

# Evoluția opiniilor în forum-uri folosind o abordare semantică de sumarizare

**Dumitru Clementin Cercel**

Universitatea Politehnica București

Bd. Splaiul Independenței, nr. 313,  
București

clementin.cercel@gmail.com

**Ștefan Trăușan-Matu**

Universitatea Politehnica București

Bd. Splaiul Independenței, nr. 313,  
București

Institutul de Cercetări pentru Inteligența  
Artificială al Academiei Române

Calea 13 Septembrie nr.13, București

stefan.trausan@cs.pub.ro

## REZUMAT

Creșterea numărului de comentarii a devenit o problemă pentru consumatorii de produse, care trebuie să citească un număr semnificativ de comentarii pentru a identifica punctele tari și punctele slabe ale unui produs. Sumarizarea opiniilor este necesară pentru a ajuta consumatorii să înțeleagă, într-un mod simplificat, o cantitate mare de opinii. Grafurile reprezintă un model adecvat pentru a reprezenta forumurile și un mod eficient de a urmări evoluția opiniilor cu privire la o anumită entitate. În acest articol noi vom prezenta un algoritm de vizualizare a evoluției opiniilor într-un forum de discuții, ce folosește o abordare de sumarizare. Nu este vorba de o evoluție în timp a opiniilor din forum, ci de o evoluție a structurii grafului, din punct de vedere al opiniilor ce caracterizează fiecare nod.

## Cuvinte cheie

Evoluție opinii, Polaritate cuvânt, Polaritate propoziție, Sumarizare opinii, Stanford Typed Dependencies.

## Clasificare ACM

H5.2. Information interfaces and presentation (e.g., HCI): Miscellaneous.

## INTRODUCERE

Resursele web cuprind două categorii de informații: fapte și opinii [23]. Faptele presupun o prezentare obiectivă a realității, în timp ce opiniile reprezintă afirmații subiective despre evenimente, produse sau persoane. Extragerea de opinii ("Opinion mining") a reprezentat un domeniu de mare interes pentru cercetători odată cu creșterea popularității Web-ului, care generează, și în continuare, cantități mari de date textuale electronice. Volumul de resurse web, cum ar fi forum-urile, blog-urile, comentariile de pe site-uri gen e-commerce, reprezintă surse de opinii care pot fi folosite atât de persoanele fizice, cât și de organizații.

Analiza sentimentelor („Sentiment analysis”) are ca scop detectarea opiniilor exprimate în resursele web și clasificarea opiniilor în funcție de sentimentul exprimat în opinii pozitive și opinii negative. Abordările propuse de cercetători pentru a rezolva această ultimă sarcină denumită clasificarea sentimentelor („sentiment

classification”) folosesc clasificarea sentimentelor la nivel de document, la nivel de propoziție și la nivel de cuvânt [3, 7].

Clasificarea propozițiilor în propoziții subiective sau obiective este numită în literatura de specialitate clasificarea subiectivității („subjectivity classification”). Subiectivitatea nu este echivalentă cu existența unei opinii [17]. Putem avea propoziții subiective care să nu exprime opinii, dar putem avea cazuri în care propozițiile obiective implică opinii. Rezultatele prezentate în [18] arată că rezultatele analizei sentimentelor se pot îmbunătăți dacă această sarcină este ulterioară cea a clasificării propozițiilor.

În acest articol vom prezenta un algoritm de vizualizare a evoluției opiniilor într-un forum de discuții, care simplifică graful cu mesajele ce conțin opinii. În graful final, fiecare nod rezultat prin agregarea opiniilor de aceeași polaritate de la nodurile vecine conține polaritatea opiniei nodului și cuvintele care dau respectiva polaritate. Soluția propusă folosește o abordare de sumarizare de opinii. Am descompus problema în următoarele subsarcini:

- filtrarea grafului de mesaje în scopul de a elimina mesajele care nu se referă la entitatea avută în vedere, respectiv a elimina mesajele care nu conțin opinii despre această entitate;
- identificarea cuvintelor de opinie într-un mesaj. Vom considera că opiniile pot fi exprimate prin cuvinte precum adjective, adverbe și verbe;
- calcularea polarității mesajului dintr-un nod al grafului;
- agregarea nodurilor (pe orizontală și pe verticală) cu aceeași polaritate din graf.

