



**UNIVERSITATEA POLITEHNICA BUCUREȘTI  
FACULTATEA DE AUTOMATICA ȘI CALCULATOARE  
CATEDRA DE AUTOMATICA**

# **TEZA DE DOCTORAT REZUMAT**

## **SISTEM AUTOMAT DE ANALIZĂ ȘI PRELUCRARE SEMANTICĂ PENTRU TEXTE SCRISE ÎN LIMBA ROMÂNĂ**

Autor: Ing. Alexandru Catalin Cosoi

Comisia de Doctorat

Presedinte	Prof.dr.ing.Adina Florea		
Conducator de doctorat	Prof.dr.ing.Valentin Sgarciu		
Referent	Prof.dr.ing.Victor Patriciu		
Referent	Prof.dr.ing.Mihaela Oprea		
Referent	Prof.dr.ing.Radu Dobrescu		

**București**

**2014**

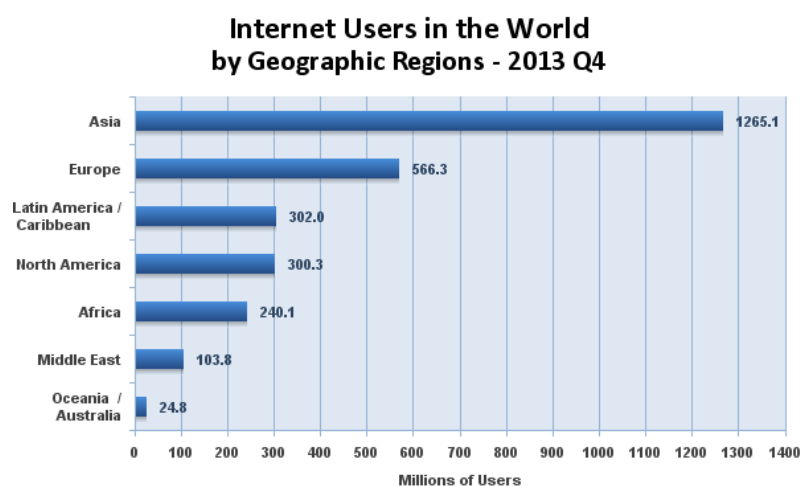
## CUPRINS

<b>INTRODUCERE .....</b>	<b>3</b>
<b>CLASIFICARE SEMANTICA .....</b>	<b>6</b>
ANALIZA SEMANTICII LATENTE.....	6
RETELE DE CONTEXT .....	9
TERM FREQUENCY - INVERSE DOCUMENT FREQUENCY.....	11
REGRESIE LOGISTICĂ .....	12
TEOREMA LUI BAYES .....	13
RETELE SEMANTICE .....	14
CONCLUZII PARȚIALE ȘI ADAPTAREA LA SPECIFICUL LIMBII ROMÂNE .....	15
<b>ANALIZA SENTIMENTELOR .....</b>	<b>16</b>
DEFINIREA OPINIILOR EXPLICITE SI IMPLICITE.....	18
CLASIFICAREA SUBIECTIVITĂȚII ȘI A SENTIMENTELOR.....	21
ANALIZA SENTIMENTELOR BAZATĂ PE CARACTERISTICI.....	21
CONCLUZII PARȚIALE .....	24
<b>SUMARIZAREA.....</b>	<b>24</b>
SEGMENTAREA .....	28
WEBSUMM .....	29
SUMARIZAREA PRIN EXTRAGEREA FRAZELOR CHEIE.....	29
SUMMARIST.....	30
INTERPRETAREA SUBIECTULUI (TOPIC-ULUI).....	31
GENERAREA SUMARULUI.....	32
CONCLUZII PARȚIALE.....	33
<b>SISTEM AUTOMAT DE MONITORIZARE MEDIA .....</b>	<b>35</b>
<b>FRAMEWORK NLP .....</b>	<b>37</b>
SIRI .....	39
CONCLUZII PARȚIALE.....	40
<b>CONCLUZII FINALE, CONTRIBUȚII ORIGINALE ȘI DIRECȚII VIITOARE DE CERCETARE.....</b>	<b>41</b>
CLASIFICAREA DOCUMENTELOR .....	41
ANALIZA SENTIMENTELOR .....	42
SISTEM AUTOMAT DE MONITORIZARE MEDIA .....	43
CONTRIBUȚII PERSONALE.....	43
DIRECȚII VIITOARE DE CERCETARE.....	44
<b>REFERINTE.....</b>	<b>45</b>

## Introducere

Importanța aplicațiilor de prelucrare a limbajului natural a crescut pe măsură ce cantitatea de date disponibile utilizatorilor s-a mărit. Acest lucru se datorează în cea mai mare parte apariției și expansiunii Internet-ului și a trecerii unui număr din ce în ce mai mare de utilizatori la conexiuni de mare viteză. La sfârșitul anului 2012 erau peste două miliarde de utilizatori de Internet, ceea ce reprezintă 34.3% din populația globală.

În mod evident, monitorizarea manuală a acestor informații este o sarcină imposibilă. Dacă de exemplu am dori să facem un rezumat despre tot ce s-a scris în presa online despre un anumit subiect dintr-o anumită zi sau de-a lungul unei săptămâni, ar trebui să extragem mai întâi toate articolele care dezbate subiectul respectiv și apoi să le citim pe fiecare în parte urmând ca apoi să putem crea rezumatul dorit. Dacă este vorba despre o căutare după un anumit cuvânt, am avea un oarecare avantaj pentru că ar trebui să citim doar articolele care conțin cuvântul respectiv (deși ne-ar scăpa toate articolele care s-ar referi la acel cuvânt prin sinonime sau metafore). Dacă însă căutarea se face după o anumită idee, atunci ar trebui să parcurgem toate articolele din toate ziarele și revistele din ziua respectivă sau din săptămâna respectivă și abia apoi să fim capabili să scriem rezumatul lor.

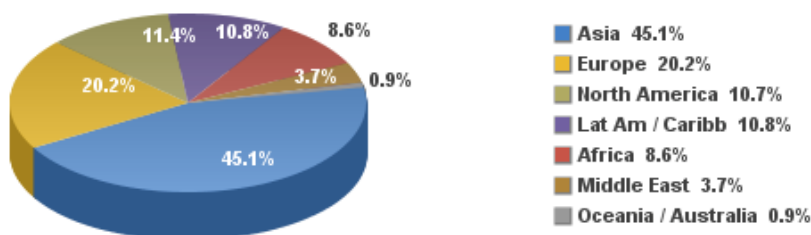


Aplicațiile procesării limbajului natural se înscriu în trei mari categorii:

1. aplicațiile bazate pe text, dintre care amintim:
  - a. clasificarea documentelor (și respectiv găsirea documentelor legate de anumite subiecte);
  - b. regăsirea informației (căutarea unor cuvinte-cheie sau concepte);
  - c. extragerea informației (legate de un anumit subiect, deci de un anumit cuvânt-cheie);
  - d. înțelegerea textelor (care presupune o analiză profundă a structurii acestora);
  - e. traducerea automată și traducerea asistată de calculator dintr-o limbă în alta;
  - f. alcătuirea de sinteze;
  - g. achiziția de cunoștințe.

2. aplicațiile bazate pe dialog, care implică comunicarea între om și mașină, aplicații cum ar fi sistemele de învățare, sistemele de interogare și răspuns la întrebări, rezolvarea problemelor, controlul (bazat pe limba vorbită) al unui calculator, etc.
3. procesarea vorbirii (Este important să facem distincția între problemele de recunoaștere a vorbirii și cele de înțelegere a limbajului. Astfel, trebuie să remarcăm încă de la început faptul că un sistem de recunoaștere a vorbirii nu folosește nici un element de înțelegere a limbajului. Recunoașterea vorbirii se ocupă numai de identificarea cuvintelor vorbite provenind de la un semnal dat, nu și de înțelegerea mesajului, adică a modului în care aceste cuvinte sunt folosite în procesul de comunicare. Pentru a deveni un sistem de înțelegere a limbajului, un dispozitiv de recunoaștere a vorbirii trebuie să furnizeze intrarea sa unui sistem de înțelegere a limbajului natural, operație care produce un așa-numit "sistem de înțelegere a limbajului vorbit". Caracteristica de bază a oricărui sistem de înțelegere este aceea că el realizează o reprezentare a înțelesului propozițiilor într-un limbaj de reprezentare, care poate fi utilizat în vederea unei procesări ulterioare).

### Internet Users in the World Distribution by World Regions - 2013 Q4



Source: Internet World Stats - [www.internetworldstats.com/stats.htm](http://www.internetworldstats.com/stats.htm)  
 Basis: 2,802,478,934 Internet users on Dec 31, 2013  
 Copyright © 2014, Miniwatts Marketing Group

Această lucrare dorește să exploreze principalele tehnici de prelucrare a limbajului natural și posibilitatea utilizării acestora pentru texte scrise în limba română. Un interes deosebit se va arata în **capitolul doi**, metodelor de clasificare semantică precum indexarea semanticii latente sau a grafurilor de context și adaptarea acestora (în funcție de necesități) pentru specificul limbii române.

În termeni matematici, clasificarea semantică este definită ca fiind partiționarea unui set de documente într-un număr de clase semantice. Aceste clase identifică o mulțime de documente care vorbesc (în cea mai mare parte) despre același subiect. De-a lungul timpului, metodele care au avut cel mai mare succes în utilizare au fost cele statistice, însă limbajul natural nu este un domeniu static, iar sinonimia și polisemia sunt elemente care pot crea probleme metodelor statistice. Un cuvânt poate avea mai multe înțelesuri în același timp iar același concept semantic poate avea mai multe reprezentări. Mai mulți algoritmi de clasificare semantică vor fi analizați în acest capitol, analizându-se rezultatele și defectele fiecăruia.

**Capitolul trei** va trata problema analizei automate a sentimentelor. Majoritatea textelor disponibile pe Internet pot fi împărțite în două mari categorii: fapte și opinii. Faptele sunt expresii obiective despre entități și evenimente și de asemenea proprietățile acestora. De cele mai multe ori, opiniile sunt subiective, expresii ce reprezintă sentimentele utilizatorilor, părerile lor despre entități, evenimente sau despre proprietățile acestora. Conceptul de "opinie" este unul foarte larg, însă capitolul trei se va concentra doar pe determinarea sentimentelor pozitive sau negative dintr-un articol. Majoritatea articolelor despre procesarea automată a textelor s-au concentrat pe regăsirea și extragerea informațiilor, căutarea inteligentă pe web, clasificarea de texte, clustering și așa mai departe, însă puțină atenție a fost îndreptată spre procesarea opiniilor. Cu toate acestea, opiniile sunt cele mai importante de fiecare dată când trebuie luată o decizie și dorim să aflăm și părerile celorlalți.

**Capitolul patru** va analiza tehnicile de sumarizare automată. Scopul unui rezumat este simplu și evident: de a facilita identificarea unui subiect din mai multe articole și selectarea materialelor cele mai relevante pentru subiectul de interes. În general, crearea unui rezumat implică o oarecare familiaritate cu subiectul dezbătut. Extragerea elementelor cheie pe care autorul articolului a dorit să le sublinieze necesită de regulă antrenament și experiență. În general, se face o distincție clară între scopul unui rezumat și tipul de rezumat folosit, însă din nefericire, terminologia în acest domeniu este încă insuficient dezvoltată pentru a permite acest lucru. O descriere mai generală împarte rezumatele în două mari subclase: rezumatele indicative, care indică subiectele dezbătute în documente și pot anunța cititorul/utilizatorul despre documentele care dezbate un anumit subiect și rezumate informative care descriu subiectele dezbătute în lucrare.

**Capitolul cinci** va prezenta partea aplicativă a acestei cercetări și de asemenea va pune la dispoziție și un set complex de resurse sintactice și semantice pentru limba română pentru viitoare cercetări. Ca și aplicație demonstrativă, va fi prezentată de asemenea și o platformă de monitorizare media ce va oferi în mod automat aceleași servicii ce acum sunt oferite de diferite agenții de media. Serviciile de monitorizare media implică monitorizarea tuturor canalelor media, împărțirea pe subiecte a articolelor, extragerea rezumatelor (a unuia sau a mai multor articole), posibilitatea de a căuta după anumite cuvinte cheie și posibilitatea măsurării opiniei generale din articolele respective (dacă este vorba despre un ton pozitiv sau negativ), lucru foarte important pentru companii sau pentru campaniile politice.

Iar ultimul capitol, **capitolul 6**, va prezenta concluziile lucrării, posibilitățile de dezvoltare ulterioară și de asemenea contribuția personală.

De asemenea, resursele sintactice din această lucrare au fost concepute în așa fel încât să poată fi folosite atât singure, cât și utilizând framework-uri deja existente, precum NLTK. Scopul acestor resurse este să ușureze munca cercetătorilor români prin oferirea de resurse necesare studiului algoritmilor de procesare a limbajului natural pentru limba română. Resursele vor fi disponibile gratuit pe internet și vor putea fi accesate de oricine dorește să continue aprofundarea acestor cunoștințe prin continuarea rezultatelor prezentate aici.

## Clasificare Semantica

În termeni matematici, clasificarea semantică de texte este definită ca fiind partiționarea unui set de documente într-un număr de clase semantice. Aceste clase identifică o mulțime de documente care vorbesc (în cea mai mare parte) despre același subiect. De-a lungul timpului, metodele care au avut cel mai mare succes în utilizare au fost cele statistice, însă limbajul natural nu este un domeniu static, iar sinonimia și polisemia sunt elemente care pot crea probleme metodelor statistice. Un cuvânt poate avea mai multe înțelesuri și în același timp același concept semantic poate avea mai multe reprezentări.

Față de sinonimie, polisemia s-a dovedit a fi cea mai problematică. De exemplu, verbul "a juca" poate apărea în numeroase contexte și poate determina apartenența textului sau documentului din care face parte în mai multe clase semantice. Pe de altă parte, dacă două cuvinte sunt sinonime, înseamnă că în interiorul aceleiași propoziții, cuvintele sunt interschimbabile fără ca propoziția să-și piardă înțelesul. Însă existența mai multor sinonime pentru același concept poate duce la diminuarea frecvențelor apariției aceluiași cuvânt, ducând implicit la modificarea ponderilor în diverse metode statistice.

În acest capitol vor fi tratate principalele tehnici de indexare și clasificare semantică cum ar fi analiza și/sau indexarea semanticii latente, TF-IDF (term frequency - inverse document frequency), analiza/indexarea probabilistică a semanticii latente, teorema lui Bayes și rețelele bayesiene, rețelele semantice, rețelele de context și deasemenea vor fi prezentate îmbunătățiri aduse acestor metode observate de-a lungul implementării și testării.

### Analiza Semanticii Latente

Analiza Semanticii Latente (ASL) reprezintă o metodă folosită pentru determinarea înțelesului cuvintelor și pasajelor prin procesarea unor baze de date de text. ASL crează o reprezentare a cuvintelor și a submulțimilor de cuvinte folosite - cum ar fi o propoziție, un paragraf sau un eseu - luat din corpusul inițial sau nou, într-un spațiu semantic multidimensional.

Deși ASL este similară cu modelul rețelelor neuronale, este bazată pe descompunerea valorilor singulare, o tehnică matematică de descompunere a unei matrice, similară analizei factorilor, și este aplicabilă corpusurilor de text ce se apropie de volumul și relevanța limbajului folosit de oameni.

Reprezentarea înțelesului pasajelor și al cuvintelor derivată din ASL este capabilă să simuleze o varietate de fenomene cognitive umane, de la construirea experimentală a unui vocabular până la categorisirea de cuvinte, corespondența dintre o frază și un cuvânt, înțelegerea unui discurs și evaluarea calității unui eseu. Câteva dintre aceste simulări vor fi prezentate mai jos.

