

Studiu comparat asupra tehnicilor de data mining utilizate în rezolvarea problemelor de regresie și clasificare

Ec. Valentin MILITARU

Catedra de Informatica Economica, A.S.E. Bucuresti

In today's economic context, the corporations and the business environment in general are producing data in enormous quantity and on a daily basis. Data mining defines the process of extracting the knowledge hidden in large volumes of data and the regression and classification tasks are perhaps the most used in such purpose.

This whitepaper presents four widely used data mining techniques (Naïve Bayes, k-Nearest Neighbor, Neural Networks and Decision Trees) and underlines twelve characteristics that differentiate these techniques. The table at the end of this material summarizes the results of this comparative analysis.

Keywords: data mining, data analysis, data mining techniques, regression, classification.

În contextul economic actual, zi de zi, corporatiile și mediul de afaceri în ansamblu produc mari volume de date care sunt utilizate în operațiuni curente. În procesul de analiză, informația care poate fi extrasă din bazele de date poate fi exploatată în continuare pentru construirea unor modele previzionale, pentru identificarea unor relații între înregistrările conținute în baza de date, pentru clasificarea acestor înregistrări sau fie și numai pentru realizarea unei descrieri a conținutului bazei de date. Tehnicile de data mining permit extragerea informațiilor și realizarea de previziuni pornind de la date istorice.

Clasificarea tehnicilor de data mining

Tehnicile de *data mining* au fost grupate în trei categorii, în funcție de tipul de probleme pe care le pot modela:

a) *Clasificarea și regresia* reprezintă cea mai largă categorie de aplicații, constând în construirea de modele pentru previzionarea apartenenței la un set de clase (clasificare) sau a unor valori (regresie). Există câteva tehnici dedicate rezolvării problemelor de clasificare și regresie, dintre care arborii decizionali, tehnica Bayes, rețelele neuronale și k-NN se bucură de o largă recunoaștere.

b) *Analiza asocierilor și succesiunilor*, denumită uneori analiză cosului de comparații; această tehnică generează modele descriptive care evidențiază reguli de corelație între atributele unui set de date.

c) *Analiza de tip cluster* este o tehnică descriptivă utilizată pentru gruparea entităților similare dintr-un set de date sau în egala măsură pentru evidențierea entităților care prezintă diferențieri substanțiale față de un grup. Tehnicile de grupare în clustere se bazează pe algoritmi din sfera rețelelor neuronale, algoritmi demografici, k-NN etc.

Pentru rezolvarea problemelor de clasificare și regresie există o serie de tehnici, iar pentru fiecare tehnică sunt disponibili mai mulți algoritmi. Diferența dintre clasificare și regresie este aceea că în primul caz, outputul previzionat este apartenența la o anumită clasă, în timp ce în al doilea caz, outputul estimează valoarea unui atribut. Regresia este utilizată în cazurile în care outputul este definit pe un domeniu foarte larg (chiar infinit) de valori și nu trebuie confundată cu noțiunea de "regresie liniară" din matematică.

De remarcat că o problemă de regresie poate fi ușor transformată într-una de clasificare și invers. În multe cazuri, instrumentele pot fi utilizate pentru rezolvarea ambelor tipuri de probleme.

● **Tehnică Bayes**¹. Mai puțin implementată în aplicațiile de explorare a datelor, tehnica Bayes naivă este o metodă de clasificare care își datorează numele ministrului britanic Thomas Bayes (1702- 1761). Teoria probabi-

¹ În limba engleză, denumirea tehnicii este însoțită de adjectivul "naiv": *Naïve-Bayes*

litatilor a lui Bayes, pe care se bazeaza si tehnica ce-i poarta numele, a fost publicata postum, abia în 1764. Bayes este o tehnica de clasificare cu potential atât predictiv, cât si descriptiv. Ea permite analiza relatiei dintre fiecare variabila independenta si variabila dependenta, prin calcularea unei probabilitati conditionate pentru fiecare din aceste relatii. Când o noua instanta se doreste a fi clasificata, predictia se realizeaza prin combinarea efectelor variabilelor independente asupra variabilei dependente.

