

#### IV. Tipurile de variabile și testele statistice

Alegerea metodei de analiză statistică pentru o anumită problemă depinde de comparația pe care vrem să o facem și de tipurile de variabile utilizate. Așadar, pentru a alege testul potrivit trebuie să ne punem două întrebări: Ce fel de date am colectat? Care este scopul nostru? Aceeași analiză o facem și citind un articol, pentru a vedea dacă testele utilizate de autorii acestuia sunt cele corecte.

Variabilele nominale sunt variabile sub formă de nume sau alte simboluri reprezentând categorii ce nu pot fi ordonate una în raport cu cealaltă, de exemplu numele, grupa sanguină, sexul, rasa, culoarea ochilor, diagnosticul etc. Atunci când o variabilă nominală nu poate lua decât două valori, ea este o variabilă dihotomică (binară, bimodală), cum ar fi sex masculin/feminin, mort/viu, fumător/nefumător, prezent/absent, normal/anormal, care a suferit efectul(*end-point*)/care nu l-a suferit etc. – variabile de tip DA/NU.

Variabilele ordinale sunt variabilele ce sunt clasificate în mai mult de două categorii și la care există o ordine naturală între categorii (de la valoarea cea mai mică la cea mai mare) - de exemplu evoluția bolii (agravat, staționar, ameliorat), stadializări (insuficiența cardiacă, TNM în cancer), scoruri etc. Chiar dacă unele variabile iau valori numerice, ele sunt considerate ordinale pentru că nu îndeplinesc condițiile celor cantitative măsurabile (luând exemplul unui scor de calitate a vieții, sau al unei scale analogice vizuale pentru durere, putem spune că un individ cu scorul 10 are o durere mai mare sau o calitate a vieții mai bună decât un individ cu scorul 8, dar nu înseamnă că diferența dintre ei este aceeași cu aceea dintre un individ cu scorul 4 și unul cu scorul 2; la fel, nu putem spune că un individ cu scorul 8 are o durere de două ori mai puternică sau o calitate a vieții de două ori mai bună decât un individ având scorul 4; din același motiv, calcularea mediei nu are nici un sens).

Variabilele cantitative (măsurabile) pot fi *continue* (variabile cu un număr potențial infinit de valori de-a lungul unui continuum: înălțimea, greutatea, TA, vârsta etc.) sau *discontinue (discrete)* (variabile descrise numai prin unități întregi ce nu pot fi măsurate în intervale mai mici decât unitatea: frecvența cardiacă, numărul de copii etc.).

În privința variabilelor cantitative, este important de văzut dacă acestea au o distribuție normală (simetrică, sub forma clopotului lui Gauss); în cazul acestei distribuții, media este egală cu mediana și cu modul, iar 95% dintre valorile pe care le poate lua variabila se află în intervalul  $media \pm$  două deviații standard. Pentru a vedea dacă o variabilă are distribuție normală, putem folosi orice program statistic și verificăm cifric dacă media, mediana și modul sunt foarte apropiate (ideal identice, dar în viața reală nu există ideal), iar media minus dublul deviației standard nu trebuie să ia valori negative; sau reprezentăm variabila sub forma unei histograme și vizual verificăm forma simetrică, de clopot.

Este important să știm dacă distribuția variabilei noastre este normală pentru că *numai variabilelor cantitative, cu distribuție normală li se pot aplica testele statistice parametrice!* De asemenea, pentru a putea aplica teste statistice parametrice trebuie ca nici dispersia (deviația standard) celor două grupuri să nu difere foarte mult. Testele parametrice sunt testele care compară mediile și deviațiile standard ale grupurilor despre care vrem să dovedim că sunt identice sau dimpotrivă, diferite, ori media și deviația standard nu au nici un sens dacă nu avem o distribuție normală.

