
Componenta Excel *Data* *Mining*. Aplicații în audit și raportări financiare

Daniel HOMOCIANU,
Universitatea „Alexandru Ioan Cuza” din Iași (UAIC),
E-mail: daniel.homocianu@feaa.uaic.ro

Dinu AIRINEI,
Universitatea „Alexandru Ioan Cuza” din Iași (UAIC),
E-mail: adinu@uaic.ro

Rezumat

Argumentele de performanță în luarea deciziilor bazate pe date economice solicită uzual un management bun al multiplelor formate de date și, de asemenea, viteză de procesare, flexibilitate, portabilitate, automatizare, putere de sugestie și ușurință de utilizare. Lucrarea vine cu idei teoretice și exemple practice în favoarea utilizării componentei Excel Data Mining pentru motivele menționate anterior. Cele mai multe exemple includ figuri legate la scenarii video construite de autori și parte a unei liste on-line interactive cu optsprezece piese. Împreună ele contribuie la înțelegerea celor mai multe cerințe care trebuie îndeplinite pentru a avea exemple valide și rezultate utile.

Cuvinte-cheie: Date economice și financiare, foi de calcul, Data Mining (DM), exemple.

Clasificare JEL: C61, D81, D83, M42

Vă rugăm să citați acest articol astfel:

Homocianu, D. and Airinei, D. (2017), The Excel Data Mining Add-in. Applications in audit and financial reports, *Audit Financiar*, vol. XV, no. 3(147)/2017, pp. 451-468, DOI: 10.20869/AUDITF/2017/147/451

Link permanent pentru acest document:

<http://dx.doi.org/10.20869/AUDITF/2017/147/451>
Data primirii articolului: 17.03.2017
Data revizuirii: 13.04.2017
Data acceptării: 20.04.2017

Introducere

Lucrarea începe de la o serie de tehnici utilizate de cele mai multe instrumente de exploatare a datelor (en. *Data Mining*) pentru volume mari de date din baze de date și prezintă avantajele utilizării foilor de calcul ca aplicații client (msdn.microsoft.com/.../dn282385.aspx). Acestea din urmă sunt foarte familiare utilizatorilor finali și au o interfață care integrează limbaje de programare sau de configurare pentru aplicații de birou precum Visual Basic pentru Aplicații (VBA) în cazul Microsoft Excel (msdn.microsoft.com/.../ee814737.aspx) și Google Apps Script pentru foi de calcul (en. Google Sheets) (developers.google.com/.../sheets). Multe funcții și facilități avansate de procesare, analiză, reprezentare și simulare bazate pe principii de interactivitate și dinamică au impact considerabil asupra capacității utilizatorilor de a percepe, interpreta, înțelege și gestiona informații complexe în diferite situații.

Conceptul de exploatare a datelor înseamnă în esență identificarea supervizată de tipare nedescoperite și relații ascunse existente în seturi foarte mari de date (searchsqlserver.techtarget.com). Inmon care este un bine cunoscut „guru” în depozite de date (computerweekly.com) a dat una dintre cele mai concise definiții ale exploatarei datelor (Inmon și Linstedt, 2014) și anume analiza unor cantități mari de date pentru a găsi tipare precum grupuri de înregistrări sau înregistrări și dependențe neobișnuite. Inițiativele de exploatare a datelor vin de obicei de la departamentele de marketing și vânzări cu amănuntul și sunt potrivite pentru organizații care au baze de date de dimensiuni foarte mari (Airinei, 2002). Acest concept este strâns legat de acela de sisteme de asistare a deciziilor (en. *Decision Support Systems - DSS*) orientate spre date și de inteligența în afaceri (en. *Business Intelligence - BI*) – în special acela pentru obiective strategice (dssresources.com/...id=174) care solicită volume uriașe de date (bi-insider.com). Deși termenul de inteligență în afaceri este cunoscut ca fiind un set de concepte și metode de îmbunătățire a luării deciziilor care a apărut în anii '90 (Howard Dresner din grupul Gartner - dssresources.com/.../dsshistory.html), cercetările publicate în literatura de specialitate indică abordări cu 15 ani mai devreme (Cleland și

King, 1975 și 1976), care conțin referințe clare la BI, planificatorii/ administratorii afacerii și manageri, respectiv la luarea deciziilor.

Conform concluziilor lui Dan Power (dssresources.com/...id=199), instrumentele de exploatare a datelor includ: raționamentul bazat pe caz, vizualizarea datelor (în special grafice, arbori și grupuri), interogări și analize vagi (en. *fuzzy*), algoritmi genetici și rețele neuronale.

Începând de acum câțiva ani suntem martorii implementării acestui concept și a modelelor asociate prezentate nu doar în aplicații dedicate sistemelor de gestiune a bazelor și depozitelor de date, dar și în module ale aplicațiilor de tip foi de calcul care le folosesc așa cum se sugerează chiar din titlul lucrării. Pare evident când ne gândim că asemenea produse dedicate au permis construirea de structuri și modele DM plecând chiar de la un singur tabel (uzual ca agregare a mai multor tabele dintr-o bază de date).

Aplicabilitatea în audit a elementelor teoretice și practice din acest articol, în special cel al performanței (Fraser, 1998) și al raportărilor financiare, este justificată plecând de la o nevoie specifică de a valoriza structurile existente de date (adesea înregistrări în tabele și tabele în baze de date) și de a obține rapid și cu cost minim rapoarte care pot prezenta informații clare despre relația de cauzalitate existentă între eficacitate (rezultate efective/estimate comparate cu cele propuse) și eficiență (resurse consumate comparate cu rezultate obținute/estimate).

Exemplele concrete din acest articol susțin anumite concluzii desprinse din analiza literaturii de specialitate, și anume:

- utilitatea abordării misiunilor de audit folosind tehnici de exploatare a datelor (Vintilescu ș.a., 2010) complementar metodelor clasice de analiză a riscurilor și intervenției la fața locului;
- consacrarea existenței unor zone posibile de integrare a procesului de exploatare a datelor cu procesele de audit grupate pe faze (Sirikulvadhana, 2002), precum planificarea, execuția, documentarea și finalizarea;
- utilizarea de exemple specifice (Wang și Yang, 2009) cum ar fi rețelele neuronale pentru evaluarea riscului, detectarea erorilor și a

fraudelor, determinarea preocupărilor pentru continuitatea activității unei firme, evaluarea dificultăților financiare, precum și realizarea de predicții de faliment și arbori de decizie pentru analiza de faliment, analiza de faliment bancar și a riscului de credit);

- progresul celor mai noi instrumente software care implementează algoritmi de exploatare a datelor și faptul că mulți utilizatori le considerau până nu demult ca nefiind foarte prietenoase (Chersan ș.a., 2013) și necesitând abilități tehnice destul de avansate.

Lucrarea își propune și eliminarea unor confuzii privind utilizarea valorilor tip ștampilă temporală (date calendaristice, părți din ea sau valori înlocuitoare) la exploatarea seriilor de timp cu date economice (de exemplu, volumul vânzărilor recunoscut ca factor de influență directă a nivelului anumitor indicatori financiari precum rezultatul din exploatare).

1. Metodologia de cercetare

Datele sursă pentru exemplele prezentate în această lucrare provin din două mostre de baze de date Microsoft. Primul set de exemple a fost creat plecând de la un fișier bază de date Access denumit *foodmart*

(sites.google.com/.../supp4excel2datamining) inițial disponibil pe CD-ul de instalare a unei versiuni anterioare de Microsoft (MS) SQL Server. Al doilea provine dintr-o bază de date mostră de tip MS SQL Server denumită

„*AdventureWorksDW2012_Data.mdf*” deja instalată și pregătită pentru utilizare într-o aplicație virtuală (y2u.be/Xs2SWtBqdzl) Windows 8.1 pe 32 biți pe care am folosit-o pentru acest articol. Această aplicație a beneficiat de licența Microsoft Imagine/fostă Dream Spark de software educațional pentru toate aplicațiile instalate în

