 & 119 & 175 & 185 & 158 & 122 & 69 & 35 & 11 & 5 \end{pmatrix}$$

Media este $\bar{x} = (2.25 \cdot 5 + 7.25 \cdot 12 + 12.25 \cdot 33 + 17.25 \cdot 71 + 22.25 \cdot 119 + 27.25 \cdot 175 + 32.25 \cdot 185 + 37.25 \cdot 158 + 42.25 \cdot 122 + 47.25 \cdot 69 + 52.25 \cdot 35 + 57.25 \cdot 11 + 62.25 \cdot 5) / 1000 = 32.18$

Mediana este $me = \frac{x_{500} + x_{501}}{2} = \frac{32.25 + 32.25}{2} = 32.25$

Dispersia este $s^2 = 114.4950$.

Abaterea standard este $s = \sqrt{s^2} = 10.7002$

3. Statistica nașterilor înregistrate lunar într-o localitate este prezentată în tabelul următor

Luna	01	02	03	04	05	06	07	08	09	10	11	12
Nr. Nașteri	8	9	13	18	15	20	24	19	12	11	6	5

a) Să se reprezinte seria de date cu ajutorul histogramei

b) Să se calculeze indicatorii seriei de date și să se interpreteze rezultatele

a) Histograma este

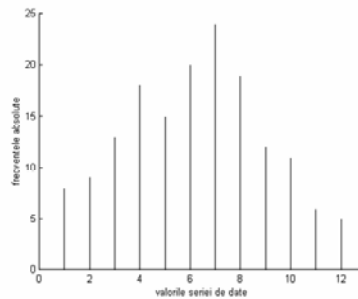


Figura 2.6 Histograma cu bare a seriei de date din exercițiul 3

b) Media este $\bar{n} = 13,3$.

Dispersia este $s^2 = 32.7222$.

Abaterea standard este $s = 5.7203$.

Coeficientul de variație este $CV = \frac{s}{\bar{n}} = 0.4290$

Deoarece coeficientul de variație este mare rezultă că media nu este un indicator reprezentativ pentru seria de date.

4. Frecvența medie a vântului pe direcțiile principale și secundare ale punctelor cardinale înregistrate la Stația meteorologică Craiova în perioada 1950-2000 este dată în tabelul următor

Direcția	N	NE	E	SE	S	SV	V	NV
Frecvența (%)	5	10	24	7	5	13	27	9

- a) Să se reprezinte seria de date cu ajutorul histogramei
 b) Să se calculeze indicatorii seriei de date și să se interpreteze rezultatele

a) Histograma este

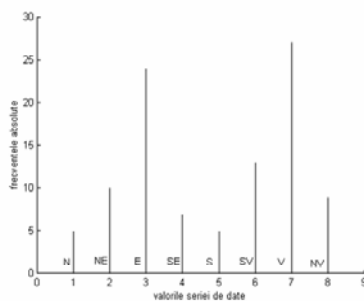


Figura 2.7 Histograma cu segmente a seriei de date din exercițiul 4

b) Media este $\bar{n} = 12.5$

Interpretarea sa este: în medie vântul a bătuț din fiecare direcție 12,5% din timp

Dispersia este $s^2 = 63$

Abaterea standard este $s = 7.9362$

Coeficientul de variație este $CV = 0.6349$

Deoarece coeficientul de variație este mare rezultă că media nu este un indicator statistic relevant.

5. Măsurătorile efectuate prin sondaj aleator asupra înălțimii a 50 de spice dintr-un lot de orz indică următoarele valori (în cm.) date în tabelul de mai jos:

Nr. crt	Înălțime	Nr. crt	Înălțime	Nr. crt	Înălțime	Nr. crt	Înălțime	Nr. crt	Înălțime
1	50,7	11	50,1	21	50,0	31	49,8	41	49,9
2	51,0	12	50,0	22	50,0	32	50,5	42	50,2
3	51,0	13	50,1	23	49,9	33	49,6	43	49,8
4	49,6	14	50,0	24	50,2	34	50,4	44	49,9
5	49,8	15	49,9	25	50,0	35	50,2	45	50,1
6	49,2	16	50,3	26	49,7	36	50,6	46	50,0
7	50,0	17	50,0	27	50,3	37	49,6	47	49,9
8	49,8	18	50,2	28	49,2	38	49,3	48	49,8
9	49,8	19	49,4	29	50,0	39	49,5	49	50,1
10	49,9	20	49,8	30	50,1	40	50,0	50	50,2

a). Să se facă gruparea datelor și să se determine frecvențele absolute și relative. Să se facă reprezentarea în batoane.

b). Să se reprezinte histograma.

c). Să se determine clase de valori de lungime 0,3, să se determine frecvențele absolute ale intervalelor și să se reprezinte histograma cu bare.

d). Să se determine valorile centrale ale claselor, media, valoarea modală, mediana, dispersia și abaterea medie pătratică.

a) Distribuția de frecvențe a seriei de date este

$$X = \begin{pmatrix} 49.2 & 49.3 & 49.4 & 49.5 & 49.6 & 49.7 & 49.8 & 49.9 & 50 & 50.1 & 50.2 & 50.3 & 50.4 & 50.5 & 50.6 & 50.7 & 50.8 & 51 \\ 1 & 1 & 1 & 1 & 2 & 1 & 6 & 6 & 10 & 5 & 6 & 2 & 1 & 1 & 2 & 1 & 1 & 2 \end{pmatrix}$$

Frecvențele relative sunt date de

$$X_{rel} = \begin{pmatrix} 49.2 & 49.3 & 49.4 & 49.5 & 49.6 & 49.7 & 49.8 & 49.9 & 50.0 & 50.1 & 50.2 & 50.3 & 50.4 & 50.5 & 50.6 & 50.7 & 50.8 & 51 \\ \frac{1}{50} & \frac{1}{50} & \frac{1}{50} & \frac{1}{50} & \frac{2}{50} & \frac{1}{50} & \frac{6}{50} & \frac{6}{50} & \frac{10}{50} & \frac{5}{50} & \frac{6}{50} & \frac{2}{50} & \frac{1}{50} & \frac{1}{50} & \frac{2}{50} & \frac{1}{50} & \frac{1}{50} & \frac{2}{50} \end{pmatrix} \text{ b)}$$

Histograma este

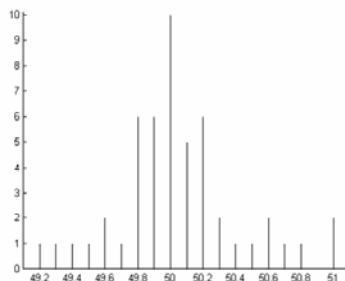


Figura 2.2 Histograma cu segmente a seriei de date din exercitiul 5

c) Clasele sunt date în tabelul următor

Clasa	Frecvența	Valoarea centrală a clasei
[49.2 49.5)	3	49.35
[49.5 49.8)	4	49.65
[49.8 50.1)	22	49.95
[50.1 50.4)	13	50.25
[50.4 50.7)	4	50.55
[50.7 51]	4	50.85

Histograma cu bare este

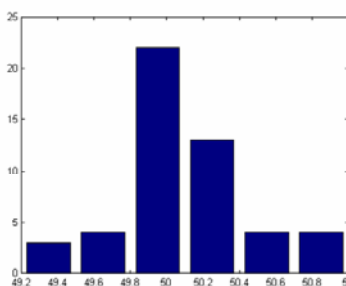


Figure 2.3 Histograma cu bare a grupării de date din exercițiul 5

Distribuția de frecvențe pentru care se calculează indicatorii este

$$X = \begin{pmatrix} 49.35 & 49.65 & 49.95 & 50.25 & 50.55 & 50.85 \\ 3 & 4 & 22 & 13 & 4 & 4 \end{pmatrix}$$

Media este $\bar{X} = 50.0880$.

