



Institutul Național de Statistică

**GHID METODOLOGIC
DE EDITARE A DATELOR**

2016

CUPRINS

I. OBIECTIVELE PRINCIPALE	3
II. DESCRIEREA ETAPELOR PROCESULUI DE EDITARE A DATELOR	4
II.1. Managementul/Organizarea cercetării statistice	4
II.2. Introducerea datelor în mediul electronic	13
II.3. Validarea datelor	16
II.4. Imputarea și ajustarea datelor	25
II.4.1. Imputarea non-răspunsurilor	25
II.4.2. Ajustarea valorilor imputate	30
III. BIBLIOGRAFIE	32
IV. ANEXE	33
Anexa 1. CORELATII PENTRU CONTROLUL DE INTEGRITATE – Ancheta forței de muncă în gospodării – 2016	34
Anexa 2. CORELAȚII DE VERIFICARE A CORECTITUDINII “DRUMURILOR” LOGICE DIN CHESTIONARUL CI – Ancheta forței de muncă în gospodării – 2016	35
Anexa 3. IMPUTAREA NON-RASPUNSURILOR pentru variabila VNET - Ancheta forței de muncă în gospodării – 2015	38
Anexa 4. Imputarea non-răspunsurilor în EU-SILC	44
Anexa 5. DESCRIEREA UNEI CERCETĂRI STATISTICE SELECTIVE - Ancheta forței de muncă în gospodării	46

I. OBIECTIVELE PRINCIPALE

În prezent, nevoia de date statistice de înaltă calitate, cu un grad ridicat de detaliere și produse într-un interval scurt de timp este în creștere, iar cererea de date cuprinde un spectru tot mai divers de domenii ale activității sociale și economice. Conform legii 226/2009 privind organizarea și funcționarea statisticii oficiale în România, Institutul Național de Statistică, împreună cu direcțiile sale teritoriale, este principala instituție a sistemului statistic național care asigură date statistice actuale, relevante și de calitate, prin programul statistic național. Principalele surse pe baza cărora se realizează colectarea datelor sunt cercetări statistice selective (prin sondaj) și exhaustive (de tipul recensămintelor). Un factor major care perturbă calitatea datelor colectate este faptul că sursele de date conțin, în mod inevitabil, erori și valori lipsă, impunând derularea unui proces continuu de editare a datelor.

Viitorul cercetărilor statistice, care prevede utilizarea extensivă a surselor administrative și a unor surse combinate de date pentru diverse cercetări ridică și mai multe probleme în ceea ce privește utilizarea unor date corecte, coerente și cât mai aproape de realitatea de zi cu zi. În aceste condiții, editarea datelor capătă valențe multiple.

Unul dintre obiectivele principale ale ghidului metodologic este de a promova o mai bună înțelegere a diferitelor metode de editare a datelor statistice rezultate din cercetările statistice, în ceea ce privește controlul, validarea și imputarea datelor produse în cadrul Institutului Național de Statistică din România.

Ghidul își propune să reprezinte și un util instrument de perfecționare a statisticienilor, prin prezentarea unor aspecte esențiale ale statisticii aplicate, a unor exemple din cercetările statistice implementate de Institutul Național de Statistică, alături de aspecte de ordin teoretic.

Prezentul ghid poate servi și pentru o viitoare standardizare a proceselor statistice în diferite domenii de interes (cercetări în întreprinderi, în gospodării etc.).

De asemenea, ghidul prezintă o scurtă trecere în revistă a metodologiilor utilizate în procesul de editare a datelor - cel mai frecvent utilizate - în alte oficii de statistică naționale din statele membre ale Uniunii Europene și din alte state ale lumii (de ex. SUA, Canada).

❖ CE ESTE EDITAREA DATELOR ȘI DE CE ESTE NECESARĂ?

Ce este editarea datelor și de ce este necesară?

Editarea datelor reprezintă ansamblul operațiunilor aplicate datelor brute (colectate în cadrul unei cercetări statistice) pentru a îmbunătăți calitatea rezultatelor (eliminarea erorilor), prin analiza acestora din prisma unui set de criterii sau judecăți de valoare.

Scopul validării/ editării este asigurarea unui nivel acceptabil de calitate a datelor finale (diseminate). Validarea/ editarea nu pot asigura, în sine, calitatea datelor ci doar îndeplinirea unui set de criterii de calitate în funcție de care datele pot fi considerate ca fiind „acceptabile”. Se poate spune deci că **validarea nu va avea ca rezultat date „perfecte” sau „adevărate” ci date plauzibile.**

Dimensiunile calității avute în vedere de procesul de validare se referă la:

- Acuratețea – se referă la diferența dintre valoarea adevărată și cea estimată a unui parametru ce trebuie măsurat. Această diferență este dată atât de erorile de sondaj și de alte erori (de măsurare, de prelucrare etc.). Erorile de sondaj sunt măsurate prin metode statistice specifice și nu pot fi corectate de procedurile de validare / editare. Alte tipuri de erori însă, ca de exemplu erorile de măsurare sau cele generate chiar de procesul de producție statistică (codificare, prelucrare etc.) pot fi identificate și corectate prin proceduri de validare / editare.
- Coerența și comparabilitatea – se referă la gradul în care datele sunt consistente ca structură, ca evoluție în timp și comparabile între zone geografice.

II. DESCRIEREA ETAPELOR PROCESULUI DE EDITARE A DATELOR

Editarea datelor este definită ca procesul care implică controlul și validarea datelor colectate pe baza cercetărilor statistice selective (anchete statistice) sau exhaustive cu scopul de a îmbunătăți calitatea rezultatelor. Acest proces se desfășoară în patru etape:

- Managementul/Organizarea cercetării statistice;
- Introducerea datelor în mediul electronic;
- Validarea datelor;
- Ajustarea datelor.

II.1. Managementul/Organizarea cercetării statistice

Managementul cercetării statistice este o activitate care se desfășoară, în diferite forme și cu diferite intensități, în fiecare etapă a acesteia.

Unele dintre activitățile de management premergătoare colectării propriu-zise a datelor constau în:

- asigurarea cadrului legal al cercetării statistice;
- verificarea oportunității și necesității desfășurării cercetării respective, precum și a eficienței (a concordanței dintre costuri și rezultate);
- asigurarea resurselor (umane, financiare și de timp) necesare desfășurării ei;
- asigurarea sustenabilității în timp a cercetării respective;
- verificarea faptului că instrumentarul proiectat permite colectarea informațiilor necesare producerii statisticilor necesare și astfel, că cercetarea își va atinge obiectivele generale și specifice etc.

Ulterior colectării datelor, managementul cercetării statistice constă în următoarele activități:

a) **verificarea integrității**¹ și

b) **controlul calității.**

¹ Un exemplu de corelații pentru controlul de integritate este prezentat în Anexa 1

Pentru cercetările statistice selective, **verificarea integrității** presupune condiții de control care verifică gradul de colectare a datelor cercetării, mai exact, dacă s-a colectat ceea ce (și cât) trebuia colectat. Verificarea integrității are o dimensiune cantitativă prin care se verifică că numărul de chestionare completate, respectiv cele aferente unităților pentru care colectarea nu s-a realizat din diverse motive, este egal cu numărul unităților selectate în eșantionul cercetării.

Dimensiunea calitativă urmărește dacă colectarea datelor s-a realizat doar de la unitățile incluse în eșantion și nu de la alte unități (în cazul în care cercetarea selectivă nu permite înlocuirea unităților nerespondente).

Un alt aspect al verificării integrității privește integritatea între chestionare. De exemplu, într-o cercetare statistică în gospodării, pentru fiecare locuință ocupată trebuie să existe cel puțin un chestionar de locuință (care are completat cel puțin rubrica privind Rezultatul interviului) completat, iar dacă există un chestionar de locuință pentru care s-au colectat datele (gospodăria a acceptat interviul), numărul de chestionare individuale completate (care au completat cel puțin rubrica privind Rezultatul interviului) trebuie să fie egal cu numărul membrilor de gospodărie prezenți sau temporar absenți din respectiva gospodărie.

Un exemplu privind verificarea integrității chestionarelor statistice este prezentat în Anexa 2.

Chestionarele trebuie să conțină o rubrică referitoare la Rezultatul interviului (sau a colectării datelor, dacă se utilizează altă metodă de colectare decât cea a interviului) prin care să se poată identifica dacă colectarea datelor a avut loc, dacă respondentul a fost inaccesibil sau a refuzat răspunsul; aceste informații potând fi utilizate în cadrul procedurilor de verificare.

Organizarea cercetării statistice include controlul calității procesului de colectare a datelor și măsuri ale impactului activității de ajustare asupra rezultatelor. Este un pas în procesul de control al calității totale, care asigură faptul că ipotezele statistice de bază ale unei cercetări statistice selective nu sunt încălcate.

❖ ETAPELE UNEI CERCETĂRI STATISTICE

Etapele și activitățile principale care trebuie desfășurate în fiecare etapă pentru realizarea unei cercetări statistice prin sondaj sunt următoarele:

I) PROIECTAREA CERCETĂRII STATISTICE are în vedere următoarele aspecte:

- a. **Definirea obiectivelor cercetării** statistice se realizează pornind de la rezultatele ce se doresc a se obține. În majoritatea cazurilor, o cercetare statistică este inițiată datorită apariției unei nevoi de informație statistică, a unui indicator sau pentru a monitoriza un anumit fenomen economic sau social;
- b. **Identificarea bazei de sondaj** celei mai adaptate la situație ținând seama de aspectele de calitate și cost. Se pot utiliza eventual mai multe baze de sondaj de proveniențe diferite pentru a se realiza baza considerată a fi cea mai adecvată;
- c. **Selecția eșantionului** presupune extragerea unui număr de unități dintr-o bază de sondaj, pe baza unor criterii pre-stabilite, de la care urmează a fi colectate datele. Volumul eșantionului se calculează în funcție de mai multe criterii, cel mai important fiind acela a gradului de reprezentativitate pe care dorim să-l aibă rezultatele finale ale cercetării statistice;;
- d. **Proiectarea chestionarului:** pornind de la obiectivele cercetării statistice, de la programul său de observare, se stabilesc variabilele ce urmează a fi colectate și, pe baza acestora,

se stabilește conținutul chestionarului, prin formularea întrebărilor, stabilirea fluxului acestora, a pertinentei lor, a duratei medii de completare etc;

- e. **Proiectarea mijloacelor de prezentare a rezultatelor cercetării statistice:** tabele de prezentare a rezultatelor, fișiere de micro-date sau baze de date on-line.

II) DESFĂȘURAREA CERCETĂRII STATISTICE ÎN TEREN are ca etape principale:

a. **Pregătirea activităților în teren** - sunt necesare activități premergătoare cum ar fi:

- **selectarea operatorilor de interviu** și a altor categorii de personal se face pe baza unor criterii precum nivelul de instruire, experiența anterioară în cercetări prin sondaj, abilitatea de a desfășura un interviu, aspectul fizic etc.;
- **instruirea personalului** în vederea cunoașterii obiectivelor cercetării, a metodelor celor mai adecvate de contactare a unităților din eșantion, a modului de desfășurare a interviurilor și de înregistrare a informațiilor în chestionare;
- **stabilirea celor mai bune practici în colectarea datelor** poate presupune de exemplu, găsirea metodelor de evitare, iar dacă acest lucru nu este posibil, de minimizare a numărului unităților care nu au participat la cercetare, pentru limitarea la maxim a **non-răspunsurilor**; de cele mai multe ori operatorii de interviu trebuie să facă vizite repetate până reușesc să contacteze unitățile pe care nu le-au găsit în prealabil sau care au refuzat inițial să participe la cercetare; dacă aceștia nu reușesc să convingă unitatea, o nouă încercare trebuie făcută de către supervisorul acestuia;
- **atribuirea eșantionului fiecărui operator de interviu** presupune identificarea unităților care trebuie intervievate de către aceștia;
- **identificarea prealabilă a unităților din eșantion** se realizează pentru a face publică în teren desfășurarea cercetării, în vederea economisirii timpului în perioada de colectare a datelor și pentru minimizarea non-răspunsurilor;
- **distribuirea materialelor cercetării către operatorii de interviu;**

b. **Colectarea informațiilor** – este etapa de care depinde în mod covârșitor calitatea rezultatelor unei cercetări statistice.

Colectarea se poate realiza prin interviuri față-în-față (metodă des utilizată în cercetările selective în gospodării din România), în care înregistrarea se realizează mai întâi prin completarea răspunsurilor în chestionarele tipărite pe hârtie (ulterior acestea fiind introduse într-o bază de date) sau prin intermediul mediului online, prin chestionare electronice (care se pot descărca automat în bazele de date). Sursa datelor: persoane fizice și/sau juridice sau surse administrative disponibile.

Pentru cercetarea statistică Intrastat, colectarea datelor de la operatorii economici se face numai în format electronic, în sistem online, pe site-ul dedicat Intrastat, sau în sistem offline, prin aplicația Intrastat pusă la dispoziție cu titlu gratuit de INS operatorilor economici.

Alte metode de colectare a datelor pot fi:

- auto-înregistrarea, în care respondentul completează el însuși răspunsurile într-un chestionar;
- prin poștă, respondentul primind și, ulterior completării, returnând chestionarele către organizația care a organizat cercetarea prin poștă;
- prin telefon, interviul desfășurându-se sub forma unei convorbiri telefonice între un operator de interviu și respondent, iar răspunsurile fiind completate fie pe un chestionar pe hârtie, fie direct, într-o bază de date;

- metode mixte, în care se combină mai multe dintre metodele clasice de colectare.

III) PRELUCRAREA DATELOR

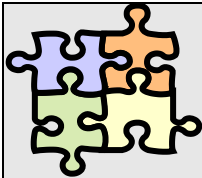
- a. Principala activitate o reprezintă **introducerea datelor din chestionare pe suport electronic și/sau preluarea automată a datelor**.
- b. **Controlul și validarea datelor** presupune realizarea unui control de integritate (daca s-au introdus toate chestionarele si toate informațiile conținute de acestea) precum și a unui control logic al datelor (controlul coerenței datelor) și corectarea erorilor constatate;
- c. În această etapă de prelucrare a datelor are loc **tratarea non-răspunsurilor**, folosind diverse metode de imputare a datelor lipsă.

IV) CALCULUL COEFICIENȚILOR DE EXTINDERE ȘI AL ESTIMATORILOR

Coeficienții de extindere sunt necesari pentru a extrapola datele obținute pe baza unui eșantion la nivelul întregii populații. Procesul de ponderare începe cu calculul ponderii de bază care, pentru o anumită unitate din eșantion, este egală cu inversul probabilității de selecție. Ulterior, aceste ponderi de bază sunt ajustate pentru compensarea non-răspunsurilor și a non-acoperirii și pentru a face estimățiile din eșantion conforme cu totalurile din populația de referință.

Non-răspunsurile totale se pot compensa prin 3 metode, și anume:

- prin ajustarea ponderilor de bază;
- prin selectarea inițială a unui eșantion mai mare care, în urma existenței non-răspunsurilor va determina atingerea unui eșantion realizat de mărimea dorită;
- prin utilizarea substituirii, adică prin înlocuirea unor unități care nu răspund cu altele care nu au fost incluse în eșantion și care sunt similare cu cele non-respondente din punctul de vedere al caracteristicilor de interes.



Caseta 1: Exemplificare din **STATISTICA SOCIALĂ**

În domeniul statisticii sociale extinderea rezultatelor obținute din anchetă se realizează pe baza coeficienților atribuiți persoanelor din gospodăriile din eșantion, care au răspuns la interviu. Pentru determinarea acestor coeficienți este necesară parcurgerea următoarelor etape:

(i) Calculul ponderilor de bază

Probabilitățile de includere ale UP, corespunzătoare primei trepte de eșantionare, au fost calculate conform unei scheme de selecție stratificată. Criteriile de stratificare utilizate au fost județ și mediu de rezidență, prin intersectarea acestora rezultând un număr de 88 straturi (în Mun. București selecția a fost realizată în mod separat pentru fiecare din cele 6 sectoare administrative). În fiecare din cele 88 de straturi au fost calculate probabilitățile de includere în prima treaptă, proporțional cu mărimea unei UP, mărime exprimată în număr de locuințe permanente, conform următoarei formule de calcul:

$$P_{1hj} = m_h \times \frac{N_{hj}}{\sum_{j=1}^{N_h} N_{hj}}$$

unde:

h = indicele stratului, h=1,...,88

j = indicele UP, j=1,...,4622

P_{1hj} = probabilitatea de includere în prima treaptă a UP j aparținând stratului h

m_h = volumul eșantionului de UP în stratul h

N_{hj} = Număr de locuințe permanente în stratul h, UP j

N_h = Număr de locuințe permanente în stratul h

În a doua treaptă, în interiorul fiecărei UP din totalul celor 792 UP incluse în prima treaptă în EMZOT'2002, au fost extrase câte 36 de locuințe pe baza unui algoritm de selecție sistematică cu start aleator. Astfel, toate locuințele compunând o anumită UP au aceeași probabilitate de includere în eșantionul trimestrial al anchetei. Probabilitatea de includere din treapta a doua a fost calculată după cum urmează:

$$P_{2hjk} = \frac{36}{N_j}$$

unde:

h = indicele stratului, h=1,...,88

j = indicele UP, j=1,...,792

k = indicele unei locuințe aparținând UP j

P_{2hjk} = probabilitatea de includere în treapta a doua a unei locuințe permanente k aparținând UP j din stratul h

N_j = numărul total de locuințe permanente în UP j

Probabilitatea generală de includere a unei locuințe k în eșantionul anchetei (PGSk), după cele 2 trepte de eșantionare, este calculată astfel:

$$PGSk = P_{1hj} * P_{2hjk}$$

Ponderea de bază a unei locuințe k, selectată în a doua treaptă de eșantionare din UP j (BWk) este, așadar, inversa probabilității generale de includere a unei locuințe k:

$$BWk = 1/PGSk$$

Ponderea de bază a unei locuințe este 'împrumutată' tuturor gospodăriilor din acea locuință.

(ii) Ajustarea non-răspunsurilor totale

Pentru a acoperi procentul gospodăriilor care refuză să participe la anchetă, se procedează la re-ponderarea unităților respondente, prin ajustarea cu inversul ratei de răspuns.

