

Identificarea entităților, citatelor și evenimentelor în știri și texte din Web-ul social în limba română

Adrian-Nicolae Zamfirescu^{1,2}, Traian Eugen Rebedea^{1,2}

¹Universitatea Politehnica din București, Facultatea de Automatică și Calculatoare, Splaiul Independenței, Nr. 313, 060042, București, România

²TeamNet International, Splaiul Independenței, Nr. 319, 060044, București, România
E-mail: zamfirescuan@yahoo.com, traian.rebedea@cs.pub.ro

Rezumat. În cadrul prelucrării limbajului natural, detectarea automată a entităților cu nume a reprezentat una dintre cele mai importante provocări, care încă nu a fost rezolvată perfect până în acest moment nici pentru toate tipurile de texte scrise în limba engleză. Mai mult, identificarea entităților a deschis calea rezolvării unor alte probleme în care sunt implicate aceste construcții lingvistice, precum identificarea citatelor și a declarațiilor făcute în general de persoane, dar și de către companii sau alte tipuri de organizații, sau a extragerii evenimentelor. Problema identificării și a clasificării automate a entităților a apărut din necesitatea de a putea observa o evoluție a statisticilor în contextul diverselor informații din textele scrise, raportate la anumite persoane publice, organizații sau alte tipuri de entități cu nume care prezintă interes în diverse domenii. În articolul de față vom încerca rezolvarea problemelor menționate anterior pentru texte scrise în limba română și provenind din diverse surse online, precum știri, articole de pe bloguri sau comentarii din rețele sociale. Întâi, vom face o scurtă trecere prin fundamentele teoretice pentru rezolvarea acestor probleme, iar apoi vom prezenta metodele de rezolvare ale problemelor bazate pe aplicarea unor algoritmi de clasificare automată, combinați cu euristici bazate pe reguli și expresii regulate. În final, vom prezenta rezultatele și eficiența diverselor metode utilizate, o comparație între acestea, precum și concluziile referitoare la problemele abordate.

Cuvinte cheie: prelucrarea limbajului natural, extragerea informațiilor, clasificare automată, entități cu nume, citate, evenimente

1. Introducere

În multe aplicații de analiză automată a textelor, identificarea entităților cu nume reprezintă un prim pas esențial în cadrul procesului de analiză întrucât extrage o parte din elementele reprezentative din cadrul textului analizat (Jurafsky și Martin, 2008). Acest articol descrie realizarea unor componente pentru detectarea automată a următoarelor elemente din cadrul articolelor

de știri, a textelor publicate pe bloguri, extrase din rețele sociale sau a comentariilor publicate pe diverse situri web în limba română:

- **Entitățile cu nume** care apar în aceste texte, prin identificarea numelor acestora, dar și prin clasificarea acestora în câteva categorii predefinite.
- **Evenimentele** descrise de texte în care aceste entități sunt implicate, cu preponderență persoanele.
- **Citate** extrase din declarațiile persoanelor, atât exprimate ca citate, interviuri sau alte forme de vorbire directă în texte, cât și cele relatate sub formă de vorbire indirectă.

Problemele menționate anterior se încadrează în domeniul Prelucrării Limbajului Natural (PLN) și au fost rezolvate cu tehnici specifice PLN și extragerii de informații (*information extraction*) din texte, care vor fi detaliate în cadrul articolului. PLN este un domeniu de cercetare multidisciplinar, care combină elemente din: inteligență artificială, învățare automată, lingvistică, statistică, antropologie și altele. Rezolvarea problemelor principale studiate în cadrul PLN, care au drept obiectiv final înțelegerea limbajului natural, pot influența multe aspecte din interacțiunea om-calculator. În mod ideal, aceasta presupune înțelegerea textului de către calculator așa cum este perceput de către mintea umană (sau într-un mod asemănător celui în care îl înțelege aceasta), preprocesarea lui și translatarea sa într-un spațiu de reprezentare care să permită unor algoritmi și euristici specifici să extragă informația dorită pentru a fi folosită de către alte aplicații.

Problema identificării și a clasificării entităților cu nume are ca punct de pornire tocmai această punte de trecere între limbajul natural, inteligibil pentru om și o formă a sa modificată, propice pentru observarea anumitor șabloane care să intuiască sau să învețe tehnici de identificare sau clasificare în cadrul textelor. În cadrul lucrării de față, metodele de clasificare a entităților se axează pe distribuirea acestora în trei categorii: **persoane**, **organizații** și **teritorii**. Acestea sunt cele mai importante categorii folosite de către majoritatea sistemelor de clasificare a entităților cu nume (Finkel, Grenager și Manning, 2005; Bird, 2006), dar la care se pot adăuga și alte entități generale (precum data calendaristică, valori monetare etc.) sau specifice (de exemplu, nume de gene sau proteine, nume de proiecte de cercetare etc.).

Odată extrase și clasificate entitățile dintr-un text, putem să descoperim alte elemente utile care pot fi aflate despre entitatea respectivă. Cum o entitate cu nume apare într-un anumit context, acesta poate fi expresia unui fapt important consumat în timp și spațiu, adică un **eveniment**. Așadar, o entitate sau un grup de mai multe entități identificate pot conduce în continuare la detectarea unui eveniment în care acestea sunt implicate. În acest articol, ne vom baza preponderent pe evenimentele în care sunt implicate persoane, atunci când acestea sunt localizate prin cel puțin o referință temporală. Definiția aceasta este un pic diferită față de abordarea uzuală a detectării evenimentelor în sistemele de extragere a informațiilor, însă în contextul aplicației dezvoltate sunt importante, în special, evenimentele în care sunt implicate entitățile cu nume. Mai mult, și în sistemele uzuale pentru detectării evenimentelor se poate observa că majoritatea acestora sunt legate de cel puțin o entitate cu nume (Yang, Pierce și Carbonell, 1998).

Cum scopul final al aplicației din care fac parte modulele prezentate în cadrul acestei lucrări îl reprezintă monitorizarea evoluției unei entități cu nume în textele publicate online în România, probabil categoria cea mai importantă de entități este aceea a persoanelor, întrucât acestea sunt cele mai frecvente în cadrul textelor analizate. Cum unul dintre scopurile lucrării este de a identifica mențiunile persoanelor în cadrul articolelor de ziare online sau al postărilor pe bloguri sau situri de socializare precum Facebook sau Twitter, un element interesant în informațiile prezentate despre o entitate îl pot reprezenta **citatele**, precum și alte tipuri de declarații sau de vorbire indirectă prezentă în textele analizate și care, momentan, pot fi legate automat de către o entitate de tip persoană.

Astfel, un al treilea obiectiv al lucrării, pe lângă identificarea entităților și a evenimentelor în care sunt implicate acestea, îl reprezintă extragerea declarațiilor și citatelor. Ele se pot găsi la nivel textual atât în vorbirea directă – prin citate din declarațiile entităților respective sau din interviuri, cât și în cea indirectă – prin relatarea autorului de articol cu privire la ce a spus o anumită persoană publică.

