

Compararea mediilor a mai mult de două grupuri (analiza varianței = ANOVA)

În paginile anterioare am vorbit despre compararea mediilor a două grupuri (exemplul cu TA ale studenților la medicina și ASE), arătând că testele statistice (Student = t) compară atât mediile, cât și varianțele, și că ele pot fi aplicate numai dacă distribuția variabilei numerice (în cazul acela TA) este normală (Gaussiană).

Sunt cazuri în care avem, însa, de comparat între ele mai mult de două grupuri - să presupunem ca vrem să vedem dacă diferă TA ale studenților tuturor institutelor din București, și atunci avem tot atâtea grupuri câte institute de învățământ superior există.

Anul trecut am predat metodologia cercetării științifice la patru serii de studenți, iar notele de la testul scris au fost cele din Figura 1 (reprezentare sub formă de *boxplot*).

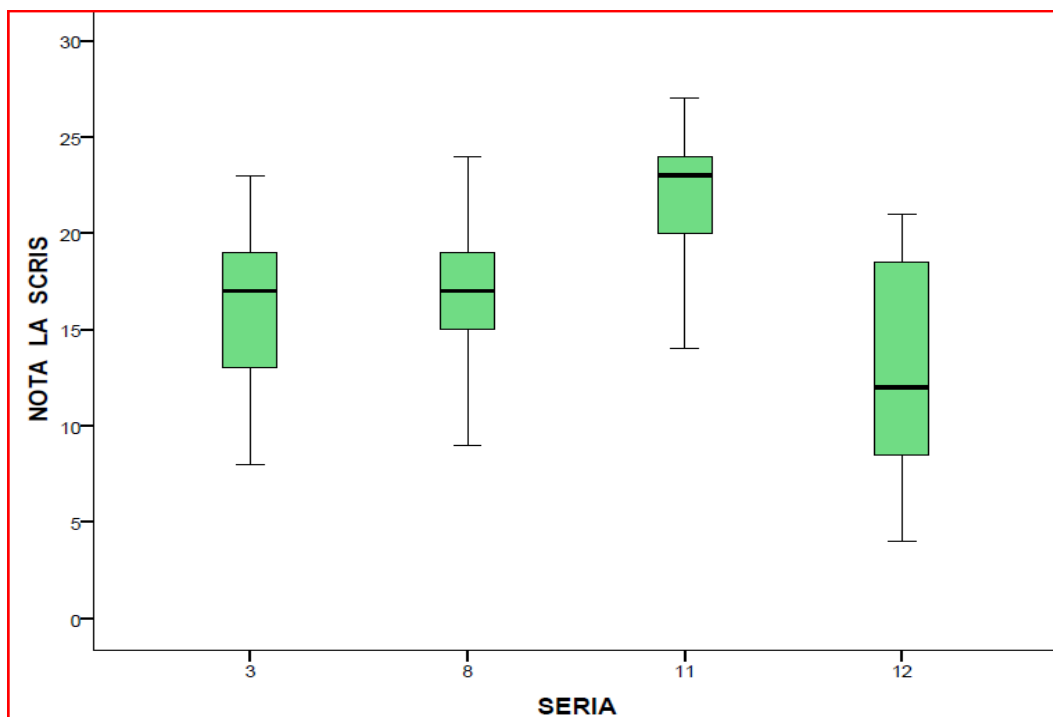


Figura 1. Notele a 4 serii de studenți la testul scris.

Vă aduc aminte dintr-un articol anterior că, la reprezentarea sub formă de *boxplot*, linia centrală reprezintă mediana, marginile orizontale ale dreptunghiului reprezintă cvartilele 25 (inferioară), respectiv 75% (superioară), iar „T”-urile reprezintă limitele. Se observă că *boxplot*-urile sunt relativ simetrice, așadar distribuția este probabil gaussiană, deci putem aplica teste statistice parametrice, respectiv ANOVA.

Baza de date rezultată va arăta ca în Figurile 2 (în programul InStat), sau 3 (în programul SPSS).

