

Marcarea sentimentelor în text: emoțiile autorului și opinia cititorului

Victoria BOBICEV, Liviu CARCEA

Technical University of Moldova

victoria.bobicev@ia.utm.md, liviu.carcea@ia.utm.md

Rezumat — Analiza automată a sentimentelor este o direcție de cercetare foarte activă în ultimii ani. Însă sentimentele sunt subiective și greu de analizat automat. Textele sunt marcate manual și apoi utilizate ca sursa pentru învățarea automată. Însă marcarea manuală la fel are multe nuanțe legate de subiectivitate. În articol se discută două tipuri de marcarea manuală și acordul între adnotatori în ambele cazuri. În final se compară rezultatele analizei automate în baza textelor marcate manual.

Cuvinte cheie — analiza sentimentelor, marcarea sentimentelor, interfața pentru marcarea, acord între adnotatori, analiza automată a sentimentelor.

I. INTRODUCERE

În ultimii decenii internetul a crescut considerabil. Volumul textelor postate online crește zilnic cât de la surse oficiale de informație atât și de la utilizatori. Informația postată de utilizatori în calitate de comentarii, discuții în bloguri și forumuri, evaluări ale produselor, filmelor sau muzicii a atras atenție sporită a cercetătorilor din sociologie, psihologie, marketing și altele.

Analiza automată a sentimentelor este o direcție de cercetare foarte activă în ultimii ani. Însă sentimentele sunt subiective și greu de analizat automat. Textele sunt marcate manual și apoi utilizate ca sursa pentru învățarea automată. Însă marcarea manuală la fel are multe nuanțe legate de subiectivitate.

În articol se discută două tipuri de marcarea manuală și acordul între adnotatori în ambele cazuri. În final se compară rezultatele analizei automate în baza textelor marcate manual.

II. LUCRĂRI PRECEDENTE

În multe cazuri este nevoie de marcarea manuală a sentimentului mesajelor online. Cu scopul de a obține o marcarea relevantă și cât mai obiectivă marcarea sentimentelor în texte se efectuează de câteva persoane independent una de alta și apoi marcările se compară. Ca regulă, marcarea coincide pentru unele texte iar pentru altele nu. Astfel, este nevoie de estimat acordul între adnotatori. Pe parcursul anilor au fost propuse mai multe metrici de acord între adnotatori pentru evaluarea acordului.

Procentajul de acord este cea mai ușoară și mai simplă măsură. Măsura este deseori criticată din cauza că nu face diferența între categorii, dar oferă totuși o aproximare de bază a acordului între adnotatori.

Însă pe parcurs au fost propuse alte măsuri de acord, mai avansate. Unele din măsurile acestea sunt Cohen's kappa (κ), Fleiss' kappa (K), Krippendorff's alpha (α) [6]. Acești coeficienți sunt numiți coeficienți ce măsoară valoarea acordului ce depășește acurd întâmplător în caz dacă adnotatorii aveau să marceze textele la întâmplare. Formula

generală pentru acești coeficienți este:

$$k = (A_o - A_e)/(1 - A_e) \quad (1)$$

unde A_o este acordul observat și A_e este acordul întâmplător în cazul în care adnotatorii aleg etichetele aleatoriu.

În cazul lui Cohen k , acordul așteptat A_e se calculează bazându-se pe faptul că alocarea aleatorie a categoriilor la elementele este guvernată de distribuții anterioare care sunt unice pentru fiecare adnotator și care sunt obținuți din distribuția lor reală.

Alpha (α) lui Krippendorff este un coeficient de acord pe baza ipotezei conform căreia acordul așteptat este calculat prin analizarea distribuției globale a marcărilor fără a lua în considerare adnotatorii care au produs aceste marcări. Această măsură poate fi utilizată în cazul marcărilor multiple și permite lipse de valori.

Adnotarea subiectivității poate fi centrată pe percepția cititorului [4] sau pe poziția autorului unui text [5]. În lucrarea curentă, ne-am propus să comparăm aceste două tipuri de marcări.