De exemplu, pentru o metodă practică pentru caracterizarea înțelesului cuvintelor, știm că ASL produce o măsură a relației cuvânt-cuvânt, cuvânt - pasaj și pasaj - pasaj care este bine corelată cu multiple fenomene umane cognitive ce implică similaritatea semantică. Aceste corelații demonstrează asemănarea dintre ceea ce extrage ASL și

modul în care reprezentarea oamenilor a înțelesului reflectă ceea ce au auzit sau au citit, și de asemenea modul în care reprezentarea înțelesului este reflectată în modul în care autorii textelor își aleg cuvintele. Ca și o consecință practică a acestei corespondențe, ASL permite aproximarea judecății umane asupra similarității dintre cuvinte și precizarea obiectivă a consecinței similarității paragrafelor, ambele elemente importante în procesarea de discursuri.

ASL este o metodă statistică automată pentru extragerea și crearea relațiilor dintre cuvinte și utilizarea acestora în pasaje de text. ASL nu este o metodă tradițională de procesare a limbajului natural sau un program bazat pe inteligență artificială; nu folosește dicționare construite de către un om, baze de date de cunoștințe, rețele semantice, gramatici, parsere sintactice, morfologii sau orice altceva din această gamă, și primește ca și date de intrare doar text împărțit în paragrafe sau propoziții.

Primul pas în aplicarea acestei metode este reprezentarea textului sub forma unei matrice în care fiecare rând reprezintă un cuvânt unic iar fiecare coloană reprezintă un pasaj (sau un document). Fiecare celulă conține frecvența cu care cuvântul de pe linia respectivă se regăsește în pasajul de pe coloana respectivă. Apoi, fiecare celulă este supusă unei transformări, prin care frecvența este ponderată printr-o funcție ce exprimă importanța cuvântului atât în pasajul analizat cât și pentru analiza unui discurs în general.

Pasul următor constă în descompunerea valorilor singulare ale matricei. Acesta este un caz particular de analiză a factorilor, sau mai bine zis, generalizarea matematică în care analiza factorilor este un caz special. În cazul SVD (singular value decomposition – descompunerea valorilor singulare - DVS), matricea inițială este descompusă în produsul a trei matrice. O matrice descrie entitățile liniare sub forma unor vectori derivați din valorile ortogonale ale factorilor, a doua matrice descrie același lucru dar pe coloane, iar cea de-a treia este o matrice diagonală ce conține valorile de scalare, astfel încât atunci când se înmulțesc cele trei matrice, se obține matricea inițială. Se poate demonstra matematic că orice matrice se poate descompune perfect utilizând nu mai mult decât cea mai mică dimensiune a matricei originale.

Descompunerea valorilor singulare (DVS) este una dintre cele mai puternice metode în algebra liniară. Probleme importante precum calculul rangului unei matrice, determinarea bazelor ortogonale pentru subspații liniare, calculul proiecțiilor ortogonale pe subspații liniare, problema celor mai mici pătrate, sunt rezolvate cel mai bine folosind DVS.

Descompunerea valorilor singulare ale unei matrice  $A \in R^{m \times n}$  este definită astfel: exista două matrice ortogonale  $U \in R^{m \times m}$  și  $V \in R^{n \times n}$  astfel încât:

$$U^T A V = \Sigma = \begin{pmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{pmatrix}$$

$$\text{unde } \Sigma_1 = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \in R^{r \times r} \text{ cu } \sigma_1 \geq \sigma_2 \geq \dots \geq 0.$$

Numerele reale pozitive  $\sigma_i$  sunt numite valori singulare ale lui A și sunt întotdeauna ordonate crescător. Coloanele  $u_j \in R^m$  ale matricei ortogonale U se numesc valori singulare la stânga ale matricei A. Coloanele matricei ortogonale V se vor numi

vectorsi singulari la dreapta ai matricei A. Calculul descompunerii valorilor singulare are la bază următorul rezultat:

Valorile singulare nenule  $\sigma_i$  ale unei matrice  $A \in \mathbb{R}$  sunt rădăcinile pătrate pozitive ale valorilor proprii nenule ale matricei simetrice pozitiv semidefinite:

$$B = A^T A$$

i.e. dacă  $\lambda_1 \geq \lambda_2 \dots \geq 0$  sunt cele r valori proprii nenule ale lui B, atunci:

$$\sigma_i = \sqrt{\lambda_i}$$

Teorema de mai sus sugerează o procedură pentru calculul valorilor singulare ale unei matrice A, date folosind algoritmul QR simetric pentru calculul valorilor proprii ale matricei  $B = A^T A$ . Această procedură nu este recomandată datorită posibilei slabe condiționări numerice a matricei B.

Aplicarea SVD/DVS în căutarea informației a fost inițial propusă de un grup de cercetători de la Bellcore, respectiv de către Deerwester și colegii săi în 1990 și numită în acest context indexare semantică latentă. ASL a fost comparată cu căutarea standard din spațiul vectorial pe mai multe colecții de documente. S-a ajuns la concluzia că ASL dă rezultate mai bune, în special în cazul căutărilor în care este necesar să se întoarcă un procent cât mai mare din documentele relevante din colecție. Pe de altă parte, din cauza falselor apariții simultane a termenilor, poate duce uneori la o precizie mai scăzută.

Rezultatele ASL depind și de colecția de documente. Într-o colecție heterogenă, documentele pot folosi termeni diferiți pentru a face referire la același subiect, iar ASL poate identifica similaritatea semantică dintre documente ce aparent nu se aseamănă. Dar într-o colecție cu un vocabular omogen, ASL este mai puțin folositoare.

Ca și aplicații ale acestei metode, se poate enumera compararea documentelor în spațiul conceptelor (clasificarea de documente, gruparea datelor, etc), căutarea de documente similare în mai multe limbi, după analiza în prealabil a traducerilor documentelor (cross language retrieval), găsirea de relații între termeni/cuvinte (polisemie și sinonimie) și de asemenea, având o mulțime de termeni ce trebuie să fie căutați, să fie transpuși într-un spațiu de concepte și apoi să poată căuta documente similare (regăsirea informației – information retrieval)

Unul din dezavantajele majore pe care le are indexarea semanticii latente este că de fiecare dată când se dorește inserarea unui document nou, baza ortogonală se schimbă și este implicit nevoie să se refacă descompunerea în valori singulare. Această operație este destul de costisitoare din punct de vedere computațional și are o complexitate calculată la  $O(m^2n + n^3)$  ceea ce presupune un efort considerabil pentru simpla operație de adăugare de documente noi. Acest lucru complică adăugarea de documente noi odată ce au fost deja învățate un set de documente, însă având în vedere că această metodă nu necesită un calcul al frecvențelor stabilit apriorii adăugării altor documente, este mult mai ușoară și mai convenabilă învățarea tuturor



documentelor în același timp. Din acest motiv, LSI este mult mai bine folosită atunci când există o separare inițială clară și concisă iar corpusul de antrenare este stabil și fix.

De asemenea, LSI/LSA a fost o perioadă destul de îndelungată o metodă denigrată deoarece nu ia în considerare ordinea în care apar cuvintele în text și implicit a relațiilor sintactice sau a morfologiei, ceea ce poate duce la diverse erori în unele situații, cu toate ca este capabilă sa extragă înțelesurile corecte în majoritatea situațiilor.

Un alt dezavantaj al acestei metode este ca a fost patentată de Telcordia, lucru ce a dus la cercetarea și dezvoltarea altor metode a.i. cercetătorii sa-și poata continua munca fără a încălca vreun patent.

## Retele de context

Așa cum a fost prezentat mai sus, un pas intermediar în LSA este crearea unei matrice de corespondență între termeni (cuvinte) și documente, care este defapt un tabel de căutare ponderată a frecvențelor tuturor cuvintelor din întreaga colecție de documente. În LSI, această matrice este interpretată ca un spațiu vectorial multidimensional.

O interpretare alternativă a acestei matrice este utilizarea ei ca și graf bipartit cu noduri termeni sau documente unde fiecare valoare nenulă în matrice corespunde unei muchii ce conectează un nod- termen de un nod-document. În acest model, fiecare document are câte o legătură către toți termenii conținuți în el (către toate cuvintele). Frecvențele ponderate din matrice corespund ponderilor atribuite muchiilor din graf. Acest graf poartă numele de graf de context (sau contextual).

Exemplu practic:

- Glacial Ice often appears blue
- Glaciers are made up of fallen snow
- Firn is an intermediate state between snow and glacial ice
- Ice shelves occur when ice sheets extend over the sea
- Glaciers and ice sheets calve icebergs into the sea
- Firn is half as dense as sea water
- Icebergs are chunks of glacial ice under water

Nod	Termen	Aparitii
<b>a</b>	Glacial ice	3
<b>b</b>	Ice	5
<b>C</b>	Glacier	2
<b>D</b>	Snow	2
<b>E</b>	Firn	2
<b>F</b>	Ice sheet	2
<b>G</b>	Sea	3
<b>H</b>	Water	2
<b>I</b>	Iceberg	2
<b>J</b>	Sheet	2

Pe baza documentelor din prima lista, se construiește tabelul de mai sus în care nodurile reprezintă termeni (notate în cazul de mai sus cu a,b,c,d...) iar documentele din care fac parte vor fi notate cu 1,2,3.... De asemenea, se vor păstra și frecvențele de apariție a fiecărui termen cu scopul stabilirii ponderilor în graful bipartit. Această matrice poate fi reprezentată apoi într-un graf de context. Pentru exemplul prezentat mai sus, graful rezultat este prezentat în figura de mai jos:



Cu toate că această metodă este destul de similară cu alte metode din IR (information retrieval) cum ar fi de exemplu grafurile de concept, grafurile (rețelele) de context nu înglobează nici un fel de informație gramaticală sau ierarhică privind relațiile dintre termeni (cuvinte). Structura sa este determinată doar din apariția listei de cuvinte în colecția de documente. Fiecare muchie din acest graf are o anumită pondere, care este direct proporțională cu alegerea ponderării globale sau locale utilizată în generarea matricii de adiacență. Singura constrângere legată de ponderarea muchiilor este că acestea trebuie să se regasească în intervalul (0,1).

Putem căuta în colecția de documente reprezentată prin acest graf prin oferirea unei “energii” nodului de început și prin oferirea posibilității ca această energie să se poată disipa în graf de la nod la nod prin conexiunile sale, urmând câteva reguli simple. Energia totală în orice nod din graf este legată atât de numărul de căi prin care se poate ajunge la acel nod, cât și prin valoarea ponderilor prin care se ajunge la el. Intuitiv, documentele care au foarte mulți termeni (foarte multe cuvinte rare) rar vor fi semantic înrudite. De asemenea, se obțin rezultate similare ca în cazul analizei semantice latente, având în vedere că o căutare după un anumit termen se poate termina sau poate trece printr-un document care nu conține acel termen, dar este suficient de înrudit cu alte documente din listă.

Deoarece energia inițială se disipă în timp ce se răspândește în rețea, (de aici și constrângerea asupra muchiilor grafului), rezultatele fiecărei căutări sunt localizate în majoritatea cazurilor într-o singură zonă din graf, lucru foarte important pentru performanța algoritmului, oferindu-i scalabilitate.

În exemplul prezentat mai sus, o căutare după cuvântul “iceberg” ar începe prin activarea nodului corespunzător acestui cuvânt, adică nodul  $i$ . Acest nod va primi energia inițială  $E$ , care este apoi distribuită nodurilor cu care este conectat, după următorul algoritm:

```

procedure = energize( energy E, node  $n_k$  )
  energy( $n_k$ ) := energy( $n_k$ ) + E
   $E' := E / \text{degree of } n_k$ 
  if (  $E' \geq T$  )
    for each  $n_j$  = node  $n_k$  in  $M_k$ 
       $E'' := E' * \epsilon_{jk}$ 
      energize(  $E''$ ,  $n_j$  )

```

Unde  $M_k$  este mulțimea tuturor nodurilor vecine lui  $n_k$ ,  $\epsilon_{jk}$  este ponderea muchiei ce unește nodurile  $n_j$  și  $n_k$ ,  $T$  este o valoare constantă și reprezintă valoarea prag de activare. Se poate observa că acest algoritm este de fapt o parcurgere în adâncime a grafului. Dacă muchiile ar fi fost sortate în ordinea crescătoare a ponderilor, parcurgerea s-ar fi putut desfășura în mod euristic optimal. Algoritmul poate fi de asemenea modificat să efectueze o parcurgere în lățime.

CNG este cu mult mai rapid decât LSI, este mult mai ușor de înțeles, adăugarea de noi documente în listă nu este atât de complicată ca la LSI și nu necesită efectuarea pașilor de descompunere în valori singulare. De asemenea, modul de reprezentare al datelor este unul foarte intuitiv și ușor de înțeles. Ca și dezavantaje putem menționa faptul că nu există încă suficient de multă documentație și de asemenea, CNG nu se descurcă la fel de ușor cu query-urile complicate.

## Term Frequency - Inverse Document Frequency

Algoritmul TF-IDF (term frequency - inverse document frequency) este folosit în general în domenii ca extragerea informațiilor sau text mining. De regulă, este vorba de o pondere determinată statistic ce evaluează importanța unui cuvânt în interiorul unui document dintr-o colecție de documente. Importanța acestuia este crescută direct proporțional cu numărul de apariții într-un document, și este ponderată de frecvența apariției acestui cuvânt în întreg corpusul de documente.

Diferite variații ale acestei metode sunt folosite de regulă de motoarele de căutare ca și metodă principală de ranking a relevanței documentelor pentru categoriile definite de către utilizator. Un exemplu de funcție de ranking (poate cel mai incipient exemplu) este prin însumarea ponderii tf-idf a fiecărui cuvânt din expresia căutată. Foarte multe modele de ranking pornesc de la acest exemplu de bază.

Să presupunem că avem un set de documente în limba română și vrem să căutăm care sunt acele documente cele mai relevante pentru expresia “vacă maro”. O metodă bună de început este eliminarea din analiză a tuturor documentelor ce nu conțin ambele cuvinte. O primă preselecție scoate din analiză documentele inutile, însă rămân în continuare suficiente documente. Următorul pas ar fi numărarea de câte ori apare fiecare cuvânt în fiecare document și însumarea acestor valori. Numărul de apariții ale unui cuvânt în interiorul unui document, va purta numele de frecvență a cuvântului. (term frequency). De asemenea, este foarte importantă eliminarea cuvintelor puțin

importante. Dacă de exemplu expresia de căutat era “acea vacă maro”, atunci el ne va crea probleme, deoarece “acea” este un cuvânt destul de frecvent. De asemenea, “acea” nu este un cuvânt cheie potrivit cu toate că este foarte des întâlnit în documente, iar “vacă” și “maro” sunt foarte potrivite, deși se regăsesc mult mai rar. De aceea, este nevoie și de un factor care să țină legătura între frecvența cuvintelor și frecvența apariției într-o colecție de documente (inverse document frequency).

Acest *inverse document frequency factor* va micșora ponderea termenilor care se regăsesc foarte frecvent în colecția de documente și o va crește acelor cuvinte care se regăsesc mai rar. Notăm cu  $tf(i,j)$  frecvența cuvântului  $i$  din documentul  $J$ .

$$tf(i,j) = \frac{n_{ij}}{\sum n_{kj}}$$

unde  $n_{ij}$  este numărul de apariții al cuvântului  $i$  în documentul  $j$  iar suma respectiva este suma tuturor frecvențelor cuvintelor din documentul  $j$ .