Un aspect important care trebuie menționat este faptul că pentru rezolvarea acestor probleme au fost folosite anumite instrumente (biblioteci și date) disponibile liber (*open-source*, respectiv *open-data*). Pentru

modulele de învățare statistică pentru clasificarea entităților cu nume au fost folosite două biblioteci (care vin cu aplicații suport) create pentru învățare automată – Weka (Hall et al., 2009) sau PLN – Mallet (McCallum, 2002). Pentru elementele de analiză lexicală și morfologică a textelor în limba română s-a pornit de la baza de date cu toate cuvintele și formele flexionate existente în principalele dicționare publicate pentru română și care este pusă la dispoziție de către DEX Online (iar datele pot fi descărcate gratuit de la <http://www.dexonline.ro>).

În continuare, acest articol va continua cu o scurtă prezentare a cercetărilor anterioare în problematica abordată, prezentând rezultate obținute pentru limba engleză. În cadrul secțiunii 3, este descrisă arhitectura aplicației care rezolvă aceste probleme și sunt detaliate tehnologiile și algoritmi folosiți pentru rezolvarea acestora. Secțiunea 4 continuă cu evidențierea rezultatelor obținute pentru fiecare problemă în parte, fiind prezentată și o comparație între abordările utilizate pentru clasificarea entităților cu nume. Lucrarea se termină cu niște observații relevante despre problematica abordată și cu concluzii despre proiectul realizat.

2. Cercetări anterioare

O entitate cu nume (sau entitate numită) reprezintă o secvență de cuvinte care exprimă elemente din lumea reală și care de obicei pot fi organizate în categorii bine definite (Aggarwal și Zhai, 2012). Problema recunoașterii entităților cu nume (REN) este parte din domeniul care se numește extragerea informațiilor (*information extraction*), care folosește atât elemente de PLN, dar și din regăsirea informațiilor. În cadrul REN, se dorește localizarea în texte a unor expresii specifice și clasificarea acestora în categorii predefinite, precum nume de persoane, organizații, teritorii, expresii ale timpului, ale cantităților, valori monetare, procente ș.a. (Aggarwal și Zhai, 2012).

Deși studiul REN în texte a început dinainte de 1990, problema a fost conceptualizată formal în cadrul celei de-a șasea conferințe de înțelegere a mesajelor (*Message Understanding Conference, MUC-6*, <http://cs.nyu.edu/faculty/grishman/muc6.html>), din 1995, ca o subproblemă din cadrul domeniului extragerii informațiilor.

Primele soluții aduse pentru problema REN s-au bazat pe aplicarea unor șabloane, reguli sau automate finite, în general create manual (Grishman,

1995). Din cauză că această abordare presupunea expertiză umană pentru elaborarea șabloanelor, precum și din cauza faptului că șabloanele create manual nu puteau acoperi toate cazurile de entități prezente în corpusurile mai mari, sistemele ulterioare au încercat să învețe automat aceste șabloane din corpusuri adnotate, folosind diverse tipuri de reguli, transductoare sau automate finite (Mikheev, Moens și Grover, 1999). Totuși, cele mai recente studii și aplicații în domeniul REN se bazează pe metode statistice de învățare automată. Una dintre primele astfel de aplicații s-a numit Nymble și descoperă entitățile cu nume folosind modele Markov ascunse (HMM) (Bikel et al., 1997). Alte modele pentru clasificarea automată utilizate cu succes în cadrul REN sunt Entropie Maximă, SVM (*Support Vector Machine*) sau CRF (*Conditional Random Fields*) (Aggarwal și Zhai, 2012).

Interesul în detectarea automată a citatelor și a vorbirii indirecte din textele scrise a urmat ca o consecință logică a descoperirii entităților cu nume. Similar tehnicilor de recunoaștere a entităților, și această problemă se reduce la două abordări:

- Folosirea expresiilor regulate și a altor tipuri de modele pentru reprezentarea șabloanelor
- Învățarea automată folosind diverse modele statistice

Prima abordare este una mai simplă. Dezvoltatorul trebuie să definească un set mai mult sau mai puțin complex de expresii regulate, care, aplicate asupra textului procesat, să realizeze o potrivire (*matching*) prin care să se extragă structurile reprezentative (vorbitor, text citat sau declarat, data etc.). Ca și în cazul REN, există și aici alternativa identificării automate a acestor șabloane (Krestel, Bergler și Witte, 2008).

Alternativa la utilizarea expresiilor regulate este antrenarea unui clasificator pe un corpus robust, adnotat manual cu declarații și vorbire indirectă. În SUA, în timpul alegerilor din 2012 a fost dezvoltat proiectul *Politics Verbatim* pentru detectarea automată a declarațiilor politicianilor în texte, iar modelul Entropie Maximă a fost tehnologia aleasă pentru a rezolva această problemă (Davies, 2012). Acest model presupune definirea unor proprietăți (*features*) care să permită calculatorului definirea unor elemente cheie pentru a determina dacă o structură de text este citat sau nu.

Un alt domeniu de interes din extragerea informației din texte îl reprezintă detectarea evenimentelor. Ceea ce definește un eveniment

depinde în mod direct de contextul în care acesta se află și din care se dorește a fi extras. Declerck (2005) afirmă că existența unui eveniment presupune inițial existența unor entități și a unei relații între ele, evenimentul fiind declanșat de o schimbare de stare survenită în contextul curent. Pornind de la aceste idei, în cadrul lucrării am considerat că un eveniment relevant în știri sau texte din rețele sociale pentru o entitate trebuie să fie legat de către o dată când acesta a avut loc și, eventual, de către o locație. Aplicațiile practice realizate pentru recunoașterea evenimentelor sunt direcționate de regulă înspre un anumit domeniu specific. Totuși, majoritatea au fost dezvoltate pentru texte din rețele sociale, în special Twitter, și folosesc metode specifice regăsirii informației, precum extragerea termenilor cei mai relevanți și identificarea entităților în texte, împreună cu clustering sau învățare supervizată pentru a detecta și grupa evenimentele în sine (Sayyadi, Hurst și Maykov, 2009).

3. Descrierea implementării – algoritmi și euristici pentru rezolvarea problemelor

Sistemul dezvoltat pentru detecția și clasificarea entităților numite a fost conceput pentru o aplicație de monitorizare a publicațiilor online în limba română, care analizează texte provenite dintr-o gamă variată de surse. Spre deosebire de majoritatea sistemelor de REN dezvoltate și evaluate în studiile anterioare, în cadrul acestui lucrări am folosit o abordare mai pragmatică folosind o metodă semi-supervizată, în locul celei supervizate clasice.