Trăsăturile părților de vorbire ce caracterizează cuvintele dintr-o propoziție determină anumite relații semantice între cuvinte. În algoritmul propus relațiile gramaticale sunt folosite pentru a determina dependența între cuvântul a cărui evoluție a opiniilor dorim să o urmărim și cuvântul ce descrie opinia despre el. Astfel, relații de dependență semnificative au fost identificate folosind parserul Stanford [16] pentru a extrage propozițiile care conțin opinii referitoare la o entitate țintă.

Expunerea continuă cu descrierea cadrului și a abordărilor importante din domeniile existente pentru rezolvarea

subsarcinilor, acestea fiind urmate de formularea problemei și prezentarea metodei propuse. Rezultatele experimentelor, discuția despre analiza erorilor și concluziile sunt expuse în ultimele trei părți prezentate în lucrare.

## STAREA ACTUALĂ A DOMENIULUI

### Clasificarea sentimentelor

Detectarea sentimentului exprimat de un cuvânt de opinie poate fi făcută de un operator uman, folosind un corpus sau cu ajutorul unui dicționar [13].

Analiza perechilor de adjective ce sunt legate prin conjuncții precum „and”, „or”, „but”, „either-or” sau „neither-nor” este abordată în [8] pentru a deduce orientarea unui adjectiv. Folosind un set de documente, algoritmul de învățare supervizat constă din următorii pași:

- extragerea tuturor conjuncțiilor dintre adjectivele setului de documente;
- folosirea unui model de regresie log-linear pentru a clasifica perechile de adjective legate prin conjuncții ca având aceeași orientare sau diferită;
- folosirea unui algoritm de grupare pentru a împărți adjectivele în două clase. Clasa pentru care adjectivele componente au frecvența medie de apariție cea mai mare în setul de documente este etichetată ca fiind pozitivă.

Acuratetea rezultatelor în [8] este de 78,08% pentru un set de 657/679 de adjective de polaritate pozitivă, respectiv negativă.

Soluția propusă de autori în [10] folosește legăturile semantice între cuvintele din WordNet [21] pentru a defini o măsură a distanței între două cuvinte, dată de cea mai scurtă cale între ele. Pe baza acestei măsuri, orientarea semantică a unui adjectiv este calculată în funcție de distanța între respectivul adjectiv și cuvintele „good” și „bad”.

Metoda propusă în [19, 20] deduce orientarea semantică a unui cuvânt din legătura sa semantică cu alte cuvinte. În acest scop, autorii definesc două mulțimi de cuvinte ce au orientare semantică în general independentă de context și cuvintele componente sunt antonime: o mulțime de cuvinte cu polaritate pozitivă  $P_{\text{words}} = \{\text{“good”, “nice”, “excellent”, “positive”, “fortunate”, “correct”, “superior”}\}$  și o altă mulțime de cuvinte cu polaritate negativă  $N_{\text{words}} = \{\text{“bad”, “nasty”, “poor”, “negative”, “unfortunate”, “wrong”, “inferior”}\}$ . Măsurile folosite pentru a calcula distanța semantică între două cuvinte sunt: PMI („informație reciprocă punctuală”) și LSA („analiza semantică latentă”). Astfel, pentru un cuvânt de intrare „word” orientarea sentimentului exprimată de el este dată de expresia:

$$SO - A(\text{word}) = \sum_{pword \in P_{\text{words}}} A(\text{word}, pword) - \sum_{nword \in N_{\text{words}}} A(\text{word}, nword)$$

### Sumarizarea opiniilor

Sumarizarea opiniilor este diferită față de problema clasică de sumarizare de text, prin faptul că prima este centrată pe trăsăturile entității asupra căreia consumatorii și-au

exprimat punctul de vedere (opinia) și pe sentimentele trezite de aceste trăsături. Scopul sumarizării de opinii este acela de a ajuta utilizatorii să înțeleagă mulțimea de opinii care se găsește pe Web într-un mod simplificat.

