

ARBORII DE DECIZIE O PROVOCARE PENTRU CLASIFICAREA ȘI LUAREA DECIZIILOR ÎN ECONOMIE

LAVINIU AURELIAN BĂDULESCU

Universitatea din Craiova; Facultatea de Automatică, Calculatoare și Electronică, Secția Inginerie Software; Craiova, str. G. Fotino , nr. 2, bl. b8, 3, 11, 200481, Dolj; mobil 0722276554;

E-mail: lavinIU_aurelian_badulescu@yahoo.com

Economists have relied on traditional models when trying to describe how a particular binary choice is related to a set of explanatory variables. There are a variety of "classifiers" like decision trees with superior attributes. Decision tree are used to select the best course of action in situations where you face uncertainty. Many business decisions fall into this category.

1. INTRODUCERE

În mod tradițional, economiștii s-au bizuit pe modele probabilistice și logice atunci când au încercat să explice cum este asociată o alegere binară particulară cu o mulțime de variabile care o descriu. Astfel de modele au o atracție larg răspândită prin aceea că ele pot fi derivate, cel puțin intuitiv, din modele economice de luări de decizii individuale, sunt ușor de estimat, iar rezultatele lor sunt ușor de interpretat.

Însă relativ de curând, specialiștii științei calculatoarelor din domeniul inteligenței artificiale și machine learning au pus în evidență alte instrumente de clasificare cu caracteristici superioare modelelor probabilistice și logice. Aceste instrumentele specifice Data Mining includ algoritmi ca: rețelele neuronale și arborii de decizie. Astfel de algoritmi induc modele puternic neliniare care sunt considerate mai flexibile decât metodele liniare tradiționale și în consecință, furnizează o potrivire mai bună la datele empirice. Astfel de tehnici Data Mining sunt mai performante în clasificarea unui obiect individual într-o categorie particulară pe baza atributelor sale individuale.[5]

Arborii de decizie sunt folosiți pentru a selecta cea mai bună direcție de acțiune în situațiile în care apare incertitudinea. Multe decizii de afaceri intră în această categorie. De exemplu, un producător trebuie să decidă cât de mare să fie stocul creat înainte de a ști precis care va fi cererea. O persoană aflată în litigiu trebuie să aleagă între o înțelegere în afara tribunalului sau riscul unui proces. Un jucător la bursă trebuie să decidă să cumpere înainte de a ști dacă ceea ce a cumpărat poate fi vândut ulterior pentru a obține un profit. În toate aceste cazuri, cel care ia decizia întâlnește un necunoscut care pare să îl facă incapabil să aleagă varianta corectă cu siguranță absolută.

Chiar dacă cel care ia decizia nu cunoaște ce efect va avea factorul necunoscut, el are de obicei niște cunoștințe despre efectele ce pot să apară și cum e mai probabil să apară fiecare efect. Această informație poate fi utilizată pentru a selecta opțiunea care este cel mai

probabil să producă rezultate favorabile. Arborii de decizie fac ușor de aplicat acest tip de analiză.[2]

Algoritmii de arbori de decizie [6] sunt utilizați la explorarea setului de date în problemele de afaceri. Aceasta se realizează adesea prin urmărirea predictorilor și a valorilor care sunt alese la fiecare ramificare a arborelui.

În Figura 1 se evidențiază câteva tipuri de aplicații economice în care se utilizează algoritmi