Experiența anterioară ne-a dovedit că două variabile pot influența decizia unei gospodării de a participa sau nu la anchetă:

- Județ;

- Mediul de rezidență (urban sau rural).

Ca urmare, tratarea non-răspunsurilor totale nu se face în mod global, pe ansamblul eșantionului, ci în mod diferențiat, pe grupe de gospodării, generate de intersecția variabilelor considerate ca variabile explicative ale non-răspunsului: județ*mediu de rezidență.

Această modalitate de tratare a non-răspunsurilor corespunde așa numitei metode a grupelor de răspuns omogen, care presupune că într-o anumită grupă din eșantion probabilitățile de răspuns sunt egale. În concluzie, pentru fiecare grupă de gospodării, obținută prin intersectarea variabilelor sus-menționate și considerată ca grupă de răspuns omogen, se calculează rata răspunsului, după cum urmează:

$$R_g = \frac{NHH_{2g}}{NHH_{1g}}$$

unde:

g = indicele grupei de răspuns omogen, g=1,..., număr de grupe generate de intersecția variabilelor județ*mediu de rezidență

NHH_{1g} = suma ponderilor de bază ale gospodăriilor eșantionate în grupa g, gospodării aparținând locuințelor eligibile pentru interviu.

NHH_{2g} = suma ponderilor de bază ale gospodăriilor respondente în grupa g.

În final, ponderea de bază a unei gospodării respondente k este ajustată cu inversul ratei de răspuns, separat pe fiecare grupă de răspuns omogen:

$$W_{adj_k} = BW_k \cdot (1/R_g)$$

Fiecare individ aparținând unei gospodării respondente primește ponderea de bază ajustată a gospodăriei.

(iii) Redresarea eșantionului și calculul ponderilor finale

Redresarea este realizată cu scopul de a îmbunătăți calitatea estimațiilor printr-o ajustare finală a ponderilor rezultate în urma pasului anterior.

Ponderile obținute în final sunt modificate astfel încât totalurile estimate din eșantion să fie egale cu totalurile în populație pentru anumite variabile. În plus, ponderile finale sunt obținute astfel încât să se îndepărteze cât mai puțin posibil de ponderile inițiale, prin minimizarea unei funcții de distanță dintre cele două ponderi, ceea ce are efect asupra îmbunătățirii preciziei estimațiilor.

Această metodă de redresare este cunoscută sub numele de calibrare, în timp ce variabilele utilizate sunt denumite variabile de calaj. Calibrarea se realizează cu ajutorul macro-ului SAS CALMAR (CALaj pe MARje), creat de către INSEE Franța.

CALMAR calculează ponderea finală utilizând o variabilă auxiliară x (pentru simplificare, presupunem aici o singură variabilă auxiliară disponibilă), ale cărei totaluri sunt cunoscute pentru întreaga populație, și variabila inițială de ponderare, astfel încât:

$$\sum_{k \in S} G(W_{final_k} / W_{adj_k}) = \min$$

sub restricția:

$$\sum_{k \in s} W_{final_k} \cdot x_k = X$$

unde:

- k = indicele unei gospodării din eşantionul disponibil s
- G = este o funcție de distanță de argument
- X = este totalul variabilei auxiliare în populație.

Spre deosebire de alte metode de redresare (de exemplu post-stratificare sau estimatorul prin raport), metodele de calibrare implementate în CALMAR urmăresc, pe lângă un calaj al ponderilor finale pe totalurile populației, și o minimizare a distanței dintre ponderile inițiale și ponderile finale.

Variabile demografice (populația pe sexe și grupe de vârstă) și variabile de localizare (populația pe regiuni și medii de rezidență) sunt utilizate în ajustarea finală. Structura populației pe variabilele menționate este cunoscută din surse externe (populația rezidentă, disponibilă de două ori pe an: 1 ianuarie și 1 iulie al anului respectiv)

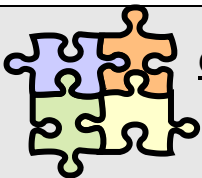
Pentru o regiune se folosesc următoarele variabile:

- Număr de persoane pe medii de rezidență (urban, rural)
- Număr de persoane pe sexe (masculin, feminin)
- Număr de persoane pe categorii de vârstă

După calibrare, totalurile estimate sunt egale cu totalurile în populație pentru fiecare din celulele obținute prin intersecția variabilelor menționate mai sus.

Uneori, destul de rar, totuși, atunci când anumite celule conțin prea puține observații în eşantion, calajul este dificil sau chiar imposibil de realizat. În aceste cazuri, se recurge la o regroupare a celulelor inițiale.

La sfârșitul acestei etape, ponderile finale, care sunt diferite de la o gospodărie la alta, sunt obținute, iar toate persoanele aparținând unei anumite gospodării primesc ponderea finală a gospodăriei.



Caseta 2: Exemplificare din **STATISTICA ÎNTREPRINDERILOR**

În cazul anchetei structurale în întreprinderi extinderea rezultatelor obținute din anchetă se realizează pe baza coeficienților atribuiți fiecărei întreprinderi din eşantion, care a răspuns la chestionar. Comparativ cu ancheta structurală în întreprinderi, în cazul altor anchete în întreprinderi coeficienții de extindere se calculează din ponderea de selecție și rata de non-răspuns. Pentru determinarea acestor coeficienți este necesară parcurgerea următoarelor etape:

- Calculul ponderii de selecție (π_{ih}) pentru fiecare unitate. Ponderea este de tipul Horvitz-Thompson calculată ca inversul probabilității de selecție.

$$\pi_{ih} = \frac{1}{p_{ih}} = \frac{N_h}{n_h}$$

unde:

p_{ih} = probabilitatea de selecție a unității i din stratul h

N_h = numărul de unități din baza de sondaj, în stratul h

n_h = numărul de unități în eșantion în stratul h

- Calculul ratei de non-răspuns. Rata de non-răspuns la nivel de strat se calculează pornind de la ipoteza conform căreia unitățile care nu au răspuns sunt similare din punct de vedere statistic cu cele care au răspuns. Se obțin coeficienții finali de extindere la nivel de strat prin raportul N_h/m_h unde N_h este numărul total de întreprinderi în stratul h al populației și m_h este numărul de întreprinderi cu date (care au răspuns) din stratul h al eșantionului.

$$c_h = \frac{n_h}{m_h}$$

unde:

n_h = numărul de unități din eșantion din stratul h

m_h = numărul unități respondente

- Urmează etapa de post-stratificare și estimare.
- În etapa de estimare se folosesc informații auxiliare din surse administrative (situații financiare etc). Indicatorii preluați din baza de calare sunt cifra de afaceri și numărul de salariați la nivel de unitate. Prin însumarea indicatorilor de interes la nivel de clasa CAEN Rev2 și clasă de mărime a întreprinderii (stabilită după numărul de salariați), din baza de calare, se obțin fișierele de calaj. Cu ajutorul pachetului software CLAN (SAS macro), se calculează coeficienții de calaj (cc_{ih}). Calcularea acestor coeficienți la nivel de unitate se face prin ajustarea coeficienților finali de extindere, ținând cont de limitele (marjele) din fișierul de calaj. Aplicând coeficienții de calaj în locul coeficienților de extindere finali, diferența dintre suma estimatelor pentru cifra de afaceri/numărul mediu de salariați, la nivel de clasa CAEN și clasă de mărime și valoarea „marginii” corespunzătoare din fișierul de calaj va tinde către 0.

Calculul coeficienților finali ($COEF_{ih}^{ext}$)

$$COEF_{ih}^{ext} = \pi_{ih} \cdot c_h \cdot CC_{ih}$$

- Calculul coeficienților de redresare (ch). Coeficienții de redresare s-au calculat la nivelul fiecărei celule de selecție a unităților primare ca inversul probabilității de răspuns. Coeficientul de redresare are rolul de compensare a unităților non-respondente în ipoteza în care aceste unități se manifestă similar cu unitățile respondente aferente stratului din care fac parte

$$c_h = \frac{1}{p_{rh}} = \frac{1}{\frac{m_h}{n_h}} = \frac{n_h}{m_h}$$

unde:

p_h = probabilitatea de răspuns din celula de selecție h

n_h = numărul de unități primare selectate în eșantion din celula de selecție h

m_h = numărul de unități primare selectate în eșantion din celula de selecție h care au răspuns la anchetă.

- Calculul coeficientului de extindere a unităților primare ($COEF_{ih}^{ext}$):

$$COEF_{ih}^{ext} = \pi_{ih} \cdot c_h ==$$

Calculul indicatorilor de calitate

Coeficientul de variație (CV) este definit ca eroarea standard ($V(\theta)$) împărțită la valoarea calculată a estimatorului (θ). CV-ul reprezintă eroarea standard în termeni relativi (procente) și cuantifică eroarea de eșantionare.

Coeficienții de variație sunt calculați la nivel de celulă utilizând procedura SAS - PROC SURVEYMEANS

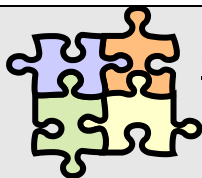
$$SE = \sqrt{\sum (\theta - \hat{\theta})^2} = \sqrt{VAR(\hat{\theta})} \quad ; \quad CV^{EST}(\theta) = \frac{\sqrt{V(\theta)}}{\theta}$$

Estimatorul trebuie să ia în considerare planul de eșantionare și trebuie să integreze efectul preciziei ajustărilor cu non-răspuns-ul, corecțiile clasificărilor eronate și informații auxiliare prin metode de calibrare etc.

V) ANALIZA DATELOR

Analiza datelor se realizează prin calculul unor indicatori medii (de exemplu, consumul mediu pe persoană dintr-un anumit produs, venitul mediu al gospodăriilor, câștigul salarial mediu pe economie etc.) totaluri de variabile repartizate în timp sau spațiu sau prin prezentarea distribuției unor variabile.

Un alt aspect care trebuie detaliat și documentat în această etapă îl reprezintă **estimarea calității rezultatelor** obținute prin calculul intervalelor de încredere și a coeficienților de variație pentru cei mai importanți indicatori calculați pe baza datelor cercetării.



Caseta 3: Exemplificare din **STATISTICA AGRICOLĂ**

În domeniul statisticii agricole, pentru Ancheta Structurală în Agricultură, estimarea calității datelor s-a realizat prin calcularea Erorii Relative Standard (ESR).

Eroarea relativă standard (RSE) s-a calculat astfel:

$$\frac{\sqrt{\hat{V}}}{\hat{Y}} \times 100,$$

Unde:

\hat{Y} =- este estimarea variabilei și

\hat{V} = este varianța estimatorului.

Întrucât, în România, Ancheta Structurală în Agricultură 2013 (ASA 2013) a fost o anchetă selectivă, cu eșantionare stratificată, metoda de estimare a varianței a ținut cont de acest lucru. Estimatorul total al lui Y a fost calculat după formula:

$$\hat{Y} = \sum_{h=1}^H \sum_{i=1}^{n_h} W_{hi} \times Y_{hi} = \sum_{h=1}^H N_h \times \bar{Y}_h,$$

Unde:

H = număr de straturi

N_h = număr (extins) al exploatațiilor agricole din stratul h

n_h = numărul exploatațiilor agricole din stratul h

Y_{hi} = valoarea variabilei Y pentru exploatația agricolă i din stratul h

$W_{hi} = \frac{N_h}{n_h}$ = factorul de extindere pentru exploatația agricolă i din stratul h

\bar{Y}_h = media valorilor Y_{hi} , $i = 1, \dots, n_h$

Estimatorul varianței lui \hat{Y} s-a calculat după formula:

$$\hat{V}(\hat{Y}) = \sum_{h=1}^H \frac{N_h \times (N_h - n_h)}{n_h \times (n_h - 1)} \times \left(\sum_{i=1}^{n_h} Y_{hi}^2 - \frac{Y_h^2}{n_h} \right)$$

VI) PUBLICAREA REZULTATELOR ȘI EVALUAREA FINALĂ

Scopul final al unei cercetări statistice este *oferirea de informații corecte și în timp util care să permită cunoașterea unui fenomen, a unui efect al unei cauze sau a unei categorii de populație*. Evaluarea finală a unei cercetări prin sondaj trebuie să evidențieze dacă aceasta și-a atins obiectivele. Se poate măsura prin gradul de satisfacție al utilizatorilor rezultatelor unei cercetări, prin numărul de publicații vândute, numărul de accesări on-line al bazelor de date care conțin estimări obținute prin sondaj etc.

II.2. Introducerea datelor în mediul electronic

Introducerea datelor într-un mediu electronic reprezintă, în cele mai multe cazuri, conversia datelor din varianta "pe hârtie" în variantă electronică. Această activitate este însoțită de mesaje de eroare care apar în cazul introducerii unor valori eronate (fie erori de înregistrare, fie erori de introducere).

Cheile de control în ceea ce privește introducerea datelor pot fi stabilite în două tipuri de abordări: **top down** sau **bottom up**.

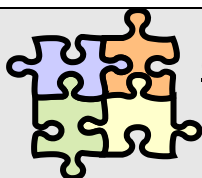
- Modalitatea **Top down** se referă la introducerea datelor fără detectarea erorilor care pot apare în momentul introducerii. Se utilizează personal care introduce datele cu viteză mare, în mod "heads down". Datele introduse în mod „heads down” sunt adesea re-verificate prin

reintroducerea chestionarului și compararea celor două variante (copii) introduse ale aceluiași chestionar.

- b. Modalitatea **Bottom up** se referă la introducerea datelor concomitent cu verificarea lor, în momentul introducerii. Modalitatea „Heads up” de introducere a datelor necesită ca personalul care introduce astfel datele să aibă cunoștințe în domeniul statistic respectiv. Introducerea datelor este mai lentă, dar revizuirea/ajustarea datelor este mai redusă ca volum deoarece inconsistențele simple între răspunsuri sunt identificate mai devreme/în faza de început a anchetei. Acest mod este eficient, în special atunci când intervievatorul sau respondentul introduce datele în timpul interviului. Acest lucru este cunoscut sub numele de **Computer Assisted Interviewing** (Intervievare asistată de calculator), care este explicată în detaliu, mai jos.

Datele pot fi capturate (introduse) prin mai multe metode automate fără a folosi introducerea tradițională a datelor. Pe măsură ce tehnologia avansează, mai multe instrumente vor deveni disponibile pentru capturarea (introducerea) datelor. Un instrument popular este interviul prin telefon cu sintetizarea vocii asistată de calculator („the touch-tone telephone key-pad with synthesized voice computer-administered interview”). Citoarele optice de caractere (OCR) pot fi folosite pentru a scana chestionarele în format electronic.

Alegerea modului de introducere a datelor precum și a metodei de ajustare a datelor au cel mai mare impact asupra tipului de personal care va fi necesar și asupra instruirii acestuia.



Caseta 4: Exemplificare din **STATISTICA ÎNTREPRINDERILOR**

Pentru cercetarea statistică Intrastat, colectarea datelor de la operatorii economici se face numai în format electronic, în sistem online, pe site-ul dedicat Intrastat, sau în sistem offline, prin aplicația Intrastat pusă la dispoziție cu titlu gratuit de INS operatorilor economici.

Există două tipuri de control al datelor colectate prin sistemul Intrastat: validarea primară și validarea secundară.

Validarea primară reprezintă validarea datelor la nivelul operatorului economici precum și validarea la nivelul INS din punctul de vedere al corectitudinii codurilor utilizate. Validarea secundară se referă atât la corectitudinea corelațiilor dintre variabile cât și a indicatorilor tip medie (ex. preț mediu).

Validarea primară:

I. La nivelul operatorului economic, în momentul introducerii datelor în declarația statistică Intrastat, sunt verificate următoarele:

- i. codurile declarate să fie valide, conform nomenclatoarelor încorporate în aplicațiile Intrastat offline și online: codurile de produs din nomenclatorul combinat, codurile de țări (de destinație / expediție sau origine), termenii de livrare, modul de transport și natura tranzacției;
- ii. masa netă, unitățile de măsură suplimentare, valoarea facturată și valoarea statistică să aibă valori pozitive.

II. La nivelul INS – Direcția Statisticii Comerțului Exterior, sunt monitorizate următoarele:

- i. Firmele prioritare* (primele 5000 de firme din punct de vedere a valorii cumulate realizate pe fiecare flux în ultimele 12 luni) – se verifică dacă firma a transmis declarația Intrastat sau este non-respondent

ii. Erorile de încărcare generate de aplicație. Aplicația Intrastat online generează pentru fiecare grup de declarații Intrastat offline încărcate în baza de date, câte un BATCH conținând erorile generate automat. Evidența acestor BATCH-uri se ține într-un registru special de BATCH-uri. Tratarea erorilor din aceste BATCH-uri se face progresiv, în Registru fiind înscris numele persoanei care a preluat spre rezolvare respectivele erori. Pentru rezolvarea fiecărei erori de încărcare (validare primară), persoanele responsabile contactează firmele prin email sau telefon pentru corectarea și retransmiterea declarațiilor Intrastat.