*Inverse document frequency factor* se definește ca logaritmul raportului dintre numărul total de documente și numărul de documente în care se regăsește cuvântul  $i$ . Bineînțeles, dacă un cuvânt nu se regăsește în nici un document, se practică împărțirea la  $1 + d(i)$ , însă este o situație care nu ar trebui să se întâmple niciodată având în vedere că lista de cuvinte face parte chiar din documentele în care se caută.

$$idf(i) = \frac{|D|}{(|d|t_i \in D|+1)}$$

Ponderea finală se obține prin înmulțirea celor doi factori:

$$tfidf_{ij} = tf_{ij} * idf_i$$

## Regresie Logistică

În statistică, regresia logistică (câteodată numit model logistic sau model logit) este utilizată pentru estimarea probabilității apariției unui eveniment, prin așezarea datelor într-o curbă logistică. Ca multe alte metode de regresie, folosește multe variabile predictoare, atât numerice cât și categoriale. De exemplu, probabilitatea ca o persoană să facă un atac de cord într-o anumită perioadă de timp ar putea fi prezisă dacă s-ar cunoaște vârsta, sexul și masa corporală a persoanei.

$$f(z) = \frac{1}{1 + e^{-z}}$$

Dacă  $z$  ar reprezenta un set de factori de risc, atunci  $f(z)$  ar reprezenta probabilitatea unui anumit rezultat, având în vedere factorii de risc  $z$ .

Această funcție primește valori între minus infinit și infinit și întoarce valori între 0 și 1 sub forma de probabilități. Dacă valoarea întoarsă este mai mare decât o anumită valoare prag, atunci se poate considera că documentul aparține de categoria B. Algoritmul necesită un vocabular, și fiecărui cuvânt din vocabular îi este asociată o pondere.

Algoritmul poate de asemenea învăța atât în timpul testării cât și înainte de testare fiind o fază specială de antrenare. La fel ca și în cazul celorlalte metode mai puțin LSA și CNG, algoritmul nu poate determina relațiile semantice între cuvinte, ci doar similaritatea între documente pe baza stasticilor aparițiilor cuvintelor fiind încadrat în categoria algoritmilor de tip bag of words.

## Teorema lui Bayes

Thomas Bayes (1702-1761) a fost un matematician britanic și în același timp un pastor prezbiterian, cunoscut a fi formulat o teoremă care îi poartă și numele: Teorema lui Bayes (adesea numită legea lui Bayes), publicată post-mortem sub forma unui eseu *Essay Towards Solving a Problem in the Doctrine of Chances* (1764), de către prietenul său Richard Price în *Philosophical Transactions of the Royal Society of London*. Teorema lui Bayes este una din teoremele fundamentale ale teoriei probabilității, care determină probabilitatea apartenenței evenimentelor și a obiectelor la o anumită grupă. În cazul filtrelor spam bazate pe teorema lui Bayes de exemplu (numite și filtre bayesiene), pentru determinarea probabilității apartenenței unui anumit mesaj la spam, sunt utilizate dicționarele create în timpul "învățării" filtrului. De regulă programul "învăță" analizând arhivele de email-uri, selectate în prealabil manual. Când dicționarele sunt create definitiv, probabilitatea apartenenței unui nou mesaj la spam este calculată prin normalizarea și însumarea probabilității fiecărui cuvânt în parte. Prin urmare, adunând informații statistice despre rata de apariție a unor diferite cuvinte și structuri în mesajele de tip spam sau în mesajele legitime, filtrul compară apoi noile mesaje cu aceste modele și le clasifică corespunzător. Formula de calcul al probabilităților în cazul filtrării antispam poate fi văzută mai jos:

$$PR(S|W) = \frac{PR(W|S) * PR(S)}{PR(W|S) * PR(S) + PR(W|H) * PR(H)}$$

În mod evident, clasificarea bayesiană se poate face pe mai mult de două categorii. În "Introducing Syntactic Features into a Bayesian Classifier" în 2008, am arătat că în funcție de numărul de categorii și de cantitatea de text din fiecare document, rezultatele obținute în urma unei clasificari bayesiane sunt în general bune sau foarte bune, însă dacă se introduc de asemenea și informațiile sintactice pentru fiecare cuvânt, clasificarea se poate îmbunătăți cu până la 10% (bineînțeles, acolo unde rata de clasificare permite acest lucru).

Sistemul care s-a folosit în lucrarea respectivă a constat într-un filtru bayesian clasic și algoritmul Link Grammar pentru determinarea părților sintactice din propozițiile analizate. Informația sintactică a fost adăugată prin aducerea cuvintelor într-o formă canonică de forma "parte sintactică\_cuvânt". Testele s-au efectuat pe un total de 25 de categorii semantice, formate din diferite domenii și subdomenii științifice și conținând câteva zeci de exemple. În mod evident, erorile de clasificare au fost mult mai mari în cazul subdomeniilor dintr-un același domeniu științific comparat cu documentele din domenii diferite, însă introducerea părților sintactice a adus îmbunătățiri în ambele situații. Ca și dezavantaje, timpul de antrenare s-a mărit iar baza de date de cuvinte s-a dublat, lucru de așteptat, având în vedere canonizarea cuvintelor.

O alta metodă de filtrare bayesiană constă în utilizarea rețelelor bayesiene. Modelul rețelelor Bayesiene, introdus de Judea Pearl [1988], pornește de la modelul probabilistic Bayesian, dar elimină numărul enorm de calcule necesare prin

considerarea caracteristicilor de modularitate și de cauzalitate ale domeniului problemei. Ideea de bază a rețelelor Bayesiene este aceea că, pentru a descrie domeniul problemei, nu este necesar să se considere probabilitățile tuturor perechilor de evenimente (fapte) posibile. Cele mai multe evenimente sunt independente între ele și interacțiunile dintre acestea nu trebuie considerate, deoarece nu există. Rețele Bayesiene, numite și rețele de încredere aparțin familiei modelelor grafice probabilistice. Modelele grafice reprezintă grafuri în care nodurile sunt variabile aleatoare și arcurile sunt condiții de propagare. Aceste structuri sunt utilizate pentru reprezentarea cunoașterii despre un anumit domeniu incert. În mod particular, fiecare nod dintr-un asemenea graf reprezintă o variabilă aleatoare, în timp ce liniile dintre noduri reprezintă dependențe probabilistice între variabilele aleatoare corespondente. Aceste dependențe condiționale din graf sunt în general estimate folosind metode statistice de calcul. Deci, rețelele Bayesiene combină principii din teoria grafurilor, teoria probabilităților, știința calculatoarelor și din statistică.

O rețea bayesiană poate fi folosită pentru predicția sau analiza problemelor din lumea reală și din sistemele naturale complexe, unde corelațiile statistice pot fi găsite între variabile sau approximate folosind opiniile experților. Aceste rețele au o vastă arie de aplicabilitate, cum ar fi ajutorul în luarea deciziilor în domenii ca medicină, inginerie, resurse naturale și managementul decizional.

Ca și aplicativitate în clasificarea semantică exemplele sunt multiple iar construcția depinde foarte mult despre ce problema este vorba. O subproblemă care poate fi rezolvată cu ajutorul rețelelor bayesiene este împărțirea documentelor într-un set de categorii prestabilit. Un alt exemplu este prezentat în unde este descris un sistem bazat pe rețele bayesiene capabil să poarte conversații și să descifreze expresii complexe. De asemenea, o dată ce sunt parcurse noțiuni precum "analiza semanticii latente" sau "grafurile de context", rețelele bayesiene reprezintă încă o alternativă ușor de înțeles pentru aceste metode.

## Retele Semantice

O rețea semantică utilizează grafuri pentru reprezentarea cunoștințelor. Conceptul de rețea semantică a fost introdus aproximativ în perioada 1965-1970. Scopul primelor rețele semantice a fost reprezentarea cuvintelor din limbajul natural. O rețea semantică are forma unui ansamblu de noduri și arce, orientate și etichetate. Nodurile reprezintă obiectele. Un obiect poate fi un concept abstract sau particular, un atribut, ș.a.m.d. Arcele sunt utilizate pentru a reprezenta legăturile care există între aceste obiecte.

Rețelele semantice sunt reprezentări geometrice utilizate pentru reprezentarea cunoștințelor și organizarea lor în noduri și conexiuni. Implementările pe calculator a rețelelor semantice au apărut mai întâi pentru diverse aplicații de inteligență artificială și traducere automată, dar chiar și implementări mai rudimentare au fost folosite foarte mult timp în filozofie, psihologie și lingvistică.

Câteva din rețelele semantice prezentate mai sus au fost dezvoltate în mod explicit să implementeze câteva din sistemele cognitive umane, însă altele au avut ca scop doar îmbunătățirea performanțelor diversilor algoritmi de procesare a limbajului natural. Însă, o dată cu apariția rețelelor, multe din soluțiile computaționale au fost apropiate

de cele care implementau sisteme cognitive. De exemplu, diferența între rețelele de definire și cele aserționale este similară cu distincția între memoria semantică și memoria episodică.

## Concluzii parțiale și Adaptarea la specificul limbii române

Majoritatea tehnicilor prezentate în acest capitol sunt în general statistice, așa că sunt destul de puține îmbunătățiri care pot fi aduse algoritmilor existenți. Rezultatele sunt în general la fel de bune dacă corpusul de antrenare este la fel de bun și la fel de complex, ceea ce ne face să concluzionăm că dacă reușim să avem un corpus de antrenare exhaustiv, atunci implicit și rezultatele vor fi pe măsură.

De asemenea, așa cum se poate observa și din exemplele de mai sus, gradul de complexitate și necesarul de resurse este direct proporțional cu numărul de cuvinte din corpusul de antrenare și în mod implicit cu dimensiunea corpusului. De asemenea, în funcție de timpul la care se face referire sau în care se petrece acțiunea oferă un grad de variație destul de larg atât verbelor cât și celorlalte cuvinte prin variația diacriticelor, a genului și a persoanei. De exemplu, variația [acest, această, acesta, aceștia, acestea, acești, aceste], cu toate că reprezintă mereu același lucru (un substitut pentru un anumit obiect/persoană), variază în funcție de număr și sex.

În funcție de problema care se dorește a fi rezolvată, se pot efectua diferite operații de preprocesare. Așa cum am văzut în “Introducing Syntactic Features into a Bayesian Classifier” inclusiv rolul sintactic al fiecărui cuvânt poate avea un rol esențial în clasificarea documentelor.

Când vorbim de clasificarea semantică și a modului în care poate fi ea aplicată, trebuie să fim foarte atenți la problema pe care trebuie să o rezolvăm și abia pe urmă să decidem ce algoritm se potrivește mai bine, sau ce corpus de antrenare este necesar. De exemplu, majoritatea algoritmilor prezentați mai sus se bazează pe un model de învățare supervizată, care în mod evident duce la rezultate mult mai bune decât în cazul unei antrenări nesupervizate, însă la cantitatea din ziua de astăzi de documentație și de informații, problema pregătirii corpusului de antrenare se complică exponențial.

Să luăm următorul exemplu: dorim să creăm un sistem automat de monitorizare media, care este capabil să:

- Preia din presa online toate articolele cu o anumită frecvență prestabilă. Statistic, frecvența optimă este din oră în oră
- Cum multe ziare vor acoperi același subiect, sistemul va trebui să fie capabil să creeze clustere de articole care tratează același topic/categorie
- Majoritatea publicațiilor oferă o preclasificare a articolelor în categorii precum economic, politic, etc, însă acest lucru este insuficient în cadrul unui sistem de monitorizare media, deoarece se dorește o filtrare la nivel de topic și nu la nivel de categorie
- Pentru fiecare topic, sistemul este capabil să determine dacă este vorba despre un topic pozitiv sau negativ. Acest lucru este foarte important în cazul review-ului de produse, al politicii, etc

- De asemenea, pentru fiecare topic în parte, pe baza articolelor din ziare diferite care ating același subiect, să se creeze rezumatul

La volatilitatea știrilor din ziua de azi, crearea unui corpus de antrenare la nivel de categorie (politic, financiar, tabloid, etc) este un lucru trivial, însă nu se poate spune același lucru și despre clasificarea la nivel de topic.

O bună metodă de creare a unui corpus de antrenare la nivel de topic s-a dovedit a fi Similaritatea Jaccard, care cu toate că pare destul de simplă, s-a dovedit în testele mele, în special pe zona clasificării articolelor din presa online, a avea rezultate foarte bune.

Distanța Jaccard dintre două documente este:  $d = \frac{|A \cap B|}{|A \cup B|}$ , unde A este mulțimea cuvintelor din primul document iar B este mulțimea cuvintelor din al doilea document. Rezultatele sunt în general foarte bune, dar în capitolul 5 vom arata diferite metode de îmbunătățire, prin adăugarea de semidistanțe Jaccard (numărul de apariții al fiecărui cuvânt, dimensiunea documentelor, etc), care ne vor duce gradul de acuratețe care în mod normal la distanța Jaccard clasică este undeva între 70-80 la sută. În cazul testelor mele, în gama 98-100 la sută. Însă și distanța Jaccard clasică poate fi folosită dacă dorim doar crearea unui corpus de antrenare.

## Analiza Sentimentelor

Majoritatea textelor disponibile pe Internet pot fi împărțite în două mari categorii: fapte și opinii. Faptele sunt expresii obiective despre entități și evenimente și de asemenea proprietățile acestora. De cele mai multe ori, opiniile sunt subiective, expresii ce reprezintă sentimentele utilizatorilor, părerile lor despre entități, evenimente sau despre proprietățile acestora. Conceptul de "opinie" este unul foarte larg, însă acest capitol se va concentra doar pe determinarea sentimentelor pozitive sau negative dintr-un articol. Majoritatea articolelor despre procesarea automată a textelor s-au concentrat pe regăsirea și extragerea informațiilor, căutarea inteligentă pe web, clasificarea de texte, clustering și așa mai departe, însă puțină atenție a fost îndreptată spre procesarea opiniilor. Cu toate acestea, opiniile sunt cele mai importante de fiecare dată când trebuie luată o decizie și dorim să aflăm și părerile celorlalți. Acest lucru este adevărat atât pentru indivizi, cât și pentru organizații.

Unul din multele motive pentru lipsa studiului opiniilor este că o mare perioadă de timp a fost foarte puțin text ce exprima opiniile oamenilor disponibil pe Internet. Înainte de utilizarea web-ului la o scară largă, în momentul în care cineva trebuia să ia o decizie, și-ar fi întrebat prietenii sau familia. Când o organizație dorește să facă acest lucru, va comanda sondaje și cercetări de piață. Dacă luăm în considerare creșterea explozivă a Internetului și a cantității de conținut generat de către utilizatori în ultimii ani, atunci putem să concluzionăm că lumea s-a schimbat. Internetul a schimbat modul în care oamenii își exprimă opiniile și părerile. Acum își pot expune părerile direct pe site-urile vânzătorilor și își pot prezenta întreaga experiență pe care au avut-o cu produsul respectiv, astfel încât următorii cumpărători să fie avizați. Această nouă lume oferă o sursă nemărginită de informații ce poate prezenta beneficii substanțiale unei game foarte largi de utilizatori. Acum nu mai este nevoie să suni un



prieten dacă ai nevoie de un sfat în vederea achiziționării unui produs ci poți avea acces la informațiile dorite imediat.

Cu toate acestea, găsirea de surse de opinii valide și monitorizarea lor pe internet încă poate cauza numeroase probleme datorită numărului foarte mare de surse, iar fiecare dintre acestea poate avea un număr foarte mare de opinii. De cele mai multe ori, opiniile sunt ascunse în discuții lungi pe blog-uri sau forum-uri. Este destul de dificil pentru o persoană să găsească surse relevante, să extragă propozițiile și frazele ce conțin opinii, să le citească, să le summarizeze și să le organizeze într-o formă utilizabilă. În mod evident, necesitatea unor sisteme ce pot extrage singure opiniile din texte și să summarizeze aceste opinii, este destul de mare. Analiza sentimentelor, cunoscută de altfel și ca "mineritul sentimentelor", s-a născut din această necesitate. Având în vedere aplicațiile multiple pe care le poate avea, atât în mediul industrial cât și în cel academic, acest domeniu a explodat în ultimii ani. La momentul actual, există cel puțin 20-30 de companii care oferă analiza sentimentelor pentru industrie.