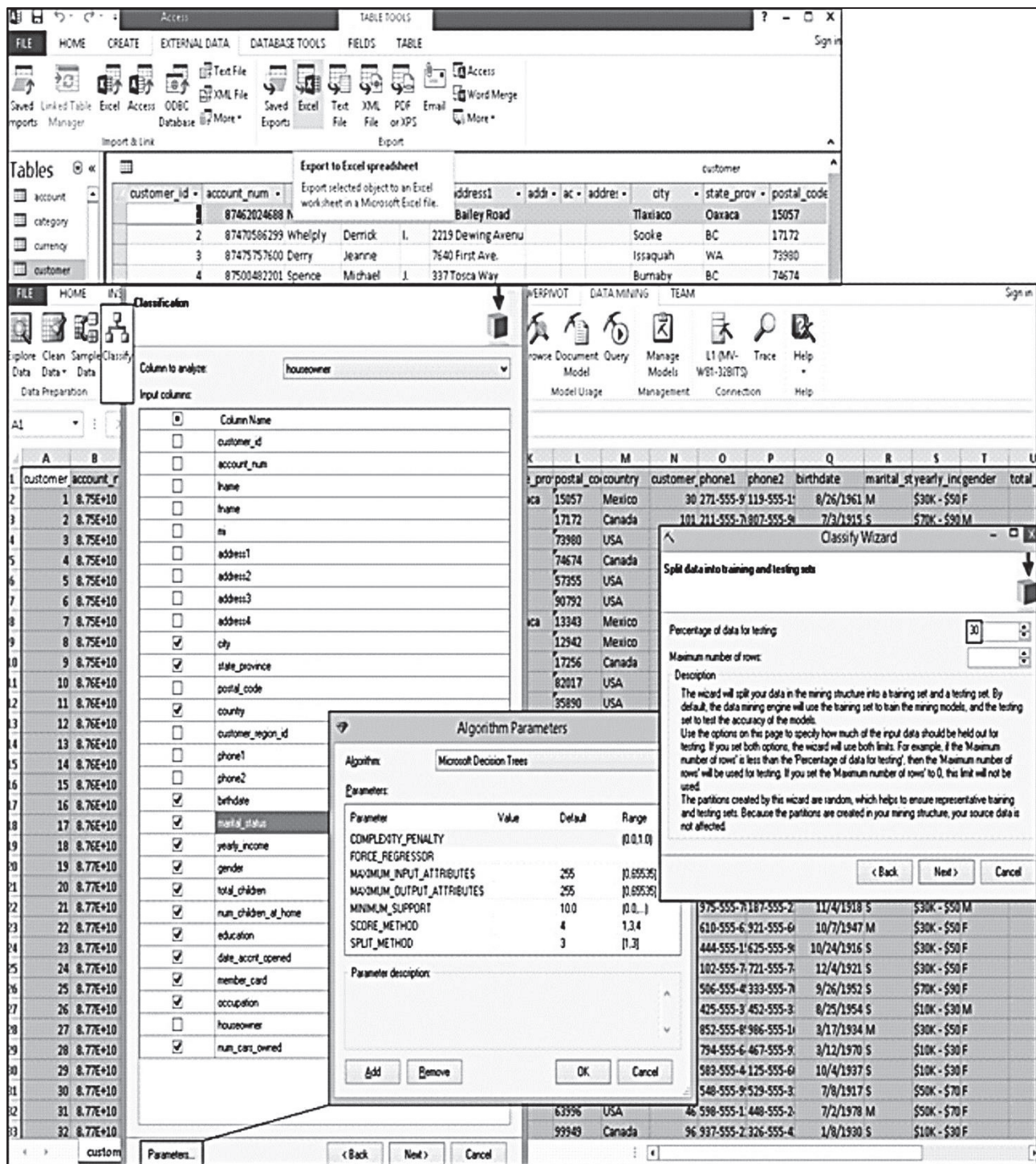
interiorul ei și a fost optimizată pentru Oracle Virtual Box. De fapt SQL Server 2012 (sau 2008) este o componentă necesară pentru instalarea Excel *Data Mining* care este detaliată în al doilea tutorial video (lista menționată mai jos). Deși servesc construirii de exemple și materiale video suport corespunzătoare (tutoriale – lista creată de autori și disponibilă la adresa goo.gl/JDDtFp), asemenea date au doar caracter instructiv în această cercetare de natură pronunțat aplicativă, asemănările cu realitatea fiind doar o coincidență.

2. De la tipare intuitive la o analiză profundă plecând de la surse simple de date precum tabellele

Primul exemplu pe care l-am ales a fost menit să clasifice prin generarea unui arbore decizional în care variabila estimată a fost una de tip categoric, având două valori posibile (proprietar de casă/*houseowner*: Da/Yes – Y sau Nu/No – N) în funcție de alte câmpuri (vezi **Figura nr. 1**), conținând informații despre clienți (un export în Excel din tabelul cu clienți (*customer*) din baza de date *foodmart*).

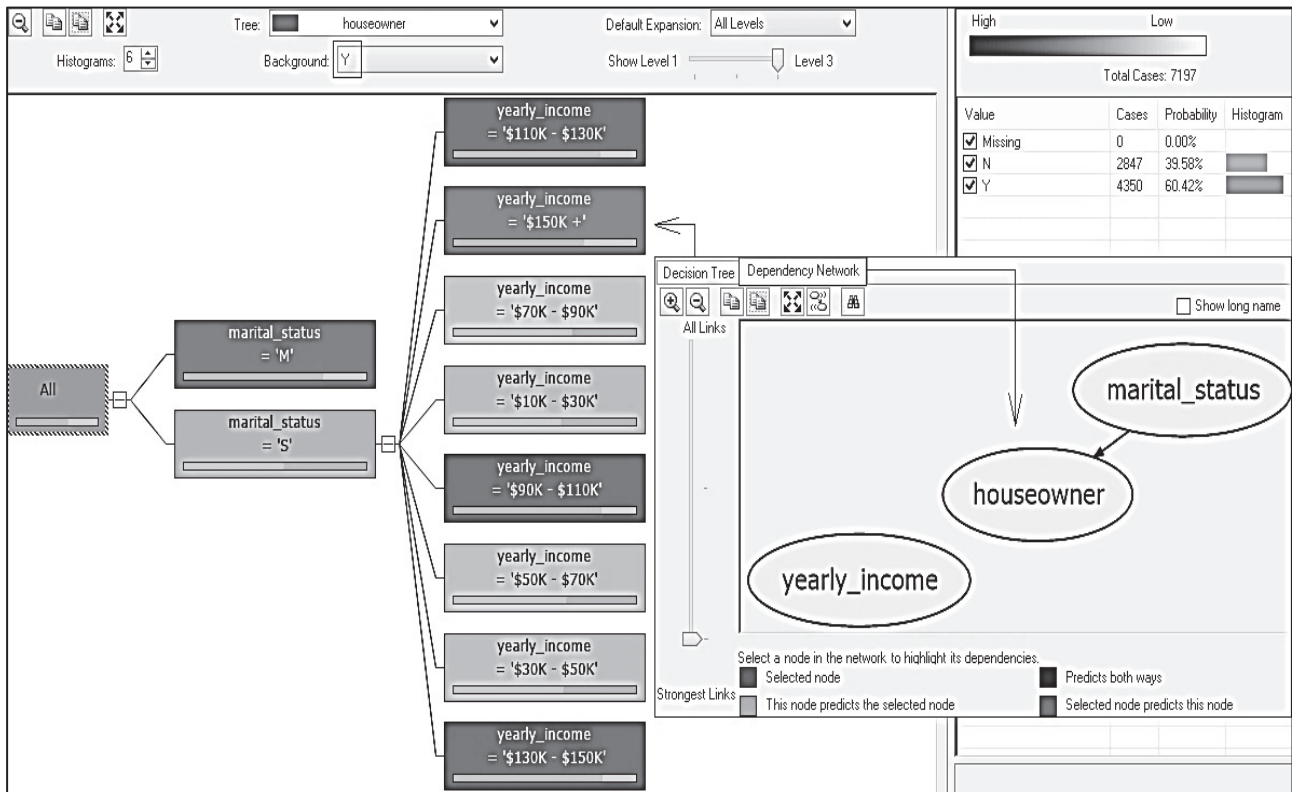
Rezultatele (**Figura nr. 2**) acestei clasificări cu reglajele implicite (algoritmul Arbori Decizionali Microsoft și 30 procente din date pentru test) sunt: (1) un arbore decizional și (2) o rețea de dependențe ce indică cele mai importante variabile ce influențează valoarea atributului *houseowner* și anume: starea civilă/*marital_status* (căsătorit/*married* – M sau necăsătorit/*single* – S) și venitul anual/*yearly_income* (opt praguri în mii/K de dolari/\$: '\$10K - \$30K', '\$30K - \$50K', '\$50K - \$70K', '\$70K - \$90K', '\$90K - \$110K', '\$110K - \$130K', '\$1300K - \$150K', '\$150K +'), în această ordine a importanței.