Limitele tehnicii. Pentru instantele care apartin setului de date utilizat la calculul probabilitatilor a priori si al celor conditionale, "predictia" atributului-obiectiv este 100% corecta. Însa pentru instante din afara setului de date de instruire, eficienta algoritmului este puternic afectata de prezenta unor probabilitati conditionate egale sau foarte aproape de zero.

O alta limita a algoritmului provine din asumtia ca între attributele independente din setul de date exista (teoretic) o independenta statistica. Aceasta asumtie sta si la originea adjectivului "naiv" din denumirea algoritmului, având în vedere ca independenta statistica de regula nu se verifica si în practica. Algoritmul este limitat din punct de vedere al inputului la date booleene sau categorice. Dincolo de efortul de preprocesare necesar pentru transformarea datelor cu caracter continuu în intervale valorice, asa cum sa mai mentionat, operatiunea este de multe ori dependentă de experienta si chiar intuitia analistului, factori subiectivi care vor marca rezultatele explorarii.

Avantajele tehnicii. Tinând seama de faptul ca pentru calculul probabilitatilor nu este nevoie decât de o singura parcurgere a setului de date, algoritmul prezinta avantajul important al unei viteze mari de construire a modelului de clasificare.

Ca avantaj semnificativ, algoritmul prezinta capacitatea de a realiza predictii din informatii parțiale. În ciuda sensibilitatii la caracteristici slab reprezentate în setul de date, pentru realizarea unei predictii algoritmul nu are obligatoriu nevoie de toate attributele independente, astfel încât cele identificate de am-

list a fi irelevante pot fi ușor eliminate din algoritm. De fapt, chiar daca nu s-ar cunoaste nimic despre attributele independente, analistul tot ar putea face o predictie (fara pretentia de a fi foarte acurata) numai pe baza probabilitatilor a priori.

Modelul obtinut prin aplicarea algoritmului are si un continut descriptiv, care poate fi util analistului. Probabilitatile conditionate aferente fiecarui atribut independent pot fi utilizate în a descrie legatura dintre acestea si atributul-obiectiv.

● **k-NN.** Tehnica (prescurtare a expresiei engleze *k-Nearest Neighbor*) este predictiva de explorare a datelor utilizata cu precadere în probleme de clasificare. Principiul care sta la baza tehnicii este relativ simplu: o instanta noua este clasificata prin analiza "proximitatii" sale (sau gradului de similitudine) cu alte instante dintr-un set de date cunoscut. k-NN este o tehnica folosita în special pentru clasificarea datelor în categorii multiple, însa poate fi aplicata inclusiv pentru previzionarea unui atribut-obiectiv de natura numerica (continua sau discreta), ca rezultat al unor dependente neliniare. Fie un set de date compus din instante care au urmatoarea structura:

- n attribute numerice independente $\{X_i, i=1, n\}$;
- m attribute booleene sau categorice independente $\{A_j, j=1, m\}$;
- un atribut-obiectiv Y, reprezentând variabila dependenta a carui valoare va trebui estimata pentru noile instante.

Pentru a previziona valoarea atributului-obiectiv al unei instante noi, algoritmul cauta în setul de date k înregistrari "apropiate" de acea instanta, pentru care se cunosc valorile lui Y. Predictia este data de media valorilor lui Y aferente "vecinilor" identificati în setul de date.

Aplicarea conceptului în practica ridica urmatoarele probleme:

- i) Prin ce metoda se stabileste relatia de vecinatate dintre doua instante? Care este metrica utilizata la calculul distantelor dintre doua instante?
- ii) Care este valoarea optima pentru k? De câte instante similare este nevoie pentru ca media atributelor-obiectiv sa se constituie

într-o predicție cu grad acceptabil de reprezentativitate?

iii) Ce metoda de calcul a mediei va fi utilizată pentru predicția lui Y ? Întrebarea este importantă în special pentru cazurile în care atributul-obiectiv nu este de natura numerică, ci booleană sau categorică.

iv) Care din atributele $\{X_i\}$ și $\{A_j\}$ sunt cu adevărat reprezentative pentru predicție? Care din aceste atribute identifică cel mai bine "vecinii" instanței de analizat?