Erorile generate automat sunt următoarele:

- Cod eroare UPL-1 Fisierul incarcat este corupt sau invalid;
- Cod eroare UPL-2 Declaratia Intrastat este pentru o firma necunoscuta in Registru Intrastat;
- Cod eroare UPL-3 Codul de judet al PDT este invalid;
- Cod eroare UPL-4 Codul de oras al PDT este invalid;
- Cod eroare UPL-5 Corelatie invalida oras / judet;
- Cod eroare UPL-6 Parte declaranta referita in fisierul incarcat nu este utilizatorul curent;
- Cod eroare UPL-7 Nu poate fi incarcata o declaratie revizuita (tip REVIZUIT), daca nu a fost incarcata in prealabil o declaratie de tip NOU;
- Cod eroare UPL-8 Nu poate fi incarcata o a doua declaratie pentru aceiasi perioada de referinta decat marcata ca REVIZUITA;
- Cod eroare UPL-9 Firma referita nu a fost asociata cu PDT-ul referit in Declaratie;
- Cod eroare UPL-10 PDT referit in declaratie nu este cunoscut in Registrul Intrastat;
- Cod eroare UPL-101 Continut invalid al declaratiei...:Cod NC invalid;
- Cod eroare UPL-102 Continut invalid al declaratiei...:Cod UMS invalid;
- Cod eroare UPL-103 Continut invalid al declaratiei...:Masa net nu poate fi 0;
- Cod eroare UPL-104 Continut invalid al declaratiei...:Cod termeni de livrare invalid;
- Cod eroare UPL-105 Continut invalid al declaratiei...:Cod natura tranzatiei A invalid;
- Cod eroare UPL-106 Continut invalid al declaratiei...:Cod natura tranzatiei B lipsa;
- Cod eroare UPL-107 Continut invalid al declaratiei...:Cod natura tranzatiei B invalid;
- Cod eroare UPL-108 Continut invalid al declaratiei...:Cod mod de transport invalid;
- Cod eroare UPL-109 Continut invalid al declaratiei...:Cod tara de origine invalid;
- Cod eroare UPL-110 Continut invalid al declaratiei...:Cod tara de expeditie invalid;
- Cod eroare UPL-111 Continut invalid al declaratiei...:Cod tara de destinatie invalid;
- Cod eroare UPL-112 Continut invalid al declaratiei...:masa neta trebuie sa fie 0)

iii. Firmele care transmit declarația Intrastat pe flux invers față de obligație – Rapoartele din Aplicația Intrastat online atenționează asupra firmelor ce au în sistem declarații Intrastat pentru alt flux decât pentru cel pentru care au obligații legale de declarare (declarare pe flux invers). Firmele sunt contactate și monitorizate astfel încât să revizuiască și să retransmită declarațiile corecte.

iv. Firmele care au greșit perioada de raportare - Rapoartele din Aplicația Intrastat online atenționează asupra firmelor ce au în sistem declarații Intrastat pentru altă perioadă de raportare decât perioada de referință la momentul respectiv. Firmele sunt contactate și monitorizate astfel încât să revizuiască și să retransmită declarațiile corecte.

v. Firmele care au transmis declarații duble - Raportul generat odată cu firmele prioritare pentru perioada de referință respectivă atenționează asupra firmelor ce au în sistem declarații Intrastat cu valori identice în luna de referință curentă față de luna de referință anterioară sau au valori identice pe cele două fluxuri (achiziții=expedieri). Firmele sunt contactate și monitorizate astfel încât să revizuiască și să retransmită declarațiile corecte.

vi. Firmele care au transmis declarația Intrastat pentru prima dată: declarațiile acestor firme nu pot fi încărcate în Aplicația Intrastat online până nu sunt completate în Registrul Intrastat datele demografice ale acestora. În consecință, înainte de încărcarea datelor pentru firmele noi, codul de identificare fiscală al acestora este transmis responsabilului de la Registrul Intrastat pentru adăugarea datelor demografice și încărcarea sa în Registrul Intrastat, utilizând funcția dedicată din Aplicația Intrastat online. După această procedură, declarația Intrastat a firmei noi este încărcată în baza de date.

Validarea secundară este prezentată în capitolul II.3 Validarea datelor.

II.3. Validarea datelor

Validarea datelor constă în trei tipuri de activități: detectarea erorilor, analiza erorilor și a datelor și corectarea erorilor.

❖ CE SUNT ERORILE, CLASIFICAREA ȘI DESCRIEREA LOR

Eroarea reprezintă o valoare a unei variabile care nu respectă (încalcă) o regulă dinaintea stabilită (condiție logică).

În cercetările prin sondaj se disting 3 categorii principale de erori:

- erori de eșantionare (sampling errors);
- erori nelegate de eșantionare (non-sampling errors) Erorile de eșantionare nu sunt generate de erori în datele colectate și fac obiectul altei analize (ESSNET pag 7);
- erori de estimare care apar în procesul de extrapolare a rezultatelor obținute la nivelul eșantionului la nivelul întregii populații și privesc în principal procesul de proiectare și implementare a eșantionului;

Erorile nelegate de eșantionare se clasifică în:²

- erori de acoperire
- erori de măsurare
- erori de procesare (inclusiv introducere, codificare, agregare)
- non-răspunsuri

După sursa care introduce eroarea putem distinge:

² Di Zio M., Fursova N., Gelsema T., GieBing S., Guarnera U., Petrauskiene J., Quensel von Kelben L., Scanu M., Bosch K.O. van der Loo M., Wlsdorfe K. - *Methodology for data validation, Essnet Validat Foundation, pag.7, 2015*

- erori ale unității observate
- erori ale operatorului de interviu
- erori ale mijloacelor de înregistrare (erori de introducere generate de operator)
- erori generate de metodă (erori de codificare sau erori introduse de programul informatic)
- erori generate de factori externi

O altă clasificare a erorilor este cea în *atenționări sau erori fatale*.

Eroarea fatală reprezintă constatarea unei situații care nu poate exista sub nici o formă în realitate (situația reală întâlnită în teren nu poate fi adevărată decât în forma definită de condiția logică).

Exemple de erori fatale:

- anul de naștere al unei persoane este 1490;
- nivel de instruire superior absolvit de o persoană în vârstă de 12 ani;

Atenționarea reprezintă constatarea unei situații care în general nu poate fi adevărată, dar care ÎN MOD EXCEPȚIONAL, poate exista în realitate (situația reală reprezintă încălcarea unei condiții logice care în general este adevărată, dar care constituie o excepție a regulii definite de condiția logică).

Exemple de atenționări:

- cazul unei persoane care are vârsta de 115 ani;
- cazul unei femei care declară că a născut 15 copii;
- un bărbat care a declarat că statutul său ocupațional este de persoană casnică;

Altă clasificare a erorilor:

- **sistematice** – sunt erori de același tip (identice) întâlnite pentru mai mulți respondenți și au de obicei drept cauză probleme în proiectarea cercetării (ambiguități în chestionar sau precizările metodologice) sau în implementarea acesteia (instruirea insuficientă a operatorilor de interviu). Deoarece mai mulți respondenți greșesc în același sens, erorile sistematice, necorectate, conduc la obținerea de estimatori deplasati.
- **aleatorii** – sunt erori care apar accidental și se datorează în principal neatenției respondentului, operatorului de interviu sau a persoanei care introduce datele.

❖ ÎN CE ETAPE APAR ERORILE

Dată fiind tipologia vastă a erorilor acestea pot apărea în toate etapele unei cercetări statistice. În continuare sunt descrise cele mai probabile erori care pot apărea, în funcție de momentul apariției lor.

În faza de proiectare a chestionarului:

- prin definirea unui salt greșit între întrebările din cuprinsul acestuia, se pot induce non-răspunsuri parțiale, prin omiterea completării unor variabile din cauza saltului greșit definit;
- prin formularea unei întrebări foarte lungi sau, dimpotrivă, a unei întrebări laconice care nu definește corect subiectul întrebării, se pot induce erori determinate de către respondent, prin neînțelegerea corectă a întrebării. Același tip de erori poate fi generat de folosirea unui limbaj prea tehnic sau complicat, fără a defini respondentului noțiunile utilizate;
- prin proiectarea unui chestionar foarte lung, se pot induce non-răspunsuri parțiale la acele întrebări aflate la sfârșitul chestionarului, datorate plictisului sau oboselii respondentului;

- prin solicitarea unor răspunsuri referitoare la evenimente care au avut loc într-un trecut îndepărtat se pot induce erori de amintire sau non-răspunsuri parțiale etc.

În faza de colectare a datelor:

În această fază erorile se pot datora atât respondenților, cât și operatorilor de interviu.

Respondenții pot induce erori prin declararea unor răspunsuri eronate la întrebările unei cercetări statistice, generate în principal de:

- incapacitatea de a înțelege întrebările;
- definirea incorectă sau ambiguă a întrebărilor;
- folosirea unui limbaj prea tehnic sau complex în cadrul chestionarului;
- imposibilitatea de a-și aminti corect anumite evenimente din trecutul prea îndepărtat;
- solicitarea de a efectua operațiuni foarte complexe pentru obținerea răspunsurilor;
- solicitarea de a răspunde la întrebări cu un caracter intim;
- din neînțelegerea caracterului confidențial acordat de sistemul statistic informațiilor colectate de la unitățile statistice și a principiului conform căruia datele statistice colectate sunt utilizate exclusiv în scopuri statistice etc.

Erorile generate de operatorii se interviu pot fi cauzate de:

- înregistrarea altor variante de răspuns din cauza modului de redactare al chestionarului (economie de spațiu);
- neînțelegerea noțiunilor utilizate, a definițiilor sau nefurnizarea unor explicații suplimentare în cazul în care respondentul le solicită;
- operatorul de interviu nu adresează întocmai respondentului întrebările din chestionar;
- introducerea unui salt „fals” pentru a urma un drum mai scurt în interiorul chestionarului;
- codificarea unor răspunsuri din propria experiență, fără a mai întreba respondentul etc.

În faza de validare a datelor

Chiar și în faza de validare a datelor se pot introduce erori. Acestea se datorează, în principal, factorului uman.

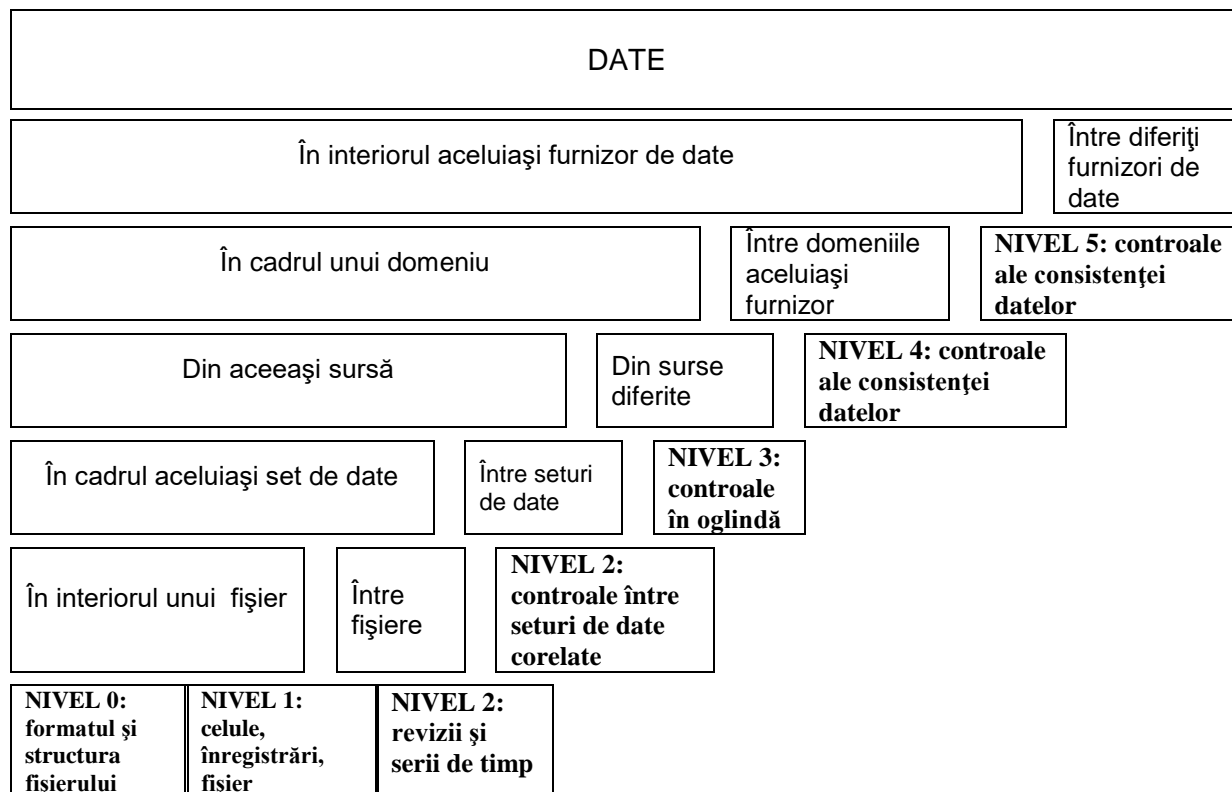
Erori generate de operatorul PC:

- erorile de tastare (introducere) sau de înțelegere greșită a răspunsurilor de pe chestionarele completate;
- erori sistematice de introducere a altor răspunsuri (decât cele completate pe chestioar) care să genereze un flux mai scurt în interiorul chestionarelor (generarea unui salt „fals”);
- alegerea primului răspuns dintr-un nomenclator pre-definit, în cazul codificării automate etc.;

Definirea greșită a unei corecții automate poate determina introducerea unor erori noi, fie prin faptul că valoarea aplicată prin corecția automată este una greșită, care intră în contradicție cu valorile altor variabile, fie din cauza faptului că se pot șterge valorile înregistrate ale unei variabile (pierderea unei informații deja existentă în baza de date).

❖ NIVELURI ALE VALIDĂRII

Deoarece validarea datelor este un amalgam de etape și proceduri care se realizează în diferite etape ale procesului de producție statistică, ea poate fi ierarhizată pe mai multe niveluri. Acestea sunt descrise într-o formă sistematică în documentul **Metodologia de validare a datelor** conform schemei următoare³:



Validarea de nivel 0 – presupune verificarea corespondenței dintre structura fișierelor și a atributelor acestora și cerințele formulate. Regulile de validare pot fi formulate ca:

- numărul de câmpuri (coloane) ale unui fișier este cel specificat
- formatul de date dintr-un câmp (colană) este cel specificat (numeric, caracter, dată etc.)

Validarea de nivel 1 – presupune verificarea consistenței informațiilor conținute într-un singur set de date, la un moment dat. Regulile de validare pot verifica în acest caz:

- la nivel de microdate:
 - dacă valorile într-un câmp (coloană) sunt valide, de exemplu:
 - numărul din coloana 4 este ne-negativ;
 - anul din coloana 2 este 2011;
 - valoarea din coloana 3 este un cod de activitate valid din Nomenclatorul Activităților Economiei Naționale

³ Di Zio M., Fursova N., Gelsema T., GieBing S., Guarnera U., Petrauskiene J., Quensel von Kelben L., Scanu M., Bosch K.O. van der Loo M., Wlsdorfe K. - *Methodology for data validation*, Essnet Validat Foundation, pag.11, 2015

- dacă valorile dintr-un câmp se încadrează într-un interval de valori plauzibile (de exemplu numărul de membri ai gospodăriei)
- dacă combinația de valori înregistrate în 2 sau mai multe câmpuri este permisă sau logică
- la nivel de macrorodate, de exemplu:
 - Total populație = Total bărbați + Total femei
 - Numărul de femei = (Total populație/2) ± 10%

Validarea de nivel 2 – presupune verificarea consistenței informațiilor dintr-un fișier cu:

- alte versiuni ale aceluiași fișier, referitoare la aceeași perioadă de timp (pentru detectarea reviziilor)
- instante ale aceluiași fișier referitoare la alte perioade de timp (pentru verificarea plauzibilității evoluțiilor în timp)
- alte fișiere corelate

Validarea de nivel 3 – presupune verificarea consistenței între date referitoare la același domeniu dar având surse diferite. În această categorie intra de exemplu statisticile în oglindă. De exemplu, exportul raportat de țara A în țara B trebuie să fie egal cu importul declarat de țara B din țara A.

Validarea de nivel 4 – presupune verificarea plauzibilității datelor referitoare la același fenomen dar provenind din domenii diferite, produse de aceeași instituție, având în vedere și diferențele de metodologie. De exemplu plauzibilitatea evoluției numărului de salariați din Ancheta forței de muncă și cel provenit din Cercetarea statistică privind costul forței de muncă în unitățile economico-sociale.

Validarea de nivel 5 – presupune verificarea plauzibilității statisticilor referitoare la un fenomen dat între diverse instituții producătoare de statistici.

Există validări de verificare:

1. **a integrității** (completitudine) – se verifică dacă s-au colectat toate datele care trebuiau colectate.
2. **a unicității a unei înregistrări** (să existe una și doar o înregistrare pentru fiecare unitate de observare sau element de detaliere);

Verificarea integrității se face: în interiorul fiecărui fișier și între fișiere.

Exemple:

Dacă există completat cel puțin un chestionar de gospodărie (indiferent de rezultatul interviului) pentru toate adresele selectate în eșantionul anchetei din fiecare centru de cercetare.

Se verifică dacă toți membrii unei gospodării enumerați în Componenta gospodăriei (sau doar cei eligibili conform unor anumite criterii) au chestionare individuale în fișierul individual.

3. de respectare a „drumurilor” și a filtrelor în cadrul unui chestionar

4. la nivel de variabilă

- de identitate (ANUL INTERVIULUI = 2016)
- de încadrare a variabilei numerice într-un interval plauzibil ($0 \leq \text{VARSTA} \leq 120$)

5. între variabilele aceluiași set de date (dacă STATUT = salariat atunci $\text{VARSTA} \geq 15$)

❖ CARACTERISTICILE UNEI PROCEDURI DE VALIDARE A DATELOR

O procedură de validare a datelor presupune definirea unui set de reguli de validare. Prin aplicarea acestui set de reguli de validare pe toate înregistrările individuale ale unităților statistice incluse într-o cercetare statistică se verifică gradul în care răspunsurile înregistrate (valorile variabilelor colectate) respectă regulile de validare. Cazurile de încălcare ale regulilor de validare constituie erori care, în funcție de tipul lor (erori fatale sau atenționări), trebuie corectate sau acceptate.