Figura nr. 1. Exemplu de export urmat de folosirea opțiunii de clasificare a componentei Excel Data Mining și configurarea atributelor de intrare



Sursa: Tutorialul video creat de autori: youtu.be/Nx9xqCX1DjY

Figura nr. 2. Exemplu de rezultat al clasificării plecând de la date dintr-un tabel cu clienți și creat folosind algoritmul Arbori Decizionali Microsoft

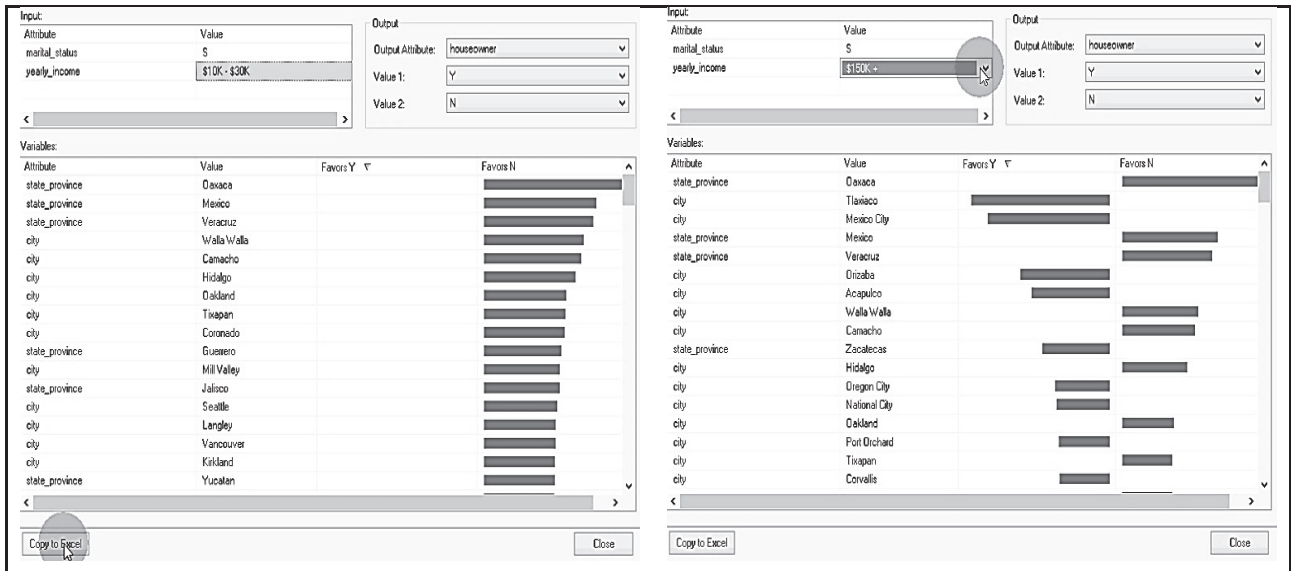


Sursa: Tutorialul video creat de autori: [y2u.be/Nx9xqCX1DjY](https://www.youtube.com/watch?v=y2u.be/Nx9xqCX1DjY)

Așa cum se observă în partea stângă a **Figurii nr. 2** ramurile ce indică o probabilitate mai mare pentru Da (Y – *houseowner*) sunt mai închise, restul fiind colorate cu o nuanță mai deschisă. Putem de asemenea, observa că variabila *houseowner* depinde esențial de variabila *marital_status* (în partea dreaptă a **Figurii nr. 2** – bara de derulare pe legăturile cele mai puternice/*Strongest Links*) și apoi de *yearly_income*

(bara pe toate legăturile/*All Links*). Aceasta se poate deduce și direct din arborele decizional în care un nod mai aproape de rădăcină exprimă un test (inf.uvc.ro) aferent atributului *marital_status*. Accesând opțiunea *marital_status*='M' (nod terminal) am obținut o probabilitate peste 74% în toate cele zece teste făcute în aceeași configurație (coloane de intrare respectiv de analizat, algoritm, procentaj date de test).

Figura nr. 3. Exemplu de analiză discriminatorie după aplicarea regresiei logistice (profs.info.uaic.ro) pentru aceleași condiții de mai sus și specificând cele două variabile de intrare cu impact major

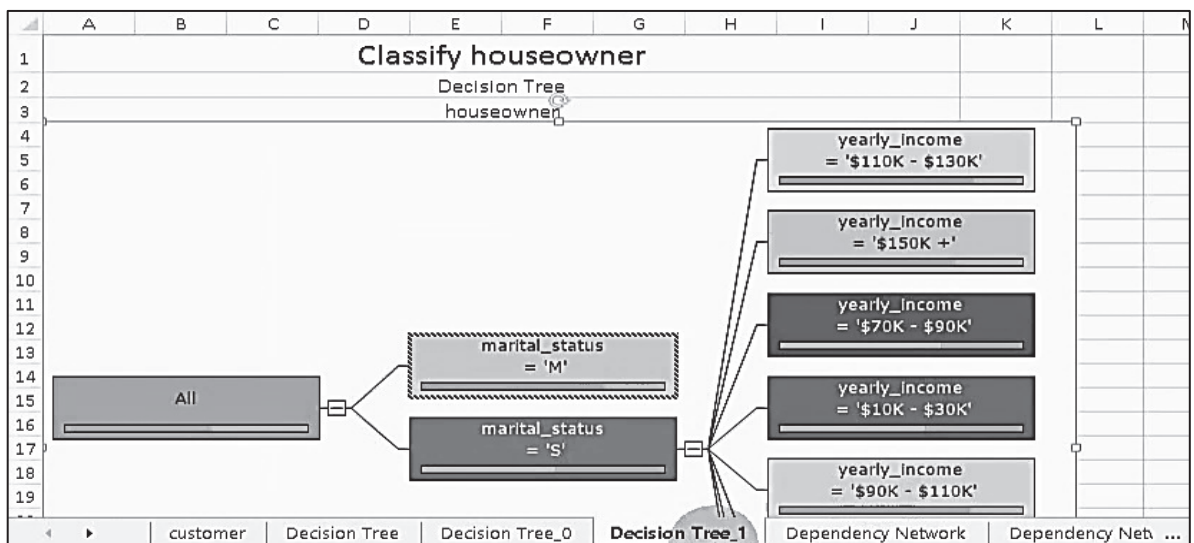


Sursa: Tutorialul video creat de autori: [y2u.be/-6jzQuyTjlo](https://www.youtube.com/watch?v=y2u.be/-6jzQuyTjlo)

În imaginile anterioare (Figura nr. 3) am încercat să punem în evidență cum anume am preconizat probabilitățile pentru clienți de a face parte din cele 2 categorii de răspuns binar (onlinecourses.science.psu.edu): proprietar de casă sau nu, în funcție de anumite variabile explicative și de valorile lor. Am realizat

analiza discriminatorie redată în mod parțial mai sus (Figura nr. 3) plecând de la un alt algoritm, și anume cel de regresie logistică, implementat de Microsoft, folosind o variație a algoritmului de tip rețele neuronale (msdn.microsoft.com/.../ms174828.aspx) care este mai ușor de instruit.

Figura nr. 4. Exemplu de rezultate ale funcționalității „Copiați în Excel” (“Copy to Excel”)



	A	B	C	D
1	Classify houseowner			
2	Neural Network			
3	houseowner			
4	Attribute	Value	Favors Y	Favors N
5	state_province	Oaxaca		
6	state_province	Mexico		
7	state_province	Veracruz		
8	city	Walla Walla		
9	city	Camacho		
10	city	Hidalgo		
11	city	Oakland		
12	city	Tixapan		
13	city	Coronado		
14	state_province	Guerrero		
15	city	Mill Valley		
16	state_province	Jalisco		
17	city	Seattle		
18	city	Langley		
19	city	Vancouver		

Sursa: Tutorialele video create de autori: [y2u.be/Nx9xqCX1DJY](https://www.youtube.com/watch?v=Nx9xqCX1DJY) și [y2u.be/6jzQuyTjlo](https://www.youtube.com/watch?v=6jzQuyTjlo)

Funcția “Copy to Excel” ne-a ajutat să trimitem rezultatele înapoi în Excel ca foi de calcul noi cu capturi de ecran (în partea stângă, **Figura nr. 4** pentru arbori decizionali) sau, mai important, seturi de date cu efecte vizuale implicând de obicei formatare condiționate realizate automat (un exemplu este analiza discriminatorie bazată pe regresie logistică – în partea dreaptă a **Figurii nr. 4**).

Plecând de la rezultatele descrise anterior (**Figurile nr. 1-4**) se pot dezvolta exemple similare care să rezolve inclusiv problema încadrării clientului (corespunzând fazei de acceptare/menținere în demersul de audit) într-una dintre cele două categorii: acceptabil/neacceptabil, plecând de la un istoric validat al unor astfel de decizii, în format tabelar care să includă și multe alte atribute descriptive (zona geografică, sectorul de activitate, numărul mediu de angajați, cifra de afaceri a clientului, evoluțiile anumitor indicatori, nivelul onorariului etc.).