Limitele tehnicii. Timpul de calcul este direct proporțional cu numărul de instanțe din setul de date. Din acest motiv, pentru seturi mari de date se impune ca în etapa de preprocesare, din setul inițial de date să se selecteze un subset de instanțe cu dimensiuni rezonabile. Algoritmul lucrează eficient în probleme de clasificare atunci când toate clasele aferente atributului-obiectiv au o reprezentare egală ca pondere în setul de date, fapt care face necesară "îmbogățirea" setului de date original. Algoritmul pe care se bazează tehnica k -NN permite doar realizarea unei estimări a valorii atributului-obiectiv, fără a produce informații suplimentare despre instanța supusă analizei, despre structura setului de date ori despre categoriile de clasificare a atributului-obiectiv.

De cele mai multe ori este dificil de stabilit ce tip de funcție estimează cel mai bine distanța dintre două instanțe. Deși din punct de vedere matematic tehnica permite calculul distanțelor și pentru atribute categorice și booleene, în astfel de cazuri metrica devine puternic influențată de transformările aplicate de analist setului de date în preprocesare. De aceea, k -NN este de preferat a fi utilizată mai mult în situațiile în care pentru toate atributele instanțelor se poate aplica aceeași funcție de distanță.

Avantajele tehnicii. Tehnica permite clasificarea în multiple clase și modelarea relațiilor neliniare dintre date (în probleme de predicție). Pentru tehnicile care necesită o etapă de învățare a carei output îl constituie un model predictiv, există riscul ca acest model să devină desuet în timp, iar predicțiile realizate în baza lui să piardă din reprezentativitate. În cazul k -NN, modelul îl constituie chiar setul

de date, care se presupune că odată supus analizei, este deja în forma sa cea mai recentă. Chiar dacă de multe ori pot apărea dificultăți în stabilirea unei metrici eficiente, algoritmul este unul dintre putinele care acceptă ca input date de natură diferită (continua, categorică, booleană etc.).

● **Retele neuronale.** Tehnica are la bază două concepte aparținând domeniului inteligenței artificiale. *Neuronul artificial* reprezintă unitatea de bază pentru prelucrarea informației în cadrul calculului neuronal. Prin analogie cu neuronul biologic, el a fost definit ca o unitate ce procesează inputuri informaționale și generează outputuri. *Reteaua neuronală artificială* reprezintă un ansamblu de neuroni artificiali, legați prin conexiuni. Retelele neuronale sunt sisteme dinamice, al căror comportament poate fi caracterizat prin urmărirea stărilor la momente diferite de timp. Starea unei rețele la un moment dat este definită de ansamblul nivelurilor de activare a neuronilor și de intensitățile conexiunilor dintre neuroni. În plus față de acești parametri ajustabili, o rețea este definită și de următorii parametri fiși: configurația conexiunilor și tipul funcțiilor de activare.

Limitele tehnicii. Retelele neuronale nu operează decât direct asupra variabilelor numerice. Drept urmare, orice variabilă non-numerică din setul de date care se dorește analizată va trebui convertită în variabilă numerică înainte de utilizarea sa în instruirea rețelei. În cazul problemelor complexe, utilizatorul este pus în situația de a rezolva un compromis, între a crește numărul de neuroni ascunși, ceea ce poate conduce la o instruire foarte lentă și a accepta o topologie mai simplă, asociată unei soluții mai puțin precise. Pentru seturi de date cu număr mare de atribute, folosirea rețelelor neuronale devine nefezabilă.