	Group A	Group B	Group C	Group D
	seria 3	seria 8	seria 11	seria 12
... 1				
... 2				
... 3				
... 4				
... 5				
... 6				
... 7				
... 8				
... 9				
... 10				
... 11				

Figura 2. Baza de date în InStat. Pe orizontală avem seriile de studenți (variabila nominală), iar în tabel vor fi trecute notele fiecărui dintre studenți (variabila numerică, ale cărei medii sunt subiectul comparației). Nu are importanță ordinea în care le scriem pe verticală.

The screenshot shows the SPSS Data Editor window titled '*Untitled1 [DataSet0] - SPSS Data'. The menu bar includes File, Edit, View, Data, and Transform. Below the menu bar, there are icons for file operations and a status bar showing '1 : nota' and '26,0'. The main data grid has two columns: 'nota' and 'seria'. The data is as follows:

	nota	seria
1	26	3
2	15	8
3	8	3
4	23	3
5	23	8
6	25	8
7	19	8
8	28	11
9	11	11
10	9	12
11	7	12
12	23	3
13	14	8
14	18	12
15	22	11
16	13	8

Figura 3. Baza de date în SPSS. În acest program, pe orizontală avem variabilele, în cazul nostru 2: nota (variabila numerică) și seria (în coloana a doua, care poate lua doar patru valori, fiind vorba de patru serii de studenți de anul IV).

În ANOVA, testăm ipoteza nulă că între seriile de studenți nu există o diferență în privința pregătirii la această materie; ipoteza alternativă este aceea că există, totuși, diferențe între serii, și la o privire sumară a graficului din figura 1, se pare că seria 11 are note ceva mai bune, iar seria 12 are note mai slabe. Rămâne să vedem dacă diferențele aparente sunt semnificative statistic.

Dacă testul Student (t) compara mediile și varianțele a două grupuri, testul ANOVA compară media și varianța totale a studenților celor 4 serii puși la un loc, cu mediile și varianțele fiecărei serii.

Dacă p rezultat în urma testului (care se mai numește și testul F) este semnificativ statistic, înseamnă că seriile de studenți nu sunt omogene, și între ele există diferențe, fără să ne spună unde se află aceste diferențe (chiar dacă noi intuim, pe baza graficului din Figura 1). Pentru a vedea

unde sunt exact aceste diferențe și dacă sunt semnificative statistic, trebuie să facem comparații utilizând testul t între seriile de studenți, luate două câte două, deci vom face un număr de „combinări de 4 serii luate câte 2” comparații, adică 6. Problema care apare aici este că noi vom face cam multe comparații (problema comparațiilor multiple), și se știe că, cu cât faci mai multe comparații, cu atât este un risc mai mare de a obține valori semnificative

statistic numai din întâmplare (ca să obții un $p=0,05$, la fiecare 100 de comparații, se obțin în medie 5 rezultate semnificative din întâmplare, iar la fiecare 20 de comparații, se obține un rezultat semnificativ din întâmplare, fără să existe o diferență reală; așadar, la 6 comparații, câte vom face noi, este foarte probabil să obținem o diferență între două serii fals pozitivă). Pentru a evita ca rezultatele noastre să fie fals pozitive, trebuie să fim mai severi în privința p -ului; **corecția Bonferroni** setează pragul de semnificație la $p = 0,05/\text{nr. de comparații} \approx 0,01$. Stabilind acest nou prag de semnificație statistică și comparând seriile două câte două, obținem $p < 0,01$ (adică diferențe semnificative statistic) numai între seria 11 și fiecare dintre celelalte serii, deci putem spune numai că seria 11 a fost semnificativ mai bună decât celelalte, nu și că seria 12 a fost semnificativ mai slabă, cum am fi putut bănui din figură. Seria a 11-a a dat testul ultima, iar observația noastră empirică (înainte de a aplica orice test statistic) a fost aceea că studenții ei au aflat, într-un mod oarecare, întrebările din test înainte de examen!

NB. Testul t (Student) nu este decât un test F, particularizat pentru numai două grupuri.