

Seminar 2

Deschideți fișierul **Seminar 2_Biostat.xls** (dedicat „Statisticii descriptive”) și salvați-l sub numele: "Nume_seminar2.xls" în directorul dedicat cursului. Fișierul conține mai multe foi (worksheets)!

Partea 1: determinarea marimilor caracteristice **tendinței centrale**

Fișierul conține o foaie (worksheet) numită "tendinta centrala" în care găsiți numărul de pacienți care s-au prezentat la UPU în luna august, pe zile.

Calculați media aritmetică și medianul (pentru numărul de pacienți) în două moduri: modul "de mână" și modul "EXCEL".

Etichetați viitoarele calcule (ex: celulele A34-A37) pentru a ști exact ce conține fiecare celulă (etichetarea se poate face pentru fiecare celulă înainte sau după ce a fost introdusă formula):

	A	B
34	suma	
35	numar date	
36	medie	
37	medie (Excel)	

38


Media aritmetică

Pentru calculul mediei aritmetice trebuie să folosim formula pentru media aritmetică:


$$\bar{x} = \frac{\sum x}{n} \quad (1)$$

I. se calculează suma pacienților din luna August:

- click în celula destinată calculului sumei tuturor valorilor (B34), apoi tastezi (ignorând ghilimelele) "=SUM("
- selectați celulele ce conțin numărul de pacienți folosind mouse-ul - veți observa referirea la celule (cell references) apărând după "=SUM(" ca parte a formulei.
- apoi apăsați ENTER și programul EXCEL va completa formula: "=SUM(B2:B32)"

Nota: Acest calcul se poate face folosind icoana  (function wizard), prin găsirea funcției SUM la capitolul "Math & Trig" și urmarea pașilor indicați.


II. se calculează numărul de date (n) cuprinse în setul de date:

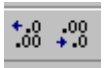
- click pe celula în care doriți să apară numărul de date (B35)
- introduceți "=COUNT(domeniu)" unde selectați domeniul (range) în același mod cum ați făcut la sumă,
- numărul de date se poate calcula și folosind icoana : se găsește funcția "COUNT" și apoi se urmează pașii indicați.

III. se calculează media aritmetică folosind formula (1):

- în următoarea celulă (B36), tastezi "="
- click pe celula B34 unde ați calculat suma: formula trebuie să se schimbe în "=B34"
- introduceți simbolul împărțit (/) folosind tastatura (keyboard)
- click pe celula B35 în care ați calculat numărul de date: formula trebuie să se schimbe în "=B34/B35"
- apăsați tasta ENTER: răspunsul va fi afișat în celulă, iar formula introdusă este vizibilă în bara de formule de sub meniu.


IV. se calculează media aritmetică folosind funcția EXCEL dedicată "=AVERAGE(range)".


- click pe următoarea celulă goală (B37).
- introduceți formula `"=AVERAGE(range)"` unde pentru domeniu (range) sunt selectate celulele cu datele a caror medie vrem să o calculăm. Funcția "Average" se poate găsi folosind icoana , apoi se urmează pașii indicați.

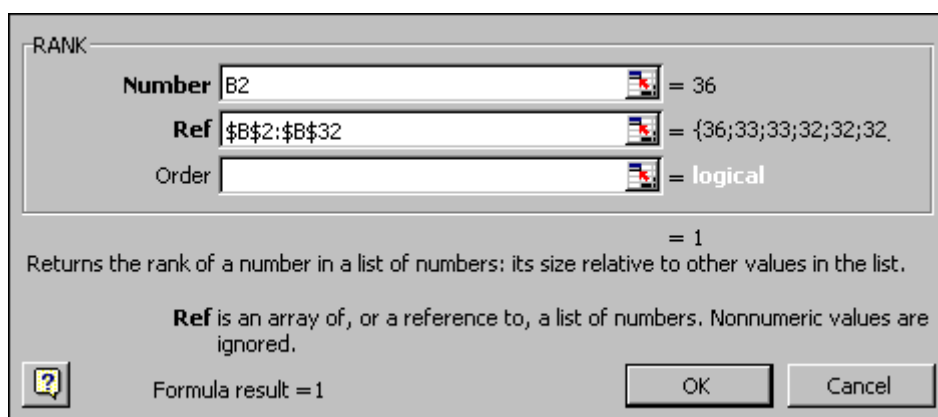
NOTA: folosiți icoanele  pentru a afișa rezultatul, din celulele B36 și B37, cu precizia dorită (setează numărul de zecimale după virgulă): întâi se activează (click) celulele care ne interesează și apoi se folosește una din icoanele indicate.