Să presupunem că vrem să vedem dacă TA este diferită la bărbații față de femeile din București; pentru aceasta, ideal ar fi să extragem la întâmplare două eșantioane, unul de femei, altul de bărbați din București și să le măsurăm TA. Cum valorile TA vor avea, probabil, o distribuție normală în cele două grupuri, atunci pentru a le compara putem folosi un test parametric, care este *testul t (Student)*. Dacă vrem să comparăm mai multe grupuri deodată (de exemplu vrem să vedem dacă TA a moldovenilor, ardelenilor sau regătenilor diferă între ele), folosim analiza varianței în sens unic (*one way ANOVA*), aplicând *testul F*, care ne va arăta dacă TA este diferită sau nu în cele trei regiuni istorice, fără a ne spune însă care este grupul care diferă de celelalte.

*Testele nonparametrice* se aplică pentru variabilele cantitative fără distribuție normală (Figura IV.1) și pentru variabilele ordinale. Să luăm un exemplu din revista Medicina Internă 2004, 1:57-59 (Dumitrașcu DL et. al), în care se compară stresul la pacienții cu dispepsie funcțională și la martori, pe baza unui chestionar. Pentru comparația scorurilor de stres s-a folosit *testul t*. În primul rând, în cazul scorurilor nu trebuie folosite teste parametrice, scorurile nefiind variabile cantitative (vezi explicația de la prezentarea variabilelor ordinale). Să ne imaginăm, totuși, că scorurile de stres sunt variabile cantitative – ca să putem aplica *testul t* ar trebui ca distribuția lor să fie normală, ori în tabelul 2 putem vedea scoruri (media și deviația standard) de genul 0,89 și 0,93; 7,5 și 5,28; 1,25 și 0,93; 0,54 și 0,60; 0,94 și 1,39 (și încă altele), din care se observă clar că distribuția nu este normală (dacă scădem din medie 2 deviații standard avem scoruri negative, care nu există în realitate), așadar pentru comparație ar fi trebuit folosit un test nonparametric (*testul Mann-Whitney U*).

Așadar, pentru variabilele cantitative care nu au o distribuție normală și pentru cele ordinale se folosesc *testele nonparametrice*. De exemplu, dacă vrem să demonstrăm că pacienții cu insuficiență cardiacă internați în spitalul X sunt mai gravi decât cei internați în spitalul Z, comparând clasa NYHA de insuficiență cardiacă între cele două grupuri.

Atunci când pacienții sunt împerecheați, folosim *testele statistice împerecheate (paired)*, parametrice sau nonparametrice. Singura împerechere perfectă se realizează atunci când împerechem pacientul cu el însuși, în comparațiile înainte-după. De exemplu, comparăm TA, sau colesterolul unor participanți la un studiu înainte de a începe tratamentul și după o lună de tratament. Variabila (TA, colesterolul) fiind continuă și cu o distribuție probabil normală, vom folosi un test parametric, și anume *testul t împerecheat*. Dacă variabila de comparat nu are o distribuție normală (de exemplu valoarea creatininei la pacienții cu insuficiență renală) sau este o variabilă ordinală (stadializarea tumorii, sau clasa NYHA a insuficienței cardiace, sau scorul durerii pe o scală analogică vizuală, înainte și după un tratament), vom folosi un test nonparametric împerecheat, care este *testul Wilcoxon*.

Echivalentul nonparametric al ANOVA (*testul F*) este *testul Kruskal-Wallis*.

Testele nonparametrice nu țin cont de valoarea efectivă a variabilei, ci de ordinea lor (*rank tests*) – care este valoarea cea mai mică, care este următoarea și așa mai departe...

În cazul variabilelor dihotomice (pentru compararea proporțiilor) se folosește *testul  $X^2$*  sau variantele sale *Yates* și mai ales *testul exact al lui Fisher* (atunci când în tabelul de contingență 2x2 avem într-una din căsuțe o valoare așteptată mai mică de 5). De exemplu atunci când vrem să comparăm proporția de pacienți care a făcut infarct în grupul tratat cu statină cu proporția de pacienți care a făcut infarct în grupul tratat cu placebo. De remarcat că în studiile terapeutice, atunci când avem de-a face cu efecte surogat studiem variabile cantitative (TA, transaminaze, clasa NYHA, fracția de ejeție,

densitatea osoasa etc.), pe când în cazul efectelor serioase avem de-a face cu variabile dihotomice (pacientul a suferit sau nu infarctul de miocard, fractura, decesul etc.).