Modul este $mo = 49.95$.

$$\text{Mediana este } me = \frac{X(25) + X(26)}{2} = \frac{49.95 + 49.95}{2} = 49.95.$$

Dispersia este $s^2 = 0.1264$.

$$\text{Abaterea medie pătratică este } s = \sqrt{s^2} = \sqrt{0.1264} = 0.3560$$

Coeficientul de variație este $CV = 0.0071 \approx 0$, deci datele sunt grupate în jurul valorii medii și media este un indicator relevant.

Exerciții propuse

1. Producția de boabe a 100 de parcele de 6 m² cultivate cu un anumit soi de grâu într-un câmp experimental este dată în tabelul:

Nr. crt	Prod.	Nr. crt	Prod.	Nr. crt	Prod.	Nr. crt	Prod.	Nr. crt	Prod.
1	2,72	21	2,97	41	3,05	61	3,11	81	3,22
2	2,76	22	2,98	42	3,05	62	3,11	82	3,23
3	2,84	23	2,98	43	3,06	63	3,13	83	3,24
4	2,85	24	2,99	44	3,06	64	3,13	84	3,24
5	2,87	25	2,99	45	3,07	65	3,13	85	3,25
6	2,87	26	3,00	46	3,07	66	3,13	86	3,25

7	2,88	27	3,01	47	3,07	67	3,14	87	3,25
8	2,90	28	3,01	48	3,07	68	3,14	88	3,25
9	2,91	29	3,01	49	3,08	69	3,14	89	3,27
10	2,93	30	3,02	50	3,08	70	3,15	90	3,28
11	2,93	31	3,02	51	3,08	71	3,15	91	3,29
12	2,93	32	3,02	52	3,09	72	3,15	92	3,29
13	2,94	33	3,03	53	3,09	73	3,16	93	3,31
14	2,94	34	3,03	54	3,09	74	3,17	94	3,31
15	2,95	35	3,04	55	3,09	75	3,17	95	3,33
16	2,95	36	3,04	56	3,10	76	3,17	96	3,34
17	2,96	37	3,04	57	3,10	77	3,19	97	3,36
18	2,96	38	3,04	58	3,10	78	3,19	98	3,37
19	2,96	39	3,04	59	3,11	79	3,21	99	3,39
20	2,97	40	3,05	60	3,11	80	3,21	100	3,41

Se cere:

- Să se facă gruparea datelor pe clase de lungime 0,05, să se întocmească histograma și să se deseneze poligonul frecvențelor.
- Să se determine valorile centrale ale claselor, media, valoarea modală, mediana, dispersia și abaterea medie pătratică.

2. Temperaturile medii înregistrate la Craiova în lunile *mai* ale anilor 1930-1979 sunt date în tabelul de mai jos:

Anul	0	1	2	3	4	5	6	7	8	9
1930...	8,1	4,0	-0,9	3,2	8,2	6,7	8,8	5,6	7,8	4,1
1940...	3,5	6,3	+0,4	4,3	3,8	6,4	6,4	8,2	5,9	0,3
1950...	5,5	6,9	-1,9	5,1	2,1	3,6	0,0	6,2	2,9	6,0
1960...	4,6	8,0	2,3	2,9	3,2	3,7	6,1	6,6	5,5	-0,1
1970...	5,2	3,6	5,5	3,0	4,9	7,7	3,1	7,2	5,8	6,3

- Să se facă gruparea în clase, de mărime 2°C cu convenția ca extremitatea dreaptă a fiecărei clase să nu aparțină clasei (ex. $[-2,0;0)$, $[0;2,0)$, $[2,0;4,0)$, ...);
- Să se completeze tabela obținută la punctul a) cu frecvențele absolute, cu frecvențele relative și cu valoarea centrală a clasei;
- Să se reprezinte histograma grupării în clase;
- Să se calculeze indicatorii distribuției de frecvențe a grupării datelor în clase și să se interpreteze rezultatele.

3. Cantitățile lunare de precipitații căzute la Craiova în lunile aprilie ale anilor 1930-1979 sunt date (în litri/m.p.) în tabelul următor

	1930...	1940...	1950...	1960...	1970...
0	55,5	92,0	24,8	39,4	64,4
1	19,6	36,5	40,0	49,4	42,5
2	17,8	33,7	40,8	75,6	16,4
3	7,8	26,9	23,5	33,7	42,6
4	89,0	42,3	52,2	62,6	74,0
5	32,7	35,4	94,3	57,9	43,8
6	22,6	16,3	31,6	65,8	47,1
7	45,3	22,8	65,3	49,5	50,2
8	57,1	37,6	51,4	8,7	31,6
9	28,1	3,9	19,3	31,9	42,7

- Să se facă gruparea în clase, de mărime 10 litri/mp.
- Să se completeze tabela obținută la punctul a) cu frecvențele absolute, cu frecvențele relative și cu valoarea centrală a clasei;
- Să se reprezinte histograma grupării în clase și să se calculeze indicatorii statistici ai grupării în clase și să se interpreteze rezultatele.

3. Studiul statistic al legăturii dintre fenomene. Aplicații.

Elementele unei populații pot avea diverse caracteristici, fiecare determinând anumite variabile aleatoare X, Y, \dots , acestea având fie un caracter determinist, fie un caracter întâmplător (stochastic) iar între ele putând exista anumite legături.

Legăturile dintre caracteristicile unei populații pot fi foarte strânse, exprimate prin funcții $y=f(x)$, numite *funcționale*.

Există însă și legături în care intervin numeroși factori sistematici și accidentali care fac ca două sau mai multe însușiri (caracteristici) să varieze în strânsă concordanță (nu însă în sens funcțional). Între acestea sunt legăturile dintre fenomenele și procesele economice care apar ca legături statistice (*stochastice*), a căror particularitate este faptul că rezultatul este determinat ca urmare a influenței unui ansamblu de factori. Legăturile statistice se manifestă, ca tendință valabilă numai la nivelul populației. Dependența de acest tip are caracter întâmplător și se numește *dependență stochastică* sau *corelație*.

În cele ce urmează vom considera fenomene descrise cu ajutorul seriilor de date (exprimate prin numere reale) sistematizate cu ajutorul distribuțiilor de frecvențe (numite și variabile statistice).

Există două aspecte ale studiului dependenței stochastice între fenomene: analiza de corelație și analiza de regresie.

Analiza de corelație studiază comportarea fiecărei variabile în funcție de valorile celorlalte variabile, precum și măsura dependenței dintre variabilele considerate. Se analizează dacă tendința ascendentă a unei variabile implică o tendință ascendentă sau descendentă la cealaltă, sau nici o tendință. Rezultatele se exprimă prin *coeficientul de corelație* sau prin *raportul de corelație*.

Analiza regresiiilor constă în determinarea *funcției de regresie* între două variabile. În ipoteza existenței unei legături între variabile se pot prognoza valorile uneia în raport cu valorile celeilalte folosind funcția de regresie.

În paragrafele următoare va fi studiată legătura directă între serii de date (care generează variabilele statistice) care descriu anumite caracteristici ale unei populații.

