

O îmbunătățire a performanțelor algoritmului KNN în sistemele de recomandare pe web

Costin-Gabriel Chiru
Universitatea „Politehnica”
Bucuresti
Splaiul Independenței nr.
313 - București, România
chirucos@gmail.com

Ștefan Trăușan-Matu
Universitatea „Politehnica”
Bucuresti
Splaiul Independenței nr.
313 - București, România
trausan@gmail.com

Traian Rebedea
Universitatea „Politehnica”
Bucuresti
Splaiul Independenței nr.
313 - București, România
trebedea@gmail.com

REZUMAT

Lucrarea de față prezintă o abordare ce poate conduce la îmbunătățirea rezultatelor întoarse de sistemele de recomandare bazate pe o implementare de tipul filtrării colaborative. Recomandările oferite de acest tip de sisteme sunt generate pe baza preferințelor utilizatorilor care au comportamente asemănătoare cu cel al utilizatorului curent. Pentru identificarea acestor utilizatori se utilizează foarte frecvent o metodă de clasificare automată bazată pe algoritmul „cei mai apropiați K vecini”. Scopul acestei lucrări este să trateze o modalitate de a prezenta și corecta unele probleme ale acestui algoritm, prin aplicarea unei interpolări asupra claselor întoarse de către varianta standard a algoritmului de clasificare, astfel încât să se îmbunătățească acuratețea metodei și, implicit, a rezultatelor furnizate de sistemele de recomandare.

Cuvinte cheie

Sisteme de recomandare, KNN, clasificare.

Clasificare ACM

Interfețe inteligente om-calculator, agenți inteligenți

INTRODUCERE

Sistemele de recomandare încearcă să furnizeze utilizatorilor propuneri de itemi (pagini web, știri, cărți, filme, muzică, poze, CD-uri etc.) care pot fi interesați pentru aceștia, plecând de la filtrarea informațiilor pe baza preferințele altor persoane, cu un profil similar. De aceea, aceste sisteme reprezintă o alternativă foarte bună pentru algoritmii de căutare, deoarece ajută utilizatorii să descopere itemi pe care nu i-ar putea descoperi singuri datorită multitudinii de opțiuni din care pot alege [3].

Dintr-o altă perspectivă, putem spune că se schimbă direcția interacțiunii om-calculator. În loc ca utilizatorul să ceară explicit un item, sistemele de recomandare îi propun acestuia un item care a fost preferat de persoane cu caracteristici asemănătoare. De exemplu, unul dintre cele mai cunoscute astfel de sisteme este cel de la amazon.com, care, plecând de la caracteristicile pe care le deduce din cărțile cumpărate sau chiar doar examinate de un utilizator pe situl web, recomandă automat cărți cumpărate de persoane care au avut interese similare.

Sistemele de recomandare compară profilul utilizatorului cu anumite caracteristici de referință, pe baza acestei comparații calculând o listă de itemi ce pot fi interesați pentru utilizator. Caracteristicile de referință pot fi de două feluri: caracteristici ce țin de contextul comportamentului

social al utilizatorului (abordare de tipul filtrării colaborative) sau caracteristici ce țin de natura itemilor (abordare bazată pe conținut) [2].

Majoritatea abordărilor de tipul filtrării colaborative folosesc metode de tipul „cei mai apropiați K vecini” (engl. „K Nearest Neighbors” - KNN). Aceasta este o metodă de învățare automată din exemple pentru aproximarea funcțiilor cu valori reale sau discrete. Acest tip de învățare stochează instanțele furnizate ca exemple pentru antrenarea sistemului și amână învățarea efectivă până când este primit un exemplu de test. Învățarea are loc de fiecare dată când se primește un nou test. Procesul constă din extragerea din memorie a instanțelor similare și clasificarea exemplului de test pe baza acestora [6].

În următoarele două secțiuni vom face o scurtă descriere a sistemelor de recomandare și a algoritmului clasic de KNN. Următoarele secțiuni descriu îmbunătățirile aduse algoritmului clasic și datele cu care s-a lucrat. Lucrarea se încheie cu prezentarea rezultatelor obținute și cu concluziile derivate din aceste rezultate.

SISTEME DE RECOMANDARE

Sistemele de recomandare încearcă să ajute utilizatorii să identifice rapid care sunt produsele și serviciile interesante dintr-o gamă de oferte extrem de variată. Numărul și complexitatea acestora fac practic imposibilă încercarea utilizatorilor de a le trece în revistă cu scopul de a lua o decizie în privința lor. Sistemele de recomandare sunt folosite de către foarte multe magazine on-line, putând să ofere recomandări despre aproape orice: începând de la lucruri simple, ca pagini web, știri, cărți, filme, muzică, poze, CD-uri, și până la itemi complecși de genul serviciilor financiare sau serviciilor de e-government.

Recomandările pot fi generate prin două metode: prin purtarea de dialoguri cu utilizatorii on-line de pe sit sau prin analiza informațiilor existente cu privire la itemii cumpărați de un individ sau de către o comunitate de utilizatori. Datorită faptului că metoda dialogului s-a demonstrat a fi inadecvată, s-a trecut la construcția de sisteme utilizând metoda analizei. Primul sistem de recomandare bazat pe această metodă a fost construit la mijlocul anilor 90 [2]. De atunci, interesul pentru astfel de sisteme a crescut extrem de mult datorită cererii mari de tehnologii personalizate de către aplicațiile de succes din domeniul comerțului electronic (ex. amazon.com).