Fleiss k este generalizarea la mai mult de doi adnotatori; acordul așteptat se calculează pe baza presupunerii că alocarea aleatorie a categoriilor la articole, de orice adnotator, este guvernată de distribuția categoriilor articolelor în lumea reală.

Analiza sentimentelor de știri este diferită de cea a altor tipuri de text. [7] au subliniat că există mai multe probleme specifice legate de acest tip de texte; una dintre ele fiind problema separării știrilor bune sau rele și sentimentelor pozitive și negative. Annotatorii tind să interpreteze greșit intenția autorului și să marceze interpretarea lor ca fiind adevăratul sentiment al textului. În [5] a fost descris experimentul de adnotare a titlurilor de știri. Acordului între adnotatori chiar și pentru aceste fragmente scurte de text, a fost relativ scăzut, sub 50%. Numai la a doua etapă de re-adnotare a aceluși citate, au reușit să ajungă la un acord de 80%, creând un set de instrucțiuni de marcarea detaliate cu exemple și explicații multiple.

[3] au adaptat metodologia lui [5] în experimentul de adnotare și adnotările făcute de primii doi autori ai lucrării

au ajuns la un acord de 76%.

[1] a afirmat, de asemenea, că în transcrierile de știri sau în articolele de știri sentimentul atașat la o declarație poate fi mult mai puțin evident. Adnotarea a fost făcută folosind Amazon Mechanical Turk (<https://www.mturk.com/>) și fiecare unitate de adnotare a obținut trei adnotări independente. Lucrarea nu a raportat un acord de adnotare, deși au menționat că aproximativ 6% din exemple au fost aruncate din cauza dezacordului adnotatorilor.

În lucrările care descriu diferite tipuri de adnotări manuale [2], [5] au fost enumerate mai multe motive ale acordului slab între adnotatori. Ambiguitatea unităților de adnotare și subiectivitatea adnotatorilor au fost una dintre principalele probleme întâmpinate în acest proces. Un acord substanțial între adnotări poate fi atins numai după mai multe iterații de adnotări de către aceiași adnotatori și cu instrucțiuni de adnotare îmbunătățite iterativ.

În [4], fiecare titlu de știri din corpus a fost adnotat de șase annotatori care nu primiseră nicio instrucțiune suplimentară și, prin urmare, și-au adnotat propriile sentimente provocate de text. Apoi au calculat adnotările medii ale titlurilor.

În [2], masteranzii au adnotat mesaje pe forum de sănătate cu șase emoții de bază, iar calculul Fleiss Kappa [6] a fost relativ scăzut: 0,26. Ei au explicat acest dezacord prin variabilitatea dintre oameni și specificul textelor corporale. Cu toate acestea, au continuat cu experimente de clasificare automată a sentimentelor care dau rezultate rezonabile (cel mai bun rezultat $F = 0,65$).

III. MARCAREA SENTIMENTELOR ÎN FORUM

Primul experiment de marcare manuală a sentimentelor care descriem în acest articol este marcare sentimentelor în forumuri dedicate problemelor sănătății. Problemele personale legate de sănătate sunt unele care preocupă majoritatea populației. Analiza sistematică a 19 studii din 1999-2009 a evidențiat mai multe motive pentru utilizarea forumurilor medicale cum ar fi căutarea informațiilor pentru a afla despre aspectele psihologice, fizice și sociale ale tratamentelor disponibile, evaluări ale tratamentelor alternative sau în căutarea unui sprijin emoțional prin comunicare anonimă, acces imediat și constant la persoane cu aceleași probleme.

Am analizat sentimentele exprimate de participanții la forumul medical de fertilizare in vitro (ivf.ca). Acest forum reunește femeile care utilizează tratamente de FIV cu speranța de a concepe. Pentru analiza empirică, am selectat 1321 de posturi care au acoperit 80 de tematici legate de FIV (de exemplu, "Peste 40 de ani și însărcinate"). Pe când în majoritatea experimentelor cu sentimente se utilizează clasificarea în sentimente pozitive și negative, în cazul nostru polaritatea aceasta era nepotrivită și am avut nevoie de alte etichete de sentimente.

