

DESCOPERIREA CUNOȘTIȚELOR DIN DATE: METODE PREDICTIVE

Cornel Lepădatu

cornel_lepdatu@biblacad.ro

Biblioteca Academiei Române - București

Rezumat. Obiectivul principal al metodelor predictive îl constituie căutarea de modele optimale pentru diferite metode de modelare: clasice (regresia multiplă, analiza discriminantă), mai puțin clasice (segmentarea) sau de instruire (rețele neuronale, agregarea de modele, mașinile cu suport vectorial). Articolul se concentrează pe prezentarea sub o formă omogenă și sintetică a celor mai frecvent utilizate metode de instruire supervizată pentru descoperirea de cunoștințe din volume (foarte) mari de date (*Big data*, *DCD*) pentru sprijinirea deciziilor în diverse domenii de aplicare. Pentru fiecare metodă au fost evidențiate, după caz, o serie de aspecte specifice esențiale pentru prospectorul de date: domeniile de aplicabilitate, semnificațiile coeficienților, puterea de discriminare a caracteristicilor, metodele de selecție a variabilelor, adecvarea modelului cu datele observate, măsurarea performanțelor, separarea estimării modelului de estimarea erorilor de previziune, controlul supra-învățării, caracterizarea și interpretarea rezultatelor, performanțele computaționale.

Cuvinte cheie: big data, descoperire cunoștințe din date (*DCD*), discriminare, instruire, modelare, previziune.

Abstract: The main objective of predictive methods is the search for optimal models for various modeling techniques: classical (multiple regression, discriminant analysis), less classic (classification and regression tree) or machine learning (neural networks, ensemble methods, support vector machines). The article focuses on an uniform and synthetic presentation for the supervised learning methods the most commonly used for knowledge discovery from (very) large amount of data (*Big data*, *KDD*) for decision support in various fields of application. For each method were highlighted, as appropriate, a number of specific issues, essential for an data prospector: the fields of the application, the significances of the coefficients, the discrimination power of the characteristics, the methods for selection of variables, the appropriateness of the model with the observed data, the performances measurement, the separation of model estimation error from the prediction estimation errors, over-learning control, characterization and interpretation of results, computational performances.

Keywords: Big data, Classification, Knowledge Discovery in Databases (*KDD*), Modeling, Prediction, Statistical learning.

1. Introducere

Mediul economic, social și politic în care se iau în prezent deciziile se caracterizează printr-o dinamică pronunțată și continuă în care tehnologiile avansate devin un determinant major al stilului de viață uman. Numărul căilor de acțiune posibile poate fi foarte mare, gradul de incertitudine poate face foarte dificilă previziunea consecințelor luării unei decizii, efectele unor erori în luarea deciziilor ar putea fi dezastruoase datorită complexității operațiilor și reacțiilor în lanț pe care aceste erori pot să le cauzeze [9, 10].

Convergența procesării informației cu tehnicile de comunicații, ilustrată elocvent mai ales prin dezvoltarea exponențială a Internet-ului, a determinat apariția unor enorme cantități de date, informații și cunoștințe reprezentate în forme din cele mai diverse. Această cantitate imensă de informații este sporită, în continuu, nu doar de dezvoltările permanente ale *web*-ului dar și de apariția agresivă a unor tehnologii emergente precum sistemele dedicate (*embeded*), sistemele mobile și respectiv sistemele omniprezente (*ubiquitous*) de prelucrare a informației [1, 2, 3, 5, 6, 7, 14, 15].

Este, deci, indiscutabil de clară necesitatea extragerii de informații și de cunoștințe, din aceste masive de date distribuite, în primul rând pentru asistarea proceselor decizionale. În acest sens, esențial este faptul că este nevoie de a reprezenta în mod explicit caracteristici importante ale informațiilor, care nu mai sunt legate de reprezentarea abstractă a conceptelor lumii reale ci, mai degrabă, de obiectivul factorilor de decizie și anume susținerea proceselor de analiză a datelor orientate către luarea deciziilor [9, 10, 12].

Menirea sistemelor suport pentru decizii este de a atenua efectul limitelor și restricțiilor factorului decizional în rezolvarea problemelor decizionale. În desfășurarea proceselor decizionale poziția centrală este ocupată de intuiția și judecata umană iar metodele utilizate se bazează pe

analiza datelor disponibile [4, 9, 10, 12, 16, 20, 21, 23].

Principalele concepte și rezultate în domeniul asistării cu mijloace informatice a activităților din procesele decizionale, care presupun analiza datelor, au provenit din prelucrarea analitică on-line (*on-line analytical processing*) și depozitarea datelor (*data warehousing*) precum și din explorarea datelor și descoperirea cunoștințelor (*data mining and knowledge discovery*) [10, 16, 17, 20, 21, 23].

2. Modelare în vederea previziunii

Descoperirea cunoștințelor din baze de date (*DCD*) vizează căutarea de informații relevante pentru previziuni, în vederea susținerii proceselor decizionale, recurgând la metode de instruire (*machine and statistical learning*) care iau în considerare specificitatea volumelor de date, mari și foarte mari.

Având la dispoziție observații asupra unei variabile explicative multidimensionale, $\mathbf{X} = \{X^j\}_{j=1}^p$, măsurată pe o mulțime de n indivizi, în funcție de absența sau prezența unei variabile de explicat Y observată în conjuncție cu \mathbf{X} , se pot distinge două tipuri de probleme, numite de instruire:

- probleme de instruire nesupervizată, absența variabilei de explicat: să se găsească o tipologie sau taxonomie a observațiilor, cum să fie acestea grupate în clase cât mai omogene dar cât mai diferite între ele;
- probleme de instruire supervizată sau de modelare, prezența variabilei de explicat: să se găsească, observându-l pe \mathbf{X} , o funcție φ susceptibilă să reproducă „cel mai bine” (conform unui criteriu ce urmează a fi definit) pe Y , $Y = \varphi(\mathbf{X}) + \varepsilon$, unde ε simbolizează eroarea sau zgomotul de măsurare.

Modelarea sau discriminarea, respectiv, deducerea unui model de previziune pentru o variabilă țintă, constituie obiectivul principal urmărit în demersul inferențial și confirmatoriu [8, 18, 19] pentru aplicațiile *DCD*. Tehnicile cele mai frecvent utilizate în atingerea acestui obiectiv au fost: modelele liniare, analiza discriminantă, rețelele neuronale, mașinile cu suport vectorial, arborii de clasificare și de regresie, agregarea modelelor [2, 8, 12, 13, 22, 24].

Principalele aspecte evolutive privind modelarea în vederea previziunii au fost [2]:

- anii 1940-1970. Statistică: o problemă, asociată cu o ipoteză discutabilă, un experiment planificat cu $p \leq 10$ variabile observate pe $n \approx 30$ indivizi, un model liniar presupus real, un test, o decizie, un răspuns, un volum de date de ordinul hectoocteților.
- anii 1970. Se generalizează primele instrumente *IT*: analiza datelor (*multivariate statistics*) explorează, fără a pretinde un model, volume mai mari de date kiloocteți.
- anii 1980. În inteligența artificială sistemele expert, considerate depășite, sunt înlocuite de instruirea rețelelor neuronale (*ANN learning*) iar în statistică sunt abordate modelele neparametrice sau funcționale, volumul de date megocteți.
- anii 1990. Prima schimbare de paradigmă: datele nu mai sunt planificate, ele sunt culese și stocate pentru activitățile curente ale companiilor dar sunt valorificate pentru sprijinirea proceselor decizionale [9], este momentul apariției marketing-ului cantitativ și al gestionării relațiilor cu clienții (*CRM*), volumul datelor gigaoteți.
- anii 2000. A doua schimbare de paradigmă: numărul p al variabilelor „explodează” ($p \gg n$, $p \approx 105$), obiectivul privind calitatea previziunii prevalează asupra realității modelului (acesta devenind „cutie neagră”), *machine learning* și *statistics* se reunesc în *statistical learning* [13, 25], modelele „bune” sunt acelea care asigură un „echilibru” între bias și varianță, respectiv, o minimizare coordonată a erorilor de aproximare (bias) și a erorii de estimare (varianță), volumul datelor teraocteți.
- anii 2010. A treia schimbare de paradigmă: în multe aplicații numărul n , de indivizi, „explodează”, bazele de date debordează (*big data*) și sunt structurate în nori (*cloud*),

mediile computaționale sunt grupate (*cluster*) dar puterea acestora tot nu este suficientă pentru complexitatea (*greed*) algoritmilor, apare o a treia noțiune privind eroarea, o optimizare indusă prin limitarea timpului de calcul sau a volumului/fluxului de date luate în considerare, decizia devine adaptivă sau secvențială, volumul datelor este de ordinul petaocteților.

- anul 2014. Comisia Europeană programează două domenii de abordare pentru *big data*: îmbunătățirea capacității companiilor de a realiza produse și servicii de date, inovatoare și multilingve și rezolvarea problemelor de cercetare (fundamentală și aplicativă cu orientare către piață) privind capacitățile analiticilor predictive de a susține scalabilitatea și capacitatea de reacție [11].

3. Estimarea calității previziunii

Performanța unui model, rezultat al unei metode de instruire, se evaluează prin capacitatea sa de previziune sau de generalizare. Măsurarea acestei performanțe este foarte importantă pentru prospectorul de date deoarece permite selecția unui model optim dintr-o familie de modele asociată metodei de învățare utilizate, ghidează alegerea metodei comparând modelele selecționate între ele și oferă o măsură a calității sau a încrederii care se poate acorda previziunii.

Estimarea calității previziunii este un element central al oricărei strategii *DCD*. În principiu, sunt avute în vedere trei tipuri de abordări: partiționarea eșantionului pentru a separa estimarea modelului de estimările erorii de previziune, penalizarea erorii de ajustare luând în considerare complexitatea modelului sau recurgerea la simulări implicând multiplicarea calculelor. Alegerea uneia dintre abordări depinde de mai mulți factori între care dimensiunea eșantionului inițial, complexitatea modelului anvizajat, varianța erorii, complexitatea algoritmilor adică volumul de calcule admisibil.

Fie $Y = \varphi(X) + \varepsilon$ modelul de estimat cu ε independent de X , $M(\varepsilon) = 0$ și $var(\varepsilon) = \sigma^2$, fie F legea lui Y în conjuncție cu X , și fie $z = \{(x_i, y_i)\}_{i=1}^n$ un eșantion. Eroarea de previziune a modelului este definită prin: $\mathcal{E}_P(z, F) = M_F [Q(Y, \tilde{\varphi}(X))]$, unde Q este o funcție de pierdere.

Dacă variabila Y de previzionat este cantitativă funcția de pierdere este, în general, pătratică $Q(y, \tilde{y}) = (y - \tilde{y})^2$ iar dacă Y este calitativă Q este un indice de misclasare $Q(y, \tilde{y}) = I_{\{y \neq \tilde{y}\}}$.

În cazul cantitativ eroarea de previziune, într-un punct x , se descompune astfel:

$$\mathcal{E}_P(x) = \sigma^2 + \text{bias}^2 + \text{varianță}.$$

Cu cât un model este mai complex, adică cu un număr mai mare de parametri, cu atât el este mai flexibil, respectiv, se poate ajusta cu atât mai bine la datele observate și deci bias-ul său va putea fi cu atât mai redus. Dar, pe de altă parte, varianța crește odată cu numărul de parametri de estimat adică odată cu complexitatea modelului. Pentru a minimiza riscul pătratic, definit mai sus, soluția este de a căuta un compromis cât mai „bun” între bias și varianță, de a accepta bias-area estimării pentru a reduce cât mai favorabil varianța.