Metodele de sumarizare de opinii pot fi împărțite astfel: metode de sumarizare bazate pe aspect și metode de sumarizare care nu sunt bazate pe aspect [11]. Prima categorie implică identificarea trăsăturilor care caracterizează un obiect (de exemplu, pentru obiectul "iPod" trăsăturile pot fi "durată baterie", "preț", "calitate sunet"), determinarea polarității opiniilor exprimate și prezentarea opiniilor sumarizate pentru fiecare trăsătură [22].

Un alt algoritm de sumarizare [12] extrage principalele subiecte din corpus prin atribuirea unei ponderi pentru fiecare cuvânt bazat pe tehnica tf-idf (term frequency - inverse document frequency), atât la nivel de paragraf, cât și la nivel de document. Sumarizarea se face per subiect și constă în extragerea de propoziții care conțin conceptele relevante obținute în pasul anterior, concepte care exprimă o opinie. Polaritatea unei propoziții se obține printr-o funcție proprie care combină sentimentele cuvintelor de opinie. Testarea algoritmului este făcută folosind un corpus de știri și blog-uri.

Ganesan et al. [6] propun un algoritm de sumarizare de opinii, bazat pe graf care elimină opiniile redundante.

Alți autori [14] propun o metodă de sumarizare folosind rating-ul global al comentariilor unui produs. Fiecare comentariu este reprezentat de o mulțime de perechi  $(w_m, w_h)$ , unde  $w_m$  reprezintă o trăsătură a entității respective și expresia  $w_h$  este opinia cu privire la trăsătură. Algoritmul este descris în trei pași:

- identificarea celor mai importante trăsături folosind PLSA („analiza semantică latentă probabilistică”) [9] și aplicarea algoritmului clasic de clustering k-means [15] pentru gruparea trăsăturilor;
- predicția rating-ului pentru fiecare trăsătură. Este calculată printr-o măsură de agregare propusă de autori folosind rating-ul frazelor care fac parte din același cluster, prin urmare având asociate aceeași trăsătură;
- extragerea frazelor reprezentative pentru fiecare trăsătură. Rating-ul unei fraze este calculat folosind două euristici: fiecare frază dintr-un comentariu are același rating ca și rating-ul global al comentariului, respectiv fiecare frază are rating-ul global al tuturor comentariilor.

### FORMULAREA PROBLEMEI

Un forum  $F$  constă dintr-o mulțime de fire de discuții:  $F = \{T_1, \dots, T_k\}$ . Un forum este caracterizat de o topică, ce reprezintă subiectul discuției. Astfel, un forum este caracterizat de setul de mesaje componente  $M$ , de setul de autori al mesajelor  $A$  și de setul de opinii pozitive și negative  $O$ , despre o entitate țintă. Firele de discuție componente ale unui forum au aceeași structură: un set de comentarii ale utilizatorilor, ca răspuns la unul din comentariile anterioare sau la subiectul discuției. Putem spune că acestea formează un graf orientat  $G = (V, E)$  unde:

- $V = \{ v_i \mid v_i = (v_i^m, v_i^a, v_{ij}^{op}, v_{ij}^{st}, v_i^t), v_i^m \in M, v_i^a \in A, v_{ij}^{op} \subset O, j \in \{0, 1, \dots, n_O\}, n_O = \text{numărul total de opinii în } v_i^m, i \in \{1, 2, \dots, n_M\}, n_M = \text{numărul total de mesaje în } M \} \cup \{v_0\}$  este setul de noduri. Nodul  $v_0$  este un nod special introdus, fiind caracterizat doar de subiectul de discuție;
- $E = \{e_1, e_2, \dots, e_{n_M}\}$  este setul de arce, astfel încât dacă  $e = (v_i, v_j) \in E$  atunci mesajul  $v_i$  este replică la mesajul  $v_j$ .

Fiecare nod din graf reprezintă un mesaj  $v_i^m$  din respectivul fir, mesaj scris de un autor  $v_i^a$  și conține opiniile  $v_{ij}^{op}$  despre entitatea respectivă. Variabila  $v_i^t$  este momentul de timp când mesajul a fost scris.