Astfel, analiza pornește de la un corpus mare de texte care este folosit pentru a detecta o primă listă de entități cu nume potențiale (candidate). Dintre acestea, sunt păstrate inițial doar acele entități potențiale care pot fi găsite în versiunea în limba română a enciclopediei Wikipedia. În acest al doilea pas, fiecare entitate candidat este automat clasificată într-una din cele trei categorii de interes. Abia după aceea, putem construi un clasificator care va fi folosit pentru etichetarea automată a restului entităților care nu sunt regăsite în Wikipedia. În plus, după ce sunt cunoscute anumite entități, acestea sunt folosite pentru extragerea citatelor și a declarațiilor rostite de către aceste entități sau a evenimentelor în care sunt implicate. Întregul proces este reprezentat grafic în Figura 1.

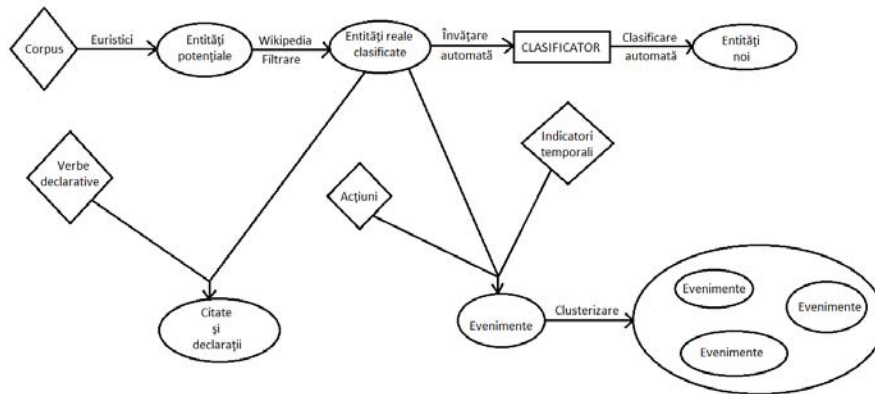


Figura 1. Diagrama funcțională a sistemului implementat pentru REN și extragerea citatelor și a evenimentelor

3.1 Detectarea și clasificarea entităților

În secțiunea curentă vom descrie pașii necesari pentru a identifica și clasifica entitățile cu nume extrase din textele apărute în articole online, în textele de bloguri sau în conținutul siturilor de socializare precum Facebook sau Twitter în limba română. În total, corpusul analizat cuprinde peste 200.000 de texte provenite din aceste surse, publicate în special în perioada iulie-decembrie 2012.

Așa cum am precizat anterior, pentru învățarea automată a șabloanelor sau a regulilor folosite pentru REN este necesar un corpus bogat, din care să se poată extrage suficiente informații. În domeniul prelucrării informației textuale, acest corpus este reprezentat de o colecție de texte similare cu cele care urmează a fi procesate pentru se extrage structurile dorite.

Lucrarea curentă urmărește detectarea și clasificarea entităților cu nume din articolele online de știri și clasificarea lor în 3 categorii: persoane, organizații și teritorii. După aceea, se construiesc liste cu aceste entități pentru a putea fi identificate și în alte tipuri de texte, precum cele preluate din rețele sociale.

Ideea premergătoare extragerii listei potențialelor entități a fost aceea a obținerii unui corpus suficient de robust astfel încât majoritatea entităților

cu nume care apar în acest corpus să aibă o frecvență de apariție suficient de mare, astfel încât identificarea acestora să poată fi făcută cu un nivel ridicat de încredere.

Colectarea corpusului a fost realizată cu ajutorul crawler-ului Apache Nutch (<http://nutch.apache.org/>), care a preluat articole de pe câteva sute de situri de știri și bloguri din România. După ce articolele au fost parcurse și analizate (parsate) și informațiile utile (precum titlul, data și textul articolului) au fost extrase, acestea au fost salvate și indexate folosind Apache Solr (<http://lucene.apache.org/solr/>) pentru a putea fi accesate rapid, inclusiv în cadrul procesului de determinare a entităților. Acest proces trece prin următoarele etape, descrise în continuare în cadrul acestei secțiuni:

1. Extragerea potențialelor entități folosind euristici
2. Filtrarea potențialelor entităților cu ajutorul Wikipedia
3. Clasificarea automată a entităților
4. Utilizarea clasificatorului pentru entitățile neconfirmate

Extragerea potențialelor entități folosind euristici

În cadrul primei etape în procesul de REN are loc detectarea și memorarea tuturor cuvintelor care încep cu majusculă. Folosind un corpus alcătuit din articole de pe situri de știri și bloguri, au fost de așteptat imperfecțiuni în redactarea lor. Deși majoritatea articolelor au și secțiuni de comentarii, iar identificarea se dorește atât la nivel de articol, cât și la nivel de utilizator care comentează articolul, extragerea inițială este realizată doar din corpul efectiv al articolului. Motivul principal al acestei alegeri este că siturile de știri și blogurile respectă, de obicei, normele de scriere a numelor proprii cu majuscule, spre deosebire de textele din comentarii sau din rețele sociale.

Așadar, pentru fiecare text de articol analizat, acesta este supus unui proces de *token*-izare, prin împărțirea sa în cuvinte și având ca delimitatori spațiile și semnele de punctuație dintre ele. Dacă un cuvânt începe cu literă mare, este reținut pentru a fi inclus în grupul potențialelor entități. Dacă următorul cuvânt care urmează după acesta începe tot cu literă mare și între cele două nu se află niciun semn de punctuație precum “.”, “,”, “?”, “!”, “:”, atunci al doilea se alătură primului și tot așa până când se încheie fraza sau a fost întâlnit un cuvânt scris cu literă mică care să închidă grupul curent

format. În acest fel se construiește o listă de *potențiale entități cu nume* din articolele analizate.

Astfel, fiecare grup de cuvinte consecutive, începând cu majuscule și despărțite exclusiv prin spații, formează un nou candidat pentru a fi o entitate cu nume. Există totuși o excepție în regula formării acestor grupuri: dacă între 2 cuvinte consecutive începând cu majusculă apar anumite cuvinte speciale (precum unele articole, prepoziții, conjuncții etc.), atunci se vor include și acestea în structura formată. Un exemplu elocvent pentru acest caz este entitatea cu nume *Cetatea de Scaun a Sucevei*.

Tot în cadrul acestei etape, este realizată prima filtrare care să ajute la restrângerea ariei de cuvinte care denotă în mod eronat o entitate cu nume: toate cuvintele începând cu majusculă, care apar la început de frază, sunt supuse unui verificări a frecvențelor de apariție în 3 contexte diferite:

- Număr de apariții ale acestor cuvinte la început de frază ($Na_{\hat{f}}$)
- Număr de apariții ale acestor cuvinte, scrise cu majusculă, în interiorul unei fraze (Na_{ifM})
- Număr de apariții ale acestor cuvinte, dar începând cu literă mică, în interiorul unei fraze (Na_{ifm})

Fiecare astfel de frecvență oferă o estimare asupra încrederii că un cuvânt sau grup de cuvinte reprezintă într-adevăr un nume de entitate. Spre exemplu, un nume de entitate nu ar trebui să apară, în mod normal, în mijlocul unei fraze, scris cu minusculă. Există totuși situații când anumite cuvinte crează ambiguitate datorită polisemiei. Un bun exemplu ar fi cuvântul *Marin*, care pe de o parte denotă un nume de persoană, iar pe de altă parte, scris cu literă mică, devine un adjectiv care sugerează o trăsătură specifică mării.