Un criteriu de estimare a erorii de previziune, care exprimă calitatea de ajustare a modelului pe eșantionul observat, este $\mathcal{E}_P = (1/n) \sum_{i=1}^n Q(y_i, \tilde{\varphi}(x_i))$.

În cazul cantitativ, acest criteriu este minimizat prin cercetarea celor mai mici pătrate, în cazul calitativ estimarea este rata de misclasare.

Modul cel mai simplu de a estima, fără bias, eroarea de previziune constă în a calcula \mathcal{E}_P pe un eșantion independent care nu a participat la estimarea modelului.

Dacă dimensiunea eșantionului este suficient de mare, se procedează la separarea eșantionului în trei părți numite de instruire, de validare și de test $z = z_{ins} \cup z_{val} \cup z_{test}$:

- $\tilde{\mathcal{E}}_P(z_{ins})$ este minimizată pentru a estima modelul;
- $\tilde{\mathcal{E}}_P(z_{val})$ servește la compararea modelelor în interiorul unei aceleiași familii pentru a-l

selecționa pe acela care minimizează această eroare;

- $\hat{\mathcal{E}}_P(z_{test})$ este utilizată pentru a compara între ele cele mai „bune” modele ale fiecărei metode considerate.

Dacă dimensiunea eșantionului este insuficientă calitatea ajustării este degradată, varianța estimării erorii poate fi importantă dar nu poate fi estimată și atunci selecția modelului se bazează pe un alt tip de estimare a erorii de previziune recurgându-se la simulare (validare încrucișată) sau la penalizare.

În situațiile în care legea F a eșantionului nu este cunoscută sau, de cele mai multe ori, atunci când nu se poate presupune că este gaussiană, evaluarea distribuției estimatorului se face prin simulare recurgându-se la tehnici de bootstrap (sau reeșantionare). Obiectivul este de a înlocui ipotezele probabilistice, nu totdeauna verificate sau chiar neverificabile, prin simulări implicând mai multe calcule. Ideea de bază a bootstrap constă în substituirea distribuției de probabilitate F aferentă eșantionului de învățare, necunoscută, cu distribuția empirică \tilde{F} obținută acordând o pondere de $1/n$ fiecărei realizări. Astfel se obține un eșantion de dimensiune n , numit eșantion bootstrap, cu legea de distribuție empirică \tilde{F} prin n extrageri aleatoare cu înlocuire dintre cele n observații inițiale. Se poate obține, fără dificultate, un număr mare de eșantioane bootstrap pe care să se calculeze estimatorul respectiv. Legea simulată a acestui estimator este o aproximare asimptotic convergentă, în ipoteze rezonabile, a legii estimatorului. Această aproximare oferă estimări ale bias-ului, ale varianței (deci ale riscului pătratic) și chiar intervalele de încredere ale estimatorului, fără vre-o ipoteză (normalitate) privind legea reală.

Fie $z^n = \{(x_i^n, y_i^n)\}_{i=1}^n$ un eșantion bootstrap al datelor:

- estimatorul plug-in al erorii de previziune, în care distribuția F este înlocuită cu distribuția empirică \tilde{F} , este definit prin: $\mathcal{E}_P(z^n, \tilde{F}) = (1/n) \sum_{i=1}^n nQ(y_i^n, \varphi_{z^n}(x_i^n))$, unde φ_{z^n} reprezintă estimarea lui φ pe z^n ;
- estimarea bootstrap a erorii medii de previziune este dată de: $\mathcal{E}_{boot} = M_{\tilde{F}} [\mathcal{E}_P(z^n, \tilde{F})] = M_{\tilde{F}}[(1/n) \sum_{i=1}^n nQ(y_i^n, \varphi_{z^n}(x_i^n))]$;
- estimarea obținută prin simulare va fi: $\hat{\mathcal{E}}_{boot} = (1/K) \sum_{\kappa=1}^K (1/n) \sum_{i=1}^n nQ(y_i^n, \varphi_{z^\kappa}(x_i^n))$.

Estimarea erorii de previziune astfel construită este, în general, biasată prin optimism deoarece, datorită simulărilor, aceleași observații apar în același timp și în estimarea modelului și în estimarea erorii. Există abordări care vizează corecția acestui bias. Estimatorul out-of-bag al erorii de previziune $\hat{\mathcal{E}}_{oob}$, inspirat din validarea încrucișată, consideră, pe de o parte, observațiile extrase în eșantionul bootstrap și, pe de altă parte, observațiile neutilizate la estimarea modelului dar reținute pentru estimarea erorii:

$$\hat{\mathcal{E}}_{oob} = (1/n) \sum_{i=1}^n 1/B_i \sum_{\kappa \in \mathcal{K}_i} nQ(y_i^n, \varphi_{z^\kappa}(x_i^n))$$

unde \mathcal{K}_i reprezintă mulțimea de indici κ ai eșantioanelor bootstrap neconținând a i -a observație după cele B simulări și B_i reprezintă numărul $|\mathcal{K}_i|$ al acestor eșantioane. B trebuie să fie suficient de mare pentru ca orice observație să poată să fie extrasă cel puțin o dată, altfel termenii cu $\mathcal{K}_i = \emptyset$ trebuiesc omiși.

4. Modele liniare

Modelele liniare urmăresc să prevadă (să explice sau să prezică) o variabilă continuă, numită variabilă de explicat (dependentă sau endogenă) cu ajutorul unor variabile numite explicative (exogene sau predictorii). În cazul în care variabilele explicative sunt continue modelul este un model de analiză a regresiei, dacă acestea sunt variabile discrete (nominale) modelul este de analiză dispersională (sau analiză de varianță) iar dacă mulțimea variabilelor exogene este mixtă modelul este de analiză de covarianță.

Analiza regresiei. În modelul de analiză a regresiei relația dintre Y și X este presupusă liniară, $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ unde: $\mathbf{y} = (y_1, y_2, \dots, y_n)'$, reprezintă vectorul observațiilor asupra variabilei dependente Y , $\mathbf{X} = \{x_{ij}, x_{i0} = 1\}_{i=1}^n \{j=0, \dots, p\}$, este matricea observațiilor asupra variabilelor explicative, $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)'$ reprezintă vectorul coeficienților iar $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$, este vectorul erorilor/reziduurilor. Pentru evaluarea coeficienților necunoscuți ai modelului, inclusiv a reziduurilor ε_i , se dispune de un sistem de n ecuații liniare având $n+p+1$ necunoscute. Sistemul admite o infinitate de soluții; o soluție posibilă $\mathbf{b} = (b_0, b_1, \dots, b_p)$ va trebui să minimizeze global mulțimea distanțelor la modelul liniar conform cu un anumit criteriu de definit; sunt aleși acei vectori \mathbf{b} care minimizează mulțimea valorilor $\{e_i\}_{i=1}^n$, unde $e_i = y_i - (b_0 + b_1x_{i1} + \dots + b_px_{ip})$.

Criteriul celor mai mici pătrate conduce la calcule algebrice simple, se pretează la interpretări geometrice clare și permite interpretări interesante, motiv pentru care se utilizează cel mai des. Estimarea funcției de regresie liniară multiplă presupune determinarea coeficienților b_0, b_1, \dots, b_p prin metoda celor mai mici pătrate pornind de la observațiile $\{y_i, x_{i0} = 1, x_{i1}, \dots, x_{ip}\}_{i=1}^n$. Se presupune că variabilele sunt centrate, ceea ce implică $b_0 = 0$; coeficienții funcției de regresie liniară multiplă sunt $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$. Căutarea lui \mathbf{y} sub forma unei combinații liniare de \mathbf{x}_i se reduce la a defini $\tilde{\mathbf{y}}$ într-un subspațiu V_X generat de variabilele explicative.

Metoda ajustării celor mai mici pătrate se reduce la aproximarea lui \mathbf{y} prin proiecția sa ortogonală $\tilde{\mathbf{y}}$, pe V_X înlocuindu-l pe \mathbf{b} . Se va obține $\tilde{\mathbf{y}} = \mathbf{X}\mathbf{b} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = P_X\mathbf{y}$, unde P_X este operatorul proiecției ortogonale pe V_X . Lungimile în R^n pot fi interpretate în termeni de dispersie deoarece $(1/n) \sum_{i=1}^n y_i^2 = (1/n) \sum_{i=1}^n (y_i - \tilde{y})^2 + (1/n) \sum_{i=1}^n \tilde{y}_i^2$, unde $(1/n) \sum_{i=1}^n y_i^2$, este dispersia totală, $(1/n) \sum_{i=1}^n (y_i - \tilde{y})^2$, este dispersia reziduală și $(1/n) \sum_{i=1}^n \tilde{y}_i^2$ reprezintă dispersia explicată.

Pentru o evaluare globală a calității aproximării se definesc: coeficientul de corelație multiplă, $R = \text{cor}(\mathbf{y}, \tilde{\mathbf{y}}) = \text{cor}(\mathbf{y}, \mathbf{X}\mathbf{b})$ și coeficientul de determinare $R^2 = \sum_{i=1}^n \tilde{y}_i^2 / \sum_{i=1}^n y_i^2$ (adică dispersia explicată împărțită la dispersia totală) sau $R^2 = \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} / \mathbf{y}'\mathbf{y}$ (în funcție de datele inițiale). Dacă $R^2 = 1$ atunci $\tilde{y}_i = y_i$ ($\forall i = 1 \div n$) adică modelul liniar ajustează perfect datele.

Prin minimizarea termenului $\sum_{i=1}^n e_i^2$ se maximizează termenul R^2 , cu alte cuvinte metoda celor mai mici pătrate determină acea combinație liniară a variabilelor explicative ce maximizează corelația cu variabila explicată \mathbf{y} .

Din punctul de vedere al prospectorului de date aspectele cele mai interesante privesc semnificațiile statistice ale coeficienților de regresie, adecvarea modelului regresiei multiple la datele observate, studiul reziduurilor, observațiilor aberante și influenței observațiilor asupra rezultatelor, stabilizarea coeficienților de regresie și tehnicile de obținere de coeficienți stabili precum și metodele de selecție a variabilelor (\mathbf{y} se „explică” doar prin $q \ll p$ predictorii) pentru a micșora numărul de predictorii, a crește viteza de calcul și a obține formule stabile cu o putere predictivă bună.

Analiza dispersională. Dacă variabilele explicative sunt discrete (nominale) regresia multiplă devine analiză dispersională. Se dispune în acest caz de n observații asupra variabilei continue Y observată în conjuncție cu cele p variabile nominale $\{X^k\}_{k=1}^p$ având, respectiv, modalitățile $\{\tau_k\}_{k=1}^p$.

Matricea variabilelor explicative, \mathbf{X} , se prezintă sub forma $[\mathbf{X}_1 \cdot \dots \cdot \mathbf{X}_k \cdot \dots \cdot \mathbf{X}_p]$ adică un tablou disjunctiv complet. Pentru fiecare submatrice \mathbf{X}_k suma coloanelor este egală cu vectorul $\mathbf{1}_n$ existând p relații liniare între coloanele lui \mathbf{X} . Sistemul de ecuații normale are o infinitate de soluții, toate soluțiile duc la același vector $\tilde{\mathbf{y}}$ care este proiecția lui \mathbf{y} pe V_X , dar coeficienții \mathbf{b} nu sunt unici. Pentru a obține o estimare unică \mathbf{b} , trebuie impuse p restricții liniare privind codificările variabilelor calitative. Cea mai des utilizată restricție este ca suma coeficienților lui \mathbf{b} , relativ la fiecare variabilă nominală, să fie nulă, aceasta revine la suprimarea unei coloane din fiecare submatrice și la înlocuirea coloanelor rămase cu diferența dintre ele și coloana suprimată. Matricea $\tilde{\mathbf{X}}$, a variabilelor explicative astfel recodate este de rang maxim, $\text{rang}(\tilde{\mathbf{X}}) = \sum_{k=1}^p (m_k - 1)$.

