

## **Analiza bivariată: cum alegem între un test parametric și unul nonparametric**

După cum am mai spus înainte, alegerea testului adecvat se face, pe de o parte, în funcție de datele pe care vrem să le colectăm (tipurile de variabile), iar pe de alta în funcție de scopul nostru.

Alegerea cea mai dificilă este în cazul variabilelor numerice (atunci când datele noastre reprezintă rezultatele unor măsurători), deoarece putem alege între două familii de teste, cele parametrice și cele nonparametrice. Primele se bazează pe presupunerea că datele provin dintr-o populație cu distribuție normală, Gaussiană, iar testele cel mai des folosite sunt testul  $t$  (Student) și analiza varianței (ANOVA).

Testele pentru aplicarea cărora nu e nevoie de nici o presupunere asupra distribuției sunt numite teste nonparametrice, iar acestea aranjează în ordine valorile variabilei și apoi face comparații între grupuri (vezi capitolul trecut). Testele Wilcoxon, Mann-Whitney și Kruskal-Wallis sunt nonparametrice (se mai numesc și *teste care nu depind de distribuție*).

Uneori, este foarte simplu să alegem între un test parametric și unul nonparametric. Vom alege, clar, un test parametric atunci când suntem siguri că eșantionul nostru se trage dintr-o populație care urmărește o distribuție Gaussiană. Vom alege, dimpotrivă, un test nonparametric, în una din următoarele situații:

- Variabila este ordinală, iar populația este clar nonGaussiană. Exemple: anul de studiu la facultate, scorul Apgar al nou-născuților, un scor analog vizual pentru durere, scale Lickert, scoruri de calitate a vieții compuse din adunarea mai multor *item*-uri etc.
- Unele valori ale variabilei sunt „înafara scalei”, adică prea mari sau prea mici pentru a putea fi măsurate. Chiar dacă populația este Gaussiană, este imposibil să folosești teste parametrice pentru astfel de date. În schimb, analiza lor cu teste nonparametrice este foarte simplă: acordăm acestor valori niște valori arbitrare, fie foarte mici, fie foarte mari, iar cum testele nonparametrice țin cont numai de ordine, nu și de valoarea în sine, nu contează dacă nu știm exact valorile.
- Variabila este cantitativă, numerică, dar știm sigur că distribuția nu este Gaussiană (există, în acest caz, posibilitatea de a obține o distribuție Gaussiană și deci a face analiza cu teste parametrice după o transformare a valorilor – logaritm, reciprocă, radical etc).
- Variabila cantitativă are distribuție Gaussiană, dar dispersia diferă mult între cele două grupuri.

De multe ori, însă, este dificil să-ți dai seama dacă distribuția valorilor unei anumite variabile este normală (Gaussiană):

- Dacă eșantionul este mare (cel puțin 100), nu trebuie decât să te uiți la distribuția valorilor și se vede clar dacă este sau nu Gaussiană, pentru a ști ce fel de test să alegi. Dacă, însă, eșantionul este mic, este dificil să-ți dai seama dacă distribuția este Gaussiană sau nu prin inspecție (histogramă), iar testele statistice (Kolmogorov-Smirnov) nu au destulă putere pentru a face diferența.
- Ne putem uita și pe date similare din alte studii. Ceea ce contează este distribuția valorilor **populației**, nu eșantionului. Pentru a vedea dacă o populație este Gaussiană, trebuie să analizăm toate datele disponibile, nu numai datele din studiul la care lucrăm.

Când nu sunt siguri, unii utilizează un test parametric (deoarece nu sunt siguri că prezumția de normalitate a fost violată), iar alții aleg un test nonparametric (deoarece nu sunt siguri că distribuția este Gaussiană). După părerea mea, a doua variantă este mai corectă.

Întrebarea este în ce măsură are importanță ce fel de test, parametric sau nonparametric, alegem, și sunt patru posibilități:

- Dacă eșantionul este mare (cel puțin 24/30 de date în fiecare grup, cifra diferă între cărțile de statistică), este mai ușor de spus dacă eșantionul provine dintr-o populație Gaussiană, dar nu are mare importanță, putem folosi orice tip de test, deoarece rezultatul este același (rezultă același p).
- Dacă eșantionul este mic, distribuția nu este Gaussiană și utilizăm un test parametric, p-ul obținut nu este corect.
- Dacă eșantionul este mic, distribuția este Gaussiană și utilizăm un test nonparametric, valorile p sunt mai mari (testele nonparametrice au putere statistică mult mai mică pe eșantioane mici).

Așadar, în cazul eșantioanelor mari nu sunt probleme. Mai mult, dacă eșantionul este  $> 100$  și nu există valori extreme, cu influență disproporționată asupra analizei statistice, se pot utiliza teste parametrice fără grijă.

Dilema apare atunci când eșantionul este mic: pe de o parte, este greu de spus dacă provin din populații cu distribuție normală, și tocmai aici acest lucru este foarte important, pentru că testele nonparametrice nu sunt puternice (pe aceleași date, cu unul parametric ai obține un p mai mic decât cu unul nonparametric), iar cele parametrice nu sunt robuste.