Pentru simplificare le vom nota

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

3.1. Analiza corelațiilor

Prin *corelație simplă* se înțelege legătura reciprocă dintre două variabile X și Y ale unei populații.

Corelațiile dintre variabile prezintă mare importanță, deoarece cunoscând variația unei însușiri putem trage concluzii asupra însușirii sau însușirilor de care aceasta este legată, fără a recurge la determinări directe.

Corelația poate fi *pozitivă*, atunci când valorile celor două variabile cresc sau descresc în același timp, sau *negativă*, atunci când valorile unei variabile cresc, iar cele ale celeilalte variabile descresc.

Metodele cele mai simple de constatare a unei corelații sunt *metoda grafică* sau *graficul de corelație (corelograma)* și *tabela de corelație*.

1.1. Metoda grafică (diagrama de împrăștiere)

Perechile de observații (x_i, y_i) , $i \in \{1, 2, \dots, n\}$ se reprezintă în planul Oxy prin punctele $M_i(x_i, y_i)$, $i \in \{1, 2, \dots, n\}$. Se obține « un nor de puncte », numit *corelogramă*. Tendința norului de puncte permite vizualizarea și stabilirea formei analitice a funcției de regresie. Corelograma arată dacă între cele două variabile există o relație și poate indica și forma legăturii prin observarea unei densități de puncte care se concentrează în jurul unei anumite curbe, care poate fi liniară sau de altă formă.

Dacă norul de puncte are forma unei elipse alungite există o legătură puternică între variabilele X și Y . Dacă norul e răspândit în interiorul unui cerc, pătrat variabilele sunt independente.

Exemplul 1: Pentru seriile de date

$$X = \{1.2, 0.8, 1.1, 3.0, 0.7, 0.8, 1.0, 0.6, 0.9, 1.1, 0.65, 0.75, 0.85, 0.95, 1.05, 1.1, 1.25\}$$

$$Y = \{10.1, 9.2, 11.0, 12.0, 9.0, 8.2, 9.35, 9.1, 10.5, 8.8, 8.5, 9.8, 8.5, 9.5, 9.5, 10, 11\}$$

norul de date este reprezentat in Figura 3.1.

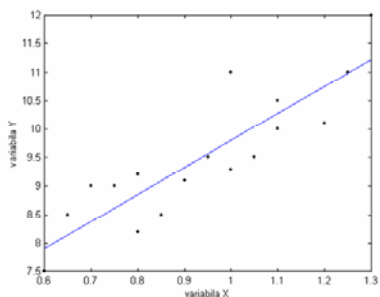


Figura 3.4 Norul de puncte al seriilor de date din Exemplul 1 și dreapta de corelație

Deoarece forma norului de puncte este apropiată de o elipsă se poate considera că seriile de date sunt puternic corelate.

Exemplul 2 Pentru seriile de date

$$X = \{1.2, 0.8, 1, 1.3, 0.7, 0.8, 1.0, 0.6, 0.9, 1.1, 0.65, 0.75, 0.85, 0.95, 1.05, 1.1, 1.25\}$$

$$Y = \{8.1, 10.2, 7.0, 10.0, 9.0, 7.2, 9.3, 8.5, 8.1, 9.5, 9.5, 8, 9.5, 7.5, 8.5, 10, 8\}$$

norul de date este reprezentat in Figura 3.2.

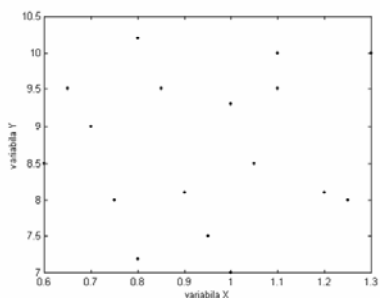


Figura 3.5 Norul de puncte al seriilor de date din Exemplul 2

Faptul că norul de puncte e răspândit în interiorul unui dreptunghi poate fi interpretat ca lipsa unei corelații între cele două variabile.

În cazul probelor cu volum mare de valori observate, pentru cercetarea legăturii dintre variabile se folosește *tabelul de corelație* care constă în gruparea pe clase a datelor de observație. În tabelul de corelație termenul $x_{i,j}$ reprezintă numărul de membrii ai populației pentru care variabila X are valoarea x_i iar variabila Y are valoarea y_j

Cu cât valorile individuale din tabelul de corelație sunt mai strâns concentrate în jurul diagonalei cu atât corelația este mai puternică. Cu cât corelația este mai puternică, cu atât valorile din tabelul de corelație sunt mai strâns concentrate în jurul unei diagonale.

Exemplul 3: ([1], pag 284) În tabelul de mai jos sunt trecute datele privind diametrul tulpinii unei plante și procentul de fibre în funcție de diametru:

$y = \text{conținut de fibre } (\%)$	$x = \text{diametrul tulpinii (mm)}$							
	2	3	4	5	6	7	8	
26	2	3	3	2				10
24	4	5	13	7	4			33
22	3	6	18	25	10	2		64
20		1	8	17	18	3		47
18			1	9	8	8	2	28
16				2	3	4	6	15
14						1	2	3
Suma x	9	15	43	62	43	18	10	200

Se observă că între cele două caracteristici există corelație pentru că valorile din tabel sunt concentrate în jurul diagonalei secundare. Corelația este negativă deoarece valorilor mai mari ale variabilei X le corespund valori mai mici ale variabilei Y , adică tendința ascendentă a lui X conduce la o tendință descendentă a lui Y .

Aceste observații intuitive reprezintă o informație primară despre corelație, descrierea ei corectă fiind realizată cu ajutorul coeficientului de corelație și al raportului de corelație.

Pentru seriile de date $X : x_1, x_2, \dots, x_n$ și $Y : y_1, y_2, \dots, y_n$ considerăm $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ și $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$.

Coeficientul de corelație (numit și coeficientul Pearson) se definește prin

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Pentru calcule directe se poate folosi formula

$$r = \frac{\sum_{i=1}^n x_i \cdot y_i - \frac{\left(\sum_{i=1}^n x_i\right) \cdot \left(\sum_{i=1}^n y_i\right)}{n}}{\sqrt{\left(\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}\right) \cdot \left(\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}\right)}}$$

Următoarele observații reprezintă elemente de bază pentru interpretarea coeficientului de corelație

- Coeficientul de corelație r este o mărime adimensională a cărei valoare absolută este subunitară, adică $|r| < 1$.
- Dacă seriile de date X și Y sunt independente atunci $r = 0$.
- Dacă coeficientul de corelație este nul, seriile statistice nu sunt în mod necesar independente, dar dependența lor nu este liniară, ea putând fi de altă natură.
- Dacă $r \approx -1$ corelația este puternică și negativă (creșterea valorilor lui X este asociată cu descreșterea valorilor lui Y)
- Dacă $r \approx +1$ corelația este puternică și pozitivă (creșterea valorilor lui X este asociată cu creșterea valorilor lui Y)

Folosirea coeficientului de corelație este recomandabilă îndeosebi atunci când legătura dintre variabile nu se abate mult de la liniaritate, iar populația studiată este de tipul distribuțiilor normale bidimensionale, adică, în cazul când datele studiate aparțin unei distribuții bidimensionale normale și relația dintre variabile este liniară coeficientul de corelație are un înțeles statistic bine definit. Dimpotrivă, dacă populația pe care o reprezintă datele nu este normală sau dacă din graficul de corelație este evident că relația dintre variabile se abate mult de la liniaritate, coeficientul de corelație r își pierde înțelesul său statistic, iar examinarea semnificației sale statistice devine lipsită de sens. Coeficientul empiric de corelație r rămâne astfel numai o mărime de calcul și nu o valoare estimativă.