Pentru exemplificare, în cazul în care se dispune de două variabile nominale A și B , numite factori, având I , respectiv, J modalități, numite nivele, analiza dispersională cu doi factori cu interacțiune se reduce la a efectua regresia variabilei \mathbf{y} cu matricea de condiție $\tilde{\mathbf{X}} = [I \cdot \tilde{\mathbf{X}}_1 \cdot \tilde{\mathbf{X}}_2 \cdot \tilde{\mathbf{X}}_{12}]$ cu $\text{rang}(\tilde{\mathbf{X}}_1) = J$; $\text{rang}(\tilde{\mathbf{X}}_2) = I$; $\text{rang}(\tilde{\mathbf{X}}_{12}) = JK$, unde $\tilde{\mathbf{X}}_1$ și $\tilde{\mathbf{X}}_2$ sunt matricile indicator reduse ale celor doi factori A și B iar $\tilde{\mathbf{X}}_{12}$ este matricea interacțiunilor corespunzând celor JK combinații ale

nivelelor lui A și B . În această situație modelul liniar devine $y = \mu I + \tilde{X}_1 \alpha + \tilde{X}_2 \beta + \tilde{X}_{12} \gamma + \varepsilon$ și deci se poate utiliza un program de regresie multiplă pentru a efectua o analiză dispersională. Procedura poate fi generalizată la modele cu mai mulți factori și nivele de interacțiune de ordin superior.

O anumită prudență se impune, totuși, din mai multe motive: este dificil de apreciat și de limitat clar natura ipotezelor testate; interacțiunile de ordin superior pot duce la „teste în lanț” delicat de interpretat; o interacțiune, mai ales de ordin superior, se poate datora prezenței unor observații ușor aberante caz în care procedura nu este robustă.

Modelele liniare generalizate extind modelele liniare clasice în două direcții: combinația liniară $a_i = b_0 x_{i0} + b_1 x_{i1} + \dots + b_p x_{ip}$ a variabilelor explicative poate fi o funcție ℓ de $M(y_i)$ (numită funcție de legătură), $a_i = \ell(M(y_i))$ în comparație cu modelele liniare obișnuite în care $a_i = M(y_i)$; legea de probabilitate a lui y poate fi și un alt membru al clasei legilor exponențiale (legile binomiale Poisson, Gamma) decât legea normală. Alegând diferite legi de probabilitate din clasa legilor exponențiale și diferite funcții de legătură pentru y , se pot obține și alte modele, printre care un loc important îl ocupă modelele *log*-liniare.

Ajustarea modelelor liniare generalizate se face prin metoda verosimilității maxime care, în cazul legii normale, coincide cu metoda celor mai mici pătrate.

5. Metode de discriminare

Metode geometrice. Metodele geometrice de analiză discriminantă, esențialmente descriptive, se bazează pe noțiunea de distanță și nu utilizează nici o noțiune probabilistă.

Se dispune de observații privind p variabile cantitative $\{X^j\}_{j=1}^p$, jucând rolul de variabile explicative și o variabilă calitativă Y cu q modalități $\{k\}_{k=1}^q$, jucând rolul de variabilă de explicat. Cele p variabile explicative X^j au fost observate pe un eșantion $\{x_i\}_{i=1}^n$, de n indivizi. Variabila nominală Y generează o partiție a celor n indivizi în q clase $\{A_k\}_{k=1}^q$.

Problema de discriminare (sau clasare) este următoarea: fiind dat un nou individ x , pe care au fost observate variabilele explicative X^j dar nu și variabila de explicat Y , se pune problema de a decide modalitatea k a lui Y (sau clasa A_k corespunzătoare) pentru x .

În context geometric, discriminarea poate fi interpretată ca o împărțire a spațiului indivizilor în regiuni, \tilde{R} , numite regiuni de decizie, fiecare regiune fiind asociată cu o clasă de indivizi. Regiunile de decizie și implicit clasele corespunzătoare, se zic separabile dacă pot fi separate prin suprafețe, \tilde{S} , numite suprafețe de decizie. Dacă suprafețele de decizie sunt hiperplane H , clasele se zic liniar separabile. Suprafețele de decizie pot fi descrise cu ajutorul unei mulțimi $G = \{g\}$ de funcții numite funcții de discriminare sau funcții de decizie. Funcția de discriminare g atașează fiecare individ x unei regiuni \tilde{R} , regiune delimitată prin intermediul unei mulțimi de suprafețe de decizie. Funcția de discriminare este instruită într-o fază de instruire când sunt stabilite clasele și suprafețele de decizie. În faza de lucru (sau decizională sau de afectare) funcției de discriminare i se prezintă date ale căror clase nu se cunosc, noii indivizi fiind asociați uneia sau alteia dintre clasele stabilite.

Pentru rezolvarea problemelor de discriminare sunt stabilite reguli de decizie (sau de afectare) și moduri de evaluare. Se disting următoarele trei cazuri de separabilitate:

1. Fiecare clasă A_k este separată de toate celelalte printr-o singură suprafață de decizie. Funcția de decizie corespunzătoare clasei A_k este $g_k(x) : R^p \rightarrow R$, $k \in [1, q]$, ecuația suprafeței de decizie ce separă clasa A_k de toate celelalte clase este: $g_k(x) = 0$.

Pentru fiecare clasă A_k , $[x \in A_k] \rightarrow [g_k(x) > 0]$. Pentru un punct x' , nou, dacă $g_k(x') > 0$ și $g_\ell(x') < 0$, $(\forall) \ell \in [1, q]$, $\ell \neq k$ atunci x' este atașat clasei A_k . Regiunea de decizie \tilde{R}_k , corespunzătoare clasei A_k , este: $\tilde{R}_k = \{x \in R^p \mid [g_k(x) > 0] \wedge [g_\ell(x) < 0], (\forall) \ell \in [1, q], \ell \neq k\}$.

2. Fiecare clasă este separată de oricare alta printr-o suprafață de decizie. Clasele sunt două câte două separabile, cele $q(q-1)/2$ suprafețe de decizie sunt generate de funcțiile $g_{k\ell}(x) : R^p \rightarrow R$ unde $g_{k\ell}(x) = -g_{\ell k}(x)$, $(\forall) x \in R^p$. Suprafața de decizie corespunzătoare claselor A_k și A_ℓ are ecuația

$g_{kl}(x) = 0$, punctele clasei A_k se află de partea pozitivă a suprafeței. Regula de decizie/afectare este: $[x \in A_k] \leftrightarrow [g_{kl}(x) > 0] (\forall) \ell \in [1, q], \ell \neq k$. Regiunea de decizie \check{R}_k corespunzătoare clasei A_k este $\check{R}_k = \{x \in R^p \mid g_{kl}(x) > 0, (\forall) \ell \in [1, q], \ell \neq k\}$.

3. Există q funcții de decizie. Regula de decizie este: $[x \in A_k] \leftrightarrow [g_k(x) > g_\ell(x)], (\forall) \ell \neq k, k \in [1, q]$. Regiunea de decizie \check{R}_k este: $\check{R}_k = \{x \in R^p \mid g_k(x) > g_\ell(x), (\forall) \ell \neq k\}, k \in [1, q]$. Suprafața de decizie dintre clasele A_k și A_ℓ este dată de ecuația: $g_k(x) = g_\ell(x), (\forall) x \in R^p, (\forall) k, \ell \in [1, q], \ell \neq k$. Obiectele clasei A_k se află de partea pozitivă a suprafeței de separare.

Pentru prospectorul de date de o mare importanță practică este cazul claselor liniar separabile. Funcțiile afine de decizie pot fi transformate în funcții liniare de decizie.

Dacă g_k este funcția liniară de decizie corespunzând clasei A_k atunci, în conformitate cu cazul 3 de separabilitate, un obiect x este atașat clasei A_k dacă $g_k(x) > g_\ell(x) (\forall) \ell \in [1, q], \ell \neq k$. În cazul 3 de separabilitate regiunile de decizie pot fi mărginite de hiperplane sau de porțiuni de hiperplane. Clasarea, prin minimizarea unei funcții criteriu, conduce la o clasă de funcții discriminante liniare. Funcția criteriu luată în considerație este distanța d de la vectorii caracteristică la prototipurile claselor. Un vector x este atașat acelei clase A_k de al cărei prototip g_k vectorul x este mai aproape, adică: $[d(x, g_k) = \min_\ell d(x, g_\ell)] \rightarrow [x \in A_k]$.

O clasificare echivalentă se obține considerând funcția de decizie $g_k: R^p \rightarrow R$ dată de formula $g_k(x) = x'g_k - (1/2)g_k'kg_k$. Funcția g_k este o funcție afină de decizie și regula de decizie devine $[g_k(x) = \max_\ell g_\ell(x)] \rightarrow [x \in A_k]$. Hiperplanul de separare este ortogonal pe dreapta ce unește prototipurile claselor, pe care o intersectează într-un punct situat la jumătatea distanței dintre prototipuri. Funcția discriminantă cu distanță minimă este adecvată pentru cazurile când punctele unei clase tind să se aglomereze în vecinătatea unui punct prototip, formând un nor (*cluster*) de puncte.

Metode probabiliste. În abordarea probabilistă, metodele de clasare sunt dedicate aspectului inferențial al analizei discriminante.

Fie (Ω, \mathcal{K}, p) un câmp de probabilitate. Probabilitatea condiționată a evenimentului $A \in \mathcal{K}$ relativ la evenimentul $B \in \mathcal{K}$ cu $p(B) > 0$, este $p_B: \mathcal{K} \rightarrow R$ cu $p_B(A) \equiv p(A|B) = p(A \cap B) / p(B)$.

Dacă $\{A_i\}_{i \in I} \subset \mathcal{K}$ formează un sistem complet de evenimente atunci are loc următoarea egalitate (*formula lui Bayes a probabilității cauzelor*):

$$p(A_i|B) = p(A_i \cap B) / p(B) = p(A_i)p(B|A_i) / (p(A_i)p(B)) = p(A_i) p(B|A_i) / \sum_i p(A_i) p(B|A_i),$$

unde $\{p(A_i)\}$ sunt probabilitățile a priori și $\{p(B|A_i)\}$ probabilitățile a posteriori.

Funcția de repartiție a variabilei aleatoare X condiționată de evenimentul $A \in \mathcal{K}$ cu $p(A) > 0$ este funcția $F_A: R \rightarrow [0, 1], F_A(x) \equiv F(x|A) = p(X \leq x|A)$.

Densitatea de repartiție a variabilei aleatoare X condiționată de evenimentul $A \in \mathcal{K}$ cu $p(A) > 0$ este funcția $f(\bullet|A): R \rightarrow R$ pentru care $F(x|A) = \int_{-\infty}^x f(t|A) dt$.

De asemenea: $f(x|A) = F'(x|A)$ aproape peste tot; $p(A|X=x) = p(A)f(x|A)/f(x)$.