Metodologia de validare a datelor⁴ elaborată de Essnet Validat Foundation – Comisia Europeană definește urătoarele caracteristici pe care trebuie să le îndeplinească o procedura de validare a datelor:

1. **completitudinea** – teoretic, o procedură trebuie să cuprindă toate condițiile posibile de validare. În practică însă, gradul de completitudine a unei proceduri de validare depinde de gradul de general de cunoaștere a fenomenului, de experiența și gradul de inteligență al persoanei care definește regulile și de gradul de repetativitate al cercetării respective (dacă este o cercetare realizată pentru prima dată sau dacă a mai fost realizată de mai multe ori în trecut. Din perspectiva acestei caracteristici un set de reguli de validare poate fi foarte restrictiv, datorită faptului că include un număr foarte mare de condiții sau pentru că acestea sunt foarte stricte, sau dimpotrivă, mai „relaxat” cuprinzând un număr mic de reguli, de regulă cele mai elementare;
2. **redundanța** – un set de proceduri de validare nu trebuie să cuprindă reguli redundante. Acestea pot interveni de regulă, prin includerea în diferite categorii de condiții de control (de salt, de verificare a valorilor aberante etc.), a două sau mai multe reguli de validare care verifică, practic, același lucru;
3. **fezabilitatea** - este caracteristica prin care un set de reguli de validare este capabil să detecteze toate erorile posibile.
4. **complexitatea** este definită de varietatea și cantitatea de informație care este necesară pentru definirea unei reguli de validare și pentru evaluarea, respectiv corecția erorilor generate de aceasta, precum și pentru găsirea soluției corecte de rezolvare a respectivelor erori.

❖ MODALITĂȚI DE VALIDARE A DATELOR

Există diverse perspective prin prisma cărora se poate realiza validarea datelor, între care cele mai importante fiind⁵ :

1. **automatizat sau manual;**

Validarea manuală a datelor poate avea loc înainte de introducerea datelor. Datele pot fi validate și pregătite / corectate înainte de introducere. Această procedură este utilizată în special în cazul introducerii datelor în modul „heads-down”.

Validarea automată a datelor poate avea loc fie în mod batch (validare pe pachete de date) fie în mod interactiv (validare individuală a datelor). Este important faptul că datele introduse în mod heads-down pot fi validate ulterior, fie în mod batch, fie în mod interactiv.

- **Validarea datelor în mod batch (validare pe pachete de date)** are loc după introducerea datelor și constă în validarea mai multor chestionare într-un singur batch (pachet). În

⁴ Di Zio M., Fursova N., Gelsema T., Gießing S., Guarnera U., Petrauskiene J., Quensel von Kelben L., Scanu M., Bosch K.O. van der Loo M., Wlisdorfe K. - *Methodology for data validation, Essnet Validat Foundation, pag.39-52,, 2015*

⁵ Di Zio M., Fursova N., Gelsema T., Gießing S., Guarnera U., Petrauskiene J., Quensel von Kelben L., Scanu M., Bosch K.O. van der Loo M., Wlisdorfe K. - *Methodology for data validation, Essnet Validat Foundation, pag.9, 2015*

general, rezulta un fisier cu mesaje de eroare. Acest fisier poate fi tiparit si utilizat pentru corectarea erorilor. Inregistrările de date pot fi impartite in doua fisiere: un fisier care conține înregistrări "bune" și un fisier conținând înregistrări de date cu erori. Fisierul continand date cu erori poate fi corectat printr-un proces interactiv.

- **Validarea interactivă a datelor** implică validarea chestionarului imediat ce corecțiile/ajustările au fost făcute. Rezultatele validării sunt vizualizate pe un terminal de afișare video iar cel care editează datele poate, fie să revizuiască/ajusteze datele fie, după caz, să ignore marcajul de eroare. Acest proces continuă până când chestionarul este considerat acceptabil de către procesul de validare automată. Apoi, un alt chestionar (următorul) va intra în procesul de validare automată. O caracteristică de dorit a software-ului de editare interactivă a datelor este aceea de a supune atenției editorului (vizualiza) doar acele chestionare care necesită corecții/ajustări.

Validarea realizată în timpul interviului asistate de calculator (Computer-Assisted Interviewing - CAI) combină validarea interactivă a datelor cu editarea interactivă a datelor în timp ce respondentul este o sursă disponibilă pentru ajustarea/corectarea datelor. Un avantaj suplimentar este acela că introducerea/capturarea datelor (keyentry) are loc în momentul interviului. Această metodă poate fi utilizată atât în timpul interviului telefonice cât și în cazul interviurilor/colectărilor realizate față în față cu respondentul, utilizând dispozitive portabile de introducere a datelor.

CAI asistă interviuatorul în formularea întrebărilor și îndrumă spre întrebările următoare pe baza răspunsurilor anterioare. Este un instrument util pentru reducerea timpului de realizare a unui interviu, ajutând interviuatorii mai puțin experimentați. CAI a fost folosită în principal în interviurile telefonice asistate de calculator (CATI) însă, dat fiind avansul tehnologiilor moderne ce au permis miniaturizarea calculatoarelor personale, interviurile asistate de computerele personale (CAPI) vor fi utilizate din ce în ce mai mult.

2. **obiectiv versus subiectiv** (bazată pe opinia și experiența experților);
3. **validare structurală versus validare de conținut**;
4. **nivelul la care se realizează validarea**: în cadrul unei înregistrări, între înregistrări, între seturi de înregistrări etc.;

Detectarea erorilor poate avea loc la mai multe nivele.

- **La nivel de linie/articol/la nivel elementar** - Validările la acest nivel sunt în general denumite "range checking" (verificare interval), din moment ce elementele sunt validate pe baza unor intervale. Exemplu: vârsta trebuie să fie > 0 și < 120 . Intervalul de verificare poate varia în funcție de straturi sau de alt identificator. Exemplu: În cazul în care stratul = "producție agricolă mare" atunci dimensiunea terenului trebuie să fie mai mare decât 500 hectare.
- **La nivel de chestionar** – Acest nivel implică verificarea corelațiilor dintre variabilele din întregul chestionar. Exemplul 1: Dacă căsătorit = 'da' atunci vârsta trebuie să fie mai mare de 16. Exemplul 2: Suma tuturor terenuri individuale trebuie sa fie egală cu dimensiunea totală a exploatației agricole.
- **Ierarhică** – Acest nivel implică verificarea elementelor din chestionarele aferente unităților care compun o altă unitate (de exemplu, chestionarele individuale ale membrilor aceleiași gospodării). Relațiile dintre datele de acest tip sunt cunoscute ca „date ierarhice”. În exemplul dat, informațiile comune la nivel de gospodărie se află pe un chestionar și informațiile despre fiecare individ din gospodărie se află pe chestionare diferite/separate. Se fac verificări pentru a asigura faptul că suma datelor individuale pentru un item/individ nu este mai mare decât totalul raportat pentru întreaga gospodărie.

Validare aplicată la nivelul tuturor chestionarelor anchetei (Across Questionnaire level) implică calculul intervalelor de validare pentru fiecare item/element din anchetă sau utilizarea

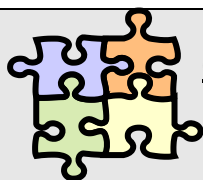
datelor istorice pentru detectarea valorilor aberante (outliers). Rutinele de analiză a datelor, care de obicei sunt rulate în momentul sintezei rezultatelor (totalizării datelor), pot fi incorporate mai ușor în validarea datelor la acest nivel. În acest fel, erorile sunt detectate suficient de devreme pentru a fi corectate în timpul procedurilor uzuale/obișnuite de corectare a datelor. Verificările la nivel de chestionare ar trebui să identifice chestionarul specific care conține date discutabile. Modificările la nivel de chestionar sunt, în general, de două tipuri: modificări statistice (**statistical edits**) și modificări macro (**macro edits**).

- **Modificările statistice (Statistical Edits)** utilizează distribuția datelor pentru detectarea posibilelor erori. Aceste proceduri utilizează datele curente din mai multe/ din toate chestionarele sau datele istorice ale unităților statistice pentru a genera limite acceptabile pentru datele anchetei curente/actuale. Valorile aberante (**Outliers**) pot fi identificate în funcție de **limitele de acceptabilitate**. **Inliers** sunt date ce se află în interiorul limitelor de acceptabilitate, dar sunt considerate ca fiind suspecte din cauza lipsei variațiilor în timp.

Variabilele aleatoare presupun un grad măsurabil de variație. Dacă valoarea este prea „consistentă”, atunci ea ar putea fi pur și simplu copiată dintr-un chestionar anterior în loc să fie raportată în chestionarul curent. De aceea, testul constă în compararea modificărilor în timp a unei unități din eșantion (comparison to the double root residual of a sample unit over time). Dacă testul eșuează, atunci variația nu este suficient de aleatoare și chestionarul ar trebui verificat. În cazul USDA-NASS acest test se aplică datelor privind greutatea animalelor sacrificate. Ipoteza este că numărul de capete de porci sacrificați nu poate varia foarte mult de la o săptămână la alta. Însă, greutatea totală a tuturor porcinelor sacrificate este o variabilă aleatoare și ar trebui să arate un grad de variație măsurabilă în fiecare săptămână.

- **Modificările macro (Macro Edits)** asunt validări ale datelor la nivele agregate. Inconsistențele sunt urmarite la nivelul înregistrărilor individuale implicate. O mare parte din activitatea curentă în acest domeniu este realizată de Leopold Granquist (1991) de la Statistica Suediei. Munca sa se bazează pe convingerea că este preferabil să determini erorile cu impact mare la nivel total și de a evita ajustările care nu au impact la nivel de total.
- Procesul de revizuire a datelor ar trebui să permită detectarea erorilor de diferite niveluri de severitate. De asemenea, ar trebui să permită decizia privind corectarea sau nu a erorii.

5. **momentul când se face validarea:** în timpul colectării datelor, în timpul introducerii în baza de date, în timpul validării, în etapa de agregare a datelor etc.;
6. **locul în procesul de producție:** input, throughput, output;
7. **tipul de regulă de validare:** egalitate, încadrarea într-un interval, condiție logică etc.



Caseta 5: Exemplificare din **STATISTICA ÎNTRERINDERILOR**

În sistemul statistic **Intrastat**, după colectarea datelor și validarea primară, este realizată validarea secundară a acestora:

I. La nivel de linie din declarație:

Sunt analizate tranzacțiile (liniile cu date) în următoarea ordine:

a. *tranzacțiile cu valori foarte mari:*

- valoare facturată sau valoare statistică > 5.000.000 LEI
- masa netă > 5.000.000 Kg
- UMS > 1.000.000

b. Codul din Nomenclatorul Combinat (NC) 27090090 – *Țiței* - majoritatea importului este din spațiul Extra-UE; importul/exportul din/în spațiul Intra-UE are o probabilitate foarte mică de realizare

c. Codurile NC – 27160000 - *Energie electrica și 27112100 - Gaz natural (prin comparație cu informațiile obținute din surse administrative)*

d. *tranzacțiile cu erori importante:*

- cod eroare 1006.1 Valoare unitară aberantă (valoare / masa netă)
- cod eroare 1006.2 Valoare unitară aberantă (valoare / cantitate)
- cod eroare 1006.3 Valoare unitară aberantă (masa netă/ cantitate)
- cod eroare 1004 Raportul dintre valoare statistică și valoare facturata <0.5 sau >1.5
- cod eroare 1009 Masa netă = 0
- cod eroare 1006.1 $c \geq 10$; eroare 1006.2 $c \geq 10$; eroare 1004, raport >4 sau raport <0.25, unde c reprezintă deviația față de medie
- cod eroare 1009; valori mai mari de 10.000 lei
- cod eroare 1006.1 $c < 10$; eroare 1006.2 $c < 10$; eroare 1004 ramase pentru valoare statistica sau facturata >100.000, unde c reprezintă deviația față de medie
- verificarea tranzacțiilor pentru care cantitatea exprimată în unitate de măsură suplimentară (UMS) = codul din NC pentru respectiva UMS

e. Gramele – rotunjire incorectă gram/Kg

f. Capitolul NC 87 pentru: masa netă/cantitate < 100 și pentru codurile pentru care cantitatea exprimată în UMS = masa netă exprimată în kg

g. Codurile NC 8901, 8902, 8904, 8905, 8906, 8801, 8802, 8805 (analiza strictă a tuturor tranzacțiilor pe aceste coduri)

h. Posibile erori legate de natura tranzacției

- natura tranzacției 4 la export
- natura tranzacției 5 la import
- natura tranzacției 6 sau 7

i. Alte erori

- valoarea statistică sau valoarea facturată = masa netă
- valoarea statistică sau valoarea facturată = UMS
- masa netă = UMS
- valoarea statistică = valoarea facturată =1

j. Alte atenționări semnalate de aplicație

- cod atenționare 1001.A Combinație invalidă: mod de transport (transport maritim) și țara de expediție fără ieșire la mare
- cod atenționare 1001.D Combinație invalidă: mod de transport (transport maritim) și țara de destinație fără ieșire la mare
- cod atenționare 1002.1 Combinație invalidă: mod de transport (propulsie proprie)

și codul NC al bunului

- cod atenționare 1002.2 Combinație invalidă: mod de transport (postal) și masa netă > 500 kg
- cod atenționare 1002.3 Combinație invalidă: mod de transport (exclusiv transport prin instalații fixe) și bun (codurile NC 27112100 și 27160000)
- cod atenționare 1003 Combinație invalidă: țara de origine și bun
- cod atenționare 1004 Raportul între valoare statistică și valoare facturată <0.5 sau >1.5 – ramase
- cod atenționare 1006.1 Valoare unitară aberantă (valoare / masa netă)
- cod atenționare 1006.2 Valoare unitară aberantă (valoare / cantitate)
- cod atenționare 1006.3 Valoare unitară aberantă (masa / cantitate)
- cod atenționare 1007.A Corelație invalidă între valoarea facturată și valoarea statistică în funcție de termeni de livrare (introduceri)
- cod atenționare 1007.D Corelație invalidă între valoarea facturată și valoarea statistică în funcție de termeni de livrare (expedieri)
- cod atenționare 1008.A Valoarea statistică este obligatorie (introduceri)
- cod atenționare 1008.D Valoarea statistică este obligatorie (expedieri)
- cod atenționare 1009 Masa netă = 0

II. Pentru toate declarațiile dintr-o lună de referință și pentru un flux (export/import): controale de credibilitate (valori aberante, încadrarea datelor transmise între anumite limite a rapoartelor preț / cantitate, preț / kg, kg / cantitate, diverse cazuri de incompatibilitate între diverse coduri – cod țară / produs, cod produs / mijloc de transport etc.);

II.4. Imputarea și ajustarea datelor

II.4.1. Imputarea non-răspunsurilor

Corectarea/Ajustarea manuală a datelor are loc atunci când selectarea unei valori mai rezonabile/bune se face de către o persoană. Aceasta poate implica scrierea corecțiilor, pentru introducere, printr-o procedură de tip batch. Corecțiile "manuale" de date pot avea loc de asemenea în mod interactiv, la fel ca în procesul "heads-up" de introducere a datelor sau validare interactivă a datelor.

Corectarea/ajustarea automată a datelor are loc ca rezultat al acțiunii computerului. O opțiune necesară în orice sistem de validare automată a datelor este posibilitatea marcării acelor niveluri pentru care aceste acțiuni să nu fie realizate (skip pentru anumite nivele). Corectarea datelor în batch produce un fișier cu înregistrări corectate (editate/imputate) și mesaje însoțitoare pentru acțiunile întreprinse de calculator în scopul ajustărilor.

Datele pot fi imputate pe baza unei game variate de metodologii, unele dintre acestea fiind mult mai ușor de programat decât altele. Cea mai simplă presupune realizarea calculelor într-un chestionar (de exemplu, obținerea unei sume lipsă în partea de jos a unei coloane).

Imputările automate se încadrează, în general, într-unul din următoarele 5 tipuri.

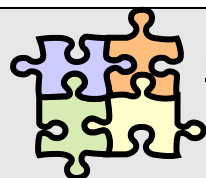
- a. **Deterministic** – in cazul in care exista o singură valoare corectă, ca și în suma lipsă în partea de jos a unei coloane de numere. O valoare este astfel determinată din alte valori aflate pe același chestionar.
- b. **Model based** – utilizarea mediilor, medianelor, ecuațiilor de regresie, etc, pentru a imputa o valoare.
- c. **Deck** – Este folosit un chestionar donator pentru a imputa valoarea lipsă.

Hot deck - un chestionar donator este găsit în aceeași cercetare ca și chestionarul cu elementul lipsă. Tehnica de căutare "**cel mai apropiat vecin**" ("**nearest neighbour**") este adesea folosită pentru a accelera căutarea unui donator de înregistrare. În această tehnică de căutare, pachetul de chestionare donatoare provine din aceeași cercetare și prezintă similitudini cu înregistrarea primitoare, în cazul în care similitudinea se bazează pe alte date din chestionar, care se corelează cu datele donate. De exemplu: dimensiunea și localizarea similară a unei ferme ar putea fi utilizată pentru donarea prețurilor la combustibili.

Cold deck - Similar ca la hot deck, cu excepția faptului că datele se găsesc într-o cercetare similară efectuată anterior.

- d. **Mixed** – In cele mai multe sisteme există, de obicei, o combinație de categorii utilizate. GEIS utilizat de Statistica din Canada (Generalized Edit și System imputare), de exemplu, folosește mai întâi o abordare deterministă. În cazul în care nu este de succes, atunci este încercată o abordare hot deck. Aceasta este urmată de o abordare model based. În cazul în care toate aceste abordări eșuează, atunci are loc o imputare manuală prin intervenție umană. Explicații mai detaliate ale metodelor a) - d) pot fi găsite în lucrarea lui Giles și Patrick, (1986).
- e. **Expert Systems** – Sistemele expert sunt doar recent aplicate pentru editarea datelor și multe cercetări sunt la început în acest domeniu. "Un sistem expert" este un program inteligent de calculator ce utilizează proceduri de cunoaștere și de inferență pentru a rezolva probleme destul de dificile pentru o expertiză umană. Fiecare sistem expert este alcătuit din două părți principale: Baza de cunoaștere și motorul de inferență. Baza de cunoaștere conține atât cunoștințe factuale cât și euristice. Cunoașterea factuală este data de elemente convenite de comun acord experții într-un anumit domeniu. Cunoașterea euristică este mai puțin riguroasă, mai experimentală și bazată pe reguli de "judecată bună" sau arta de a "ghici"/"presupune", într-un domeniu. O reprezentare utilizat pe scară largă pentru baza de cunoaștere este regula sau IF / THEN. Partea IF afișează un set de condiții într-o anumită combinație logică. O dată ce partea IF a regulii este îndeplinită, partea THEN poate fi încheiata sau problema rezolvată. Sistemele expert cu cunoștințe reprezentate în formă de regulă se numesc "**sisteme bazate pe reguli**" (**rule based systems**), (Magnas, 1989). Motorul de inferență face inferențe prin determinarea regulilor ce sunt satisfăcute prin fapte, prin ordonarea regulilor satisfăcute și executarea regulii cu cea mai mare prioritate.

