Probabilități și Statistică Probabilități și Statistică

# Probabilități și Statistică - Curs 8

Olariu E. Florentin
Probabilități și Statistică

### Table of contents

1 Statistică	
Probabilități și Statistică Probab <b>Introducere</b>	
Vocabular	
2 Statistică descripti	vă Probabilități și Statistică
Variabile Variabile	
ProbabReprezentări gra	ficebabilități și Statistică
Măsuri ale tendi	nței centrale
Proba <b>Media</b> i Statistică	
Probabilita Mediana că	
Proba <b>Mòdul</b> i Statistică	
Compararea dife	ritelor măsuri Probabilități și Statistică Probabilități și Statistică
	bilit <b>ății</b> babilități și Statistică Probabilități și Statistică

Probabilită Valori aberante

3 Bibliografie

# Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

- Rădăcinile cuvântului Statistică sunt latine: Status (latina veche) care înseamnă means stat (politic), Statista (italiană) înseamnă politician.
- La mijlocul secolului al XVII-lea într-o Universitate Germană a fost folosit pentru prima oară cuvântul statistik cu sensul de ştiinţă politică a statelor: analiza datelor privind statele.
- În Marea Britanie la sfârșitul secolului XVIII termenul de statistică a fost introdus cu un înțeles similar: știința statelor (sau aritmetica politică).
- Utilizarea statisticii fără a o numi în mod expres datează de la începutul civilizației umane: forme incipiente de recensământ al populației, sistematizarea datelor geografice şi economice etc.

- Primul studiu statistic este considerat în general a fi cel care a pus bazele demografiei: în 1662 doi englezi au introdus metode statistice cum ar fi tabelele speranței de viață şi probabilitățile de supraviețuire la diferite vârste.
  - Abia în secolul XIX rezultatele din teoria probabilităților au început a fi folosite în raționamentul statistic.
- Bazele matematice ale statisticii s-au consolidat datorită rezul
   tatelor profunde obținute în teoria probabilităților din sec
   olul anterior.

   Probabilităților din sec-
  - Începând cu secolul XX au fost dezvoltate noi metode și teorii, iar o influențămajoră asupra statisticii a avut-o dezvoltarea informaticii.

"Statistics has become the universal language of the sciences."

Elementary Statistics, R. Johnson, P. Kuby

O problemă tipică de statistică este format din

- unul sau mai multe experimente aleatoare din efectuarea
- o metodă de extragere a informației din date si de interpretare a rezultatelor.

Modul în care informați aeste procesată și interpretată dă naștere la doua ramuri ale statisticii ca știință

- Statistica descriptivă colectează, prezintă și descrie datele (de multe ori în formă grafică).
- Statistica inferențială folosind datele deja colectate ia de-Proba cizii relativ la populația în cauză.

#### Introducere

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

### Definition 1 atistică

Statistica este știința colectării, descrierii, interpretării datelor și luării de decizii pe baza acestor date.

Cele două ramuri ale statisticii sunt și cei doi pași dintr-un studiu statistici:

Probabilități și Statistică
Probabilități și Statistică

- statistica descriptivă are rolul de a
  - sintetiza, aduna și reprezenta datele;
- aranja informaţia, pregătind-o pentru luarea deciziilor;
  - statistica inferențială are drept scop
    - o luarea deciziilor pe baza datelor strânse; pabilităti și Statistică
- estimarea parametrilor (cum sunt media, dispersia etc);
  - Proba o verificarea ipotezelor statistice.

- Statistica își are propriul limbaj, dincolo de împărțirea în descriptivă și inferențială.
- Cel mai important concept în statistică este acela de *populație*: colecția completă (exhaustivă) a obiectelor care prezintă interes pentru cel care face studiul.
- Exemple de populații: mulţimea studenţilor din Iaşi, mulţimea românilor analfabeţi, mulţimea dozelor de cola produse întroba o luna întro fabrică, mulţimea furtunilor tropicale din 2015.

### Definition 2 distica

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

• O populație poate fi finită (dacă poate fi teoretic listată) sau infinită (populația cutremurelor de pământ din zona Vrancea).

