

NOI GENERAȚII DE SISTEME DE RECOMANDARE ȘI ALGORITMI PENTRU IMPLEMENTAREA LOR ÎN COMERȚUL ELECTRONIC

*Gheorghe MIȘCOI, Elena BĂDĂRĂU**

Universitatea Liberă Internațională din Moldova,

**Institutul de Relații Internaționale din Moldova*

În articol a fost efectuată o analiză a problemelor de bază ale sistemelor de recomandare, a metodelor de elaborare a noi sisteme de recomandare, mai adecvate proceselor reale, și a algoritmilor efectivi de realizare a lor.

Cuvinte-cheie: *sisteme de recomandare, procese semimarkoviene, priorități, flux Poisson, modelare matematică, algoritm Bayes, algoritmul aproximărilor succesive.*

NEXT-GENERATION SYSTEMS OF RECOMMENDATION AND ALGORITHMS FOR IMPLEMENTATION IN E-COMMERCE

In the article was given an analysis of the basic problems of recommendation systems, methods of developing new systems for recommendation, more adequately to real processes and effective algorithms of their realization.

Keywords: *recommendation systems, semimarkovien processes, priorities, Poisson flow, mathematical modeling, Bayes algorithm, algorithm of successive approximations.*

Introducere: Internetul contemporan oferă utilizatorului o cantitate enormă de informație diversificată, în care este tot mai dificil de a te orienta și naviga, iar sistemele clasice de căutare și sistematizare a informației nu mai verifică necesitățile utilizatorului. Ca urmare, au apărut așa-numitele Sisteme de Recomandare (SR), orientate spre furnizarea de informații la maxim posibil verificând interesele utilizatorului. SR reprezintă soft-uri prin intermediul cărora se încearcă de a elabora pronosticuri referitor la bunurile și/sau serviciile (cărți, filme, algoritmi, hoteluri, restaurante etc.) pe care le-ar dori utilizatorul având anumite informații despre profilul lui [3, 8].

SR au devenit obiectul atenției științifice începând cu mijlocul anilor '90 ai secolului XX, de la primele lucrări cu privire la filtrația colaborativă. În ultimul deceniu SR au fost dezvoltate atât sub aspect teoretic, cât și aplicativ, dar în acest domeniu rămân multe probleme, a căror soluționare promite vaste aplicări practice ce ar permite utilizatorului să se lupte cu volumul mare de informație și să-l înarmeze cu instrumente de elaborare a recomandărilor personificate. Exemple de aplicații ale SR pot servi sistemele de recomandare a cărților, CD-urilor și altor mărfuri pe Amazon.com, a filmelor pe Netflix și MovieLens, a noutăților pe VERSIFI Technologies etc.

Rezultatele cercetării și analiza: Perfecționarea în continuare a SR este necesară pentru o utilizare mai vastă și mai efectivă în situații reale pentru elaborarea de recomandări ce se referă, de exemplu, la călătorii, servicii financiare pentru investitori, mărfuri în magazine procurate prin intermediul cardului „inteligent” al consumatorului, diferite servicii. Devin necesare metode noi pentru urmărirea comportamentului consumatorului și informația propusă lui, modele performante de elaborare a recomandărilor ce includ informația referitor la contextul de procurare, utilizarea unui sistem multifactorial de estimare, reguli speciale de prioritate, modele de comportament virtual, precum și folosirea sau adoptarea diferitelor compartimente ale matematicii contemporane, ceea ce ar permite elaborarea unor SR mai flexibile și mai puțin sofisticate.

În acest context se înscrie și ipoteza despre evoluția piețelor. Nu demult era important pe cât e de calificativă marfa. Piețele mai complexe consideră important nu atât „valorile reale”, cât cum sunt ele percepute de consumatori. Acum începe o perioadă a „piețelor noi”, când primordial devine nu „valoarea” produsului sau chiar „percepția” lui, ci ce este cunoscut despre percepția acestui bun de alți consumatori. Motorul piețelor moderne constă în abilitarea consumatorului de a adera la consumul colectiv sau a se distanța de el; devine important câți oameni și cum percep produsul dat. Astfel, elaborarea și dezvoltarea SR, capabile să dea recomandări în „timp real”, să efectueze pronosticuri, devine stringentă pe piața tehnologiilor informaționale.

Pronosticul în SR se elaborează în baza datelor despre utilizator, care se obțin prin metode explicite (cereri de evaluare a obiectului de către consumator după o scară diferențială, cereri de aranjare a unui grup

de bunuri de la cel mai bun la cel mai rău, prezentarea consumatorului a două obiecte cu întrebarea care este mai bun, elaborarea unei liste de preferințe) și implicite (urmărirea vizualizării consumatorului internet-magazinului sau a diferitelor baze de date, înregistrarea comportamentului consumatorului on-line, conținutul calculatorului etc.). În rezultatul datelor colectate SR determină un set de recomandări pentru utilizatorul concret.

Formal, problema elaborării recomandărilor poate fi reprezentată astfel: fie C – un grup (set) de consumatori, S – un grup (set) de bunuri propuse. Grupul S poate fi foarte voluminos, atingând în unele domenii sute de mii sau milioane de unități de bunuri (cărți, CD etc.), grupul C la fel poate fi extrem de mare. Fie U – funcția de utilitate ce descrie utilitatea bunului s pentru c , $U : C \times S \rightarrow R$, unde R – numărul de bunuri comandate. Astfel, pentru fiecare consumator $c \in C$ dorim să alegem așa un bun $s \in S$ care ar fi cât mai util consumatorului.

În SR utilitatea unui bun este determinată pe cât i-a plăcut unui consumator acest bun. În funcție de aplicație, utilitatea u poate fi determinată de consumator (de regulă, în sistemele bazate pe evaluările consumatorului) sau poate fi calculată automat, similar funcției de utilitate, bazate pe beneficiu.

