

## Sample rectifying by post-stratification and calibration

Lect. Ileana Gabriela NICULESCU-ARON  
Catedra de Statistică și Previziune Economică, ASE București

*Insuring the representativity of the sample is a main concern for the person who organizes a survey. The methods for rectifying the samples have constantly improved. Gradually, by using certain algorithms, on which efficient software programs are based, a transition from generalized post-stratification to generalized calibration was made.*

**Keywords:** sample, post-stratification, calibration.

**S**tabilirea planului de sondaj și a procedurilor de selecție a unităților în eșantion sunt foarte importante și se realizează cu multă rigurozitate.

Întocmirea chestionarului pe baza obiectivelor stabilite la începutul studiului presupune și definirea clară a variabilelor. În studiile realizate pe baza sondajelor statistice distingem următoarele tipuri de variabile:

- **variabile de interes.** Sunt cele care trebuie estimate și țin de obiectivele sondajului. În cazul anchetelor asupra forței de muncă variabile de interes pot fi de exemplu statutul după participarea la activitatea economică (ocupat, șomer, inactiv) sau durata efectivă a săptămânii de lucru.

- **variabile auxiliare sau de identificare.** Sunt cele pe care le cunoaștem din alte surse și nu formează obiectul studiului prezent dar sunt utilizate în stabilirea planului de sondaj și ulterior în ameliorarea extrapolării.

În ciuda eforturilor de a asigura reprezentativitatea eșantionului, în anchetele de mare amploare, realizate la nivel național, de cele mai multe ori estimatorii variabilelor auxiliare sunt deplasați față de parametrii populației. Datorită existenței corelației dintre aceste variabile și variabilele de interes putem presupune că și estimatorii acestora din urmă vor avea aceeași problemă.

Plecând de la rezultatele brute furnizate de respondenți se estimează cât mai bine posibil, prin extrapolare, situația populației. Această metodă constă în atribuirea unei anumite ponderi sau coeficient de extrapolare fiecărui respondent după cum el reprezintă o fracțiune din populație. Pentru a extrapola trebuie să dăm fiecărui respondent ponderea sa inițială

reprezentată de inversul probabilității cu care a fost selecționat ( $d_k = \frac{1}{\pi_k}$   $k=1, \dots, n_r$ ,  $n_r$

fiind numărul de respondenți). Acest lucru nu este posibil deoarece<sup>97</sup>:

- unitatea selectată este gospodăria, nimic nu ne garantează că un individ din 500 va face parte din eșantion;

- structura populației după variabilele auxiliare (vârstă, sex medii de rezidență etc.) se va regăsi în eșantion cu o oarecare aproximație;

- mai mult ca sigur că, datorită non-răspunsurilor totale, numărul respondenților va fi inferior efectivului eșantionului selectat inițial.

Datorită acestor cauze este necesară modificarea coeficienților de extrapolare inițiali având drept scop asigurarea preciziei și coerenței valorilor extrapolate.

Până nu demult redresarea eșantioanelor în vederea extrapolării se realiza prin post-stratificarea realizată pe baza variabilelor auxiliare.

Variabilele clasice utilizate de obicei în vederea post-stratificării eșantioanelor din sondajele realizate la nivel național sunt:

- **REG:** regiunea de dezvoltare în care locuiește respondentul. România este împărțită în 8 regiuni: Nord Est, Sud Est, Sud, Sud Vest, Vest, Nord Vest, Centru, București.

- **VÂRSTA:** grupa de vârstă căreia îi aparține respondentul. În ancheta asupra forței de muncă se utilizează următoarele 6 grupe: 15-24 ani, 25-34 ani, 35-44 ani, 45-54 ani, 55-64

<sup>97</sup> Luminet D., Vanderhoeft C., Une méthode de calibration appliquée aux statistiques de l'emploi, Carrefour de l'Economie 2003/7-8A, Bruxelles, pg. 4

ani și peste 65 ani.

o **MREZ**: mediul de rezidență Urban sau Rural.

o **SEX**: masculin sau feminin.