• Din cauza dimensiunilor mari ale unei populații studiul statistic se concentrează asupra unei porțiuni mai mici a populației. Acesta este un *eşantion* care constă din indivizi selectați din populație.

# Definition 3

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

- Când se alege o populație sau un eșantion pentru studiu, interesează o anumită trăsătură a indivizilor.
  - Astfel de trăsături (atribute) pot fi: înălţimea, volumul, magnitudinea pe scara Richter, vârsta, presiunea sângelui, culoarea ochilor, suprafata etc.

# Definition 4

O variabilă sau un atribut este o caracteristică a indivizilor din populație sau eșantion.

• După ce alegem un eșantion trebuie să măsurăm valorile unuia sau mai multor variabile asociate. Acestea sunt *datele*, ele pot fi numere reale, întregi, cuvinte, litere etc.

# Definition 5

Datele sunt valorile variabilei colectate de la fiecare individ din eşantion.

• O populație este descrisă numeric de *parametri* (medie, disperse, deviație standard); parametrii sunt în centrul unui probal studiu statistic.

# Terminologia statistică

## Definition 6

Un parametru este o valoare numerică care privește întreaga populație.

Probabilități și Statistică

Probabilități și Statistică

- Dacă populația este foarte mare (ceea ce se întâmplă adesea) un parametru anume nu poate fi calculat.
- O soluție este de a calcula parametrul doar pentru un eșantion al populației. Aceasta este o *statistică*.
- Pentru orice parametru și fiecare eșantion există o statistică

  Probabilitări și Statistică

  Probabilitări și Statistică

  Probabilitări și Statistică

# **Definition** 7 atistica

O statistică este un parametru calculat pentru un eșantion în locul întregii populații.

### Folosirea terminologiei

- pbabilități și Statistică Probabilități și Statistică
  Probabilități și Statistică Probabilități și Statistică
  pbabilități și Statistică Probabilități și Statistică
- Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică
- Populație: mulțimea studenților din primul an din Iași.
  - Eşantion: studenţii din primul an de la FII.
- Variabilă/atribut: dimensiunea vocabularului lor curent.
- Date: 4200, 3520, 1800, ... dimensiunile vocabularului pentru fiecare student din primul an de la FII.
  - Parametru: media dimensiunii vocabularului studenţilor din primul an din Iaşi
- Statistică: media dimensiunii vocabularului studenților din
  - babilități și Statistică Probabilități și Statistică Probabilități Probabilități și Statistică Probabilități și Statistică Probabilități Babilități și Statistică Probabilități și Statistică Probabilități

## Tipuri de variabile

obabilități și Statistică Probabilități și Statistică

- Clasificarea variabilelor împarte atributele în cantitative sau calitative. Astfel există
- Variabile care oferă o *informație calitativă*, cum ar fi culoarea ochilor studenților, genurile literare ale cărților dintro bibliotecă (ficțiune, știință, literatură motivațională etc), tipul de personalitate ale peroanelor dintr-o comunitate (sanguin, coleric, melancolic sau flegmatic), nivelul de satisfacție a clienților unui magazin etc.
  - Variabile care dau o *informație cantitativă*; spre exemplu: înălțimea studenților, greutatea lor, suma de bani pe care un student o cheltuie pe cărți într-un an școlar ș. a.

## Tipuri de variabile

Probabilități și Statistică Probabilități Probabilită Probabilități Probabilită Probabilită

### **Definition** 8 distication

O variabilă calitativă (sau categorică) este o variabilă care descrie un individ dintr-o populație (conform unor categorii). O variabilă cantitativă este o variabilă care măsoară un individ dintr-o populație.

- Probabilităti și Statistică Pro • Variabilele calitative pot fi *nominale* sau *ordinale*.
  - Variabilele nominale sunt: culoarea ochilor, tipul de personalitate, numele membrilor unei comunități etc.
  - Exampe de variabile ordinale: nivelul de satisfacție a clienților, nivelul de educație (liceal, post liceal, universitar, post universitar, doctoral) etc.