Clasificarea sentimentelor sau clasificarea sentimentelor la nivel de document încearcă să stabilească o părere generală asupra sentimentelor ce reies din text cu privire la subiectul dezbătut. Dacă de exemplu ne gândim la evaluarea unui produs, o astfel de aplicație ar determina în mod automat dacă este o evaluare pozitivă sau negativă. Celălalt pas se referă la identificarea unei opinii la nivel de propoziție și stabilirea tipului de opinie, respectiv negativă sau pozitivă.

Analiza sentimentelor bazată pe caracteristici: acest model caută obiectele asupra cărora sunt exprimate sentimentele dintr-o propoziție, urmând ca apoi să se determine dacă sunt pozitive sau negative. Un obiect poate fi reprezentat printr-un produs, un serviciu, un individ, o organizație, un eveniment, un subiect, etc. De exemplu, în cazul unei evaluări de produs, se poate determina la ce componente ale produsului se face referire în cadrul evaluării și dacă exprimă o părere pozitivă sau negativă. De exemplu, propoziția "durata de viață a bateriei acestei camere este prea scurtă" este un comentariu ce se referă la durata de viață a bateriei unei camere, iar comentariul este unul negativ. Acest gen de aplicații sunt necesare în foarte multe situații: de exemplu, ca o companie să-și îmbunătățească produsele, trebuie să fie la curent cu toate evaluările produselor curente.

Evaluarea unui obiect se poate face prin două metode principale: direct sau prin comparație. Analiza directă sau opiniile directe oferă în mod direct sentimentele negative sau pozitive asupra subiectului dezbătut, fără a fi necesară menționarea altor obiecte. Analiza prin comparație înseamnă comparația obiectului analizat cu alte obiecte similare, eventual concurente. De exemplu, următoarea propoziție reprezintă o opinie directă: "calitatea acestei camere este slabă", iar următoarea propoziție reprezintă o opinie pozitivă exprimată prin comparație: "Calitatea pozelor este mai bună decât a celor făcute cu camera x". Este evidentă necesitatea determinării acestor tipuri de propoziții și situarea produsului analizat în comparație cu concurența.

Așa cum este ușor de înțeles de ce opiniile și sentimentele sunt foarte importante de monitorizat, la fel de ușor de înțeles este și de ce este foarte important pentru companii să răspândească opinii pozitive în propriul interes, create sau nu de oameni. Acest fenomen se numește "opinion spam" și de regulă se traduce prin oferirea de comentarii negative produselor concurenței și oferirea de opinii pozitive produselor

proprii. Acest gen de opinii sunt inutile pentru un evaluator și este foarte importantă eliminarea acestora din lista opiniilor ce vor ajunge la utilizator.

Sentiment Analysis (Analiza sentimentelor) sau Opinion Mining (mineritul opiniilor) este studiul computațional al opiniilor, sentimentelor și emoțiilor exprimate într-un text. De exemplu, în următoarea evaluare, se prezintă un set de păreri personale despre iPhone:

*“Mi-am cumpărat noul iPhone acum câteva zile. Este un telefon foarte bun. Ecranul tactil este foarte incitant. De altfel, inclusiv calitatea vocii este foarte bună. Cu toate că bateria nu rezistă foarte mult, nu mă deranjează foarte tare lucrul acesta. Cu toate acestea, mama s-a suparat pe mine ca nu i-am spus înainte sa-l cumpăr. I s-a părut scump și voia să ne întoarcem la magazin sa-l returnăm.”*

Putem observa că există atât păreri pozitive cât și negative. De asemenea, se poate remarca că fiecare propoziție face referire la un anumit obiect despre care este enunțată o opinie. Se mai poate observa și că în textul de mai sus, majoritatea opiniilor sunt pozitive, iar cele care sunt negative sunt atribuite altei persoane.

În general, opiniile pot fi exprimate în legătură cu orice (un produs, un serviciu, un individ, etc). Fiecare obiect poate avea la rândul său un set de atribute (sau proprietăți, iar dacă obiectul respectiv este compus din mai multe componente, atunci fiecare componentă poate avea la rândul său un set de atribute sau proprietăți. În mod evident, această relație poate continua. Se poate concluziona că un obiect poate fi descompus într-un arbore de componente, subcomponente și proprietăți ale acestora. O definiție corectă ar fi: "un obiect este o entitate care poate fi un produs, o persoană, un eveniment, o organizație, sau un subiect. Obiectului îi este asociată o pereche (T,A) unde T este o ierarhie de componente, subcomponente, s.a..m.d iar A este un set de atribute." În exemplul de mai sus, telefonul este obiectul iar câteva din componentele sale sunt bateria și ecranul. De asemenea, setul adiacent de atribute este format din calitatea vocii, mărime și greutate. Bateria are propriul set de atribute: durata de viață și dimensiunile. Astfel, un obiect poate fi reprezentat sub formă arborescentă ce are ca rădăcină obiectul în sine. Fiecare nod care nu este rădăcină este ori o componentă ori o subcomponentă a obiectului iar fiecare muchie este relația dintre acestea. De asemenea, fiecare nod are asociat un set de atribute. O opinie se poate lega de orice componentă din acest arbore.

## Definirea opiniilor explicite si implicite

O opinie explicită referitoare la caracteristica  $f$  este o opinie exprimată vis a vis de  $f$  într-o propoziție subiectivă. O opinie implicită referitoare la caracteristica  $f$  reprezintă o opinie asupra  $f$  subînțeleasă în cadrul unei propoziții obiective. Propoziția următoare ilustrează exemplul clasic al unei opinii pozitive explicite: *Calitatea sunetului pentru telefonul acesta este extraordinara*. Comparativ, această propoziție ne oferă un exemplul unei opinii negative implicite: “Difuzorul telefonului s-a defectat in doua zile”. Astfel, deși această propoziție prezintă o realitate obiectivă, implicit indică și o opinie negativă referitoare la difuzorul telefonului.

Definirea propoziției în care sunt exprimate opinii: O propoziție în care sunt exprimate opinii este o propoziție care exprimă în mod explicit sau implicit opinii pozitive sau negative. Aceasta poate fi o propoziție de natură subiectivă sau obiectivă.

Asa cum se poate observa, noțiunea de propoziție subiectivă și cea de propoziție în care sunt exprimate opinii nu se suprapun perfect, deși în multe cazuri, propozițiile în care sunt exprimate opinii reprezintă de fapt o subcategorie a propozițiilor subiective.

Abordările utilizate pentru clasificarea acestor două tipuri de propoziții este totodată similară. Metoda prin care se determină dacă o propoziție este de natură subiectivă sau obiectivă poartă denumirea de clasificare a subiectivității.

În mod evident și documentele pot fi evaluate în ce măsură prezintă sau nu opinii. Până acum, am considerat din start documentele ca fiind unele în care sunt prezentate opinii. În practică însă trebuie stabilită și această caracteristică a documentului. De exemplu, forumurile de discuții dispun și de numeroase pasaje cu întrebări și răspunsuri care nu conțin opinii. Firul logic pe care îl vom aplica mai departe consideră că un document conține propoziții în care sunt exprimate opinii, atunci și documentul respectiv poate fi considerat drept unul care prezintă această caracteristică. Această definiție nu este totuși exhaustivă, deoarece există și cazuri precum o știre, în care sunt preluate opiniile altor persoane însă numai cu titlu de informație. Nu ar fi în întregime corect să catalogăm documentul respectiv drept unul subiectiv.

Cel mai corect în cazul de față ar fi să conchidem că documentul respectiv prezintă într-o manieră obiectivă o serie de opinii care aparțin unei alte persoane, alta decât autorul documentului. O definiție mai clară a tipologiei documentelor poate fi realizată pe baza intenției autorului, respectiv dacă acesta intenționează sau nu să exprime o serie de opinii proprii prin intermediul textului. Evaluările de produse corespund acestei categorii de documente. Dacă o propoziție conține sau nu opinii personale este mai ușor de stabilit, în timp ce într-un text obișnuit se regăsesc atât propoziții care conțin opinii cât și altele fără.

Prin acest lucru se dovedește cât de necesară este aplicarea unei forme de sumarizare asupra rezultatelor obținute în urma procesului de căutare (“mining”), deoarece valoarea adăugată oferită de o simplă listare a opiniilor identificate și a elementelor acestora este minimă pentru utilizator. În vederea atingerii obiectivului principal și anume acela de a transforma limbajul natural nestructurat în date structurate, ne putem folosi cu succes de evidențierea rezultatelor sub forma celor 5 elemente esențiale. Datele astfel descompuse pot fi păstrate în baze de date ce permit interogări ulterioare. În cazul bazelor de date astfel construite se pot aplica tehnici de vizualizare pentru a putea beneficia de rezultate în formatul cel mai elocvent extragerii de informații relevante, care de obicei se numesc sumări structurate și prezintă datele sub forma de grafice.

Sumarizarea structurată bazată pe opinii reprezintă o metodă simplă de a produce o sumarizare la nivel de caracteristici ale opiniilor la nivelul unui obiect sau a unor obiecte aflate în competiție/opoziție.

Presupunând că dorim să realizăm sumarizarea unei evaluări pentru un telefon mobil, denumit în continuare Telefon Mobil 1, rezumatul va arăta așa cum este ilustrat în tabelul de mai sus. În cadrul figurii, cuvântul “telefon” reprezintă telefonul în sine

(rădăcina reprezentării ierarhice a obiectului). 125 de evaluari erau de natură pozitivă, în timp ce doar 7 erau de natură negativă cu privire la performanța telefonului.

Calitatea sunetului și mărimea telefonului sunt două caracteristici ale obiectului la care s-a făcut referire, prima dintre ele primind 120 de evaluări pozitive și doar 8 negative. Propozițiile care conțin evaluări venite de la o singură persoană pot fi corelate cu alte propoziții sau chiar cu un text de evaluare, care oferă comentarii pozitive sau negative vis a vis de o caracteristică a produsului. Folosind acest beneficiu, clienții pot afla relativ ușor care este părerea clienților actuali față de produsul în cauză. Dacă persoana care citește evaluarea este interesată într-o masura mai mare de evaluarea unei anumite caracteristici ale produsului, atunci ea poate urmări link-ul din cadrul evaluării individuale pentru a afla mai multe păreri și de la alți clienți. Un astfel de rezumat al informației poate fi vizualizat folosind graficele de tip bara.

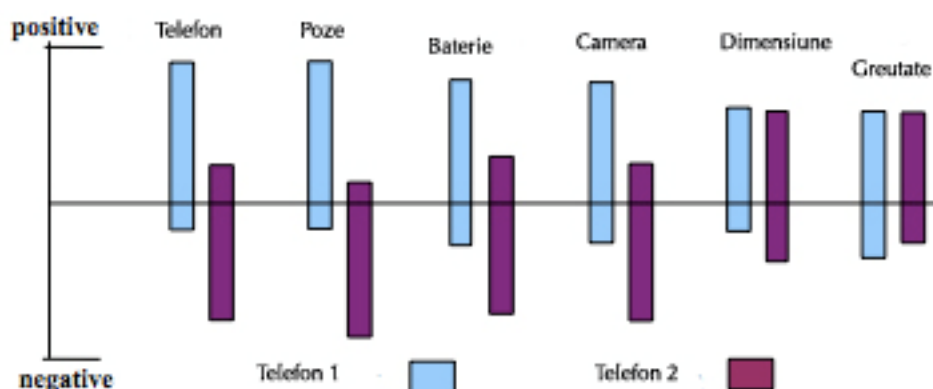


Figura de mai sus ilustrează un rezumat al opiniilor în cadrul unei evaluări pentru un telefon mobil. Fiecare bară situată deasupra axei X mediane indică numărul de opinii pozitive înregistrate de către caracteristica respectivă a produsului, iar bara situată dedesubtul axei X ne arată câte opinii negative s-au înregistrat cu privire la aceeași caracteristică. Cu siguranță, putem folosi și alte metode de vizualizare, cum ar fi prezentarea procentului de opinii pozitive, cuantumul opiniilor negative fiind calculat drept unitatea minus procentul de opinii pozitive.

Față de telefon în sine, majoritatea persoanelor și-au exprimat o părere pozitivă față de primul dintre produse și una de natură negativă față de cel de-al doilea telefon. După cum se poate observa, vizualizarea permite utilizatorilor să observe cum performează un telefon față de un competitor de-a lungul unei liste de caracteristici. În afară de aceste tehnici de vizualizare prezentate sunt posibile multe alte tipuri particularizate în funcție de obiectivele urmărite sau de contextul în care se doresc a fi folosite. Totodată, și sumarizarea care nu vizează texte în care sunt exprimate opinii poate fi foarte folositoare, după cum se va putea observa din exemplele următoare.

Sumarizarea bazată pe aparițiile unei anumite caracteristici ("feature buzz") se referă la frecvența relativă a mențiunilor cu privire la o anumită caracteristică a obiectului. Această tehnică poate indica managementului unei companii care sunt cu adevărat aspectele importante care îi preocupă pe clienții lor. De exemplu, în cadrul unui studiu asupra mediului online bancar, caracteristica cea mai des menționată de către clienți a fost securitatea tranzacțiilor. Sumarizarea bazată pe aparițiile unui anumit obiect

(“object buzz”) vizează cuantificarea frecvenței de menționare a diferite produse sau mărci concurente. Această metodă este utilă deoarece oferă informații privind gradul de popularitate al produselor sau mărcilor respective la nivelul pieței în care activează.

De asemenea, având în vedere că momentul când s-a realizat exprimarea opiniei este cuantificat ca fiind unul dintre cele cinci elemente de bază, o altă analiză care poate fi implementată se referă la monitorizarea în timp a schimbărilor survenite la nivelul fiecărui aspect. Monitorizarea trendurilor (“trend tracking”) poate fi obținută prin adăugarea dimensiunii temporale în cadrul analizelor de mai sus.

Raportarea ce include și monitorizarea trendurilor este extrem de utilă pentru utilizatorii finali, deoarece le permite să fie la curent cu schimbările survenite în timp. Toate aceste rezumate pot fi realizate și vizualizate extrem de ușor deoarece sunt rezultatele interogării unei baze de date fără a mai fi nevoie de acțiuni suplimentare, dovedind încă o dată valoarea adusă de structurarea informației care conține opinii în funcție de cele cinci elemente prestabilite.

## **Clasificarea subiectivității și a sentimentelor**

Atunci când discutăm despre analiza sentimentelor, clasificarea este procesul la care se va face referire cu precădere. Astfel, clasificarea sentimentelor vizează crearea unei diferențieri la nivelul documentelor în care sunt exprimate opinii în funcție de tipul opiniei: pozitivă sau negativă. Această activitate este cunoscută și sub numele de clasificare a sentimentelor la nivelul documentelor, deoarece ia în considerare întreg documentul drept unitate de bază pentru analiză. Cercetarea existentă în domeniu pornește de la premisa ca documentul este deja considerat ca fiind unul în care sunt exprimate opinii. În mod firesc, același procedeu de clasificare a sentimentelor poate fi aplicat și la nivelul propozițiilor, însă în această situație literatura de specialitate nu consideră ca fiecare propoziție drept una în care sunt exprimate opinii. Metoda prin care se stabilește dacă într-o propoziție sunt exprimate sau nu opinii, poartă denumirea de clasificare a subiectivității, pentru că ulterior, propozițiile despre care se consideră că exprimă opinii să fie clasificate în funcție de tipul opiniei exprimate, pozitivă sau negativă, folosind metoda de clasificare a sentimentelor la nivelul propoziției.

