

## Methodological aspects of data quality modeling

Otilia PÎRLOG  
Ministerul Apărării Naționale

*Even though a fairly new field of economical science, data quality provides approaches that are both comprehensive and thorough. In order to state quality requirements and evaluate the final results, it is required to define the axiomatic system and introduce the objective and subjective aspects of the analyzed concept. The quality levels of data are also defined: attribute characteristics, dimensions and category. Thus, by investigating the quality levels, the category to which the deficiencies belong to is relieved. Improving the quality is associated with systems or processes belonging to the data category levels, the effort is focused towards the determination of the changes in the data acquisition and / or processing. A very important aspect in this field is represented by the data usage premises, the ultimate measure of the quality being the level of fulfillment of the end user's requirements.*

**Keywords:** data, quality, deficiency, usage, attribute, characteristic, dimension, category.

### Atributele de calitate ale datelor

Proiectarea conceptuală a colecțiilor de date vizează problematica aplicației informatice, transpusă la nivelul entităților și a formalizării relațiilor dintre acestea. Maniera de abordare tradițională a proiectării colecțiilor de date nu recurge la incorporarea explicită a aspectelor de calitate.

Următoarele problematice ale utilizării datelor conduc la concentrarea atenției asupra calității acestora:

- folosirea nu numai de către aplicația pentru care au fost inițial proiectate;
- procesarea împreună cu date provenite din alte colecții;
- punerea datelor la dispoziția unor utilizatori noi, nefamiliarizați cu acestea.

Una din definițiile calității datelor unanim acceptată este conformitatea la cerințele exprimate de utilizatori. Din punct de vedere operațional, aceasta este definită prin intermediul parametrilor și indicatorilor corespunzători.

Un *indicator* este o dimensiune a calității da-

telor prin care se furnizează informații obiective despre aceasta. Astfel de indicatori pot fi: sursa, timpul de creare, metoda de colectare etc.

Scopul modelării calității datelor este acela de a permite utilizatorilor să-și definească cerințele privind calitatea datelor și să stabilească indicatorii de calitate care sunt adecvați pentru o aplicație dată.

*Valoarea unui indicator* este o caracteristică măsurabilă a datelor stocate.

*Valoarea unui parametru* reprezintă valoarea determinată pentru un parametru de calitate, bazată pe interpretarea valorii indicatorilor aferenți.

Translatarea valorii indicatorilor în valori parametrice este subiectivă și se realizează prin funcții utilizator, expresia acestora fiind dependentă de context.

Un *atribut* este un termen bidimensional, care pune în corelație parametrii și indicatorii calitativi. Relația dintre atribute, parametri și indicatori este prezentată în Figura 1.

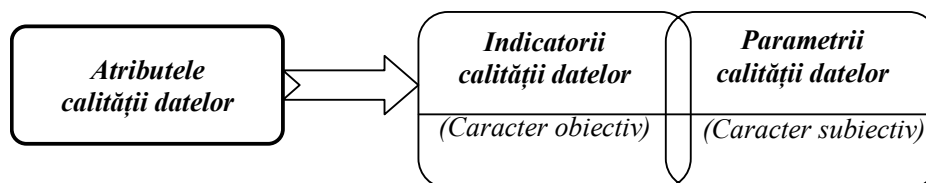


Fig.1. Definierea atributelor de calitate ale datelor

*Cerințele calității datelor* specifică acei indicatori a căror valoare este necesar să se înscrie între anumite limite specificate, astfel încât să se asigure punerea la dispoziția utilizatorilor a unor date cu un nivel calitativ acceptabil.

*Administratorul datelor* este o persoană (sau sistem) a cărei responsabilitate este asigurarea îndeplinirii de către datele stocate a cerințelor exprimate.

În general, utilizatori diferiți au cerințe diferite de calitate, iar calitatea unor date diferite poate varia între limite extrem de largi. De aici decurg și cele două categorii de premise de studiu referitoare la calitatea datelor, respectiv modelarea și utilizarea datelor.