### Tipuri de variabile

Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

### Probabilități și Statistică **Definition**i **9**atistică

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

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

• Variabilele cantitative pot fi discrete sau continue. Cele două tipuri pot fi distinse astfel: unele numără iar celelalte măsoară.

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică babilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică babilități și Statistică Probabilități și Statistică Probabilități și Statistică

- O variabilă discretă de obicei numără: numărul de credite ale unui student, numărul de pagini ale unei cărţi etc; câteodată o asemenea variabilă sumează puncte/note care nu pot fi continue.
- O variabilă continuă măsoară: volumul, înălţimea, viteza, presiunea etc.

# Definition 10

O variabilă discretă este o variabilă care are un număr finit sau infinit dar numărabil de valori; o astfel de variabilă poate avea valori corespunzând unor puncte izolate de pe un interval real.

O variabilă contină este o variabilă care are un număr infinit și nenumărabil de valori; o astfel de variabilă poate avea, de obicei, orice valoare dintr-un interval real, incluzând orice valoare posibilă dintre orice două valori.

## Reprezentări grafice

- pabilități și Statistică Probabilități și Statistică
- O primă formă de explorare a datelor este utilizarea reprezentărilor grafice care pot revela un comportament sistematic (un şablon) al variabilei.
- Tipul de reprezentare grafică depinde în mod normal de tipul
  variabilei.
  - Pentru date calitative reprezentările grafice folosite sunt pie
     charts și bar graphs.
- Pentru datele cantitative scopul reprezentărilor grafice este

## Reprezentări grafice - date calitative

Probabilități și Statistică Probabilități și Statistică

- Datele calitative sunt mai întâi transformate în frecvențe.
- Frecvenţa unei observaţii (o valoare a unei variabile) este numărul de repetări ale acelei observaţii în eşantion.
- Frecvenţa relativă a unei observaţii este raportul dintre frecvenţa observaţiei respective şi numărul total de observaţii (dimensiunea eşantionului).
- Distribuţia frecvenţelor unei variabile calitative este familia tuturor perechilor formate din observaţie şi frecvenţa sa corepubation spunzătoare.

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistic obabilități și Statistică Probabilități și Statistică obabilități și Statistică Probabilități și Statistică

# Reprezentări grafice - date cantitative

Probabilități și Statistică

relative sau gruparea datelor pentru a regăsi distribuția frecvențelo

- datele sunt *grupate* în clase (sau *bins*) care sunt uzual intervale cu aceeași lungime; clasele nu trebuie sa se acopere.
- o regulă pentru determinarea lungimii claselor:  $1+\log n/\log 2$  unde n este dimensiunea eşantionului.
- apoi datele sunt *sortate* pe clase: se determină numărul observațiilor din fiecare clasă acestea sunt frecvențele.
  - suma frecvențelor este dimensiunea eșantionului (n).
  - frecvențele relative se pot afla împărțind frecvențele la n.

# Reprezentări grafice - date cantitative

Probabilități și Statistică Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

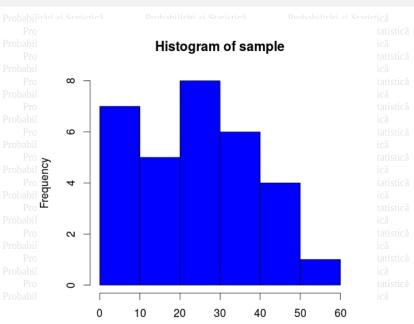
- Cea mai utilizată metodă de reprezentare grafică a datelor cantitative este *histogram*.
- O altă metodă la îndemână pentru eşantioanele relativ mici este stem-and-leaf.

Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

Probabilitāti şi Statisticā
Probabilitāti şi Statistică

# Reprezentări grafice - histograma



### Datele

- Probabilități și Statistică Probabilități și Statistică
  - Când privim reprezentarea grafică a datelor din eşantion ne putem pune următoarele întrebări.
- Care sunt valorile centrale/medii?
- Cât de mult sunt împrăștiate aceste date în jurul valorilor Probabilităti și Statistică Probabilități și Statistică
- Pro Care este froma distribuţiei?
  - Există valori care nu se potrivesc cu imaginea generală a distribuţiei?