Fiecare membru din mulțimea C poate fi descris prin intermediul unui profil individual, care conține diverse caracteristici ale lui: vârsta, sexul, venitul, starea familială etc. În cel mai simplu caz profilul poate conține un singur parametru – **ID**–consumatorului. Sau, de exemplu, în SR a filmelor, unde S – totalitatea filmelor, fiecare film poate fi reprezentat nu doar prin **ID**-ul său, dar și prin denumire, gen, regizor sau actorii principali.

Cea mai importantă problemă a SR constă în faptul că utilitatea u , de regulă, nu este determinată pe întregul set $C \times S$. În SR utilitatea este reprezentată, de regulă, prin puncte și inițial definită doar pentru bunurile deja apreciate de consumator.

Extrapolarea de la estimările cunoscute la cele necunoscute are loc, de regulă, astfel: 1) prin alegerea unor reguli euristice ce determină funcția de utilitate și justifică comportamentul ei euristic sau 2) determinarea funcției de utilitate ce optimizează comportamentul parametrilor ei dați (abaterea medie pătratică).

Cum numai estimările necunoscute au fost analizate, consumatorul primește ca recomandare bunurile c cu cele mai înalte aprecieri [8]. SR pot determina sau totalitatea N de bunuri (cele mai potrivite) pentru consumator sau totalitatea consumatorilor relevanți bunului dat. Estimările noi pentru bunurile neestimate încă pot fi făcute prin diferite metode, de exemplu: euristice, teoria aproximării, luarea deciziilor.

SR diferă în abordarea analizei estimărilor. Clasificarea modernă a SR le împarte în următoarele categorii în funcție de modul cum sunt făcute recomandările:

- **Recomandări de content:** consumatorul primește recomandări de bunuri similare cu cele alese anterior;
- **Recomandări colaborative:** consumatorul primește recomandări de bunuri selectate anterior de persoane cu gusturi și cerințe similare;
- **Recomandări ce combină cele două metode anterioare.**

Metodici de content. În SR de content concluzia referitor la utilitatea $u(c, s)$ se ia reieșind din utilitatea $u(c, s_i)$, atribuită de consumatorul c bunului $s_i \in S$ similar bunului s . Mai formal, fie **Content (s)** profilul bunului s , adică totalitatea proprietăților bunului s . Acest set în continuare se folosește pentru a analiza utilitatea bunului pentru consumator. Cel mai frecvent, SR se utilizează pentru a recomanda bunuri ce conțin informație textuală. Conținutul acestor sisteme se descrie folosind cuvinte-cheie.

„Importanța” sau („informativitatea”) cuvântului k_i în documentul d_j se determină prin ponderea lui, w_{ij} , care poate fi determinată prin mai multe metode. Una din cele mai cunoscute metode de măsurare a ponderii cuvintelor în sistemele de căutare este așa-numita metoda frecvenței directe/indirecte, a cărei esență se reduce la: fie N – cantitatea de documente ce pot fi recomandate utilizatorului, cuvântul-cheie k_j se întâlnește în n_i documente; în plus, f_{ij} – va indica de câte ori cuvântul-cheie k_i se întâlnește în documentul d_j .

Atunci $TF_{i,j}$ – frecvența cuvântului-cheie k_i în documentul d_j , se determină prin formula

$$TF_{i,j} = \frac{f_{i,j}}{\max_z f_{z,j}}, \quad (1)$$

unde maximul se ia după toate frecvențele $f_{z,j}$ ale tuturor cuvintelor-cheie k_z , ce se întâlnesc în domeniul d_j . Cu toate acestea, cuvintele-cheie ce se întâlnesc frecvent nu ajută suficient pentru distingerea unui document

relevant de unul nerelevant. De aceea, măsurarea frecvenței inverse a cuvântului (IDF_i) este adesea utilizată de rând cu $TF_{i,j}$. Frecvența inversă a cuvântului-cheie k_i se definește de obicei prin

$$IDF_i = \log \frac{N}{n_i}, \quad (2)$$

ponderea cuvântului-cheie k_i în documentul d_j – prin

$$w_{i,j} = TF_{i,j} \cdot IDF_i, \quad (3)$$

iar conținutul documentului d_j – prin $Content(d_j) = (w_{1j}, w_{2j}, \dots, w_{kj})$.

Fie *Content Profil (c)* – profilul consumatorului c ce conține informația despre gusturile și preferințele acestui consumator. Așa tip de profiluri se creează ca rezultat al analizei conținutului (proprietăților) bunurilor deja estimate de consumatori și se bazează pe analiza cuvintelor-cheie. De exemplu, *Conținutul Profilului (c)* poate fi definit ca un vector de liste $(w_{c1}, w_{c2}, \dots, w_{ck})$, unde fiecare listă w_{ci} determină importanța cuvântului-cheie k_i și poate fi calculată pe baza vectorilor conținutului estimat individual.