O altă metodă de imputare a datelor în cazul non-răspunsurilor totale și/sau parțiale este preluarea din surse administrative disponibile.



Caseta 6: Exemplificare din **STATISTICA SOCIALĂ**

Un exemplu de cercetare statistică în care s-a realizat imputarea totală a unor înregistrări îl reprezintă Recensământul Populației și al Locuințelor din 20 octombrie 2011.

În urma procesului de prelucrare a formularelor individuale, sub-înregistrarea constatată în etapa prelucrării rezultatelor provizorii ale RPL 2011 s-a confirmat. Prin urmare, s-au aplicat

metode pentru asigurarea completitudinii datelor de recensământ, folosindu-se metoda colectării indirecte din surse administrative și metode statistice de imputare a înregistrărilor.

Numărul persoanelor nerecenzate în teren la Recensământul Populației și al Locuințelor din 20 octombrie 2011 (RPL 2011) și identificate în sursele administrative de date a fost 1.183 mii. Pentru aceste persoane s-au imputat total înregistrări individuale și respectiv, de gospodărie, locuință și clădire (dacă acestea nu existau deja în baza de date).

Sursele de date administrative identificate ca având informații utile pentru definitivarea rezultatelor RPL 2011 sunt cele cuprinse în:

- Registrul Național de Evidența Persoanei - RNEP – gestionat de Direcția pentru Evidența Persoanelor și Administrarea Bazelor de Date
- Declarația privind obligațiile de plată a contribuțiilor sociale, impozitul pe venit și evidența nominală a persoanelor asigurate - D112 – gestionată de Agenția Națională de Administrare Fiscală
- Registrul de Evidență a Salariaților - IM – gestionat de Inspekția Muncii
- Baza de date CNPP – gestionată de Casa Națională de Pensii Publice
- Baza de date CNAS – gestionată de Casa Națională de Asigurări de Sănătate
- Declarația de înregistrare fiscală/Declarație de mențiuni pentru persoanele fizice care desfășoară activități economice în mod independent sau exercită profesii libere - D70 – gestionată de Agenția Națională de Administrare Fiscală
- Registrul de evidență a beneficiarilor de alocație de stat pentru copii, alocație de susținere a familiei și de ajutorul minim garantat – gestionat de Agenția de Plăți și Inspekție Socială
- Baza de date a elevilor înscriși în anul școlar 2011-2012 – gestionată de Ministerul Educației Naționale.

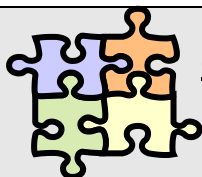
Procedura de colectare indirectă din sursele administrative a avut ca punct de plecare compararea înregistrărilor din baza de date a RPL 2011 (informațiile înregistrate în formulare individuale completate de recenzori în perioada de colectare în teren) cu înregistrările existente în baza de date de la Evidența Populației. Setul de înregistrări găsite la Evidența Populației care lipseau din baza de date a RPL 2011 (din toate formularele⁶ - P, PPI sau TP) au fost căutate în alte surse administrative aferente lunii octombrie 2011 și lunilor precedente și ulterioare din același an, în conformitate cu conceptul reședinței pe teritoriul României pentru o perioadă de cel puțin 12 luni, concept aplicat la RPL 2011 pentru măsurarea populației rezidente (stabile). Doar pentru persoanele identificate în sursele administrative utilizate, pentru care existau dovezi clare că au fost pe teritoriul României în perioada recensământului și în cea mai mare parte a anului 2011, s-a adăugat (s-a imputat) o înregistrare în baza de date a RPL 2011, pentru care s-au completat ulterior identificatori și valori pentru variabilele de recensământ. În acest fel, toate înregistrările obținute prin colectare indirectă au avut aceeași structură cu restul înregistrărilor obținute prin interviuarea persoanelor în perioada de colectare a datelor în teren și s-au referit la aceeași perioadă de referință, permițând agregarea informației pentru întreaga populație rezidentă (stabilă) a țării, indiferent dacă aceasta a fost interviuată de recenzori sau nu.

Pentru adulții identificați (mamă și tată) pentru care s-au preluat informațiile din sursele administrative de mai sus, precum și pentru adulții (mamă și tată) care au fost recenzați la RPL 2011, au fost căutați și identificați minorii în Registrul Național de Evidența Persoanei.

⁶ P – Persoană (prezentă sau temporar absentă); PPI – Persoană plecată pentru o perioadă îndelungată din gospodărie (în țară sau în străinătate); TP – persoană temporar prezentă.

Ponderea persoanelor adulte pentru care s-au colectat informațiile indirect din sursele administrative a fost de: 64,9% din declarația D112, 0,5% din registrul de evidență a salariaților, 2,1% din declarația D070, 16,1% din Registrul de evidență a beneficiarilor de alocație de stat pentru copii, alocație de susținere a familiei și de ajutor minim garantat, 4,6% din baza de date a persoanelor asigurate la sistemul public de asigurări de sănătate și 11,8% din combinații de mai multe surse.

În Anexa 3 este prezentat un exemplu privind condițiile de control pentru imputarea non-răspunsurilor variabilei VNET – 2015 din Ancheta asupra forței de muncă în gospodării.



Caseta 7: Exemplificare din **STATISTICA AGRICOLĂ**

Tratarea non-răspunsurilor în Ancheta Structurală în Agricultură (ASA) 2013 se realizează diferit, în funcție de informațiile înregistrate în cadrul capitolului privind codul de completitudine.

Se presupune că avem următoarele variante de răspuns în cadrul codului de completitudine:

- a) Interviu complet
- b) Exploatație agricolă desființată (sau comasată cu altă exploatație)
- c) Exploatație agricolă temporar fără activitate
- d) Interviu refuzat
- e) Exploatație agricolă neidentificată
- f) Exploatație agricolă necontactată
- g) Alte situații

Baza de eșantionare a fost mai întâi stratificată după criteriile stabilite inițial (de ex.: județ, regiune, clase de mărime, statut juridic etc.). Din această bază a fost extras eșantionul prin metoda alocării proporționale în cadrul fiecărui strat.

În cele ce urmează este descrisă modalitatea efectivă de tratare a non-răspunsurilor.

I) Asupra exploatațiilor agricole desființate sau comasate cu altă exploatație (de ex.: având codul de completitudine = b), sau pentru exploatațiile agricole temporar fără activitate (de ex: cod completitudine = c), se acționează în două moduri:

1. Nu se efectuează nici o operațiune (imputare sau ajustare).

Se recurge la această metodă, presupunând că valorile indicatorilor statistici aferenți acestor unități vor dispărea după desființarea acestora, sau se vor însuma în cadrul altor unități, iar orice altă acțiune de imputare sau ajustare a acestora ar duce la creșterea în mod artificial și nejustificat a valorii totale a indicatorilor (în cazul variabilelor numerice, însumabile în special).

2a. Se realizează procedura de imputare a acestor unități, prin înlocuirea valorii fiecărui indicator statistic ale acestora cu valorile indicatorilor similari, de la o unitate având aceleași caracteristici (de ex: din același strat) cu unitatea imputată.

2b. Se realizează procedura de imputare a acestor unități, prin înlocuirea valorii fiecărui indicator statistic cu media valorilor indicatorilor similari ale tuturor unităților din același strat cu unitatea imputată.

II) În cazul exploatațiilor agricole cu interviu refuzat (cod completitudine = d), exploatațiilor agricole mutate la o adresă necunoscută (cod completitudine = e), pentru exploatațiile agricole cu interviu nerealizat (cod completitudine = f), precum și pentru celelalte situații se aplică metoda ajustării coeficienților de extindere pentru fiecare strat în parte.

Rezultatul ajustării va fi următorul:

$$K_{if} = \frac{N_{i2}}{n_{i2}},$$

Unde:

$$N_{i2} = N_{i1} - X_1$$

$$n_{i2} = n_{i1} - X_2$$

K_{if} = coeficientul final de extindere al stratului „i”;

N_{i2} = numărul final de unități din baza de eșantionare, aparținând stratului „i”, după eliminarea exploatațiilor agricole desființate sau comasate cu altă exploatație;

n_{i2} = numărul final de unități din eșantion, aparținând stratului „i”, după eliminarea exploatațiilor agricole desființate sau comasate cu altă exploatație

X_1 = exploatații agricole desființate (sau comasate cu altă exploatație) și cele temporar fără activitate

X_2 = reprezintă toate exploatațiile agricole considerate ca non răspunsuri

Coeficientul inițial de extindere:

$$K_{i1} = \frac{N_{i1}}{n_{i1}}$$

Coeficientul final de extindere:

$$K_{if} = K_{i1} * C_{ir}$$

$$C_{ir} = \frac{n_{i1}}{n_{i1} - n_{ix}},$$

$$K_{if} = \frac{N_{i1}}{n_{i1}} * \frac{n_{i1}}{n_{i1} - n_{ix}}, \quad \Rightarrow \quad K_{if} = \frac{N_{i1}}{n_{i1} - n_{ix}}$$

unde:

K_{i1} = coeficientul inițial de extindere al stratului „i”;

N_{i1} = numărul inițial de unități din cadrul de eșantionare, aparținând stratului „i”;

n_{i1} = numărul inițial de unități din eșantion, aparținând stratului „i”;

K_{if} = coeficientul final de extindere, aparținând stratului „i”;

C_{ir} = coeficientul final de ajustare, aparținând stratului „i”;

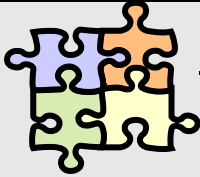
n_{ix} = numărul tuturor non-răspunsurilor (fără exploatațiile agricole desființate sau comasate cu altă exploatație și cele temporar fără activitate), aparținând stratului „i”;

II.4.2. Ajustarea valorilor imputate

Această activitate constă în verificarea consistenței valorilor imputate.

În cele mai multe cazuri, regulile de editare nu sunt luate în considerare de metodele de imputare. Ca o consecință, înregistrările imputate sunt, în general, incompatibile cu regulile de editare.

Această problemă este în prezent rezolvată prin introducerea unei etape de ajustare în care sunt aduse modificări la valorile imputate astfel încât înregistrările să satisfacă toate regulile de editare și ajustările sunt cât mai mici posibil. Această problemă este formulată ca o problemă de programare liniară, minimizând suma valorilor absolute ale ajustărilor sub constrângerea că imputările ajustate rezultate satisfac regulile de editare. Un algoritm pentru a rezolva această problemă este pusă în aplicare în SLICE.



Caseta 8: Exemplificare din **STATISTICA ÎNTREPRINDERILOR**

Pentru cercetarea statistică Intrastat, se fac următoarele estimări și imputări de date:

I. Estimarea datelor pentru firmele non-raspondente și pentru firmele aflate sub pragurile de raportare, utilizând datele fiscale (TVA/VIES):

- utilizarea informațiilor fiscale (TVA/VIES) pentru estimările Intrastat implică faptul ca datele fiscale și datele Intrastat sa fie aduse la aceeași “sfera de cuprindere” - conținutul ambelor seturi de date - este similar din punct de vedere metodologic, dar nu egal. Prin urmare, mai întâi se calculează totalurile pentru datele de TVA/VIES și pentru datele Intrastat pentru acele societăți care au declarat în sistemul statistic Intrastat, în momentul estimării. Se realizează diverse ajustări pentru a aduce cele două seturi de date la un nivel comparabil. Astfel:
 - Din datele Intrastat se elimină operațiunile de comerț intracomunitar cu natura tranzacției 3 (tranzacții care implică transferul de proprietate dar fără beneficii), 4 (operațiuni în scopul prelucrării pe bază de contract), 5 (operațiuni ce succed prelucrării pe bază de contract) și 8 (furnizarea materialelor și echipamentelor de construcție).
 - Din datele de TVA se elimină: comerțul triunghiular, serviciile, plățile în avans, declaratiile TVA si/sau VIES gresite, declaratiile TVA si/sau VIES facute pe alt cui, declaratiile TVA si/sau VIES pentru plinuri combustibil pentru camioane, declaratiile TVA si/sau VIES pentru softul la comanda, declaratiile TVA si/sau VIES pentru valoarea certificatelor de emisii poluanti, declaratiile TVA si/sau VIES pentru licențele achiziționate pe Internet etc.
- pe baza acestor seturi de date “purificate”, se va calcula un raport care va fi folosit pentru a transforma valorile de TVA/VIES ale societăților care nu au declarat Intrastat (firme non-raspondente sau firme aflate sub pragurile de raportare) în valorile Intrastat (valori care ar fi trebuit declarate).
- Raportul se calculează pe baza totalurilor declarațiilor de TVA/VIES “purificate” și declarațiile Intrastat “purificate”, numai pentru o categorie de societăți, și anume pentru acelea care au declarat atât în deconturile de TVA/declarațiile recapitulative VIES-ul cât și în declarațiile statistice Intrastat.
- Acest raport este necesar deoarece, chiar și după “purificare”, cele două seturi de date produc încă valori diferite pentru categoriile menționate de societăți.

- Totalurile de TVA/VIES pentru societățile care nu au declarat Intrastat sunt multiplicare cu acest raport și rezulta estimările de TVA/VIES. Nu sunt folosite valorile din TVA/VIES pentru care $\left| \frac{TVA - VIES}{VIES} \right| * 100 \geq 50\%$ (se presupune ca datele de TVA sau cele de VIES sunt eronat completate).

II. Estimare valoare statistică

În România numai operatorii economici peste un anumit prag, sunt obligați să furnizeze valoarea statistică a tranzacțiilor, împreună cu valoarea facturată. Din această motiv trebuie estimată valoarea statistică pentru firmele sub acest prag. Pentru construirea matricii de estimare a valorii statistice vom lua în considerare următoarele variabile: condițiile de livrare și țara de destinație / expediție.

Pe baza datelor istorice declarate, referitoare la valoarea statistică și valoarea facturată, se compoză această matrice cu media rapoartelor pe flux, condiții de livrare și țară de destinație / expediție. Aceste rapoarte vor fi utilizate pentru estimare a valorii statistice pe baza valorii facturate:

Valoare statistică = Valoarea facturata * Medie raport_{conditie de livrare/tara de destinație sau expediție} .

❖ INDICATORI DE MĂSURARE A PERFORMANȚELOR PROCEDURII DE EDITARE

Calitatea sau eficiența unei proceduri de editare poate fi măsurată cu ajutorul unor indicatori de performanță care evidențiază, în principal, cât de eronate au fost datele brute și cât de multe operațiuni s-au efectuat asupra acestora pentru a se obține datele finale corecte.

În *Metodologia de validare a datelor* realizată de către Essnet Validat Foundation sunt prezentați următorii indicatori de performanță ai unei proceduri de editare⁷:

1. Numărul de înregistrări pentru care se înregistrează erori;
2. Numărul minim de variabile care trebuie schimbate pentru ca înregistrările să îndeplinească un set de reguli de validare;
3. Numărul înregistrărilor care au respectat o anumită regulă de validare;
4. Numărul înregistrărilor care au încălcat o anumită regulă de validare;
5. Distribuția înregistrărilor care au încălcat una, două sau k reguli de validare;
6. Numărul de înregistrări care îndeplinesc/ încălcă/ au non-răspuns parțial pentru fiecare variabilă;
7. Raportul dintre numărul înregistrărilor lipsă și numărul înregistrărilor greșite;
8. Diferența dintre valorile indicatorilor 1-6 din etapa curentă și cea anterioară de validare.

⁷ Di Zio M., Fursova N., Gelsema T., Gießing S., Guarnera U., Petruskiene J., Quensel von Kelben L., Scanu M., Bosch K.O. van der Loo M., Wlsdorfe K. - *Methodology for data validation*, Essnet Validat Foundation, pag.53, 2015

III. BIBLIOGRAFIE

- Bethlehem, J. G. (2009), *Applied Survey Methods*. Wiley Series in Survey Methodology, John Wiley
- Fellegi, I. P., Holt, D., (1976), *A Systematic Approach to Automatic Edit and Imputation*. Journal of the American Statistical Association, 71, 17—35.
- Granquist, L., (1995), *Improving the Traditional Editing Process*. In: *Business Survey Methods* (eds. Cox, Binder, Chinnappa, Christianson and Kott), John Wiley & Sons, New York, pp. 385-401
- Hartwig, P. (2009), *How to Use Edit Staff Debriefings in Questionnaire Design*. Paper presented at the
- Hoogland, J., van der Loo, M., Pannekoek, J., and Scholtus, S. (2011), *Data Editing: Detection and Correction of Errors*. Methods Series Theme, Statistics Netherlands, The Hague.
- Lindgren, K. (2012), *The Use of Evaluation Data Sets when Implementing Selective Editing*. Working Paper, UN/ECE Work Session on Statistical Data Editing, Oslo.
- Norberg, A. (2011), *The Edit*. Working Paper, UN/ECE Work Session on Statistical Data Editing,
- Norberg, A. (2012), *Tree Analysis – A Method for Constructing Edit Groups*. Working Paper, UN/ECE Work Session on Statistical Data Editing, Oslo.
- Pannekoek, J., (2009), *Research on edit and imputation methodology: the throughput programme*, Statistics Netherlands, ISSN: 1572-0314
- Scholtus, S., (2008), *Algorithms for detecting and resolving obvious inconsistencies in business survey data*. DMV-2008-04-25-SSHS, Discussion paper, CBS.
- https://ec.europa.eu/eurostat/cros/sites/crosportal/files/Statistical%20Data%20Editing-05-M-Manual%20Editing%20v1.0_5.pdf
- http://epp.eurostat.ec.europa.eu/portal/page/portal/quality/documents/RPM_EDIMBUS.pdf
- <http://www.cbs.nl/nr/rdonlyres/693e4b18-9322-4ac2-99fd-db61f03637b2/0/200818x10pub.pdf>
- https://ec.europa.eu/eurostat/cros/sites/crosportal/files/Statistical%20Data%20Editing-05-M-Manual%20Editing%20v1.0_5.pdf
- <https://books.google.ro/books?id=WmX-0yU7a5EC&pg=PA400&lpg=PA400&dq=netherlands+Guide+for+data+editing&source=bl&ots=CA0kiX1S8H&sig=4xLNMRCJZFdi4WZAzScmMhXV8CM&hl=en&sa=X&ved=0ahUKEwjQw-Ltw8fLahXm73IKHRx8DowQ6AEIKDAC#v=onepage&q=netherlands%20Guide%20for%20data%20editing&f=false>

IV. ANEXE

Anexa 1. CORELATII PENTRU CONTROLUL DE INTEGRITATE – Ancheta forței de muncă în gospodării – 2016

Controlul între fișierele CL și CI se face la nivel de COD CENTRU, COD LOCUINTA, NR. CL și NR. PERSOANA.
Controlul între fișierele CL, CI și LG se face la nivel de COD CENTRU și COD LOCUINTA.