Medianul

Pentru a găsi medianul, mai întâi trebuie ordonate datele, iar apoi trebuie găsită data care are ordinul $(n - 1) / 2$.

I. pentru a ordona datele folosiți funcția "RANK" care se găsește cu ajutorul icoanei  (function wizard).

- etichetați coloana în care veți calcula ordinul datelor: tasteați "Ordin" în prima celulă din dreapta etichetei "Cluj" (C1).
- selectați celula situată la dreapta primei valori din coloana cu număr de pacienți (C2)
- click pe icoana 
- selectați funcția RANK din subsetul "statistics"
- click pe câmpul "Number" și apoi folosiți mouse-ul pentru a selecta prima valoare din coloana corespunzătoare numărului de pacienți (B2).
- click pe câmpul "Ref", apoi folosiți mouse-ul pentru a selecta toate valorile din domeniul dorit (de la B2 la B32).
- apăsați tasta "F4" pentru a face trimitere absolută la domeniul de celule (B2:B32).
- dacă se dorește ordonare ascendentă se pune orice cifră diferită de 0 în câmpul "Order". Dacă nu se pune nimic, sau se pune "0", se va face ordonare descendentă.
-



- click "OK" - celula trebuie să indice ordinul primei valori din setul de date dedicat numărului de pacienți.

- la final, copiați celula cu formula în celulele goale de deasupra - click pe mica cruciuliță neagră din partea de jos-stânga a celulei cu formula și trageți-o în jos până când este acoperit tot setul de date.

NOTA: Prin setarea, în prima celulă în care se introduce formula, a referinței absolute la domeniu de celule (B2:B32), partea "Ref" din formulă rămâne constantă în toate celulele de mai jos, permițându-ne să raportăm fiecare celulă la același grup de referință (ex. numărul de pacienți). Click pe orice celulă din coloana C pentru a vedea cum s-a schimbat conținutul formulei.

II. Pentru a determina medianul trebuie cunoscut ordinul median:

- o formula "COUNT" ne dă valoarea lui n (ar trebui să se obțină $n = 31$ - numărul de date din set)
- o ordinul median este dat de formula: $(n - 1) / 2$.
- o Efectuând calculul se obține $30 / 2 = 15$
- o medianul este acea valoare care are ordinul 15, deci medianul este 29

III. Acum vom determina medianul folosind funcția "MEDIAN":

- o click pe celula dedesubt celei în care s-a calculat media aritmetică.
- o introduceți formula "`=MEDIAN(range)`", unde "range" este domeniul de celule în care se află datele de interes.
- o nu uitați să etichetați acest rezultat înainte de a salva fișierul.

Modul

Calculul "de mână" a modului este dificil în EXCEL, deci este mai simplu să se folosească funcția "MODE".

- În aceeași manieră a funcțiilor "MEDIAN" și "AVERAGE", se calculează "modul" datelor din coloana B (fie prin tastarea formulei "`=MODE(range)`", fie prin folosirea icoanei funcție.

Partea 2. Crearea histogramelor de frecvență

Programul EXCEL are o unealtă folosită dar limitată pentru a crea histograme, pe care le vom folosi pentru a explora distribuția în frecvență a datelor.

Pentru a crea histograma numărului de pacienți urmați indicațiile următoare.

- Pentru început trebuie definit un set de categorii ("bins") pentru histograma de frecvențe. În acest scop trebuie determinate valoarea minimă (MIN) și maximă (MAX) din setul de date.

- o într-o zonă goală a foii (ex. celula F1), se introduce textul "numar" - acesta este eticheta pentru setul de categorii (bins).

- o în celula imediat dedesubt (F2), introduceți numărul corespunzător minimului setului de date - acesta este cea mai joasă categorie.

- o în celula F3, introduceți formula "`=F2+2`" (unde "F2" este celula referință pentru cea mai joasă categorie).

- o copiați formula în jos pe următoarele celule până când ajungeți la valoarea maximă a setului de date - acum ar trebui să aveți o listă cu categorii între 18 și 36.

- o *NOTA: în acest caz s-a folosit o variație a categoriei (bin size) cu 2, aleasă aleator. O regulă generală (pentru determinarea variației ideale de la o categorie la alta) este aceea că numărul de categorii (bins) trebuie să fie $1,0 + 3,3 \log(n)$, unde n este mărimea eșantionului; mărimea unei categorii este domeniul $(1,0 + 3,3 \log(n))$, sau în EXCEL (referința la celulă este arbitrară): $=(\text{max}(A1:A10)-\text{min}(A1:A10))/(1+3.3*\log(\text{count}(A1:A10)))$*

- În continuare folosiți "unealta" Histogram din Analysis Toolpak al programului EXCEL pentru a crea o histogramă.