Prin urmare, sunt folosite următoarele euristici pentru o filtrare preliminară a listei de potențiale entități cu nume:

- $Na_{ifM} = 0$ și $Na_{ifm} > 0$
- $Na_{ifM} > 0$, dar $Na_{ifm} / Na_{ifM} > 10$
- $Na_{\hat{f}} + Na_{ifM} + Na_{ifm} < 0.05\% * N_{docs}$, unde N_{docs} reprezintă numărul total de documente din corpus.

Pe baza procesării rezultatelor numerice obținute se formează o primă tranșă de cuvinte detectate ca necorespunzătoare pentru a fi candidate în

continuare ca entități cu nume. Astfel sunt eliminate multe dintre apariții, care constituie în majoritatea cazurilor conjuncții, prepoziții sau substantive comune la început de frază.

Pentru restul de expresii potențiale se declanșează procesul de grupare a lor, având drept principale obiective realizarea unei clasament al celor mai des întâlnite entități în corpusul analizat și gruparea tuturor formelor lexicale folosite pentru a desemna aceeași entitate. În cadrul acestei etape este folosită baza de date DEX Online pentru a aduce la forma de bază cuvintele declinate (de ex. *Băncii Centrale a României* va ajunge *Banca Centrală a României*).

Filtrarea potențialelor entități cu ajutorul Wikipedia

Cu toate ajustările făcute, lista de entități, împreună cu formele lor lexicale și frecvențele de apariție, conține în continuare reziduuri, ca urmare a unei structuri imperfecte a textelor din articolele de știri.

Primul pas efectiv în detectarea unei entități reale din lista de entități potențiale (și totodată a clasificării ei) este făcut cu ajutorul variantei în limba română a enciclopediei Wikipedia. Datorită baze de date foarte mari pe care o pune la dispoziție gratuit, Wikipedia a fost folosită ca sursă pentru căutarea entităților cu nume identificate. De altfel, una dintre noile metode din domeniul REN de identificare a entităților este legarea (*cross-linking*) la Wikipedia, proces denumit și *wikification*. În acest caz, o entitate nu este neapărat descoperită, cât mai degrabă asociată cu o referință existentă deja.

În cazul de față, entitățile căutate pe Wikipedia și care sunt regăsite în conținutul acesteia devin confirmate ca entități reale. Mai mult decât atât, arhitectura Wikipedia face posibilă și clasificarea lor, încă din această etapă, facilitând totodată formarea corpusului adnotat ce va fi folosit în pasul următor, pentru învățarea supervizată.

Clasificarea presupune distribuirea entităților într-una din cele trei categorii: persoane, organizații sau teritorii. Pentru a clasifica entitățile care au pagini asociate pe Wikipedia, sunt folosite două metode:

- În primă fază s-a încercat găsirea casetei de informații tipice pentru fiecare tip de entitate în parte, dar care există doar pentru unele pagini. Astfel, pentru fiecare din cele trei clase, au fost asociați vectori de cuvinte-cheie (de ex. data nașterii pentru persoane) care să indice categoria. În cazul unei potriviri de șablon, entității îi este asociată clasa respectivă.

- În lipsa casetei de informații, este analizată prima frază din articol. Cum Wikipedia are un format specific al primului paragraf care descrie succint entitatea, și în acest caz este folosită o euristică ce folosește șabloane și expresii regulate. Spre exemplu, în cazul persoanelor, multe dintre descrieri încep cu secvența de caractere “(n. “, imediat după numele persoanei. Alt caz apare când după numele entității urmează un predicat nominal, al cărui nume predicativ este un cuvânt care indică meseria persoanei; în acest caz, cuvântul este căutat pentru potrivire în lista completă de meserii din România, extrasă din alte surse online.

Clasificarea automată a entităților

În urma “wikificării”, au fost găsite referințe pentru majoritatea entităților candidat – în special, pentru cele foarte frecvente, care au fost imediat și clasificate. Rămâne totuși un număr considerabil de potențiale entități neconfirmate. În acest moment, datorită adnotării cu clase a unui număr semnificativ de entități, este posibilă crearea unui corpus de antrenare pentru un clasificator automat. Acesta poate fi folosit pentru a identifica și clasifica restul entităților candidat, precum și a entităților noi.

Dintr-o perspectivă practică, un aspect important din spatele algoritmilor de învățare supervizată care vor fi testați este identificarea atributelor (*feature-uri*) semnificative, care descriu documentele în punctele esențiale pentru a putea permite procesul de inferență și generalizare. Practic, combinația de atribute cea mai potrivită pentru clasificare se determină prin aplicarea unei funcții-indicator asupra setului de documente – funcția de învățare. Așa cum a fost menționat în pasul precedent, clasificarea primelor entități identificate a fost făcută cu ajutorul articolelor corespunzătoare paginilor de pe Wikipedia. Așadar, aceste texte pot constitui și corpusul de antrenare utilizat pentru a extrage atributele necesare.

Ne propunem deci să determinăm acele atribute care să fie relevante pentru categoria din care face parte fiecare entitate în parte. La prima privire, întregul articol referitor la o entitate ar trebui să fie util pentru discriminare între clase, întrucât se concentrează pe entitatea respectivă. Deci o primă abordare presupune folosirea întregului text al articolului pentru definirea setului de antrenare.

Din punctul de vedere al acestei abordări, atributele care vor fi puse la dispoziție algoritmului de clasificare sunt chiar cuvintele din fiecare articol, eliminând cuvintele de stop și aplicând lemantizare. Această abordare ar trebui să fie utilă deoarece fiecărei entități îi sunt asociate o multitudine de cuvinte care vor dobândi anumite ponderi în raport cu clasa din care face parte entitatea. Totuși, în cadrul etapei de testare prin crosvalidare în 10 runde a clasificatorului rezultat, după ce acesta a fost antrenat cu 800 de entități distribuite aproximativ egal în cele trei categorii, s-a constatat obținerea unei acurateți de doar 55% folosind modelul Entropie Maximă din cadrul Mallet.

După analiza acestor rezultate, am ajuns la concluzia că deși fiecare categorie a avut asociate multe atribute provenite din cuvintele extrase din articolele de pe Wikipedia, acest lucru a avut totuși un dezavantaj major prin faptul că multe cuvinte s-au dovedit a fi irelevante pentru clasele analizate. De fapt, studiind mai atent un astfel de articol, se observă că multe dintre fraze se referă la lucruri, fapte, persoane aferente entității în cauză, iar acest lucru crează confuzie în cadrul procesului de învățare.