Când vrem să vedem cum (și dacă) variază o variabilă cantitativă în funcție de o altă variabilă cantitativă, așadar vrem să vedem în ce măsură două variabile cantitative se corelează, calculăm *coeficientul de corelație al lui Pearson* ( $r$ ). De exemplu, putem vedea dacă vârsta se corelează cu VSH (adică VSH crește odată cu vârsta).

Dacă variabilele cantitative nu au o distribuție normală, sau sunt ordinale (de exemplu, corelația dintre fracția de ejeție și clasa NYHA a insuficienței cardiace stângi, sau dintre valoarea transaminazelor și cea a scorului necroinflamator găsit la biopsia hepatică) utilizăm echivalentul nonparametric al coeficientului Pearson, care este *coeficientul de corelație Spearman*.

Dacă, în cazul a două variabile care se corelează, putem spune care variabilă o determină pe cealaltă și/sau vrem să calculăm valoarea unei variabile știind-o pe cealaltă, utilizăm *regresia* lineară (de exemplu, știind valoarea ALAT, putem prezice scorul necro-inflamator de la biopsie, sau știind înălțimea prezicem valoarea VEMS, sau știind doza de captopril pe care o administrăm prezicem cu cât va scădea TA).

Variabilele cantitative pot fi transformate oricând în variabile ordinale sau dihotomice (de exemplu valorile colesterolului în quartile, sau în colesterol normal/crescut). În baza noastră de date este indicat să trecem (și pentru aceasta să culegem) variabilele noastre ca atare, pentru că apoi putem să le transformăm oricând în ordinale sau dihotomice, pe când invers nu vom putea niciodată (de exemplu introducem în baza de date anemie DA/NU și apoi descoperim că ar fi fost mai bine să avem chiar valorile hemoglobinei!).

Partea cea mai dificilă este alegerea între testele parametrice și cele neparametrice. Alegem clar un test nonparametric în trei situații: 1. efectul este o variabilă ordinală și populația este clar non-Gaussiană (de exemplu notele studenților, scorul Apgar, scala vizuală analogică pentru durere etc.); 2. efectul este o variabilă cantitativă și suntem siguri că nu are o distribuție gaussiană în populație (în acest caz o putem aduce la o distribuție normală prin transformare: logaritmul, reciproca, rădăcina pătrată – din punct de vedere matematic este corect, mai puțin din punct de vedere biologic); și 3. efectul este o variabilă cantitativă cu distribuție gaussiană, dar dispersia (deviația standard) este mult diferită între grupurile de comparat.

Deseori alegerea este dificilă. Când avem cazuri puține, este greu de spus dacă distribuția este Gaussiană, iar testele speciale pentru verificarea normalității (Kolmogorov-Smirnov) au putere mică. De fapt, ceea ce contează este distribuția la nivelul populației, și nu la nivelul eșantionului nostru, iar informații despre distribuția valorilor unei variabile în populație trebuie căutate în literatură! (Este bine de reținut că în natură, distribuțiile non-gaussiene sunt frecvente, iar acest fapt este valabil îndeosebi în cazul valorilor biologice).

Când nu știm dacă distribuția este normală, alegerea tipului de test depinde de mărimea eșantionului: dacă eșantionul este mare (cel puțin 24/30 de date în fiecare grup), este mai ușor de spus dacă eșantionul provine dintr-o populație Gaussiană, dar nu are mare importanță, putem folosi orice tip de test, rezultatul va fi același. Problema apare dacă eșantionul este mic, când este greu de spus dacă populația este gaussiană, dar tocmai atunci este foarte important: testele nonparametrice nu sunt puternice, iar cele parametrice nu sunt robuste.