- o În bara de meniu EXCEL, click pe "Tools" apoi pe "Data Analysis..."

NOTA: Dacă nu găsiți opțiunea "Data Analysis...", va fi nevoie să o instalați folosind opțiunea "Add Ins..." din meniul "Tools". Cereți ajutor dacă nu sunteți sigur ce trebuie făcut.

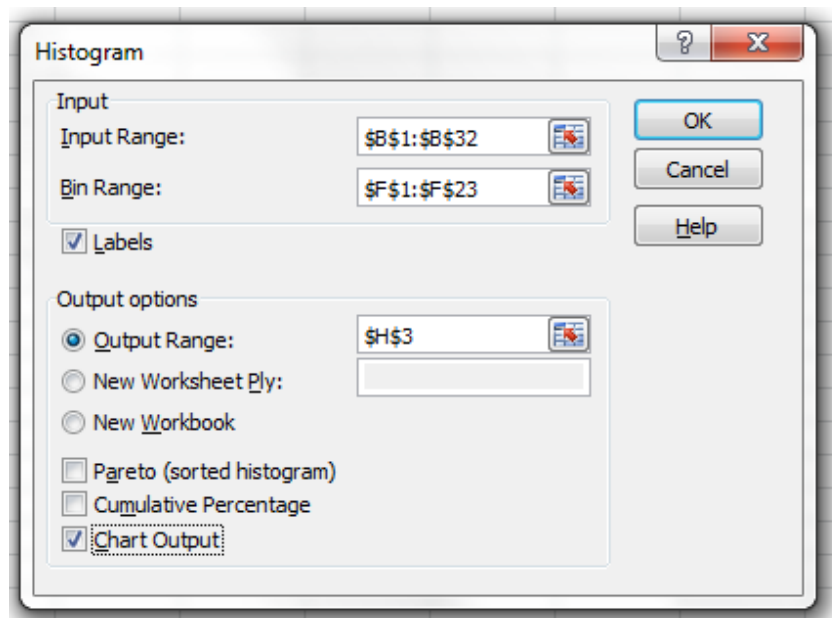
- o Selectați "Histogram" din opțiunile afișate în fereastră "Data Analysis", click OK și se va deschide fereastra Histogram.

- o Completați următoarele câmpuri, făcând click în câmp și apoi folosiți mouse-ul pentru a selecta domeniul potrivit din foaia de lucru:

- Input Range: datele "nr pacienți", *incluzând titlul*
- Bin Range: categoriile pe care tocmai le-ați creat, *incluzând titlul*.
- Output Range: o celulă din partea de sus stânga a unei părți goale din foaia de lucru (ex. H3)

- o Click pentru a selecta opțiunile "Labels" și "Chart Output".

- o Când fereastra "Histogram" va arata ca imaginea de jos, click "OK":



EXCEL va produce:

- un tabel în care sunt trecute categoriile (sau bins) ce le-ați definit, și va lista frecvențele de apariție a datelor în fiecare categorie.
- un grafic ce reprezintă histograma frecvențelor pentru setul de date.

Petreceți ceva timp pentru a vă juca cu graficul, pentru a vedea cum se poate schimba modul în care arată:

- click pe grafic pentru a-l selecta, mutați-l undeva unde îl puteți vedea și apoi ajustați-i mărimea la o dimensiune rezonabilă făcând click și apoi trageți-l de unul din colțuri.
- click-dreapta oriunde pe grafic pentru a schimba opțiunile generale ale graficului.
- dacă doriți puteți să editați graficul:
 - click pe orice etichetă pentru a edita un text - ex. schimbați titlul din "Histogram" în "UPU - August"
 - click pe legendă și deplasați-o în altă poziție, sau chiar stergeți-o de vreme ce aveți un singur set de date.
 - dublu-click pe orice parte a graficului pentru a schimba opțiunile: ex. pe axa OY (axa verticală) pentru a schimba scala, culoarea etc.

Partea 3: Determinarea marimilor ce descriu variabilitatea unui set de date

Pentru a determina dispersia unui set de date vom explora funcțiile grafice ale programului EXCEL, precum și capacitățile programului EXCEL de a calcula măsurile dispersiei.


Graficele create în programul EXCEL sunt legate de setul de date, orice modificare adusă setului de date este vizualizată imediat prin modificarea graficului.

I Vom folosi același set de date (numarul de pacienti tratați la UPU în luna august).