TIPURI DE APLICAȚII ECONOMICE

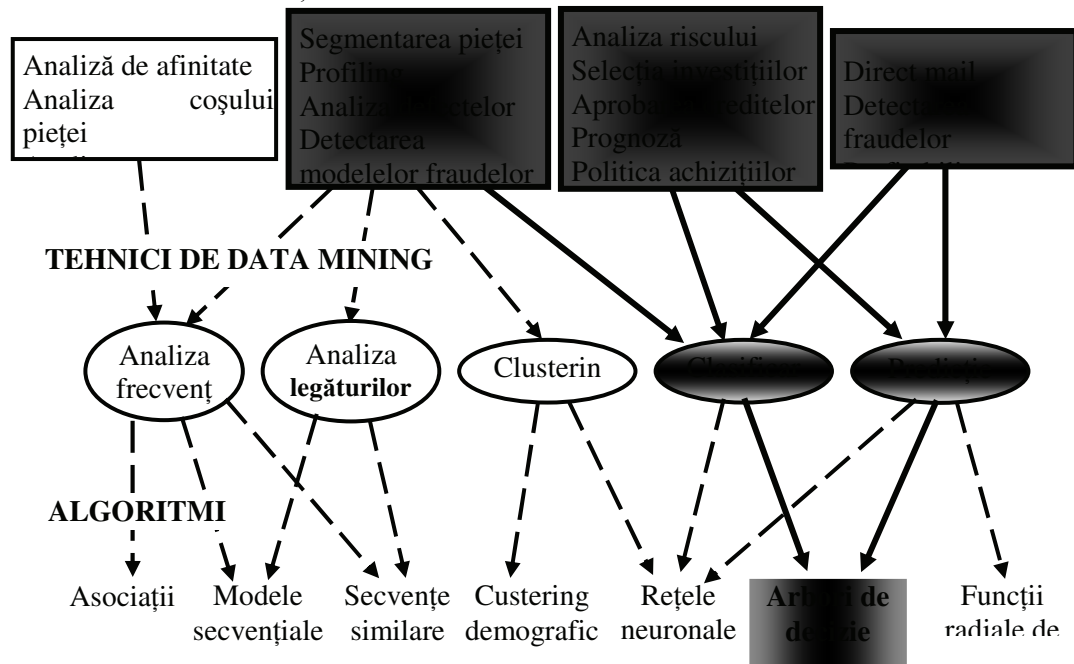


Figura 1.

de Data Mining [3]. Se constată că *algoritmii de arbori de decizie* se folosesc în cazul a două tehnici de Data Mining, respectiv *clasificare* și *predicție*. În același timp aceste două tehnici au utilitate în aproape toate aplicațiile economice pentru care putem folosi Data Mining. Acest fapt evidențiază puterea deosebită a algoritmilor de arbori de decizie pentru aplicațiile economice ce se subsumează aplicării tehnicilor Data Mining.

2. CONSTRUIREA UNUI ARBORE DE DECIZIE DE CALITATE

Să considerăm un arbore de decizie cu M frunze. Acest arbore de decizie corespunde la descompunerea spațiului caracteristic în M subregiuni nesuprapuse E^1, \dots, E^M , astfel încât subregiunea E^S corespunde la frunza numărul S (vezi Figura 2). E^S este definit ca produsul cartezian

$$E^S = E_1^S \times E_2^S \times \dots \times E_n^S, \quad (1)$$

unde E_j^S este proiecția lui E^S pe caracteristica numărul j . E_j^S este obținută în următorul mod. Dacă caracteristica X_j nu este situată pe drumul de la rădăcină la frunza numărul S , atunci E_j^S coincide cu un domeniu al definițiilor caracteristicii X_j . Altfel, E_j^S este egal cu intersecția

tuturor subregiunilor caracteristicii X_j , care a fost întâlnită pe drumul de la rădăcină la frunza numărul S .

Fie o mulțime a unor observații experimentale

$$Data = (x^i, y^i), \quad i = 1, \dots, N. \quad (2)$$

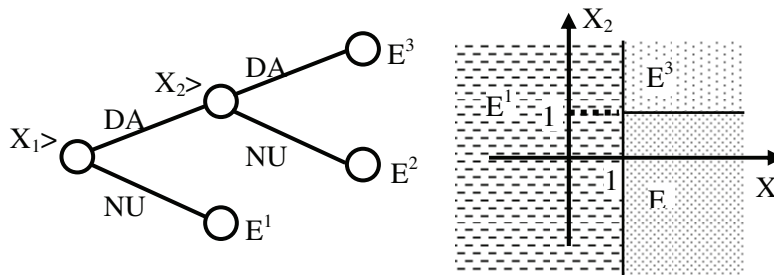
Fiecare din aceste observații aparține (cu privire la X) unora dintre subdomeniile considerate, adică $x^i \in E^S$. Vom nota mulțimea observațiilor aparținând lui E^S prin $Data^S$ și numărul de observații îl vom nota cu N_i^S . Fie N_i^S notația pentru numărul de observații din $Data^S$ aparținând clasei numărul i .

Pentru fiecare scop specific al analizei statistice, există un număr mare de variante diferite de arbori de decizie. Ne punem întrebări de forma: „Care arbore este cel mai bun?” și „Cum să-l găsim?”.

Pentru a răspunde la prima întrebare, vom considera diferite moduri de definire a parametrilor care descriu calitatea unui arbore. Teoretic, putem considera eroarea așteptată a predicției ca parametru de bază. Totuși, această valoare poate fi definită numai dacă cunoaștem legea de distribuție probabilistică a variabilelor examinate. De obicei în practică această lege este necunoscută. În consecință, putem estima calitatea numai aproximativ, utilizând mulțimea observațiilor pe care le avem la dispoziție.

Să presupunem că există un arbore de decizie și un eșantion de obiecte de dimensiune N . Este posibil să alegem două tipuri de bază de parametrii pentru descrierea calității unui arbore. Primul tip este reprezentat de parametrii de acuratețe și al doilea de parametrii de complexitate ai arborelui.