Fiecare opinie din  $v_{ij}^{op}$  este caracterizată de tăria sentimentului  $v_{ij}^{st}$  pe care îl exprimă. Intensitatea opiniilor este exprimată printr-o funcție  $s_{op} : O \rightarrow (-1, 1) \setminus \{0\}$ , unde  $s_{op}(o) \in (0, 1)$  dacă opinia este pozitivă și  $s_{op}(o) \in (-1, 0)$  dacă opinia este negativă. Asociem fiecărui mesaj din graf initial un scor  $msgScore(v_i^m)$ , dat de suma dintre puterea sentimentelor exprimate de opiniile pozitive și opiniile negative, ce sunt extrase din mesaj împărțit la numărul total de opinii din mesaj:

$msgScore(v_i^m) : M \rightarrow (-1, 1), n_0^k + n_0^l = n_0$

$$msgScore(v_i^m) = \frac{\sum_{k=1}^{n_0^k} s_{op}(v_{ik}^{op}) + \sum_{l=1}^{n_0^l} s_{op}(v_{il}^{op})}{n_0}$$

**Definiție (agregare noduri pe verticală):** Două noduri din graf  $v_i = (v_i^m, v_i^a, v_{ij}^{op}, v_{ij}^{st}, v_i^t)$  și  $v_k = (v_k^m, v_k^a, v_{kl}^{op}, v_{kl}^{st}, v_k^t)$  vor fuziona dacă opiniile  $v_{ij}^{op}$  și  $v_{kl}^{op}$  au aceeași polaritate (pozitivă sau negativă) și  $(v_i, v_k) \in E$ , rezultatul fiind un singur nod  $v_i' = (v_{rs}^m, v_{rs}^a, v_{rs}^{op}, v_{rs}^{st}, v_{rs}^t)$ , caracterizat prin  $v_{rs}^m = v_i^m \cup v_k^m, v_{rs}^a = v_i^a \cup v_k^a, v_{rs}^{op} = v_{ij}^{op} \cup v_{kl}^{op}, v_{rs}^t = v_i^t \cup v_k^t$  și  $v_{rs}^{st} = v_i^{st} \cup v_k^{st}$ .

**Definiție (agregare noduri pe orizontală):** Două noduri din graf  $v_i = (v_i^m, v_i^a, v_{ij}^{op}, v_{ij}^{st}, v_i^t)$  și  $v_k = (v_k^m, v_k^a, v_{kl}^{op}, v_{kl}^{st}, v_k^t)$  vor fuziona dacă opiniile  $v_{ij}^{op}$  and  $v_{kl}^{op}$  au aceeași polaritate (pozitivă sau negativă) și  $\exists v_k^l \in V$  astfel încât  $(v_i, v_k^l) \in E$  și  $(v_k, v_k^l) \in E$  rezultatul fiind un singur nod  $v_i' = (v_{rs}^m, v_{rs}^a, v_{rs}^{op}, v_{rs}^{st}, v_{rs}^t)$ , caracterizat prin  $v_{rs}^m = v_i^m \cup v_k^m, v_{rs}^a = v_i^a \cup v_k^a, v_{rs}^{op} = v_{ij}^{op} \cup v_{kl}^{op}, v_{rs}^t = v_i^t \cup v_k^t$  și  $v_{rs}^{st} = v_i^{st} \cup v_k^{st}$ .

Pentru fiecare nod rezultat prin agregare (pe verticală sau pe orizontală) se va calcula valoarea normalizată:

$$msgScore(v_i') = \frac{msgScore(v_i) + msgScore(v_k)}{2}$$

Vizualizarea opiniilor într-un forum de discuții se poate face în mod simplificat într-un graf final  $G' = (V', E')$  unde:

- $V' = \{ v'_i \mid v'_i = (M_i, A_i, O_i, S_i, T_i), M_i \subset M, A_i \subset A, O_i \subset O, i \in \{1, 2, \dots, n_M\}, n_M < n_M \}$ . Un nod din graf este un ansamblu de mesaje  $M_i$  scrise de un set de autori  $A_i$ , la momentele de timp  $T_i$ . Opiniile  $O_i$  despre o entitate considerată având tăria sentimentelor  $S_i$ , corespund la aceste mesaje.

- $E' = \text{setul de arce, astfel încât dacă } e' = (v'_i, v'_j) \in E'$  atunci setul de mesaje corespunzătoare nodului  $v'_i$  sunt replică la mesajele nodului  $v'_j$ .