### Tendinţa centrală

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică robabilități și Statistică Probabilități și Statistică robabilități și Statistică Probabilități și Statistică

- Tendinţa centrală sau centrul distribuţiei este centrul (abstract) al datelor. Toate măsurile tendinţei centrale sunt legate într-un fel sau altul de noţiunea de medie.
  - Diferite moduri de a defini tendința centrală:
- Probabilități și Statistică Probabilități Probabilități Probabilități Probabilită Probabilită
- Probabilită Numărul care minimizează suma tuturor deviațiilor absolute.
  - Numărul care minimizează suma tuturor deviațiilor la pătrat.
  - o Cea mai frecventă valoare.

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică
robabilități și Statistică
robabilități și Statistică
Probabilități și Statistică
robabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică
Probabilități și Statistică
Probabilități și Statistică

• Să presupunem că valorile din eșantion sunt  $x_1, x_2, \ldots, x_n$ .

### **Definition 11**

Media de selecție sau media eșantionului este media aritmetică a tuturor datelor din eșantion:

$$\overline{x}_n = rac{x_1 + x_2 + \cdots + x_n}{n}.$$

- Formula mediei pentru întreaga populație este în esență iden-Probabilităti și Statistică Probabilități și Statistică Probabilități și Statistică
- Media populaţiei se notează cu μ.

- În limbajul teoriei probabilităților media populației este media unei variabile aleatoare, X, ale cărei valori sunt sunt cele probable indivizilor din populație; deci $M[X] = \mu$ .
- Pro Media de selecție este o statistică care estimează media populației. Si Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică
- Să presupunem că  $X_1, X_2, \ldots, X_n$  sunt variabilele din spatele fiecărui individ al eșantionului, iar  $x_i$  este doar o valoare a variabilei  $X_i$ .
- Atunci  $X_i$  este o variabilă aleatoare cu aceeași distribuție ca a a lui X. Mai mult, variablele  $(X_i)_{1\leqslant i\leqslant n}$  sunt independente  $\hat{I}$  în ansamblu.
- Aceste observații conduc la faptul ca media de selecție poate fi văzută ca o variabilă aleatoare, iar media aritmetic calculată pentru un eșantion este una dintre posibilele valori ale ei (fiecare eșantion dă o altă valoare a mediei de selecție).

Dacă media de selecție este o variabilă aleatoare, îi putem
 calcula media:

Probabilități și Statistică 
$$M[\overline{x}_n] = M[\overline{x}_1 + X_2 + \cdots + X_n]$$
 pabilități și Statistică  $M[\overline{x}_n] = M[\overline{x}_1 + X_2 + \cdots + X_n]$  probabilități și Statistică Probabil

- Media mediei de selecție este media populației.
- O astfel de statistică se numește *estimator nedeplasat* al parametrului corespunzător.
- Media de selecție este un estimator nedeplasat pentru media

   populației.

  Probabilităt și Statistică

  Probabilităt și Statistică

  Probabilităt și Statistică

  Probabilităt și Statistică

  Probabilităt și Statistică

- Pro• Formula din definiția anterioară este valabilă pentru date negrupate. În acest caz toate datele din eşantion contribue direct la calulul mediei de selecție.
  - Pentru date grupate se folosește o formulă cu ponderi:

Probabilităti
$$M = \frac{m_i * f_i}{Probabilități şi States $_i$ Probabilități şi States $_i$$$

- unde  $m_i$  este mijlocul intervalului clasei i, iar  $f_i$  este numărul de observații care aparțin clasei i.
- În această formulă observațiile nu contribuie direct la calculul mediei; cu toate acestea este o formulă preferată în cazul datelor grupate pentru eșantioane mari fiind mai ușor de calculat.