Pentru a realiza o post-stratificare după cele patru variabile auxiliare fiecare respondent va fi clasat într-una din cele  $8 \times 6 \times 2 \times 2 = 192$  celule numite post-straturi. În fiecare celulă  $h$  ( $h=1, \dots, 192$ ) vom avea un număr  $n_h$  de respondenți iar  $\sum_{h=1}^{192} n_h = n$  unde  $n$  este numărul

total de respondenți.

Pe baza informațiilor din alte surse putem determina distribuția pentru populația României cu vârsta 15 ani și peste ( $N$ ) după cele patru variabile: REG, VÂRSTA, MREZ, SEX reprezentată de frecvențele  $N_h$  unde  $\sum_{h=1}^{192} N_h = N$ .

Frecvențele relative ale eșantionului sunt date de raportul  $\frac{n_h}{n}$  iar ale populației de  $\frac{N_h}{N}$ .

Din diferite motive vor exista celule pentru care  $\frac{n_h}{n} \neq \frac{N_h}{N}$  eșantionul respondenților ne-reprezentând fidel populația României cu vârsta de 15 ani și peste, anumite celule (post-straturi) fiind sub sau supraestimate. „În sens strict matematic am putea spune că eșantionul nu este reprezentativ pentru populația din care a fost extras. Totuși noțiunea (ne)reprezentativ este în general utilizată într-un sens mai puțin strict. În acest caz se poate pretinde că eșantionul este reprezentativ pentru populație dacă toate frecvențele  $n_h$  sunt nenule și suficient de mari”<sup>98</sup>. În felul acesta fiecare combinație dintre regiune, vârstă, mediu de rezidență și sex este suficient reprezentată.

Existența unor diferențe semnificative statistice între frecvențele relative ale eșantionului de respondenți și populație poate avea consecințe grave asupra calității estimatorilor.

Am considerat  $y$  o variabilă de interes (de

exemplu durata săptămânii de lucru. Pe baza anchetei asupra forței de muncă urmează să se estimeze durata medie efectivă a săptămânii de lucru. Această variabilă nu este auxiliară și nu dispunem de informații în ceea ce o privește din alte surse. Calculând durata medie efectivă a săptămânii de lucru ( $\bar{y}$ ) ca o medie neponderată pe baza datelor din eșantion vom obține un estimator deplasat pentru media populației din motivele de sub sau supra-reprezentare prezentate anterior.

Estimatorul timpului total de muncă ( $N\bar{y}$ ) care va prezenta aceleași neajunsuri se poate scrie sub forma:  $N\bar{y} = N \frac{Y_T}{n}$  unde  $Y_T$  reprezintă timpul total de muncă din eșantion.

Din această relație se deduce că pentru a trece de la totalul eșantionului la totalul populației trebuie să atribuim fiecărui respondent același coeficient de ponderare  $\frac{N}{n}$ . Se poate spune că, în medie, fiecare respondent al anchetei reprezintă  $\frac{N}{n}$  din populația României cu vârste de 15 ani și peste.

Utilizând post-stratificarea această pondere nu se va efectua uniform. Fiecare respondent din post-stratul  $h$  va primi un coeficient de ponderare de forma  $\frac{N_h}{n_h}$  cu  $h=1, \dots, 192$ . În aceste condiții timpul total de muncă din populație  $Y_T$  va fi estimat pe baza relației:  $\hat{Y}_T = \sum_{h=1}^{192} N_h \bar{y}_h$ , iar durata medie efectivă a săptămânii de lucru se va estima pe baza relației:

$$\hat{\bar{y}} = \sum_{h=1}^{192} N_h \bar{y}_h / N.$$

Practic, fiecare respondent din stratul  $h$  reprezintă  $\frac{N_h}{n_h}$  indivizi ce îndeplinesc aceleași caracteristici.