Una din cele mai des utilizate metode folosite în sistemele de recomandare este filtrarea colaborativă [4, 5], care

adoptă, în general, o metodă de KNN aplicată unei matrice de evaluări. În filtrarea colaborativă, opiniile prietenilor și rapoartele ce compară diferiți itemi reprezintă principalele elemente care influențează decizia utilizatorilor de a cumpăra unul din acei itemi [2]. Alte metode utilizate frecvent în sistemele de recomandare sunt metodele bazate pe conținut (pe ceea ce a cumpărat respectivul utilizator în trecut) și respectiv cele bazate pe cunoștințe (pentru cumpărături complexe care necesită atât informații suplimentare despre context cât și modalități de interacțiune cu sistemul) [2].

ALGORITMUL KNN

KNN este cea mai simplă metodă de clasificare bazată pe instanțe. Algoritmul presupune că toate exemplele sunt plasate într-un spațiu n -dimensional. Fiecare instanță de test este clasificată în funcție de cele mai apropiate k exemple din memorie, fără a încerca să se determine funcția de distribuție pe baza căreia instanțele sunt plasate în spațiul n -dimensional [6].

Pentru a găsi și extrage din memorie aceste exemple, este nevoie de o metodă de a compara cât de similare sunt două instanțe. În cazul în care se lucrează cu valori reale, se pot aplica metrici de genul distanței Euclidiene, distanței normate sau funcției cosinus între instanțe. Pentru situațiile în care valorile sunt discrete, se calculează numărul caracteristicilor distincte între cele două instanțe. În lucrarea curentă, ne vom concentra doar asupra ultimei variante, cazul valorilor reale fiind similar, unica diferență fiind modul de calcul al similarității dintre două instanțe. După definirea unei metrici pentru determinarea similarității dintre două instanțe, următorul pas este să se calculeze distanța dintre exemplul de test și toate exemplele de învățare aflate în memorie. În ultimul pas se determină pentru fiecare instanță de test care sunt cele mai apropiate k instanțe din setul din memorie și pe baza acestor instanțe similare se decide care este cea mai probabilă clasificare a acesteia. Acest lucru se realizează prin votarea între cele mai similare k instanțe.

Din modul de calcul al clasificării instanțelor de test se observă faptul că algoritmul KNN utilizează o aproximare diferită pentru fiecare instanță de test în parte. Acest lucru oferă un mare avantaj în cazul în care funcția de distribuție este foarte complexă, dar poate fi descrisă printr-o colecție de aproximări locale mai puțin complicate.

ÎMBUNĂȚĂIREA ALGORITMULUI CLASIC DE KNN

Algoritmul KNN este o metodă inductivă de inferență foarte eficientă pentru multe probleme practice. Acest algoritm este rezistent la zgomotul ce poate să apară în datele de antrenament și este foarte eficient în cazul în care i se furnizează un set de date de antrenament suficient de mare. În ciuda rezultatelor bune obținute folosind KNN, acest algoritm are totuși câteva probleme, pe care ne propunem să le corectăm prin îmbunătățirile propuse. Una din cele mai mari probleme ale algoritmului KNN este faptul că nu oferă informații despre ce valoare ar trebui să fie aleasă pentru „ k ” astfel încât să se obțină cele mai bune rezultate. O altă problemă a algoritmului KNN este aceea a ponderii pe care o au instanțele în calculul rezultatului final. Se pune problema dacă toate instanțele din setul de antrenare trebuie să aibă aceeași pondere. Ultima problemă considerată în lucrarea de față a fost

aceea a distanței medii dintre instanța de test și instanțele din setul de antrenament. Se pune întrebarea dacă se pot obține rezultate mai bune în cazul în care se folosește distanța ca modalitate de ponderare a rezultatelor obținute.

Pentru a elimina problemele identificate la algoritmul KNN, s-a construit un sistem care să poată aborda fiecare dintre probleme în parte.

Acest sistem ia toate valorile lui „ k ” începând de la valoarea 1 și până la o limită maximă și calculează care este clasificarea instanței de test pentru fiecare valoare a lui „ k ” în intervalul specificat și pentru datele de antrenare. Această limită poate fi impusă de utilizator sau poate fi aleasă astfel încât să fie egală cu numărul total de instanțe de antrenare furnizate sistemului. În continuare, există două variante de utilizare a rezultatelor întoarse de aplicările algoritmului cu diferite valori ale lui „ k ”:

- prima variantă, în care nu se folosește nici o pondere pentru valorile întoarse (toate aplicările au pondere egală)
- a doua variantă, în care valorile întoarse de algoritmul KNN sunt ponderate cu valorile folosite pentru „ k ”.

Aceste posibilități au fost luate în calcul pentru a putea adresa problema ponderării, adică pentru a vedea dacă rezultatele obținute trebuiesc (sau nu) ponderate în vreun fel în funcție de valoarea lui „ k ”.