3. Cumularea datelor istorice și folosirea câmpurilor descriptive din tabelele bazei de date

Rapoartele dinamice și interactive care răspund la multe nevoi informaționale și cu care suntem atât de

familiarizați dar și cele statice mai vechi, ca și fotografiile ale informațiilor la momente precise ce generează mai multe întrebări decât răspunsuri (Rasmussen ș.a., 2002), pot utiliza atât date curente cât și istorice. Prima categorie este reprezentată de date din sistemele de procesare a tranzacțiilor (TPS) care se referă în mod obișnuit la anul curent, în timp ce a doua categorie presupune în esență date ce implică o perioadă mai mare de timp. Proportia utilizării celor două categorii depinde în mod esențial de nevoile decizionale (la nivel operațional, tactic sau strategic). Pentru minimizarea redundanței și a dependenței datelor sau din rațiuni de spațiu de stocare și de nevoie de viteză de scriere (deshpande.mit.edu) schema unei surse de date tradiționale de tip relațional este gândită în mod uzual ca fiind formată din mai multe tabele obținute prin aplicarea principiilor normalizării (w3schools.in). Mai mult, din cauza unor motive suplimentare de performanță (nevoi de viteză de citire respectiv de scriere) datele istorice trebuie separate de cele curente. Ambele categorii includ în esență înregistrări din tabele cu tranzacții (de exemplu cheltuieli, vânzări, examene etc.) diferența fiind dată de valoarea ștampilei temporale. Aceasta explică de ce aceste tabele încărcate doar cu date istorice sunt redenumite cu un indicativ de timp, arhivate și separate de restul sistemului tranzacțional pentru a-i îmbunătăți

performanța operațională (curentă). Când este nevoie de volume mari de date istorice pentru analize bazate pe interogări ad-hoc, sistemele trebuie să procedeze invers prin agregarea într-un singur tabel (sursă pentru un tabel de fapte într-un depozit de date) a tuturor înregistrărilor din arhivele istorice ale tabelelor cu tran-

zații (de același tip ca cel rezultat). În cele mai multe cazuri, aceasta generează avantajul unui potențial crescut de identificare de tipare, dar vine și cu dificultăți legate de punerea laolaltă a datelor într-un format comun și consistent, în special atunci când aplicațiile și structura sursei de date s-au modificat și ele în timp.

Figura nr. 5. Acumularea de date tranzacționale și descriptive folosind două interogări SQL MS Access în cascadă

The screenshot displays the Microsoft Access interface. At the top, a SQL query window shows the following code:

```
SELECT * INTO inventory_fact
FROM (SELECT * FROM inventory_fact_1997
UNION SELECT * FROM inventory_fact_1998
ORDER BY time_id)
```

A dialog box titled "Microsoft Access" is open, displaying the message: "You are about to paste 11352 row(s) into a new table. Once you click Yes, you can't use the Undo command to reverse the changes. Are you sure you want to create a new table with the selected records?" with "Yes" and "No" buttons.

Below the dialog, a table structure for "inventory_fact" is shown:

Field Name	Data Type
inventory_id	AutoNumber
product_id	Number
time_id	Number

The main window shows a database relationship diagram with tables: time_by_day, product, product_class, store, warehouse, warehouse_class, and inventory_fact. The inventory_fact table is linked to product_class, store, and warehouse. A query window titled "query4inventory(ext)" contains a complex SQL query joining multiple tables.

At the bottom, a field list table is visible:

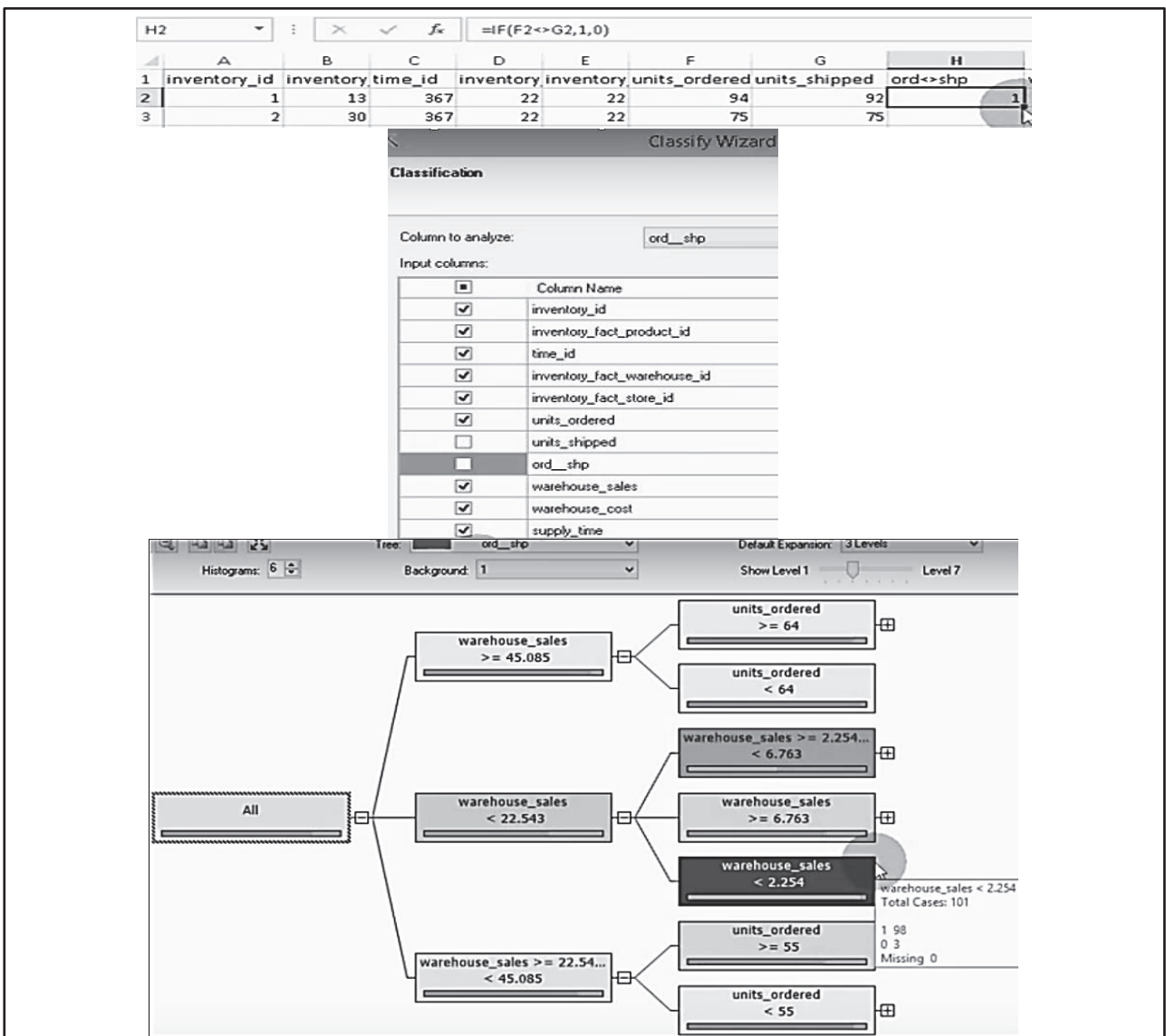
Field:	inventory_fact.*	product.*	product_class.*	warehouse.*	warehouse_class.*	store.*	the_date
Table:	inventory_fact	product	product_class	warehouse	warehouse_class	store	time_by_day
Sort:							
Show:	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Criteria:							