În sistemele de conținut funcția de utilitate, de regulă, se definește ca

$$u(c, s) = \text{punctaj}(\text{conținutul profilului } (c), \text{Conținutul } (s)) \quad (4)$$

În utilizarea sistemului de căutare a informației de mai sus pentru recomandări de web-sait-uri, URL, mesaje ale rețelei Usenet, *Conținutul Profil (c)* pentru consumatorul c și *Conținutul (s)* a documentului s pot fi reprezentați ca vectori ai frecvenței directe/ indirecte a ponderii cuvintelor-cheie. În plus, în literatura axată pe problemele de găsire a informației, funcția de utilitate $u(c, s)$ de obicei se reprezintă ca un produs al unor algoritmi euristici, exprimați de vectorii \vec{w}_c și \vec{w}_s ,

$$u(c, s) = \frac{\vec{w}_c \cdot \vec{w}_s}{|\vec{w}_c| \cdot |\vec{w}_s|} = \frac{\sum_{i=1}^k w_{i,c} \cdot w_{i,s}}{\sqrt{\sum_{i=1}^k w_{i,c}^2} \sqrt{\sum_{i=1}^k w_{i,s}^2}}, \quad (5)$$

unde k – numărul general de cuvinte-cheie în sistem. De exemplu, dacă consumatorul c citește un număr mare de articole on-line pe bioinformatică, SR de conținut vor putea recomanda acestui utilizator și alte articole din acest domeniu, deoarece aceste articole conțin mai mulți termeni speciali (ghenom, secvențiere, proteomică etc.) decât în articole cu altă tematică și, prin urmare, *Conținutul Profil (c)* descris de vectorul \vec{w}_c va reprezenta așa termeni k_i cu o pondere w_{ic} mai mare.

În afară de metodele euristice tradiționale, bazate pe principiile de căutare a informației, există și alte metode-conținut. Aceste metode sunt caracterizate prin aceea că ele au la bază datele anterioare obținute cu ajutorul analizei statistice sau învățare automată. De exemplu, pe totalitatea paginilor web, estimate de consumator ca „utile” sau „inutile”, se definește clasificatorul Bayes pentru a clasifica paginile neestimate. Acest clasificator trebuie să prezică probabilitatea că pagina p_j aparține clasei C_i (adică, este importantă sau nu) reieșind din cuvintele-cheie $k_{1,j}, k_{2,j}, \dots, k_{n,j}$ de pe această pagină:

$$P(C_i / k_{1,j} \& k_{2,j} \& \dots \& k_{n,j}) \quad (6)$$

În [2], autorii reiese din presupunerea că toate cuvintele-cheie sunt independente și, prin urmare, probabilitatea (6) este proporțională cu

$$P(C_i) \prod_x P(k_{x,j} | C_i) \quad (7)$$

În plus, probabilitățile $P(k_{x,j} / C_i)$ și $P(C_i)$ pot fi estimate reieșind din datele deja existente. În așa caz, pentru fiecare pagină p_j probabilitatea (6) se calculează pentru fiecare clasă C_i și pagina p_j se atribuie la clasa C_i cu cea mai mare probabilitate.

Metodele colaborative. SR colaborative (sau sistemele de filtrare colaborativă) încearcă să prezică utilitatea bunului pentru un consumator aparte, reieșind din estimările făcute de alți consumatori. Mai formal, utilitatea $u(c, s)$ a bunului s pentru consumatorul c se estimează prin utilitățile $u(c_j, s)$ atribuite bunului s de acei consumatori c_j , care „seamănă” cu consumatorul c .

Conform [11], algoritmi pentru filtrația colaborativă pot fi divizați în două clase mari: amnestici (**memory-based**) și de model (**model-based**).

Algoritmi amnestici, bazându-se pe ipoteze euristice, fac pronosticuri referitor la estimările consumatorului reieșind din toate estimările precedente date de acest consumator.

Valoarea estimării necunoscute $r_{c,s}$ pentru consumatorul c și bunul s se calculează, de regulă, din totalitatea estimărilor bunului s date de alți consumatori. De exemplu:

$$r_{c,s} = \frac{1}{N} \sum_{c' \in C^*} r_{c',s} \quad (8)$$

$$r_{c,s} = k \sum_{c' \in C^*} \text{sim}(c, c') \cdot r_{c',s} \quad (9)$$

$$r_{c,s} = \bar{r}_c + k \sum_{c' \in C^*} \text{sim}(c, c') \cdot (r_{c',s} - \bar{r}_c), \quad (10)$$

unde: C^* – totalitatea a N consumatori cu gusturi similare consumatorului c , care au estimat deja bunul s ;
 k – coeficient de normare. De regulă,

$$k = \frac{1}{\sum_{c' \in C^*} |\text{sim}(c, c')|}.$$

Estimarea medie \bar{r}_c a consumatorului c se definește ca

$$r_c = \frac{\sum_{s \in S_c} r_{c,s}}{|S_c|}; \quad (11)$$

$\text{sim}(c, c')$ – măsură liniară a similitudinii consumatorilor c și c' .

În cel mai simplu caz, însumarea poate fi o medie simplă (8). Totodată, de obicei, pentru însumare se ia suma ponderată (9). Măsura liniară a similitudinii de regulă servește ca măsură de distanță dintre c și c' și se utilizează pentru măsurarea ponderii, adică cu cât sunt mai asemănători în gusturi c și c' cu atât mai mare va fi ponderea estimării $r_{c',s}$ în pronosticul $r_{c,s}$. Trebuie de reținut că $\text{sim}(x, u)$ reprezintă un artefact euristic, necesar pentru a determina totalitatea „celor mai apropiați clienți” pentru fiecare consumator.

În cazul aplicării (9) trebuie să se ia în considerare faptul că diferiți consumatori pot folosi diferite scări de evaluare. Pentru a ocoli această restricție putem utiliza suma ponderată *ajustată* (11). În așa caz, nu se consideră valorile absolute ale estimărilor, suma ponderată ia în considerare devierile acestor valori de la estimarea medie a consumatorului respectiv. Altă metodă de a evita diferențele în utilizarea scării de evaluare constă în filtrarea, bazată pe preferințe: se supun pronosticului preferințele utilizatorilor și nu valorile absolute ale estimărilor. În SR colaborative se utilizează diferite metode de calcul al similarității $\text{sim}(c, c')$ a consumatorilor. În majoritatea acestor metode similaritatea a doi consumatori se bazează pe faptul ce evaluări au făcut ei unora și acelorași bunuri.