- Ne întoarcem acum la definiția inițială (pentru date negrupate) a mediei de selecție.
- Variații mici în suma de la numărător nu modifică prea mult probabilită la variații mici ale probabilită datelor.
- Pro Valorile aberante sau extreme pot avea o influență mare asupra mediei; introducând o valoare foarte mare sau foarte prică media se poate schimba foarte mult.
- Media este o funcție *liniară* (la fel ca media unei variabile Probabilităt și Statistică Probabilităt și Statistică
- Deviațiile de la medie sunt  $(x_i \overline{x}_n)$ ; suma lor este zero:

$$\sum_{i}(x_i-\overline{x}_n)=0.$$

(Definiția variațională) Se poate arăta că media este numărul
 M care minimizează suma deviațiilor la pătrat:

$$\sum (x_i - M)^2$$
.

• Există şi alte tipuri de medie în afară de cea artimetică (A): media geometrică (G), media armonică (H).

abilități și Statistică Probabilități și Statistică Pro
$$n$$
bilități și Statistică Probabilități și Statistică Probabilități și Statistică abilități și Statistică Probabilități Probabilită Probabilită Probabilități Probabilită Probabilită Probabilită Probabilită Probabilită Probabilită Probabilită Probabilită

- Să presupunem că o maşină parcurge distanţa dintre două orașe de patru ori cu vitezele 80km/h, 90km/h, 60km/h, and 120km/h, respectiv. Care a fost viteza sa medie?
- Folosim media aritmetică obţinem 87.5 km/h; dar media adecvată aici este cea armonică: 82.3km/h.

• Mediana este o *statistică ordonată*; calculul unei astfel de statistici presupune ordonarea crescătoare a datelor din eşan-

### Definition 12

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

- Mediana împarte datele din eşantion în două jumătăţi: o jumătate conţine datele mai mari sau egale decât mediana, iar cealaltă jumătate le conţine pe cele mai mici sau egale.
- Valoarea medianei este o observaţie sau media a doua observaţii (pentru eşantioane de dimensiune pară).
- Ca statistică mediana este mult mai puţin influenţată de existenţa valorilor aberante.

Probabilități și Statistică Probabilități și Statist

### Definition 13

Mòdul este observația cea mai frecventă din eșantion.

- Pentru date grupate se alege mai întâi clasa cu cea mai mare frecvenţă clasa modală. Fie i indexul acestei clase, ai marginea stângă a intervalului corespunzător şi L lungimea comună a intervalelor.
- Atunci mòdul poate fi calculat folosind formula

Probabilități și Statistică 
$$mod = a_i + \frac{1}{f_i} \frac{L * (f_i - f_{i-1})}{(f_i - f_{i-1})} + \frac{1}{f_i - f_{i-1}}$$
 e și Statistică Probabilități Probabilită Probabilități Probabilită Probabilități Probabilită Pro

• Antimòdul este cea mai puţin frecventă observaţie.

### Compararea diferitelor măsuri

bbabilități și Statistică Probabilități și Statistică

- Mai stabile la valorile aberante sunt mediana şi mòdul.
- Media incorporează toate valorile şi nu poate fi calculată, în cazul datelor grupate pentru distribuţii deschise (primul, sau ultimul interval deschis).
- Pro Mediana și mòdul nu sunt funcții liniare. robabilități și Statistică
- Pro Mòdul este calculat mai laes pentru date grupate. Statistică
  - Pentru distribuții asimetrice mòdul oferă cea mai reală imagine asupra tendinței centrale.

## Compararea diferitelor măsuri

- Dacă eșantionul conține date foarte mari sau foarte mici mediana este măsura preferată mediei - stabilitatea o face mai reprezentativă.
- Pentru distribuţii simetrice cele trei măsuri sunt aproape

   Probabilităti si Statistica

   Probabilităti si Statistica

   Probabilităti si Statistica

   Probabilităti si Statistica
- Forma distribuţiei poate fi legată ade relaţia dintre medie şi
   mediană; forma poate fi
  - asimetrică spre stânga dacă  $\overline{x}_n < Me$ ;
- ullet simetrică dacă  $\overline{x}_n=Me;$ 
  - asimetrică spre dreapta dacă  $\overline{x}_n > Me$ ;

• Relativ la măsurile tendinței centrale există *măsuri de poz-*iţie care sunt statistici ordonate ca și mediana.

### **Definition 14**

Cvartilele sunt valori care împart domeniul (ordonat al) observațiilor în patru segmente egale.