Analiza preventivă a conținutului discuțiilor pe forum a reafirmat faptul că majoritatea postărilor prezentau schimbul de experiențe personale, furnizarea de informații sau sfaturi, expresii de recunoștință / prietenie, chat, cereri de informații și expresii de susținere [8].

Astfel, începând cu mai multe sentimente posibile, am clasificat în final textele în *încurajare*, *gratitudine*,

confuzie, *fapte + încurajare* și *fapte*. Fiecare post în discuțiile selectate pentru marcare a fost marcat de două persoane. Nici una din persoanele acestea nu avea informație despre deciziile altor adnotatori. Textele în care adnotatorii nu au fost de acord cu o etichetă de clasă au fost considerați *ambigue*. Aceasta marcare este descrisă în [9] în mai multe detalii. Adnotatorii au lucrat cu textele date de două ori: prima etapă fiind analiza textelor din punct de vedere a tuturor sentimentelor posibile și a doua etapă atașarea unei etichete din cele stabilite la finalul primei etape. Respectiv, acordul obținut este comparativ bun ($Fleiss\ Kappa = 0.73$). Astfel de acord este considerat înalt [6]. Statistica rezultatelor marcării este prezentată în tabelul 1.

TABEL I. PROCENTAJUL TEXTELOR AMBIGUE ÎN MARCAREA SENTIMENTELOR ÎN POSTĂRILE UTILIZATORILOR

Categorie	Numarul de posturi	Procentajul
<i>Incurajare</i>	206	27%
<i>Gratitudine</i>	88	12%
<i>Confuzie</i>	48	6%
<i>Fapte + Incurajare</i>	73	10%
<i>Fapte</i>	187	25%
Ambigue	150	20%
Total	752	100%

În [10] este descrisă continuarea experimentului de marcare a sentimentelor pe același forum cu același set de sentimente. O interfață web a fost creată pentru adnotatori care oferea același set de sentimente: *încurajare*, *gratitudine*, *confuzie* și *fapte*.

Însă adnotatorilor a fost dată și posibilitatea de a adăuga orice altă emoție care ei observă în textele propuse spre analiză.

Ca rezultat, adnotatorii au adăugat o mulțime de sentimente ca: susținere, îngrijorare, incertitudine, compasiune, speranță, optimism, dispreț, îngrijorare, tristețe, bucurie, fericire, dezamăgire, frustrare.

În cazul forumurilor textele participanților conțin exprimările sentimentelor sale în abundență. Scopul autorilor textelor este de a împărtăși emoțiile și sentimentele sale. Astfel, începând marcarea am optat pentru modelul bazat pe dispoziția autorul textului și am cerut adnotatorilor să analizeze sentimentul exprimat în text de autorul lui. Acest model este firesc pentru marcare textului generat de utilizatori (User Generated Content). În majoritatea cazurilor utilizatorii nu ascund opiniile și sentimentele sale ci le exprimă liber și activ.

Textele marcate manual au fost utilizate în experimentele de clasificare automată. În primul pas au fost stabilite etichetele de sentimente prin alegerea sentimentului atașat de cel puțin doi adnotatori. Apoi au fost efectuate experimentele de clasificare automată pe baza etichetelor date. Rezultatul obținut a fost satisfăcător, F-measure în jur de 0,65.

IV. MARCAREA SENTIMENTELOR ÎN NOUTĂȚI

Absolut altă situație este în cazul noutăților postate online de ziare și agenții profesionale. Cum a fost deja menționat, în cazul acesta autorii textului îl prezintă în forma maximal neutră însă așteptând să trezească reacție

emoțională în cititor.