Metoda corelațională și metoda asemănării liniare au o răspândire largă. Fie $S_{x,y}$ – mulțimea bunurilor estimate atât de consumatorul x , cât și de consumatorul y :

$$S_{x,y} = \{s \in S \mid r_{x,s} \neq \emptyset \ \& \ r_{y,s} \neq \emptyset\}.$$

În SR colaborative, $S_{x,y}$ se utilizează ca un rezultat intermediar la calculul așa-numiților „vecini apropiați” ai consumatorului x și frecvent este determinat prin calculul direct al intersecțiilor S_x și S_y . Metoda corelațională pentru calculul similarității folosește coeficientul de corelație Pearson:

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{x,y}} (r_{x,s} - \bar{r}_x)(r_{y,s} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{x,y}} (r_{x,s} - \bar{r}_x)^2 \sum_{s \in S_{x,y}} (r_{y,s} - \bar{r}_y)^2}} \quad (12)$$

În metoda asemănării liniare, ambii consumatori x și y pot fi reprezentați ca vectori ai spațiului m -dimensional, unde $m=|S_{x,y}|$. Atunci similitudinea a doi vectori poate fi determinată ca cosinusul unghiului dintre ei:

$$sim(x, y) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| \cdot |\vec{y}|} = \frac{\sum_{s \in S_{x,y}} r_{x,s} \cdot r_{y,s}}{\sqrt{\sum_{s \in S_{x,y}} r_{x,s}^2} \sqrt{\sum_{s \in S_{x,y}} r_{y,s}^2}}, \quad (13)$$

unde $\vec{x} \cdot \vec{y}$ – produsul scalar a doi vectori.

Atragem atenția că diferite SR pot folosi diferite abordări pentru un mod cât mai efectiv de determinare a similitudinii consumatorilor și analizei estimărilor obținute. Conceptul general acceptat constă într-un calcul preliminar al corespondențelor între toți utilizatorii sistemului și revizuirea lor periodică (așa necesitate nu apare frecvent, deoarece comunitatea utilizatorilor „vecini” este destul de constantă și nu se schimbă radical în perioade mici de timp).

Menționăm că atât metodele colaborative, cât și cele de content folosesc aceeași măsură a similitudinii liniare (cosinusul unghiului dintre vectori), fapt recomandat de literatura de căutare a informației. Cu toate acestea, SR de content utilizează măsura similitudinii liniare între vectorii ponderii frecvenței/frecvenței inverse, iar în SR colaborative similitudinea liniară se măsoară între vectorii estimărilor efectuate, specifice consumatorului.

Abordarea probabilistă a filtrației colaborative. În [2] se propune o abordare probabilistă a filtrației colaborative, unde estimările necunoscute se calculează astfel:

$$r_{c,s} = E(r_{c,s}) = \sum_{i=0}^n i \cdot P(r_{c,s} = i | r_{c,s'}, s' \in S_c), \quad (14)$$

unde expresia probabilistică înseamnă că consumatorul c va da o notă sau alta bunului s , reieșind din estimările precedente, iar evaluările sunt numere întregi, de la 0 la n . Pentru a aprecia această probabilitate, în [2] se propun două modele probabilistice alternative: *modele cluster* și *rețele Bayes*.

În primul model consumatorii sunt grupați în clustere. Reieșind din apartenența consumatorului unui sau altui cluster, evaluările lui sunt considerate independente, adică se poate aplica modelul clasic Bayes. Numărul de clustere și parametrii modelului sunt determinați de baza de date.

Al doilea model reprezintă fiecare bun ca unul din nodurile rețelei Bayes, în care poziția fiecărui nod corespunde estimării presupuse a bunului. Structura rețelei și probabilitățile unor sau altor rezultate depind de baza de date.

În literatura de specialitate întâlnim și alte abordări. Modelul statistic al filtrării colaborative a fost analizat în [10], unde au fost comparați diferiți algoritmi de evaluare a parametrilor modelului (clasterizarea k -dimensională și selecția Gibbs). Shani [9] aplică modelul soluțiilor succesive pentru a elabora recomandări și propune aplicarea lanțurilor Markov în același scop. Kumar [6] utilizează un model probabilistic simplu, pentru a demonstra că filtrația colaborativă lucrează și în cazul unei cantități relativ mici de informație despre consumator și că în unele situații, în lipsă de informație suficientă, algoritmi simpli de filtrare colaborativă sunt la fel de efectivi, ca și cei mai buni algoritmi bazați pe analiza utilității. În [12] se propune folosirea abordărilor probabilistice în calitate de metodă de reuniune a metodelor de model și a celor anamnestice. În particular, se propune

- abordarea studiului activ ce formează modele probabilistice pentru preferințele fiecărui consumator;
- în modelele complexe să se refere la profilurile păstrate ale consumatorilor, pentru a pronostica recomandări.

Metoda propusă dezvoltă ideile algoritmilor tradiționali anamnestici. Alte modele probabilistice de filtrare colaborativă includ analiza Bayes, modelul relativist de probabilitate, modelul regresiei liniare, modelul entropiei maxime.

Unele SR folosesc metode hibride ce combină abordările colaborative și de content, ceea ce permite de a evita restricțiile specifice fiecărui SR aparte. Reieșind din modul de a combina mecanismele de content și colaborativitate, metodele hibrid pot fi grupate în următoarele clase:

- realizarea aparte a algoritmilor de content și colaborativi și reuniunea pronosticurilor lor;
- încorporarea unor reguli de content în metodică colaborativă;

- încorporarea unor reguli colaborative în metodica de content;
- elaborarea unui model comun ce cuprinde regulile ambelor metodici.