Determinarea numărului de neuroni ascunși, pentru probleme complexe de clasificare, nu se poate face decât experimental, ceea ce pe de o parte crește substanțial timpul alocat căutării modelului optim de clasificare, iar pe de altă parte lasă calitatea rezultatelor analizei să depindă de nivelul de experiență al utilizatorului. Absența componentei descriptive

într-un model generat de o retea neuronală face ca evoluția modelului în etapa de instruire să fie lipsită de transparentă pentru utilizator. Datorită acestei caracteristici, tehnica este deseori comparată cu o "cutie neagră". Totuși, cea mai supărătoare caracteristică a rețelelor neuronale este timpul îndelungat necesar pentru o bună instruire, fapt corelat cu necesitatea existenței unui număr relativ mare de instanțe în setul de instruire.

Avantajele tehnicii. Reteaua odată instruită poate realiza predicții rapide pentru instanțe noi. Această caracteristică face ca rețelele neuronale să fie utilizate cu succes în probleme care necesită răspuns în timp real. Până în prezent, rețelele neuronale reprezintă metoda cea mai eficientă de modelare a unor relații neliniare. Mai mult, aplicațiile de până acum au demonstrat aplicabilitatea acestei tehnici în domenii dificil de modelat, precum vederea electronică sau recunoașterea vocală. Spre deosebire de celelalte tehnici de data mining, rețelele neuronale nu restricționează output-ul la un singur atribut. Folosind o arhitectură de rețea potrivită se pot obține predicții simultane pentru mai multe variabile, ceea ce poate însemna o eficientizare semnificativă a proceselor de explorare a datelor.

● **Arbori decizionali.** Arborele decizional este o tehnică de explorare a datelor cu potențial atât predictiv, cât și descriptiv. Denumirea sa provine de la aceea că rezultatul se prezintă utilizatorului sub forma unui graf de tip arbore. Output-ul major al unui model bazat pe arbori decizionali este arborele în sine. Procesul de instruire care creează arborele este numit inductie. Inducția presupune, ca și în cazul rețelelor neuronale, parcurgerea de câteva ori a setului de date de instruire, cu deosebirea că în cazul arborilor, timpul de instruire și implicit numărul de baleieri ale setului de date este mult mai mic decât la rețelele neuronale. Mai precis, numărul de parcurgeri ale setului de instruire este egal cu numărul de niveluri în arbore.

Limitele tehnicii. Majoritatea algoritmilor nu folosesc întregul set de date indicat de utilizator pentru inductie. Pentru acești algoritmi, construirea arborelui presupune transferul instantelor din setul de date de instruire în me-

moria RAM. Dimensiunea limitată a memoriei face ca programul să transfere în RAM numai un subset de date, selectat aleator. În consecință, gradul de reprezentativitate al modelului construit este determinat de capacitatea aplicației de a selecta un subset reprezentativ pentru întreg setul de inductie. O critică adusă frecvent arborilor decizionali este aceea că algoritmi de inductie nu iau în considerare la momentul splitării efectul pe care respectiva separare o are asupra viitoarelor splitări. În plus, toate separările se fac secvențial, ceea ce determină dependența fiecărei splitări de cele precedente.

Avantajele tehnicii. Majoritatea algoritmilor care construiesc arbori decizionali pot fi aplicați fără restricții legate de tipul datelor. Deși variabila dependentă trebuie să fie de natură numerică (în cazul problemelor de regresie) sau categorică (în cazul problemelor de clasificare), pentru majoritatea algoritmilor, variabilele independente pot lua valori în orice domeniu. Tehnica se caracterizează prin capacitate de prelucrare a unor seturi de date cu număr mare de atribute. Există situații în care o instanță poate fi descrisă printr-un număr relativ mare de atribute, de ordinul sutelor sau chiar miilor. În astfel de situații, explorarea prin tehnica arborilor decizionali reprezintă singura alternativă, cei mai mulți algoritmi fiind capabili să trateze seturi de date cu peste 1000 de coloane.