COD REFUZ	DESCRIERE
RI1:	CENTR și/sau LOC și/sau CL sunt în CI și nu sunt în CL
RI2:	SIL din CL diferita de col 1 (SIL) din LG
RI3:	CL din chestionar CL diferit de col 2 (CL) din LG
RI4:	NTG din CL (I4) diferit de col 3 (NTG) din LG
RI5:	PA CL (I1 cap.2) diferit de cel din col 4 din LG
RI6:	Numărul de persoane cu SEX = 1 și PREZ = 1...9 din CL diferit de col 5 (MASCT) din LG în cadrul aceleiasi locuinte
RI7:	Numărul de persoane cu SEX = 2 și PREZ = 1...9 din CL diferit de col 6 (FEMT) din LG în cadrul aceleiasi locuinte
RI8:	Numarul de persoane cu SEX = 1 și vârsta >= 15 ani și PREZ = 1...9 din CL diferit de col 7 (MASC) din LG în cadrul aceleiasi locuinte
RI9:	Numarul de persoane cu SEX = 2 și vârsta >= 15 ani și PREZ = 1...9 din CL diferit de col 8 (FEM) din LG în cadrul aceleiasi locuinte
RI10:	CENTR și/sau LOC din CL 1 nu este în LG
RI11:	CENTR și/sau LOC din LG nu se află în CL 1
RI12:	NRP cu CI nu se afla în componenta gospodariei (Cap.2) la locuinta și centrul respectiv
RI13:	NRP din CL (Cap.2) cu vârsta >= 15 ani și PREZ = 1...6 nu are CI
RI14:	Persoana NRP din CI are PREZ = 7...9 în CL 2 (sit. incorectă)
RI15:	Să existe toate județele încărcate pentru LG, CL 1, CL 2 și CI.
RI16:	La nivelul fiecărui județ numărul de centre trebuie să fie cel specificat în listă, pentru LG, CL 1, CL 2 și CI
RI18:	Gospodăriile din CL 2 în componența cărora sunt persoane cu vârsta >= 15 ani și PREZ = 1...6 trebuie să apară obligatoriu în CI.
RI19:	Toate locuințele din CL 1 cu PA = 1 trebuie să apară obligatoriu în CL 2.

**Anexa 2. CORELAȚII DE VERIFICARE A CORECTITUDINII “DRUMURILOR”
LOGICE DIN CHESTIONARUL CI – Ancheta forței de muncă în gospodării –
2016**

Dacă	Atunci	și	Altfel	Cod
LUCRM (CI.1) = 1	STAP (CI.5) ≠ 0	ABST(CI.2).SITSP(CI.3).LUCREX(CI.4)=0	R	C7.1
LUCRM (CI.1) = 2	ABST (CI.2) ≠ 0		R	C7.2
ABST (CI.2) = 1.....9	STAP (CI.5) ≠ 0	SITSP (CI.3), LUCREX (CI.4) = 0	R	C7.3
ABST (CI.2) = 10	SITSP (CI.5) ≠ 0		R	C7.4
SITSP (CI.3) = 1, 2.....7	LUCREX (CI.4) ≠		R	C7.6
LUCREX (CI.4) = 1	STAP (CI.5) ≠ 0		R	C7.7
LUCREX (CI.4) = 2	LUCRU (CI. 67)	STAP (CI.5), AGRACT (CI.6), AGRPROP (CI.7),	R	C7.8
STAP (CI.5) = 1	STAPSPEC (CI.8)	AGRACT (CI.6), AGRPROP (CI.7), AGRVZ	R	C7.9
STAP (CI.5) = 2	NPERS (CI. 19)	AGRACT (CI.6), AGRPROP (CI.7), AGRVZ	R	C7.10
STAP (CI.5) = 3,4	AGRACT (CI.6)		R	C7.11
STAP(CI.5) = 5,6	ACT (CI. 20) ≠ 0	AGRACT (CI.6), AGRPROP (CI.7), AGRVZ	R	C7.12
AGRACT (CI.6) = 1,2...7	AGRPROP (CI.7)		R	C7.177
AGRACT (CI.6) = 8	PERS (CI.18) ≠ 0	AGRPROP (CI.7), AGRVZ (CI.8), AGRCONST	R	C7.178
AGRPROP (CI.7) = 1,2	AGRVZ (CI.8) ≠		R	C7.179
AGRPROP (CI.7) = 3	PERS (CI.18) ≠ 0	AGRVZ (CI.8), AGRCONST (CI.9), AGRCONSA	R	C7.180
AGRVZ (CI.8) = 1,2	AGRCONST		R	C7.181
AGRCONST (CI.9) = 1	PERS (CI.18) ≠ 0	AGRCONSA (CI.10), STAPSPEC (CI.11),	R	C7.182
AGRCONST (CI.9) = 2	AGRCONSA		R	C7.183
AGRCONSA (CI.10) = 1,2	PERS (CI.18) ≠ 0	STAPSPEC (CI.11), ANGSAL (CI.12),	R	C7.184
STAPSPEC (CI.11) = 1,2	ANGSAL (CI.12)		R	C7.185
ANGSAL (CI.12) = 1,2,3,4	CONTR (CI.16)	MOTEMP (CI.13), STAPLA (CI.14), ANGT	R	C7.186
ANGSAL (CI.12) = 5,6,7,8,9	MOTEMP		R	C7.187
MOTEMP (CI.13) = 1,2	STAPLA (CI.14)		R	C7.188
MOTEMP (CI.13)= 3,4,5	ANGT (CI.15) ≠	STAPLA (CI.14) = 0	R	C7.18
STAPLA (CI.14) = 1,2	ANGT (CI.15) ≠		R	C7.189
ANGT (CI. 15) = 1.....8	CONTR (CI. 16)		R	C7.19
CONTR (CI. 16) = 1,2	PROP (CI. 17) ≠		R	C7.20
PROP (CI. 17) = 1..... 5	PERS (CI. 18) ≠		R	C7.25
PERS (CI. 18) = 1	NPERS (CI. 19)		R	C7.21
PERS (CI. 18) = 2.....7	ACT (CI. 20) ≠ 0	NPERS (CI. 19) = 0	R	C7.22
DURE (CI. 32) ≠ 0	OSUPLIM (CI.		R	C7.23
NPERS (CI. 19) ≠ 0	ACT (CI. 20) ≠ 0		R	C7.24
ACT (CI. 20) ≠ 0	TARA1 (CI. 21)		R	C7.26
TARA1 (CI. 21) = 1	IUD (CI. 21) ≠ 0	TARA2 (CI. 21)= 0 ∩ REG (CI. 21) = 0	R	C7.27
TARA2 (CI. 21) ≠ 0		TARA1 (CI. 21)= 0 ∩ IUD (CI. 21) = 0		C7.29
TARA1 (CI.21) ≠ 0 ∪	OCUP (CI.22) ≠		R	C7.28
OCUP ≠ 0	COORD (CI. 23)		R	C7.30
COORD (CI. 23) = 1,2,3	LUNIN, ANIN		R	C7.31
LUNIN, ANIN (CI. 24) ≠ 0	MODG (CI. 25)		R	C7.32
MODG (CI. 25) = 1,2	PROG (CI. 26) ≠		R	C7.33
PROG (CI. 26) = 1	EDURO (CI. 29)	MOPARM(CI. 27), CAUZAPP (CI. 28) = 0	R	C7.34
PROG (CI. 26) = 2	MOPARM (CI.		R	C7.35
MOPARM (CI. 27) =	EDURO (CI. 29)		R	C7.36
MOPARM (CI. 27) = 3,4,5	CAUZAPP (CI.		R	C7.37
CAUZAPP (CI. 28) ≠ 0	EDURO (CI. 29)		R	C7.38
EDURO (CI. 29) = 1	DURO (CI. 31) ≠	TOMP (CI. 30) = 0	R	C7.39
EDURO (CI. 29) = 2	TOMP(CI. 30) ≠		R	C7.40
TOMP(CI.30) ≠ 0	DURE (CI. 32) >	DURO (CI. 31) = 0	R	C7.41
DURO (CI. 31) ≠ 0	DURE (CI. 32) >		R	C7.42
OSUPLIM (CI. 33) ≠ 0 și	MOTMU (CI.		R	C7.43
OSUPLIM (CI. 33) ≠ 0 și	LUCDOM (CI.	MOTMU (CI.34), MODUREM (CI. 35),	R	C7.44
OSUPLIM (CI.33) ≠ 0 și	MODUREM (CI.	MOTMU (CI. 34) = 0	R	C7.45
OSUPLIM (CI. 33) ≠ 0 și	MOTNLUC (CI.	MOTMU (CI.34), MODUREM (CI. 35),	R	C7.46

Dacă	Atunci	și	Altfel	Cod
MOTMU (CI. 34) ≠ 0	LUCDOM (CI.	MODUREM (CI. 35), COSMI (CI. 36),	R	C7.47
MODUREM (CI. 35) =	LUCDOM (CI.	COSMI (CI. 36), MOTNLUC (CI. 37), COSMII	R	C7.48
MODUREM (CI. 35) = 4,5,6	COSMI (CI. 36)		R	C7.49
COSMI (CI. 36) ≠ 0	LUCDOM (CI.	MOTNLUC (CI. 37), COSMII (CI.38), LUNSF	R	C7.50
MOTNLUC (CI. 37) = 1,2,3	LUCDOM (CI.	COSMII (CI. 38), LUNSF (CI.39), ANSF (CI.	R	C7.51
MOTNLUC (CI. 37) =	COSMII (CI. 38)		R	C7.52
MOTNLUC (CI. 37) = 6	LUNSF, ANSF	COSMII (CI. 38) = 0	R	C7.53
COSMII (CI. 38) ≠ 0	LUNSF, ANSF		R	C7.54
LUNSF, ANSF (CI. 39) ≠ 0	REV (CI. 40) ≠ 0		R	C7.55
REV (CI. 40) = 1	ABSDURT (CI.		R	C7.56
REV (CI. 40) = 2	CAUTNOCM	ABSDURT (CI. 41), ABSPL (CI. 42), ABSCONT	R	C7.57
ABSDURT (CI. 41) = 1	LUCDOM (CI.	ABSPL (CI. 42), ABSCONT (CI. 43) = 0	R	C7.58
ABSDURT (CI. 41) = 2 și STAP	ABSPL (CI. 42) ≠		R	C7.59
ABSDURT (CI. 41) = 2 și STAP	ABSCONT (CI.	ABSPL (CI. 42) = 0	R	C7.60
ABSDURT (CI. 41) = 2 și STAP	CAUTNOCM	ABSPL (CI. 42), ABSCONT (CI. 43), LUCDOM	R	C7.61
ABSPL (CI. 42) = 1	LUCDOM (CI.	ABSCONT (CI. 43) = 0	R	C7.62
ABSPL (CI. 42) = 2,3	CAUTNOCM	ABSCONT (CI. 43), LUCDOM (CI. 44), LUCSE	R	C7.63
ABSCONT (CI. 43) = 1,2,3,4	LUCDOM (CI.		R	C7.64
ABSCONT (CI. 43) = 5	CAUTNOCM	LUCDOM (CI. 44), LUCSE (CI. 45), LUCNO	R	C7.65
LUCDOM (CI. 44) ≠ 0	LUCSE (CI. 45) ≠		R	C7.66
LUCSE (CI.45) ≠ 0	LUCNO (CI. 46)		R	C7.67
LUCNO (CI. 46) ≠ 0	LUCSAM (CI.		R	C7.68
LUCSAM (CI.42-47) ≠ 0	LUCDUM (CI.43		R	C7.69
LUCDUM (CI.43-48) ≠ 0	SCH (CI.44-49)		R	C7.70
SCH (CI. 49) ≠ 0	TWE (CI.50) ≠ 0		R	C7.71
TWE (CI.50) ≠ 0	TWI (CI.51) ≠ 0		R	C7.190
TWI (CI.51) ≠ 0	TWLOC (CI.52)		R	C7.191
TWLOC (CI.52) ≠ 0	TWDUR (CI.53)		R	C7.192
TWDUR (CI.53) ≠ 0	TWDE (CI.54) ≠		R	C7.193
TWDE (CI.54) = 1,2...5	TWDI (CI.55) ≠	TWDI (CI.55), TWS (CI.56)=0	R	C7.194
TWDE (CI.54) = 6	ASEC (CI.57) ≠		R	C7.195
TWDI (CI.55) ≠ 0	TWS (CI.56) ≠ 0		R	C7.196
TWS (CI.56) ≠ 0	ASEC (CI.57) ≠		R	C7.197
ASEC (CI.45-57) = 1	STAPS (CI. 58) ≠		R	C7.72
ASEC (CI. 57) = 2	DOROREM (CI.	STAPS (CI. 58), PROPS (CI. 59), ACTS (CI.	R	C7.73
STAPS (CI. 58) = 1	PROPS (CI. 59)		R	C7.74
STAPS (CI. 58) = 2,...,6	ACTS (CI. 60) ≠	PROPS (CI. 59) = 0	R	C7.75
PROPS (CI. 59) ≠ 0	ACTS (CI. 60) ≠		R	C7.76
ACTS (CI. 60) ≠ 0	OCUPS (CI. 61)		R	C7.77
OCUPS ≠ 0	DURES (CI. 62)		R	C7.78
DURES (CI. 62) ≠ 0	DOROREM (CI.		R	C7.79
DOROREM (CI. 63) ≠ 0	ORT (CI.64) ≠ 0		R	C7.80
ORT (CI. 64) ≠ 0	CAUTALT (CI.		R	C7.81
CAUTALT (CI. 65) = 1	MOCAUTM (CI.		R	C7.82
CAUTALT (CI. 65) = 2	DISP (CI. 100) ≠	MOCAUTM (CI.66), LUCRU (CI. 67),	R	C7.83
MOCAUTM (CI. 66) ≠ 0	LOCM (CI. 79) ≠	LUCRU (CI. 67), MONLUCRUM (CI. 68),	R	C7.84
LUCRU (CI. 67) = 1	MONLUCRUM		R	C7.85
LUCRU (CI. 67) = 2	CAUTNOCM	INLUCRU (CI. 69), ANII (CI. 70), LUNU (CI.	R	C7.86
MONLUCRUM (CI. 68) ≠ 0	INLUCRU (CI.		R	C7.87
INLUCRU (CI. 69) = 1	ANII (CI. 70) ≠ 0		R	C7.88
INLUCRU (CI. 69) = 2	LUNU, ANU (CI.	ANII (CI. 70) = 0	R	C7.89
ANII (CI. 70) ≠ 0	CAUTNOCM	LUNU (CI. 71), ANU (CI. 71), STAPU (CI. 72),	R	C7.90
LUNU, ANU (CI. 71) ≠ 0	STAPU (CI. 72)		R	C7.91
STAPU (CI. 72) = 1	PROPU (CI. 73)		R	C7.92
STAPU (CI. 72) = 2,...,6	ACTU (CI. 74) ≠	PROPU (CI. 73)= 0	R	C7.93
PROPU (CI. 73) ≠ 0	ACTU (CI. 74) ≠		R	C7.94
ACTU (CI. 74) ≠ 0	OCUPU (CI. 75)		R	C7.95
OCUPU (CI. 75) ≠ 0	CAUTNOCM		R	C7.96