A doilea experiment a fost dedicat marcării sentimentelor în noutăți scrise în limba rusă și ucraineană [12]. Textele pentru procesare au fost extrase de pe două surse de știri: ucrainene (<https://tsn.ua/>) și rusești (<http://censor.net.ua/>). Inițial corpusul ucrainean conținea 5817 texte de știri, iar rus conținea 10194 texte de știri. O mare parte a textelor a fost neutră deoarece textele conțineau știri, nu comentariile utilizatorilor. Acestea au fost difuzate în forma comparativ neutră. Chiar dacă atitudinea autorului a fost prezentă într-un text, ea a fost exprimată implicit, fără cuvinte ce exprimă direct sentimente. Astfel, am selectat textele care conțineau cel puțin un cuvânt din lexiconul afectiv. Cu scopul acesta au fost utilizate lexicoanele descrise în [11]. După filtrarea au fost lăsate 2018 de texte ruse și 2133 ucrainene pentru marcarea manuală.

Adnotarea a fost realizată de studenții Institutului Politehnic Kharkiv prin interfața online. În cazul dat au fost propuse doar trei variante de etichete: pozitiv, negativ, neutru.

Aproximativ 40 de studenți au participat la adnotări ale fiecărei părți, ruși și ucraineni generând în 7248 adnotări pentru texte ruse și 6733 adnotări pentru limba ucraineană. Întrebarea la care au răspuns adnotatorii a fost: Ce sentimente evocă acest text? De asemenea, aceștia aveau posibilitatea de a adăuga comentarii dacă considerau necesar. În final, fiecare text a fost adnotat de 2 până la 5 elevi folosind trei etichete: pozitiv, negativ, neutru. Numărul mediu de adnotatori pe text a fost de 3,59 pentru limba rusă și de 3,2 pentru textele ucrainene. Acordul inter-adnotator calculat pe aceste texte a fost extrem de scăzut: Fleiss Kappa = 0,14 pentru textele ucrainene și Fleiss Kappa = 0,24 pentru rusă.

TABEL II. PROCENTAJUL TEXTELOR AMBIGUE ÎN MARCAREA SENTIMENTELOR ÎN NOUTĂȚI

Categorie	Procentajul de texte ambigue ucrainene	Procentajul de texte ambigue ruse
<i>Society</i>	14% (63 din 436)	16% (62 din 384)
<i>Politics</i>	21% (103 din 487)	13% (65 din 510)
<i>Incidents</i>	12% (58 din 486)	13% (30 din 239)
<i>Sport</i>	14% (58 din 426)	29% (92 din 413)
<i>Economics</i>	31% (68 din 218)	18% (43 din 245)
Total	19% (350 din 1849)	15% (292 din 1994)

În cazul textelor ucraineni, doi adnotatori au selectat o etichetă și un adnotator a selectat altă etichetă pentru 2/3 dintre texte. Chiar dacă în acest caz putem folosi eticheta selectată de doi adnotatori ca fiind cea corectă, distribuția voturilor este de fapt identică cu distribuția "din întâmplare", iar Fleiss Kappa pentru o astfel de adnotare este egală cu 0. În afară de aceasta, o mare parte din texte a fost marcată cu toate trei etichete de trei persoane. Aceste texte au fost absolut ambigue și nu era posibil de definit eticheta lor. Statisticile privind rezultatele marcării textelor de știri sunt prezentate în tabelul II.

Textele absolut ambigue au fost eliminate din seturile utilizate în experimentele de clasificare automată iar textele care au fost marcate cu aceeași etichetă de două persoane au fost păstrate în set cu eticheta dată.

Experimentele de clasificare automată în trei clase: pozitiv, negativ și neutru au demonstrat un rezultat comparativ bun: F-measure în jur de 0,73.

V. PROBLEMELE MARCĂRII

În această secțiune, analizăm și comparăm aceste două experimente și rezultatele lor.

Datele comparative pentru aceste două experimente sunt prezentate în tabelul III.