Graful final rezultat va fi caracterizat de următoarea proprietate importantă: orice cale de la nodul  $v_0$  la un nod frunză este o alternanță de noduri de sentimente de polaritate pozitivă și negativă.

## SISTEMUL PROPUȘ

Algoritmul propus primește ca intrare corpul de forumuri și entitatea pentru care dorim să urmărim evoluția opiniilor. Entitatea poate fi un substantiv comun sau propriu. Pentru entitățile care constau din substantive comune, algoritmul va ține cont atât de entitatea în sine, cât și de sinonimele sale. Este folosit în acest scop WordNet pentru a obține o listă de sinonime pentru un substantiv.

Folosind parserul Stanford vom obține, pentru un mesaj din graf, relațiile de dependență care evidențiază perechile compuse din entitatea avută în vedere și cuvântul care conține opinia despre entitate. Dependențele între cuvinte sunt relații gramaticale binare între un cuvânt conducător („governor”) și un cuvânt dependent (de exemplu, pentru dependența „nsubj (interesting, movie)” adjectivul “interesting” este cuvântul conducător, iar substantivul “movie” este partea de vorbire dependentă). Parserul Stanford poate detecta 55 de relații gramaticale care pot fi prezente într-o propoziție. De asemenea, obținem în acest pas părțile de vorbire pentru fiecare cuvânt (de exemplu, JJ – adjectiv; VBN – verb participiu trecut). Vom lua în considerare următoarele relații de dependență:

- Nsubj („nominal subject”). Am considerat cazul în care cuvântul conducător al relației este orice adjectiv și cuvântul dependent este entitatea țintă. De exemplu, în propoziția: “The movie is interesting”, reprezentarea dependențelor Stanford sunt: det (movie, The), nsubj (interesting, movie), cop (interesting, is), root (ROOT, interesting);
- Acomp („adjectival complement”). Cuvântul conducător al relației acomp este un verb la orice timp și partea dependentă este un adjectiv. De asemenea, verbul este în relație Nsubj cu entitatea urmărită. Verbul prezintă starea, acțiunea entității, în timp ce adjectivul exprimă opinia cu privire la această funcție. Pentru propoziția: “She looks very beautiful.” relațiile de dependență sunt: nsubj (looks, She), root (ROOT, looks), advmod (beautiful, very), acomp (looks, beautiful);
- Advmod („Adverbial Modifier”). Am considerat cazurile în care cuvântul conducător al relației advcomp este un verb la orice timp sau adjectiv și partea dependentă este un adverb.
- Neg („negation modifier”), preconj(preconjunct), conj\_nor (conjunct negation). În aceeași etapă considerăm adverbe („not” and „neither nor”) care schimbă polaritatea de opinie. În cazul în care adjectivul este determinat de o negație, polaritatea adjectivului este schimbată. De exemplu, în propoziția: “The movie is not interesting.”, relațiile de dependență sunt: det (movie, The), nsubj (interesting, movie), cop (interesting, is), neg (interesting, not),

root (ROOT, interesting), în timp ce în propoziția: “The movie is neither amazing nor funny.”, dependențele sunt: det (movie, The), nsubj (amazing, movie), nsubj (funny, movie), cop (amazing, is), preconj (amazing, neither), root (ROOT, amazing), conj\_nor (amazing, funny).

Cu privire la cuvintele de opinie găsite, determinăm polaritatea și tăria polarității lor folosind SentiWordNet [4]. SentiWordNet este o resursă lexicală în care pentru fiecare mulțime de sinonime („synset”) s din WordNet este asociat trei scoruri numerice: Obj(s) (indică gradul de obiectivitate al cuvântului), Pos(s) (indică gradul de pozitivitate al cuvântului), și Neg(s) (indică gradul de negativitate al cuvântului). Cele trei scoruri variază de la 0 la 1 și suma lor este egală cu 1. Valorile SentiWordNet au fost obținute prin utilizarea unui algoritm de clasificare semi-supervizată și analiza definițiilor conceptelor („glosses”) ce sunt asociate cu mulțimile de sinonime SentiWordNet.