Calitatea estimatorilor obținuți prin post-stratificare depinde de măsura în care variabilele auxiliare regiune, vârstă, mediu de rezidență și sex explică variația variabilelor de interes. Este clar că cele patru variabile sunt explicative pentru un număr foarte mare de variabile de interes totuși, în cazul în care se dorește o analiză detaliată sau estimarea unor

<sup>98</sup> Luminet D., Vanderhoeft C., Une méthode de calibration appliquée aux statistiques de l'emploi, Carrefour de l'Economie 2003/7-8A, Bruxelles, pg. 13.

variabile de interes specifice, ele nu mai sunt suficiente. Acest fapt a determinat abandonarea tehnicii clasice de post-stratificare. Modelul de post-stratificare prezentat anterior poate fi prezentat succint sub forma: REG x VÂRSTĂ x MREZ x SEX.

Pentru aceasta, pornind de la modelul inițial REG x VÂRSTĂ x MREZ x SEX putem adăuga alte variabile ajungând la un model detaliat de forma: REG x VÂRSTĂ x MREZ x SEX x X<sub>1</sub> x X<sub>2</sub>..... Noile variabile introduse vor permite ameliorarea estimațiilor. În acest model variabilele utilizate nu se vor mai numi variabile de post-stratificare ci **variabile de calibrare**.

În momentul în care numărul variabilelor de calibrare este foarte mare nu se vor putea calcula izolat ponderile pentru fiecare celulă rezultată din încrucișarea variabilelor. Literatura de specialitate propune diferiți algoritmi pentru determinarea unei soluții.

Într-un model de post-stratificare este posibil, ca pentru fiecare post-strat h să se definească o ecuație exprimată ca sumă a greutateilor (ce trebuie calculate) acordate respondenților din post-stratul h și a căror sumă trebuie să corespundă cu efectivul populației N<sub>h</sub>: astfel:

$$\sum_{k \in h} w_k = N_h \text{ unde } w_h \text{ este greutatea acordată respondentului } h.$$

Pentru fiecare post-strat stabilim o astfel de ecuație în final obținând un sistem de h ecuații liniare ce ar trebui rezolvat. Rezolvarea lui duce la mai multe soluții deoarece, pe de o parte acest sistem cuprinde mai multe necunoscute decât ecuații iar pe de altă parte fiecare individ k nu este cuprins într-o singură ecuație. O soluție particulară este obținută impunând ca toți indivizii din același post-strat h să aibă aceeași greutate w<sub>h</sub> astfel încât fiecare va fi tratat în aceeași manieră. Ecuația post-stratului h poate fi formulată astfel:

$$n_h w_k = N_h \text{ deci } w_k = N_h / n_h$$

Deoarece sistemul de ecuații de calibrare are mai multe soluții putem alegea cea soluție care modifică cel mai puțin **coeficienții de extrapolare inițiali**. Practic vom căuta acei coeficienți de redresare w<sub>k</sub> care să verifice sistemul de ecuații de calibrare și în același timp să fie cât mai apropiați posibil de coefi-

cienții de extrapolare ce rezultă direct din planul de eșantionare ( $d_k = 1/\pi_k$  k=1,...n<sub>r</sub>).

Problema generală a calibrării constă în ajustarea coeficienților de extrapolare inițiali  $d_k = 1/\pi_k$  și obținerea unor coeficienți de

redresare calibrați de forma :  $w_k = g_k d_k$  unde g<sub>k</sub> reprezintă factorul de ajustare.

Este vorba de următoarea problemă de optimizare:

$$\begin{cases} \sum_{k=1}^{n_r} d_k G\left(\frac{w_k}{d_k}\right) \text{ minim} \\ \sum_{k=1}^{n_r} w_k x_{kj} = T_j (j=1, \dots, m) \text{ restricțiile de calibrare} \end{cases}$$

unde:

○ m reprezintă numărul de variabile de calibrare iar x<sub>kj</sub> este valoarea variabilei de calibrare j pentru respondentul k

○ T<sub>j</sub> reprezintă totalul populației pentru variabila de calibrare j.

○ G este funcția distanțelor definită pe vecinătatea convexă a lui 1 și care verifică următoarele condiții:

$$G \geq 0;$$

G este strict convexă;

G este de două ori continuu derivabilă;

$$G(1) = 0$$

$$G'(1) = 0;$$

$$G''(1) = 1.$$

Teorema funcțiilor implicite afirmă că funcția reciprocă F a lui G', definită și continuu derivabilă în vecinătatea lui 0 satisface condițiile F(0)=1 și F'(0)=1. F se numește **funcție de calibrare**.<sup>99</sup>

Pentru simplificarea explicației presupun că problema de optimizare are doar o soluție și o singură restricție și anume:

$$\sum_{k=1}^{n_r} x_k w_k = T \text{ unde } x_k \text{ este variabila de cali-}$$

brare iar T reprezintă totalul populației pentru variabila de calibrare.