## **Analiza sentimentelor bazată pe caracteristici**

Deși clasificarea textelor în care sunt exprimate opinii la nivelul documentului sau a propoziției se dovedește a fi utilă în multe cazuri, aceasta clasificare nu oferă gradul de detaliere necesar în cazul unor alte aplicații. Caracterizarea unui document drept unul în care sunt exprimate opinii pozitive vis a vis de un obiect evaluat, nu înseamnă că autorul documentului are o părere pozitivă despre toate trăsăturile obiectului respectiv. Aceași regulă aplicându-se și în cazul în care opinia generală desprinsă la nivelul textului este una negativă. În general, în cazul unui text în care sunt exprimate opinii cu privire la un obiect, persoană, serviciu, etc, autorul include atât referințe pozitive sau negative cu privire la trăsăturile obiectului, deși opinia desprinsă la nivel general poate să fie ori pozitivă, ori negativă.

Pe de altă parte, atunci când vine vorba de identificarea obiectului evaluat, această informație este de o reală valoare și în cazul textelor informale, așa cum sunt cele disponibile online. Deși în general atunci când se face evaluarea unui produs se cunoaște obiectul la care se face referire, de-a lungul evaluării există și posibilitatea de a se face referi la alte obiecte pentru comparații, ceea ce îngreunează relativ procesul de identificare a tuturor obiectelor relevante pentru textul respectiv.

În ceea ce privește extragerea cu exactitate a perioadei în care a fost redactat textul, această informație este mai ușor de reperat în cazul informațiilor postate online, unde orice postare este înregistrată și se menționează ora exactă la care aceasta a devenit disponibilă publicului, însă destul de dificilă în cazul altor tipuri de documente. Modalitatea prin care sunt extrase aceste trei tipuri de informații la nivelul unui text, respectiv persoana care emite opinia, obiectul sau obiectele la care se fac referire și perioada la care s-a expus opinia, poartă denumirea de Named Entity Recognition (NER).

Una dintre resursele ce va fi disponibilă pe site este o listă, în format XML, ce va conține un set de obiecte, și pentru fiecare obiect, un set de criterii ce îl definesc. Exemplu:

```
<title>Obiecte și componente</title>
<continut>

<obiect>
  <nume>telefon</nume>
  <caracteristica>ecran</caracteristica>
  <caracteristica>baterie</caracteristica>
  <caracteristica>greutate</caracteristica>
  <caracteristica>dimensiune</caracteristica>
</obiect>

<obiect>
  <nume>baterie</nume>
  <caracteristica>greutate</caracteristica>
  <caracteristica>durată de viață</caracteristica>
</obiect>

<obiect>
  <nume>laptop</nume>
  <caracteristica>ecran|display</caracteristica>
  <caracteristica>greutate</caracteristica>
  <caracteristica>fiabilitate</caracteristica>
  <caracteristica>dimensiuni</caracteristica>
  <caracteristica>RAM</caracteristica>
  <caracteristica>hard|hard-disk|HDD</caracteristica>
  <caracteristica>procesor</caracteristica>
  <caracteristica>culoare</caracteristica>
  <caracteristica>design</caracteristica>
</obiect>

</continut>
```

Un format care ar fi mult mai util pentru procesare online ar fi pastrarea într-un format mult mai rapid și ușor de procesat de către limbajele ce se folosesc pentru programarea web. Un format mult mai recomandat ar fi păstrarea în format JSON, ceea ce înseamnă că fișierul va arata astfel:

```
{"telefon":{"cnt":4,"p1":"ecran","p2":"baterie","p3":"greutate",
"p4":"dimensiune"}}
{"baterie":{"cnt":2,"p1":"greutate","p2":"durată de viață"}}
{"laptop":{"cnt":9,"p1":"ecran|display","p2":"greutate","p3":
"fiabilitate","p4":"dimensiuni","p5":"RAM","p6":"hard|hard-disk|HDD",
"p7":"procesor","p8":"culoare","p9":"design"}}
```

De asemenea, având în vedere că se lucrează cu cantități foarte mari de date, este atât dificil cât și nepractic să păstrăm aceste informații într-un fișier text, așa că una dintre propuneri ar fi păstrarea acestor informații într-o bază de date. Una din bazele de date foarte recente aparute în industrie și care se împacă nativ formatul JSON, este MongoDB.

MongoDB este capabilă să memoreze cantități impresionante de informație și datorită modului diferit de reprezentare a datelor, este de asemenea capabilă să răspundă foarte rapid la cereri.

De exemplu, crearea unei colecții în care să păstrăm informații ca cele prezentate mai sus, se face folosind următoarele comenzi.

```
mongo
use mydb
db.elements.save({"telefon":{"cnt":4,"p1":"ecran","p2":"baterie",
"p3":"greutate","p4":"dimensiune"}});
show collections
```

În mod evident, fiecare element al unui obiect poate are anumite proprietăți care sunt pozitive și altele care sunt negative. De asemenea, nu există o corelație foarte clară între elementele unor obiecte diferite și atributele lor. Un atribut asociat unui element și care exprimă o opinie pozitivă poate exprima o opinie negativă pentru chiar același element dar al unui alt obiect și în mod evident pentru elemente diferite.

Dacă avem la dispoziție ambele resurse, atât descompunerea în elemente pentru fiecare obiect, cât și lista de proprietăți pozitive și negative pentru fiecare element, atunci putem face o aproximare destul de bună a orientării opiniei din documentul analizat.

Aceste resurse pot fi păstrate separat, ca două resurse separate, sau pot fi păstrate împreună, având la dispoziție astfel o resursă completă pentru analiza sentimentelor. Această resursă va fi disponibilă pe site.

## Concluzii parțiale

Am discutat mai sus despre analiza sentimentelor bazate pe sentimentele asociate fiecărei componente sau subcomponente ale obiectului analizat. Mai este de făcut o singură observație: în momentul în care facem această analiză, este foarte important să ne uităm și la zona geografică, deoarece se pot observa diferențe de gusturi. Cel mai elocvent exemplu este vânzarea de telefoane iPhone 6 și iPhone 6 plus: dacă în Statele Unite și în Europa iPhone 6 a stat de aproape 6 ori mai bine decât iPhone 6 plus, însă nu același lucru se poate spune și despre Asia, unde iPhone 6 plus a fost un succes.

Motivul este foarte simplu: dacă în Statele Unite și în Europa review-urile au fost oarecum spre negativ datorită ecranului foarte mare comparativ cu toate generațiile precedente de iPhone, în Asia, acest lucru a fost mai degrabă apreciat.

Dacă ar fi să facem această analiză la nivel mondial și să coroborăm rezultatele cu cele de vânzări, atunci am putea concluziona că iPhone 6 este net superior față de iPhone 6 plus, iar ca principal argument ar fi dimensiunea foarte incomodă a telefonului și faptul că este mai degrabă o tabletă decât un telefon și poate în unele situații inclusiv prețul poate fi un argument. Spre exemplu, dacă un telefon iPhone valorează în Statele Unite suma de X USD, atunci în România acesta va valora X EUR, iar în UK X GBP și având în vedere diferența de curs valutar, atunci diferența de preț este chiar una substanțială.

Putem concluziona că analiza sentimentelor la nivel mondial poate avea rezultate eronate sau poate chiar controversate. Ca să-l cităm pe Mark Twain: "There are three kinds of lies: lies, damned lies and statistics" ceea ce face foarte grea analiza sentimentelor la nivel mondial fără să ne uităm și la zona geografică.

## Sumarizarea

O dată cu dezvoltarea rapidă a Internet-ului, din ce în ce mai multă informație este disponibilă on-line. Această explozie de date a evoluat însă într-o supraîncărcare cu informații. Nu există timp suficient pentru a putea parcurge toate materialele disponibile și de regulă deciziile critice trebuie făcute pe baza datelor disponibile la acel moment de timp. Ca și răspuns la această problemă, numărul de cercetători interesați de sumarizare a crescut considerabil.

O dată cu creșterea acestui interes, se poate observa de asemenea și o creștere a bugetului dedicat cercetării acestui domeniu. Statele Unite (DARPA - Defense Advanced Research Projects Agency), Uniunea Europeană și țările din zona Pacificului au identificat sumarizarea textului ca fiind un domeniu de cercetare critic și încep să investească în dezvoltarea sa. De asemenea, se poate remarca un interes și din partea sectorului privat, în special industria telecom care a dezvoltat o serie de tehnici în acest sens: Prosum - BT - data mining pentru explorarea unor colecții mari de date, Context - Oracle - filtre web pentru găsirea documentelor online, Inxight - AltaVista Discovery.

Așa cum am arătat și în introducere, nivelul comunicațiilor în interiorul unei companii a crescut considerabil în ultimii ani, iar șansele să crească și mai mult în următorii ani



sunt foarte mari. Pe lângă această creștere a volumului comunicațiilor, se poate observa un număr copleșitor și în creștere de potențiale surse de informare, majoritatea fiind în formă scrisă, ca de exemplu, rețelele de socializare.

Resursele și timpul limitat forțează managerii să-și limiteze documentarea la doar un mic subset din toată informația disponibilă. Acest lucru poate afecta performanța și calitatea deciziilor pe care aceștia le iau. De exemplu, în companiile de succes se investesc mai mulți bani în monitorizarea presei și a altor surse de documentare față de companiile cu performanțe mai slabe. Accesul la informație în timp util poate aduce multe beneficii pentru deciziile luate în cadrul unei companii.

Scopul unui rezumat este simplu și evident: de a facilita identificarea unui subiect din mai multe articole și selectarea materialelor cele mai relevante pentru subiectul de interes. În general, crearea unui rezumat implică o oarecare familiaritate cu subiectul dezbătut. Extragerea elementelor cheie pe care autorul articolului a dorit să le sublinieze necesită de regulă antrenament și experiență. În general, se face o distincție clară între scopul unui rezumat și tipul de rezumat folosit, însă din nefericire, terminologia în acest domeniu este încă insuficient dezvoltată pentru a permite acest lucru. O descriere mai generală împarte rezumatele în două mari subclase: rezumatele indicative, care indică subiectele dezbătute în documente și pot anunța cititorul/utilizatorul despre documentele care dezbate un anumit subiect și rezumatele informative care descriu subiectele dezbătute în lucrare.

Rezumatul manual este întotdeauna influențat de trecutul persoanei care îndeplinește această sarcină, de personalitatea acestuia și de dispoziția sa. Părerile personale ale persoanei care extrage rezumatul și interesele sale sunt de asemenea responsabile pentru subiectivismul rezumatului rezultat. Putem extrapola această situație și putem afirma că un articol poate avea mai multe rezumate în funcție de persoana care l-a extras și dispoziția acestuia la acel moment de timp.

Un alt factor ce influențează cercetarea în domeniul sumarizării automate este lipsa unei măsurători obiective a calității unui abstract sau a unui rezumat. Cu toate acestea, Morris et al afirmă că o metodă destul de obiectivă poate fi măsurarea înțelegerii textului inițial comparat cu varianta rezumată a acestuia.

Având în vedere interesul pentru acest domeniu, la momentul actual se dorește eliminarea subiectivității umane și găsirea de algoritmi performanți și capabili să efectueze această sarcină riguros și consecvent.

Imediat după ce s-a început cercetarea în domeniul traducerii automate, s-a remarcat un interes crescând pentru posibilitatea extragerii automate de rezumate ale documentelor. Accentul s-a pus mai mult pe rezumatele indicative, adică pe rezumatele ce îi conferea cercetătorului un indiciu referitor la ce documente sunt mai importante decât altele pentru un anumit subiect decât pe rezumatele informative, scopul fiind acela de a oferi un rezumat ce poate substitui citirea unui document.

O altă ipoteză investigată a fost că extrasul unui document (selectarea propozițiilor cheie din document) ce ar putea fi considerat de asemenea rezumatul documentului. Cu toate ca această ipoteză a fost recunoscută o perioadă îndelungată de timp drept extract automat, în mod informal, este denumit tot rezumat automat.

Orice cercetare în domeniul sumarizării automate trebuie să înceapă cu abordările clasice ale lui Luhn și Edmundson, trebuie să parcurgă abordările bazate pe corpus/ținut și să exploateze structura discursurilor, să treacă prin abordările ce utilizează o bază de date de cunoștințe, să exploreze diverse metode de evaluare și în mod evident să sublinieze problemele nerezolvate.

Metodele originale ale lui Luhn atribuiau ponderi fiecărei propoziții din document în funcție de diferite criterii statistice, în special fiind vorba despre o funcție legată de frecvența cuvintelor ce apar într-o propoziție. Cuvintele foarte comune sunt eliminate (de, pe, care, unde, sus...etc) cu toate că au cea mai mare frecvență atât în document cât și în fiecare propoziție. Eliminarea se face pe baza unui dicționar iar pentru restul cuvintelor au asignate ponderi în funcție de frecvența acestora, stabilindu-se astfel o funcție ce exprimă înțelesul (subiectul, semnificația) propozițiilor.

Extrasele ce foloseau aceste metode au oferit rezultate suficient de bune încât să încurajeze aprofundarea cercetării în acest domeniu. În același timp, o metodă pur statistică pentru determinarea rezumatelor unui document s-a dovedit neadecvată și implicit, au apărut și alte metode de-a lungul timpului.

În decursul ultimilor ani au fost prezentate o multitudine de abordări pentru sumarizarea automată, însă acestea se pot împărți în trei mari categorii, în funcție de nivelul de procesare: abordări de suprafață, abordări la nivel de entitate și abordări la nivel de discurs.

Abordările la nivel de suprafață au tendința să reprezinte informația printr-un set de euristici ce vizează anumite caracteristici. Aceste caracteristici sunt apoi combinate cu ajutorul unui grup de funcții într-un set finit de variabile predefinite din care se pot extrage anumite concluzii mult mai ușor.

Aceste euristici pot urmări:

- termenii tematici - prezența unor termeni importanți determinați pe bază statistică
- localizarea - poziția în text, poziția în paragraf, etc
- fundalul - prezența unor termeni din titlu și subtitlu în interiorul documentului
- cuvintele de evidențiere sau fraze de evidențiere: în concluzie, cercetarea noastră a evidențiat, important, bonus, în particular, etc

Abordările la nivel de entitate contruiesc o reprezentare internă pentru text, modelând entitățile din text și legăturile dintre acestea. Aceste abordări urmăresc să găsească tipare între legăturile dintre entități și se folosesc de teoria grafurilor pentru a determina ceea ce este important. Relațiile dintre entități includ:

- similaritatea - suprapunerea de dicționar
- proximitatea - distanța dintre unitățile de text
- înrudirea - cuvinte legate prin apariția acestora în expresii comune
- legăturile de dicționar - sinonimie, hipernimie, etc
- coreferențierea - propozițiile care fac o referențiere la un anumit subiect
- relațiile logice - aprobarea, contradicția, consistența, etc
- relațiile sintactice

- relațiile bazate pe înțeles

Abordările la nivel de discurs modelează întreaga structură a textului în relație cu informația transmisă. Aceasta structură poate include:

- formatul documentului
- legaturile dintre subiectele de discuție pe măsură ce sunt descoperite
- structura retorică a textului

Crearea unui rezumat poate fi descompusă în 3 părți principale: analizarea textului primit la intrare, transformarea acestuia într-o reprezentare pe scurt și sintetizarea acestuia într-un rezumat corespunzător.