TABEL III. COMPARAREA EXPERIMENTELOR DE MARCARE A SENTIMENTELOR

Categorie	Primul experiment	Al doilea experiment
<i>Tipul textelor</i>	postările forumului despre problemele sănătății	textele noutăților de pe site-urile de știri
<i>tipul marcării</i>	centrată pe dispoziția autorului	centrată pe percepția cititorului
<i>întrebarea către adnotatori</i>	ce sentimente exprimă autorul textului	ce sentiment evocă în cititor textul dat
<i>înțelegerea între adnotatori</i>	Fleiss Kappa = 0.73 pentru doi adnotatori și Fleiss Kappa = 0.46 pentru trei	Fleiss Kappa = 0.14 pentru texte ucraineni și Fleiss Kappa = 0.24 pentru ruse
<i>rezultatele experimentelor de clasificare automată</i>	F-measure în jur de 0,65	F-measure în jur de 0,73

După cum observăm din tabelul III înțelegerea între adnotatori în cazul postărilor utilizatorilor este mai bună decât în cazul noutăților. Lucrul acesta poate fi explicat cu câteva cauze.

În cazul noutăților adnotatorii marcau sentimentele sale evocate de textul noutății și acestea puteau să fie absolute opuse. Cum a fost menționat și în alte cercetări, oamenii reacționează foarte diferit la aceeași informație. De exemplu, un anunț sportiv raportează despre un meci între două echipe de fotbal și rezultatul lui în care prima echipă a câștigat. Dacă adnotatorul este fanatul echipei date, acesta marchează textul ca pozitiv. În caz opus, dacă adnotatorul este fanatul echipei care a pierdut acesta marchează textul ca negativ și dacă adnotatorul nu este fanatul fotbalului și nici nu cunoaște echipele menționate, acesta marchează textul ca neutru; textul într-adevăr fiind doar constatația faptului fără orice discuție pe marginea lui. La fel se întâmplă și cu alte tipuri de noutăți; adnotatorii marcau sentimentele sale care erau diferite în dependență de pozițiile sale față de faptele descrise în text. De exemplu, în politică ca și în sport fiecare are părerea sa ce este corect și ce nu și reacționează destul de emoționant la noutățile politice.

În cazul marcării posturilor adnotatorilor scopul de bază a adnotatorilor a fost definit altfel. Lor li se cerea de indicat sentimentele autorilor textelor care erau exprimate destul de deschis, ca de exemplu: "Mă gândesc că, dacă sunt așa stresat, oricum, nu este un moment bun" sau "Ma bucur pentru tine! Arăți mai relaxat". Astfel, pentru adnotatori este mai simplu de indicat același sentiment și de ajuns la o înțelegere mai bună. Putem constata că adnotatorii

identificau sentimentele autorului în baza textului scris destul de uniform. Una din problemele întâmpinate este faptul că textele date sunt despre problemele sănătății și în cazuri când autorul textului descria problemele adnotatorii au dificultăți cu marcarea. Dintr-un punct de vedere acestea sunt descrieri cu aspect negativ dar din alt punct de vedere autorul textului poate reda informația dată în forma neutră, doar ca o descriere pentru alții însoțind cu informații despre tratament sau medicamente utilizate.

Însă, trebuie de spus că experimentele de detectare automată a etichetei au demonstrat că rezultatele mai bune au fost obținute pentru textele de noutăți decât pentru postările forumurilor. Aici la fel pot fi propuse câteva explicații posibile.

În primul rând, trebuie de reținut că în cazul forumurilor a fost efectuată clasificarea în patru clase, ce complică sarcina și duce la rezultate mai slabe decât în cazul clasificării doar în sentiment pozitiv și negativ.

A doua cauză poate fi mărimea textelor. În cazul noutăților textele au fost comparative mici (100 – 150 cuvinte) și, ce este esențial, conțineau informație doar pe un subiect. Postările forumurilor au fost destul de lungi (150 – 300 cuvinte), mai lungi decât sunt postările obișnuite ale utilizatorilor și în cazul postărilor mai lungi conțineau câteva părți cu tonalitatea diferită. Un post putea să se înceapă cu cuvinte de mulțumire interlocutorului precedent pentru susținere, apoi să conțină descrierea stării autorului care putea să aibă tonalitatea negativă și să termine cu enunțuri de speranță și urări de bine pentru alți interlocutori. Astfel, textul în sine deja putea să fie ambiguu și să creeze dificultăți în procesul clasificării.