**Analiza multivariabilă (sau multivariată)** este o unealtă statistică prin care se determină contribuția fiecăruia dintre mai mulți factori la apariția unui efect. De exemplu, există o mulțime de factori asociați cu apariția bolii coronariene (fumatul, obezitatea, sedentarismul, diabetul, hipercolesterolemia, hipertensiunea) – numiți *factori de risc, variabile independente, sau variabile explicative*. Analiza multivariabilă ne permite să determinăm contribuția independentă a fiecăruia dintre acești factori de risc la apariția bolii coronariene (numită *efect* sau *variabilă dependentă*).

În studiile observaționale, fiind mai mulți factori de risc, nu știm care dintre ei este adevărat sau în ce măsură asocierea aparentă dintre un factor de risc și efect nu este datorată de fapt altora. Să presupunem că în baza noastră de date avem toate informațiile și variabilele privind pacienții, și vom testa dacă există o asociere între variabilele fumat și boala coronariană, fără a ține cont de vreo altă variabilă. Neexistând randomizare (nu putem pune pacienții, prin tragere la sorți, să fumeze sau nu), chiar dacă în analiza univariată (bivariată, după alții) găsim o asociere între fumat și apariția coronaropatiei, aceasta reprezintă o dovadă prea slabă pentru a o considera cauzală. Poate că fumătorii fac mai degrabă coronaropatie pentru că sunt mai frecvent bărbați și/sau sunt săraci și/sau au mai degrabă un stil de viață nesănătos în alte privințe, care sunt adevărații factori de risc. Cu alte cuvinte, relația dintre fumat și coronaropatie poate fi *confundată* de aceste alte variabile.

Confuzia apare atunci când o asociere aparentă dintre un factor de risc și un efect este afectată de relația unei a treia variabile cu factorul de risc și cu efectul. Pentru ca o variabilă să fie un *factor de confuzie*, aceasta trebuie să se asocieze atât cu factorul de risc, cât și cu efectul.

Sexul masculin și sedentarismul pot fi factori de confuzie, deoarece sunt asociate atât cu fumatul, cât și cu boala coronariană. Prin analiza multivariabilă, putem demonstra că și după ajustarea pentru sexul masculin și sedentarism, fumatul are o relație independentă cu boala coronariană. (De altfel acest cuvânt, “ajustat” care apare într-un articol ne spune de fiecare dată că la rezultatul prezentat s-a ajuns printr-o analiză multivariabilă.) Să presupunem că în analiza univariată, boala coronariană este asociată cu consumul de cafea (riscul relativ=5, cu semnificație statistică). Bănuim, însă, că există un factor de confuzie, și anume fumatul, care este asociat cu consumul de cafea (cei care beau cafea, în general fumează) și cu efectul (am descoperit asocierea fumat-coronaropatie). Pentru a verifica această ipoteză, facem o analiză multivariată simplă, în care introducem ca variabile independente atât cafeaua, cât și fumatul, și vom vedea efectul fiecăreia dintre ele asupra apariției coronaropatiei. Dacă riscul de a face boală coronariană se menține semnificativ statistic la cei care consumă cafea, înseamnă că aceasta reprezintă un factor de risc independent pentru boala coronariană. Dacă însă, în analiza multivariată, riscul dispare, înseamnă că într-adevăr, relația cafea-coronaropatie a fost confundată de relația adevărată, fumat-coronaropatie. Bineînțeles că în analiza multivariată vom găsi o relație semnificativă statistic fumat-coronaropatie.

Deși teoretic se poate face distincția între asocierea independentă și confuzie, o variabilă poate avea în același timp un efect independent și să fie un factor de confuzie: de exemplu sărăcia este un factor de confuzie între fumat și coronaropatie (cei săraci fumează mai mult și fac mai des boala coronariană), dar sărăcia are de asemenea și un efect independent asupra apariției bolii coronariene (după ajustarea pentru fumat, colesterolemie și alți factori de risc, aceasta rămâne totuși asociată semnificativ cu apariția bolii).