Sumarizarea mai poate fi catalogată și în funcție de metoda folosită: (1) folosind procesarea limbajului natural pentru construcția unei hărți de cunoștințe a unui document pe baza căreia se poate genera un extract; și (2) utilizarea unui algoritm pentru extragerea celor mai importante propoziții și fraze din document. Procesarea limbajului natural necesită cunoștințe sintactice, semantice și pragmatice. Cunoștințele sintactice conțin regulile de structurare al unui limbaj iar cunoștințele semantice și pragmatice sunt specifice domeniului analizat. Din nefericire, utilizarea singulară a sintacticii poate produce uneori interpretări eronate datorită expresiilor combinate și a metaforelor.

Sumarizarea automată acoperă mai multe sub-probleme, ce vor fi dezbătute de asemenea în acest capitol. Pentru limba română, primele încercări de sumarizare automată datează din anii 70 prin lucrările lui Erika Nistor și Eliza Roman și utilizau algoritmul lui Luhn. Testele realizate de aceștia au constatat în procesarea a 75 de pagini de text, reprezentând 21 de articole sau capitole scurte din cărți din mai multe domenii (știință medicală, matematică, istorie și muzică).

Lucrările prezintă rezultatele experimentelor și discută dificultățile întâlnite în timpul experimentelor, în special datorate limbii române. Metoda lui Luhn aplicată pentru limba română implica la momentul respectiv împărțirea fiecărei propoziții în segmente despărțite de termeni de legătură (termeni comuni sau cuvinte-stop), la nu mai mult de 4 termeni ne semnificativi distanță. Luhn puncta fiecare segment utilizând pătratul numărului de elemente semnificative despartite de termenii de legătură. Acest tip de segmentare este puțin semantic comparat cu alte metode de segmentare a textului. Nistor și Roman au observat necesitatea unei liste de cuvinte stop, însă lipsa resurselor lingvistice în format electronic le-au creat destule probleme în anii 70.

În articolele ulterioare Roman și Nistor încearcă împărțirea textului în propoziții nucleu. Aceste transformări urmăresc recomandările lui Chomsky de a împărți frazele în propoziții nucleu. În timp ce Chomsky a găsit 9 categorii de propoziții nucleu, Nistor și Roman au încercat un studiu asemănător pentru limba română și au aplicat aceste rezultate pentru sumarizarea textelor în limba română, descoperind astfel 8 categorii. Lăsând la o parte aportul indubitabil al acestei cercetări, evaluarea calității acestor rezultate (atât teoretice cât și practice) cu privire la limba română rămâne la latitudinea lingviștilor români.

O dată ce textul a fost împărțit în propoziții nucleu, având 8 forme de baza stabilite de către autori, abstractul textului poate fi apoi contruit în mod automat. În mod evident, primele întrebări care apar sunt: care sunt cele mai eficiente cuvinte cheie? cele mai frecvente propoziții nucleu? care sunt cele mai frecvente cuvinte? Autorii ajung la concluzia că acele cuvinte cu cea mai mare frecvență sunt deasemenea și subiectele propozițiilor nucleu cu cea mai mare frecvență. Acest lucru poate fi tradus prin faptul că metoda lui Luhn poate fi combinată cu această transformare și se poate obține astfel un rezumat din cele mai frecvente propoziții nucleu.

În următoarele lucrări, cei doi autori își largesc noțiunea de rezumat, ceea ce reprezintă o realizare majoră pentru perioada respectivă. Astfel, Nistor și Roman (1979) încearcă să identifice acele propoziții dintr-un text care se referă la un subiect apriori specificat de către utilizator. Așa cum scriu și cei doi autori, două idei se pot evidenția: (1) utilizarea de mulțimi fuzzy și funcția acestora de a măsura similaritatea semantica dintre documente/text și cuvinte cheie și (2) utilizarea dicționarului disponibil sau construit special pentru acest scop.

Concluzia principală a acestui studiu este ca relația dintre un set de documente și rezumatele acestora nu mai este biunivocă. Poate exista un singur document și mai multe rezumate care indică către el. Acesta este unul dintre motivele pentru care autorii își denumesc metoda "rezumare dinamică". La momentul actual, sunt puține referințe în literatura de specialitate care se referă la sumarizarea în limba română.

Începem explorarea acestui domeniu prin prezentarea mai întâi a tehnicilor existente de sumarizare, și inclusiv a diverselor tehnici de preprocesare a textelor a priori sumarizării și vom încheia cu câteva aplicații practice pentru limba română. Este destul de greu să ne imaginăm viața de zi cu zi fără sumarizare. Titlurile știrilor sunt de regulă un rezumat al întâmplărilor descrise în conținut. De multe ori, prima frază este deja suficientă pentru a ne face o idee despre întregul document. În cazul articolelor științifice, abstractul sau rezumatul care se completează la început prezintă foarte pe scurt subiectul dezbătut și concluziile care se vor regăsi la final.

## Segmentarea

Segmentarea textului reprezintă procesul de împărțire a textului sursă în paragrafe și propoziții. Vom numi un document procesat în acest fel "DocSent" provenit din cuvintele englezești *document* și *sentence*. Această împărțire în paragrafe și propoziții este realizată printr-o parcurgere secvențială a textului, când se încearcă detectarea semnelor de punctuație care pot arăta sfârșitul unei propoziții. Acestea pot fi semne simple, precum punctul (.), semnele întrebării (?) și exclamării (!), sau semne compuse, precum elipsa (...) sau combinații de semne ale întrebării și exclamării (?!, !?). O problemă care poate apărea la acest nivel o reprezintă abrevierile, care de obicei se termină cu punct. În cazul în care algoritmul nu este pregătit pentru această situație, el va considera că propoziția se sfârșește după abreviere (Mr.) sau chiar în interiorul ei (U.S.A.). Pentru a depăși acest inconvenient, se folosește un corpus de abrevieri disponibil on-line, care a servit la verificarea existenței acestora în text. Astfel, ele vor fi substituite astfel încât să nu existe confuzii cu privire la sfârșitul real al propozițiilor. După efectuarea împărțirii în propoziții, abrevierile înlocuite anterior vor fi readuse la forma inițială, textul fiind păstrat în acest fel nealterat.

## WebSumm

WebSumm funcționează în două moduri: (1) ca un instrument de filtrare, ordonare și navigare printre titlurile și subtitlurile textelor întoarse de motoarele de căutare, și (2) ca un instrument de navigare în adâncime ce folosește extragerea de text și sumarizarea pentru a analiza informația din setul de documente.

Algoritmul descris de Mani și Bloedorn poate fi împărțit în trei pași:

Mai întâi se identifică toate asemănările și diferențele între toate documentele implicate, acțiune ce se poate efectua utilizând o serie multiplă de algoritmi. Unul din avantajele determinării similarităților și diferențelor dintre aceste documente, este faptul că similaritățile ne arată ce este comun în fiecare document corelat cu subiectul dezbătut în toate documentele, diferențele ne arată ceea ce este unic în fiecare document. Dacă presupunem ca documentele sunt ordonate în ordine cronologică, diferențele din ultimul document pot indica ce este nou într-un set de documente ce dezbate același subiect. În mod evident, atât similaritățile cât și diferențele pot ajuta foarte mult la construirea unui rezumat.

Se poate apoi determina care sunt cei mai importanți termeni în cazul unei căutări, acest lucru pornind de la ipoteza că localizarea și ordinea liniară au un rol foarte mare în determinarea termenilor importanți. Această sarcină poate fi îndeplinită utilizând un algoritm de ierarhizare (ranking) pe paragrafe. Cu toate acestea, alegerea dimensiunii unei ferestre cu scopul alegerii paragrafelor potrivite poate fi o problemă, chiar dacă sunt folosite ferestre de dimensiune fixă sau ferestre cu dimensiune variabilă în funcție de discurs. Dacă dimensiunea ferestrei este prea mică, atunci observatorul ar putea obține o mulțime mai mare de ferestre mici care oferă puțină informație în sine, dar împreună pot fi foarte relevante. Dacă însă dimensiunea ferestrei este prea mare, atunci e posibil să conțină prea multă informație irelevantă. În loc să se folosească ferestre de dimensiune fixă sau să se utilizeze un efort mult prea mare pentru a utiliza ferestre de dimensiune dinamică, Mani și Bloedorn recomandă utilizarea unui algoritm ce conferă diferite ponderi elementelor din text (unde termenii respectivi corespund nodurilor din graf). De asemenea, identificarea acestor regiuni din documente relevante pentru căutarea efectuată, reduce foarte mult spațiul în care se caută asemănări sau diferențe.

Se explorează de asemenea un model de reprezentare a textului care ia în considerare și gradul de conectare dintre termeni. Această reprezentare explorează un model de conectivitate unde puterea legăturilor (în general bazată pe similaritate) între diferite bucăți de text este utilizată pentru identificarea caracteristicilor importante într-unul sau mai multe documente.

## Sumarizarea prin extragerea frazelor cheie

Probabil cea mai întâlnită metodă de sumarizare este prin extragerea celor mai importante paragrafe sau fraze din text și așezarea acestora într-o anumită ordine. Voi prezenta în continuare câțiva din acești algoritmi și voi discuta eficacitatea, avantajele și dezavantajele acestora. În mod normal, sumarizarea unuia sau a mai multor documente conține mai mulți pași:

- Înțelegerea conținutului documentului sau a documentelor
- Identificarea celor mai importante componente ale articolului sau ale articolelor
- Scrierea sumarului

Având în vedere varietatea informației, este mult mai avantajoasă implementarea unei tehnici de sumarizare independentă de domeniul cu care se lucrează, însă automatizarea pașilor 1 și 3 este la momentul actual într-un stadiu incipient. De aceea, cel mai frecvent tip de sumarizare la momentul actual este extrasul. Extrasele sunt în general create cu ajutorul euristicilor bazate pe o analiză statistică a aparițiilor unor anumite cuvinte, dorindu-se astfel selectarea acelor fraze sau propoziții care au cel mai mare potențial pentru rezumarea documentului. De-a lungul timpului, au fost prezentate diverse metode pentru selectarea și extragerea celor mai concludente propoziții.

## Summarist

Summarist este una din încercările ce a dorit dezvoltarea unui sistem capabil să extragă rezumatul unui document, urmărind definiția: sumarizare = identificarea subiectului/topic-ului, interpretare și generare. Trebuie să dezbaterem fiecare punct în parte, deoarece conțin o multitudine de subalgoritmi, fiecare antrenați pe un corpus mare de date de antrenare.

Înainte de crearea oricărui rezumat, trebuie să determinăm care este subiectul central dezbătut în document. Apoi acest subiect/topic trebuie interpretat și apoi se poate genera rezumatul textului.

- *identificarea* - scopul acestui pas este eliminarea subiectelor inutile pentru crearea unui rezumat al textului. De regulă, se pleacă de la presupunerea că orice text dezbate mai multe subiecte, ceea ce introduce două ipoteze: (1) utilizatorul va dori extragerea celor mai importante  $N$  subiecte din document sau (2) extragerea celor mai importante subiecte (pentru el) din text.
- *interpretarea* - o dată ce subiectele principale au fost identificate, acestea pot fi listate pentru prezentarea unui extract. Cu toate acestea, în sumarele extrase manual, este nevoie și de un pas de interpretare cu scopul obținerii a unei compactări mai mari. Într-unul din studiile sale, Daniel Marcu a observat că lucrându-se cu un corpus de 10 documente extrase din ziare și cu 14 analiști, se poate obține un factor de compresie de 2.76. Deasemenea, a observat că un extract este de regulă de trei ori mai lung decât un sumar extras manual, calculul făcându-se pe numărul de cuvinte folosit. Aceste rezultate indică nevoia unei sumarizări ulterioare extragerii informațiilor importante, cu scopul eliminării informațiilor care se repetă, reformularea propozițiilor a.i. să conțină mai multă informație condensată, și de asemenea, generalizarea rezumatelor.
- *generalizarea* - are ca scop reformularea extracturilor într-un text coerent, dens și scurt. Dacă se sare această fază, atunci rezumatul va fi compus din simpla extragere a unor fraze cheie din textul inițial.

Chiar dacă un rezumat extras manual nu se va compara niciodată cu un extras automat, în mod special datorându-se faptului că pentru crearea unui rezumat corect este totuși nevoie de înțelegerea adevăratului sens al documentului analizat, nu se poate totuși contesta necesitatea rezumatelor automate.

Probabil una dintre cele mai frecvente aplicații pentru sumarizarea automată sunt aplicațiile de extragere a informației. Începând cu 1950, mulți cercetători și-au îndreptat atenția spre acest domeniu, încercând să categorisească textele pe categorii, și să extragă esența documentelor analizate.

Abordările de bază ale metodelor de extragere a informațiilor (information retrieval) au totuși un set de limitări. De multe ori, orice depășea nivelul cuvintelor era privit cu neîncredere. De aceea, pasul doi din summarist este atât de important: fără înțelegerea deplină a subiectelor dezbătute și posibilitatea reformulării lor, rezumatele nu vor fi mai mult decât propoziții din text extrase și concatenate alături de seturi de cuvinte care se repetă cu o frecvență mai mare.

Cu toate că metodele la nivel de cuvânt au fost destul de bine analizate, există în continuare câteva aspecte care împiedică aceste metode de a fi perfect:

- *Sinonimie* - un concept poate fi reprezentat prin mai multe cuvinte
- *polisemie* - un cuvânt poate avea mai multe înțelesuri
- *expresii* - o expresie poate avea un alt înțeles decât cuvintele din care este compusă
- *dependența dintre termeni* - majoritatea tehnicilor probabilistice privesc cuvintele ca fiind independente, fără a lua în considerare legăturile dintre ele

De exemplu, Summarist nu calculează doar frecvența cuvintelor (o tehnică frecventă în IR (information retrieval) ci calculează și frecvența conceptelor (utilizând WordNet și alte resurse) a.i. poate opera la un nivel mai adânc decât cel al cuvintelor. La momentul actual, partea de interpretare este încă la un nivel incipient, deoarece majoritatea atenției a fost îndreptată spre extragerea conceptelor.

Fiecare modul este conceput să lucreze independent, deși este posibil ca unele module să aștepte informații de la altele. În faza de preprocesare, fiecare cuvânt este scris pe câte o linie. Apoi, fiecare modul deschide fișierul de preprocesare și adaugă ce informație este considerată ca fiind relevantă pentru fiecare cuvânt, fiind vorba de multe ori de valori numerice, însă poate fi vorba și de cuvinte (stringuri). Fișierul poate fi deschis și inspectat la orice moment de timp. La sfârșitul fiecărei faze, un modul de integrare combină scorurile precedente într-un scor final, ce se adaugă de asemenea în același fișier.

## Interpretarea Subiectului (Topic-ului)

Al doilea pas în sumarizare este interpretarea subiectelor dezbătute în articol. În acest pas, unul, două, sau mai multe topici sunt combinate într-una (sau mai multe) concepte unificatoare. Acest pas poate fi unul simplu, ca de exemplu combinarea unor cuvinte ca "roată", "lanț", "pedală", "lumină", "șa" într-un singur concept/cuvânt "bicicletă". În mod evident, procesul nu este atât de simplu precum reiese din

exemplul de mai sus, însă nu este o sarcină imposibilă. Cu toate acestea, este unul dintre cele mai dificile probleme din cadrul sumarizării automate, motivul principal fiind faptul că această fuziune de concepte necesită cunoștințe suplimentare, care de regulă nu sunt incluse explicit în articol.

Există diverse metode ce pot fi folosite la acest pas, iar majoritatea implică un set de indicatori ce au ca și corespondent un concept mai general.