### **Premise privind modelarea calității**

Modelarea calității datelor este o extensie a metodologiilor tradiționale de modelare a datelor. Diferența dintre acestea constă în faptul că în timp ce modelarea datelor se adresează structurii și semanticii datelor, modelarea calității datelor vizează aspectele structurale și semantice aferente calității datelor.

#### Premisa a.1. Relația dintre atributele aplicației și ale calității

Nu este întotdeauna necesar ca atributele aplicației și ale calității să fie distincte. Pot fi identificate două domenii diferite de activitate din acest punct de vedere: utilizarea datelor și administrarea calității.

Astfel, valoarea unui indicator poate să fie folosită în contextul procesării datelor pentru obținerea unor informații livrabile la utilizator sau care să fie trecute într-o nouă etapă de procesare. Aceeași valoare poate fi interpretată printr-o funcție dedicată pentru stabilirea nivelului de dezvoltare al unui parametru calitativ.

Unele informații cuprind detalii referitoare la procesul de prelucrare al datelor. Această categorie reprezintă numai indicatori de calitate.

#### Premisa a.2. Atributele de calitate nu sunt ortogonale

Calitatea datelor poate fi reliefată prin parametri extrem de diferiți, funcție de contextul în care acestea sunt analizate. Unii parametri sunt interdependenți, necesitățile de interpre-

tare impunând reliefaarea unor aspecte intercorelate dintr-un anumit unghi de referință. Exemplul cel mai elocvent pentru acest caz îl constituie interdependența dintre oportunitatea și volatilitatea datelor.

#### Premisa a.3. Heterogenitatea și ierarhizarea calitativă a datelor furnizate

Calitatea datelor aparținând unor colecții, entități, atribute sau valori distincte poate fi diferită. De aici necesitatea menținerii parametrilor de calitate în limitele admisibile stabilite funcție de prelucrările ulterioare ale datelor și a gradului de interdependență al acestora.

#### Premisa a.4. Indicatori de calitate recursivi

O premiză esențială în studiul calității este stabilirea setului de indicatori de calitate relevanți. Apare astfel necesitatea definirii *meta-indicatorilor de calitate* ca un obiectiv de bază în dezvoltarea unei perspective calitative în analizele ulterioare.

### **Premise privind utilizarea datelor**

Utilizatorii datelor au opinii și cerințe diferite referitoare la calitatea datelor. Utilizatorul unui sistem informatic dat poate cunoaște calitatea datelor acestuia. Prin exportarea acestui sistem la alți utilizatori sau prin combinarea cu informații de calități diferite, calitatea datelor poate deveni necunoscută, rezultând necesități diferite privind atributele de calitate acoperitoare pentru utilizatorii și domeniile aplicației.

#### Premisa b.1. Specificitatea atributelor funcție de utilizator

Parametrii și indicatorii de calitate pot diferi de la un utilizator la altul. Astfel, pentru un manager parametrul de calitate critic pentru un raport de cercetare îl poate constitui costul, în timp ce pentru un comerciant credibilitatea și oportunitatea sunt mult mai importante.

#### Premisa b.2. Utilizatorii au standarde de calitate diferite

Nivelurile acceptabile de calitate a datelor pot diferi în limite foarte largi de la un utilizator la altul. Funcție de natura activității desfășurate și a datelor vehiculate, un același nivel al atributelor de calitate poate fi considerat acceptabil sau nu.

Cel mai adesea se face referire în acest sens la oportunitatea datelor. Timpul de întârziere de la preluarea datelor în sistemul informatic și până la punerea lor la dispoziția beneficiarilor poate să reprezinte un parametru critic pentru un anumit utilizator sau conex pentru altul.

Premisa b.3. Atribute și standarde diferite pentru un singur utilizator

Funcție de contextul folosirii datelor, același utilizator poate avea cerințe diferite privind colecțiile, entitățile, atributele sau valorile propriu-zise ale datelor. Prin definirea cu claritate a acestora se oferă posibilitatea concentrării și distribuirii efortului conform priorităților stabilite.