În termenii teoriei statistice a deciziei problema de discriminare se formulează astfel:

Dându-se,

- m grupe (sau populații), $\{\Pi_k\}_{k=1}^m$, specificate prin distribuțiile lor de probabilitate $p_k(x) = p(X=x|x \in \Pi_k)$ cu $k = 1 \div m$;
- m probabilități a priori $\{q_k\}_{k=1}^m$, ca un individ să provină din populațiile Π_k , formând un sistem complet de probabilități ($\sum_{k=1}^m q_k = 1$);
- $\mathcal{E} \subset R^p$, spațiul observațiilor asupra a p variabile aleatoare

$X = \{X^j\}_{j=1}^p$ (predictori);

- $\{C(j|k)\}_{k,j=1}^m$, costurile erorilor de clasare (costul clasării unui individ, provenind din populația Π_k în populația Π_j , $j \neq k$);

să se găsească o partiție \mathcal{P} a spațiului \mathcal{E} astfel încât $\sum_{k=1}^m q_k \sum_{j=1, j \neq k}^m C(j|k) \rho(j|k, \mathcal{P})$ să fie minimă, unde:

$$\mathcal{P} = \{ \check{R}_k \}_{k=1}^m, \cup_{k=1}^m \check{R}_k = \mathcal{E}, \check{R}_k \cap \check{R}_j = \emptyset \quad (\forall) k, j = 1 \div m, j \neq k$$

$$\{ \rho(j|k, \mathcal{P}) = \int_{\check{R}_j} \rho_k(\mathbf{x}) d\mathbf{x} \}_{j=1}^m \}_{k=1, k \neq j}^m \text{ sunt probabilitățile de eroare pentru partiția } \mathcal{P}.$$

Regula Bayes pentru distribuții cunoscute. Se presupun cunoscute probabilitățile a priori $\{q_k\}_{k=1}^m$ și distribuțiile de probabilitate $\{\rho_k\}_{k=1}^m$.

Fie $Y = \{k\}_{k=1}^m$ mulțimea etichetelor claselor și fie $\rho_Y(\ell) = \sum_{k=1}^m q_k \delta_k(\ell)$ distribuția de probabilitate pe Y , unde $\delta_k(\ell)$ este funcția Dirac ($\delta_k(\ell) = 1$ dacă $\ell = k$ și $\delta_k(\ell) = 0$ în rest):

- se numește plasator o funcție $c : \mathcal{E} \rightarrow Y$ ce estimează clasa lui \mathbf{x} , după ce $\mathbf{x} \in \mathcal{E}$ a fost observat $c(\mathbf{x}) = \ell$;
- probabilitatea de misclasare pentru clasa k este: $pmc(k) = \rho[\{c(\mathbf{x}) \neq k \mid \{\mathbf{x} \in \Pi_k\}\}]$;
- funcția de pierdere discretă pentru plasatorul c față de clasa k este notată cu $fpd(c(\mathbf{x}), k)$;
- riscul funcțional al plasatorului c este

$$rf(c) \equiv M_\mu[fpd(c(\mathbf{x}), k)] = \sum_{j=1}^m q_j pmc(j) = \sum_{j=1}^m q_j \sum_{k=1, k \neq j}^m \int_{\check{R}_j} \rho_k(\mathbf{x}) d\mathbf{x}$$

deoarece, distribuția de probabilitate pe $\mathcal{E} \times Y$ este, din construcție, $\mu(\mathbf{x}, k) = q_k \rho_{\ell(\mathbf{x})}(\mathbf{x})$ unde cu $\ell(\mathbf{x}) \in Y$ s-a notat clasa lui \mathbf{x} ;

- dacă se consideră costurile misclasării $\{C(j|k)\}_{k,j=1}^m$ egale cu 1 atunci un plasator va fi optim dacă minimizează $rf(c) = \sum_{k=1}^m q_k \sum_{j=1, j \neq k}^m C(j|k) \rho(j|k, \mathcal{P})$ adică exact funcționala din enunțul problemei de clasare;
- dacă $X = \mathbf{x}$ probabilitatea a posteriori a clasei k este $\rho(k|\mathbf{x}) = q_k \rho_k(\mathbf{x}) / \sum_{k=1}^m q_k \rho_k(\mathbf{x})$.

Partiția lui \mathcal{E} care minimizează riscul funcțional $rf(c)$ este $\mathcal{P} = \{ \check{R}_k \}_{k=1}^m$, unde regiunile de decizie $\check{R}_k = \{ \mathbf{x} \in \mathcal{E} \mid \sum_{j=1, j \neq k}^m q_j \rho_j(\mathbf{x}) \leq \sum_{j=1, j \neq \ell}^m q_j \rho_j(\mathbf{x}), (\forall) \ell \in [1, m], \ell \neq k \}$ sunt numite regiuni de decizie Bayes și se înscriu în cazul 3 de separabilitate.

Dacă $\rho(j|\mathbf{x}) = \max_{1 \leq k \leq m} \rho(k|\mathbf{x})$ atunci plasatorul care minimizează riscul funcțional este notat cu $c_B(\mathbf{x})$. Plasatorul $c_B(\mathbf{x})$ se numește plasator Bayes, riscul funcțional pe care acesta îl minimizează se numește risc (sau eroare) Bayes, iar partiția \mathcal{P} , care determină și este determinată de plasatorul Bayes, se numește procedură de discriminare (sau clasare) bayesiană. Rezultatul fundamental al analizei discriminante probabiliste clasice este: *dacă $\rho(\rho_j(\mathbf{x})/\rho_\ell(\mathbf{x})) = b \mid \mathbf{x} \in \Pi_k = 0, (\forall) j, k, \ell = 1 \div m, \ell \neq j$ și $0 \leq b < \infty$, atunci clasa procedurilor bayesiene este minimală și completă.*

Regula Bayes pentru distribuții cunoscute permite să se construiască o procedură de clasare cu proprietăți de optimalitate dar aplicabilitatea practică directă este redusă deoarece, în realitate, cel puțin distribuțiile $\{\rho_k\}_{k=1}^m$ nu se cunosc.

Regula de decizie Bayes cu parametrii cunoscuți. Se consideră $m = 2$, cazul a două populații normale, multidimensionale $\{\Pi_k\}_{k=1}^2$ caracterizate de densitățile de probabilitate:

$\rho_k(\mathbf{x}) = (1 / ((2\pi)^{p/2} |V|^{1/2})) \exp[(-1/2)(\mathbf{x} - \boldsymbol{\mu}_k)' V^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)]$, adică $[X \in \Pi_k] \rightarrow [X \sim N(\boldsymbol{\mu}_k, V)]$, unde $\boldsymbol{\mu}_k \in M_{p \times 1}(R)$ este vectorul medie și $V \in M_{p \times p}(R)$ este matricea de varianță-covarianță:

- regiunea de clasificare în Π_1 , și anume \check{R}_1 , este mulțimea punctelor $\mathbf{x} \in R^p$ pentru care raportul densităților $\rho_1(\mathbf{x}) / \rho_2(\mathbf{x}) \geq c$, cu c o constantă convenabil aleasă. Condiția de definire a lui \check{R}_1 revine la: $f(\mathbf{x}) \equiv \mathbf{x}' V^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + (-1/2)(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)' V^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \geq \ln c$;

- dacă $\{\Pi_k\}_{k=1}^2$ sunt populații multidimensionale, normal distribuite, de medie μ_i și cu matricea V de varianță-covarianță, comună, atunci cele mai bune regiuni de clasificare sunt date de: $\check{R}_1 \equiv f(x) \geq \ln c$ și $\check{R}_2 \equiv f(x) < \ln c$;
- dacă probabilitățile a priori q_1 și q_2 sunt cunoscute, atunci constanta c este dată de relația $c = q_2 C(1|2) / q_1 C(2|1)$. Dacă $q_1 = q_2$ și $C(1|2) = C(2|1)$ atunci suprafața de separare a celor două regiuni este hiperplanul $H: (\mathbf{g}_1 - \mathbf{g}_2)'(\mathbf{x} - (1/2)(\mathbf{g}_1 + \mathbf{g}_2)) = 0$, unde $\mathbf{g}_k = V^{-1}\mu_k$ este prototipul populației Π_k iar clasificatorul obținut este un clasificator cu distanță minimă;
- dacă probabilitățile a priori nu sunt cunoscute atunci constanta $C = \ln c$ va fi aleasă astfel încât $C(1|2)(1 - \Phi((C + (1/2)\alpha) / \sqrt{\alpha})) = C(2|1)(\Phi((C - (1/2)\alpha) / \sqrt{\alpha}))$, unde $C(k|j)$ sunt cele două costuri ale misclasării, $\alpha = (\mu_1 - \mu_2)'V^{-1}(\mu_1 - \mu_2)$ este distanța Mahalanobis dintre cele două populații iar $\Phi(x) = \int_{-\infty}^x (1/\sqrt{2\pi})e^{-t^2/2} dt$ cu $\varphi(t) = -(t^2/2)$, este funcția de repartiție a variabilei aleatoare Gauss-Laplace.

Regula de decizie Bayes cu parametri necunoscuți. În cazul în care probabilitățile a priori nu sunt cunoscute, se generează o clasă de proceduri admisibile pe bază de estimății.

Dacă $\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{n_i}^{(i)} \in N(\mu_i, V)$, $i \in \{1, 2\}$ sunt două selecții bernoulliene atunci estimatorii $\bar{\mathbf{x}}_i = (1/n_i) \sum_{j=1}^{n_i} \mathbf{x}_j^{(i)}$ și $((n_1 - 1) + (n_2 - 1))\mathbf{S} = (n_1 + n_2 - 2)\mathbf{S} = \sum_{i=1}^2 \sum_{j=1}^{n_i} (\mathbf{x}_j^{(i)} - \bar{\mathbf{x}}_i)(\mathbf{x}_j^{(i)} - \bar{\mathbf{x}}_i)'$ sunt estimatori nedeplasați, de verosimilitate maximă, ai lui μ_i , și V . Pentru selecții suficient de mari folosirea estimățiilor în locul valorilor exacte implică erori mici.

Substituind parametri estimați în relațiile de definiție ale regiunilor de decizie se obține: $\check{R}_1 \equiv f(x) \geq \ln c$ și $\check{R}_2 \equiv \check{f}(x) < \ln c$, unde $\check{f}(x) = \mathbf{x}'\mathbf{S}^{-1}(\bar{\mathbf{x}}^{(1)} - \bar{\mathbf{x}}^{(2)}) - (1/2)(\bar{\mathbf{x}}^{(1)} + \bar{\mathbf{x}}^{(2)})'\mathbf{S}^{-1}(\bar{\mathbf{x}}^{(1)} - \bar{\mathbf{x}}^{(2)})$.

Dacă se dorește clasificarea selecțiilor reunite ca un tot se utilizează următorii estimatori:

$$n = n_1 + n_2, \quad \bar{\mathbf{x}} = (1/n) \sum_{j=1}^n \mathbf{x}_j, \quad [\mathbf{x}_j \in \Pi_1] \vee [\mathbf{x}_j \in \Pi_2]; \quad (n_1 + n_2 + n - 3)\bar{\mathbf{S}} = \mathbf{S} + \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'$$

$$\check{R}_1 \equiv (\bar{\mathbf{x}} - (1/2)(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2))'\mathbf{S}^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \geq c.$$

Prospectorul de date poate obține diverse particularizări ale regiunilor de decizie Bayes pentru diverse valori privind numărul m de populații și numărul p de variabile sau pentru diverși estimatori de verosimilitate maximă definiți în cadrul unor ipoteze compozite.

Estimare bayesiană. În abordările anterioare (frecventiste) s-a presupus o selecție aleatoare dintr-o populație având densitatea de probabilitate $f(\mathbf{x}; \theta)$ cu $\mathbf{x} \in X$ și $\theta \in \Theta$. O procedură de inferență frecventistă depinde de funcția de verosimilitate $L(\theta) = \prod_{i=1}^n f(x_i; \theta)$, unde θ este necunoscut dar fixat.

În demersul bayesian se presupune, a priori, că parametrul necunoscut θ este o variabilă aleatoare având o distribuție de probabilitate proprie pe spațiul Θ al parametrilor, notată $h(\theta)$ și numită distribuția a priorică a lui θ , $f(\mathbf{x}; \theta)$ devenind $f(\mathbf{x}|\theta)$. Distribuția a priorică este, în cazul ideal, fixată înainte de începerea culegerii selecției bernoulliene.