Algoritmi de construire a arborilor decizionali necesită un număr redus de parcurgeri a setului de date utilizat în inductie. Consecința directă a acestei caracteristici funcționale este rapiditatea procesului de inductie și aplicarea eficientă asupra seturilor mari de date.

Forma outputului permite nu numai realizarea de previziuni și clasificări, ci și descrierea relațiilor existente între variabilele independente și variabila dependentă. În plus, forma grafică a outputului facilitează analiza relațiilor. Există aplicații care permit reprezentarea arborelui sub forma unui set de reguli care, pentru arbori de dimensiuni mari, este mai ușor de înțeles.

Concluzie

În funcție de cazul concret, anumite tehnici

de data mining sunt mai eficiente decât altele, existând chiar situații în care pentru rezolvarea problemei nu există decât o unică opțiune (de exemplu arborii decizionali sunt singura alternativă viabilă pentru analiza seturilor de date cu număr mare de variabile, rețelele neuronale reprezintă unică soluție pentru probleme în care outputul are o formă vectorială etc.). Tabelul 1 concentrează caracteris-

ticile prezentate pe larg în acest material, pentru fiecare din cele patru tehnici de data mining tratate urmărindu-se un set de 12 caracteristici. S-a notat cu semnul “+” situația în care tehnica satisface criteriul curent și cu “-” situația contrară. De departe, arborii decizionali prezintă cele mai multe avantaje, în timp ce rețelele neuronale prezintă gradul cel mai mic de flexibilitate.

Tabelul 1 - Analiza comparativă a celor patru tehnici de data mining

Criteriul de comparație	Denumirea tehnicii			
	Naive - Bayes	k-NN	Rețele neuronale	Arbori decizionali
1. Rapiditate în etapa de instruire	+	-	-	+
2. Rapiditate în aplicarea modelului	+	+	+	+
3. Instruire eficientă pe seturi largi de date	+	-	-	+
4. Operare eficientă pe seturi de date cu nr. mare de atribute	-	-	-	+
5. Capacitate de generare a unor outputuri complexe (mai multe atribute simultan)	-	-	+	-
6. Capacitate de generare a unor outputuri de natură vizuală	-	-	-	+
7. Output cu potențial descriptiv	+	-	-	+
8. Utilizare în probleme de predicție	-	+	+	-
9. Utilizare în probleme de clasificare	+	+	+	+
10. Nu comportă restricții legate de tipul datelor de input	-	+	-	+
11. Soluția (modelul) nu depinde de experiența utilizatorului	+	+	-	+
12. Transparența modelului față de utilizator	+	-	-	+
Total	(+7) (-5)	(+5) (-7)	(+4) (-8)	(+11) (-1)

Bibliografie

1. *** - *Critical Features of High performance Decision Trees* - Salford Systems, San Diego (U.S.A.), 2001 (www.salford-systems.com)
2. *** - *IBM's Data Mining Technology – Data Management Solutions* - International Business Machines Corporation, 1996 (www.ibm.com)
3. Bodea, Constanta – Nicoleta – *Inteligenta Artificiala si Sisteme Expert* – Editura Infoc, Bucuresti, 1998
4. Brand, Estelle; Gerritsen, Rob – *Classification and Regression* – DBMS Magazine, Data mining Solutions Supplement, Feb. 1998 (www.dbmsmag.com)

5. Brand, Estelle; Gerritsen, Rob – *Decision Trees* – DBMS Magazine, Data mining Solutions Supplement, Feb. 1998 (www.dbmsmag.com)
6. Brand, Estelle; Gerritsen, Rob – *Naive-Bayes and Nearest Neighbor* – DBMS Magazine, Data mining Solutions Supplement, Feb. 1998 (www.dbmsmag.com)
7. Brand, Estelle; Gerritsen, Rob – *Neural Networks* – DBMS Magazine, Data mining Solutions Supplement, Feb. 1998 (www.dbmsmag.com)