Reconsiderând modul de alegere al setului de antrenare, informațiile păstrate pentru fiecare entitate sunt reduse pentru a fi mai relevante pentru entitatea respectivă (și, implicit, pentru clasa din care aceasta face parte). Cum aproape întotdeauna prima frază dintr-un articol Wikipedia se referă la entitatea respectivă, este normal ca și cuvintele constituente ale frazei să fie relevante și descriptive pentru acea entitate.

Ca o abordare intermediară între cele două strategii, am decis să considerăm doar primul paragraf din fiecare articol. Antrenând din nou corpusul, observăm o creștere a acurateții până la 73%. Cu toate acestea, rezultatele obținute nu sunt acceptabile pentru a putea fi folosite în practică, deci trebuie căutate alte îmbunătățiri.

În mod ideal, ar trebui să extragem acei termeni care să se raporteze sau să determine exact persoana, instituția sau locația referită de către fiecare articol. La nivel de text, acest lucru poate fi făcut analizând morfo-sintactic propozițiile constituente ale articolelor. Multe cuvinte nu descriu în mod direct entitatea. Astfel, această abordare se axează pe extragerea strict a atributelor care se leagă de entitate. Din fiecare articol asociat unei entități sunt identificate adjectivele din vecinătatea entității, verbele predicative ce reflectă o acțiune săvârșită de către entitate și numele predicative din predicatul nominale asociate entității.

În acest fel, numărul de atribute se reduce considerabil, însă de această dată devine mult mai specific, mai orientat spre descrierea entității. Folosind aceste noi atribute se obține o îmbunătățire substanțială, obținând o acuratețe de 85% folosind același clasificator Entropie Maximă implementat în Mallet. Această abordare a fost testată cu mai mulți algoritmi de învățare supervizată: SMO, Bayes Naiv și Entropie Maximă (Mohri, Rostamizadeh și Talwalkar, 2012) – primii doi implementați în Weka, ultimul în Mallet, iar acuratețea medie pentru cele trei clase variază în intervalul 80-90%. Rezultatele comparației sunt prezentate în secțiunea următoare.

Utilizarea clasificatorului pentru entitățile neconfirmate

În urma etapei de antrenare a clasificatorului, au fost testate mai multe abordări de alegere a atributelor. Mai departe, clasificatorul rulat pe corpul de antrenare poate fi folosit pentru a clasifica entitățile care nu au fost găsite pe Wikipedia sau pe cele găsite, dar neclasificate din cauză că nu s-au încadrat în șabloanele definite anterior.

Pentru această etapă, este important să poată fi colectat un corpus suficient de mare de fraze ce includ entitățile din care vor fi extrase atributele. O problemă în această abordare este lipsa unui corpus consistent de texte pentru entitățile care apar foarte rar. Tocmai de aceea, pentru entitățile care apar rar, este stringentă obținerea tuturor caracteristicilor posibile existente în text, raportate la respectivele entități. Acest lucru presupune metode avansate de analiză sintactică a textelor (pentru extragerea co-referințelor, arbori de parsare etc.). Dacă în etapa de antrenare, obținerea atributelor se realiza extrăgând din text cuvinte apropiate, care sunt evident legate de către entitate (un adjectiv alăturat entității, un substantiv care apare ca nume predicativ al unui predicat nominal reprezentând o acțiune săvârșită de entitate etc.), de această dată este necesară o analiză aprofundată a întregului context în care apar menționate entitățile ce se doresc clasificate.

Pentru acest lucru, se folosesc parsere sintactice special implementate și dedicate procesului de stabilire a relațiilor dintre cuvintele unui text. Analizând relațiile ce includ entitatea observată, se extrag cuvintele care fac referire la aceasta și se adaugă ca atribute pentru setul de test, în scopul

clasificării. Această etapă este în curs de implementare și nu vor fi raportate rezultate obținute în cadrul acestui articol.

3.2 Extragerea citatelor și a declarațiilor

Declarațiile unei persoane reprezintă o sursă permanentă de analiză și de generare a noi subiecte de discuție. În acest moment, în cadrul aplicației dezvoltate, acestea sunt detectate folosind șabloane și expresii regulate. Pe de o parte, această soluție este foarte rapidă și conduce la un procent de identificare ridicat. Pe de altă parte, în acest moment nu există un corpus adnotat manual de declarații ca să poată fi aplicată clasificarea automată. Rezultatele satisfăcătoare apar în urma observării folosirii unor șabloane specifice exprimării jurnalistice. În continuare, vor fi prezentate câteva tipuri de exprimări ce anticipează (marchează) declarația unei persoane:

- “Doamnă Ministru, așa cum ați declarat...”, a declarat Victor Ponta.
- Boc a declarat că pentru el este important să termine cursa.
- Preda a mai spus că nimeni nu trebuie să țină de scaun...
- “Voi depune plângere penală împotriva lui Ponta...” a spus Dan Diaconescu.
- Ioan Oltean a declarat vineri, într-o conferință de presă, că România...

După cum se observă și din exemplele anterioare, o declarație poate însemna, la nivel de text, un extras direct din cuvintele persoanei (un citat) sau poate fi expresia vorbirii indirecte.

Indiferent de tipul declarației, un aspect esențial pe baza căruia se începe căutarea este identificarea unor cuvinte-cheie, care să indice existența unei astfel de declarații. În urma analizei modului în care sunt scrise articolele și a limbajului folosit de către redactori, se poate observa că aceste cuvinte-cheie sunt reprezentate de anumite verbe care indică faptul că urmează o afirmație sau o declarație a unei terțe persoane. Acestea se numesc verbe de zicere (Bălășoiu, 2004), iar printre acestea se află: *a spune, a declara, a zice, a afirma, a preciza, a anunța* ș.a.

Așadar, ne vom folosi de existența acestor verbe declarative pentru a localiza paragraful din care se poate extrage un citat/o declarație. Abuzând ușor de terminologie, numim un citat spusele unei persoane afișate ca atare în articol (de obicei încadrate în ghilimele) – vorbirea directă, iar o

declarație va fi exprimată prin relatarea la persoana a treia de către autor a afirmației – vorbirea indirectă.

În cazul citatelor, se observă că există, de obicei, simboluri specifice care delimitează vorbirea directă (ghilimele, un paragraf etc.), precum și faptul că entitatea cu nume apare după verbul declarativ în textele de știre. Identificând caracterele delimitatoare ce se află în apropierea verbului declarativ, citatul poate fi extras cu ușurință.

Și al doilea caz, cel al declarațiilor sau al vorbirii indirecte, este caracterizat de anumite cuvinte-cheie. O declarație este precedată într-o majoritate covârșitoare de cazuri de construcții specifice (de exemplu, conjuncția “că”) care urmează verbului declarativ. De cele mai multe ori, aceste construcții specifice sunt situate la nivelul textului imediat după verbul care anunță declarația.