Dacă	Atunci	si	Altfel	Cod
CAUTNOCM (Cl. 76) = 1	LOCM (Cl. 79) ≠	MONCAUTM (Cl. 77), CAUZANE (Cl. 78) = 0	R	C7.97
CAUTNOCM (Cl. 76) = 2	MONCAUTM		R	C7.98
MONCAUTM (Cl.77) = 1, 2.	LOCM (Cl. 79) ≠	CAUZANE (Cl. 78) = 0	R	C7.99
MONCAUTM (Cl.77) = 5, 6	MET1 (Cl.82)	CAUZANE (Cl.78), LOCM (Cl.79), PROGCM	R	C7.100
MONCAUTM (Cl.77) = 7	ILUNC(Cl. 81) ≠	CAUZANE (Cl.78), LOCM (Cl.79), PROGCM	R	C7.101
MONCAUTM (Cl. 77)= 8,9,10	CAUZANE (Cl.		R	C7.102
MONCAUTM (Cl. 77) =		CAUZANE (Cl.78), LOCM (Cl. 79), PROGCM	R	C7.103
CAUZANE (Cl.78) = 1,2,3,4	MET1 (Cl.82)	LOCM (Cl.79), PROGCM (Cl. 80), LUNC (Cl.	R	C7.105
LOCM (Cl.79) = 1	PROGCM		R	C7.106
LOCM (Cl.79) = 2	ILUNC (Cl.81) ≠	PROGCM (Cl.80) = 0	R	C7.107
PROGCM (Cl.80) ≠ 0	ILUNC (Cl.81) ≠		R	C7.108
DURAC (Cl.81) = 1	SITANT (Cl. 97)	MET1 (Cl.82), MET2 (Cl.82), MET3 (Cl.82),	R	C7.109
LUNC (Cl. 81) ≠ 0	ANC(Cl. 81) ≠ 0	DURAC (Cl.81) = 0	R	C7.110
LUNC, ANC, LUNFC, ANFC	MET1 (Cl.82)		R	C7.111
MET 18 (Cl. 96) = 1,2	SITANT (Cl. 97)		R	C7.113
SITANT (Cl.97) ≠ 0	CAUT (Cl. 98) ≠		R	C7.115
CAUT (Cl.98) ≠ 0	DOR (Cl. 99) ≠ 0		R	C7.116
DOR (Cl.99) = 1, 2	DISP (Cl. 100) ≠		R	C7.117
LUNFC (Cl.81) ≠ 0	ANEC (Cl.81) ≠	DURAC (Cl.81) = 0	R	C7.118
DISP (Cl.100) = 1	MOREF (Cl.	MONDISP (Cl. 101) = 0	R	C7.119
DISP (Cl.100) = 2	MONDISP		R	C7.120
MONDISP (Cl.101) = 1.....7	OFMS (Cl. 103)	MOREF (Cl.102) = 0	R	C7.121
MOREF (Cl.102) = 1.....11	OFMS (Cl.103)		R	C7.122
OFMS (Cl.103) = 1	ALOCM (Cl.		R	C7.123
OFMS (Cl.103) = 2	STAPP (Cl. 105)	ALOCM (Cl.104) = 0	R	C7.124
ALOCM (Cl. 104) = 1,2	STAPP (Cl.105)		R	C7.125
STAPP (Cl.105) ≠ 0	STAPPAP ≠ 0		R	C7.126
STAPPAP (Cl.118) = 1 și STAP	STAPAP (Cl.		R	C7.142
STAPPAP (Cl.118) = 2..... 7 și	NETA (Cl.121) ≠	STAPPAP (Cl.119), ACTAP (Cl. 120) = 0	R	C7.143
STAPPAP (Cl.119) ≠ 0	ACTAP (Cl.120)		R	C7.144
STAPPAP (Cl.118) = 1 și STAP	STAPAP	NETA (Cl.121), INETA (Cl. 122),	R	C7.145
STAPPAP (Cl. 118) = 2... 7 și	DIFICIL (Cl. 124)	STAPPAP (Cl.119), ACTAP (Cl.120), NETA	R	C7.146
RASP (Cl. 125) = 1.....3	ZI, LUN, DURM		R	C7.149
OSUPLIM (Cl. 33) = 1	OSUPLIM (Cl.		R	C7.150
OSUPLIM (Cl. 33) 2 ≠ 0	OSUPLIM (Cl.		R	C7.151
OSUPLIM (Cl. 33) 3 ≠ ' '	MOTMU (Cl.		R	C7.152
INETA (Cl.122) = 35, 36	DIFICIL (Cl.124)	LINTR (Cl.123) = 0	R	C7.154
LINTR (Cl.123) ≠ 0	DIFICIL (Cl.112)		R	C7.155
LINTR (Cl.123) = 2	LINTR (Cl.123)		R	C7.156
NETA (Cl.121) = 1	NETA (Cl.121)		R	C7.157
NETA (Cl.121) 1 ≠ 0	LINTR (Cl.123)	INETA (Cl. 122) = 0	R	C7.158
NETA (Cl.121) = 2	INETA (Cl.122)		R	C7.159
INETA (Cl.122) = 1..... 34	LINTR (Cl.123)		R	C7.161
DIFICIL (Cl.124) = 1.....5	RASP (Cl.125) ≠		R	C7.162
MET1 (Cl.82) = 1 ∪ MET2	MET5 (Cl. 83) ≠		R	C7.163
MET 5 (Cl. 83) = 1,2	MET 6 (Cl.84) ≠		R	C7.164
MET 6 (Cl.84) = 1,2	MET 7 (Cl. 85) ≠		R	C7.165
MET 7 (Cl.85) = 1,2	MET 8 (Cl. 86) ≠		R	C7.166
MET 8 (Cl.86) = 1,2	MET 9 (Cl. 87) ≠		R	C7.167
MET 9 (Cl.87) = 1,2	MET 10 (Cl. 88)		R	C7.168
MET 10 (Cl.88) = 1,2	MET 11 (Cl. 89)		R	C7.169
MET 11 (Cl.89) = 1,2	MET 12 (Cl. 90)		R	C7.170
MET 12 (Cl.90) = 1,2	MET 13 (Cl. 91)		R	C7.171
MET 13 (Cl.91) = 1,2	MET 14 (Cl. 92)		R	C7.172
MET 14 (Cl.92) = 1,2	MET 15 (Cl. 93)		R	C7.173
MET 15 (Cl. 93) = 1,2	MET 16 (Cl. 94)		R	C7.174
MET 16 (Cl.94) = 1,2	MET 17 (Cl.95)		R	C7.175
MET 17 (Cl.95) = 1,2	MET 18 (Cl. 96)		R	C7.176

Anexa 3. IMPUTAREA NON-RASPUNSURILOR pentru variabila VNET - Ancheta forței de muncă în gospodării – 2015

Scop: Toți salariații trebuie să aibă o valoare corectă (≥ 724 RON) la VNET.

Înregistrările din fisier pot fi clasificate în:

INDIFERENTE: $CAT \neq A \cup STAP \neq 1 \cup [CAT=A \cap STAP=1 \cap MOTNLUC=4 \cap (VNET \geq 724)]$ adică inactivi sau non-salariați sau salariați dar în concediu pentru creștere a copilului care au declarat un salariu corect

DONORI: $CAT=A \cap STAP=1 \cap VNET \geq 724 \cap MOTNLUC \neq 4$ adică salariați pentru care am o valoare corectă la VNET și care nu sunt în concediu pentru creștere a copilului⁸

PRIMITORI: $CAT=A \cap STAP=1 \cap [0 \leq VNET < 724]$ adică salariați pentru care VNET este incorect

Reuniunea celor 3 categorii = total înregistrări

Pentru înregistrările care sunt DONORI sau INDIFERENȚI nu se face nimic. Pentru înregistrările care sunt PRIMITORI (adică ar trebui să aibă valori pentru VNET dar nu au sau valorile sunt incorecte) se caută între înregistrările care au o valoare corectă pentru VNET (DONORI) un "geamăn" de la care să se preia datele pentru câștiguri.

"Geamănul" este căutat folosind mai întâi criteriul cel mai strict (aceiași sex, aceeași grupa de vârstă, același nivel de instruire, aceeași activitate, aceeași ocupație, aceeași regiune) și care are același program de lucru (PROG) și care nu și-a mai donat anterior valoarea unui PRIMITOR. Dacă se găsește un "geamăn", PRIMITORUL preia datele acestuia, prin procedura de imputare 1 și se trece la următoarea înregistrare. Dacă nu găsește, se continuă cautarea folosind același criteriu dar făcând abstracție de programul de lucru (PROG). Dacă se găsește un "geamăn", PRIMITORUL preia datele acestuia, prin procedura de imputare 2 și se trece la următoarea înregistrare. Dacă nu găsește, se continuă cautarea folosind un criteriu mai relaxat (aceiași sex, aceeași grupa de vârstă, aceeași instruire – doar pe cele 3 categorii superior, mediu, scăzut, aceeași activitate, aceeași ocupație) și care are același program de lucru (PROG) și care nu și-a mai donat anterior valoarea unui PRIMITOR. Dacă se găsește un "geamăn", PRIMITORUL preia datele acestuia, prin procedura de imputare 1. Se procedează în felul acesta, aplicând criteriile din ce în ce mai relaxate și cautând prima dată printre DONORII care îndeplinesc criteriul și au același program de lucru și abia apoi (dacă nu găsește) făcând abstracție de programul de lucru, până când printre DONORI se găsește un "geamăn" care nu și-a mai donat anterior valoarea unui PRIMITOR și de la care se pot prelua datele pentru câștiguri

Criteriile de căutare sunt ierarhizate, de la 1 la 17, criteriul 1 fiind cel mai strict. Fiecare dintre criteriile 2, 3, ... este mai relaxat decât cel anterior, în sensul că fie se renunță la precizie pentru una dintre variabilele după care se face căutarea (ex. la criteriul 3 "nivelul de instruire" devine "instruire" – superior/mediu/scăzut, la criteriul 5 "grupa minoră de ocupații" devine "grupa majoră de ocupații" ș.a.m.d) fie se renunță la una dintre variabilele după care se face potrivire (ex. la criteriul 3 "regiunea" nu mai apare, la criteriul 9 nu mai apare "grupa de ocupații" ș.a.m.d). Criteriul 17 cel mai relaxat (practic, la aplicarea criteriului 17 este imposibil să nu se găsească un DONOR care nu și-a mai donat anterior valoarea unui PRIMITOR și de la care să se poată prelua datele).

Criteriile folosite pentru a identifica un "geamăn":

CRITERIUL 1:

SEX

GRUPA DE VARSTA_1

NIVEL DE INSTRUIRE

ACTIVITATE

GRUPA MINORA DE OCUPATII

REGIUNE

CRITERIUL 2:

⁸ Aceste persoane, fiind în absența de lungă durată de la locul de muncă au declarat ultimul salariu primit (posibil cu mai mult de 1 an în urmă) → nu sunt folosite de DONORI

SEX
GRUPA DE VARSTA_1
INSTRUIRE
ACTIVITATE
GRUPA MINORA DE OCUPATII
REGIUNE
CRITERIUL 3:
SEX
GRUPA DE VARSTA_1
NIVEL DE INSTRUIRE
ACTIVITATE
GRUPA MINORA DE OCUPATII
CRITERIUL 4:
SEX
GRUPA DE VARSTA_1
INSTRUIRE
ACTIVITATE
GRUPA MINORA DE OCUPATII
CRITERIUL 5:
SEX
GRUPA DE VARSTA_1
NIVEL DE INSTRUIRE
ACTIVITATE
GRUPA MAJORA DE OCUPATII
CRITERIUL 6:
SEX
GRUPA DE VARSTA_1
INSTRUIRE
ACTIVITATE
GRUPA MAJORA DE OCUPATII
CRITERIUL 7:
SEX
GRUPA DE VARSTA_1
NIVEL DE INSTRUIRE
GRUPA DE ACTIVITATE
GRUPA MAJORA DE OCUPATII
CRITERIUL 8:
SEX
GRUPA DE VARSTA_1
INSTRUIRE
GRUPA DE ACTIVITATE
GRUPA MAJORA DE OCUPATII
CRITERIUL 9:
SEX
GRUPA DE VARSTA_1
NIVEL DE INSTRUIRE
ACTIVITATE
CRITERIUL 10:
SEX
GRUPA DE VARSTA_1

INSTRUIRE

ACTIVITATE

CRITERIUL 11:

SEX

GRUPA DE VARSTA_1

NIVEL DE INSTRUIRE

GRUPA DE ACTIVITATE

CRITERIUL 12:

SEX

GRUPA DE VARSTA_1

INSTRUIRE

GRUPA DE ACTIVITATE

CRITERIUL 13:

SEX

GRUPA DE VARSTA_1

NIVEL DE INSTRUIRE

CRITERIUL 14:

SEX

GRUPA DE VARSTA_1

INSTRUIRE

CRITERIUL 15:

SEX

GRUPA DE VARSTA_1

CRITERIUL 16

SEX

GRUPA DE VARSTA_2

CRITERIUL 17:

SEX

GRUPA DE VARSTA_3

unde:

- GRUPA DE VARSTA_1 → VARSTA

15-19 ani

20-24 ani

25-29 ani

30-34 ani

35-39 ani

40-44 ani

45-49 ani

50-54 ani

55-59 ani

60-64 ani

65 ani si peste

GRUPA DE VARSTA_2 → VARSTA:

15-24 ani

25-34 ani

35-44 ani

45-54 ani

55-64 ani

65 ani si peste

- GRUPA DE VARSTA_3 → VARSTA

15-24 ani

25-54 ani

55 ani si peste

NIVELUL DE INSTRUIRE → NIVS

INSTRUIRE → NIVS

superior → 1, 2, 3, 4

mediu → 5, 6, 7, 8, 9, 10

scazut → 11, 12, 13, 14, 15, 16

ACTIVITATE → ACT

GRUPA DE ACTIVITATE → ACT - grupele din publicație (atentie – vezi grupele corespunzatoare CAEN Rev 2)

GRUPA MINORA DE OCUPATII → OCUP - primele 3 caractere

GRUPA MAJORA DE OCUPATII → OCUP - primul caracter (grupele din publicație)

Pentru a avea o evidență a înregistrărilor asle căror valori au fost imputate, a celor care au fost folosite ca donori si a criteriilor de potrivire folosite, pentru fiecare înregistrare se vor adauga 3 noi variabile:

- IMPUTAT - se macheaza in cazul în care pentru înregistrarea în cauza valorile pentru castiguri au fost imputate
– adica inregistrarea este PRIMITOR si s-a realizat efectiv imputarea.

- DONOR - se macheaza in cazul în care înregistrarea în cauza este DONOR și a donat efectiv valori

-CRITERIU – numai pentru PRIMITORI pentru care s-a realizat efectiv imputarea- se atribuie valori de la 1 la 17 corezpunzator criteriului (1, 2, 3...,17) care s-a folosit pentru realizarea efectiva a imputarii.

Modul de lucru:

Pentru fiecare înregistrare i din fisier (i de la 1 la n, unde n este numărul de înregistrari din fisier)

1. dacă i este DONOR sau INDIFERENT → nimic de făcut și treci la urmatoarea înregistrare

1. dacă este i este PRIMITOR → există o înregistrare j care este DONOR și $CRITERIUL1j=CRITERIUL1i$ și $PROGj=PROGi$ și nu și-a mai donat anterior valoarea ?

2. dacă da → PROCEDURA DE IMPUTARE_1 și treci la urmatoarea înregistrare

2. dacă nu → există o înregistrare j care este DONOR și $CRITERIUL1j=CRITERIUL1i$ și nu și-a mai donat anterior valoarea ?

3. dacă da → PROCEDURA DE IMPUTARE_2 și treci la urmatoarea înregistrare

3. dacă nu → există o înregistrare j care este DONOR și $CRITERIUL2j=CRITERIUL2i$ și $PROGj=PROGi$ și nu și-a mai donat anterior valoarea ?

4. dacă da → PROCEDURA DE IMPUTARE_1 și treci la urmatoarea înregistrare

4. dacă nu → există o înregistrare j care este DONOR și $CRITERIUL2j=CRITERIUL2i$ și nu și-a mai donat anterior valoarea ?

5. dacă da → PROCEDURA DE IMPUTARE_2 și treci la urmatoarea înregistrare

5. dacă nu → există o înregistrare j care este DONOR și $CRITERIUL3j=CRITERIUL3i$ și $PROGj=PROGi$ și nu și-a mai donat anterior valoarea ?

6. dacă da → PROCEDURA DE IMPUTARE_1 și treci la urmatoarea înregistrare

6. dacă nu → există o înregistrare j care este DONOR și $CRITERIUL3j=CRITERIUL3i$ și nu și-a mai donat anterior valoarea ?

7. dacă da → PROCEDURA DE IMPUTARE_2 și treci la urmatoarea înregistrare

7. dacă nu → există o înregistrare j care este DONOR și $CRITERIUL4j=CRITERIUL4i$ și $PROGj=PROGi$ și nu și-a mai donat anterior valoarea?

8. dacă da → PROCEDURA DE IMPUTARE_1 și treci la urmatoarea înregistrare

8. dacă nu → există o înregistrare j care este DONOR și $CRITERIUL4j=CRITERIUL4i$ și nu și-a mai donat anterior valoarea?

9. dacă da → PROCEDURA DE IMPUTARE_2 și treci la urmatoarea înregistrare

9. dacă nu → există o înregistrare j care este DONOR și CRITERIUL5j=CRITERIUL5i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
10. dacă exista → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
10. dacă nu → există o înregistrare j care este DONOR și CRITERIUL5j=CRITERIUL5i și nu și-a mai donat anterior valoarea?
11. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
11. dacă nu → există o înregistrare j care este DONOR și CRITERIUL6j=CRITERIUL6i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
12. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
12. dacă nu → există o înregistrare j care este DONOR și CRITERIUL6j=CRITERIUL6i și nu și-a mai donat anterior valoarea?
13. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
13. dacă nu → există o înregistrare j care este DONOR și CRITERIUL7j=CRITERIUL7i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
14. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
14. dacă nu → există o înregistrare j care este DONOR și CRITERIUL7j=CRITERIUL7i și nu și-a mai donat anterior valoarea?
15. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
15. dacă nu → există o înregistrare j care este DONOR și CRITERIUL8j=CRITERIUL8i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
16. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
16. dacă nu → există o înregistrare j care este DONOR și CRITERIUL8j=CRITERIUL8i și nu și-a mai donat anterior valoarea?
17. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
17. dacă nu → există o înregistrare j care este DONOR și CRITERIUL9j=CRITERIUL9i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
18. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
18. dacă nu → există o înregistrare j care este DONOR și CRITERIUL9j=CRITERIUL9i și nu și-a mai donat anterior valoarea?
19. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
19. dacă nu → există o înregistrare j care este DONOR și CRITERIUL10j=CRITERIUL10i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
20. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
20. dacă nu → există o înregistrare j care este DONOR și CRITERIUL10j=CRITERIUL10i și nu și-a mai donat anterior valoarea?
21. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
21. dacă nu → există o înregistrare j care este DONOR și CRITERIUL11j=CRITERIUL11i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
22. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
22. dacă nu → există o înregistrare j care este DONOR și CRITERIUL11j=CRITERIUL11i și nu și-a mai donat anterior valoarea?
23. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
23. dacă nu → există o înregistrare j care este DONOR și CRITERIUL12j=CRITERIUL12i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
24. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
24. dacă nu → există o înregistrare j care este DONOR și CRITERIUL12j=CRITERIUL12i și nu și-a mai donat anterior valoarea?
25. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare

- 25 . dacă nu → există o înregistrare j care este DONOR și CRITERIUL13j=CRITERIUL13i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
26. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
26. dacă nu → există o înregistrare j care este DONOR și CRITERIUL13j=CRITERIUL13i și nu și-a mai donat anterior valoarea?
27. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
- 27 . dacă nu → există o înregistrare j care este DONOR și CRITERIUL14j=CRITERIUL14i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
28. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
28. dacă nu → există o înregistrare j care este DONOR și CRITERIUL14j=CRITERIUL14i și nu și-a mai donat anterior valoarea?
29. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
- 29 . dacă nu → există o înregistrare j care este DONOR și CRITERIUL15j=CRITERIUL15i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
30. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
30. dacă nu → există o înregistrare j care este DONOR și CRITERIUL15j=CRITERIUL15i și nu și-a mai donat anterior valoarea?
31. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
- 31 . dacă nu → există o înregistrare j care este DONOR și CRITERIUL16j=CRITERIUL16i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
32. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
32. dacă nu → există o înregistrare j care este DONOR și CRITERIUL16j=CRITERIUL16i și nu și-a mai donat anterior valoarea?
33. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
- 33 . dacă nu → există o înregistrare j care este DONOR și CRITERIUL17j=CRITERIUL17i și PROGj=PROGi și nu și-a mai donat anterior valoarea?
34. dacă da → PROCEDURA DE IMPUTARE_1 și treci la următoarea înregistrare
34. dacă nu → există o înregistrare j care este DONOR și CRITERIUL17j=CRITERIUL17i și nu și-a mai donat anterior valoarea?
35. dacă da → PROCEDURA DE IMPUTARE_2 și treci la următoarea înregistrare
- 35 . dacă nu → nimic de făcut și treci la următoarea înregistrare

OBS. Dacă la o căutare se găsesc mai multe înregistrări care satisfac criteriile de cautare, datele se preiau de la prima înregistrare găsită.