Pragul de încredere pentru interpretarea coeficientului de corelație este definit prin

$$PI = |r| \cdot \sqrt{n-1}.$$

Se consideră că legătura dintre variabile este sufficient de probabilă dacă $PI \geq 3$.

Pentru prezentarea corelației între două fenomene se procedează astfel

- se realizează diagrama de împrăștiere a norului de puncte și se observă în mod empiric dacă datele sunt corelate.

- Dacă variabilele sunt corelate și corelația e aproape liniară (norul de puncte se află în interiorul unei elipse alungite) se calculează coeficientul de corelație și pragul de încredere și se interpretează rezultatele.

Exemplul 4 ([1], pag 288) În urma efectuării a 8 măsurători asupra două caracteristici X și Y ale unei populații, s-au găsit valorile date în tabelul de mai jos:

Proba	1	2	3	4	5	6	7	8
$X: x_i$	26,9	26,3	23,6	24,8	29,1	19,6	17,9	19,5
$Y: y_i$	54,0	52,2	55,5	57,1	54,3	63,2	70,1	70,2

Să se determine coeficientul de corelație al variabilelor X și Y .

Norul de puncte corespunzător seriilor de date X și Y este reprezentat în Figura 3.3.

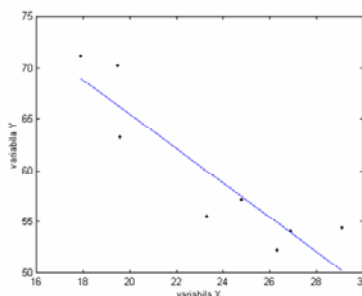


Figura 3.6 Norul de puncte di Exemplul 4 și dreapta de corelație

Configurația norului de puncte indică o corelație liniară negativă.

Pentru calculul coeficientului de corelație așezăm datele în tabelul de mai jos, pe primele două coloane. Celelalte coloane se completează folosind datele problemei.

x_i (cm)	y_i (%)	x_i^2	y_i^2	$x_i y_i$
26,9	54,0	723,61	2916,00	1452,60
26,3	52,2	691,69	2724,84	1372,86
23,6	55,5	556,96	3080,25	1309,80
24,8	57,1	615,04	3260,41	1416,08
29,1	54,3	846,81	2948,48	1580,13
19,6	63,2	384,16	3994,24	1238,72
17,9	70,1	320,41	4914,01	1254,79
19,5	70,2	380,25	4928,04	1368,90
$\sum x_i = 187,7$	$\sum y_i = 476,6$	$\sum x_i^2 = 4518,93$	$\sum y_i^2 = 28766,28$	$\sum x_i y_i = 10993,88$
$\bar{x} = 23,46$	$\bar{y} = 59,57$	$\overline{x^2} = 564,86$	$\overline{y^2} = 3595,78$	$\overline{x \cdot y} = 1374,23$

Efectuând calculele necesare obținem:

$$\bar{x} \cdot \bar{y} = 139751,22, \quad \bar{x}^2 = 550,3716, \quad \bar{y}^2 = 3548,5849;$$

$$s_x^2 = \frac{n}{n-1} \left[(\overline{x^2}) - (\bar{x})^2 \right] = \frac{8}{7} (564,86 - 550,37) = 16,56;$$

$$s_y^2 = \frac{n}{n-1} \left[(\overline{y^2}) - (\bar{y})^2 \right] = \frac{8}{7} (3595,78 - 3548,58) = 53,94;$$

$$s_{xy} = \frac{n}{n-1} (\overline{x \cdot y} - \bar{x} \cdot \bar{y}) = \frac{8}{7} (1374,23 - 1397,51) = -26,61;$$

Rezultă $r = \frac{s_{xy}}{s_x \cdot s_y} = -0,89$, ceea ce indică o corelație negativă puternică

Pragul de încredere este $PI = 0,89 * \sqrt{7} = 2,3547 < 3$, deci coeficientul de corelație nu este relevant.

Exemplul 5 (Exemplul 3 continuat) În tabelul de mai jos sunt trecute datele privind diametrul tulpinii unei plante și procentul de fibre în funcție de diametru:

$y = \text{conținut de fibre } (\%)$	$x = \text{diametrul tulpinii (mm)}$							
	2	3	4	5	6	7	8	
26	2	3	3	2				10
24	4	5	13	7	4			33
22	3	6	18	25	10	2		64
20		1	8	17	18	3		47
18			1	9	8	8	2	28
16				2	3	4	6	15
14						1	2	3
Suma x	9	15	43	62	43	18	10	200

Cele două variabile statistice pentru care se cere coeficientul de corelație sunt

$$X = \begin{pmatrix} 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 9 & 15 & 43 & 62 & 43 & 18 & 10 \end{pmatrix} \text{ și } Y = \begin{pmatrix} 26 & 24 & 22 & 20 & 18 & 16 & 14 \\ 10 & 33 & 64 & 47 & 28 & 15 & 3 \end{pmatrix}$$

$$\sum x_i = 2 \cdot 9 + 3 \cdot 15 + 4 \cdot 43 + 5 \cdot 62 + 6 \cdot 43 + 7 \cdot 18 + 8 \cdot 10 = 1009$$

$$\sum y_i = 26 \cdot 10 + 24 \cdot 33 + 22 \cdot 64 + 20 \cdot 47 + 18 \cdot 28 + 16 \cdot 15 + 14 \cdot 3 = 4186$$

$$\sum x_i^2 = 2^2 \cdot 9 + 3^2 \cdot 15 + 4^2 \cdot 43 + 5^2 \cdot 62 + 6^2 \cdot 43 + 7^2 \cdot 18 + 8^2 \cdot 10 = 5479$$

$$\sum y_i^2 = 26^2 \cdot 10 + 24^2 \cdot 33 + 22^2 \cdot 64 + 20^2 \cdot 47 + 18^2 \cdot 28 + 16^2 \cdot 15 + 14^2 \cdot 3 = 89044$$

$$\begin{aligned} \sum x_i \cdot y_i &= 2 \cdot 26 \cdot 2 + 3 \cdot 26 \cdot 3 + 4 \cdot 26 \cdot 3 + 5 \cdot 26 \cdot 2 + 2 \cdot 24 \cdot 4 + 3 \cdot 24 \cdot 5 + 4 \cdot 24 \cdot 13 + 5 \cdot 24 \cdot 7 + 6 \cdot 24 \cdot 4 \\ &+ 2 \cdot 22 \cdot 3 + 3 \cdot 22 \cdot 6 + 4 \cdot 22 \cdot 18 + 5 \cdot 22 \cdot 25 + 6 \cdot 22 \cdot 10 + 7 \cdot 22 \cdot 2 + 3 \cdot 20 \cdot 1 + 4 \cdot 20 \cdot 8 + 5 \cdot 20 \cdot 17 + \\ &+ 6 \cdot 20 \cdot 18 + 7 \cdot 20 \cdot 3 + 4 \cdot 18 \cdot 1 + 5 \cdot 18 \cdot 9 + 6 \cdot 18 \cdot 8 + 7 \cdot 18 \cdot 8 + 18 \cdot 2 + 5 \cdot 16 \cdot 2 + 6 \cdot 16 \cdot 3 + \\ &+ 7 \cdot 16 \cdot 4 + 8 \cdot 16 \cdot 6 + 7 \cdot 14 \cdot 1 + 8 \cdot 14 \cdot 2 = 20624 \end{aligned}$$

$$\text{Coeficientul de corelație este } r = \frac{20624 - \frac{1009 \cdot 4186}{200}}{\sqrt{\left(5479 - \frac{1009^2}{200}\right) \cdot \left(89044 - \frac{4186^2}{200}\right)}} = -0.6630$$

Coeficientul indică o corelație negativă (confirmând observația intuitivă asupra norului de puncte) puternică (deoarece $|r| > 0.5$).