Pentru obținerea polarității unui mesaj este nevoie de a calcula scorul său folosind formulele din prezentarea problemei. Polaritățile mesajelor pot fi reprezentate ca o funcție msgPolarity :  $M \rightarrow \{-1, 0, +1\}$ . Funcția are valoarea -1 când opinia mesajului este negativă, valoarea 0 când opinia este neutră și +1 când opinia este pozitivă. Astfel, mesajul unui nod al grafului poate fi clasificat în mesaj de polaritate pozitivă sau negativă.

Algoritmul ce identifică polaritatea opiniei unui mesaj din graf este prezentat mai jos:

*Date intrare: entitatea considerată e, un mesaj  $v_i^m$  din forum, n relații de dependență de forma relație<sub>i</sub>( $g_i, d_i$ ) între cuvintele mesajului  $v_i^m$*

*Date ieșire: opinie mesaj forum*

```

1. scor_opinie_msg <- 0
2. for i = 1 to n
3.   if relației = nsubj și gi = adjectiv
4.     și di = e
5.       scor_opinie_msg +=intensitate
6.       opinie gi obținută folosind
7.       SentiWordNet
8.   end if
9.   if relației = acomp și gi = verb și
10.    di = adjectiv și există relațiej(gj,
11.    dj) astfel încât relațiej = nsubj și
12.    gj = verb și dj = e și gi = gj
13.      scor_opinie_msg += intensitate
14.      opinie di obținută folosind
15.      SentiWordNet
16.   end if
17.   if relației = advmod și (gi = verb sau
18.    gi = adjectiv) și di = adverb și
19.    există relațiej(gj, dj) astfel încât
20.    relațiej = nsubj și gj = adverb și
21.    dj = e și di = gj
22.      scor_opinie_msg += intensitate
23.      opinie di obținută folosind
24.      SentiWordNet
25.   end if
26.   if relației = neg și gi = adjectiv și
27.    există relațiej(gj, dj) astfel încât

```

```

28.   relațiej = nsubj și gj = adjectiv și
29.   dj = e și gi = gj
30.     scor_opinie_msg += (-1) *
31.     intensitate opinie gi obținută
32.     folosind SentiWordNet
33.   end if
34.   if relației = conj_nor și gi =
35.   adjectiv și di = adjectiv și există
36.   relațiej(gj, dj) astfel încât relațiej
37.   = nsubj și gj = adjectiv și dj = e și
38.   gi = gj și există relațiek(gk, dk)
39.   astfel încât relațiek = nsubj și
40.   gk = adjectiv și dk = e și di = gk
41.     scor_opinie_msg += (-1) *
42.     (intensitate opinie gi +
43.     intensitate opinie gk obținute
44.     folosind SentiWordNet)
45.   end if
46. end for
47. if scor_opinie_msg > 0
48.   polaritate pozitivă mesaj
49. end if
50. if scor_opinie_msg < 0
51.   polaritate negativă mesaj
52. end if
53. if scor_opinie_msg = 0
54.   polaritate neutră mesaj
55. end if

```

Primul pas al algoritmului constă în eliminarea nodurilor care nu se referă la entitatea avută în vedere, precum și cele care nu conțin opinii despre entitate. Celelalte noduri din graf nu vor fi modificate. Aplicarea celor două definiții de reducere a nodurilor va determina o serie de grafuri intermediare, în care nodurile reprezintă mesaje ce conțin opinii despre entitate. După aplicarea unei reduceri vom restabili legăturile dintre nodurile afectate. Graful final, obținut când cele două definiții nu mai pot fi aplicate, este folosit pentru vizualizarea evoluției opiniilor din forumul inițial.

## REZULTATE EXPERIMENTALE

Corpusul pe care l-am folosit la evaluarea sistemului propus este o colecție de 5 forum-uri din domeniul politicii [5] extrase de pe site-ul The Huffington Post. Întrucât articolul este un studiu preliminar în acest subiect am ales să folosim o mulțime de forum-uri de dimensiune redusă iar etichetarea mesajelor cu opinii necesită un volum mai redus de timp.

Fiecare mesaj din forumurile considerate a fost adnotat manual dacă conține sau nu opinii despre entitatea urmărită. Atât în algoritmul propus cât și în etichetarea manuală a forumurilor suntem interesați de propozițiile care conțin opinii la modul general despre o entitate. Nu vom căuta și opinii despre trăsături ale sale. De asemenea, nu vom trata cazurile când opiniile sunt exprimate prin informație non-textuală (de exemplu emoticoane).