- dacă $f(\mathbf{x}|\theta)h(\theta)$, distribuția comună a lui \mathbf{x} și θ , și $m(\mathbf{x}) = \int_{\Theta} f(\mathbf{x}|\theta)h(\theta)d\theta$, distribuția marginală a lui \mathbf{x} , sunt cunoscute, atunci distribuția lui θ condiționată de evenimentul $\mathbf{X}=\mathbf{x}$ sau distribuția a posteriori a lui θ , este: $h(\theta|\mathbf{x})=h(\theta|X=\mathbf{x}) = f(\mathbf{x}|\theta)h(\theta)/m(\mathbf{x})$, $m(\mathbf{x})>0$, $\mathbf{x} \in \mathcal{E}$, $\theta \in \Theta$.
- dacă $\theta \sim N(\mathbf{m}, \mathbf{S})$ și $\mathbf{x} \sim N(\theta, V)$, atunci $h(\theta|\mathbf{x})$ este densitatea de probabilitate a unei $N(\mu, C)$ cu $\mu = \mathbf{S}(\mathbf{S} + V)^{-1}\mathbf{x} + V(\mathbf{S} + V)^{-1}\mathbf{m}$ și $C = V(\mathbf{S} + V)^{-1}\mathbf{S}$;
- dacă $\theta \sim N(\tau, \sigma^2_0)$ și $\mathbf{x} \sim N(\theta, \sigma^2_1)$, atunci densitatea a posteriori a lui θ este: $N(\mu, \sigma^2)$, unde $\mu = (x/\sigma^2_1 + \tau/\sigma^2_0) / (1/\sigma^2_0 + 1/\sigma^2_1)$ și $\sigma^2 = (\sigma^2_0 \sigma^2_1) / (\sigma^2_0 + \sigma^2_1) = (1/\sigma^2_0 + 1/\sigma^2_1)^{-1}$.

Pentru variabila aleatoare X , cu densitatea de probabilitate $f(x, \theta)$, funcția $T: \Omega \rightarrow R$ se numește statistică suficientă pentru $\theta \leftrightarrow f(x|T(x) = t, \theta) = f(x|T(x) = t) (\forall) t \in \Delta \subseteq R$, adică dacă și numai dacă densitatea de probabilitate condiționată a lui X este independentă de θ .

Fie $X = (x_1, \dots, x_n)$ o selecție bernoulliană asupra unei variabile aleatoare ce depinde de θ și fie $\tilde{\theta} \equiv \tilde{\theta}(T)$ un estimator a lui θ :

- funcția de pierdere, estimând θ prin $\tilde{\theta}$, este: $L^b(\theta, \tilde{\theta}) \equiv L^b(\theta, \tilde{\theta}(T)) = (\tilde{\theta}(T) - \theta)^2$;
- riscul funcțional este: $R^b(\theta, \tilde{\theta}) = M[L^b(\theta, \tilde{\theta})] = \int_{\Delta} L^b(\theta, \tilde{\theta}(t)) f(t|\theta) dt$;
- se numește risc bayesian: $r^b(\theta, \tilde{\theta}) = \int_{\Theta} R^b(\theta, \tilde{\theta}) h(\theta) d\theta$;
- se numește estimator bayesian: $r^b(\theta, \tilde{\theta}^b) = \inf_{\tilde{\theta} \in B} r^b(\theta, \tilde{\theta})$, $\tilde{\theta}^b \in B$, unde B este clasa estimatorilor pentru care riscul bayesian este finit;
- în cazul funcției de pierdere „suma pătratelor erorilor” estimatorul bayesian este dat de relația: $\tilde{\theta}^b(t) = \int_{\Theta} \theta h(\theta|t) d\theta \equiv M[\theta|T(x) = t]$, adică media distribuției à posteriori $h(\theta|t)$ pentru toate valorile posibile observate, $t \in \Delta$.

Fie x_1, \dots, x_n variabile aleatoare independente și identic repartizate $N(\theta, \sigma^2)$, cu θ necunoscut și $\sigma_1 > 0$ dat și fie statistica $T = (1/n)\sum_{i=1}^n x_i$, care este suficientă pentru θ :

- dacă distribuția a priori a lui θ pe spațiul $\Theta = R$ este $N(\tau, \sigma_0^2)$ cu $\tau, \sigma_0 \in R$ dați și $\sigma_0 > 0$, distribuția à posteriori a lui θ , condiționată de observațiile x_1, \dots, x_n , este $N(\mu, \sigma^2)$, unde $\mu = ((n\sigma_0^2)/(n\sigma_0^2 + \sigma_1^2))T(x) + ((\sigma_1^2)/(n\sigma_0^2 + \sigma_1^2))\tau$ și $\sigma^2 = (\sigma_0^2 \sigma_1^2)/(n\sigma_0^2 + \sigma_1^2)$;
- dacă $\sigma_0 = 0$, atunci $\mu = \tau$ indiferent de observațiile efectuate;
- dacă $\sigma_0 > \sigma_1$ rezultă $\mu \approx \bar{x}$, cunoașterea mediei à priorice τ este de importanță redusă;
- raportul $a = \sigma_1^2 / \sigma_0^2$ măsoară încrederea à priori că τ este o estimare corectă a mediei;
- dacă $a < \infty$ atunci $\lim_{n \rightarrow \infty} \mu = \lim_{n \rightarrow \infty} \bar{x}$;
- dacă dispersia inițială este mică, media estimată tinde să rămână în apropierea mediei inițiale τ chiar dacă media empirică \bar{x} diferă considerabil de aceasta;
- dacă raportul a este mic, atunci media și dispersia à priori au doar o influență redusă asupra estimării parametrilor care sunt determinați aproape exclusiv din datele empirice;
- dacă $T(x) = t$, atunci estimatorul Bayes al mediei unei variabile aleatoare $N(\mu, \sigma^2)$ este:
 $\tilde{\theta}(t) = \tilde{\theta}_B = (nt/\sigma_1^2 + nt/\sigma_0^2) (1/\sigma_1^2 + 1/\sigma_0^2)^{-1}$;
- pentru cazul multidimensional: $\tilde{\theta}_B = \mathbf{S}(\mathbf{S} + (1/n)\mathbf{V})^{-1} \mathbf{t} + (1/n)\mathbf{V}(\mathbf{S} + (1/n)\mathbf{V})^{-1} \mathbf{m}$.

Fie $X = (x_1, \dots, x_n)$ o selecție bernoulliană din populațiile Π_1 și Π_2 :

- dacă $X \in \Pi_k$, atunci densitatea de probabilitate este $f_k(x|\theta)$, $\theta \in \theta_k$ și densitatea à priorică este $h_k(\theta)$, $k = \{1, 2\}$;
- dacă q_1 și q_2 sunt probabilitățile à priori ale populațiilor Π_1 , și Π_2 , probabilitățile à posteriori sunt: $\mathcal{P}(\Pi_k|x) = m_k(x)q_k / (m_1(x)q_1 + m_2(x)q_2)$, unde $m_k(x) = \int_{\Theta_k} f_k(x|\theta) h_k(\theta) d\theta$, este densitatea de probabilitate marginală a lui x condiționat de faptul că provine din Π_k ;
- procedura bayesiană de discriminare este:

$$\mathbf{x} \in \begin{cases} \Pi_1 & \text{dacă } \mathcal{P}(\Pi_1|x) / \mathcal{P}(\Pi_2|x) = (q_1/q_2)B_{12}(\mathbf{x}) \geq 1 \\ \Pi_2 & \text{în caz contrar} \end{cases}$$

unde $B_{12}(\mathbf{x}) = m_1(\mathbf{x})/m_2(\mathbf{x})$, fiind cunoscut ca factorul Bayes al populației Π_1 versus Π_2 .

6. Mașini cu suport vectorial

Mașinile cu suport vectorial reprezintă o clasă de algoritmi de instruire destinați, inițial, problemelor de discriminare adică de predicție unei variabile calitative. Ulterior, algoritmi au fost generalizați pentru a prezice o variabilă cantitativă adică de a găsi o funcție de discriminare (sau

clasificator) a cărei capacitate de generalizare (sau calitate a predicției) să fie cea mai mare posibilă. Abordarea s-a concentrat pe proprietățile de generalizare (sau de previziune) ale modelului controlându-i complexitatea, mai precis, integrând în estimare numărul de parametri, în acest caz numărul de vectori suport. Ideea de bază al mașinilor cu suport vectorial a fost de a reduce problema discriminării la o problemă, liniară, de căutare a unui hiperplan optimal:

- prin definirea hiperplanului optimal ca soluție a unei probleme de optimizare cu restricții, în care funcția obiectiv se exprimă numai cu ajutorul produselor scalare între vectori iar numărul de restricții „active” (vectorii suport) controlează complexitatea modelului;
- prin căutarea unor suprafețe de separare neliniare;
- prin introducerea unei funcții nucleu în produsul scalar inducând implicit o transformare neliniară a datelor către un spațiu hilbertian, intermediar, de dimensiune mai mare și în care este rezolvată problema liniară.

Fie Y variabila de explicat, fie $X = \{X^j\}_{j=1}^p$ o variabilă explicativă sau de predicție multidimensională și fie φ un model pentru Y :

- $X = \{X^j\}_{j=1}^p$ este o variabilă cu valori într-o mulțime $\mathcal{E} \subset R^p$, $\mathbf{x} = (x_j)_{j=1}^p \in \mathcal{E}$;
- Y este dicotomică, $\varphi : \mathcal{E} \rightarrow B$, $\varphi(\mathbf{x}) \in B = \{-1, 1\}$;
- $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^n$ este un eșantion de mărime n și de lege F necunoscută;
- obiectivul este de a construi o estimare $\tilde{\varphi} : \mathcal{E} \rightarrow \{-1, 1\}$ a lui φ astfel încât probabilitatea $\rho(\varphi(X) \neq Y)$ să fie minimă.

Problema revine la a căuta o frontieră de decizie în spațiul \mathcal{E} pentru valorile lui X și la a găsi un compromis între complexitatea acestei frontiere, respectiv, capacitatea de ajustare a modelului, și calitățile de generalizare (sau de previziune) ale modelului.

Demersul constă în a găsi o funcție reală f al cărui semn să ofere previziunea: $\tilde{\varphi} = \text{sign}(f)$. Eroarea de previziune se exprimă prin cantitatea: $\rho(\varphi(X) \neq Y) = \rho(Yf(X) \leq 0)$. Valoarea absolută a acestei cantități, $|Yf(X)|$, furnizează o indicație privind încrederea care poate fi acordată rezultatului clasării. Se spune că $Yf(X)$ este marja lui f în (X, Y) .

Primul pas este de a transforma valorile lui X , adică obiectele din \mathcal{E} , prin funcția $\Phi : \mathcal{E} \rightarrow \tilde{H}$ cu valori într-un spațiu \tilde{H} , intermediar, înzestrat cu un produs scalar. Această transformare, fundamentală pentru abordarea SVM, ia în considerare eventuala neliniaritate a problemei de rezolvat și conduce la rezolvarea unei separări liniare. În cazul în care Φ este funcția identitate (adică în cazul liniar), atunci când separarea este posibilă, dintre toate hiperplanele, soluții de separare a observațiilor, se alege acela care este situat „cel mai departe” de toate exemplele, adică de marjă maximală.