Astfel, pentru fiecare sarcină de analiză a sentimentelor trebuie de luat în seamă mai multe aspecte: mărimea textelor, uniformitatea lor, tematica lor, numărul și caracterul sentimentelor. Un aspect important este scopul cercetării care determină ce anume trebuie să fie evidențiat în text analizat. Aici sunt posibile variante multiple: prezența sau absența subiectivității în text; sentimentul autorului textului; cauza sentimentului dat sau ținta lui; percepția și sentimentele cititorului și elementele textului care le evocă.

VI. CONCLUZIE

În lucrarea dată am comparat și am analizat două experimente de marcarea manuală a sentimentelor în textele online. În cadrul fiecărui experiment s-a realizat marcarea manuală a sentimentelor, calculul acordului între adnotatori, selectarea etichetelor pentru textele date și apoi experimentele de detectare automată a sentimentelor în baza marcării manuale.

Au fost evidențiate și discutate problemele și dificultățile marcării manuale și comparate rezultatele clasificării automate.

REFERINȚE

- [1] Joseph G. Ellis, Brendan Jou, and Shih Fu Chang. 2014. Why We Watch the News: A Dataset for Exploring Sentiment in Broadcast Video News. In the Proceedings of the 16th International Conference on Multimodal Interaction (ICMI 2014), pages 104-111.
- [2] Soumia Melzi, Amine Abdaoui, Jerome Azé, Sandra Bringay, Pascal Poncelet, and Florence Galtier. 2014. Patient's rationale: Patient Knowledge re-trieval from health forums. In the Proceedings of the Sixth International Conference on eHealth, Telemedicine, and Social Medicine, pages 140-145.
- [3] Patrik F. Bakken, Terje A. Bratlie, Cristina Marco and Jon Atle Gulla. 2016. Political News Sentiment Analysis for Under-resourced Languages. In Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, pages 2989–2996.
- [4] Strapparava, C. and R. Mihalcea. Semeval-2007 task 14: Affective text. Proceedings of the 2008 ACM symposium on Applied computing, 2008.
- [5] Balahur, A. and R. Steinberger. Rethinking Sentiment Analysis in the News: from Theory to Practice and back. Proceedings of the 1st Workshop on Opinion Mining and Sentiment Analysis, 2009.
- [6] Ron Artstein, Massimo Poesio. 2008 Inter-coder agreement for computational linguistics. Computational Linguistics Journal Volume 34 Issue 4, pages 555-596. doi: 10.1162/coli.07-034-R2
- [7] Alexandra Balahur, Ralf Steinberger, Mijail A. Kabadjov, Vanni Zavarella, Erik Van der Goot, Matina Halkia, Bruno Pouliquen, and Jenya Belyaeva. 2013. Sentiment Analysis in the News. In the Proceedings of the 7th International Conference on Language Resources and Evaluation (LREC'2010), pages 2216-2220.
- [8] Malik S. and N. Coulson. Coping with infertility online: an examination of self-help mechanisms in an online infertility support group. Patient Educ Couns, vol. 81, no. 2, pp. 315–318, 2010
- [9] Marina Sokolova and Victoria Bobicev. 2013. What Sentiments Can Be Found in Medical Forums? In Galia Angelova, Kalina Bontcheva, Ruslan Mitkov: Recent Advances in Natural Language Processing, RANLP 2013, Bulgaria.
- [10] Victoria Bobicev and Marina Sokolova. Inter-Annotator Agreement in Sentiment Analysis: Machine Learning Perspective. Recent Advances in Natural Language Processing, RANLP 2017, Bulgaria, 2017.
- [11] Olessia Koltsova, Svetlana Alexeeva, and Sergei Koltcov. 2016. An Opinion Word Lexicon and a Training Dataset for Russian Sentiment Analysis of Social Media. In the Proceedings of the International Conference "Dialogue 2016", pages 277-287.
- [12] Olga Kanishcheva and Victoria Bobicev. Good News vs. Bad News: What are they talking about? Recent Advances in Natural Language Processing, RANLP 2017, Bulgaria.