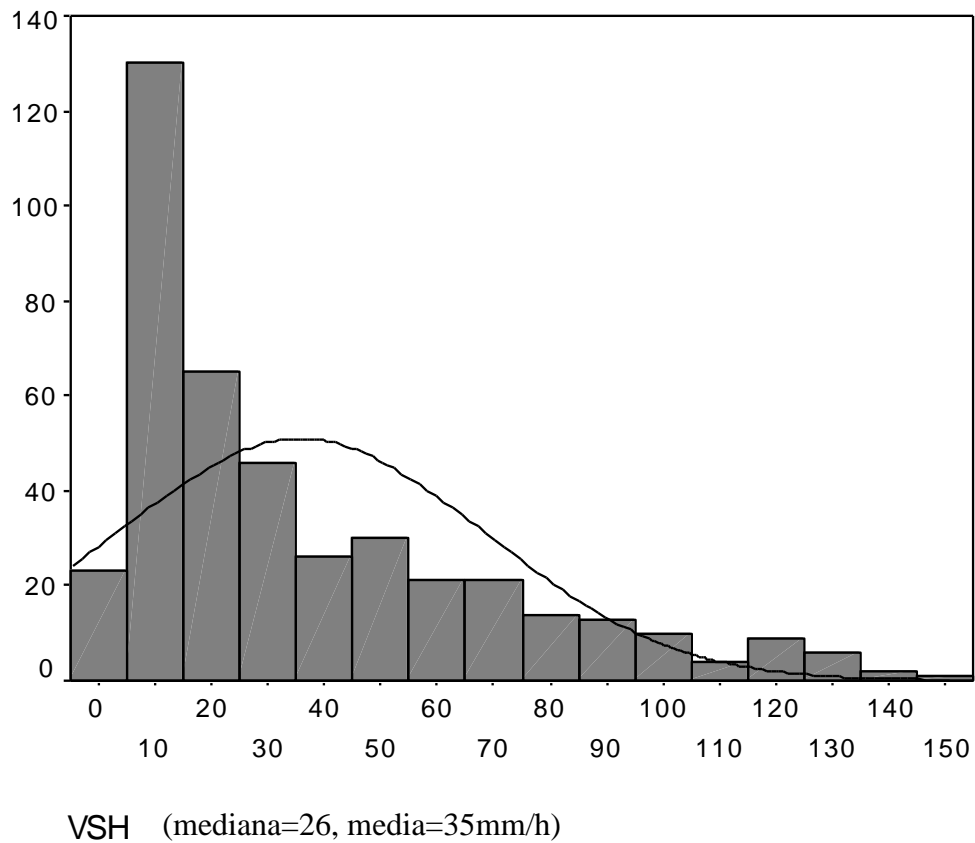
În afara analizei multivariabile, pentru eliminarea confuziei mai poate fi utilizată *analiza stratificată*, prin care se cercetează rolul unui factor de risc în apariția unui efect, în timp ce se ține cealaltă variabilă constantă. Astfel, în exemplul în care cafeaua era asociată în analiza univariată coronaropatiei, putem reface această analiză separat, la fumători și nefumători, și dacă asocierea rămâne în picioare în cele două grupuri, înseamnă că efectul cafelei în apariția coronaropatiei este independent de fumat; dimpotrivă, dacă asocierea cafea-boală dispăre, înseamnă că fumatul a fost un factor de confuzie care a determinat apariția unei false relații între cafea și boala coronariană.

Ne putem folosi de stratificare atunci când există două sau trei variabile potențiale factori de confuzie; atunci însă când acestea sunt mai multe, stratificarea ar crea zeci de grupuri în care investigatorul ar trebui să determine relația dintre variabile, iar numărul de pacienți din fiecare grup ar fi din ce în ce mai mic, pe măsură ce progresăm cu stratificarea și s-ar pierde puterea statistică.