1. Calculați manual deviația standard a numărului de pacienti.

- Pentru a face acest calcul trebuie să determinați câteva valori intermediare de care aveți nevoie (abaterea, suma patratelor abaterilor, etc.). Nu uitați să puneți etichete acestor valori!
- Amintiți-vă formulele și funcțiile folosite în practica precedentă, precum și tehnicile de trimitere la celule: "cell referencing".
- În plus, veți avea nevoie de următoarele funcții:
 - =SUM(A1:A10) pentru a calcula suma datelor ce se află în celulele A1-A10;
 - =SQRT(A1) pentru a calcula radicalul valorii din celula A1;
 - =A1^2 pentru a calcula pătratul valorii din celula A1.

(în exemplele de mai sus trimiterea la celule este arbitrară – pentru a obține rezultatul corect, va trebui să folosiți o trimitere la celule adaptată setului de date folosit).


2. Calculați deviația standard pentru setul de date "**Numar pacienti**" folosind următoarele două funcții EXCEL prin activarea icoanei .

- =STDEV(A1:A10)
- =STDEVP(A1:A10)
-

3. Comparați răspunsurile obținute pentru pașii 1 și 2. De ce credeți că sunt două valori diferite pentru pasul 2 și care este mai apropiat de valoarea reală?

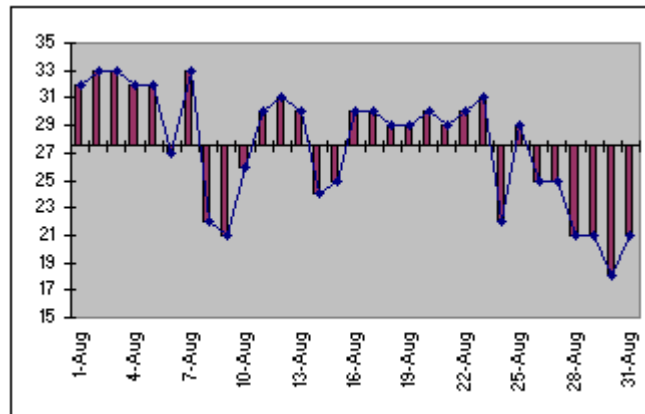
Indiciu: citiți descrierea celor două funcții în fereastra "EXCEL help".

4. Folosiți "Chart Wizard" pentru a reprezenta grafic setul de date:

- Folosiți mouse-ul pentru a selecta numarul de pacienti din setul de date.
- Deschideți "chart wizard" făcând click pe icoana  din bara de meniu și creați un grafic tip coloane (column chart).
- Adăugați un titlu, precum și numele celor două axe, click next.
- Click pe opțiunea "as object în worksheet 1" pentru a crea graficul în același worksheet cu datele, și apoi click "Finish".
- Acum aveți o reprezentare grafică elementară a datelor, reprezentare ce poate fi îmbunătățită în diferite moduri:
 - Adăugați informațiile din prima coloană (datele din luna August pentru care s-au înregistrat acele valori) pe axa x făcând click-dreapta cu mouse-ul așezat pe suprafața graficului și selectând "Source Data".
 - Selectați fereastra "Series" și introduceți trimiterea la domeniul celulelor din prima coloană în căsuța "Category (x) axis labels".
 - Schimbați formatul datelor: click-dreapta pe axa x, selectați "Format Axis", apoi "Number". Alegeți ce variantă doriți.

5. Pentru a avea o mai bună ilustrare a modului în care deviația standard derivă din date, vom schimba tipul graficului.

- Selectați "Source Data", apoi selectați "Add a new series", apoi vom selecta aceleași date (nr. pacienti) după ce casuța "Values" a fost activată
- Click-dreapta pe una din coloanele noii serii pe care ați adăugat-o, apoi selectați "Chart Type":
 - în "Custom Types", selectați graficul "line-column".
- Acum formatați axa y astfel încât axa x să o intersecteze la media aritmetică a numărului de pacienti (calculată anterior):
 - click-dreapta pe axa y a graficului, alegeți "Format Axis", "Scale" și introduceți în căsuța "Category (X) Axis Crosses At" media aritmetică.
- La final graficul vostru trebuie să arate ca cel de jos.
- În continuare puteți să modificați caracteristicile graficului, încercând cât mai multe opțiuni, pentru a dobândi abilități diverse de operare a graficelor. De exemplu puteți adăuga liniaturi orizontale sau verticale (gridlines), puteți să faceți fundalul graficului alb, sau altă culoare, puteți schimba culoarea barelor, puteți schimba domeniul axei y, etc. O procedură utilă uneori este mutarea etichetelor corespunzătoare marcajelor (tick-mark labels) folosind meniul de formatare al axei x.
- Indiciu: dacă doriți să printați numai graficul, click pe grafic pentru a-l selecta, apoi alegeți "File" și "Print" din meniul EXCEL și selectați opțiunea "Selected Chart" în secțiunea "Print What".