Pentru identificarea subiectelor (topic-urilor), se pornește în general de la premiza că cu cât un cuvânt este folosit mai frecvent în text, cu atât este mai important pentru textul analizat. Cu toate că această metodă oferă rezultate destul de bune, indiferent de domeniul analizat, nu se poate descurca însă cu sinonime, expresii complexe sau polisemie. Printr-o listă de concepte ce pot fi exprimate printr-o multitudine de cuvinte se pot depăși aceste dificultăți. Deasemenea, Summarist calculează frecvența conceptelor și nu doar a cuvintelor din text, folosindu-se pentru generalizare (determinarea conceptelor părinte) WordNet-ul. Dacă însă nu se găsesc informațiile necesare în WordNet, atunci se numără cuvintele.

Algoritmul este destul de simplu: se numără de câte ori apare fiecare cuvânt în text, iar numărul de apariții este adăugat la numărul de apariții al conceptului reprezentat prin acel cuvânt. Procesul se poate termina o dată ce s-au trecut prin toate cuvintele din document, sau în orice moment în care considerăm că un anumit concept este suficient de bine reprezentat. Acest algoritm va alege cele mai generale concepte prezente în articol.

Rezultate similare se pot obține și dacă se folosesc semnăturile de topic-uri.

## Generarea Sumarului

Ultimul pas în sumarizarea automată este generarea sumarului, ce poate fi rezolvat prin mai multe metode, de la simpla tipărire a propozițiilor importante din text, până la construirea de propoziții complexe ce includ cât mai multă informație din textul original.

Summarist conține 3 module de generare a rezumatului, în funcție de cerințele utilizatorului. În cazul în care este dorit un extract, se poate elimina complet partea de generare. Termenii și/sau propozițiile extrase în faza de extragere a topic-ului pot fi listați direct ca și rezumat al textului inițial. Cu toate că ar putea părea incoerent, acest extract ar putea conține suficiente informații pentru ca un utilizator să-si dea seama de relevanța acestui articol.

De multe ori, nu este nevoie de generarea unui rezumat. Simpla listare a subiectelor dezbătute în textul analizat poate fi suficientă pentru utilizator. Summarist este capabil să efectueze și această sarcină.

Summarist conține un generator rudimentar de rezumate, ce poate concatena propoziții și poate obține propoziții simple din elementele extrase. Urmărind legăturile între cuvinte și între concepte și fiind capabil să formeze propoziții enunțiative sau cauzale.



Cu toate că este departe de a fi perfect, Summarist se regăsește printre cele mai apropiate metode de sumarizare existente. Rezultatele întoarse pot varia de la "foarte bune" până la "acceptabile", în funcție de tipul rezumatului dorit și de pretențiile utilizatorului.

De asemenea, Summarist poate fi folosit ori pentru extragerea rezumatelor din texte, dar poate fi folosit și pentru diferite sarcini intermediare, precum extragerea cuvintelor cheie, identificarea subiectelor și determinarea conceptelor principale din text.

## Concluzii Parțiale

Având în vedere ca în capitolul 5 tratăm presa din România, o dată ce ne uităm suficient de mult timp în trecut, putem observa cum majoritatea publicațiilor au adoptat un stil tabloid. Dacă așa cum am văzut până acum, titlul putea fi de multe ori considerat ca un sumar al articolului, în ziua de azi, lucrurile nu mai stau așa.

Să luăm ca și exemplu următorul articol din ziarul Libertatea, publicat pe data de 22 Noiembrie 2014:

Titlu: DEZVĂLUIRI din INTIMITATE. CUM își ÎMPART TREBURILE CASNICE Klaus Iohannis și SOȚIA SA, o GOSPODINĂ DESĂVÂRȘITĂ

Conținut: "Marius Vecerdea, antrenorul de tenis al lui Klaus Iohannis, din 2009, a dezvăluit cum își împart treburile casnice președintele ales și soția sa, care se pare că este o gospodină desăvârșită.

Klaus Iohannis, președintele ales al României va avea sâmbătă seară ocazia să se bucure de rezultatul alegerilor prezidențiale, acasă, la Sibiu, unde va da o petrecere. Evenimentul va avea loc într-un cadru restrâns, președintelui Klaus Iohannis alăturându-i-se prieteni foarte apropiați cu care își petrece de obicei concediile, potrivit oradesibiu.ro.

"Îmi este imposibil să vă spun exact cum se numește ce gătește Prima Doamnă. Oricum, e genială în bucătăria franceză. Gătește o rață de te lingi pe degete. Apoi, e tare la dulciuri. Nicăieri n-am mâncat bezele și tarte cu fructe mai bune ca ale ei. Are mașină de înghețată acasă și, dacă vă vine să credeți, își coace singură chiar și pâinea. Totul e ecologic. Iar la final vine șocul. Iohannis spală vasele. Am rămas mască. M-am oferit să-l ajut, dar n-a vrut să audă. E treaba lui dintotdeauna", povestea Vecerdea."

Link către știre: <http://www.libertatea.ro/detalii/articol/klaus-iohannis-treburile-casnice-sotie-spalat-vase-515624.html>

Așa cum putem observa, există o oarecare legătură între titlul ales și conținutul articolului, însă de asemenea, putem observa cât de exagerat este și cât de mult denaturează situația. Având în vedere situația în care se afla presa în momentul de față (nu doar în România ci chiar la nivel internațional), după ce s-a trecut printr-o scădere considerabilă a bugetelor și multe publicații tipărite chiar s-au închis, de multe ori titlurile sunt denaturate a.î. să șocheze și să atragă cititori.

În cazul Summarist ce folosește OPP (optimal position policy) - statistic s-a observat că propoziția imediat următoare sub titlu este cea care are cea mai mare valoare pentru un rezumat al textului. Acest lucru se poate remarca și în articolul de mai sus, însă folosirea cuvântului "dezvăluiri" (care sunt prezentate apoi în fraza nr. 3) are în continuare ca scop ținerea cititorului în suspans.

Pe de altă parte, să luăm în continuare exemplul următor:

Titlu: Antivirusul românesc Bitdefender va proteja infrastructura IT a guvernului din Danemarca

Continut: Soluțiile de securitate Bitdefender au fost selectate de Agenția de Modernizare a Administrației Publice din Danemarca pentru protejarea infrastructurii IT a instituțiilor din cadrul guvernului danez, a anunțat producătorul român de soluții de securitate.

Urmare a selecției realizate de agenția daneză, specialiștii Bitdefender implementează soluțiile Bitdefender GravityZone Security for Virtualized Environments (SVE), GravityZone Security for Endpoints și Cloud Security for Endpoints (Small Office Security) în ministere și alte instituții publice din Danemarca.

"Înainte de a fi selectați de agenția daneză, aveam deja municipalități dar și companii private din Danemarca în portofoliul de clienți, ce aleseră soluțiile Bitdefender atât datorită securității de top oferite cât și datorită soluțiilor unice dedicate mediilor virtualizate, ce reduc consumul de resurse și costurile asociate mentenanței", a declarat directorul de vânzări pentru Danemarca și țările Scandinave în cadrul Bitdefender, Brian Ziegler, într-un comunicat.

Potrivit contractului, Bitdefender furnizează Guvernului Danez protecție împotriva virușilor, tentativelor de phishing, amerințărilor de tip ransomware și împotriva altor amenințări ce vizează terminale de tip PC/laptop indiferent de sistemul de operare, pentru toate tipurile de servere, precum și pentru terminale mobile precum Android și iPhone.

Pachetul cuprinde și o soluție pentru email precum și o soluție de scanare a traficului în rețea.

Securitatea infrastructurilor IT virtualizate de mari dimensiuni, aparținând statului danez e realizată cu soluția Bitdefender GravityZone Security for Virtualized Environments (SVE), soluție compatibilă cu toate tipurile de hipervizoare (tehnologie ce stă la baza virtualizării), și care rezolvă provocările tradiționale legate de mediul cloud, oferind avantaje semnificative în termeni de securitate, consum redus de resurse și un timp mai redus dedicat mentenanței.

Din consola centrală a GravityZone, instituțiile ce implementează soluția în cadrul Guvernului Danez pot gestiona toate dispozitivele mobile sau fixe și toate tipurile de clienți virtuali și servere.

Brandul Bitdefender a fost lansat în 2001, prin redenumirea antivirusului AVX, produs și comercializat de Softwin. În 2007, divizia Bitdefender s-a desprins de grupul de firme Softwin.

Sursa: <http://www.zf.ro/business-hi-tech/antivirusul-romanesc-bitdefender-va-proteja-infrastructura-it-a-guvernului-din-danemarca-13603207>

În continuare putem observa că primul paragraf este cel mai important și acest lucru se poate observa deoarece majoritatea articolelor respecta piramida inversă.

În teoria și practica jurnalistică, “piramida inversată” este o tehnică de redactare, aplicată în special știrilor, prin care informațiile esențiale sunt prezentate la începutul articolului. Celelalte detalii sunt redade descrescător, în ordinea importanței.

Piramida inversată este formată din:

- Lead – (prima frază/paragraful inițial, introducere, capul știrii) conține mesajul esențial și are rolul principal de a incita cititorul la lectură;
- Paragrafe de susținere – conțin date care explică și aprofundează introducerea. Ajută la situarea evenimentelor în context. Conțin o serie de date secundare care întregesc imaginea faptului. Background-ul își are locul tot în această zonă;
- Final – este aproape la fel de important ca lead-ul. El trebuie să fie memorabil, să fixeze în mintea cititorului informația esențială.

Adevăratul scop și potențial al piramidei inversate este captarea atenției publicului încă din primele secunde de lectură. Pe online studiile au demonstrat că publicul scanează informația. Nici pe print lucrurile nu stau foarte diferit.

În principiu, publicul n-are răbdare sau n-are timp să citească tot articolul ca să afle despre ce e vorba. De regulă, cititorul parcurge titlurile, fotografiile, explicațiile foto, infografiile, apoi începe să citească. Dacă prima frază nu i-a trezit interesul, sunt mari șanse să nu citească în continuare.

Acum dacă ne uităm la cele doua articole, unul cu tentă tabloidă iar celălalt într-o tentă mai serioasă fiind vorba de un comunicat de presă, prima și ultima propoziție poate sumariza destul de bine ideea de bază din articol.

## **Sistem Automat de Monitorizare Media**

Partea aplicativă a acestei lucrări este formată din două componente. Prima parte este formată dintr-un website ce va pune la dispoziție atât majoritatea aplicațiilor prezentate în această lucrare (atât cod sursă cât și executabil), cât și un set de resurse lingvistice ce pot fi folosite pentru dezvoltarea ulterioară de aplicații. Aceste resurse sunt disponibile gratuit și pot fi utilizate de către oricine dorește să continue cercetarea în acest domeniu sau are nevoie de resurse lingvistice pentru diverse acțiuni de procesare automată a limbajului natural.

Partea a doua este o aplicație ce dorește să demonstreze utilitatea acestor algoritmi, prin implementarea unui sistem automat de monitorizare media.

Deasemenea, partea de monitorizare media nu își dorește să fie un sistem de monitorizare media în adevăratul sens al cuvântului, ci să ofere câteva din funcționalitățile de bază într-un mod complet automatizat și gratuit. Serviciile de acest gen sunt oferite acum contra cost de o multitudine de firme (în general de către agențiile de PR), însă de cele mai multe ori selecția și analiza se face manual prin citirea individuală a fiecărui articol și împărțirea pe categoriile de interes - în general fiind vorba de o clasificare pe bază de cuvinte cheie. Deasemenea, la cerere, se pot crea rezumate, astfel încât clientul să nu fie nevoit să citească toată presa din ziua respectivă, ci doar un rezumat ce conține strict informațiile necesare.

Sistemul prezentat va extrage toate articolele din presa online (la momentul actual, peste 53 de publicații), le va grupa în funcție de subiectul dezbătut, le va ordona în funcție de popularitate, va extrage cuvintele cheie, va măsura pozitivitatea articolului și de asemenea va crea rezumate atât multidocument cât și single-document.

Sistemul automat de monitorizare media folosește un set extins de algoritmi, atât din punctul de vedere al algoritmilor de procesare a limbajului natural cât și al procesării în timp real al textelor de pe internet:

Pentru a oferi o imagine de ansamblu a ce se întâmplă în media, știrile trebuie extrase suficient de frecvent pentru a fi la curent cu tot ceea ce se întâmplă (fiind vorba de publicații online, se poate întâmpla ca un articol să fie completat pe parcurs, pe măsură ce apar informații noi), dar suficient de rar pentru a nu supra-aglomera sistemul. De asemenea, având în vedere că anumite știri vor fi acoperite de către toate publicațiile, trebuie utilizat un algoritm de clasificare semantică care să fie capabil să grupeze documentele în funcție de subiectul dezbătut în timp real. Un subiect de maxim interes va fi descris de majoritatea publicațiilor, așa că de multe ori se dorește așezarea știrilor într-o ordine cronologică, pentru a se putea determina ce publicație a fost prima care a publicat știrea respectivă și care a fost raportul gradual al celorlalte publicații de-a lungul investigării subiectului dezbătut.

Extragerea articolelor dintr-un set de publicații online, cu toate că este o sarcină aparent simplă, ascunde suficiente sarcini încât să nu fie un pas ușor de ignorat.

Varianta folosită în implementarea noastră conține următorii pași:

- abonarea la feed-urile RSS al tuturor acestor publicații
- parsarea feed-ului atunci când se primește ceva nou
- descărcarea articolelor și verificarea dacă este vorba despre un update sau un articol nou. Având în vedere că majoritatea ziarelor de pe internet sunt gratuite, în general fiecare pagină va conține cel puțin un spațiu publicitar ce trebuie eliminat din analiza textului. Acest lucru se face prin diverși algoritmi euristici și poate diferi de la ziar la ziar.
- inserarea articolelor în baza de date spre analiză

O dată ce avem articolele extrase și grupate în funcție de subiectul dezbătut, se poate determina dacă este un articol ce are o tentă pozitivă sau negativă. În mod evident, algoritmul de analiză a sentimentelor are o tentă generală aici, deoarece va lucra pe un set extrem de diferit de documente, de la știri ce prezintă decese sau dezastre naturale până la știri ce prezintă variații de preț la carburant sau la diferite produse.

Deasemenea, un alt subiect de interes este sumarizarea documentelor. Având în vedere că este vorba despre un sistem de monitorizare media, ne interesează în primul rând sumarizarea multi-document, astfel încât un utilizator al acestui sistem să nu fie nevoit să parcurgă toate articolele pentru a-și face o părere despre subiectul respectiv.

În continuare vom prezenta toate resursele disponibile pentru limba română.

## Framework NLP

Așa cum am precizat încă de la început, unul din scopurile principale ale lucrării a fost oferirea unui set de resurse sintactice și semantice pentru limba română, atât cod sursă și librării cât și fișiere de date.

- în primul rând se pot regasi implementări pentru toți algoritmi descriși în lucrare
- diferite resurse tehnice pentru pre-procesarea textelor scrise în limba română (e.g. transformarea unui text scris cu diacritice într-un text ușor digerabil de algoritmi dezvoltați în această lucrare
- întreaga arhitectura (inclusiv codul sursă) al aplicației de monitorizare media

Ca și resurse lingvistice, se pot găsi următoarele elemente:

- expresii de concluzionare - acestea se pot dovedi utile atât pentru algoritmi de sumarizare, cât și pentru algoritmi dedicați determinării subiectelor dezbătute în articolele analizate. Multe lucrări de specialitate se folosesc de aceste expresii de concluzionare pentru efectuarea diferitelor task-uri, însă la momentul actual există foarte puține astfel de resurse pentru limba română. De aici și inițiativa de a pregăti un set de resurse pentru oricine dorește să continue sau să îmbunătățească cercetarea în acest domeniu, în special pentru limba română. Exemple de expresii de concluzionare: "în concluzie", "pe scurt", "această lucrare", etc.
- acronime - în analiza unui text, este foarte important să se poată determina că este vorba despre același lucru atunci când se folosește ori titlul complet al unei expresii și acronimul său. De asemenea, sunt foarte puține resurse ce conțin acronime dedicate pentru România
- abrevieri - în cazul știrilor publicate în presă sau în articolele științifice nu este cazul, dar se poate întâmpla ca într-unul din textele analizate să fie nevoie de înțelegerea abrevierilor pentru a se putea lucra mai ușor.