Pragul de încredere este $PI = 0.6630 \cdot \sqrt{199} = 9,3527 > 3$, deci coeficientul de corelație este un indicator relevant.

3.2. Analiza regresiiilor

În general punctele din norul de puncte asociat seriilor de date nu se găsesc toate pe graficul unei funcții $y = f(x)$, ci sunt mai mult sau mai puțin împrăștiate.

Folosind metoda celor mai mici pătrate se poate determina totuși o funcție față de graficul căreia suma abaterilor valorilor individuale să fie minime. Aceasta este **funcția de regresie**.

Scopul construirii funcției de regresie este prognoza valorilor unei variabile folosind valorile celeilalte variabile.

Regresia liniară

Vom considera cazul când punctele corespunzătoare unei serii statistice sunt dispuse aproximativ după o dreaptă, adică variabilele sunt liniar corelate ($r \approx 1$ sau $r \approx -1$). În acest caz legătura cea mai simplă este

cea liniară în care unei creșteri a lui x (care este considerată variabila “predictor”) îi corespunde o creștere sau o scădere proporțională a lui y (care este considerată variabila “răspuns”).

Această relație se numește *regresia liniară* și este dată de ecuația

$$y = \alpha \cdot x + \beta.$$

numită *ecuația dreptei de regresie*.

Coeficienții dreptei de regresie se calculează folosind relațiile

$$\alpha = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i \cdot y_i - \frac{\left(\sum_{i=1}^n x_i\right) \cdot \left(\sum_{i=1}^n y_i\right)}{n}}{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}}$$

$$\beta = \bar{y} - \alpha \cdot \bar{x} = \frac{1}{n} \cdot \left(\sum_{i=1}^n y_i - \alpha \cdot \sum_{i=1}^n x_i \right).$$

Regresia liniară se poate folosi dacă sunt îndeplinite următoarele ipoteze:

- valorile variabilei dependente Y trebuie să aibă o repartiție normală
- Y și X trebuie să aibă dispersia (sau abaterea standard) asemănătoare
- Legătura dintre variabile trebuie să fie liniară (verificare empirică, pe baza norului de puncte care trebuie să aibă o formă alungită)

Din ecuație de regresie se pot determina valorile lui Y dacă se știu valorile lui X .

Estimatorul dispersiei lui Y în jurul dreptei de regresie este

$$s^2 = \frac{1}{n-2} \cdot \left[\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n} - \frac{\sum_{i=1}^n x_i \cdot y_i - \frac{\left(\sum_{i=1}^n x_i\right) \cdot \left(\sum_{i=1}^n y_i\right)}{n}}{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}} \right]$$

Exemplul 6 Producția de struguri obținută într-o fermă în mai mulți ani și numărul de zile însorite observate de-a lungul anilor sunt înregistrate în tabelul următor. Pe baza datelor din table să se precizeze dacă cele două serii de date sunt corelate.

Producția de struguri/ha x_i	Numărul de zile însorite y_i	x_i^2	y_i^2	$x_i \cdot y_i$
1.2	101	1.44	10201	1.2120
0.8	92	0.64	8464	0.7360
1	110	1.0	12100	1.1000
1.3	120	1.69	14400	1.5600
0.7	90	0.49	8100	0.6300
0.8	82	0.64	6724	0.6560
1.0	93	1.0	8649	0.9300
0.6	75	0.36	5625	0.4500
0.9	91	0.81	8281	0.8190
1.1	105	1.21	11025	1.1550

$\sum x_i = 9.4$	$\sum y_i = 959$	$\sum x_i^2 = 9.28$	$\sum y_i^2 = 93569$	$\sum x_i \cdot y_i = 924.80$
------------------	------------------	---------------------	----------------------	-------------------------------

Norul de puncte corespunzător seriilor de date este prezentat în figura 3.4.

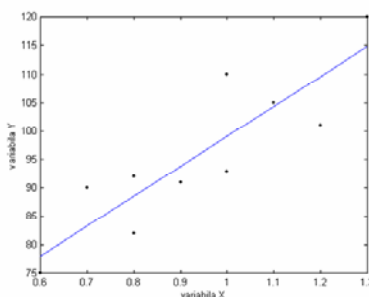


Figura 3.7 Norul de puncte și dreapta de regresie a serieilor de date din Exemplul 6

Coefficientul de corelație este

$$r = \frac{924.80 - \frac{9.4 \cdot 959}{10}}{\sqrt{\left(9.28 - \frac{9.4^2}{10}\right) \cdot \left(93569 - \frac{959^2}{10}\right)}} = \frac{23.34}{\sqrt{0.444 \cdot 1600.9}} = 0.8754.$$

Coefficienții dreptei de regresie sunt dați de

$$\alpha = \frac{924.80 - \frac{9.4 \cdot 959}{10}}{9.28 - \frac{9.4^2}{10}} = 118.39 \text{ și } \beta = 959 - 118.39 \cdot 9.4 = -153.86.$$

Dreapta de regresie (desenată în figura 4) are ecuația

$$Y = 118.39 \cdot X - 153.86.$$

Indicatorul dispersiei lui Y în jurul dreptei de regresie este

$$s^2 = 1548.3224.$$

Interpretarea rezultatului:

- coeficientul de corelație este pozitiv, deci o tendință ascendentă a variabilei “x” antrenează o tendință ascendentă a variabilei “y”
- coeficientul de corelație este apropiat de 1, deci corelația este puternică.
- Pragul de încredere este $PI = r \cdot \sqrt{10-1} = 2.6262 < 3$, deci numărul de date nu este suficient de mare pentru a asigura faptul ca e semnificativ coeficientul de corelație. Aceasta observație este confirmat de faptul că este mare coeficientul de dispersie al lui Y în jurul dreptei de regresie.
- Dreapta de regresie se va folosi cu precauție pentru prognoze, deoarece nu reprezintă o estimare precisă a dependenței dintre seriile de date.

Exemplul 7 Cantitatea de nutreț folosit și numărul de animale crescute în 14 ferme sunt prezentate în tabelul următor. Pr baza datelor din table să se precizeze dacă există corelații între cele două aspecte ale activității fermei.

	Cantitatea de nutreț x_i	Numărul de animale y_i	x_i^2	y_i^2	$x_i \cdot y_i$
1	6.4	380	42.25	144400	2470.0
2	5.2	200	27.04	40000	1040.0
3	0.4	15	0.16	225	6.0
4	1.7	50	2.89	2500	85.0
5	1.9	40	3.61	1600	76.0
6	2.4	40	5.76	1600	96.0
7	3.2	41	10.24	1681	131.2
8	4.7	18	22.09	324	84.6

9	10.1	210	102.01	44100	2121.0
10	12.5	190	156.25	36100	2375.0
11	13.1	200	171.61	40000	2620.0
12	5.5	55	30.25	3025	302.5
13	2.5	38	6.25	1444	95.0
14	1.5	20	2.25	400	30.0
	$\sum x_i = 71.2$	$\sum y_i = 1497$	$\sum x_i^2 = 582.66$	$\sum y_i^2 = 317399$	$\sum x_i \cdot y_i = 11532.3$

In Figura 3.5 este prezentat norul de puncte asociat seriilor de date și dreapta de regresie.