În afara acestor două variante, se mai dă posibilitatea utilizatorului de a alege dacă dorește sau nu să folosească distanța medie dintre instanța de test și cele de antrenare considerate în aplicarea respectivă a algoritmului KNN, pentru a pondera în plus rezultatele obținute.

Alegerea lui „ k ”

Așa cum s-a specificat în secțiunile anterioare, alegerea valorii lui „ k ” pentru care se va aplica algoritmul KNN ridică numeroase probleme.

Valorile mici ale lui „ k ” permit identificarea mai eficientă a structurilor de dimensiuni mici din spațiul problemei, dar pot conduce la crearea unui model prea complex (engl. „overfitting”). În plus, dacă se consideră valori mici pentru „ k ”, atunci rezultatele sunt influențate în mod sensibil de datele incorecte (engl. „noisy data”).

Dacă se aleg valori mari pentru „ k ”, se evită problema zgomotului din datele de antrenare și, de asemenea, se pot estima mai bine probabilitățile în cazul în care se lucrează cu valori discrete.

În schimb, pe lângă avantajele oferite de o astfel de alegere apar și 2 dezavantaje majore. O primă problemă este că nu se pot aplica valori mari pentru „ k ” în cazul în care se lucrează cu seturi mici de date de antrenare. Cel de-al doilea dezavantaj este că nu se cunoaște cât de mult se poate mări valoarea lui „ k ” fără a scădea performanțele – mărirea valorii lui „ k ” înseamnă de fapt simplificarea modelului din spațiul problemei și dacă acest model se simplifică prea mult, este posibil să nu mai obținem o soluție corectă.

Cercetătorii din domeniu au arătat că nu există o valoare general valabilă pentru „ k ” astfel încât obținerea rezultatelor optime să fie garantată și că această valoare

depinde de fiecare problemă în parte. De asemenea, aceștia au ajuns la concluzia că valoarea lui „k” trebuie aleasă din intervalul 1 – 30 .

Abordarea noastră încearcă să elimine aceste probleme prin evitarea alegerii între diferite valori ale lui „k”, algoritmul KNN rulându-se cu toate valorile posibile mai mici decât o valoare limită. Cu alte cuvinte, în loc să se aplice algoritmul KNN pentru o singură valoare a lui „k”, acesta se aplică de mai multe ori având valori diferite pentru „k”. În acest fel, se încearcă eliminarea problemelor apărute din alegerea unei valori prea mici sau prea mari pentru „k”.

Ponderarea

Problema ponderării constă în a determina dacă datele de antrenament trebuie să aibă ponderi diferite în rezultatele furnizate. Aceasta presupune să se dea o importanță mai mare (tradusă printr-o pondere mai mare) instanțelor care sunt mai aproape de instanța de test și o importanță mai mică (tradusă printr-o pondere mai mică) instanțelor care sunt departe de aceasta.

Datorită îmbunătățirilor prezentate anterior, această problema se rezolva de la sine. Abordarea noastră rezolvă în mod implicit această problemă prin considerarea multiplă a instanțelor ce sunt mai aproape de instanța de test. Cu cât o instanță este mai aproape de instanța de test, cu atât valoarea acesteia va fi luată în considerare în mai multe aplicări ale algoritmului KNN având valori diferite ale lui „k”. În acest fel, instanțele aflate în apropierea instanței de test capătă ponderi mai mari ca celelalte în final, când se decide care va fi valoarea instanței de test.

Chiar dacă această problemă este rezolvată implicit, o altă problemă apare: valorile întoarse de aplicările diferite ale algoritmului KNN trebuiesc ponderate în raport cu numărul de vecini considerați sau nu? Ceea ce încercăm să exprimăm prin această întrebare este posibilitatea ca valorile întoarse de aplicările algoritmului KNN pentru diferite valori ale lui „k” să aibă importanță diferită în votarea finală. Din această cauză, sistemul permite cele două variante: valorile întoarse nu sunt ponderate în nici un fel sau ele sunt ponderate cu valorile folosite pentru „k”. Dând aceasta posibilitate, am încercat să vedem dacă clasificările obținute folosind valori mici ale lui „k” ar trebui să aibă mai mare importanță (deci și ponderi mai mari) în raport cu cele obținute utilizând valori mai mari pentru „k”.

Distanța

O ultima problemă este aceea a identificării în ce măsură considerarea distanței în ponderarea rezultatelor este utilă. Ideea care a generat posibilitatea considerării acesteia ca factor de ponderare a fost că în cazul în care se obține o valoare mai mică a distanței medii între instanța de test și cele de antrenament avem de-a face cu o structură mai compactă decât în cazul în care această valoare este mare, și de aceea structura ar putea fi mai aproape de realitate. Această idee este mai ușor de înțeles urmărind exemplul următor în care valoarea lui „k” este 3 și în care avem două posibilități pentru instanțele de test:

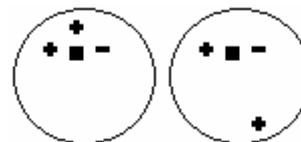


Figura 1 Exemplu 3NN

În ambele figuri instanța de test (reprezentată printr-un pătrat) va fi clasificată ca fiind pozitivă, dar încrederea cu care se face această clasificare nu este aceeași pentru cele două figuri. Dacă se calculează distanță medie în ambele cazuri, și se ponderează clasificarea returnată de algoritmul 3NN cu această distanță se observă că avem mai multă încredere în primul caz decât în al doilea.