Cu produsul scalar al spațiului \tilde{H} , un hiperplan H este definit prin ecuația $\langle \mathbf{w}, \mathbf{x} \rangle + b = 0$, unde \mathbf{w} este un vector ortogonal pe hiperplan, $\mathbf{w} \perp H$, iar semnul funcției $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$ arată de care parte a hiperplanului este situat punctul \mathbf{x} de explicat:

- un punct \mathbf{x} este bine clasat $\leftrightarrow yf(\mathbf{x}) \geq 1$;
- un hiperplan $H \equiv (\mathbf{w}, b)$ este un separator dacă: $y_i f(\mathbf{x}_i) \geq 1 \ (\forall) i \in [1, n]$;
- distanța de la un punct \mathbf{x} la (\mathbf{w}, b) este: $d(\mathbf{x}) = |\langle \mathbf{w}, \mathbf{x} \rangle + b| / \|\mathbf{w}\| = |f(\mathbf{x})| / \|\mathbf{w}\|$;
- marja hiperplanului are valoarea $2 / \|\mathbf{w}\|^2$.

Căutarea hiperplanului separator de marjă maximală revine la rezolvarea problemei (primare) de optimizare cu restricții: $(1/2) \min_{\mathbf{w}} \|\mathbf{w}\|^2 \cdot y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \ (\forall) i$. Problema duală se obține prin introducerea multiplicatorilor Lagrange. Soluția este furnizată de un punct \mathbf{w}^* , b^* , λ^* al lagranjianului $L(\mathbf{w}, b, \lambda) = (1/2) \|\mathbf{w}\|^2 - \sum_{i=1}^n \lambda_i (y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1)$, punctul \mathbf{w}^* verificând condițiile $\lambda_i^* [y_i (\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) - 1] = 0 \ (\forall) i \in [1, n]$.

Vectorii suport sunt vectorii \mathbf{x}_i pentru care restricția este activă (cele mai apropiate de hiperplan) adică verifică: $y_i(\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) = 1$. Condițiile de anulare a derivatelor parțiale permit exprimarea formulei duale a lagranjianului: $W(\lambda) = \sum_{i=1}^n \lambda_i - (1/2) \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$. Pentru a găsi punctul de șa, se maximizează $W(\lambda)$, $\lambda_i \geq 0$ ($\forall i \in [1, n]$). Rezolvarea acestei probleme de optimizare pătratică, de dimensiune n , furnizează ecuația hiperplanului optimal:

$$\sum_{i=1}^n \lambda_i^* y_i = \langle \mathbf{x}, \mathbf{x}_i \rangle + b^* = 0 \text{ cu } b_0 = -(1/2)(\langle \mathbf{w}^*, \mathbf{SV}_{clasa+1} \rangle + \langle \mathbf{w}^*, \mathbf{SV}_{clasa-1} \rangle).$$

Pentru o nouă observație \mathbf{x} prezentată modelului, este suficient să se vadă semnul expresiei $f(\mathbf{x}) = \sum_{i=1}^n \lambda_i^* y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b^*$ pentru a ști în care semi-spațiu se află \mathbf{x} și deci ce clasă i se va atribui.

Dacă observațiile nu sunt separabile printr-un hiperplan atunci se recurge la o „relaxare” a restricțiilor introducându-se termenii de eroare, ξ_i , $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq 1 - \xi_i$ ($\forall i \in [1, n]$), care controlează depășirile. Modelul va oferi un răspuns greșit pentru un vector \mathbf{x}_i dacă valoarea termenului de eroare corespunzător este mai mare decât 1, $\xi_i > 1$. Introducând o penalizare π pentru încălcarea restricțiilor, problema de minimizare se reformulează în felul următor:

$$\min(1/2) \|\mathbf{w}\|^2 + \pi \sum_{i=1}^n \xi_i \text{ cu } y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq 1 - \xi_i, (\forall i \in [1, n]).$$

Problema se formulează în aceeași formă duală ca și în cazul separabilității cu o singură diferență: coeficienții λ_i sunt mărginiți de constanta π de control a penalizării. Din punctul de vedere al prospectorului de date parametrul π , care controlează penalizarea, trebuie „bine” ales fiind parametrul care reprezintă compromisul între o bună ajustare și o bună generalizare. Cu cât el este mai mare cu atât importanța atribuită ajustării modelului este mai puternică.

Observațiile făcute în mulțimea \mathcal{E} sunt transformate prin aplicația neliniară $\Phi : \mathcal{E} \rightarrow \tilde{H}$, spațiul \tilde{H} fiind de dimensiune mai mare și înzestrat cu un produs scalar. Formularea problemei de minimizare și soluția sa $f(\mathbf{x}) = \sum_{i=1}^n \lambda_i^* y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b^*$ implică numai elementele \mathbf{x} și \mathbf{x}' , prin intermediul produsului scalar $\langle \mathbf{x}, \mathbf{x}' \rangle$. Prin urmare, nu ar mai fi necesară explicitarea transformării Φ , ceea ce de multe ori este imposibil, cu condiția de a dispune de o exprimare a produselor scalare în \tilde{H} cu ajutorul unei funcții $\kappa : \mathcal{E} \times \mathcal{E} \rightarrow R$, simetrică, numită nucleu (*kernel*), astfel încât: $\kappa(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$. Convenabil ales, nucleul permite materializarea unei noțiuni de „proximitate”, adaptată problemei de discriminare și structurii sale de date. Pentru construirea de funcții nucleu se recurge la combinații ale unor nuclee simple liniare $\kappa(\mathbf{x}', \mathbf{x}'') = \langle \mathbf{x}', \mathbf{x}'' \rangle$, polinomiale $\kappa(\mathbf{x}', \mathbf{x}'') = (c + \langle \mathbf{x}', \mathbf{x}'' \rangle)^d$ sau gaussiene $\kappa(\mathbf{x}', \mathbf{x}'') = e^{-\alpha \langle \mathbf{x}', \mathbf{x}'' \rangle}$, unde $\alpha \langle \mathbf{x}', \mathbf{x}'' \rangle = \|\mathbf{x}' - \mathbf{x}''\|^2 / 2\sigma^2$, pentru a se obține nuclee mai complexe (multidimensionale) asociate cu situația întâlnită. Pentru prospectorul de date, o mare flexibilitate în definirea nucleelor, care să permită definirea unor noțiuni adecvate de similitudine, conferă mai multă eficacitate acestei abordări cu condiția, desigur, de a construi și a testa un nucleu „bun”. Rezultă, din nou, importanța unei evaluări corecte a erorilor de previziune, de exemplu, prin validare încrucișată.

7. Metode conexioniste

O rețea neuronală este asocierea într-un graf, mai mult sau mai puțin complex, a neuronilor formali. Neuronul formal este un model al neuronului biologic, se caracterizează prin: stări interne, $s \in S$, semnale de intrare $\{x_i\}_{i=1}^p$, funcția de tranziție a stărilor $s = h(x_1, \dots, x_p) = f(\beta_0 + \sum_{j=1}^p \beta_j x_j)$. Valorile coeficienților $\{\beta_j\}_{j=0}^p$ sunt estimate într-o fază de instruire și constituie „memoria” sau „cunoașterea distribuită” a rețelei, coeficientul β_0 este numit bias al neuronului. Rețelele neuronale sunt caracterizate prin organizarea grafului (în straturi), prin numărul de neuroni și prin tipul neuronilor, respectiv, funcțiile lor de tranziție. Perceptronul multistrat este o rețea formată din straturi succesive de neuroni formali; stratul este un set de neuroni fără nici-o legătură între ei; stratul de intrare citește semnalele $\{x_j\}_{j=1}^p$ de intrare și conține câte un neuron pentru fiecare intrare x_j ; unul sau mai multe straturi ascunse participă la transfer, un neuron al unui strat ascuns este conectat la intrare cu fiecare dintre neuronii stratului precedent și la ieșire cu fiecare neuron al stratului următor; stratul de ieșire furnizează răspunsul sistemului. Un perceptron multistrat realizează o transformare $y = \varphi(x_1, \dots, x_p; \boldsymbol{\beta})$ unde $\boldsymbol{\beta}$ este vectorul conținând parametrii β_{jkl} corespunzători intrării j a neuronului k din stratul l ; stratul de intrare ($l = 0$) nu este parametrizat

pentru că nu face altceva decât să distribuie intrările în neuronii din stratul următor.

Intrările rețelei $\{x_i\}_{i=1}^p$, sunt variabilele explicative ale modelului, ieșirea y este variabila de explicat (dependentă sau țintă) iar β , vectorul ponderilor intrărilor în fiecare neuron al rețelei, reprezintă parametrii de estimat în urma unui proces de instruire.

Pentru un eșantion de instruire $\{(x^1_i, \dots, x^p_i; y_i)\}_{i=1}^n$ construit din n observații asupra a p variabile explicative $\{X^j\}_{j=1}^p$ și a unei variabile de explicat Y , instruirea constă în estimarea vectorului de parametri β rezolvând o problemă a celor mai mici pătrate:

$$\tilde{\beta} = \min_b Q(\mathbf{b}), \quad Q(\mathbf{b}) = (1/n) \sum_{i=1}^n (y_i - \varphi(x^1_i, \dots, x^p_i; \mathbf{b}))^2.$$

Algoritmul de optimizare cel mai utilizat este un algoritm de retropropagare (propagare inversă) a gradientului bazat pe faptul că în orice punct \mathbf{b} vectorul gradient al lui Q este orientat în direcția de creștere a erorii și deci pentru a-l descrește pe Q este suficientă o deplasare în sens contrar. Pornind de la erorile observate pe ieșiri, formula retropropagării erorii furnizează expresia erorii atribuite fiecărei intrări, de la stratul de ieșire către stratul de intrare. Proprietățile acestui algoritm implică o convergență aproape sigură, probabilitatea de atingere a unei precizii dorite (fixate a priori) tinde către 1 atunci când dimensiunea eșantionului de instruire tinde către infinit.

În practică, prospectorul de date se confruntă cu o serie de opțiuni privind, în principal, controlul suprainsurii: alegerea unor parametri (limitarea numărului de neuroni, limitarea duratei de instruire, creșterea coeficientului de penalizare a normei parametrilor); alegerea modului de estimare a erorii (pe eșantionul de test sau validare încrucișată).

8. Metoda segmentării

Metoda segmentării este o metodă complementară de rezolvare a problemelor de discriminare și de regresie prin împărțirea progresivă a eșantionului de observații într-un arbore de decizie binară.

Fie \mathbf{y} variabila privilegiată, discretă, cu q modalități, $\{k\}_{k=1}^q$, care este explicată prin variabilele, cantitative sau calitative, $\{X^j\}_{j=1}^p$, și fie $\{\mathbf{x}_i; y_i\}_{i=1}^n \equiv \{\{x^j_i\}_{j=1}^p; y_i\}_{i=1}^n$ eșantionul observațiilor, unde $y_i \in \{k\}_{k=1}^q$.

Metoda de segmentare constă, mai întâi, în a căuta variabila X^j care, explică „cel mai bine” variabila \mathbf{y} și definește o împărțire a eșantionului în două submulțimi de indivizi, numite segmente sau noduri. Apoi, se reiterează procedeul căutându-se cea mai bună variabilă în interiorul fiecăruia dintre cele două segmente definite, ș.a.m.d. Prin împărțirea succesivă a eșantionului în câte două submulțimi rezultă un arbore de decizie binară în care se disting: segmente intermediare, segmente terminale, ramuri ale unui segment, arborele binar complet, A_{max} , și subarbori. Efectuarea diviziunii unui nod se face astfel încât cele două segmente descendente să fie mai omogene decât nodul părinte și cât mai diferite între ele față de variabilă.

Fazele de construire ale arborelui de decizie binară sunt: stabilirea, pentru fiecare nod, a mulțimii diviziunilor admisibile; definirea unui criteriu de selecționare a „cele mai bune” diviziuni a fiecărui nod; definirea unei reguli care să permită declararea unui nod ca terminal sau intermediar; afectarea fiecărui nod terminal unei clase; estimarea riscului de misclasare.