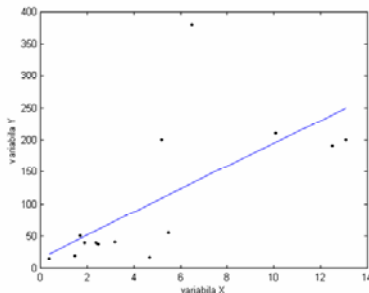


Figura 3. 8 Norul de puncte și dreapta de regresie asociate seriilor de date din Exemplul 7
Coeficientul de corelație este

$$r = \frac{11532.3 - \frac{71.2 \cdot 1497}{14}}{\sqrt{\left(585.66 - \frac{71.2^2}{14}\right) \cdot \left(317339 - \frac{1497^2}{14}\right)}} = 0.6653$$

Dreapta de regresie (desenată în figura 5) are ecuația

$$Y = 17.7685 \cdot X + 16.5626.$$

Interpretarea rezultatelor

- coeficientul de corelație este pozitiv, deci o tendință ascendentă a variabilei “X” antrenează o tendință ascendentă a variabilei “Y”
- coeficientul de corelație nu este apropiat de 0, deci cele două variabile ar putea fi corelate.
- pragul de încredere este $PI = r \cdot \sqrt{14-1} = 2.3984 < 3$, deci numărul de date nu este suficient de mare pentru a asigura faptul ca e semnificativ coeficientul de corelație.
- Dreapta de regresie se va folosi cu precauție pentru prognoze, deoarece nu reprezintă o estimare precisă a dependenței dintre seriile de date.

Exemplul 8 L_E reprezintă limita de elasticitate a tulpinei unei plante iar L_R reprezintă limita sa de ruptură. Stiind că raportul $X = \frac{L_E}{L_R}$ este strâns legat de conținutul în fibre al tulpinii, notat Y , să se analizeze corelația obținută între cei doi parametri pe un eșantion de 79 plante, date prezentate în tabelul de mai jos. Numerele întregi din interiorul tabelului reprezintă frecvența de apariție în cele 79 probe a perechilor (X, Y) corespunzătoare.

X\Y	0.5	0.6	0.7	0.8	0.9	Distribuție marginală pentru X
0.5	0	2	0	0	8	10
0.6	0	4	2	9	0	15
0.7	2	12	3	1	0	18
0.8	21	14	0	0	0	35
0.9	1	0	0	0	0	1
Distribuție marginală pentru Y	24	32	5	10	8	79

Variabilele statistice pentru care se studiază corelația sunt

$X = \begin{pmatrix} 0.5 & 0.6 & 0.7 & 0.8 & 0.9 \\ 10 & 15 & 18 & 35 & 1 \end{pmatrix}$ și $Y = \begin{pmatrix} 0.5 & 0.6 & 0.7 & 0.8 & 0.9 \\ 24 & 32 & 5 & 10 & 8 \end{pmatrix}$, pentru care tabelul interdependențelor este prezentat anterior.

Pentru calculul coeficientului de corelație sunt necesare următoarele rezultate:

$$\sum x_i = 10 \cdot 0.5 + 15 \cdot 0.6 + 18 \cdot 0.7 + 35 \cdot 0.8 + 1 \cdot 0.9 = 49.9$$

$$\sum y_i = 24 \cdot 0.5 + 32 \cdot 0.6 + 5 \cdot 0.7 + 10 \cdot 0.8 + 8 \cdot 0.9 = 55.5$$

$$\sum y_i^2 = 24 \cdot 0.5^2 + 32 \cdot 0.6^2 + 5 \cdot 0.7^2 + 10 \cdot 0.8^2 + 8 \cdot 0.9^2 = 39.93$$

$$\sum x_i \cdot y_i = 0.5 \cdot 0.6 \cdot 2 + 0.5 \cdot 0.9 \cdot 8 + 0.6 \cdot 0.6 \cdot 4 + 0.6 \cdot 0.7 \cdot 2 + 0.6 \cdot 0.8 \cdot 9 + 0.7 \cdot 0.5 \cdot 2 + 0.7 \cdot 0.6 \cdot 12 + 0.7 \cdot 0.7 \cdot 3 + 0.7 \cdot 0.8 \cdot 1 + 0.8 \cdot 0.5 \cdot 21 + 0.8 \cdot 0.6 \cdot 14 + 0.9 \cdot 0.5 \cdot 1 = 34.1400$$

Coeficientul de corelație este $r = -0.8194$

Interpretarea rezultatelor:

- coeficientul de corelație nu este apropiat de 0, deci variabilele sunt corelate.
- coeficientul de corelație este negativ, deci valorilor mari ale lui X le corespund multe valori mici ale lui Y (se confirmă prin poziționarea datelor în table)
- pragul de încredere este $PI = 0.8194 \cdot \sqrt{77} = 7.1902 > 3$, deci indicatorul de corelație este semnificativ și poate fi folosit în studiul corelației dintre variabilele X și Y
- în acest caz dreapta de regresie dă informații semnificative asupra valorilor lui Y, dacă se cunosc valorile lui X.

Dreapta de corelație are ecuația

$$Y = -0.6885 \cdot X + 17.9713.$$

Folosind această ecuație de regresie putem determina (aproximativ) valorile lui Y.

De exemplu, dacă $X = 0.75$ rezultă că $Y = -0.6885 \cdot 0.75 + 17.9713 = 17.4549$.

Exerciții propuse

1. Pentru a stabili în ce măsură depinde producția de tulpini de perioada de vegetație a diferitelor soiuri de cânepă de fibre, s-au realizat observații asupra cinci soiuri de cânepă foarte diferite ca perioadă de vegetație. Datele sunt prezentate în tabelul următor (pe orizontală este prezentată perioada de vegetație –în zile- și pe verticală este prezentată producția de tulpini – în q/ha) pentru cinci ani de producție.

	55 zile	70 zile	85 zile	100 zile	115 zile	130 zile
An 1	12	18	25	39	48	64
An 2	10	20	27	36	46	57
An 3	14	24	30	34	44	66
An 4	15	22	29	40	54	59
An 5	13	19	26	37	52	62

Să se precizeze dacă cele două caracteristici ale producției (perioada de vegetație și producția obținută) sunt corelate. $X = 90$ zile

Pe baza ecuației de regresie să se precizeze valoarea aproximativă a producției Y (per ha) dacă perioada de vegetație ar fi $X = 90$ zile.

Pentru stabilirea acțiunii azotului asupra conținutului de fibre din tulpinile plantelor de cânepă au fost efectuate măsurători în patru ani consecutiv la patru ferme asupra plantelor produse. Rezultatele sunt prezentate în tabelul următor. Pe orizontală este prezentă cantitatea de sulfat de amoniu folosită ca îngrășământ în cele cinci ferme (în kg/ha) iar pe verticală este prezentat conținutul de fibre din tulpini (în procente). Pe baza datelor prezentate în tabel să se precizeze dacă între cantitatea de îngrășământ folosită și conținutul de fibre ale tulpinilor există corelație.