Pentru soluționarea problemei de minimizare ținând cont de restricție (extreme cu legături)

<sup>99</sup> Luminet D., L'enquête sur les Forces de travail: calibrage et autres développements, Statistics Belgium, Working Paper nr. 8, pg. 32

se recurge la funcția Lagrange definită astfel:

$$L(w_{k,\lambda}) = \sum_{k=1}^{n_r} d_k G\left(\frac{w_k}{d_k}\right) - \lambda \left( \sum_{k=1}^{n_r} x_k w_k - T \right)$$

unde  $\lambda$  este valoarea multiplicatorului Lagrange.

Pentru a afla valorile minime anulăm derivatele parțiale și obținem următorul sistem:

$$\begin{cases} \sum_{k=1}^{n_r} x_k w_k = T \\ G'\left(\frac{w_k}{d_k}\right) = \lambda x_k \Rightarrow w_k = d_k F(\lambda(x_k)) \end{cases}$$

Introducând  $w_k$  în prima relație obținem :

$$\sum_{k=1}^{n_r} x_k d_k F(\lambda(x_k)) = T$$

Pe baza acestei relații putem afla valoarea multiplicatorului Lagrange și ulterior valorile  $w_k$ .

Avem posibilitatea de a alege un criteriu de apropiere între coeficienții de extrapolare inițiali și coeficienții de redresare calibrați, deci de a alege o formă a funcției distanțelor  $G$  și implicit a funcției de calibrare  $F$ .

Calibrarea generalizată a devenit în momentul de față o metodă cunoscută și puternică în mediile de specialitate deoarece, utilizând informații auxiliare din diferite surse, reușește să îmbunătățească estimațiile sondajelor prin creșterea preciziei estimatorilor obținuți. Mai mult decât atât, calibrarea este utilizată pentru corectarea deplasării produse de non-răspunsurile totale.

Rezolvarea matematică a problemei de minimizare cu restricții prezentată necesită un volum mare de calcule cu un nivel ridicat de dificultate. Literatura de specialitate propune diferiți algoritmi pentru determinarea soluției ce stau la baza programelor software pe baza cărora se realizează calibrarea eșantioanelor. În continuare voi prezenta succint aceste programe:

**Generalized Estimation System (GES).** Este utilizat de Statistics Canada și este realizat sub programul SAS. Are la bază estimarea pe baza regresiei generalizate (GREG) descrisă de Stårndal, Swensson și Wretman în „Model Assisted Survey Sampling”, 1992. Această

metodă acoperă o clasă de estimatori calibrați care cuprinde cei mai utilizați estimatori. Cu toate acestea metoda de calibrare generalizată introdusă de Deville și Stårndal în 1992 este mai cuprinzătoare.

GES este utilizat împreună cu un alt soft realizat sub SAS, GSAM (Generalised Sampling System). Ambele programe acoperă diverse tipuri de sondaj simple sau complexe.

Avantajul GES este deci integrarea calibrării, a estimării pentru totaluri, medii, proporții, rapoarte și variația estimatorilor chiar dacă numai pe baza metodei GREG.

**BASCULA** a fost realizată în Delphi pentru Windows 95 de către Nieuwenbroek în 1997 și este utilizat în Olanda. Ca și GES se bazează pe metoda regresiei generalizate (GREG). Variația estimatorilor se bazează pe tehnica reeșantionării (balanced repeated sampling – BRR)

O particularitate a programului BASCULA este modul în care factorul de ajustare  $g_k$  este limitat în metoda liniară. Contrar procedurii utilizate de CALMAR și g-CALIB factorii de ajustare nu sunt trunchiați ci mai degrabă redimensionați printr-o procedură iterativă. Aceasta este considerată o tehnică de netezire limitată față de tehnica trunchierii iterative.