Anexa 4. Imputarea non-răspunsurilor în EU-SILC

Problemele apar în special în determinarea venitului total al gospodăriei din cauza lipsei informațiilor asupra unor componente de venit, dar pot să apară și când nu există toate informațiile pentru toți membrii dintr-o gospodărie.

Există două motive majore pentru imputarea datelor:

- din punct de vedere *statistic* se impune imputarea pentru minimizarea abaterii medii pătratice a estimărilor din anchetă, în particular pentru deplasarea (bias) de non-răspuns când lipsa datelor nu este întâmplătoare;
- din punct de vedere *practic* se impune imputarea atunci când nu există coerență între rezultate din analize diferite.

Lipsa datelor în EU-SILC

Erori legate de acoperire și de selectare a eșantionului

Acestea apar atunci când unitățile din populația țintă nu sunt reprezentative în planul de eșantionare sau când probabilitățile de selectare a unităților sunt distorsionate. Câteva corecții pot fi posibile pe baza informațiilor externe planului de eșantionare. Astfel de ajustări sunt numite post-stratificări, calibrări, analiză comparativă etc.

Unitatea non-răspuns

Se referă la absența informațiilor pentru toate unitățile (gospodării și/ sau persoane) selectate din eșantion. De regulă, impactul unităților non-răspuns este redus prin atașarea ponderilor adecvate la cazurile care au răspuns.

Unitatea non-răspuns parțială

În EU-SILC există două nivele de analiză: gospodării și persoane. Analiza care implică distribuția unităților la oricare dintre cele două nivele, non-răspunsul poate fi tratat prin ponderare. O caracteristică aparte a EU-SILC este aceea că un număr de variabile la nivel de gospodărie nu sunt colectate în mod direct, ci sunt construite prin agregarea informațiilor culese la nivel de persoană, de fiecare membru al gospodăriei. Termenul de unitate non-răspuns parțială este introdus pentru a descrie situația în care doar câțiva, nu toți membrii gospodăriei selectați în anchetă au fost anchetați cu succes. Există două posibilități pentru a rezolva această problemă:

- ajustarea ponderii eșantionului persoanelor anchetate din gospodărie în scopul de a compensa lipsa celorlalți membrii din gospodărie;
- construirea variabilelor pentru fiecare persoană neanchetată din gospodărie prin imputare.

Item non-răspuns

Se referă la situația în care unitățile din eșantion au fost anchetate cu succes, însă nu au fost obținute toate informațiile. În anumite situații, când incidența de non-răspuns este infimă este de preferat ca aceste cazuri să fie ignorate și să se facă analiza doar pentru cazurile care au informații complete.

Tratarea datelor lipsă

Date lipsă	Măsuri de compensare	
Erori legate de acoperire și selectare a eșantionului	Bench-marking (analiză comparativă), post-stratificare, calibrare	
Unitatea non-răspuns	Bench-marking (analiză comparativă), post-stratificare, calibrare	Ponderare
Unitatea non-răspuns parțială	Ponderare	Imputare
Item non-răspuns	Imputare	Ponderare
Variabile țintă (componentele venit brut)		Modelare

Construirea variabilelor țintă de venit

În EU-SILC componentele de venit trebuie să fie în forma brută. Componentele care sunt disponibile (colectate) în formă netă trebuie să fie convertite în formatul brut cerut.

În conformitate cu Regulamentul (CE) nr. 1177/2003 al Parlamentului European și al Consiliului din 16 iunie 2003 privind Statisticile comunitare referitoare la venit și la condițiile de viață (EU-SILC), Statele Membre transmit Comisiei (Eurostat), sub formă de fișiere de microdate, datele transversale și longitudinale ponderate, verificate integral, editate și imputate în funcție de venit.

Regulamentul (CE) nr. 1982/2003 al Comisiei din 21 octombrie 2003 de punere în aplicare a Regulamentului (CE) nr. 1177/2003 al Parlamentului European și al Consiliului privind statisticile comunitare referitoare la venit și la condițiile de viață (EU-SILC), în ceea ce privește regulile de eșantionare și urmărire prevede:

1. În cazul în care non-răspunsul la variabilele venitului la nivelul componentei are drept rezultat date lipsă, se aplică metodele de imputare statistică corespunzătoare.
2. În cazul în care o variabilă a venitului brut la nivelul componentei nu este culeasă direct, se aplică metode de imputare statistică și/sau modelare corespunzătoare pentru a obține variabilele țintă necesare.
3. În cazul în care apare un non-răspuns la un chestionar individual într-o gospodărie eșantion, se utilizează proceduri statistice de ponderare și/sau imputare corespunzătoare pentru estimarea venitului total al gospodăriei.
4. Factorii de ponderare se calculează după cum este necesar pentru a se lua în considerare probabilitatea de selectare a unităților, non-răspunsurile și, dacă este cazul, pentru a se ajusta eșantionul la datele externe referitoare la distribuția gospodăriilor și a persoanelor în populația țintă, de exemplu după sex, vârstă (grupe de vârstă de cinci ani), dimensiune și compoziție a gospodăriei și regiune (nivel NUTS II) sau referitoare la datele privind venitul din alte surse naționale, în cazul în care Statele Membre în cauză consideră că datele externe respective sunt suficient de fiabile.
5. Statele Membre furnizează Comisiei (Eurostat) toate informațiile necesare privind organizarea și metodologia anchetei și, în special, indică criteriile adoptate în alegerea planului de eșantionare și a mărimii eșantionului.

Anexa 5. DESCRIEREA UNEI CERCETĂRI STATISTICE SELECTIVE - Ancheta forței de muncă în gospodării

1. Obiectivele anchetei

Obiectivul principal al Anchetei forței de muncă în gospodării (AMIGO) îl constituie asigurarea informațiilor necesare pentru evaluarea situației existente pe piața forței de muncă din România, măsurarea dimensiunilor și evoluției fenomenelor de ocupare, șomaj și inactivitate.

Începând cu anul 1996, ancheta forței de muncă în gospodării se realizează trimestrial, ca o cercetare continuă, permițând astfel obținerea de date conjuncturale asupra mărimii și structurii ofertei de forță de muncă și evidențierea fenomenelor cu caracter sezonier care se manifestă pe piața forței de muncă. Metodologia anchetei este armonizată cu standardele Uniunii Europene, respectiv cu ancheta europeană Labour Force Survey (LFS) iar rezultatele sunt comparabile cu cele din statele membre.

Cercetarea statistică se realizează în conformitate cu normele europene, respectiv cu Regulamentul (CE) nr. 577/98 al Consiliului privind organizarea unei anchete prin sondaj asupra forței de muncă din Comunitate, cu modificările ulterioare, Regulamentul (UE) nr. 545/2014 al Parlamentului European și al Consiliului de modificare a Regulamentului (CE) nr. 577/98 al Consiliului privind organizarea unei anchete prin sondaj asupra forței de muncă din Comunitate, Regulamentul (CE) nr. 377/2008 al Comisiei de punere în aplicare a Regulamentului (CE) nr. 577/98., în ceea ce privește codificarea utilizată pentru transmiterea datelor începând cu anul 2009 și utilizarea unui subeșantion pentru culegerea datelor referitoare la variabilele structurale și definirea trimestrelor de referință, cu modificările ulterioare.

2. Sfera de cuprindere

Ancheta se desfășoară pe întreg cuprinsul țării.

Doar gospodăriile individuale din locuințele permanente sunt intervievate. Unitățile de locuit în comun (cămine de bătrâni, de handicapați, cămine muncitorești, sanatorii etc.) și persoanele care locuiesc permanent în astfel de unități nu sunt cuprinse în anchetă. De asemenea nu sunt cuprinse în anchetă locuințele sezoniere.

Fac obiectul anchetei persoanele rezidente - temporar sau permanent - în România, membre ale gospodăriilor din locuințele selectate. Se consideră membri ai gospodăriei și persoanele plecate din localitate pentru o perioadă mai mare de 6 luni, care se află în țară sau străinătate⁹, dacă acestea păstrează legături familiale cu gospodăria din care fac parte, precum: militarii în termen, elevii și studenții plecați la studii, persoanele plecate la lucru, deținuții și arestații, persoanele spitalizate sau aflate temporar în sanatorii pentru tratament sau recuperare.

3. Unitatea de observare

Unitatea de observare este persoana.

4. Periodicitatea și perioada de referință

Ancheta se realizează cu periodicitate trimestrială, iar rezultatele se prezintă trimestrial și anual conform regulamentelor europene.

⁹⁾ Dacă durata absenței din țară este mai mică de 12 luni.

5. Metoda și perioada de înregistrare a datelor

Datele sunt colectate prin metoda interviului față-în-față. Înregistrarea informațiilor în chestionarele anchetei se realizează prin interviuarea persoanelor de 15 ani și peste, la domiciliul gospodăriilor din locuințele cercetate.

Interviurile sunt repartizate uniform de-a lungul trimestrului, ancheta realizându-se ca o cercetare continuă.

6. Planul de sondaj

Planul de sondaj folosit pentru ancheta AMIGO este un plan de sondaj în două trepte: construirea, în **prima treaptă**, a Eșantionului Multifuncțional de Zone Teritoriale (eșantionul “master” EMZOT’); în **a doua treaptă**, au fost selectate sistematic, din EMZOT, clusteri (grup de trei locuințe), eșantionul final aferent unui trimestru constând în 28080 locuințe¹⁰. Toate gospodăriile dintr-o locuință sunt incluse.

Eșantionul este reprezentativ la nivelul țării și pe regiuni. Reprezentativitatea se referă atât la structura gospodăriilor, cât și la distribuția populației pe medii, sexe și grupe de vârstă. Diferențele de structură, care apar datorită situației din teren la momentul realizării anchetei, sunt anulate prin aplicarea unor procedee de ajustare în funcție de rata de non-răspuns pe medii de rezidență și în funcție de distribuția populației pe medii, sexe și grupe de vârstă, distribuții obținute din surse exhaustive de cercetare demografică.

Eșantionul este construit pe baza unui procedeu de înnoire parțială a eșantionului trimestrial (“schema rotațională 2-2-2”), care are ca principiu de bază următoarea tehnică: o locuință este cercetată două trimestre consecutive, este scoasă temporar din cercetare următoarele două trimestre, este reintrodusă în cercetare următoarele două trimestre, apoi este scoasă definitiv din cercetare. Așadar, o locuință este gestionată 6 trimestre, fiind interviuată de 4 ori. Conform schemei aplicate, în fiecare trimestru se păstrează în eșantion 50% din locuințele trimestrului anterior, 25% din locuințele cercetate cu două trimestre în urmă, iar restul de 25% sunt locuințe absolut noi. Astfel se păstrează și o acoperire de 50% a eșantioanelor de la un trimestru al unui an la același trimestru al anului precedent.

6.1. Baza de sondaj

În lipsa unor registre adecvate (registru de locuințe, registru al populației etc), Ancheta forței de muncă în gospodării se bazează pe utilizarea unui eșantion master, ceea ce impune aplicarea unor planuri de sondaj multistadiale. Eșantionului Multifuncțional de Zone Teritoriale (eșantionul “master” EMZOT) este constituit din 780¹⁰ centre de cercetare (unități primare de eșantionare), repartizate în toate județele și sectoarele Municipiului București. Începând cu anul 2015 se utilizează eșantionul master EMZOT realizat pe baza datelor de la Recensământului Populației și Locuințelor din anul 2011.

6.2. Unitatea de selecție

Unitatea primară de eșantionare, corespunzătoare primei trepte de eșantionare (selecția eșantionului master), a fost un grup de secții de recensământ.

Unitatea secundară de eșantionare, corespunzătoare celei de a doua trepte de eșantionare (selecția eșantionului anchetei), a fost locuința.

6.3. Variabilele de stratificare

¹⁰ Eșantionul anchetei a cuprins 18036 locuințe (din 501 centre de cercetare) pe trimestru – până în anul 2003 și 28080 locuințe (din 780 centre de cercetare) pe trimestru – în perioada 2004-2014.

EMZOT este un eşantion stratificat. Criteriile de stratificare au fost județul și mediul de rezidență, obținându-se 88 de straturi.

6.4. Mărimea și alocarea eşantionului

Volumul eşantionului este de 28080 locuințe/trimestru (112320 locuințe pe an).

Conform metodologiei adoptate pentru anchetele în gospodării realizate de către INS-România, o locuință poate conține una sau mai multe (cazuri rare, totuși) gospodării. Toate gospodăriile aparținând locuințelor selectate, precum și toate persoanele de 15 ani și peste, aparținând gospodăriilor din locuințele selectate, sunt anchetate.

În aceste condiții, volumul eşantionului de gospodării, ca și volumul eşantionului de persoane, sunt variabile aleatoare, depinzând de eşantionul de locuințe selectat în treapta a doua.

6.5. Extragerea eşantionului

Pentru extragerea unităților primare, în interiorul fiecărui strat s-a utilizat metoda de extragere echilibrată, folosind macro SAS Cube.

6.6. Extinderea rezultatelor

Extinderea rezultatelor obținute din anchetă, la nivelul întregii țări, se realizează pe baza ponderilor atribuite persoanelor din gospodăriile care fac parte din eşantion și au răspuns la interviu. Pentru determinarea acestor coeficienți se parcurg următoarele etape:

- calculul ponderilor de bază: ponderea de bază atribuită unei locuințe reprezintă inversul probabilității generale de includere a locuinței în eşantionul anchetei; toate gospodăriile dintr-o locuință "împrumută" ponderea de bază a locuinței respective;
- tratarea non-răspunsurilor totale: se realizează cu ajutorul metodei grupelor de răspuns omogen, utilizând ca variabile explicative județul și mediul de rezidență; în această etapă, ponderile de bază ale gospodăriilor repondente sunt ajustate cu inversul ratei de răspuns;
- redresarea eşantionului și calculul ponderilor finale: redresarea este realizată în scopul de a îmbunătăți calitatea estimațiilor printr-o ajustare finală a ponderilor în etapa precedentă: metoda de redresare folosită este cunoscută sub numele de calibrare. Calibrarea se realizează la nivel de regiune de dezvoltare utilizând ca variabile populația pe sexe și grupe de vârstă, mediul de rezidență și numărul total de gospodării. Utilizarea acestei metode conduce la creșterea gradului de precizie al estimațiilor.

7. Chestionarul anchetei

*Informațiile sunt colectate pe chestionare identice pe întreaga perioadă a anului. Colectarea datelor se face utilizând trei chestionare statistice și anume: **CL** - chestionarul locuinței; **CI** - chestionarul individual.*

8. Clasificări utilizate

Ocupația: se definește și codifică conform Clasificării Ocupațiilor din România (COR 2008) armonizată cu Clasificarea internațională standard a ocupațiilor ISCO-08¹¹⁾.

Activitatea: se definește și codifică conform Clasificării Activităților din Economia Națională¹²⁾ (CAEN Rev.2) armonizată cu clasificarea europeană a activităților economice NACE Rev.2.

¹¹⁾ Anterior anului 2011 s-a utilizat clasificarea COR armonizată cu clasificarea internațională standard a ocupațiilor ISCO-COM (88).

Profil teritorial: se codifică conform criteriilor Regulamentului 1059/2003, privind stabilirea unei clasificări comune a unităților teritoriale statistice și a reglementărilor europene emise de EUROSTAT, corespunzătoare *Nomenclatorului Unităților Teritoriale pentru Statistică (NUTS)*.

Statutul profesional: se definește și codifică în conformitate cu clasificarea internațională **ICSE-93**.

Nivelul de instruire: gruparea datelor după nivelul de instruire absolvit s-a făcut având în vedere: nivelul de învățământ absolvit, corespondența între nivelurile de învățământ (stabilite conform legislației naționale) și nivelurile de educație definite conform Clasificării Internaționale Standard a Educației în vigoare la data respectivă ¹³⁾.

În publicații și alte medii de diseminare, datele privind nivelul de instruire absolvit pot fi grupate și sub forma: **scăzut:** fără școală absolvită, primar, gimnazial; **mediu:** liceal (ciclul superior sau inferior), profesional, complementar sau de ucenici, postliceal de specialitate sau tehnic de maiștri; **superior:** universitar de scurtă durată (colegii universitare, secții de subingineri/conducători arhitecți ale instituțiilor de învățământ superior) și de lungă durată (învățământ universitar de lungă durată, licență și masterat), postuniversitar, doctorat, postdoctorat.

¹²⁾ Anterior anului 2011, datele au fost colectate și diseminate astfel: anterior anului 2003 - conform CAEN armonizat cu NACE Rev.1, 2003-2007 conform CAEN Rev.1 - armonizat cu NACE Rev.1.1; în anul 2008 datele au fost colectate și diseminate în dublă clasificare CAEN Rev.1 și CAEN Rev 2

¹³⁾ ISCED97 – până în anul 2013 (inclusiv) și ISCED 2011 începând cu anul 2014