Datorită acestei situații, s-a realizat un test care să testeze influența distanței în clasificarea instanței de test. Astfel, pentru fiecare „k” în parte, s-a calculat valoarea medie a distanței între instanța de test și instanțele de antrenament considerate.

DATELE

Pentru efectuarea testelor s-au folosit cinci seturi diferite de date, dintre care trei (Connect-4, Clasificarea ciupercilor și X și 0) au fost preluate din baza de date a UCI Machine Learning Repository [1]. Al patrulea set de date a fost preluat de pe situl Universității din Antwerp de pe pagina cursului de tehnologia limbajului („Language Technology” [7]). Ultimul set de date a fost preluat de pe situl Universității din Texas la Dallas de pe pagina cursului de învățare automată (engl. „Machine Learning”) [8]. În continuare vom prezenta fiecare din aceste seturi de date:

Jocul Connect-4

Acest set de date conține toate pozițiile legale din jocul connect-4 pe o tabla de 6 x 7 pătrățele, în care nici unul dintre jucători nu a câștigat încă și în care următoarea mutare nu este forțată [9]. Setul de date conține 67.557 instanțe, fiecare având câte 42 atribute, corespunzătoare celor 42 de pătrățele disponibile în joc. Aceste atribute pot să ia una din valorile X, 0 sau B ele semnificând că jucătorul X respectiv 0 au ocupat poziția respectivă sau că nici unul din cei doi nu a ocupat pătrățelele și acesta este gol (engl. blank). Exemplele nu conțin valori lipsă și sunt clasificate din perspectiva jucătorului X, și anume: victorie (engl. win), înfrângere (engl. loss) sau egalitate (engl. draw). Din acest set, pentru experimente s-au folosit doar 300 de instanțe, împărțite în 200 instanțe pentru antrenare și 100 pentru test.

Clasificarea ciupercilor

Acest set de date conține o descriere a unor exemple de ciuperci făcând parte din 23 de specii din familiile Agaricus și Lepiota [10]. Fiecare exemplu este clasificat ca fiind o ciupercă otrăvitoare, bună de mâncat sau suspectă. Setul de date are 8.124 instanțe, fiecare având câte 22 de atribute. Aceste atribute au ca valori clasificările obținute pe baza criteriului ce definește atributul. Setul conține 2.480 de valori lipsă în coloanele atributelor și 124 în coloana clasei specificând faptul că ciupercile respective nu au putut fi clasificate exact. Din acest set, s-au reținut pentru teste numai instanțele ce au valori în coloana rezultatelor (8.000 exemple), acestea

fiind împărțite în 7.000 instanțe de antrenare și 1.000 instanțe de test.

Jocul X și 0

Acest set de date codifică toate configurațiile posibile ale unei table de joc de 3 x 3 pătrățele pentru jocul „X și 0” în care se presupune că jucătorul X este primul care joacă [11]. Setul conține 958 de instanțe, fiecare având câte 9 atribute corespunzătoare celor 9 pătrățele de pe tabla de joc. La fel ca la jocul connect-4 prezentat anterior, și aici atributele pot avea valoare X, 0 sau B având aceeași semnificație ca mai sus. Exemplele nu conțin valori lipsă iar clasificarea folosită este binară, reprezentând capacitatea jucătorului x de a obține victoria: pozitiv (victorie) sau negativ (înfrângere sau egalitate). Din acest set se folosesc doar 250 de instanțe, grupate în setul de antrenament (200 instanțe) și setul de test (50 instanțe).

Formarea diminutivelor în limba olandeză

Acest set de date descrie modul de formare a diminutivelor în limba olandeză [12]. Setul de date conține 3.949 de instanțe, fiecare având câte 12 atribute. Atributele au ca valori clasificările diminutivelor realizate pe diferite criterii, iar clasificarea instanțelor se poate face în cinci clase distincte. Setul de date este complet specificat (nu are valori lipsă) și a fost împărțit într-o parte de antrenare de 2.999 de instanțe și o parte de test de 950 de instanțe.

Setul 5

Acest set de date a fost creat în scopuri didactice și descrie o clasificare binară bazată pe șase atribute care sunt și ele tot binare [8]. Setul de date conține 800 de instanțe de antrenare și 203 instanțe de test.

PREZENTAREA REZULTATELOR

Pentru a evalua îmbunătățirile propuse, s-au construit două sisteme: unul care calculează clasa unei instanțe de test folosind algoritmul clasic de KNN și altul care folosește îmbunătățirile algoritmului KNN propuse în secțiunile anterioare. Cele două sisteme au fost aplicate în paralel asupra fiecăruia din cele 5 seturi de date prezentate anterior.