Există însă și cazuri particulare care trebuie tratate. Spre exemplu, între conjuncția “că” și verbul declarativ poate apărea o sintagmă care să ofere un plus de informații asupra declarației. Pentru a stabili validitatea unor astfel de situații, grupurile necunoscute de cuvinte sunt tratate morfo-sintactic, pentru a nu depista verbe care să anuleze detectarea citatului sau a declarației.

3.3 Detectarea evenimentelor

Un alt aspect care poate fi exploatat în domeniul prelucrării textului atunci când sunt implicate entitățile cu nume este extragerea evenimentelor la care acestea au participat. Domeniul semantic al percepției unui eveniment este unul vast și poate fi mult detaliat și extrapolat. În cadrul abordării curente, am considerat că un eveniment implică săvârșirea unei acțiuni. La nivelul analizei textuale, o acțiune este exprimată efectiv printr-un predicat, adică un cuvânt sau un grup de cuvinte care să reprezinte un verb sau un grup nominal.

Până a ajunge la verb, trebuie precizată condiția primordială pusă pentru a considera existența unui eveniment, și anume plasarea sa în timp. Cu alte cuvinte, în textul analizat trebuie localizate coordonate temporale care să precizeze într-o anumită măsură momentul la care evenimentul are (sau a avut) loc. Această idee este folosită în special pentru a colecta informații cât

mai specifice despre eveniment. Tocmai de aceea, precizarea timpului în care se declanșează o anumită acțiune crește probabilitatea ca aceasta să fie nucleul unui eveniment.

Problema localizării expresiilor temporale în texte ar putea fi din nou un subiect de dezbătut cu privire la ce fel de metode pot fi folosite pentru detectarea lor. Cum acestea sunt un tip specific de entități, se pot aplica aceleași metode ca cele folosite pentru REN. Însă datorită specificului expresiilor temporale, vom apela la expresiile regulate pentru identificarea lor, pentru că în acest context sunt rapide și foarte eficiente.

Apariția verbului împreună cu marcarea timpului nu sunt însă suficiente. Pentru a putea face o conexiune cu entitățile implicate în eveniment, acestea trebuie să fie inițial detectate în fraza respectivă, iar predicatul exprimat prin verbul găsit să refere o astfel de entitate în arborele de parsare.

După ce avem o listă de astfel de evenimente extrase din texte, este util să descoperim evenimentele care formează un centru de interes mai mare în opinia publică. Deci, ne dorim să grupăm aceste evenimente în funcție de anumite criterii pentru a putea observa în ce zone se concentrează atenția oamenilor din mass-media cu privire la persoane publice, locuri, instituții.

Pentru aceasta am apelat la folosirea algoritmilor de clustering pentru texte (Manning și Schütze, 1999). Clustering-ul presupune gruparea unui set de obiecte în așa fel încât obiectele din același grup au un grad de similaritate mai mare între ele decât au cu cele ce formează alte grupuri.

Ceea ce rămâne de definit este similaritatea între două documente. În cazul textelor, abordarea uzuală implică extragerea cuvintelor și calcularea ponderilor TF-IDF (*term frequency-inverse document frequency*) ale acestora pentru întregul set de documente. Aceste ponderi vor fi folosite pentru reprezentarea în spațiul vectorial. Mai exact, fiecare eveniment (frază) din spațiul limbajului natural va avea asociat un vector de astfel de ponderi, care va fi reprezentarea sa în spațiul vectorial asociat.

Spațiul TF-IDF produce o statistică numerică ce reflectă importanța unui cuvânt într-un document care face parte dintr-un corpus. Valoarea TF-IDF crește proporțional cu numărul de apariții ale unui cuvânt într-un document, dar este invers proporțională cu frecvența cuvântului în tot corpusul, lucru care ajută la penalizarea cuvintelor care sunt mai comune.

Trecând paragrafele din spațiul limbajului natural în cel vectorial, în continuare aceste numere trebuie folosite pentru a determina gradul de similaritate dintre evenimente. Măsura folosită pentru acest lucru este

similaritatea cosinus (Manning și Schutze, 1999). Aceasta presupune determinarea cosinusului unghiului dintre cei doi vectori reprezentând două documente diferite. Valoarea obținută este pozitivă, subunitară și are următoarele semnificații:

- O valoare egală cu 0 implică faptul că cele 2 documente sunt complet diferite (nu au nici un cuvânt în comun).
- O valoare egală cu 1 semnifică faptul că documentele sunt identice (conțin aceleași cuvinte, chiar dacă de un număr diferit de ori).

Prima observație care s-a putut face este faptul că această abordare inițială, de a utiliza fiecare cuvânt din frază, nu oferă scoruri satisfăcătoare.

Având în vedere faptul că un eveniment are la bază un fapt săvârșit, o acțiune, ne dorim ca această acțiune să fie un punct de echivalență între două evenimente. De asemenea, atunci când am determinat un eveniment, la baza căutării sale a stat existența cel puțin a unei entități în fraza ce descrie evenimentul. Deci un alt coeficient comun este prezența aceluiași entități în cele două evenimente comparate. Pe baza acestei logici, se aplică înainte o “curățare” a corpusului, prin filtrarea cuvintelor irelevante. Practic, se vor păstra doar cuvintele începând cu majusculă care descriu componente de entități cu nume și cele care sunt verbe (descriind practic acțiunile) în frazele care sunt extrase pentru fiecare eveniment în parte.

4. Rezultate

Rezultatele obținute pentru fiecare dintre cele trei probleme tratate în cadrul acestui articol vor fi evidențiate prin comparații, statistici sau exemple în cadrul acestei secțiuni.

4.1 Detectarea și clasificarea entităților

După cum am menționat și în secțiunea anterioară, am folosit trei metode diferite de învățare supervizată pentru clasificarea entităților cu nume, pentru a observa diferențele de performanță. Cele trei metode sunt reprezentate de către: SMO (*Sequential Minimal Optimization*, o variantă de SVM) și Bayes naiv, implementate în cadrul Weka, precum și modelul Entropie Maximă din Mallet. Valorile acurateței obținute în urma validării

clasificatoarelor automate prin crosvalidare folosind acești algoritmi sunt prezentate în Tabelul 1.

Tabelul 1. Comparație între metodele de învățare supervizată pentru clasificarea entităților cu nume

Algoritmi Categorii	SMO	Bayes Naiv	Entropie Maximă
Persoane	0.789	0.888	0.784
Organizații	0.804	0.750	0.888
Teritorii	0.939	0.777	0.965

După cum se observă, rezultatele sunt similare în cazul algoritmilor SMO și Entropie Maximă. Acuratețea ușor mai mică apare în cazul algoritmului Naive Bayes, generând o valoare ușor mai scăzută decât a celorlalți doi, însă are o viteză de antrenare și de clasificare mai mare.