- Prima cvartilă,  $Q_1$ , este o valoare astfel în cât 25% dintre observații sunt cel mult egale cu  $Q_1$  și cel mult 75% sunt mai mari sau egale.
- A treia cvartilă,  $Q_3$ , este o valoare astfel în cât 75% dintre observații sunt cel mult egale cu  $Q_3$  și cel mult 25% sunt mai mari sau egale.

- A doua cvartilă,  $Q_2$ , este o valoare astfel în cât 50% dintre observații sunt cel mult egale cu  $Q_2$  și cel mult 50% sunt mai mari sau egale. Din acest motiv a doua cvartilă este egală cu mediana:  $Me = Q_2$ .
  - Cvartilele au proprietăți similare cu cele ale medianei. Cea mai importantă fiind aceea că sunt stabile în prezenţa valorilor aberante.
  - Statistici ordonate similare sunt: decilele, percentilele etc.

    Toate aceste statistici împart datele ordonate în subeşantioane egale.
  - De exemplu există nouă decile care împart datele sortate în zece părți egale, fiecare parte reprezentând 10% din eşantion.

- După determinarea "centrului" datelor studiul statistic continuă cu analiza *împrăștierii* sau a *variabilității* datelor
- Valorile din eşantion pot să difere mult între ele şi faţă de valoarea "centrală".
  - Măsura în care valoare "centrală"/medie este reprezentativă pentru întreg eşantionul depinde de variabilitatea (sau dispersia) datelor.
- Eşantionul are variabilitate mare dacă există valori foarte mari sau foarte mici față de valoarea medie.
- Deoarece avem două moduri importante de a măsura tendința centrală (media și mediana) vom avea două metode de a măsura împrăștierea.

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statist

### Definition 15

Domeniul este diferența dintre cea mai mică și cea mai mare valoare din eșantion.

range = max - min.

- Deoarece definiția aceasta se bazează doar pe valorile extreme, dacă minimul sau maximul este foarte mare respectiv foarte mic, domeniul nu este reprezentativ pentru variabilitatea datelor.
  - Se observă că valorile aberante au o influență directă asupra domeniului. Probabilităti și Statistică Probabilități și Statistică

## Dispersia eşantionului

pabilități și Statistică Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

# Începem cu măsurile variabilității legate de medie.

- Deviațiile față de medie sunt  $(x_i \overline{x}_n)$ .
- O deviație  $(x_i \overline{x}_n)$  este pozitivă (negativă) când  $x_i$  este mai mare (mai mică) decât media de selecție.
  - Pentru a descrie o valoare medie a deviaţiilor s-ar putea utiliza media aritmetică a acestor deviaţii. Dar suma acestor deviaţii fiind zero, o astfel de medie este nulă.
- Putem îndepărta acest efect ridicând la pătrat deviațiile și utilizând o medie pătratică în locul uneia aritmetice.

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică obabilități și Statistică Probabilități și Statistică obabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică pabilități și Statistică
Probabilități și Statistică

Probabilități și Statistică
Probabilități și Statistică
Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

#### **Definition 16**

Dispersia eşantionului,  $s^2$ , n fiind dimensiunea eşantionului, este:

$$s^2=rac{\sum\limits_{i=1}^n(x_i-\overline{x}_n)^2}{n-1}$$
 ,

- Dispersia eșantionului este nenegativă și este zero dacăși numai dacă valorile sunt toate identice.
- Dispersia eşantionului este statistica asociată dispersiei popProbal ulației, notată cu  $\sigma^2$ . Dispersia estatistică probabilităti și Statistică probabilităti probabilitătică probabilităti probabilitătică probabilităti probabilităti probabilităti probabilită prob

- Morivul pentru care se utilizează (n-1) ca numitor în definiția dispersiei eșantionului este acela că astfel se obține un estimator nedeplasat.
  - Media dispersiei eşantionului (văzută ca o variabilă aleatoare)
     este