	Ferma 1: 0	Ferma 2: 150	Ferma 3: 300	Ferma 4: 450
An 1	19.0	21.8	22.1	21.8
An 2	18.1	22.5	23.0	22.7
An 3	18.9	20.6	22.6	22.4
An 4	19.8	22.0	23.1	20.8

4. Elemente de teoria probabilităților aplicate în biologie și agricultură

4.1. Elemente de analiză combinatorică

Analiza combinatorică se ocupă cu numărarea anumitor grupări ce se pot realiza cu elementele unei mulțimi finite.

Prin cardinalul unei mulțimi finite $A = \{a_1, a_2, \dots, a_n\}$ se înțelege numărul “ n ” al elementelor sale. Se notează $\text{card}(A) = n$.

O grupare care permută elementele mulțimii A este formată din toate elementele mulțimii. Două permutări diferă prin ordinea în care sunt scrise elementele. Din punct de vedere matematic, o **permutare** a mulțimii A este o bijecție de la A la A .

Numărul permutărilor lui A este

$$P_n = 1 \cdot 2 \cdot 3 \cdot \dots \cdot n \stackrel{\text{notatie}}{=} n!$$

O submulțime ordonată de k elemente ale lui A se numește **aranjament de ordin k** . Numărul de aranjamente de ordin k ale unei mulțimi cu n elemente este

$$A_n^k = n \cdot (n-1) \cdot \dots \cdot (n-k+1) = \frac{n!}{(n-k)!}.$$

El reprezintă numărul aplicațiilor injective ale mulțimii $\{1, 2, \dots, k\}$ în A .

Submulțimile de câte k elemente ale lui A care nu sunt ordonate se numesc **combinări de ordin k** .

Numărul acestor combinări este

$$C_n^k = \frac{A_n^k}{P_k} = \frac{n!}{k!(n-k)!}.$$

Principalele proprietăți ale combinărilor sunt:

1. $\sum_{k=0}^n C_n^k = 2^n$
2. $C_n^k = C_n^{n-k}$
3. $C_n^k = C_{n-1}^k + C_{n-1}^{k-1}$ pentru orice $1 \leq k \leq n-1$ (formula lui Pascal)
4. $(a+b)^n = \sum_{k=0}^n C_n^k \cdot a^k \cdot b^{n-k}$ (binomul lui Newton)

Observație C_n^k se mai numește și **coeficient binomial**, datorită formulei de dezvoltare a binomului lui Newton.

Dacă n_1, n_2, \dots, n_k sunt numere naturale și $n_1 + n_2 + \dots + n_k = n$ se definește **coeficientul multinomial** prin

$$C_n^{n_1, n_2, \dots, n_k} = \frac{n!}{n_1! \cdot n_2! \cdot \dots \cdot n_k!}.$$

El are următoarea interpretare: Dacă mulțimea A conține n elemente, atunci există $C_n^{n_1, n_2, \dots, n_k}$ partiții ordonate diferite $\{A_1, A_2, \dots, A_k\}$ ale lui A astfel încât fiecare A_i să conțină n_i elemente, $i = 1, 2, \dots, k$.

$C_n^{n_1, n_2, \dots, n_k}$ se numește coeficient multinomial pentru că are loc relația

$$(x_1 + x_2 + \dots + x_k)^n = \sum_{n_1 + \dots + n_k = n} C_n^{n_1, n_2, \dots, n_k} \cdot x_1^{n_1} \cdot x_2^{n_2} \cdot \dots \cdot x_k^{n_k}.$$

În încheiere amintim principiul (regula) produsului:

Dacă o operațiune O_1 poate fi efectuată în n_1 moduri diferite, operațiunea O_2 poate fi executată în n_2 moduri diferite, etc., operațiunea O_k poate fi executată în n_k moduri diferite, atunci cele k operațiuni pot fi executate una după alta în $n_1 \cdot n_2 \cdot \dots \cdot n_k$ moduri diferite.

4.2. Introducere euristică în teoria probabilităților

In cele ce urmează prezentăm într-o formă simplă noțiuni de bază ale teoriei probabilităților, pornind de la definiția euristică a probabilității de realizare a unui eveniment.

Datele cu care operează teoria probabilităților sunt obținute prin observații asupra evenimentelor necontrolate din natură, societate, fie ca rezultat al experimentelor controlate.

Noțiunile primare în teoria probabilităților sunt cele de eveniment într-un experiment aleator și de probabilitate a evenimentului.

Definim un experiment ca fiind procesul prin care efectuăm o observație sau o măsurătoare.

Experiențele care pot avea rezultate diferite în funcție de o serie de circumstanțe întâmplătoare și rezultatele nu pot fi cunoscute înaintea realizării experimentului se numesc experiențe aleatoare.

Rezultatul unui experiment aleator se numește **realizare**. Colecția tuturor realizărilor acoperă orice posibilitate (adică este exhaustivă) și nici o realizare nu se suprapune peste alta (realizările sunt exclusive).

O colecție de realizări se numește **eveniment**, iar mulțimea tuturor realizărilor formează evenimentul sigur. **Evenimentul sigur** se produce cu certitudine la orice efectuare a experimentului.

Evenimentul care nu se produce ori de câte ori repetăm experiența se numește **eveniment imposibil**.

Evenimentul sigur va fi notat cu X , evenimentul imposibil cu \emptyset , iar evenimentele particulare cu A, B, C, \dots

Evenimentele compuse se obțin folosind operații cu evenimentele simple:

- evenimentul $A \cup B$ se realizează dacă se realizează A sau se realizează B .
- evenimentul $A \cap B$ se realizează dacă se realizează și A și B .
- evenimentul $A - B$ se realizează dacă se realizează A și nu se realizează B .

Unui eveniment A în corespunde evenimentul contrar, notat $C_X(A)$, a cărui producere înseamnă nerealizarea lui A .

Analogia între evenimentele compuse și teoria mulțimilor este evidentă, un eveniment fiind asociat unei submulțimi a lui X .

Probabilitatea unui eveniment A , notată $P(A) \in [0,1]$ reprezintă șansa pe care o are evenimentul de a se produce.

Dacă experimentul aleator are un număr finit de realizări și acestea sunt egal probabile (adică nu există un motiv ca o realizare să se producă mai frecvent decât alta) atunci se definește probabilitatea unui eveniment ca raportul dintre numărul cazurilor favorabile și numărul cazurilor posibile, adică

$$P(A) = \frac{\text{numar cazuri favorabile lui } A}{\text{numar cazuri posibile}}.$$

Observație: Dacă experimental aleator are un număr finit de realizări ce nu sunt egal probabile, nu există o modalitate teoretică ce permite calculul probabilității cu acuratețe absolută.

Exemplu: Experimentul aleator clasic este aruncarea unui zar cubic, realizat din material omogen. Realizările posibile ale experimentului sunt apariția feței cu numărul 1, 2, 3, 4, 5, 6. Evenimentele de apariție al feței cu nr "k" se numesc evenimente elementare. Evenimentul sigur este apariția unei fețe și este asociat mulțimii $X = \{1,2,3,4,5,6\}$. Alte evenimente sunt reprezentate simbolic prin mulțimi. Spre exemplu apariția unei fețe pare este reprezentată de mulțimea $A = \{2,4,6\}$.

Probabilitatea de realizarea a lui A este $P(A) = \frac{3}{6}$.

Dacă zarul nu e cubic sau nu este bine centrat, atunci probabilitatea de apariție a unei fețe nu este $1/6$.