Abordarea [1] sugerează implementarea analizei Bayes cu aplicarea lanțurilor Markov și a metodei Monte-Carlo. În procesul analizei, conform metodei Monte-Carlo calculatorul folosește o procedură de generare a numerelor aleatoare pentru imitarea datelor selecției generale studiate. Modelarea prin ecuații structurale alege eșantioane din selecția generală conform indicațiilor consumatorului, apoi efectuează următoarele acțiuni:

- simulează un eșantion aleatoriu din selecția generală;
- efectuează analiza rezultatelor;
- salvează rezultatele.

După un număr mare de iterații, rezultatele păstrate reflectă bine repartiția reală a statisticii de selecție.

Metoda Monte-Carlo permite obținerea informației și în cazurile când teoria repartițiilor selective este inutilă. În particular, în [1] informația referitor la profilurile consumatorilor se folosește într-un model statistic unic ce studiază estimările necunoscute r_{ij} pentru consumatorul i și bunul j .

$$r_{ij} = x_{ij} \mu + z_i \gamma_j + \omega_j \lambda_i + e_{ij}, \quad (15)$$

unde: $e_{ij} \sim N(0, \sigma^2)$, $\lambda_i \sim N(0, \Lambda)$, $\gamma_j \sim N(0, \Gamma)$ – variabile aleatoare ce iau în considerare erorile, necunoscute de izvoarele heterogene cu repartiții normale indicate; $i = \overline{1, I}$, $j = \overline{1, J}$ reprezintă, respectiv, consumatorii și bunurile; x_{ij} – matricea ce conține informația despre utilizatori și bunuri; z_i – vectorul caracteristicilor consumatorului; ω_j – vectorul caracteristicilor bunurilor. Parametrii necunoscuți ai acestui model – μ , σ^2 , Λ și Γ se determină după estimările deja primite utilizând lanțurile Markov și metoda Monte-Carlo. Astfel, în [1] se folosesc parametrii consumatorului $\{z_i\}$ ce formează profilul lui, parametrii bunului $\{\omega_j\}$ ce formează profilul bunului și interacțiunea lor $\{x_{ij}\}$ pentru analiza estimării bunului.

În plus, pe lângă metodele tradiționale de elaborare a profilului de consum (cum ar fi bazarea pe cuvintele-cheie și informația demografică), în ultimul timp au apărut metode noi ce au la bază prelucrarea automată a textelor (data-mining), analiza comportamentului de rețea etc.

Aceste metode permit să se ia în considerare interesele și preferințele consumatorilor și, ca rezultat, să extindă profilul consumatorului. Asemenea metode se aplică și pentru a extinde profilurile tradiționale ale descrierii bunurilor, care până nu demult se limitau la descrierea cu ajutorul cuvintelor-cheie.

Odată ce profilurile bunului și ale consumatorului au fost descrise, în caz general funcția de precizie a estimărilor poate fi definită prin aceste profile și estimările anterioare în modul următor. Fie profilul consumatorului i se definește ca un vector de p atribute, $\vec{c} = (a_{i1}, a_{i2}, \dots, a_{ip})$, profilul bunului j la fel se definește ca un vector de p atribute, $\vec{s}_j = (b_{j1}, b_{j2}, \dots, b_{jp})$. În mod deliberat nu definim valorile atributelor a_{ij} , b_{ki} , deoarece acestea pot fi diferite în diferite aplicații (de exemplu, numere, categorii, reguli, traiectorii etc.). Fie \vec{c} – vectorul tuturor consumatorilor, $\vec{c} = (\vec{c}_1, \vec{c}_2, \dots, \vec{c}_m)$, iar \vec{s} – vectorul tuturor bunurilor, $\vec{s} = (\vec{s}_1, \vec{s}_2, \dots, \vec{s}_n)$. Atunci, în general, procedura de precizie a estimării se reprezintă prin

$$r'_{ij} = \begin{cases} r_{ij}, & \text{dacă } r_{ij} \neq \emptyset \\ u_{ij}(R, \vec{c}, \vec{s}), & \text{dacă } r_{ij} = \emptyset \end{cases} \quad (16)$$

Această funcție prezice valoarea estimării r'_{ij} și o exprimă prin estimări deja cunoscute și profilurile \vec{c} și \vec{s} . Pentru a determina funcția de utilitate u_{ij} , putem folosi diverse metode: euristice, clasificarea după cel mai apropiat „vecin”, regresii, rețele neuronale, aproximări pe porțiuni, metode relaționale etc.

Majoritatea SR colaborative folosesc algoritmul probabilistic Bayes sau algoritmul SVD (or SVD++). Dar, ambii algoritmi necesită un eșantion voluminos pentru studiu. În [3] a fost efectuată o analiză comparativă și propus un algoritm alternativ ce nu cere eșantioane voluminoase. Compararea a fost făcută după criteriile de corectitudine a recomandărilor și operativitatea.

Algoritmul Bayes. Teorema Bayes [13] – una din teoremele de bază ale probabilității clasice, definește probabilitatea unui eveniment în condițiile unei informații incomplete. Probabilitatea condiționată a evenimentului A în condiția realizării evenimentului B se notează $P(A/B)$, conform [13]:

$$P(A/B) = \frac{P(A \cap B)}{P(B)},$$

unde $P(A \cap B)$ – probabilitatea realizării simultane a evenimentului A și B, $P(B)$ ($P(A)$) – probabilitatea realizării evenimentului B(A). Ușor de observat, că

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A) \quad (17)$$

Din (17) rezultă

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad (18)$$

Vom considera B – datele inițiale (informația cunoscută), A – parametrii modelului supus studiului. De exemplu, B – datele referitor la ratingul bunurilor, determinat de consumatori, A – factorii supuși studiului pentru consumatori și bunuri. Fiecare probabilitate are sensul său: $P(B|A)$ – repartiția probabilităților parametrilor modelului, după ce au fost luate în considerare datele inițiale (probabilități a posteriori). $P(A|B)$ – așa-numita verosimilitate. Neajunsul principal al algoritmului Bayes este cerința unor eșantioane voluminoase pentru analiză, și, în situația când probabilitatea condițională inițială este egală cu zero, probabilitatea pronosticată la fel va fi egală cu zero. Din aceste motive, la moment sunt preferabili alți algoritmi.

