

Reprezentarea datelor clinice în vederea prelucrării statistice.

A.D. Corlan, 2014.

Orice studiu clinic presupune extragerea unui set sintetic de date în vederea prelucrării statistice.

În cazul cel mai simplu, acest set constă într-un tabel, cu câte un rând pentru fiecare **caz** (fiecare pacient) și câte o coloană pentru fiecare **variabilă** (cum ar fi vârsta sau prezența unui diagnostic). În cazurile mai complicate, pot fi mai multe tabele cu legături între ele, sau alte structuri de date.

Unele dintre coloane sunt destinate prelucrării statistice computerizate a datelor, în vreme ce altele sunt simple comentarii.

În cele de mai jos, vom numi 'variabile' doar acele coloane care sunt destinate prelucrării computerizate.

De obicei, realizarea unui studiu include următorii pași:

- A:** Se stabilesc obiectivele studiului, se estimează pe baza unor rezultate anterioare din literatură de câte cazuri ar putea fi nevoie, se evaluează fezabilitatea.
- B:** Se stabilește protocolul studiului, inclusiv "capul de tabel", adică se aleg și se documentează variabilele corespunzătoare capului de tabel.

- C:** Se execută protocolul studiului, din care rezultă un set de documente primare (foi de observație, buletine de analiză, înregistrări, imagini, etc) pentru fiecare caz.
- D:** Se întocmește tabelul sintetic, completând valorile variabilelor în fiecare caz.
- E:** Se face o analiză computerizată a datelor din tabel, rezultând un număr de statistici. Prin **statistici** înțelegem aici numere calculate pe baza valorilor variabilelor, care caracterizează întreaga populație¹.
- F:** Se interpretează statisticile rezultate din E și se întocmește un raport asupra acestora, care apoi se sintetizează într-o lucrare.

În faza E se dovedește însă, deseori, că statisticile necesare atingerii obiectivelor din faza A nu se pot determina din cauza unor probleme tehnice ivite în fazele B sau D, adică la stabilirea și documentarea variabilelor și la completarea valorilor.

Este atunci necesară corectarea capului de tabel și reluarea fazei D, uneori și C (ceea ce poate fi imposibil sau foarte costisitor), ceea ce poate duce la întârzieri substanțiale.

Scopul acestui document este să enumere cele mai frecvente erori și probleme de acest tip, pur tehnic, în speranța că astfel se vor preveni sau reduce întârzierile respective.

Cea mai utilă recomandare generală este: documentați în scris, în faza B, capul de tabel cu: (i) numele fiecărei variabile; (ii) unitatea de

¹În termeni tehnici: parametrii estimați ai distribuției probabilității variabilei.

măsură; (iii) intervalul sau enumerarea valorilor pe care le poate lua; (iv) descrierea, narativă, a semnificației variabilei.

Discutați apoi această documentație cu statisticianul, înainte de a începe colectarea datelor. Dacă nu ați făcut această documentare și discutare, faceți-o acum.

TOATE VALORILE VARIABILELOR PRELUCRABILE SUNT DE FAPT NUMERE

Pentru a putea face prelucrarea statistică, calculatorul trebuie mai întâi să traducă automat fiecare valoare a unei variabile într-un număr.

Calculatorul nu poate interpreta, în acest scop, decât tipuri foarte limitate de expresii (formulări, date introduse), și anume:

- (1) Valorile unei variabile pentru fiecare caz sunt date chiar prin numere, întregi sau fracționare (cu virgulă). În acest caz trebuie stabilit la documentarea capului de tabel ce fel de număr (întreg sau fracționar), care e valoarea minimă sau maximă, și ce reprezintă el. Unitatea de măsură trebuie să fie aceeași pentru toate valorile unei variabile.
- (2) Valorile sunt simboluri—cuvinte sau acronime dintr-un set predefinit, de exemplu 'bărbat/femeie'—care sunt ulterior transformate în numere întregi în cursul prelucrării. Setul predefinit și semnificația fiecărei valori trebuie stabilite la documentare. Este esențial ca fiecare semnificație să fie

reprezentată prin exact același simbol. De exemplu, nu putem reprezenta sexul bărbătesc, în aceeași coloană, o dată prin 'B' și altă dată prin 'M' sau 'm', ci trebuie să alegem de la început un singur simbol.