**CALMAR (Calage sur Marges).** Acest program a fost propus de membrii INSEE (Institut National de la Statistique des Etudes Economiques – Franța). CALMAR este realizat sub SAS și are la bază metoda de calibrare generalizată introdusă de Deville și Stårndal în 1992 concentrată pe calcularea coeficienților de calibrare  $w_h$  și a factorilor de ajustare  $g_h$ . Un instrument central al metodei îl constituie funcția distanțelor  $G$ . Din punct de vedere practic, CALMAR este considerat superior sistemelor prezentate anterior deoarece permite utilizatorilor să limiteze în mai multe moduri flexibilitatea coeficienților de redresare  $w_k$ . Din punct de vedere teoretic, deoarece se bazează pe metoda calibrării generalizate, sfera estimatorilor calibrați este mult mai cuprinzătoare.

Și în România Institutul Național de Statistică utilizează acest pachet program pentru calcularea coeficienților de ponderare în vederea creșterii gradului de precizie al estima-

țiilor și pentru tratarea non-răspunsurilor totale atât în Ancheta asupra forței de muncă AMIGO cât și în celelalte anchete în care unitatea de selecție este gospodăria sau întreprinderea.

**g-CALIB.** A fost introdus de către INS (Institut National de Statistique) – Belgia, având la bază pachetul statistic SPSS. Prima versiune a acestui program a fost realizată de Vanderhoeft și a devenit un instrument performant, aplicabil în situații diverse și capabil să rezolve probleme complicate de redresare a eșantioanelor. Ca și CALMAR, are la bază metoda de calibrare generalizată introdusă de Deville și Stärndal în 1992 concentrată pe calcularea coeficienților de calibrare  $w_h$  și a factorilor de ajustare  $g_h$ .

Pachetele program **g-CALIB și CALMAR** sunt comparabile din punct de vedere al fundamentării teoretice și al performanțelor. Cu toate acestea, din anumite puncte de vedere CALMAR este privit ca fiind în prezent cel mai bun soft în acest domeniu din următoarele motive:

- o Interfața este foarte prietenoasă nefiind necesar ca utilizatorul să fie un expert în metoda calibrării generalizate;
- o CALMAR este astfel conceput ca variabilele cantitative și calitative de calibrare să fie transformate automat într-o matrice proiectată de program. Acest lucru reduce foarte mult munca de pregătire a fișierelor input de către utilizator. Totuși, pentru calibrarea variabilelor cantitative utilizatorul CALMAR trebuie să realizeze o transformare a acestora astfel încât să obțină un format standard pentru fișierul de intrare.
- o Detectarea și raportarea erorilor se realizează într-un mod mai precis și mai eficient decât celelalte pachete program similare.

## Concluzii

Pentru ca estimatorii variabilelor de interes obținuți în urma prelucrării datelor din eșantion să nu fie deplasați față de parametrii populației este recomandat să fie aplicată una din metodele de ameliorare a eșantioanelor prezentate în această lucrare..

Deși calibrarea generalizată este în momentul de față o metodă puternică pentru că reușește să îmbunătățească estimațiile sondajelor prin creșterea preciziei estimatorilor obținuți corectând într-o oarecare măsură și deplasările produse de non-răspunsurile totale, ea nu este suficient cunoscută în România. Unul dintre motive este și costul foarte ridicat al programelor software pe baza cărora se poate implementa această metodă. Consider că orice institut ce are pretenția realizării unui eșantion reprezentativ la nivel național ar trebui să poată realiza calibrarea acestuia.

## Bibliografie

- Luminet D., Vanderhoeft C., Une méthode de calibrage appliquée aux statistiques de l'emploi, Carrefour de l'Economie 2003/7-8A, Bruxelles
- Luminet D., Vanderhoeft C., Une méthode de calibrage appliquée aux statistiques de l'emploi, Carrefour de l'Economie 2003/7-8A, Bruxelles
- Droesbeke J-J., Théorie et Pratiques de l'échantillonnage, Formation des statisticiens européens, Eurostat, Support de cours