**Dimensiunile și categoriile calității datelor**

În urma studiului efectuat de Richard Wang, Diane Strong și Lisa Guarascio [Wang96] au fost identificate principalele atribute de calitate a datelor. Cercetarea a fost continuată prin analiza factorială a unui număr de 118 variabile, delimitându-se cincisprezece dimensiuni ale calității datelor.

În timp ce atributele susțin nivelul cel mai scăzut pentru care problemele datelor pot fi identificate și înțelese, dimensiunile facilitează înțelegerea la cel mai înalt nivel. Altfel spus, atributele oferă mecanismul cu nivelul cel mai scăzut începând de la care problemele datelor devin evidente, în timp ce dimensiunile susțin condițiile prin care se previne apariția problemelor datelor. Aceasta implică evaluarea și identificarea problemelor datelor prin analizarea atributelor și apoi stabilirea domeniului deficiențelor constatate, prin gruparea atributelor funcție de dimensiunile de apartenență. Oprirea analizei și extragerea concluziilor la nivelul atributelor împiedică reliefaarea problemelor de bază, sistematice.

Funcție de contextul utilizării și al scopului inițial al proiectării datelor, nivelul de importanță acordat diverselor dimensiuni ale acestora poate varia extrem de mult. Un alt aspect important îl reprezintă dependența calității datelor de procesele care le generează, trans-

formă și modifică.

Deși datele pot avea probleme de calitate care aparțin uneia dintre dimensiuni în timp ce sunt satisfăcătoare pentru altă dimensiune, o singură cauză de bază poate afecta dimensiuni multiple. Prin gruparea atributelor în dimensiuni s-a facilitat structurarea problemelor calității datelor. Wang, Strong și Guarascio au observat că dimensiunile formează câteva familii, sau categorii, așa cum se prezintă în Tabelul 1. Fiecare din aceste categorii denotă existența unei omisiuni a procesului sau a unei defecțiuni în procesarea curentă referitoare la un anumit aspect major al calității datelor.

Categoriile reprezintă al treilea nivel de înțelegere al problemelor calității datelor. Pot fi depistate astfel lacunele de procesare a căror apariție sau existență a fost facilitată de condițiile asigurate la nivelul anterior. Conducerea analizei în acest caz permite identificarea și interpretarea cauzelor de bază legate de procesele de achiziționare și manipulare a datelor.

*Calitatea intrinsecă* denotă faptul că datele au valoare prin însăși existența lor. Aici se include nu numai acuratețea și credibilitatea, care sunt evidente pentru specialiștii din domeniu, dar și obiectivitatea și reputația. Astfel se sugerează că, în mod contrar accepțiunii tradiționale, consumatorii de date privesc de asemenea obiectivitatea și reputația drept părți integrante ale calității intrinseci a datelor.

*Contextul* în care vor fi folosite datele este funcție de timp și de consumatorii acestora. Consumatorul de date trebuie să specifice ce tip de lucrare va executa și parametrii contextuali corespunzători. Astfel, în scopul de a se furniza informații oportune și de acuratețe, cercetătorii din domeniul aviației militare au recunoscut necesitatea incorporării explicite a calității contextuale a datelor în sistemele informatice: localizare, rezoluție, condiții atmosferice, tip de obiectiv.

Tabelul 1. Categoriile calității datelor

Categoria	Dimensiunile	Deficiențe de proces
Intrinsecă	Acuratețe, Obiectivitate, Credibilitate, Reputație	La crearea datelor care corespund valorilor actuale sau adevărate.

Contextuală	Valoare Adăugată, Relevanță, Oportunitate, Completitudine, Cantitate adecvată de date	Cu privire la obținerea datelor pertinente necesităților utilizatorului.
Reprezentățională	Interpretabilitate, Ușurința înțelegerii, Consistența reprezentării, Conciziunea reprezentării	Referitoare la furnizarea unor date inteligibile și clare, compatibile cu datele arhivate.
Raportată la abordabilitate	Accesibilitate, Securitatea accesului	Pentru asigurarea disponibilității rapide și obținerea datelor.