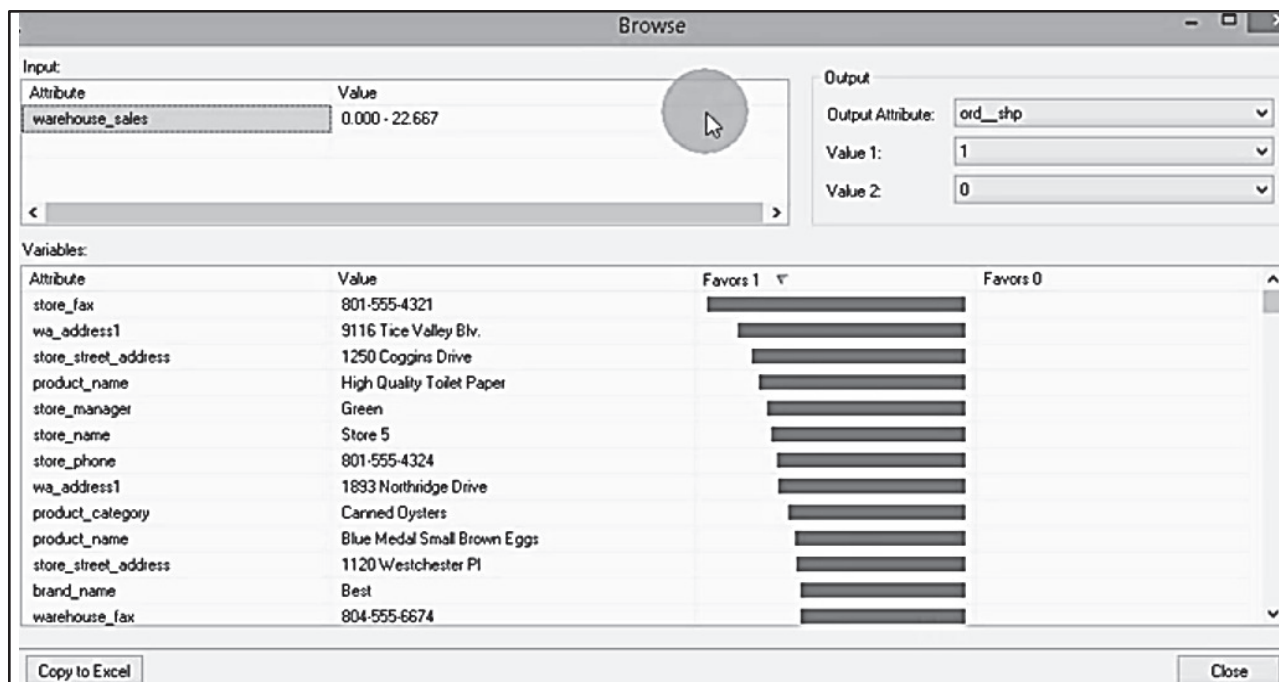
Sursa: Tutorialul video creat de autori: [y2u.be/kTuYLuav3Eo](https://www.youtube.com/watch?v=y2u.be/kTuYLuav3Eo)

Figura nr. 5 prezintă un exemplu de acumulare de date de inventar (aplicații inclusiv în auditul transportului de marfă) în doi pași majori, ce corespund celor două interogări SQL în Microsoft Access. Prima variantă are la bază cumulara (clauza UNION) înregistrărilor din două tabele cu tranzacții de același tip ce corespund doar celor doi ani (1997 și 1998), respectiv adăugarea unei coloane id necesare (*inventory_id* cu valori generate automat – de tip *AutoNumber*) în tabelul persistent rezultat (clauza INTO). Cea de a doua variantă are la bază extragerea temporară a valorilor

câmpurilor descriptive din toate tabelele aflate în legătură sau potrivite pentru o legătură (Figura nr. 5 – clauza INNER JOIN) cu cel care rezultă din prima interogare de mai sus, și anume *inventory_fact*. În acest caz datele tabelare rezultate, constând în al doilea set de doar 11.352 înregistrări, nu vor intra într-un tabel persistent al bazei de date (un fel de denormalizare - searchoracle.techtarget.com) ce este necesar în alte situații pentru economisirea de timp în detrimentul spațiului de stocare și va servi pentru export extern (Excel) imediat după executarea/rularea interogării în sine.

Figura nr. 6. Rezultatele folosirii consecutive a 2 modele Data Mining - câmp țintă derivat cu 2 valori posibile





Sursa: Tutorialele video create de autori: y2u.be/4nOMMRoC2BU și y2u.be/wce_aoTTsbw

Mai mult, din motive legate de viteza de proiectare am ales toate câmpurile sursă fără selectarea lor explicită, dar indicându-le prin folosirea celui mai flexibil caracter de căutare și anume asteriscul / "*" după denumirea tabelului (Figura nr. 5), atât în modul SQL cât și în cel de proiectare asistată/design (techrepublic.com). Din aceleași motive prezentate mai sus, noua coloană derivată necesară pentru analize (Figura nr. 6 - atribut de ieșire pentru ambele modele: clasificare-sus și regresie logistică-jos) a fost definită ulterior direct în Excel folosind funcția IF (partea superioară a fig.6, ord<>shp ca 1 sau 0, adică unitățile comandate/units_ordered și unitățile expediate/units_shipped sunt diferite, respectiv egale).

4. Reguli de asociere pentru identificarea de modele comportamentale

În teoria și practica depozitelor de date și a modelării multidimensionale exemplele de mai jos amintesc de schema „fulg de nea”, însemnând că

sursa unei dimensiuni (perspectivă de analiză bazată pe coloane descriptive, organizate în ierarhii) nu este reprezentată doar de un singur tabel, ci de mai multe tabele aflate în legătură (relații de timp „unu la mai multe”: de exemplu categoria de produs/product category, subcategoria/product subcategory și produsul/product – Figura nr. 7) și capabile să asigure suportul unei analize pentru mai mult decât un singur câmp descriptiv pe dimensiune. Pentru a putea aplica algoritmul regulilor de asociere, în cazul de mai jos am avut nevoie și de valori repetitive pentru câmpul Numărul comenzii de vânzare/SalesOrderNumber de asociat la diferite categorii, subcategorii sau nume de produse.

Motivul principal pentru acumularea acestor date descriptive din multiplele tabele ale bazei de date din exemplul de mai sus (Figura nr. 7) este de a determina reguli de asociere în comportamentul de cumpărare de tipul „Dacă achiziționez produsul X, voi cumpăra și produsul Y.” (tabelul sursă FactInternetSales) precum și cele mai importante dependențe (Figura nr. 8).

Figura nr. 7. Acumularea ambelor tipuri de date: tranzacționale despre vânzări și descriptive despre produse folosind o singură interogare MS SQL Server

The image shows two overlapping windows. The top window is a Notepad file named 'sql.txt' containing an SQL query. The bottom window is a SQL Server Enterprise Edition query window showing the same query and its results in a table format.

```

SELECT [SalesOrderNumber], [EnglishProductCategoryName],
[EnglishProductSubcategoryName], [EnglishProductName], [UnitPrice]
FROM [dbo].[FactInternetSales], [dbo].[DimProductCategory], [dbo].
[DimProductSubcategory], [dbo].[DimProduct]
WHERE [dbo].[FactInternetSales].[ProductKey]=[dbo].[DimProduct].[ProductKey] AND
[dbo].[DimProduct].[ProductSubcategoryKey]=[dbo].[DimProductSubcategory].
[ProductSubcategoryKey] AND [dbo].[DimProductSubcategory].[ProductCategoryKey]=
[dbo].[DimProductCategory].[ProductCategoryKey]
    
```

SalesOrderNumber	EnglishProductCategoryName	EnglishProductSubcategoryName	EnglishProductName	UnitPrice	
1	SO43697	Bikes	Road Bikes	Road-150 Red, 62	3578.27
2	SO43698	Bikes	Mountain Bikes	Mountain-100 Silver, 44	3399.99
3	SO43699	Bikes	Mountain Bikes	Mountain-100 Silver, 44	3399.99
4	SO43700	Bikes	Road Bikes	Road-650 Black, 62	699.0982
5	SO43701	Bikes	Mountain Bikes	Mountain-100 Silver, 44	3399.99
6	SO43702	Bikes	Road Bikes	Road-150 Red, 44	3578.27
7	SO43703	Bikes	Road Bikes	Road-150 Red, 62	3578.27
8	SO43704	Bikes	Mountain Bikes	Mountain-100 Black, 48	3374.99
9	SO43705	Bikes	Mountain Bikes	Mountain-100 Silver, 38	3399.99
10	SO43706	Bikes	Road Bikes	Road-150 Red, 48	3578.27
11	SO43707	Bikes	Road Bikes	Road-150 Red, 48	3578.27
12	SO43708	Bikes	Road Bikes	Road-650 Red, 52	699.0982
13	SO43709	Bikes	Road Bikes	Road-150 Red, 52	3578.27
14	SO43710	Bikes	Road Bikes	Road-150 Red, 56	3578.27
15	SO43711	Bikes	Road Bikes	Road-150 Red, 48	3578.27