1. *Parametrii de acuratețe* ai arborelui sunt definiți cu ajutorul eșantionului și caracterizează cât de bine sunt împărțite obiectele din clase diferite (în cazul problemei de



Figura

recunoaștere sau clasificare), sau cât de înaltă este eroarea de predicție (în cazul problemei analizei de regresie).

Numărul relativ (frecvența) greșelilor reprezintă partea obiectelor considerate incorect de un arbore ca făcând parte din altă clasă:

$$\hat{P}_{err} = \frac{N_{err}}{N} \quad (3)$$

$$N_{err} = \sum_{S=1}^M \sum_{\substack{i=1 \\ i \neq \hat{Y}(S)}}^K N_i^S \quad (4)$$

unde K este numărul claselor.

Varianța relativă pentru un arbore de regresie poate fi calculată de următoarea formulă:

$$d_{om} = \frac{d_{oc}}{d_0} \quad (5)$$

$$d_{oc} = \frac{1}{N} \sum_{S=1}^M \sum_{i \in Data^S} \left(\hat{Y}(S) - y^i \right)^2 \quad \text{varianța reziduală} \quad (6)$$

$$d_0 = \frac{1}{N} \sum_{i=1}^N \left(y^i - \bar{y} \right)^2 \quad \text{varianța inițială} \quad (7)$$

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y^i \quad (8)$$

2. *Parametrii de complexitate* caracterizează forma arborelui și nu depind de eșantion. De exemplu, parametrii de complexitate ai unui arbore sunt numărul de frunze ale arborelui, numărul de noduri interne ale sale și lungimea maximă a unei căi de la rădăcină la o frunză. De asemenea, este posibil să utilizăm lungimea unei căi externe care este definită ca numărul de ramuri ce trebuie adăugate arborelui până la obținerea unui arbore complet.

3. CONCLUZII

Parametrii de complexitate și acuratețe sunt interconectați: un arbore mai complex, în general, este mai exact, mai corect, adică vom avea acuratețe maximă pentru arborele la care fiecare frunză corespunde la un singur obiect. Însă este de preferat un arbore mai puțin complex. El exprimă dorința de a obține un model mai simplu asupra fenomenului economic cercetat și facilitează interpretarea ulterioară, explicarea modelului. În plus, din cercetările teoretice rezultă că în cazul unui eșantion de dimensiune mică, în comparație cu numărul caracteristicilor, un arbore prea complex devine instabil, adică furnizează o eroare mai mare pentru noi observații.

Pe de altă parte, este clar că un arbore foarte simplu nu va permite obținerea de rezultate bune pentru predicție. Astfel, la alegerea celui mai bun arbore de decizie, trebuie să facă un compromis între parametrii de acuratețe și cei de complexitate.

Pentru a obține o astfel de variantă de compromis, este posibil să utilizăm, de exemplu, următorul criteriu pentru calitatea arborelui:

$$Q = p + \alpha M \quad (9)$$

unde p este un parametru de acuratețe, iar α este un parametru dat. Cel mai bun arbore corespunde la valoarea minimă a criteriului dat. Mai există o variantă utilizată și anume aceea în care este specificată valoarea maxim admisibilă a complexității unui arbore simultan cu căutarea celei mai precise variante de arbore de decizie.[4]

Arborii de decizie reprezintă toate opțiunile și consecințele potențiale într-o manieră care îl face ușor de înțeles și de comunicat situația cu care te întâlnești. Deplasarea de la stânga la dreapta în arbore, este în general o deplasare înainte în timp. Ceea ce se poate întâmpla ca rezultat al alegerii fiecărei opțiuni este ilustrat sub forma unor ramuri adiționale.[1]

BIBLIOGRAFIE

1. ***, „Decision Diagrams”, Vanguard Software Corporation, <http://www.vanguardsw.com/>.
2. ***, „Decision Tree Basics” , Vanguard Software Corporation, <http://www.vanguardsw.com/>.
3. Baragoin, C., Andersen, C., M., Bayerl, S., Bent, G., Lee, J., Schommer, C., „Mining Your Own Business in Retail Using DB2 Intelligent Miner for Data”, International Technical Support Organization, IBM Corporation, RedBook, San Jose, California, 2001, pp. 28-29.
4. Bryant, R., Samaranayake, V., A., Klinkenberg, R., Wilhite, A., „Alternative Models of Choice: the decision to work”, <http://cas.uah.edu/wilhite/papers/alcohol/alcoholchoice.html>.
5. Kamber, M., Winstone, L., Gong, W., Cheng, S., Han, J., „Generalization and Decision Tree Induction: Efficient Classification in Data Mining”, Proc. of 1997 Int'l Workshop on Research Issues on Data Engineering (RIDE'97), Birmingham, England, 1997, p. 111.