*Calitatea reprezentățională* include aspecte legate de formatul datelor (reprezentare concisă și consistentă) și înțelesul datelor (interpretabilitate și ușor de înțeles). Din punct de vedere al beneficiarului, datele sunt de calitate nu numai dacă sunt reprezentate concis și consistent, dar posedă și caracteristici de interpretabilitate și ușor de înțeles.

#### *Calitatea raportată la abordabilitate*

Datorită restricțiilor fizice ale accesării datelor (copii hard în loc de date on-line), în etapa de început a studiului calității datelor s-a recurs la tratarea distinctă a accesibilității. Schimbarea radicală a condițiilor de utilizare și a performanțelor calculatoarelor au modificat opiniile utilizatorilor, în sensul considerării accesibilității ca fiind un aspect important al calității datelor.

#### **Analiza cerințelor de calitate**

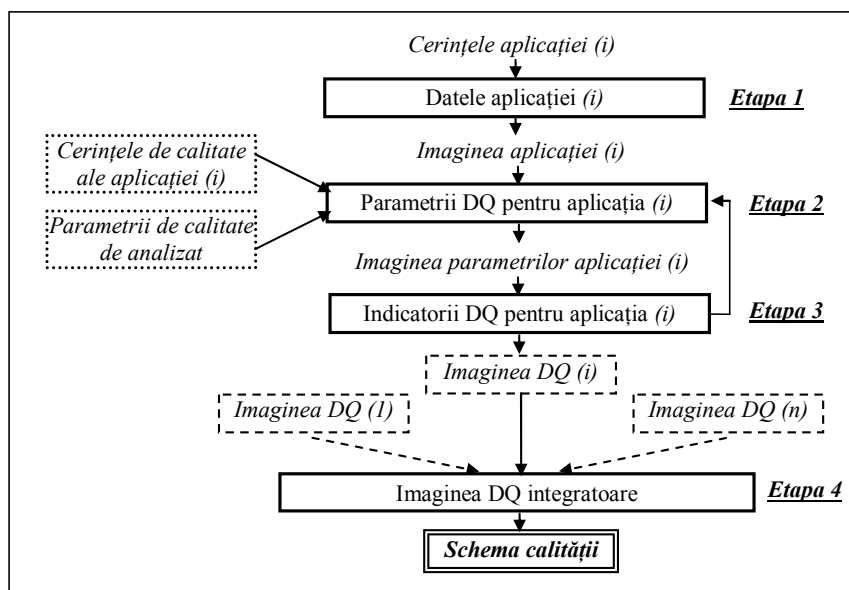
Analiza cerințelor de calitate, *modelarea datelor*, reprezintă un proces similar cu cel al analizei tradiționale, dar cu concentrarea efortului pe aspectele calitative ale datelor. Se poate stabili astfel o paralelă între ciclul

de viață al datelor și metodologia de determinare a cerințelor.

Rezultatul final al modelării calității datelor, *schema calității*, documentează atât cerințele de date ale aplicației cât și problemele calității datelor considerate importante de către echipa de proiectare. Diagrama din Figura 2 ilustrează această metodologie. Pentru fiecare etapă de analiză sunt evidențiate intrarea, procesarea și ieșirea corespunzătoare.

Deoarece una din restricțiile majore ale asigurării calității datelor impune stabilirea numai a acelor date care sunt strict necesare pentru aplicația curentă, prima etapă nu reprezintă activitatea tradițională de determinare a cerințelor de date pentru o anumită aplicație informatică.

Integrarea indicatorilor nu reprezintă o simplă comasare a subseturilor identificate pentru aplicațiile particulare. Prin examinarea determinărilor reciproce ale acestora, se elimină redundanțele și inconsistențele rezultate inițial.