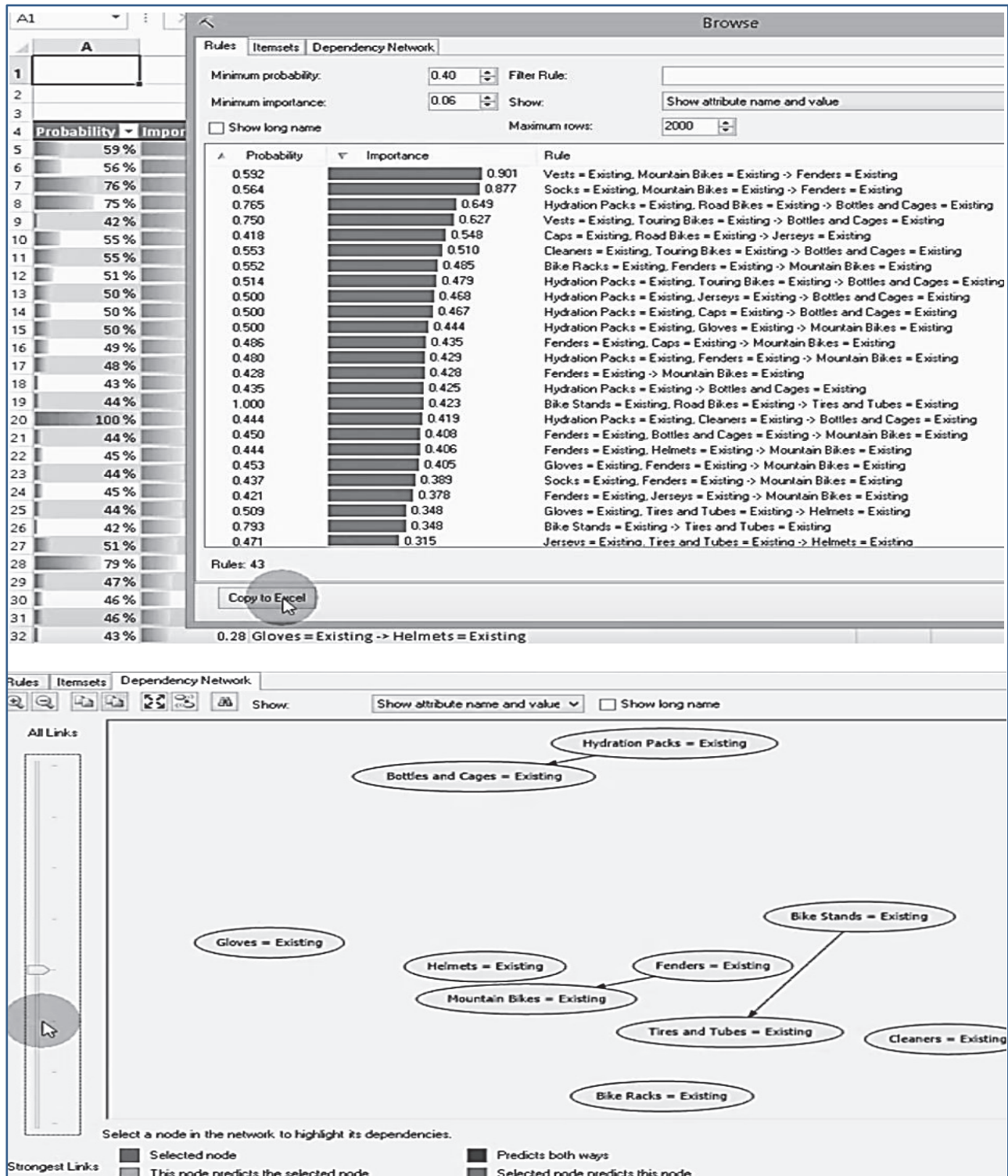
Query executed successfully. MV-W81-32BITS (11.0 SP1) mv-w81-32bits\admin (56) AdventureWorksDW2012 00:00:00 60398 rows

Sursa: Tutorialul video creat de autori: [y2u.be/2rW2wK77HD8](https://www.youtube.com/watch?v=y2u.be/2rW2wK77HD8)

Rezultatele din **Figurile nr. 7 și 8** ne pot face să înțelegem de ce aplicațiile algoritmilor de identificare de reguli de asociere pot contribui în audit și la detectarea și prevenirea fraudelor. Cu titlu de exemplu, dacă setul de date de intrare ar avea atribute precum: Identificator daună, Tip produs de asigurare, Nume asigurat, Asigurator, Nume agent constatator și

Soluționare (daună totală sau parțială și respingere), iar algoritmul ar identifica „DACĂ asigurare Casco, Asigurat X și Agent constatator Y ATUNCI daună totală” ca asociere cu probabilitate și importanță mari, aceasta nu ar însemna neapărat o alarmă de fraudă, dar ar merita măcar efortul de a face investigații suplimentare.

Figura nr. 8. Rezultatele aplicării algoritmilor cu reguli de asociere Microsoft (opțiunea asociază/asociate)



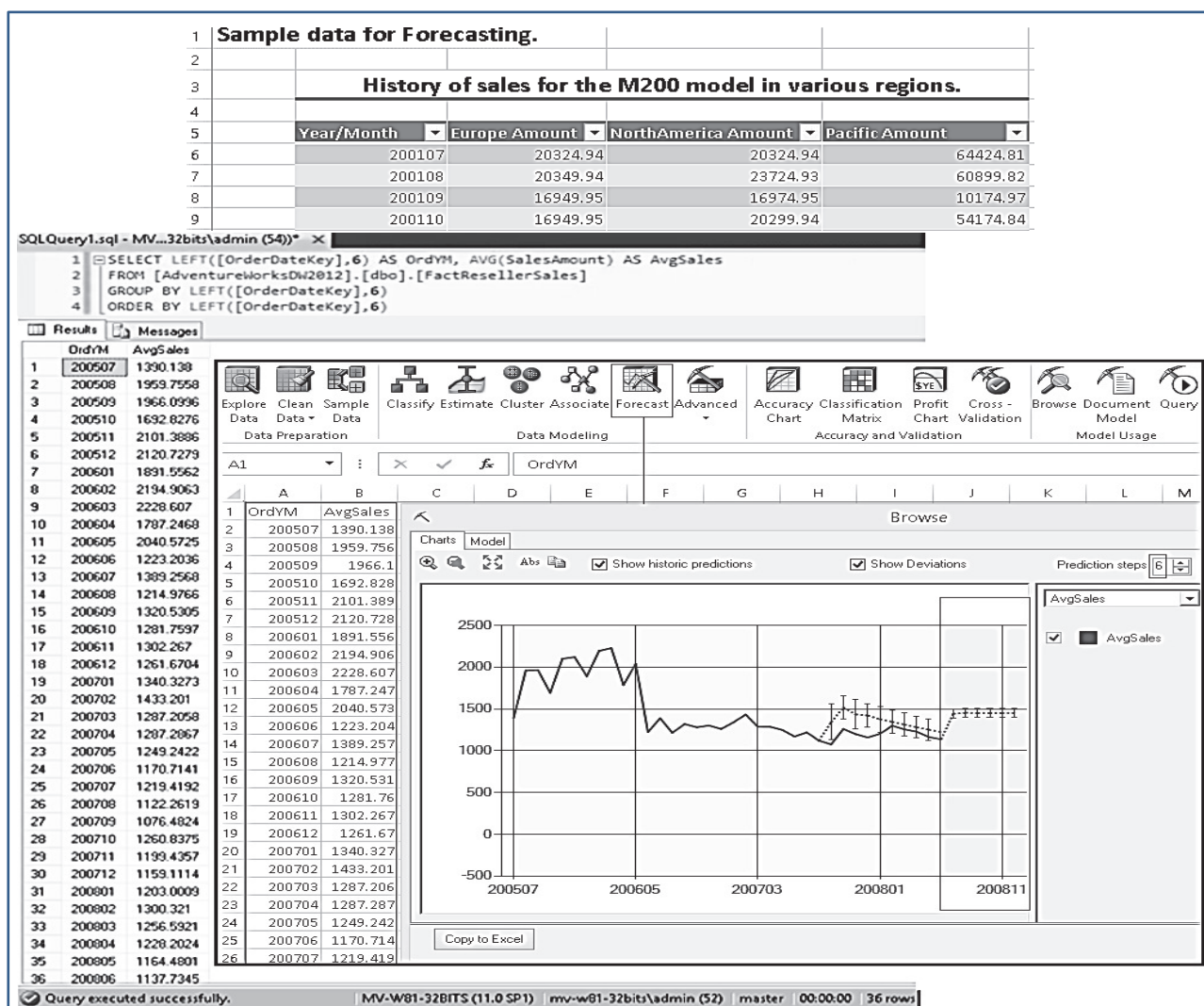
Sursa: Tutorialul video creat de autori: [y2u.be/3_8E01hnSD0](https://www.youtube.com/watch?v=y2u.be/3_8E01hnSD0)

5. Previziuni plecând de la date istorice agregate

Pentru mai multe date istorice decât cele din exemplul precedent (Figura nr. 5) am luat în considerare crearea unui scenariu special de previziune mai apropiat de realitate. Am pornit de la zero cu un nou exemplu care solicită date pe 36 de luni, din patru ani calendaristici, de data aceasta folosind o interogare simplă SQL pentru un singur tabel, dar cu clauze ORDER BY și GROUP BY

pentru a obține rezultate sortate și valori agregate, cum ar fi: sume, valori medii, număr total de apariții, număr de apariții pentru o condiție specificată etc. În cazul nostru acestea au fost medii lunare pe mai mulți ani, combinate într-un singur câmp numeric derivat, prin trecerea de la stânga la dreapta în ordinea specifică: de la ani la luni, corespunzând celei de la unități mai mari la unități mai mici (Figura nr. 9 – exact ca în mostra de date a Microsoft care este furnizată la instalarea componentei Data Mining).