- (3) Valorile sunt momente în timp, cum ar fi data nașterii sau data unei intervenții, care sunt introduse cu an/luna/zi/ora/minut, dar în vederea prelucrării vor fi transformate de calculator tot în numere simple—de pildă număr de secunde scurse de la o dată unică din trecutul mai îndepărtat. Trebuie precizat, la documentarea capului de tabel, ce reprezintă fiecare dată și cu ce acuratețe este măsurată (un an, o zi, ora și minut, etc).

Lipsa unei valori poate fi înregistrată printr-un spațiu sau, preferabil, printr-un simbol specific, cum este 'NA'.

Un exemplu de eroare frecventă în această privință ar fi o coloană cu titlul "diagnostic" sau "boli.asociate" care este completată cu un scurt text, eventual cu prescurtări, de genul "ICC, DZ.II". O astfel de coloană va fi tratată ca un comentariu și informația pe care o conține nu va putea fi folosită în scop statistic, decât dacă, mai întâi, este recodificată sub forma unor variabile obișnuite—de exemplu: ce tip de diabet are pacientul (coloana cu titlul DZ și cu valorile 0 reprezentând lipsa bolii, 1—diabet de tip I, 2—diabet de tip II).

Dacă datele disponibile în documentele primare—de exemplu, buletine de analiză—sunt cu diferite unități de măsură, cercetătorul care introduce datele trebuie să facă o conversie a acestor unități în unitatea stabilită pentru variabila respectivă.

În cazul că acest pas este dificil, se poate face o a doua variabilă, cu altă unitate de măsură, urmând să se realizeze conversia la prelucrarea datelor. De exemplu, glicemia *à jeun* în prima zi de internare, poate fi reprezentată în două coloane alăturate, GLIC0MGDL și GLIC0MMDL. În cazul în care este disponibilă în mg/dl se scrie în prima coloană, iar în cea de a doua se notează NA, iar în cazul în care valoarea este disponibilă în mmol/dl, se scrie în a doua coloană, iar în prima se notează NA.

UNELE VARIABILE SE POT CALCULA DIN ALTELE

O dată ce datele introduse au fost traduse în numere, se pot calcula automat valori pentru bazate pe aceste numere, din care să rezulte noi variabile.

Un exemplu, în continuarea celui de mai sus cu glicemia, este calculul glicemiei în mg/dl, GLIC0 (care nu este completată sau măcar prezentă în tabelul inițial) plecând de la GLIC0MGDL și GLIC0MMDL. Calculatorul va folosi o formulă (introdusă de statistician) care va lua în seamă care dintre valori este prezentă și va face conversia dacă este necesar.

Un alt exemplu este calculul vârstei când sunt prezente data nașterii și o data curentă, a unei faze din protocol; sau calculul indicelui de masă corporală sau a suprafeței corporale estimate plecând de la greutatea corporală și înălțimea pacientului.

Aceste calcule nu trebuie făcute de către experimentator înainte de introducerea datelor deoarece consumă timp și pot introduce erori. Este însă important să se stabilească de la început, în măsura posibilului, ce variabile se estimează a se calcula, pentru a verifica dacă toate datele necesare sunt prevăzute în tabel.

Este recomandat ca, pentru numele de variabile, se se folosească acronime scurte, care încep cu o literă și conțin litere, cifre, și în mod excepțional caractere de despărțire cum este ‘_’ (liniuța de subliniere).