Pentru varianta îmbunătățită de KNN s-au considerat mai multe funcții de ponderare posibile – $\sin(x)$, $\exp(-x)$, $\text{sum}(x)$ – pentru a vedea care este cea mai bună alegere. De asemenea au fost considerate combinații ale lui „k” cu distanța medie a exemplilor de antrenament față de instanța de test.

S-a observat că dintre toate aceste funcții, cele mai bune rezultate au fost obținute în cazul în care se folosește votarea între clasele returnate de aplicările algoritmului KNN pentru diferite valori ale lui „k”, fără a considera nici un fel de funcție de ponderare între termenii implicați în votare.

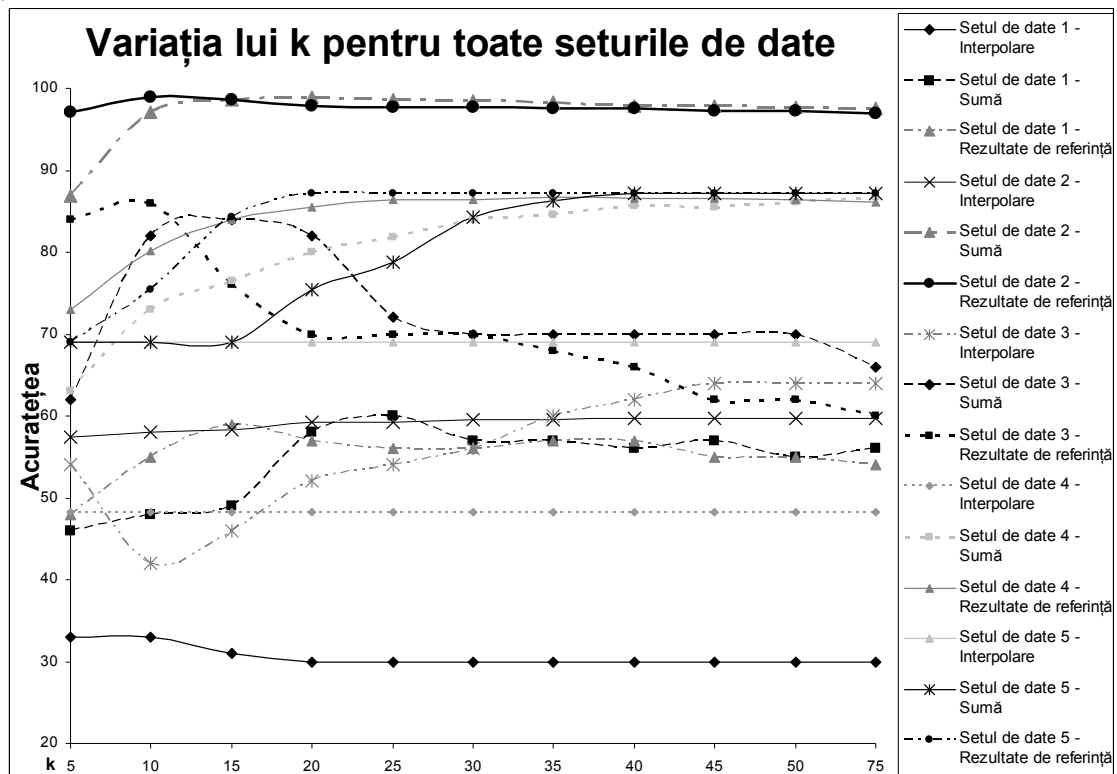


Figura 2 Variația lui „k” pentru toate seturile de date

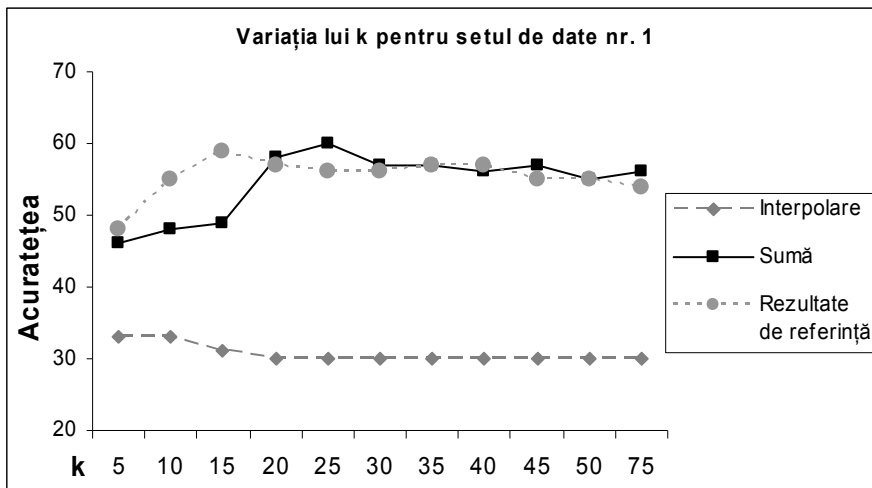


Figura 3 Variația lui „k” pentru setul de date numărul 1

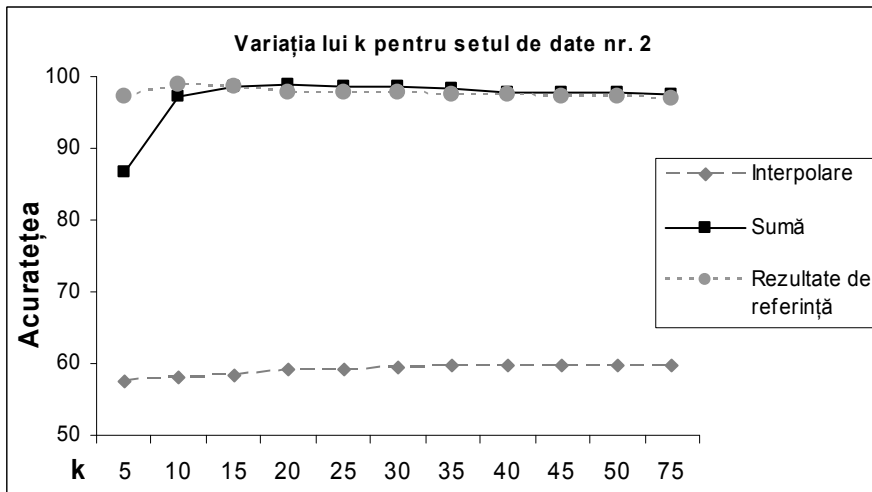


Figura 4 Variația lui „k” pentru setul de date numărul 2