Algoritmul SVD. Se consideră o matrice R, ale cărei elemente sunt rating-uri (like-uri, cumpărături realizate etc.), pe care consumatorii (liniile matricei) au alocat bunurilor (coloanele matricei). De regulă, un consumator nu va estima un număr semnificativ de bunuri, de aceea matricea R va fi rarefiată puternic. Pentru așa matrice de regulă se folosește descompunerea singulară [3] în forma $R=UDV^T$. R este o matrice de dimensiuni mari $M \times N$, dar cu rangul mic și poate fi descompusă în produs de matrice $N \times f$ și $f \times M$, astfel micșorând numărul parametrilor de la $N \times M$ la $(N+M)f$.

Proprietatea de bază a SVD constă în aproximarea optimă, dacă în matricea D de lăsat primele f elemente diagonale, iar cele rămase de egalat cu zero.

$$XUDV^T = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \sigma_k \end{pmatrix} V^T \approx U \begin{pmatrix} \sigma_1 & \dots & 0 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \sigma_f & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & 0 & \dots & 0 \end{pmatrix} V^T \quad (19)$$

În matricea diagonală D, ce se află în mijlocul descompunerii singulare, elementele sunt ordonate: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k$, așa că egalarea cu zero a ultimelor elemente înseamnă a anula cele mai mici elemente, iar f se alege reieșind din valorile singulare ale matricei. În cazul SR se obține că fiecare consumator se reprezintă ca un vector din f factori U_i , iar fiecare bun – un vector din f factori V_j ; apoi, pentru a prognoza rating-ul atribuit de consumatorul i bunului j, considerăm produsul lor scalar $U_i V_j = U_i V_j$. Problema se formulează astfel: având estimările cunoscute ale bunurilor, să prognozăm cât de bine va estima un consumator nou fiecare bun. Se introduc așa-numiți predictorii de bază $b_{i,a}$ care se obțin din predictorii de bază ai fiecărui consumator aparte b_i , a unor bunuri aparte b_a , și din rating-ul mediu general a bazei μ : $b_{i,a} = \mu + b_i + b_a$, unde μ – rating-ul mediu general al bazei, b_j – rating-ul mediu al consumatorului i, b_i – rating-ul mediu al consumatorului fiecărui bun a. Se cere de găsit așa μ , b_i , b_a pentru predictorii de bază, pentru care $b_{i,a}$ cel mai

bine aproximează rating-urile date, apoi pot fi adăugați factorii. Dat fiind faptul că predictorii de bază au fost modificați, resturile sunt comparabile între ele, factorii se determină ca:

$$r_{i,a}^* = \mu + b_i + b_a + v_a^T u_i, \quad (20)$$

unde v_a – vectorul factorilor ce reprezintă bunul \mathbf{a} , u_i – vectorul factorilor ce reprezintă consumatorul \mathbf{i} . Acum problema poate fi formulată mai exact: trebuie de găsit cei mai buni predictorii, care aproximează $r_{i,a}^*$, adică de minimizat eroarea:

$$L(\mu, b_i, b_a, v_a, u_i) = \sum_{(i,a) \in D} (r_{i,a} - r_{i,a}^*)^2 = \sum_{(i,a) \in D} (r_{i,a} - \mu - b_a - v_a^T u_i)^2. \quad (21)$$

Funcția L se minimizează prin metoda gradientului.

Algoritm RNSA (The Refined Neighbor Selection Algorithm). Cu ajutorul algoritmului de clasterizare a K -mediilor se formează K clastere, fiecare constând din clienți ce au preferințe similare. Inițial se alege un consumator arbitrar, iar în calitate de punct inițial al centrului clusterului – k . Atunci fiecărui consumator i se atribuie un cluster, astfel încât distanța dintre el și k să fie maximă. În calitate de distanță poate fi utilizat coeficientul de corelație Pearson:

$$w_{u,a} = \frac{\sum_{i \in I} (r_{a,i} - \bar{r}_a)(r_{u,i} - \bar{r}_u)}{\sqrt{\sum_{i \in I} (r_{a,i} - \bar{r}_a)^2 \sum_{i \in I} (r_{u,i} - \bar{r}_u)^2}}, \quad (22)$$

unde $w_{u,a}$ – indică cât de aproape sunt consumatorii \mathbf{u} și \mathbf{a} , I – mulțimea bunurilor estimate atât de \mathbf{u} , cât și de \mathbf{a} ; $r_{a,i}$, $r_{u,i}$ – evaluările bunului i de consumatorul \mathbf{a} , respectiv \mathbf{u} , \bar{r}_a , \bar{r}_u – evaluările medii date de \mathbf{a} și \mathbf{u} .

După finisarea clasterizării se alege clusterul cu cel mai înalt coeficient de corelație Pearson. Se calculează pronosticul pentru toți consumatorii în acest cluster. Etapele principale ale algoritmului RNSA sunt:

- Intrare: clientul-test t , datele inițiale S
- Ieșire: vecini
 1. Crearea a K clastere din S prin metoda clasterizării.
 2. De găsit cel mai bun cluster C pentru t .
 3. De adăugat t în C și de considerat ca v .
 4. De adăugat v în lista vecinilor.
 5. Dacă lista vecinilor este suficientă, întoarcem lista vecinilor. În caz contrar, extragem v din C și pornim un proces de căutare în lățime. Asemănarea consumatorului se verifică conform (22), dacă ea primește valori într-un diapazon fixat, atunci v – consumator și se trece la pasul 4.