6. Incercați să răspundeți la următoarele întrebări:

- În ce mod este graficul, pe care l-ați creat, similar cu calculul de mână (longhand) a deviației standard la pasul 1?
- Care este media aritmetică a tuturor valorilor (\bar{x}) pe care le-ați calculat? De ce?
- De ce calculul deviației standard necesită ridicarea la pătrat a acestor valori?
- Ce diferență ar fi dacă s-ar folosi valoarea absolută a diferențelor ($|x - \bar{x}|$) în loc de pătratul lor pentru a calcula deviația standard? (Acest calcul ne dă media aritmetică a deviațiilor - folosiți funcția EXCEL "=AVEDEV()" și comparați rezultatul cu deviația standard.
- Încercați să descrieți în cuvintele voastre ce reprezintă deviația standard. *Indiciu:* citiți explicațiile din "EXCEL help"

•

7. Calculați deviația cvartală pentru setul de date:

- Folosiți funcția QUARTILE pentru a calcula valorile medii ale cvartalelor de jos LQ (sau Q1) și de sus UQ (sau Q3) corespunzătoare setului de date.
- Introduceți propria voastră formulă pentru a calcula deviația cvartală: $(Q3 - Q1)/2$.
- Etichetați celulele pentru a ști mai tarziu ce reprezintă (ex. "Q1", "Q3" and "QD")

•

8. Calculați domeniul de dispersie:

- determinați în două moduri valorile maxime și minime ale setului de date folosind:
 - funcția QUARTILE
 - funcțiile MAX și MIN
- introduceți o formulă proprie pentru a calcula domeniul de dispersie, folosind celulele ce conțin maximul și minimul setului de date.

•

9. Calculați coeficientul de variație (CV) pentru setul de date creând propria formulă în EXCEL. Formula pentru calculul coeficientului de variație o găsiți în partea de teorie.

10. Folosiți funcțiile "SKEW" și "KURT" pentru a calcula parametrii de asimetrie și de formă corespunzători datelor.

Tema:

P1. Pentru evaluarea efectului anticonceptionalelor asupra tensiunii arteriale sistolice, s-a realizat un studiu pe un esantion de 10 persoane, care a furnizat urmatoarele date:

TAS (inainte de administrare)	TAS (dupa administrare)
115	128
112	115
107	106
119	128
115	122
138	145
126	132
105	109
104	102
115	117

- Calculati media diferentei tensiunilor arteriale sistolice obtinute inainte si dupa administrarea de anticonceptionale.
- Calculati varianta si deviatia standard a TAS (inainte si dupa) precum si a diferentei.
- Calculati mediana diferentei TAS.
- Stabiliti gradul de omogenitate al esantionului din punctul de vedere al nivelului Tensiunii arteriale sistolice inainte de administrare.

P2: Urmatoarele date reprezinta varsta inbolnavirii (in ani) de o boala "A" in 25 cazuri de aparitie a acestei boli (selectate aleator):

39, 50, 26, 45, 71, 51, 33, 40, 40, 51, 63, 55, 36, 57, 41, 61, 47, 44, 48, 59, 42, 47, 53, 54, 47.

- Calculati cu o zecimala urmatoarele statistici: mediana, modul, media aritmetica, domeniul de dispersie, deviatia cvartala, varianta, deviatia standard, coeficientul de variatie
- Cate din observatii cad in afara urmatoarelor intervale: $\{\bar{x} \pm 1 \cdot s\}$, $\{\bar{x} \pm 2 \cdot s\}$
- Determinati nivelul de omogenitate al esantionului din punctul de vedere al varsei de imbolnavire.

P3: Calculați următoarele valori pentru datele ce reprezintă numărul mediu de pacienți ce au avut boala "A" la cabinetele de familie din orașul "B"

- medianul
- media aritmetică
- deviația standard
- deviația cvartală
- clasa modală (folosiți "histogram wizard" pentru a construi un tabel al frecvențelor, folosind o mărime pentru intervalul dintre frecvențe "BIN" de 35 și pornind de la 0)
- coeficientul de variație (creați propria voastră formulă).
- parametrul de asimetrie (skewness).

2358, 4980, 4140, 8506, 4472, 3155, 4070, 2364, 6069, 5578, 9693, 10721, 2044, 12166, 4643, 7153, 15193, 6405, 5122, 5760, 5354