## Dispersia eşantionului

Pr $n^2(n-1)$ ristică

## Dispersia eşantionului

**n**robabilități și Statistică  $[X_i^{ ext{P}} X_i^{ ext{P}}]$ bilități și Statistică Probabilități și Statistic ati si Statistică i < jProbabilități și Statistică Probabi<u>lit</u>ăți și Statistică  $n^2(n-1)$ Probabilităti și Statistic i < j<del>obabilitătis</del>i Statistică Proba<del>bil</del>it<del>ăti și S</del>i  $n^2(n-1)$ P<sub>2</sub>obabilități și Statistică Probabilități și Statistica Probabilități și Statistică

O formulă mai simplă (exercițiu) pentru dispersia eșantionu
 Probabilități și Statistică
 Probabilități și Statistică
 Probabilități și Statistică

$$s^2 = rac{n}{n}\sum_{i=1}^n x_i - \left(\sum_{i=1}^n x_i
ight)^2}{n(n-1)}.$$

### **Definition 17**

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

- Deviația standard a eșantionului este un estimator deplasat al deviației standard a populației,  $\sigma$ .
- Se poate arăta că media deviației standard a eșantionului este mai mică decât cea a populației,  $M[s] < \sigma$ .

#### Sumarul celor cinci numere

Probabilități și Statistică Probabilități și Statistică

Continuăm cu măsuri ale împrăștierii legate de mediană. Mai întâi sumarul celor cinci numere.

### Definition 18

Sumarul celor cinci numere este compus din

- 1 min, cea mai mică valoare din eșantion;
- Q<sub>1</sub>, prima cvartilă;
- 3 Me, mediana;
- $\bigcirc Q_3$ , a treia cvartilă;
- 5 max, cea mai mare valoare din eşantion.

- obabilități și Statistică Probabilități și Stat
  - Probabilități și Statistica Probabilități și Statistică
    Probabilități și Statistică Probabilități și Statistică
    Probabilități și Statistică Probabilități și Statistică
- O metodă grafică de a reprezenta sumarul celor cinci numere:
- Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

### Definition 19

Cvartila medie este valoarea de mijloc dintre prima și cea de-a treia cvartilă:

$$midq=rac{Q_1+Q_3}{2}.$$

Domeniu intercvartilic este diferența dintre prima și cea de-a treia cvartilă:

$$IQR = Q_3 - Q_1.$$

#### Valori aberante

obabilități și Statistică Probabilități și Statistică

- Valorile aberante sunt acele valori din eşantion care pot fi considerate prea mici sau prea mari faţă de "tabloul" general al eşantionului.
- Evident, valorile aberante sunt legate de variabilitatea datelor. În mod obiņuit aceste valori vin din erori de măsură, dar pot avea şi cauze naturale.
  - Câteodată aceste valori aberante (dacă sunt datorate măsurilor) pot fi eliminate din eşantion înainte de orice altă analiză statistică.
  - Vom avea două reguli de determinare a valorilor aberante, deoarece şi variabilitatea datelor se măsoară în două feluri.

#### Valori aberante

Probabilități și Statistică Probabilități și Statistică

- Prima regulă este legată de medie. Pot fi considerate valori aberante acele valori ale eşantionului care nu aparțin intervalului  $(\overline{x}_n 2s, \overline{x}_n + 2s)$ .
- A doua regulă se numește regula 1.5 \* IQR și spune că o valoare este aberantă dacă nu aparține intervalului ( $Q_1 1.5 * IQR$ ,  $Q_3 + 1.5 * IQR$ ).

Probabilități și Statistică Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

Probabil **Sfârșit**ică Probabilităti și Statistică

# Bibliography

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică

- Proba
  - Freedman, D., R. Pisani, R. Purves, *Statistics*, W. W. Norton & Company, 4th edition, 2007.
- Johnson, R., P. Kuby, *Elementary Statistics*, Brooks/Cole, Cengage Learning, 11th edition, 2012.
- Shao, J., Mathematical Statistics, Springer Verlag, 1998.
- Spiegel, M. R., L. J. Stephens, *Theory and Problems of Statistics*, Schaum's Outline Series, McGraw Hill, 3rd edition, 1999.

Probabilități și Statistică Probabilități și Statistica Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică Probabilități și Statistică