Figura nr. 9. Agregarea de date istorice (stilul ștampei temporale din mostra Microsoft) cu o interogare SQL Server (clauza GROUP BY) și o copiere simplă a rezultatelor ce vor fi folosite pentru previziune (componenta Excel Data Mining)



Sursa: Tutorialele video create de autori: y2u.be/RjTGWROD0TI și y2u.be/qHJ3Zm3JBT4

După pașii descriși mai sus (Figura nr. 9) și alte câteva operații de prelucrare (Figura nr. 10) vom ajunge la un set de date potrivit pentru previziuni implementate folosind algoritmi de serii de timp Microsoft ca o combinație a algoritmilor ARIMA (medie mobilă auto-

regresivă și integrată - optimizată pentru creșterea acurateței în predicțiile pe termen lung) și ARTXP (arbori de auto-regresie cu predicție încrucișată optimizați pentru estimarea următoarei valori probabile într-o serie de timp - msdn.microsoft.com/.../bb677216.aspx).

Figura nr. 10. Derivarea și explicarea ștampilelor temporare corecte ca date complete stocate intern (Excel) drept numere în surse de date în format potrivit ca suport pentru previziuni nedistorsionate

1	OrdYM	AvgSales													
2	200507	1390.138	7/31/2005												

1	% from 1st end of month (date)	AvgSales	OrdYM	AvgSales	% from 1st	
2	100.00 %	38564	1390.138	200507	1390.138	100.0000 %
3	100.08 %	38595	1959.756	200508	1959.756	100.0005 %
4	100.16 %	38625	1966.1	200509	1966.1	100.0010 %
5	100.24 %	38656	1692.828	200510	1692.828	100.0015 %
6	100.32 %	38686	2101.389	200511	2101.389	100.0020 %
7	100.40 %	38717	2120.728	200512	2120.728	100.0025 %
8	100.48 %	38748	1891.556	200601	1891.556	100.00469 %
9	100.55 %	38776	2194.906	200602	2194.906	100.00474 %
10	100.63 %	38807	2228.607	200603	2228.607	100.00479 %

Sursa: Tutorialul video creat de autori: y2u.be/e0SkDwG9mNY

Ne-am gândit, de asemenea, la derivarea automată a etichetelor ștampilelor temporale corecte (formatul LL/ZZ/AAAA tradus într-un număr întreg – Figura nr. 10) și am prezentat mai multe detalii despre comportamentul lor comparativ la obținerea de funcții de trend și previziunea rezultatelor cu această componentă Excel Data Mining (ultimele trei tutoriale din lista menționată anterior).

6. Suport pentru interogarea de modele Data Mining persistente

În primul rând, persistent în acest context, se referă la un model definit astfel încât să fie procesat și stocat pe server (SQL Server Analysis Services/servicii de analiză – modul diferit de motorul de baze de date/Database

Engine) și disponibil pentru interogare (Figura nr. 11).

Componenta Data Mining din Excel oferă multe avantaje față de utilizarea directă a SQL Server Analysis Service. Printre altele, putem menționa aici: viteza de utilizare a mediului tabelar și a setului de formule din Excel, posibilitatea multor exporturi/importuri în/din foi de calcul tabelar plecând de la diferite formate de baze de date și de a implica indirect multiple tabele sursă folosind limbajul structurat de interogare (SQL), posibilitățile de exploatare a structurilor și modelelor obținute: direct (opțiunea copy to Excel), cu interogări (extensia DMX a SQL – Figurile nr. 12 și 13) sau programatic?? (Figura nr. 12). Ultimele două sunt condiționate de activarea persistenței la definirea modelelor (opțiunea de folosire a unui model temporar/use temporary model nebitată – Figura nr. 13 comparativ cu Figurile nr. 1 și 2).

Figura nr. 11. Exemple de interogări Data Mining eXtensions (DMX) pentru previziunea vânzărilor (SQL Server Analysis Services) pe baza unui model DM de serii de timp plecând de la o sursă de date în format greșit (ștampilă temporală text: 200815 / a 15-a lună în 2008)

The screenshot displays the SQL Server Analysis Services (SSAS) interface. At the top, a DMX query is shown in the query editor:

```

1 SELECT
2 PredictTimeSeries([forecast_AVG_sales_model].[AvgSales],3) AS PredAvgSales
3 FROM [forecast_AVG_sales_model]
    
```

The results pane shows a table with two columns: \$TIME and AvgSales. The data rows are:

\$TIME	AvgSales
200807	1443.24228516...
200808	1454.63464828...
200809	1450.83187874...

The Object Explorer shows the server structure, with the 'forecast_AVG_sales_model' selected under Mining Models.

The bottom part of the screenshot shows a second DMX query:

```

1 SELECT
2 PredictTimeSeries([forecast_AVG_sales_model].[AvgSales],9)
3 FROM [forecast_AVG_sales_model]
    
```

The results pane shows a table with two columns: \$TIME and AvgSales. The data rows are:

\$TIME	AvgSales
200811	1452.83271562
200812	1453.50726198
200813	1453.77358263
200814	1454.14759460
200815	1454.43486230

The year '200815' is highlighted in the results table, indicating the forecast for the 15th month of 2008.

Sursa: Tutorialul video creat de autori: <http://y2u.be/qHJ3Zm3JBT4>

Figura nr. 12. Exemplu brut de interogare programatică a unui model Data Mining bine definit folosind o interogare DMX în Visual Basic (.NET) precedată de testarea a mare parte din ea în SQL Server Analysis Services

The screenshot displays a Visual Studio environment with a VB.NET application. The code in `Form1_Load` establishes a connection to the `DMAddinsDB` database on the `MV-W81-32BITS` server. It executes a DMX query to predict time series data for `AvgSales`. The results are shown in a message box, with one example being `7/31/2008: 1449.64`. Below the code, the `Object Explorer` shows the server structure, and the `SQL Server Enterprise Manager` displays a time series chart for `AvgSales` from 2005 to 2008. A `Results` window at the bottom shows the output of the query.

PredAvgSales.\$TIME	PredAvgSales.AvgSales
7/31/2008 12:00:00 AM	1449.64321772663
8/31/2008 12:00:00 AM	1460.4264515793
9/30/2008 12:00:00 AM	1456.07251908031

Sursa: Proiecția autorilor obținută după încercări de dezvoltare cu VB și SQL Server

Figura nr. 13. Exemplu de interogare de predicție DMX (proprietar casă/houseowner) pe baza unui model de clasificare persistent cu arbori de decizie (vizualizarea generică a conținutului în fundal)

The screenshot displays the SQL Server Enterprise Miner interface. On the left, a tree view shows the mining model structure. The main area is divided into three panes:

- Node Details:** A table showing metadata for the 'DT_m_HO' model. Key fields include 'MODEL_CATALOG' (DMAddinsDB), 'MODEL_NAME' (DT_m_HO), 'ATTRIBUTE_NAME' (houseowner), 'NODE_NAME' (00000000r0107), 'NODE_UNIQUE_NAME' (00000000r0107), 'NODE_TYPE' (4 (Distribution)), 'NODE_GUID', 'NODE_CAPTION' (yearly income = '\$150K +'), 'CHILDREN_CARDINALITY' (0), 'PARENT_UNIQUE_NAME' (00000000r01), 'NODE_DESCRIPTION' (marital_status = 'S' and yearly_income = '\$150K +'), 'NODE_RULE' (XML representation of the decision rule), 'MARGINAL_RULE', 'NODE_PROBABILITY' (0.0094483812699736), 'MARGINAL_PROBABILITY' (0.019036954087346), and 'NODE_DISTRIBUTION' (table with columns: ATTRIBUTE_NAME, ATTRIBUTE_VALUE, SUPPORT, PROBABILITY, VARIANCE, VALUETYPE).
- Query Editor:** Contains the following DMX query:


```

      1 SELECT
      2 [houseowner],
      3 PredictProbability([houseowner], 'Y') AS [HouseOwner = Yes],
      4 PredictProbability([houseowner], 'N') AS [HouseOwner = No]
      5 FROM [DT_m_HO]
      6 NATURAL PREDICTION JOIN
      7 (SELECT 'S' AS [marital_status],
      8 '$150K +' AS [yearly_income]) AS t
      
```
- Results:** A grid showing the output of the query. The columns are 'houseowner', 'HouseOwner = Yes', and 'HouseOwner = No'. The results show a predicted probability of approximately 0.806 for 'Y' and 0.1938 for 'N'.