Inițial, există un singur segment conținând toți indivizii \mathbf{x}_i , $i = 1 \div n$. Sunt examinate, secvențial, toate variabilele explicative X^j , $j = 1 \div p$. În funcție de natura fiecărei variabile X^j (continuă sau discretă) se definesc toate diviziunile posibile. O diviziune posibilă este admisibilă dacă segmentele descendente sunt nevide. Dintre toate diviziunile admisibile \tilde{d}_m , unde m reprezintă a m -a diviziune (sau a m -a valoare ordonată a variabilei din eșantion), este selecționată diviziunea \tilde{d}^j „cea mai bună” în sensul unui criteriu de impuritate. Astfel, pentru fiecare din cele p variabile, se obține diviziunea optimă „locală” \tilde{d}^j și, în final, din cele p diviziuni se va reține diviziunea \tilde{d} , care va furniza cele două segmente „cele mai caracteristice” vis-à-vis de \mathbf{y} . Procedeul se aplică iterativ fiecărui segment descendent obținut și se oprește când toate segmentele sunt declarate terminale. Afectarea unui individ nou se face prin „coborârea” lui pe ramurile arborelui.

Fie $\rho(r|a)$ probabilitatea condiționată de apartenență la grupul G_r , $r \in \{1, 2, \dots, q\}$ a mulțimii

observațiilor din nodul a . Impuritatea unui nod, a , este o funcție nenegativă de $\{\rho(r|a)\}_{r=1}^q$, care verifică următoarele condiții: este maximală când probabilitățile de apartenență la diferite grupuri sunt egale între ele: $\rho(r|a) = 1/q, (\forall)r$; este nulă dacă nodul conține observații aparținând unui singur grup: $\rho(r|a) = 1$ și $\rho(s|a) = 0, (\forall)s, s \neq r$; este o funcție simetrică de probabilitățile $\rho(r|a)$. Funcțiile de impuritate, cele mai des utilizate, sunt: $i(a) = -\sum_{r=1}^q \rho(r|a) \ln(\rho(r|a))$, funcție derivată din noțiunea de entropie Shannon și indicele de diversitate Gini $j(a) = -\sum_{r \neq s} \rho(r|a) \rho(s|a)$.

Fie ∂ o diviziune admisibilă care împarte/divide nodul a în segmentele t_s și t_d cu probabilitățile: $\rho_s \equiv \rho(t_s|a) = \rho(t_s)/\rho(a)$ și respectiv $\rho_d \equiv \rho(t_d|a) = \rho(t_d)/\rho(a)$. Reducerea impurității nodului a datorată diviziunii ∂ este definită prin expresia: $\Delta_i(\partial, a) = i(a) - \rho_s i(t_s) - \rho_d i(t_d)$. Orice diviziune, ∂ , a unui nod, a , duce la o reducere pozitivă sau nulă a impurității. Cea mai „bună” diviziune este $\tilde{\partial} = \text{argmax}_{m \in \partial_j} \Delta_i(\partial_m, t)$ adică aceea pentru care reducerea impurității este maximă, unde ∂ este mulțimea diviziunilor admisibile ale variabilei X^j . Pe mulțimea $\{X^j\}_{j=1}^p$, a variabilelor explicative, diviziunea nodului t este efectuată cu ajutorul variabilei X^j care asigură $\tilde{\partial} = \max_{1 \leq j \leq p} \{\tilde{\partial}^j\}$.

În procesul de construire a lui A_{max} este posibil ca toate nodurile terminale, a , ale arborelui curent, A , să fie afectate unuia din cele q grupuri (sau clase). Fiecărei erori de clasare i se asociază un preț de misclasare $\gamma(s/r), s, r = 1, 2, \dots, q$ costul misclasării fiind $\sum_{r=1}^q \gamma(s/r) \rho(r|a)$.

Un nod a va fi asignat acelei clase \tilde{s} pentru care $\tilde{s} = \min_{1 \leq s \leq q} \sum_{r=1}^q \gamma(s/r) \rho(r|a)$. Dacă minimum este atins pentru cel puțin două clase atunci nodul este afectat arbitrar uneia dintre aceste clase. Dacă $\gamma(s/r) = 1, (\forall)s \neq r$ și $\gamma(s/s) = 0, (\forall)s$, atunci nodul va fi asignat clasei cu cei mai mulți reprezentanți în ea.

Costul misclasării unei observații aparținând nodului a este: $c(a) = \min_{1 \leq s \leq q} \sum_{r=1}^q \gamma(s/r) \rho(r|a)$. Costul misclasării datorat nodului a , este $C(a) = c(a) \rho(a)$, unde $\rho(a)$ este probabilitatea nodului.

Riscul erorii de afectare datorat arborelui A : $rea(A) = \sum_{a \in \hat{A}} C(a) = \sum_s \sum_{a \in \hat{A}(s)} \sum_r \gamma(s/r) \rho(r|a) \rho_r = \sum_s \sum_r \gamma(s/r) (n_{sr}/n)$, unde \hat{A} este mulțimea nodurilor terminale ale lui A , $\hat{A}(s)$ este mulțimea nodurilor terminale ale lui A asignate clasei s , ρ_r este probabilitatea a priori ca un nod să provină din clasa r , n_{sr} este numărul de indivizi din clasa r clasăți în clasa $s, s \neq r$.

Un subarbore al lui A_{max} este optimal („cel mai bun”) dacă numărul de segmente terminale conținute și riscul erorii de afectare sunt minime și, în plus, furnizează o estimatie corectă a erorii teoretice de clasare. Pentru selecția subarborelui optimal se împarte eșantionul inițial într-un eșantion de instruire și un eșantion de testare. Pornind de la eșantionul de instruire se construiește arborele A_{max} . Operația de „tundere” a arborelui A_{max} constă în construirea unui șir optimal $A_H, \dots, A_h, \dots, A_1$ de subarbori incluși, unde A_H este A_{max} , A_h este subarborele cu h segmente terminale, A_1 este eșantionul total. Fiecare subarbore A_h din acest șir este optimal, în sensul că eroarea aparentă a subarborelui este minimală printre toți subarborii având același număr de segmente terminale, $ea(A_h) = \min_{A \in S_h} ea(A)$, unde S_h este mulțimea subarborelor lui A_{max} cu h segmente terminale. Se selectează din șirul de arbori optimali subarborele \tilde{A} cu eroarea teoretică minimă dată de relația: $et(\tilde{A}) = \min_{1 \leq h \leq H} et(A_h)$. Eroarea teoretică se estimează după formula $\tilde{et}(A) = \sum_{t \in A} \tilde{R}_t$, unde avem: $\tilde{R}_t = (\tilde{n}_t / \tilde{n}) \times \tilde{s}_t^2$, \tilde{n} este volumul eșantionului de test, \tilde{n}_t este numărul de indivizi din eșantionul de test aparținând segmentului t , \tilde{s}_t^2 este dispersia de selecție a variabilei y în interiorul segmentului t , $\tilde{s}_t^2 = (1/\tilde{n}_t) \sum_{i=1}^{|\tilde{t}|} (y_i - \tilde{y})^2$, \tilde{y} este media de selecție în interiorul segmentului t , $|\tilde{t}| = \text{card}(\tilde{t})$.

Deși cea mai bună diviziune, $\tilde{\partial}$, a unui nod este cea care asigură cea mai mare reducere a dispersiei reziduale (sau a impurității), prin trecerea de la acel nod la segmentele descendente, prospectorul de date poate utiliza și alte diviziuni (echi-reductive, echi-divizante), aproximativ la fel de bune, dar foarte importante la nivelul interpretării.

9. Metode de agregare a modelelor

Agregarea (sau combinarea) unui număr mare de modele permite ameliorarea ajustării modelelor evitându-se, totodată, supraajustarea acestora și se bazează pe două tipuri de strategii de

agregare: aleatoare (*bagging*) și adaptive (*boosting*).

Strategii aleatoare. Principiul bagging-ului se bazează pe faptul că medierea previziunilor mai multor modele independente permite reducerea varianței și deci reducerea erorii de previziune.

Fie Y variabila de explicat, cantitativă sau calitativă cu modalitățile $\tau = 1 \div q$, fie $X = \{X^j\}_{j=1}^p$ variabilele explicative, fie $\varphi(X)$ un model funcție de X și fie $z = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ un eșantion de lege F . Speranța, $M_F(\tilde{\varphi}_z)$, a unui estimator $\tilde{\varphi}_z$ definit pe eșantionul z , este un estimator fără bias, de varianță nulă.

Se consideră K eșantioane independente, notate $\{z_\kappa\}_{\kappa=1}^K$, și se construiește familia de modele $\{\varphi_{z_\kappa}\}_{\kappa=1}^K$.

Estimarea medie va fi:

$$\tilde{\varphi}_K(\bullet) = \begin{cases} M_F(\tilde{\varphi}_{z_\kappa}) = (1/K) \sum_{\kappa=1}^K \tilde{\varphi}_{z_\kappa}(\bullet), & \text{dacă variabila de explicat } Y \text{ este cantitativă} \\ \operatorname{argmax}_{1 \leq \tau \leq q} \{ \kappa \mid \tilde{\varphi}_{z_\kappa}(\bullet) = \tau, \kappa = 1 \div K \}, & \text{dacă } Y \text{ este calitativă} \end{cases}$$

În primul caz, estimarea medie este media rezultatelor obținute pentru modelele asociate fiecărui eșantion. În al doilea caz, a fost constituit un „comitet de modele” pentru a vota și a alege răspunsul cel mai probabil. Când modelul returnează probabilități, asociate cu fiecare modalitate τ sau cu fiecare arbore de decizie, se calculează mediile acestor probabilități.

Practic, cele K eșantioane independente, z_κ , ar necesita, în general, prea multe date și ele sunt înlocuite prin K eșantioane bootstrap, z_κ^* , obținute, fiecare, prin n extrageri cu înlocuire conform legii empirice \tilde{F} . În fiecare iterație κ ($\kappa = 1 \div K$), se extrage eșantionul bootstrap, z_κ^* , și se calculează $\tilde{\varphi}_{z_\kappa^*}(\mathbf{x})$ pe acest eșantion. În final, după cum variabila de explicat Y este cantitativă sau calitativă, estimarea medie este sau media estimărilor sau rezultatul votului.

Păduri aleatoare. Pentru metoda segmentării o îmbunătățire a bagging-ului se poate obține prin adăugarea unei randomizări. Obiectivul este de mări independența arborilor de agregare prin intervenția hazardului în alegerea variabilelor implicate în modele. În fiecare iterație κ ($\kappa = 1 \div K$): se extrage un eșantion bootstrap z_κ^* și se estimează un arbore pe z_κ^* prin randomizarea variabilelor (căutarea fiecărui nod optimal este precedată de selecția aleatoare a unei submulțimi de $q \leq p$ predictorii). În final, $\tilde{\varphi}_K(\mathbf{x}) = (1/K) \sum_{\kappa=1}^K \tilde{\varphi}_{z_\kappa^*}(\mathbf{x})$ sau $\tilde{\varphi}_K(\mathbf{x}) =$ rezultatul votului. Față de bagging, în cazul „pădurilor aleatoare” (*Random Forest*), strategia de tăiere poate fi mai simplă limitându-se la arbori de mărimi, q , relativ reduse (chiar triviale: $q = 2$). Într-adevăr, doar cu bagging arborii limitați la o singură ramificație riscă să fie foarte asemănători (puternic corelați) implicând, aceleași, câteva variabile care apar ca fiind cele mai explicative. În fiecare etapă de construcție a unui arbore, selectarea aleatoare a unui număr redus de predictorii potențiali crește semnificativ variabilitatea având în mod necesar alte variabile. Fiecare model de bază este în mod evident mai puțin eficient dar agregarea duce în cele din urmă la rezultate bune. Numărul de variabile extrase aleator nu este un parametru sensibil fapt pentru care Breiman (2001) sugerează alegerea implicită $q = p$. Evaluarea iterativă a erorii out-of-bag previne o eventuală supraajustare dacă aceasta tinde să se degradeze. Ca la toate modelele construite prin agregare (sau „cutie neagră”), pentru prospectorul de date nu există nici o interpretare directă. Informațiile relevante sunt obținute prin calcul și prin reprezentarea grafică a unor indici, proporționali cu importanța fiecărei variabile din modelul agregat adică cu participarea acesteia la regresie sau discriminare. Aceste informații sunt cu atât mai utile cu cât variabilele sunt mai numeroase. Pentru a evalua importanța unei variabile prospectorul de date utilizează criterii precum: frecvența cu care apare fiecare variabilă în arborii pădurii, *MDA* (*Mean Decrease Accuracy*) sau *MDG* (*Mean Decrease Gini*).