La pasul 5 algoritmul se oprește, dacă numărul consumatorilor în căutare în lățime este mai mare decât un careva număr fixat. Acest număr se determină prin diferite experimente. Formula de calcul al pronosticului

$$P_{a,j} = A(\bar{r}_{a,i}) + \frac{\sum_k (w_{a,k} (r_{k,j} - A(\bar{r}_{k,j})))}{\sum_k |w_{a,k}|}, \quad (23)$$

unde $P_{a,j}$ – pronosticul estimării, $w_{a,k}$ – gradul de apropiere între utilizatorii a și k , $A(\bar{r}_{a,j})$, $A(\bar{r}_{k,j})$ – estimările medii date de utilizatorii a și k respectiv.

Algoritmul RNSA se recomandă în cazul unor date inițiale insuficiente și permite obținerea rezultatelor mai exacte. În cazul eșantioanelor mari se recomandă utilizarea algoritmului SVD.

SR perspective bazate pe metodologii stocastice. Majoritatea metodelor de elaborare a SR poartă un caracter determinist, deși natura oricărui SR în mod obiectiv este stocastică. În adevăr, cerințele consumatorilor parvin aleatoriu, orice prelucrare durează un timp aleatoriu. Astfel, este necesară dezvoltarea metodelor de elaborare a SR bazate pe teoria probabilităților, teoria proceselor aleatoare, teoria așteptării, statisticii matematice etc. Nu se are în vedere aplicarea formulilor Bayes, care permit simplu de a recalcula probabilitățile a priori, după ce a avut loc un eveniment, ceea ce se folosește azi. La fel nu se iau în considerare SR

construite cu utilizarea lanțurilor Markov [13]. Indiscutabil, la o etapă anumită aceste metode și-au îndeplinit rolul și au contribuit esențial la dezvoltarea SR, acum pot fi aplicate și alte metode ale matematicii contemporane.

În SR pronosticul se elaborează în baza datelor referitor la consumator, care se obțin atât explicit, cât și implicit. Această informație generează o multitudine de date noi, pe care SR contemporane le ignorează, deși informația suplimentară, fiind prelucrată și formatizată, poate schimba cardinal SR și să aducă la elaborarea a noi SR, mai flexibile, mai adecvate proceselor reale. Această afirmație poate fi ilustrată [4] în baza exemplului de formalizare a fluxului cererilor consumatorilor.

Fie t_0, t_1, \dots, t_n , momente de timp (aleatoare) în care se înregistrează cererile consumatorilor, $z_k = t_k - t_{k-1}$, $k=1, 2, \dots, n, \dots$, $-z_k$ denotă intervalul de timp între două cereri consecutive și reprezintă o variabilă aleatoare continuă, $F_k(t)$ – funcția de repartiție a variabilei aleatoare z_k , $F_k(t) = P(z_k < t)$. Astfel, setul de funcții $F_k(t)$ formează fluxul cererilor consumatorilor. Se consideră că fluxul de cereri este definit, dacă pentru orice n este definită repartiția vectorului aleatoriu (z_1, z_2, \dots, z_n) . Cu alte cuvinte, fluxul este definit, dacă sunt date repartițiile $F_i(t)$, $i = \overline{1, n}$.

Pentru a reflecta această informație în SR, este necesară informația referitor la natura concretă a fluxului: dacă variabilele aleatoare z_k sunt independente în totalitate și identic repartizate, atunci $F_1(t) = F_2(t) = \dots = F_n(t) = F(t)$ (în așa caz, fluxul se numește flux recurent și este suficient de bine studiat). Foarte comod în utilizare este cazul $F(t) = 1 - e^{-\lambda t}$, unde λ – constantă, $\lambda > 0$ – așa-numitul flux Poisson. Pentru fluxurile Poisson se poate a priori prognoza numărul de cereri pe o perioadă fixă de timp. Conform [4], probabilitățile $P_n(t)$ în intervalul de timp $[0, t]$ sunt egale cu

$$P_n(t) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}.$$

Menționăm că aceste probabilități se calculează simplu și că și alte proprietăți ale fluxului Poisson pot fi considerate în viitoarele SR.

În SR poate fi considerată și dependența de trafic. Fie B – timpul de prelucrare a unei comenzi, ρ – coeficientul traficului, atunci

$$\rho = \frac{M(B)}{M(Z_k)},$$

unde $M(B) = \int_0^{\infty} t dB(t)$, $B(t) = P(B < t)$, $M(z_k) = \int_0^{\infty} t d(A_k(t))$, $A_k(t) = P(z_k < t)$. Este cunoscut [4] că

dacă $\rho < 1$, trafic nu există, iar pentru $\rho \geq 1$, SR vor lucra în condiții de trafic intens. Acest exemplu arată cum pot fi considerate informații suplimentare în elaborarea a noi modele de SR mai flexibile și mai adecvate proceselor reale.