Sursa: Proiecția autorilor obținută după încercări de dezvoltare cu SQL Server

Acest ultim avantaj ne amintește că generarea programatică (Airinei și Homocianu, 2009) de tablouri de bord și tabele de scoruri Excel folosind reprezentări sugestive, indicatori de alertă și formatați dinamice cu suport pentru BI s-a simplificat mult începând cu versiunea 2007 a pachetului Microsoft Office. Combinarea acestora cu abilitatea de a determina programatic modele comportamentale și de a genera valori previzionate plecând de la instrumente de înaltă performanță și ușor de folosit precum această

componentă *Data Mining*, disponibilă pentru Office 2010, 2013 și 2016, promite mult în termeni de productivitate. Toate aceste progrese au fost definite după mai mulți ani de utilizare de tehnologii dedicate și acum bine-cunoscute (de exemplu, SQL Server testat de autori de la începutul anilor 2000).

Pentru produsele de tip foi de calcul tabelar (dssresources.com/.../sshistory.html) precum: VisiCalc, Lotus 1-2-3, Microsoft Excel, Microsoft Works Spreadsheet, Sun Open Office Spreadsheets, Polaris

Office Sheet și Google Sheets experiența medie a utilizatorilor finali este de până la zeci de ani. Mai mult, ușurința de utilizare a acestor aplicații chiar și doar ca instrumente de interfață pentru conectarea la datele din baze de date și depozite de date și afișarea lor a fost un motiv obiectiv pentru a continua cu testarea componentei *Data Mining* care a condus la conceperea acestui articol.

Folosind o modalitate de raportare care se identifică cu o secvență de pași care împrumută numele de la cele optsprezece tutoriale suport și, de asemenea, unele tehnici anterior definite și anume: E2P4CAFR (Homocianu, 2015), ACCORD/CADRE (Homocianu și Airinei, 2014) și S-DOT (Homocianu și Airinei, 2014) se poate ajunge în etape, dar cu un număr minim de pași de urmat, la anumite reprezentări care sunt dinamice, interactive, sugestive, bazate pe cauzalitate și înrădăcinate în realitatea curentă și în istoria definită de datele stocate în sursele de date ale organizației.

Concluzii

Putem concluziona că posibilitățile componentei Excel *Data Mining* sunt peste așteptările unui analist de afaceri, oferind avantajul integrării tiparelor de clasificare identificate, a regulilor de asociere și a predicțiilor cu suportul pentru conectivitate la formate variate de date. Validările de date, reprezentările grafice avansate, referențieri geografice, formatările condiționate automate și indicatorii cheie de performanță (KPI), tabelele și graficele de tip pivot și *power pivot*, rezolvarea automată de probleme de optimizare (*solver*) și limbajul DAX (expresii de analiză a datelor/*Data Analysis eXpressions*) împreună cu limbajul tradițional de formule

sporesc șansele definirii de tablouri de bord fundamentate pe simulări, analize și modele *Data Mining* cu adevărat utile pentru personalul de audit interesat de monitorizarea performanței.

Sperăm că am identificat multe motivații reale pentru alegerea acestei componente Microsoft pentru pachetul Office ca un instrument *Data Mining* aproape în timp real, dincolo de multe alte recomandări disponibile în literatura și practica de specialitate.

Dincolo de exemple substanțiale de lucru cu aplicații software bine cunoscute, disponibile pentru o gamă largă de utilizatori, care asigură metode avansate de analiză, interogare și reprezentare a datelor curente și istorice specifice instrumentelor suport pentru *Data Mining* și *Business Intelligence*, lucrarea furnizează și o scurtă descriere teoretică necesară înțelegerii unei modalități rapide de generare de rapoarte complexe și dinamice precum tablourile de bord fundamentate pe analize și modele *Data Mining* plecând în special de la date despre vânzări și date financiare.

Tutorialele video dezvoltate de autori, integrate într-o listă, la care s-a făcut referire în mod succesiv în această lucrare, demonstrează încercările de îmbogățire a modului de raportare anterior menționat și de asigurare a minimizării numărului de pași necesari atunci când se încearcă implementarea de exemple similare.

În ansamblu, articolul încearcă să transmită prin exemple clare anumite caracteristici dorite cum ar fi: viteză, simplitate, capacitate de sinteză, transparență, flexibilitate și disponibilitate în raportarea fundamentată pe exploatarea datelor, ca elemente cheie de performanță în pregătirea situațiilor financiare și sprijinirea activităților de audit.

BIBLIOGRAFIE

1. Airinei, D. (2002), *Depozite de date*, Editura Polirom, Iași.
2. Airinei, D. și Homocianu, D. (2009), The Geographical Dimension of DSS Applications, *Analele Științifice ale Universității „Alexandru Ioan Cuza” din Iași*, Tome LVI, pp. 637-642.
3. Chersan, I.C., Carp, M. și Mironiuc, M. (2013), Data mining – o provocare pentru auditorii financiari, *Audit Financiar*, vol. XI, nr. 10, pp. 57-64.
4. Cleland, D.I. și King, W.R. (1975), Competitive Business Intelligence Systems, *Business Horizons Journal*, vol. 18, nr. 6, pp. 19-28, DOI 10.1016/0007-6813(75)90036-1.
5. Fraser, L.E. (1998), Public Sector Audit - Business Integration and Causal Analysis, *Quality Audit Conference*, February 26-27, 1998, Louisville, KY, vol. 7.
6. Homocianu, D. (2015), Excel Power Pivot's Applications in Audit and Financial Reports, *Audit Financiar*, vol. XIII, nr. 11, pp. 127-138.
7. Homocianu, D. și Airinei, D. (2014), Business Intelligence facilities with applications in audit and

- financial reporting, *Audit Financiar*, vol. XII, nr. 9, pp. 17-29.
8. Homocianu, D. și Airinei, D. (2014), Consolidating source data in audit reports, *Audit Financiar*, vol. XII, nr. 8, pp. 10-19.
 9. Inmon, W.H. și Linstedt, D. (2014), *Data architecture: A primer for the data scientist. Big Data, Data Warehouse and Data Vault*, Morgan Kaufmann, MA.
 10. Rasmussen, N.H., Goldy, P.S., Solli, P.O. (2002), *Financial business intelligence – trends, technology, software selection and implementation*, John Wiley and Sons, Inc., New York, pp. 98-99.
 11. Sirikulvadhana, S. (2002), *Data mining as a financial auditing tool (thesis)*, Swedish School of Economics and Business Administration, pp. 49-57, disponibil online la adresa: <https://pdfs.semanticscholar.org/2612/f764664796f911e9ff9a79b7bb9de84bf16c.pdf>, accesat pe data de 15.03.2017.
 12. Vintilescu Belciug, A., Crețu, D. și Gegea, C. (2010), Utilizarea tehnicilor de data mining ca metodă complementară în audit, *Audit Financiar*, vol. VIII, nr. 7, pp. 30-35.
 13. Wang, J. și Yang, J.G.S. (2009), Data Mining Techniques for Auditing Attest Function and Fraud Detection, *Journal of Forensic & Investigative Accounting*, vol. 1, no. 1, pp. 1-24.
 14. <http://bi-insider.com/business-intelligence/operational-bi-vs-strategic-bi/>
 15. <http://dssresources.com/faq/index.php?action=artikel&id=174>
 16. <http://dssresources.com/faq/index.php?action=artikel&id=199>
 17. <http://dssresources.com/history/dsshistory.html>
 18. <http://dssresources.com/history/sshistory.html>
 19. <http://inf.ucv.ro/documents/rstoean/5.%20Arbori%20de%20decizie.pdf>
 20. http://profs.info.uaic.ro/~val/statistica/StatWork_10.pdf
 21. <http://searchoracle.techtarget.com/tip/Optimizing-database-performance-part-2-Denormalization-and-clustering>
 22. <http://searchsqlserver.techtarget.com/definition/data-mining>
 23. <http://www.computerweekly.com/tip/Inmon-vs-Kimball-Which-approach-is-suitable-for-your-data-warehouse>
 24. <http://www.techrepublic.com/article/10-tips-for-using-wildcard-characters-in-microsoft-access-criteria-expressions/>
 25. <http://www.w3schools.in/dbms/database-normalization/>
 26. <http://y2u.be/Xs2SWtBqdzl>
 27. <https://deshpande.mit.edu/portfolio/project/hybrid-dbms-optimized-read-intensive-applications>
 28. <https://developers.google.com/apps-script/guides/sheets>
 29. <https://msdn.microsoft.com/en-us/library/bb677216.aspx>
 30. <https://msdn.microsoft.com/en-us/library/dn282385.aspx>
 31. <https://msdn.microsoft.com/en-us/library/ms174828.aspx>
 32. <https://msdn.microsoft.com/en-us/library/office/ee814737.aspx>
 33. <https://onlinecourses.science.psu.edu/stat504/node/149>
 34. <https://sites.google.com/site/supp4excel2dataminig2017af/d>