În unele situații realizarea unui eveniment este condiționată de realizarea prealabilă a altui eveniment. Ideea care conduce la definiția **probabilității condiționate** este următoarea: știm că evenimentul B s-a produs, deci cazurile posibile pentru $A \cap B$ sunt cazurile favorabile pentru B , adică

$$P_B(A) = \frac{\text{nr cazuri favorabile si pentru } A \text{ si pentru } B}{\text{nr cazuri favorabile pentru } B} = \frac{\frac{\text{nr cazuri favorabile si pentru } A \text{ si pentru } B}{\text{nr cazuri posibile}}}{\frac{\text{nr cazuri favorabile pentru } B}{\text{nr cazuri posibile}}}$$

Probabilitatea unui eveniment A , condiționată de evenimentul B , cu $P(B) \neq 0$, se definește prin

$$P_B(A) = \frac{P(A \cap B)}{P(B)}$$

Două evenimente se numesc independente dacă $P(A \cap B) = P(A) \cdot P(B)$.

Dacă două evenimente sunt independente atunci realizarea unuia nu influențează realizarea celuilalt eveniment, adică $P_B(A) = P(A)$

Pornind de la definiția probabilității, se pot demonstra următoarele proprietăți:

Propoziția 1.

1. $0 \leq P(A) \leq 1$, $P(X) = 1$ și $P(\emptyset) = 0$
2. $P(C_X(A)) = 1 - P(A)$
3. Dacă $A \subset B$ atunci $P(A) \leq P(B)$
4. Dacă A și B sunt două evenimente și $A \cap B = \emptyset$, atunci $P(A \cup B) = P(A) + P(B)$
5. Dacă A și B sunt două evenimente, atunci $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
6. Dacă $A_1 \cup A_2 \cup \dots \cup A_n = X$ și $A_i \cap A_j = \emptyset$ pentru $i \neq j$ atunci

$$P(B) = \sum_{i=1}^n P(A_i) \cdot P_{A_i}(B) \quad (\text{formula probabilității totale})$$

$$P_B(A_k) = \frac{P(A_k) \cdot P_{A_k}(B)}{\sum_{i=1}^n P(A_i) \cdot P_{A_i}(B)} \quad (\text{formula lui Bayes})$$

4.3 Aplicații în biologie

4.3.1 Într-un organism există genotipurile AA, Aa, aa. Părinții transmit către urmași fiecare câte o singură genă. Se presupune că populația parentală este suficient de mare încât încrucișarea să se facă la întâmplare și că proporțiile genotipurilor sunt respectiv α , 2β , respectiv γ , cu $\alpha > 0$, $\beta > 0$, $\gamma > 0$ și $\alpha + 2\beta + \gamma = 1$. De asemenea se presupune că probabilitatea ca un părinte să transmită o genă este $1/2$.

Să se precizeze proporțiile genotipurilor după prima generație și după a doua generație. Să se interpreteze rezultatele.

În prima generație pot să apară tipurile AA, Aa, aa. Pentru fiecare tip tabelul de calcul al probabilităților este prezentat mai jos.

a) pentru tipul AA

Pentru transmiterea genotipului AA este obligatoriu ca cel puțin o genă A să apară în genotipul fiecărui părinte.

Tipul mascul	Tipul femel	Probabilitatea formării cuplului	Probabilitatea transmiterii genotipului AA	Probabilitatea existenței genotipului AA la urmași
AA	AA	$\alpha \cdot \alpha$	$1 \cdot 1 = 1$	α^2
AA	Aa	$\alpha \cdot 2\beta$	$1 \cdot \frac{1}{2} = \frac{1}{2}$	$2\alpha\beta$
aA	AA	$\alpha \cdot 2\beta$	$1 \cdot \frac{1}{2} = \frac{1}{2}$	$2\alpha\beta$
aa	Aa	$2\beta \cdot 2\beta$	$\frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$	β^2

Deci probabilitatea existenței unor urmași de tipul AA la prima generație este

$$P_1(AA) = \alpha^2 + \alpha\beta + \alpha\beta + \beta^2 = (\alpha + \beta)^2 \quad (1)$$

b) Raționând la fel se obține probabilitatea existenței unor urmași de tip aa la prima generație este

$$P_1(aa) = (\beta + \gamma)^2 \quad (2)$$

c) pentru tipul generic Aa (sau aA) tabelul probabilităților este

Tipul mascul	Tipul femel	Probabilitatea formării cuplului	Probabilitatea transmiterii genotipului Aa	Probabilitatea existenței genotipului Aa la urmași
AA	Aa	$\alpha \cdot 2\beta$	$1 \cdot \frac{1}{2} = \frac{1}{2}$	$\alpha\beta$
Aa	AA	$2\beta \cdot \alpha$	$\frac{1}{2} \cdot 1 = \frac{1}{2}$	$\alpha\beta$
AA	aa	$\alpha \cdot \gamma$	$1 \cdot 1 = 1$	$\alpha\gamma$
aa	AA	$\gamma \cdot \alpha$	$1 \cdot 1 = 1$	$\alpha\gamma$
Aa	aa	$2\beta \cdot \gamma$	$\frac{1}{2} \cdot 1 = \frac{1}{2}$	$\beta\gamma$
aa	Aa	$\gamma \cdot 2\beta$	$1 \cdot \frac{1}{2} = \frac{1}{2}$	$\beta\gamma$
Aa	Aa	$2\beta \cdot 2\beta$	$2 \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{2}$	$2\beta^2$

Probabilitatea existenței unor urmași de tip Aa la prima generație este

$$P_1(Aa) = 2\alpha\beta + 2\alpha\gamma + 2\beta\gamma + 2\beta^2 = 2(\alpha + \beta) \cdot (\beta + \gamma) \quad (3)$$

Faptul că rezultatele sunt corecte este reflectat și de relația

$$P_1(AA) + P_1(Aa) + P_1(aa) = (\alpha + \beta)^2 + 2(\alpha + \beta) \cdot (\beta + \gamma) + (\beta + \gamma)^2 = (\alpha + 2\beta + \gamma)^2 = 1$$

La a doua generație probabilitățile vor fi:

$$P_2(AA) = [(\alpha + \beta)^2 + (\alpha + \beta)(\beta + \gamma)]^2 = (\alpha + \beta)^2 (\alpha + 2\beta + \gamma)^2 = (\alpha + \beta)^2 = P_1(AA)$$

$$P_2(aa) = [(\beta + \gamma)^2 + (\alpha + \beta)(\beta + \gamma)]^2 = (\beta + \gamma)^2 (\alpha + 2\beta + \gamma)^2 = (\beta + \gamma)^2 = P_1(aa)$$

$$P_2(Aa) = 2(\alpha + \beta) \cdot (\beta + \gamma) = P_1(Aa)$$

Prin urmare, de la prima generație încolo probabilitățile de menținere a genotipurilor sunt aceleași. Se spune că procesul evolutiv este stohastic stabil.

Bibliografie

1. Bălan V., Matematici Superioare Aplicate, Editura Universitaria, Craiova, 2007
2. Petrișor E., Probabilități și statistică, Editura Politehnica, Timișoara 2005
3. Cristea M., Genetica ecologică și evoluția, Editura Ceres, București, 1991
4. Ștefănescu D.T., Călin G., Genetica și cancerul : (Elemente de genetică și patologie moleculară), Editura Didactică și Pedagogică București, 1996
5. Raicu P. (coordonator), Biologie : Genetică și evoluționism : Manual pentru clasa a XII-a, Editura Didactică și Pedagogică București, 1998
6. Biji E. M. (coordonator), Statistica managerială a agentului economic din agricultură, Editura Ceres, București, 1998
7. Howitt D. Cramer D. Introducere în SPSS, Editura Polirom, 2006