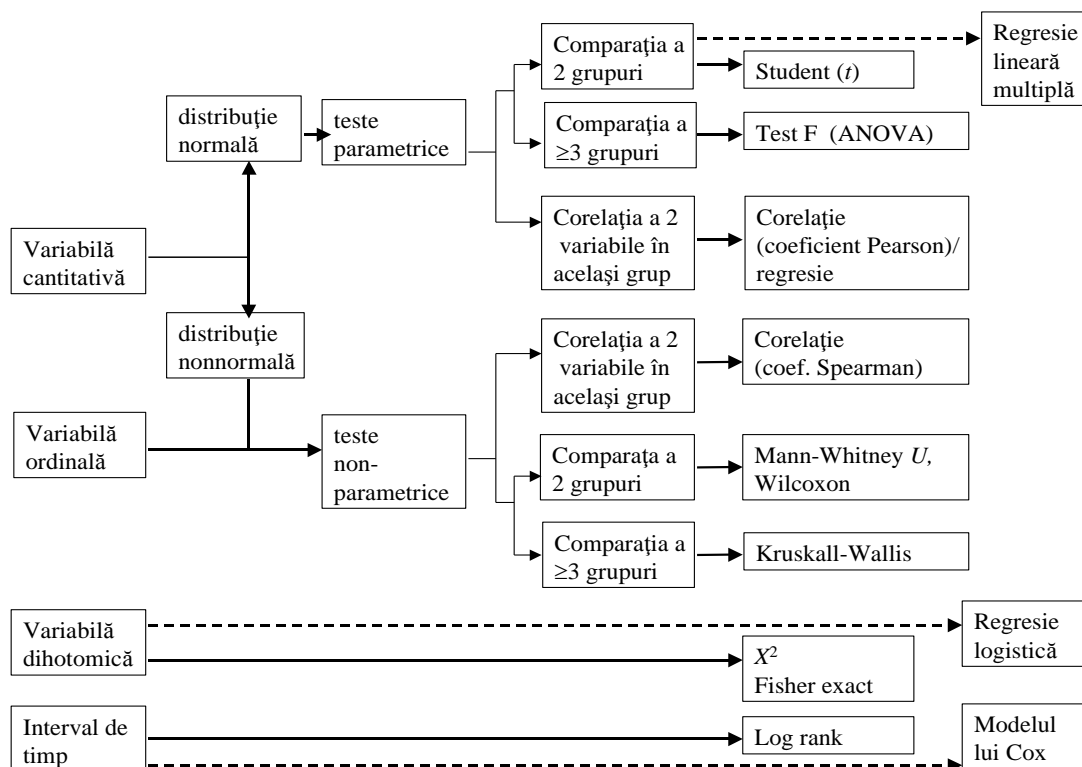
Indiferent dacă folosim stratificarea sau analiza multivariată, nu trebuie să uităm că putem stratifica sau ajusta doar pentru variabilele pe care le cunoaștem, ori există mulți factori de confuzie necunoscuți și deci nemăsurați, care ne pândesc la toate colțurile!

**Tipurile de analiză multivariabilă** sunt trei, în funcție de variabila dependentă (efectul): atunci când variabila dependentă este continuă se utilizează **regresia lineară multiplă**, dacă aceasta este dihotomică se utilizează **regresia logistică**, iar când este reprezentată prin durata de timp până la apariția unui eveniment (“supraviețuirea”), se folosește **analiza hazardului proporțional (modelul lui Cox)**.

**Figura IV.1.** Exemplu de distribuție non-normală: distribuția VSH într-un studiu (histograma) în comparație cu distribuția normală (curba lui Gauss). Se observă asimetria distribuției VSH datorită existenței unor pacienți cu VSH foarte mare, care trag media spre dreapta, în timp ce mediana nu este influențată.



**Figura IV.2.** Algoritmul utilizării testelor statistice în funcție de variabile (exemple în text).  
 ( - - - - - = analiză multivariabilă)



C Băicuș. Medicina bazată pe dovezi. Cum  
înțelegem studiile.

**Ed**  
**Medicală,**  
**2007**

---