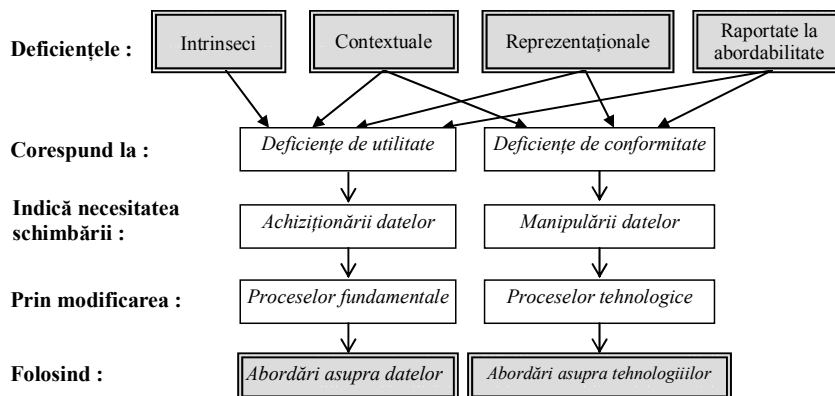


**Fig.2.** Procesul modelării calității datelor

**Abordarea deficiențelor**

Îmbunătățirea calității datelor este asociată cu sistemele sau procesele ale căror deficiențe au fost reliefate la nivelul categoriilor din ierarhia calității datelor. Aceasta deoarece investigarea dimensiunilor scoate în evidență cărei categorii îi corespund deficiențele identificate. Efortul ultim este îndreptat spre realizarea schimbărilor corespunzătoare în pro-

cesele de achiziționare și / sau manipulare ale datelor. Acestea le sunt puse în corespondență metode performante, de achiziționare eficientă a datelor, sau soluții tehnologice, cu impact asupra procesării datelor. Relațiile dintre deficiențele reliefate la nivelul categoriilor calității datelor și abordările corespunzătoare de îmbunătățire sunt prezentate în Figura 3.



**Fig.3.** Relația dintre deficiențe și abordările de îmbunătățire

Așa cum se arată în Figura 3, există două categorii de bază de îmbunătățiri, care corespund necesității schimbării proceselor fundamentale și a schimbării proceselor tehnologice.

Abordările asupra datelor constau în analiza economică, statistică, modelarea și simularea datelor primare, ingineria cerințelor precum și modalități specifice care vizează diminuarea efectelor negative datorate acțiunii anacronice a factorilor umani. Abordările tehnologice sunt impuse ca urmare a instalării de echipamente noi, implementării unor produse informatice sau dezvoltării de noi interfețe. Pentru a îmbunătăți utilitatea datelor este necesar să se implementeze soluții asupra datelor, deoarece schimbările tehnologice nu afectează procesele de achiziție. Deși este evident că utilitatea datelor nu poate fi îmbunătățită independent de procesele care le generează, abordările tradiționale pentru rezolvarea problemelor calității datelor au vizat aplicarea strictă a unor soluții tehnologice.

**Bibliografie**

[Bech00] - Bechtel Hanford Inc. Procedure: *Data Quality Objectives*, BHI-EE-01, Environmental Investigations Procedure: Proce-

dure Number 1.2; Revision 3; September 30, 1999.

[Dill00] - Dillman, D. A.: *Mail and Internet Surveys: The Total Design Method*, New York: John Wiley & Sons, 2000.

[EPA00] - EPA: *Guidance for the Data Quality Objectives Process*; Office of Environmental Information, U.S. EPA, Washington D.C. 20460, August 2000.

[Han02] - Han, J., Kamber, M.: *Data Mining: Concepts and Techniques*, Publisher: Wiley, John & Sons, Incorporated, March 2002.

[Ivan99] - Ivan, I., Noșca, Gh., Tcaciuc, S., Pîrlog, O., Căciulă, R.: *Calitatea datelor*, Infolec, București, 1999.

[Redm01] - Redman, T.C., Daugherty, M., Daugherty, M.: *Data Quality: The Field Guide*, Publisher: Elsevier Science & Technology Books, January 2001.

[Wang96] - Wang, R., Strong, D. and L.Guarascio: "Beyond accuracy: What data quality means to data consumer", *Journal of Management Information Systems*, Spring 1996.

[Wang01] - Wang, R.Y., Ziad, M., Lee, Y.W.: *Data Quality*, Kluwer, 2001.