Strategii adaptive. Boosting-ul adoptă același principiu general ca și bagging-ul: construirea unei familii de modele care să fie agregate prin o medie ponderată a estimărilor sau a unui vot. El diferă net de bagging în ce privește modul de construire a familiei care, de această dată, este recurent: fiecare model este o versiune adaptivă a precedentului acordând, în momentul estimării următoare, o pondere mai mare observațiilor prost ajustate sau prost previzionate. Intuitiv, acest algoritm își concentrează eforturile asupra observațiilor celor mai dificil de ajustat

astfel încât combinarea ansamblului de modele permite evitarea supraajustării.

Pentru exemplificare se consideră problema de discriminare în două clase și fie δ funcția de discriminare cu valori în $\{-1, 1\}$. Pentru estimarea primului model ponderile w_i ale fiecărei observații sunt inițializate la $1/n$, în continuare aceste ponderi evoluează la fiecare iterație adică pentru fiecare nouă estimare. Importanța, w_i , a unei observații rămâne neschimbată dacă observația este bine clasată, dacă observația nu este bine clasată w_i crește proporțional cu deficitul de ajustare al modelului. Agregarea finală a previziunilor, $\sum_{\kappa=1}^K c_{\kappa} \delta_{\kappa}(\mathbf{x})$, este o combinație ponderată a calităților de ajustare ale fiecărui model. Valoarea absolută a sa, numită marje, este proporțională cu încrederea care poate fi acordată semnelui său care furnizează rezultatul previziunii.

Fie $\mathbf{z} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ un eșantion și \mathbf{x} individul de previzionat. Se inițializează \mathbf{w}_1 , vectorul de ponderi: $w_{1,i} = 1/n$. În fiecare iterație κ ($\kappa = 1 \div K$): mai întâi se estimează δ_{κ} pe eșantionul \mathbf{z}_{κ} (\mathbf{z} ponderat cu \mathbf{w}_{κ}); apoi se consideră vectorul $\mathbf{Q}_{\kappa} = \{Q_{\kappa,i}\}_{i=1}^n$, unde $Q_{\kappa,i}$ este un indice de misclasare ($Q_{\kappa,i} = 1$ dacă $\delta_{\kappa}(x_i) \neq y_i$ și $Q_{\kappa,i} = 0$ dacă $\delta_{\kappa}(x_i) = y_i$) cu ajutorul căruia se estimează eroarea de previziune $\tilde{\epsilon}_P = (\sum_{i=1}^n w_i Q_{\kappa,i}) / (\sum_{i=1}^n w_i)$; se calculează $c_{\kappa} = \log((1 - \tilde{\epsilon}_P) / \tilde{\epsilon}_P)$; se calculează noile ponderi: $w_{\kappa+1,i} = w_{\kappa,i} \exp[-c_{\kappa} Q_{\kappa,i}]$, $i = 1 \div n$. În final, rezultatul este: $\tilde{\phi}_K(\mathbf{x}) = \text{sign}[\sum_{\kappa=1}^K c_{\kappa} \delta_{\kappa}(\mathbf{x})]$.

Principiile bagging-ului sau boosting-ului se pot aplica la orice metodă de modelare dar nu sunt interesante și nu reduc sensibil eroarea de previziune decât în cazul modelelor instabile deci, mai degrabă, neliniare. Astfel, pentru prospectorul de date, utilizarea acestor algoritmi nu are nici un sens cu regresia multiliniară sau cu analiza discriminantă. Ei pot fi foarte utili în asociere cu arborii binari ca modele de bază.

10. Concluzii

Practica de a obține din date cunoștințe valoroase și utile pentru susținerea activităților decizionale, denumită tot mai frecvent *data science*, este în continuă și rapidă dezvoltare pentru a face față provocărilor de prelucrare a seturilor uriașe de date (structurate, nestructurate sau semi-structurate generate de dispozitive inteligente, telefoane mobile, web, mass-media sau rețele sociale), *big data*.

Informatica decizională utilizează statistica descriptivă, pentru date cu mare densitate în informație, pentru a măsura fenomene, a detecta tendințe, etc. în timp ce *big data* utilizează statistica inferențială, pentru date cu slabă densitate în informație, ale căror volume, foarte mari, permit inferențe ale legilor conferindu-le capacități predictive (cu limitele acestor inferențe).

Pentru prospectarea datelor și interpretarea rezultatelor *data scientist*, specializat de obicei pe un anumit domeniu (marketing, medicină, securitate, fraudă, finanțe, etc.), se bazează pe expertize din statistică, instruire, optimizare, procesare de semnale, regăsire de informații sau procesare a limbajului natural.

Având o pregătire de bază în matematică și statistică, noul *data scientist* poate privi cu seninătate sosirea valului sau tsunami-ului *Big Data*.

Activitatea informatică din amonte (permanent reînnoită de evoluția rapidă a tehnologiilor) este importantă, pentru a stoca datele și a face executabile metodele dar, conceptual, matematica necesară modelelor respective a luat deja în considerare mărimi și dimensiuni infinite în spații hilbertiene.

Înzestrat cu acest „instrumentar” durabil, *data scientist* poate deci aborda și susține, cu șanse de succes, cercetările emergente.

BIBLIOGRAFIE

1. **BANCIU, D.; COARDOȘ, D.; LEPĂDATU, C-I.; LEPĂDATU, C.:** Enhancement of the Retrospective National Bibliography of the Romanian Book through the Application of the Informational Technologies, Proceedings of BIBLIO 2011 „Innovation en

- bibliotheque/Innovation within libraries”, Editura Universității Transilvania din Brașov, 2011, pp. 131-142.
2. **BESSE, P.; LAURENT, B.:** Apprentissage Statistique: modélisation, prévision et data mining, Institut National des Sciences Appliquées de Toulouse, 2014, 159 p.
 3. **CIUREA, C.; DUMITRESCU, G.; LEPĂDATU, C.:** The impact analysis of implementing virtual exhibitions for mobile devices on the access to national cultural heritage, Proceedings of 2nd International Conference Economic Scientific Research – Theoretical, Empirical and Practical Approaches, ESPERA 2014, Bucharest, Romania.
 4. **COARDOȘ, D.; COARDOȘ, V.; LEPĂDATU, C-I.; LEPĂDATU, C.:** Support Systems for Libraries Based on Business Intelligence Tools, 2008 IEEE International Conference on Intelligent Computer Communication and Processing - Digital Libraries Workshop, Cluj Napoca, August 2008.
 5. **COARDOȘ, D.; COARDOȘ, V.; LEPĂDATU, C-I.; LEPĂDATU, C.:** Integrated On-line System for Management of the National Retrospective Bibliography – SIMBNR, 2009 IEEE International Conference on Intelligent Computer Communication and Processing - Workshop on Digital Libraries, e-Content Management and e-Learning”, Cluj Napoca, August 2009.
 6. **DUMITRESCU, G.; FILIP, F.-G.; IONIȚĂ, A.; LEPĂDATU, C.:** Open Source Eminescu’s Manuscripts: A Digitization Experiment, Studies in Informatics and Control, 19(1), 2010, pp. 79-84.
 7. **DUMITRESCU, G.; LEPĂDATU, C.; CIUREA C.:** Creating Virtual Exhibitions for Educational and Cultural Development, INFOREC Publishing House, Informatica Economică Journal, 2014, 18(1), pp. 102-110.
 8. **ENĂCHESCU, D.:** Data Mining: metode și aplicații, Edit. Academiei Române, 2009, 277 p.
 9. **FAYYAD, U.; PIATETSKY-SHAPIRO, G.; SMYTH, P.:** From Data Mining to Knowledge Discovery in Databases, AAAI, AI Magazine, 17 (3), 1996, pp. 37-54.
 10. **FILIP, F.-G.:** Decizie asistată de calculator: decizii, decidenți - metode de bază și instrumente informatice asociate, Ed. a 2-a, București, Editura Tehnică, 2005, 376 p.
 11. **FILIP, F.-G. HERERA-VIEDMA, E.:** Big Data in the European Union, National Academy of Engineering (NAE), SUA, Winter Bridge: A Global View of Big Data, 2014, 44(4), pp. 33-37.
 12. **HAN, J.; KAMBER M.; PEI, J.:** Data Mining: Concepts and Techniques, Third Ed., Elsevier, 2011, 703 p.
 13. **HASTIE, T.; TIBSHIRANI, R., FRIEDMAN, J.:** The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2nd Edition, Springer-Verlag New York, 2009, 745 p.
 14. **IONIȚĂ, A.; LEPĂDATU, C.; DUMITRESCU, G.:** Digital Cultural Landscape Content, Hernik, Jozef (ed.) Cultural Landscape – Across Disciplines, Oficyna Wydawnicza BRANTA, Kracow, Poland, 2009, pp. 255-277.
 15. **LEPĂDATU, C.:** De la descriere bibliografică la web semantic, Academica, 2006, XVI (185-186/48-49), pp 78-81 și XVI (188/51), pp. 42-85.
 16. **LEPĂDATU, C.:** Support Systems for Knowledge Culture based on Solution and Tools from the Field of Business Intelligence – SSCBI, Proceedings of the Workshop IST – Multidisciplinary Approaches, Bucharest, Romania, 2006, pp. 7-12.
 17. **LEPĂDATU, C.:** Acquisition Policy of a Library and Data Mining Techniques, Studies in informatics and control, 16(4), 2007, pp. 413-420.
 18. **LEPĂDATU, C.:** Explorarea datelor și descoperirea cunoștințelor - probleme, obiective și strategii, Revista Română de Informatică și Automatică, 2012, 22(4), pp. 5-14.

19. **LEPĂDATU, C.:** Metode exploratorii multidimensionale, Revista Română de Informatică și Automatică, 23(1), 2013, pp. 14-30.
20. **LEPĂDATU, C.:** Sisteme suport pentru decizii și bibliomining, Revista Română de Informatică și Automatică, 24(2), 2014, pp. 17-30.
21. **LEPĂDATU, C.:** Sisteme suport pentru decizii de bibliotecă, Revista Română de Informatică și Automatică, 24(3), 2014, pp. 5-17.
22. **MAIMON, O. ROKACH, L. (EDS.):** Data Mining and Knowledge Discovery Handbook, 2nd Ed., Springer New York Dordrecht Heidelberg London, 2010, 1306 p.
23. **NICULESCU, C.; LEPĂDATU, C.; ȘTEFĂNESCU, D.:** SSCBI - A Teleworking Environment of Support Systems for Knowledge Culture. In the CD REV 2007 Proceedings of the International Conference Remote Engineering Virtual Instrumentation, Porto, Portugal, iunie 2007.
24. **TUFFÉRY, S.:** Modélisation Predictive et Apprentissage Statistique avec R, TECHNIP, 2015, 415 p.
25. **VAPNIK, V. N.:** Statistical learning theory, Wiley-Interscience, 1998, 768 p.