Entitatea avută în vedere în experimente este "Barack Obama". Tabelul 1 arată numărul de noduri ale grafului după eliminarea mesajelor cu opinii, care nu se referă la entitatea vizată (notată ca fiind condiția 1), eliminarea nodurilor care conțin informații fără opinie despre entitate (notată ca fiind condiția 2), precum și agregarea de noduri

care includ opinii de aceeași polaritate (notată ca fiind condiția 3). În tabelul 2 vom prezenta precizia și acuratețea între adnotarea manuală de opinii pe graful inițial și opiniile calculate de algoritm pentru graful obținut după aplicarea condiției 2.

Tabelul 1. Rezultate dimensiuni grafuri obținute

	Număr de noduri graf inițial	Număr de noduri aplicând cond. 1	Număr de noduri aplicând cond. 2	Număr de noduri aplicând cond. 3
Forum 1	100	24	9	4
Forum 2	100	11	5	3
Forum 3	100	18	6	3
Forum 4	100	13	5	2
Forum 5	100	25	7	4

Tabelul 2. Performanțele analizei sentimentelor în forum-urile considerate

	Precizie	Acuratețe
Forum 1	0.54	0.52
Forum 2	0.66	0.64
Forum 3	0.63	0.54
Forum 4	0.75	0.63
Forum 5	0.45	0.55

## DISCUȚIE

Analiza erorilor procesului de extragere de opinii este influențată de mai mulți factori:

- polaritatea cuvintelor care exprimă sentimente depinde de domeniul și de contextul în care sunt utilizate [17]. Într-un context specific, adjectivele pot avea un sens pozitiv, dar într-un alt context dat, sensul poate fi negativ sau neutru. Aceeași expresie poate avea diferite polarități ale sentimentului exprimat în diferite domenii. De exemplu, în ceea ce privește expresia „citește cartea!”, într-un domeniu despre cărți sensul este favorabil, dar într-un domeniu despre filme sensul este negativ. Chiar și în același domeniu, depinde de trăsătura urmărită (de exemplu: expresia „durata mare a bateriei telefonului” comparativ cu expresia „timpul de incarcare mare al telefonului”);
- în aceeași propoziție pot fi opinii multiple care exprimă aceeași polaritate, dar pot fi prezente și opinii multiple de polarități diferite [17];
- detectarea opiniilor implicite, date de expresii retorice, metafore (de exemplu: „Mâncarea a fost ca o piatră”), dubla negație. Un cuvânt pozitiv poate avea un sens negativ într-un context metaforic sau ironic [2].

## CONCLUZII ȘI DIRECȚII VIITOARE

Sistemele de extragere de opinii au fost dezvoltate îndeosebi pentru comentariile produselor aflate pe site-uri de comerț electronic, fiind mai puține în cazul discuțiilor interactive. Pentru a identifica într-un text opiniile și a le clasifica în: pozitive, negative sau neutre este necesar să se dezvolte metode complexe.

În viitor dorim să îmbunătățim rezultatele metodei propuse pentru vizualizarea evoluției opiniilor într-un forum, prin propunerea unui algoritm îmbunătățit de detectare a polarității pentru o propoziție în care sunt exprimate opinii de sentimente diferite, precum și a unui algoritm de identificare a opiniilor implicite, luând în calcul contextul în care apar opiniile. De asemenea, dorim să efectuăm o analiză a extragerii de opinii din perspectiva modelului polifonic [1] deoarece este posibil ca unele opinii să se influențeze reciproc.

## MULȚUMIRI

Adresăm mulțumiri profesorului Julien Velcin, Universitățile Lumière Lyon2, pentru ajutorul său în realizarea acestei metode.