Crt.	Persoane	Apariții	Organizații	Apariții	Teritorii	Apariții
1	Traian Băsescu	11870	Oltchim	4274	România	17307
2	Victor Ponta	3623	Apple	637	București	5531
3	Dan Diaconescu	2465	Google	398	Londra	1805
4	Crin Antonescu	1799	Samsung	215	Germania	1216
5	Adrian Năstase	1117	Renault	129	Franța	1148
6	Emil Boc	822	Academia Română	100	Cluj	1122
7	Vasile Blaga	537	Microsoft	97	Uniunea Europeană	1115
8	Mona Pivniceru	533	Petrom	93	Statele Unite	1044
9	Barack Obama	531	Rompetrol	90	Marea Britanie	1013
10	Ioan Rus	486	Nokia	84	Bruxelles	962

Figura 2. Topul celor mai frecvente entități cu nume într-un corpus format din 25.000 de știri

În toate situațiile de mai sus, procentele au fost stabilite printr-un proces de validare cunoscut sub numele de *10-fold cross-validation* – semnificând faptul că se împarte setul de documente în 10% set de testare, prin care se verifică corectitudinea clasificatorului format, și 90% set de antrenare.

Pentru a reflecta identificarea entităților extrase într-un corpus format din 25.000 de știri colectate în perioada iulie-decembrie 2012, în Figura 2 este reprezentat topul frecvenței entităților cu nume. Sunt afișate cele mai frecvente 10 entități, grupate în cele trei categorii, împreună cu numărul lor de apariții și sortate descrescător după apariții.

4.2 Extragerea citatelor și a declarațiilor

Citatele și declarațiile au fost extrase dintr-un set de aproximativ 50.000 de documente text, în special de pe siturile de știri și bloguri, din aceeași perioadă analizată. După cum se observă în Figura 3, predomină în top persoanele politice sau cele cu funcții publice. Deci în opinia publică subiectul politic predomină, precum și declarațiile făcute de persoanele din acest peisaj.

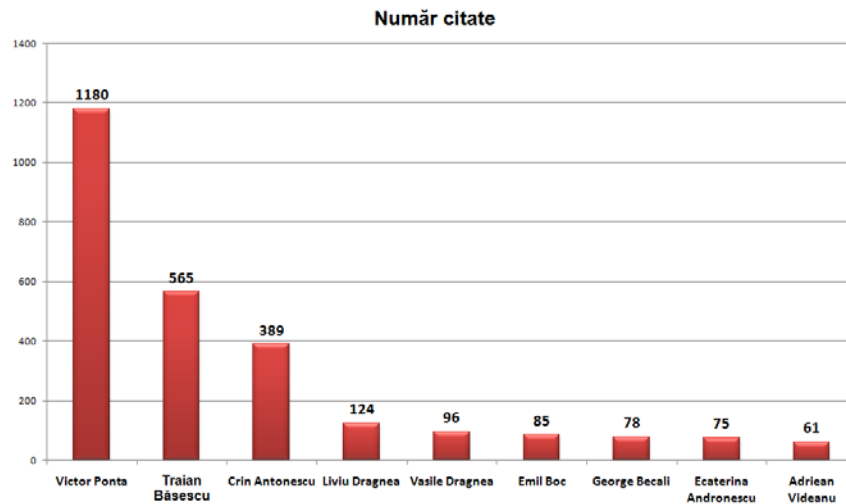


Figura 3. Entitățile cu cele mai multe declarații și citate extrase din 50.000 de articole de pe siturile de știri și de pe bloguri

4.3 Detectarea evenimentelor

Un eveniment în sine este reprezentat textual printr-o frază. Având în vedere că dintr-un set de 5000 de documente au fost detectate 2235 de evenimente,

pentru a putea oferi statistici relevante care să se limiteze în spațiul articolului, vom afișa în continuare cele mai relevante clustere de câte două evenimente, împreună cu scorul similarității aferent, extrase dintr-un subset de 100 de documente:

- “Duckadam a câștigat cu Steaua Cupa Campionilor Europeni în 1986, la Sevilla, când a apărut 4penalty-uri, intrând în Cartea Recordurilor.”
- “Helmut Duckadam are această suferință încă din 1986, la doar două luni după ce a câștigat cu Steaua Cupa Campionilor Europeni la Sevilla, intrând în Cartea Recordurilor pentru că a apărut patru penalty-uri.”
→ Scor: **0.7**
- “TVR Info și-a suspendat emisia pe 15 august, la ora 23.”
- “TVR Cultural și-a suspendat emisia pe 15 septembrie, de la ora 23.”
→ Scor: **0.63**
- În rechizitoriul întocmit de procurorii DNA se arată că, în perioada mai-septembrie 2009, Dan Diaconescu l-ar fi amenințat în mod repetat, atât în mod direct, în cadrul emisiunii Dan Diaconescu Direct din 21 iulie 2009, difuzată de postul de televiziune OTV, cât și indirect, prin intermediul lui Doru Pârv, pe primarul unei comune din județul Arad, pentru a-l determina să le dea suma totală de 200.000 de euro.
- De asemenea, în cursul lunii aprilie 2005, Dan Diaconescu l-ar fi amenințat de mai multe ori, atât în mod direct, în cadrul emisiunii Dan Diaconescu Direct din seara zilei de 20 aprilie 2005, cât și indirect, prin intermediul lui Ghezea Mitruș, realizatorul emisiunii Semnal de alarmă difuzată pe același post de televiziune, pe omul de afaceri Paul Petru Țârdea, pentru a-l determina să-i dea suma totală de 100.000 euro.
→ Scor: **0.4**

4. Descrierea aplicației de monitorizare

Sistemul de recunoaștere și clasificare a entităților numite descris în această lucrare este integrat și folosit în cadrul unei aplicații de monitorizare a publicațiilor online în limba română dezvoltat de către compania TeamNet

International. În acest moment, sistemul indexează aproape un milion de elemente text diverse incluzând știri, articole de pe bloguri, comentarii și statusuri din rețele sociale, precum și texte din forumuri de discuții. Sistemul de etichetare cu entități numite este folosit pentru a descoperi aparițiile, citatele și evenimentele în care sunt implicate entitățile menționate în textele monitorizate.

În prezent, sunt monitorizate peste 10.000 de entități cu nume diferite, însă sperăm ca această listă să fie extinsă la o iterație ulterioară a sistemului descris în această lucrare, folosind un corpus mai mare pentru etapa de învățare semi-supervizată. Mai mult, pentru fiecare entitate cu nume, aplicația reține mai multe forme lexicale (de exemplu, “CCR”, “Curtea Constituțională a României”, “CC a României”), iar acestea sunt folosite pentru REN în texte. După cum am menționat și în capitolele anterioare, fiecare formă lexicală este, de asemenea, identificată dacă apare cu sau fără diacritice, în forma normală sau flexionată.

Aplicația permite și vizualizarea interactivă a informațiilor despre entități, putând compara numărul de apariții în textele monitorizate ale mai multor entități numite diferite, precum și numărul de citatelor și ale evenimentelor în care acestea sunt implicate. Această aplicație este încă în curs de dezvoltare, iar la finalizarea sa rezultatele vor fi publicate într-un articol ulterior.