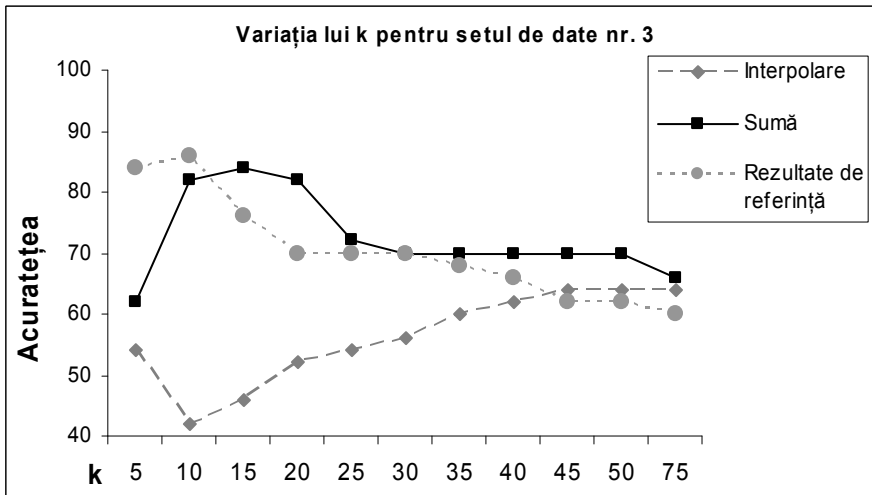


Figura 5 Variația lui „k” pentru setul de date numărul 3

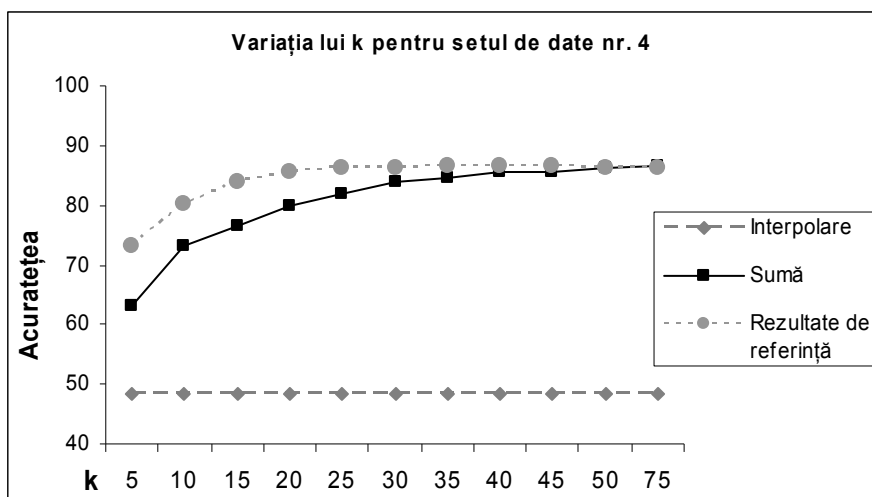


Figura 6 Variația lui „k” pentru setul de date numărul 4

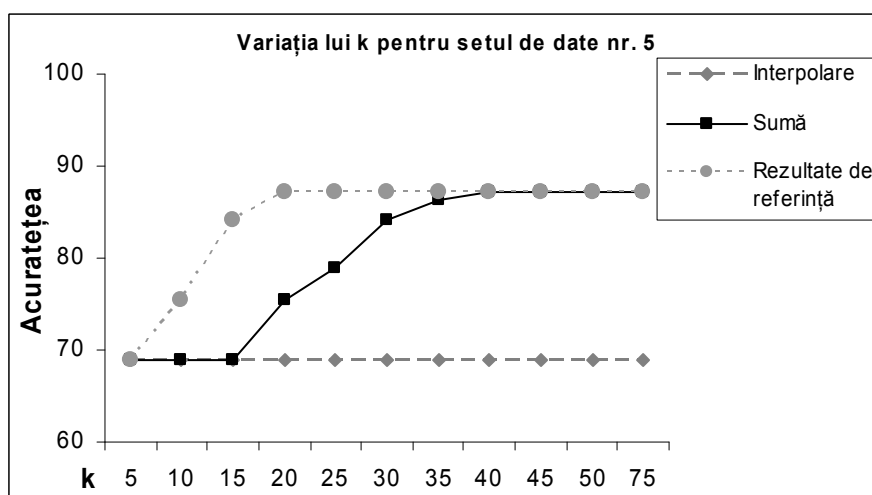


Figura 7 Variația lui „k” pentru setul de date numărul 5

Funcția considerată pentru interpolare în graficele prezentate mai sus a fost $\sin(1/(K \cdot \text{Avg}(D)))$. În graficul de mai jos prezentăm alte funcții de ponderare încercate. Pentru testarea diferitelor funcții de interpolare, am

reprezentat în grafic doar rezultatele obținute pe setul de date numărul 1 (Jocul Connect-4), pentru celelalte seturi de date curbele fiind similare.

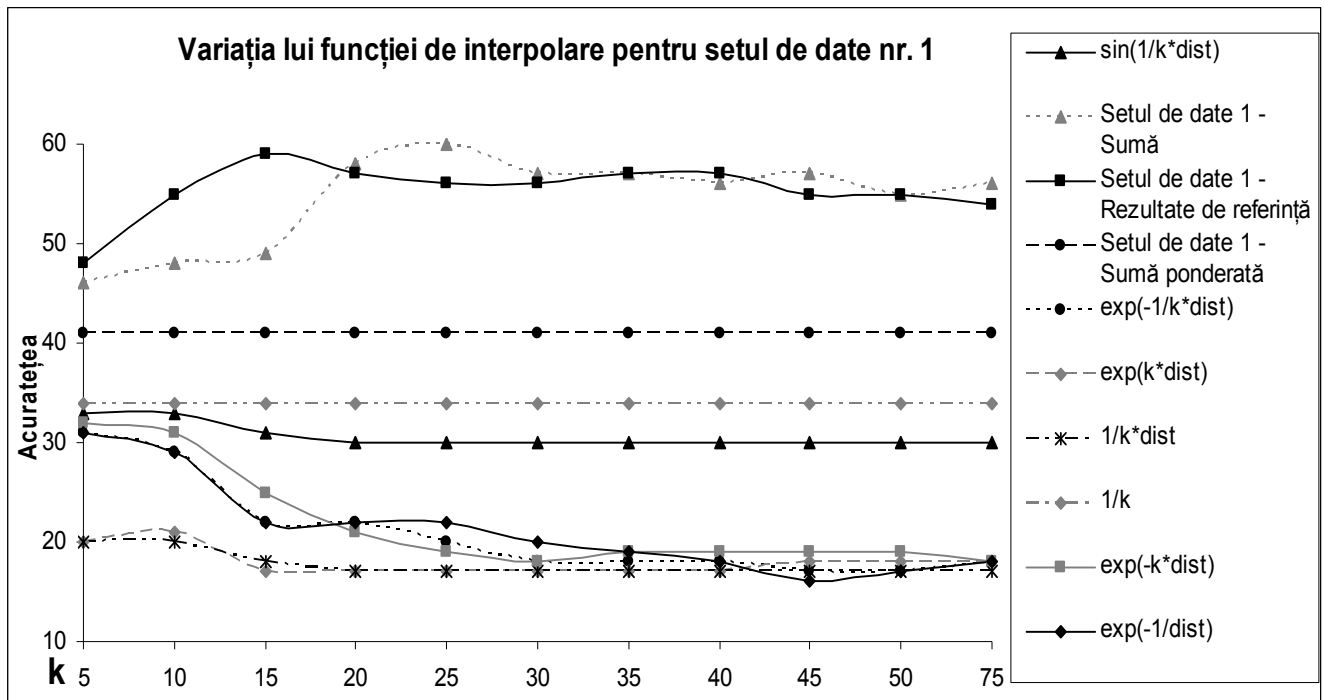


Figura 8 Variația funcțiilor de interpolare pentru setul de date numărul 1

CONCLUZII

În urma cu 100 de ani, Henry Ford a revoluționat industria manufacturieră prin introducerea producției de masă (fabricarea în serie a mai multor produse identice). În prezent, producția de masă a devenit un model învechit, companiile trebuind să ofere produse și servicii care să corespundă nevoilor cumpărătorilor [2]. În acest context, dezvoltarea de tehnologii avansate de recomandare ce pot fi aplicate unor produse și servicii configurabile astfel încât nevoile utilizatorilor să fie satisfăcute cât mai bine și mai ușor reprezintă un nou pas înainte. Unul dintre cei mai utilizați algoritmi în construcția de sisteme de recomandare este algoritmul KNN și orice îmbunătățire ce se poate aduce acestui algoritm poate atrage după sine o îmbunătățire a recomandărilor generate de sistemele care se bazează pe el. În cele ce urmează vom prezenta concluziile autorilor relativ la experimentele efectuate cu scopul de a îmbunătăți performanțele algoritmului KNN.