Rezultatele prezentate în acest articol au fost obținute cu sprijinul Ministerului Muncii, Familiei și Protecției Sociale prin Programul Operațional Sectorial Dezvoltarea Resurselor Umane 2007-2013, Contract nr. POSDRU/107/1.5/S/76903

## REFERINȚE

1. M.M. Bakhtin, *Speech Genres & Other Late Essays*, University of Texas Press, Austin, 1986.
2. D. Davidov, O. Tsur, A. Rappoport, *Semi-supervised Recognition of Sarcastic Sentences in Twitter and Amazon*, 14th Conference on Computational Natural Language Learning, Uppsala, Sweden, 2010.
3. X. Ding, B. Liu, and P. S. Yu, *A holistic lexicon-based approach to opinion mining*, In *Proceedings of the 2008 International Conference on Web Search and Data Mining*, Palo Alto, California, USA, 2008.
4. A. Esuli, and F. Sebastiani, *SENTIWORDNET: A publicly available lexical resource for Opinion Mining*, LREC, 2005.
5. M. Forestier, J. Velcin, D. A. Zighed, *Extracting Social Networks to Understand Interaction*, *International Conference on Advances in Social Network Analysis and Mining (ASONAM'11)*, vol. (2011), pp. 213-219, 2011.
6. K. Ganesan, C. Zhai, and J. Han, *Opinion: a graph-based approach to abstractive summarization of highly redundant opinions*, In *Proceedings of the 23rd International Conference on Computational Linguistics. COLING '10*. Association for Computational Linguistics, Stroudsburg, PA, USA, 2010, pp. 340-348.
7. A. B. Goldberg and X. Zhu, *Seeing stars when there aren't many stars: graph-based semi-supervised learning for sentiment categorization*, In *Proceedings of the First Workshop on Graph Based Methods for Natural Language Processing*, 2006.
8. V. Hatzivassiloglou and K. R. McKeown, *Predicting the semantic orientation of adjectives*. In *Proceedings of ACL 1997*, pp.174-181, 1997.
9. T. Hofmann, *Probabilistic latent semantic indexing*, In *Proceedings of SIGIR '99*, pp. 50-57, 1999.
10. J. Kamps and M. Marx, *Words with attitude*. In *Proceedings of the First International Conference on Global WordNet*, pp.332-341, 2002.

11. H.D. Kim, K. Ganesan, P. Sondhi, and C. Zhai, *Comprehensive Review Of Opinion Summarization*, 2011.
12. L.-W. Ku, Y.-T. Liang, and H.-H. Chen, Opinion extraction, summarization and tracking in news and blog corpora, In *AAAI Symposium on Computational Approaches to Analysing Weblogs (AAAI-CAAW)*, pp. 100–107, 2006.
13. B. Liu, M Liu B, Tutorial on sentiment analysis, based on Chapter 11 of the book “Web Data Mining – Exploring Hyperlinks, Contents and Usage Data”, 2007.
14. Y. Lu, D. Huizhong, W. Hongning, and Z. ChengXiang, Exploiting Structured Ontology to Organize Scattered Online Opinions, In *Proceedings of International Conference on Computational Linguistics (COLING-2010)*, 2010.
15. J. B. MacQueen, Some methods for classification and analysis of multivariate observations, In *Proceedings of the fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, pp. 281-297. University of California Press, 1967.
16. M.-C. de Marnee, and C.D. Manning, *Stanford Typed Dependencies Manual*, <http://nlp.stanford.edu/software/dependencies/manual.pdf>, 2008.
17. B. Pang, and L. Lee, Opinion mining and sentiment analysis, *Foundations and Trends in Information Retrieval*, vol. 2, pp. 1–135, 2008.
18. B. Pang and L. Lee. A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts. In *Proceedings of ACL*, 2004.
19. P.D. Turney. Thumbs up or thumbs down: Semantic orientation applied to unsupervised classification of reviews. In *Proceedings of ACL 2002*, pp.417-424, 2002
20. P.D, Turney, and M.L. Littman, Measuring praise and criticism: inference of semantic orientation from association, *ACM TOIS* 21(4), 2003.
21. C. Fellbaum, *WordNet: an Electronic Lexical Database*, MIT Press, 1998.
22. B. Liu, M. Hu, and J. Cheng. Opinion observer: Analyzing and comparing opinions on the web, In *Proceedings of International Conference on World Wide Web (WWW-2005)*, 2005.
23. B. Liu, *Opinion Mining*, In *Proceedings of the 17th International World Wide Web Conference*, Beijing, 2008.