5. Concluzii

Prelucrarea limbajului natural constituie una dintre cele mai complexe și interesante ramuri ale inteligenței artificiale, întrucât presupune un proces de tranziție între perceperea limbajului scris de către creierul uman și învățarea conotațiilor acestui limbaj de către calculator. Problema identificării și a clasificării entităților a reprezentat de la început o provocare pentru cercetătorii care au lansat ideea detectării automate a acestora în texte.

Această provocare a condus mai departe la noi centre de interes. Extragerea citatelor din text, detectarea și extragerea evenimentelor reprezintă o parte din activitățile în continuă dezvoltare din universul PLN și al extragerii informațiilor din texte.

Articolul de față a prezentat o abordare a problemei recunoașterii și a clasificării entităților, a detectării automate a citatelor și a extragerii evenimentelor din text. Combinând atât modalități clasice de folosire a regulilor și a expresiilor regulate, cât și tehnici moderne de învățare automată și statistică, s-a putut observa pe parcurs cum evoluția rezultatelor a depins de alegerea inteligentă și contextual potrivită a euristiciilor și a metodelor de clasificare folosite. În plus, probabil acesta este cel mai complex studiu pentru REN în limba română pentru articole de știri, din bloguri sau texte din rețelele sociale. Deși rezultatele obținute sunt rezonabile pentru a fi integrate în aplicații practice (85% pentru clasificarea entităților cu nume), considerăm că o rafinare a lor mai poate fi făcută în viitorul apropiat pentru a obține performanțe și mai bune.

Studiile anterioare ale REN în limba română au atins rezultate mai bune, însă pentru corpusuri de validare de dimensiuni mult mai reduse și pentru categorii de texte mai puțin variate, în special pentru texte jurnalistice. De exemplu, Pastra et al. (2002) raportează o precizie de 88% pentru persoane, 92% pentru locații și 95% pentru organizații, însă considerăm că este folosit un corpus destul de mic și insuficient variat: “1MB de texte” din ziarul “Amprenta” (Pastra et al.). Pe de altă parte, un alt studiu efectuat pentru REN în limba română nu specifică detalii despre rezultatele obținute sau corpusul folosit pentru antrenare și validare (Tufis et al. 2008). Din acest motiv, considerăm că rezultatele din această lucrare aduc o valoare față de toate studiile anterioare, în special datorită volumului mare de texte folosite, precum și abordării semi-supervizate pentru rezolvarea problemei recunoașterii și clasificării entităților numite.

Referințe

- Aggarwal, C.C., Zhai, C. *Mining Text Data*. Springer, 2012.
- Agrawal, R., Mannila, H., Srikant, R., Toivonen, H. Fast Discovery of Association Rules. *Advances in knowledge discovery and data mining*, American Association for Artificial Intelligence, Menlo Park, CA, 1996.
- Bălășoiu, C. *Discursul raportat în textele dialectale românești*, Ed. Univ. București, București, 2004.
- Bikel, D.M., Miller, S., Schwartz, R., Weischedel, R. Nymble: a high-performance learning name-finder. *Proceedings of the fifth conference on Applied natural language processing (ANLC '97)*, Association for Computational Linguistics, Stroudsburg, USA, pp. 194-201, 1997.

- Bird, S. NLTK: the Natural Language Toolkit. *Proceedings of the COLING/ACL on Interactive presentation sessions (COLING-ACL '06)*. Association for Computational Linguistics, Stroudsburg, PA, USA, pp. 69-72, 2006.
- Brodely, C.E., Friedl, M.A. Identifying and Eliminating Mislabeled Training Instances. *Proceedings of the thirteenth national conference on Artificial intelligence - Volume 1 (AAAI'96)*, Vol. 1. AAAI Press, pp. 799-805, 1996.
- Davis, C. *Using machine learning to extract quotes from text*. CIR Labs, available online at <http://cironline.org/blog/post/using-machine-learning-extract-quotes-text-3687>, 2012.
- Declerck, T. *Automatic event extraction from text on the base of linguistics and semantic annotation*. Language Technology Lab, 2005.
- Finkel, J.R., Grenager, T., Manning, C. Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling. *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL 2005)*, pp. 363-370, 2005.
- Foroutan, I., Sklansky, J. Feature Selection for Automatic Classification of Non-Gaussian Data. *IEEE Transactions on Systems, Man and Cybernetics 17 (2)*, 1987.
- Grishman, R. The NYU system for MUC-6 or where's the syntax?. *Proceedings of the 6th conference on Message understanding (MUC6 '95)*, Association for Computational Linguistics, Stroudsburg, USA, pp. 167-175, 1995.
- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., Witten, I. H. The WEKA data mining software: an update. *SIGKDD Explor. Newsl.*, 11(1), 2009.
- Jurafsky, D., Martin, J.H. *Speech and Language Processing*. Pearson Prentice Hall, 2008.
- Krestel, R., Bergler, S., Witte, R. Minding the Source: Automatic Tagging of Reported Speech in Newspaper Articles. *Proceedings of the Sixth International Language Resources and Evaluation (LREC 2008)*, European Language Resources Association, Marrakech, Morocco, 2008.
- Manning, C.D., Schütze, H. *Foundation of Statistical Natural Language Processing*. MIT Press, 1999.
- Mikheev, A., Moens, M., Grover, C. Named Entity recognition without gazetteers. *Proceedings of the ninth conference on European chapter of the Association for Computational Linguistics (EACL '99)*, Association for Computational Linguistics, Stroudsburg, USA, pp. 1-8, 1999.
- McCallum, A.K. *MALLET: A Machine Learning for Language Toolkit*. <http://mallet.cs.umass.edu>, 2002.
- Mohri, M., Rostamizadeh, A., Talwalkar, A. *Foundations of Machine Learning*. MIT Press, 2012.
- Pastra, K., Maynard, D., Hamza, O., Cunningham, H., Wilks, Y. How feasible is the reuse of grammars for Named Entity Recognition. *Proceedings of the 3rd Conference on Language Resources and Evaluation - LREC 2002*, ELRA - European Language

- Ressources Association, 2008.
- Sayyadi, H., Hurst, M., Maykov, A. Event Detection and Tracking in Social Streams. *Proceedings of International Conference on Weblogs and Social Media (ICWSM 2009)*, AAAI, 2009.
- Smith, M.R., Martinez, T. Improving Classification Accuracy by Identifying and Removing Instances that Should Be Misclassified. *Proceeding of International Joint Conference on Neural Networks*, 2011.
- Tufiş, D., Ion, R., Ceaşu, A., Ştefănescu, D. RACAI's Linguistic Web Services. *Proceedings of the 6th Language Resources and Evaluation Conference - LREC 2008*, Marrakech, Morocco, ELRA - European Language Resources Association, 2008.
- Yang, Y., Pierce, T., Carbonell, J.. A study of retrospective and on-line event detection. *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR '98)*, ACM, New York, USA, pp. 28-36, 1998.