Este dezirabilă evitarea numelor lungi, descriptive, mai ales conținând spații, care crează tot felul de dificultăți și confuzii ulterior.

VARIABLELE REPREZINTĂ MĂSURABILE

Rezultatul publicat al cercetării va fi o metodă predictivă prin care, pe baza unor măsurători, să se poată estima probabilitatea de apariție a unor evenimente clinice. De exemplu, pe baza măsurătorilor hemoglobinei glicozilate se prezice incidența neuropatiei diabetice.

Prin măsurătoare se înțelege o determinare obiectivă a unei valori. Esențială pentru aplicabilitatea rezultatelor este reproductibilitatea independentă a măsurătorii.

De exemplu, părerea subiectivă a clinicianului privind tipul de germen care produce o infecție, chiar dacă—să presupunem—bazată pe o vastă experiență și frecventă confirmată prin investigații independente, nu este o măsurabilă, pentru că nu este reproductibilă

obiectiv în mod independent. Cititorul rezultatelor studiului poate avea și el păreri personale privind tipul de germen dintr-o infecție, dar nu ne putem aștepta ca aceste păreri să aibă aceeași valoare predictivă—pentru că depind de experiența lui diferită—ceea ce face rezultatul obținut inaplicabil.

STUDIILE TREBUIE SĂ FIE CÂT MAI UȘOR DE COMPARAT CU ALTE STUDII

După publicare, rezultatele vor fi valorificate în continuare, în alte studii, de exemplu în meta-analize. Oricum, ele vor trebui confruntate cu rezultatele altor studii. Din acest motiv, este important ca, dacă există deja în literatură niște practici în alegerea variabilelor într-un anumit tip de studiu, să se urmeze aceste practici dacă se poate.

EVITAREA PIERDERII DE INFORMAȚIE

Variabilele alese reprezintă întotdeauna niște reprezentări simplificate ale realității. Această simplificare înseamnă că o parte din informația prezentă în documentele primare se pierde.

De exemplu, în documentul primar (buletin de analiză) se găsește valoarea glicemiei *à jeun*. Aceasta este deja o reprezentare simplificată a fenomenului mai complex al dinamicii glicemiei în cursul zilei, funcție de alimentație, efort și alți parametri.

Experimentatorul poate face eroarea de a reprezenta această valoare sub forma unei variabile simbolice: 1: hipoglicemic, 2: normoglicemic, 3: hiperglicemic alegând mai mult sau mai puțin arbitrar niște praguri.

Această reprezentare este o măsurabilă (presupunând că pragurile au fost bine documentate și respectate) și se poate prelucra statistic, dar o mare parte din informația prezentă în valoare exactă a glicemiei se pierde.

Este mai simplu în acest caz să se reprezinte informația primară (valoarea glicemiei) urmând ca, dacă din diverse motive—cum ar fi comparația cu datele altui studiu—dorim ulterior să folosim niște praguri să facem asta automat, printr-o variabilă calculată din valoarea primară.

Putem ilustra superioritatea reprezentării valorii primare prin faptul că ne permite comparația cu studii diferite care folosesc seturi diferite de praguri. Dacă alegem noi niște praguri de la început nu mai putem face această comparație decât cu studii care folosesc același set.

În general, este bine să introducem în tabel datele așa cum sunt disponibile în documentele primare, urmând ca orice prelucrare să se facă ulterior, automat.

INDEPENDENȚA VARIABILELOR

Uneori se face eroarea de a se introduce variabile care se pot calcula unele din altele, ceea ce reprezintă o irosire de efort.

De exemplu, un tabel ar putea avea două câmpuri: externat viu (EXTV) și externat decedat (EXTD), fiecare luând valorile 'da' sau 'nu'. Însă de fiecare dată când EXTV=da, EXTD=nu (evident) și invers. În astfel de situații trebuie păstrate numai minimum de variabile, celelalte putând fi, la nevoie, calculate.