Graficele prezentate în secțiunea anterioară permit realizarea unei clasificări a metodelor descrise din punctul de vedere al acurateței. Astfel, cea mai mare acuratețe se obține fie folosind algoritmul clasic de KNN, fie suma clasificărilor rezultate din aplicarea algoritmului KNN pentru diferite valori ale lui „k”. Pe următorul nivel de acuratețe s-a situat o variantă de ponderare în funcție de valoarea lui „k” a sumei clasificărilor rezultate din aplicarea algoritmului KNN pentru diferite valori ale lui „k”.

Valori mai mici ale acurateței au fost obținute în momentul în care s-a încercat combinarea a două mărimi de ponderare (valoarea lui „k” și distanța medie față de instanța de test). Acest comportament este explicat de faptul că în momentul în care cele două mărimi sunt combinate, se combină și efectele pe care acestea le au asupra rezultatelor.

Cea mai slabă acuratețe a fost obținută de funcția $\exp(k \cdot \text{dist})$ deoarece aceasta este o funcție construită greșit din punct de vedere conceptual, dar care a fost intenționat folosită pentru a se vedea rezultatele obținute.

Testele referitoare la ponderare au demonstrat că această abordare nu este foarte utilă, deoarece în acest fel se diminuează mult prea mult importanța instanțelor mai puțin similare cu cea de test. Importanța acestora este deja micșorată prin considerarea lor de mai puține ori decât cele apropiate de instanța de test, așa că o nouă ponderare nu mai este necesară.

Testele efectuate pentru a identifica necesitatea considerării ca valoare de ponderare a distanței medii între instanța de test și cele de antrenare au arătat că situația descrisă în exemplul din Figura 1 este un caz particular și că în majoritatea cazurilor, considerarea acesteia are ca efect micșorarea acurateței. De aceea, considerăm că este mai bine să se evite această modalitate de ponderare.

Din rezultatele obținute se pot trage câteva concluzii importante:

rezultatele obținute folosind metoda votării între clasele returnate de aplicările algoritmului KNN pentru diferite valori ale lui „k” au fost mai bune decât cele returnate de algoritmul clasic de KNN pentru valori mari ale lui „k”;

cea mai bună soluție este să nu se folosească ponderări ale rezultatelor claselor întoarse de către algoritmul KNN, deoarece fiecare astfel de funcție de ponderare introduce efecte negative;

în cazul în care se combină mai multe funcții de ponderare se observa o scădere mai mare a acurateței comparativ cu situațiile în care acestea se folosesc individual;

ponderarea claselor întoarse de aplicările algoritmului KNN (cu diferite valori ale lui „k”) folosind ca măsură de

ponderare valoarea lui „k” este mai puțin dăunătoare pentru acuratețe decât ponderarea folosind distanța medie a exemplilor de antrenament față de instanța de test;

atunci când valoarea lui „k” devine prea mare, se încearcă construcția unui sistem mult prea simplu ce nu mai corespunde realității.

MULȚUMIRI

Mulțumim UC Irvine Machine Learning Repository pentru punerea la dispoziție a seturilor de date folosite la testarea variantei de îmbunătățire a algoritmului KNN propusă în lucrarea de față.

Cercetările prezentate în această lucrare au fost parțial finanțate din proiectul European FP7 STREP LTfLL și din proiectul național CNCSIS K-Teams.

REFERINȚE

1. Asuncion and D.J. Newman, (2007). UCI Machine Learning Repository, Irvine, CA: University of California, School of Information and Computer Science:
<http://www.ics.uci.edu/~mlearn/MLRepository.html>
2. Felfernig, A., Friedrich, G., Schmidt-Thieme, L.: Recommender Systems. In IEEE Intelligent Systems Special Issue on Recommender Systems, Vol. 22(3) (2007)
3. Herlocker, J. L., Konstan, J. A., Terveen, L. G., and Riedl, J. T. Evaluating collaborative filtering recommender systems. In: ACM Trans. Information Systems, vol. 22, no. 1, pp. 5-53 (2004)
4. T. Hofmann. Latent Semantic Models for Collaborative Filtering. ACM Trans. Information Systems, vol. 22, no. 1, 2004, pp. 89–115.
5. J. Konstan et al. GroupLens: Applying Collaborative Filtering to Usenet News. Comm.ACM, vol. 40, no. 3, 1997, pp. 77–87.
6. T. Mitchell. Machine Learning. McGraw Hill, 1997
7. Situl cursului de tehnologia limbajului al Universității din Antwerp: <http://ilk.uvt.nl/~antalb/ltua/week2.htm>
8. Situl cursului de învățare automată (CS6375) al Universității din Texas la Dallas:
<http://www.hlt.utdallas.edu/~vince/cs6375/assignments/hw1/>
9. <http://archive.ics.uci.edu/ml/datasets/Connect-4>
10. <http://archive.ics.uci.edu/ml/datasets/Mushroom>
11. <http://archive.ics.uci.edu/ml/datasets/Tic-Tac-Toe+Endgame>
12. <http://ilk.uvt.nl/~antalb/ltua/dimin.data.arff>