Modele de SR cu pronostic virtual. În [15] a fost studiat un SR bazat pe aparatul transformării Fourier. Transformarea Fourier este bine cunoscută și are aplicații în fizica matematică, mecanică etc. În opinia noastră, pentru elaborarea SR mai productive se vor utiliza transformarea Laplace și Laplace-Stieltjes, care se determină astfel:

$$\overline{\alpha}(s) = \int_0^{\infty} e^{-sx} A(x) dx, \quad (I)$$

$$\alpha(s) = \int_0^{\infty} e^{-sx} dA(x) \quad (II)$$

În expresia (II) integrala se înțelege în sensul lui Stieltjes, iar funcția $A(t)$ – ca funcție de repartiție a unei variabile aleatoare. Una din proprietățile principale ale acestor transformări este descrisă de teoremele Tauber. Fie $P(z, t) = \sum_k k P_k(t)$ – funcția generatoare de repartiții, $P_k(t)$ – numărul de comenzi în momentul

de timp t , $\rho(z, t) = \int_0^{\infty} e^{-st} P(z, t) dt$ – transformata Laplace a funcției $P(z, t)$. Vom nota cu $P(z)$ funcția

generatoare staționară. Una din afirmațiile teoremei Tauber constă în următoarele: în condițiile unui regim staționar (dacă coeficientul de încărcare ρ este mai mic ca unu)

$$P(z) = \lim_{s \rightarrow 0} \rho(z, s),$$

ultima egalitate stabilește legătura dintre repartiția numărului de cereri în orice moment de timp și numărul de cereri în regim deja stabilit de funcționarea sistemului. Această dependență va ajuta la fundamentarea și va extinde domeniul de extrapolare a soluțiilor sistemelor de recomandare. În plus, prin integrala Stieltjes putem atribui un sens probabilistic transformatei Laplace-Stieltjes. În teoria sistemelor de deservire aceasta metodă este cunoscută sub denumirea „metoda catastrofelor” [14] și poate fi utilizată pentru determinarea strategiei SR.

Modele de SR cu priorități incluse. Este normal ca în anumite condiții unele cereri sau clase de cereri să pozeze anumite privilegii în deservire. Măsura acestor privilegii se numește prioritate. Prioritatea poate fi introdusă atât din exteriorul sistemului, cât și din interiorul lui, în dependență de anumite circumstanțe. Mai mult, pot fi considerate priorități dinamice – priorități ce se schimbă în dependență de starea sistemului. Dezvoltarea SR cu priorități poate fi justificată și de faptul că în [4] a fost demonstrat că din toate modelele SR legate de deservirea cererilor, deservirea cu priorități este optimală. O clasă largă de priorități – discipline generalizate de prioritate – a fost introdusă în [14]. Unele discipline noi de prioritate dinamică de clasă DD (Discretional Dynamical) au fost recent studiate în [5].

Concluzii

Avantajul disciplinelor de prioritate generalizate constă în faptul că modelele de deservire înzestrate cu asemenea priorități permit a considera diferite pierderi de timp necesare pentru a schimba modul de deservire, a restabili informația, a îndeplini unele operații auxiliare etc. ce au loc în sistemele reale.

Cel mai important avantaj al SR cu priorități constă în semimarkovitatea disciplinelor de prioritate generalizate, ceea ce înseamnă că ele sunt invariante în raport cu tipul funcțiilor de repartiție. Aceasta le face extrem de atractive și importante pentru diverse aplicații, deoarece ele permit modelarea unui spațiu larg de probleme, inclusiv evoluția diferitelor moduri de funcționare a sistemelor reale.

Bibliografie:

1. ANSARI, A., ESSEGAIER, S. and KOHLI, R. „Internet Recommendations Systems”. In: *J. Marketing Research*, 2000, p.363-375.
2. BREESE, J.S., HECKERMAN, D. and KADIE, C. „Empirical Analysis of Predictive Algorithms for Collaborative Filtering”. In: *Proc. 14th Conf. Uncertainty in Artificial Intelligence*, 1998.
3. GEDIMINAS, A., TUZHILIN, A. Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions IEEE. In: *Transactions on Knowledge and Data Engineering*, 2005, vol.17, no.6.
4. MISHKOY, Gh. *Textbook on Queueing Analysis*. ULIM, 2008.
5. MISHKOY, Gh., MITEV, L. Performance Characteristics for DD Priority Discipline with Semi-Markov switching. In: *Communications in Computer and Information Science*. Springer, 2014, p.2204-2018.
6. KUMAR, R., RAGHAVAN, P., RAJAGOPALAN S., and TOMKINS, A. Recommendation Systems: A Probabilistic Analysis. In: *J. Computer and System Sciences*, 2001, vol.63, no.1, p.42-61.
7. PAZZANI, M. and BILLSUS, D. Learning and Revising User Profiles: The Identification of Interesting Web Sites. In: *Machine Learning*, 1997, vol.27, p.313-331.
8. RICCI, F. *Systems Handbook*. Springer, 2011.
9. SHANI, G., BRAFMAN, R. and HECKERMAN, D. An MDP-Based Recommender System. In: *Proc. 18th Conf. Uncertainty in Artificial Intelligence*, 2002.
10. UNGAR, L.H. and FOSTER, D.P. *Clustering Methods for Collaborative Filtering*. Proc. Recommender Systems, Papers from 1998 Workshop, Technical Report WS-98-08, 1998.
11. WHITTLE, P. *Networks: Optimization and Evolution*. Cambridge Univ. Press, 2007.
12. YU, K., SCHWAIGHOFER, A., TRESP, V., XU, X. and KRIEGEL, H.-P. Probabilistic Memory-Based Collaborative Filtering. In: *IEEE Trans. Knowledge and Data Eng.*, 2004, vol.16, no.1, p.56-69.
13. МИШКОЙ, Г.К. *Вероятность и Математическая Статистика. Курс лекций*. Кишинэу: Elan Poligraf, 2012.
14. МИШКОЙ, Г. *Обобщенные приоритетные системы*. Академия Наук Молдовы, 2009.
15. ШТЕПИКОВ, Д.Г. и др. *Алгоритмы рекомендаций, основанные на Фурье-анализе профилей пользователей*. Санкт-Петербург: Телематика, 2009.
16. <http://www.machinlearning.ru/wiki/>

Prezentat la 01.06.2015