Politica	Economie	Societate	Extern	Sport	Monden	Toate
<b>Principalele stiri ale zilei</b>						<b>Linkuri utile</b>
<b>Cele mai populare</b> Cele mai noi						
3 articole pe	<b>C. V. Tudor refuza evacuarea sediului central al PRM si sustine ca Basescu este in ...</b> "Romania Libera acum 7 ore 28 minute - Presedintele PRM, C. V. Tudor discuta, la sediul central al PRM, cu executorul judecatoresc venit sa aplice decizia de evacuare din sediu si cu avocatul partii care a obtinut hotararea In justitie, europarlamentarul declarand ca nu va parasii cladirea si ca vinovat de situatie este seful statului, po ... <a href="#">Mai multe detalii</a>					In continuare sunt prezentate cateva din linkurile ce au ajutat la dezvoltarea acestui website si de asemenea cateva resurse suplimentare <a href="#">Teza de doctorat</a>
3 articole pe	<b>Patru acuzati In dosarul Petrom Service, eliberati din arest.</b> "Romania Libera acum 16 ore 15 minute - Mihai Sorin, fost director general al Petrom Service, si alte trei persoane cu functii de conducere In PSV Company si In alte firme vor fi cercetati In libertate In dosarul In care sunt acuzati de spalare de bani, In urma deciziei de joi a Curtii de Apel Bucuresti, anunta Mediafax.Curtea de Apel Buc ... <a href="#">Mai multe detalii</a>					
2 articole pe	<b>Rompan: Pretul painii va creste cu peste 15% din luna februarie.</b> "Romania Libera acum 3 ore 27 minute - Painea se va scumpi cu peste 15% Incepanad cu luna februarie a anului viitor, din cauza importurilor de grau, dupa ce In ultimele patru luni pretul a crescut cu 10%, a declarat Aurel Popescu, presedintele patronatului Rompan citat de Mediafax."Cresterea de pret la grau de la 120 de euro pe tona In iu ... <a href="#">Mai multe detalii</a>					
2 articole pe	<b>UE interzice biberonele de plastic pe baza de bisfenol A.</b> "Romania Libera acum 4 ore 54 minute - tarile Uniunii Europene (UE) au hotarat joi sa interzica, din primavara lui 2011, productia, iar ulterior comercializarea, In Europa, a biberonelelor din plastic care contin bisfenol A (BPA), un compus chimic controversat, utilizat In fabricarea plasticului de uz alimentar, relateaza AFP, citata de M ... <a href="#">Mai multe detalii</a>					
2 articole pe	<b>Masacrul de la Katin a fost o "crima" ordonata de Iosif Stalin, afirma Duma rusa.</b> "Romania Libera acum 5 ore 5 minute - Camera inferioara a Parlamentului rus, Duma de stat, a adoptat vineri, In prima lectura, o declaratie care recunoaste masacrarea de catre NKVD a mii de ofiteri polonezi, In 1940, la Katin, drept o "crima" ordonata de Iosif Stalin, potrivit presei ruse, citata de AFP si de Mediafax."Documentele publi ... <a href="#">Mai multe detalii</a>					
2 articole pe	<b>Gigi Becali, atac verbal dur la adresa presedintelui FRF.</b> "Romania Libera acum 3 ore 27 minute - Finantatorul FC Steaua, Gigi Becali, l-a atacat verbal cu duritate pe presedintele FRF, Mircea Sandu, dupa ce acesta a afirmat ca patronii din Liga I sunt "circari" si "creditori", declarand despre seful federatiei, printre altele, ca este "vagabond", "nenorocit" si "nebun", informeaza Mediafax."Pai ... <a href="#">Mai multe detalii</a>					

In interfața de bază ofera următoarele funcționalități:

- Clustering semantic pe stiriile din ultimile șapte zile și ordonarea acestora după numărul de elemente din fiecare categorie
- listă cu linkuri utile despre algoritmi folosiți în acest proiect
- Posibilitatea selectării doar a tipurilor de știri de interes (e.g. Politică, Economie, etc)
- Pentru fiecare cluster, se pot afișa toate știrile din clusterul respectiv
- Butonul `\textit{Mai Mult}`, returnează prima știre din cluster, din punct de vedere cronologic

Se poate face de asemenea și o selecție după cele mai populare sau cele mai recente știri.

Pentru calculul distanței dintre două texte se pot folosi o multitudine de algoritmi, însă cel mai rapid s-a dovedit a fi o combinație între mai mulți algoritmi rapizi: distanța Jacard pe toate cuvintele documentului (din care excludem totuși cuvintele uzuale, distanța Jacard doar pentru cuvintele care încep cu majusculă, și un parametru euristic ce depinde de diferență de numărul de cuvinte.

O listă incipientă de cuvinte uzuale poate fi regăsită mai jos, însă pe site poate fi găsită o lista completa: *este cele din despre deci deodata când care celor ceea cei celalalt doamna domnul doamnei are tot domnului mi-aș multe multa mi-l mult n-a*

*asta chiar dacă dar după pentru în însuși doar aștia între cât îl însă într-o într-un între își și cine ce cum unde când a acest îmi prin acesta acestuia lui un unde pe sus jos acestea acestora aceștea număr numărul unu una primul prima de o unu doi trei patru cinci șase șapte opt nouă zece una doua treia patra cincea șasea șaptea opta noua zecea oarecare mie să astfel așa.*

Extragerea textului corect din presă a fost un efort în sine, deoarece trebuiesc eliminate toate tag-urile HTML, toate scripturile javascript și orice altceva mai poate influența textul final și implicit rezultatele finale. Pentru determinarea celor mai recente articole din presă m-am folosit de feed-urile RSS ale fiecărei publicații, iar următorul script extrage linkurile către toate articolele noi primind ca și date de intrare link-ul către feed-ul RSS.

Alte resurse dezvoltate și adaptate pentru limba română vor fi disponibile pe website și de asemenea tool-uri noi vor fi adăugate în mod frecvent. De asemenea, cele două script-uri prezentate mai sus rulează din februarie 2010, adunând o cantitate foarte mare de date.

## Siri

Există foarte multe probleme încă nerezolvate în procesarea limbajului natural, însă în același timp există și foarte multe probleme care au fost rezolvate. Un exemplu de succes este Siri, asistentul digital introdus de compania Apple pentru toate dispozitivele mai noi de iPhone 4s. Sistemul a fost initial dezvoltat de Nuance și ulterior cumpărat de Apple. Prima versiune de Siri era capabilă să facă rezervări la diferite restaurante, să cumpere bilete la film sau să comande un taxi la locația curentă a utilizatorului și totul doar prin oferirea de comenzi vocale în limbaj natural. Ultima versiune este însă capabilă de mult mai multe acțiuni: poate interacționa cu aplicațiile de vreme, mesagerie, email, calendar, contacte, muzica, notițe, browser și hărți. De asemenea, Siri funcționează în următoarele limbi: Engleză, Franceză, Germană, Japoneză, Italiană, Mandarină, Koreană și Cantoneză.

În mod evident, în spatele Siri există un conglomerat de tehnologii din sfera procesării limbajului natural. Dacă în primă fază este vorba doar de transformarea din voce în text, în pașii următori se aplică o mare parte din tehnologiile prezentate în această lucrare.

Ca și mod de funcționare, în momentul în care utilizatorul transmite o întrebare, Siri înregistrează cerința utilizatorului în format audio și o trimite la serverele sale, unde din fișier audio devine o întrebare în format text.

Pasul următor este înțelegerea sensului întrebării. De exemplu, știm foarte clar că orice informație poate fi prezentată în mai multe feluri. "Aș vrea un croissant", "Există cumva vreo brutărie prin apropiere?" sau "Mi-ar plăcea niște produse de patiserie franțuzească", toate converg în aceeași direcție. În mod normal, pasul următor ar fi fost împărțirea propoziției în părțile de vorbire corespunzătoare. Însă Siri nu funcționează așa. În loc să modeleze concepte lingvistice, sistemul folosit a fost modelarea obiectelor reale. De exemplu, la cererea "As vrea să văd un thriller", Siri ar identifica imediat "thriller" ca și gen de film și ar începe să caute într-o listă de filme, în loc să caute cum este subiectul conectat de verb.

Siri este capabilă să mapeze conținutul unei întrebări într-un domeniu de posibile acțiuni și apoi să o aleagă pe cea mai probabilă, bazat pe "înțelegerea ei" asupra legăturilor dintre concepte reale. De exemplu, Siri știe că un restaurant trebuie să aibe o adresă, o notă, un tip de bucătărie și o gamă de preț. De asemenea, inclusiv legat de partea de transformare din audio în text. Creatorii Siri au explicat asta printr-un exercițiu simplu de imaginație. Teoretic, ar trebui să ne imaginăm un angajat în cadrul unui hotel. O întrebare "closest coffee shop" poate suna la fel de ușor cu "closest call Felicia" dacă hotelul este foarte aglomerat și sunt foarte multe persoane care discută în același timp. Însă în același timp, având în vedere contextul prezentat, este mult mai probabil ca "closest" să se refere la un loc decât la o persoană și este de asemenea mult mai probabil ca angajatul să fie întrebat despre recomandări de cină decât rugat să contacteze o persoană. La fel funcționează și Siri, care a fost concepută mai degrabă ca un asistent personal decât o persoană cu care ar trebui să ai conversații complexe.

## Concluzii Parțiale

Așa cum am văzut de-a lungul tehnicilor explorate, putem observa o schimbare de paradigmă. De-a lungul zecilor de ani prețuți în studierea procesării limbajului natural, majoritatea tehnicilor dezvoltate au fost concepute pentru a compensa lipsa puterii de procesare de la momentul respectiv. Tehnici care păreau atunci evidente și foarte ușor de implementat nu puteau fi folosite cu adevărat în practică datorită complexității algoritmilor, a cantității uriașe de date necesare și a puterii de procesare. Însă în ziua de azi, având în vedere progresul uimitor în materie de hardware și al cloud computing-ului, putem observa cum tehnici foarte simple (poate chiar vechi) sau cel puțin suficient de intuitive pot fi folosite cu o rată mult mai mare de succes decât algoritmi foarte complecși.

În ziua de azi, posibilitățile tehnice sunt mult mai avansate, ceea ce permite utilizarea unor algoritmi atât intuitivi cât și foarte intenși din punct de vedere procesare și în același timp să obținem un timp de execuție foarte bun. Tehnologiile din spatele Siri sunt în general euristice și se bazează pe cantități foarte mari de date și pe un data-center dedicat. Însă în același timp, rezultatele sunt statistic întotdeauna 77 la sută corecte iar timpul de procesare este insesizabil.

În platforma de pe site, se găsesc atât implementări adaptate pentru limba română a algoritmilor existenți, cât și implementări pentru algoritmi mai simplii, cu rezultate mai bune, însă care necesită resurse hardware mai complexe.

De asemenea, o parte din o mare parte din algoritmi folosiți în platforma de sumarizare oferă rezultate atât de bune deoarece infrastructura pe care rulează este suficientă pentru cerințele proiectului.



## Concluzii finale, contribuții originale și direcții viitoare de cercetare

Așa cum am văzut de-a lungul tehnicilor explorate, putem observa o schimbare de paradigmă. De-a lungul zecilor de ani prețeuți în studierea procesării limbajului natural, majoritatea tehnicilor dezvoltate au fost concepute pentru a compensa lipsa puterii de procesare de la momentul respectiv. Tehnici care păreau atunci evidente și foarte ușor de implementat nu puteau fi folosite cu adevărat în practică datorită complexității algoritmilor, a cantității uriașe de date necesare și a puterii de procesare. Însă în ziua de azi, având în vedere progresul uimitor în materie de hardware și al cloud computing-ului, putem observa cum tehnici foarte simple (poate chiar și foarte vechi) sau cel puțin suficient de intuitive, pot fi folosite cu o rată mult mai mare de succes decât algoritmi foarte complecși.

În ziua de azi, posibilitățile tehnice sunt mult mai avansate, ceea ce permite utilizarea unor algoritmi atât intuitivi cât și foarte intenși din punct de vedere procesare și în același timp să obținem un timp de execuție foarte bun. Așa cum am văzut în capitolul 5, tehnologiile din spatele Siri sunt în general euristice și se bazează pe cantități foarte mari de date și pe un data-center dedicat. Însă în același timp, rezultatele sunt statistic întotdeauna 77 la sută corecte iar timpul de procesare este insesizabil.

În platforma de pe site, se găsesc atât implementări adaptate pentru limba română a algoritmilor existenți, cât și implementări pentru algoritmi mai simplii, cu rezultate mai bune, însă care necesită resurse hardware mai complexe.

De asemenea, o parte din o mare parte din algoritmi folosiți în platforma de sumarizare oferă rezultate atât de bune deoarece infrastructura pe care rulează este suficientă pentru cerințele proiectului.

### Clasificarea documentelor

Majoritatea tehnicilor prezentate în capitolul 2 sunt în general statistice, așa că sunt destul de puține îmbunătățiri care pot fi aduse algoritmilor existenți. Rezultatele sunt în general la fel de bune dacă corpusul de antrenare este la fel de bun și la fel de complex, ceea ce ne face să concluzionăm că dacă reușim să avem un corpus de antrenare exhaustiv, atunci implicit și rezultatele vor fi pe măsură.

De asemenea, așa cum se poate observa și din exemplele de mai sus, gradul de complexitate și necesarul de resurse este direct proporțional cu numărul de cuvinte din corpusul de antrenare și în mod implicit cu dimensiunea corpusului. De asemenea, în funcție de timpul la care se face referire sau în care se petrece acțiunea oferă un grad de variație destul de larg atât verbelor cât și celorlalte cuvinte prin variația diacriticelor, a genului și a persoanei. De exemplu, variația [acest, această, acesta, aceștia, acestea, acești, aceste], cu toate că reprezintă mereu același lucru (un substitut pentru un anumit obiect/persoană), variază în funcție de număr și sex.

În funcție de problema care se dorește a fi rezolvată, se pot efectua diferite operații de preprocesare. Așa cum am văzut în inclusiv rolul sintactic al fiecărui cuvânt poate avea un rol esențial în clasificarea documentelor.

Când vorbim de clasificarea semantică și a modului în care poate fi ea aplicată, trebuie să fim foarte atenți la problema pe care trebuie să o rezolvăm și abia pe urmă să decidem ce algoritm se potrivește mai bine, sau ce corpus de antrenare este necesar. De exemplu, majoritatea algoritmilor prezentați mai sus se bazează pe un model de învățare supervizată, care în mod evident duce la rezultate mult mai bune decât în cazul unei antrenări nesupervizate, însă la cantitatea din ziua de astăzi de documentație și de informații, problema pregătirii corpusului de antrenare se complică exponențial.

Să luăm următorul exemplu: dorim să creăm un sistem automat de monitorizare media, care este capabil să:

- Preia din presa online toate articolele cu o anumită frecvență prestabilită. Statistic, frecvența optimă este din oră în oră
- Cum multe ziare vor acoperi același subiect, sistemul va trebui să fie capabil să creeze clustere de articole care tratează același topic/categorie
- Majoritatea publicațiilor oferă o preclasificare a articolelor în categorii precum economic, politic, etc, însă acest lucru este insuficient în cadrul unui sistem de monitorizare media, deoarece se dorește o filtrare la nivel de topic și nu la nivel de categorie
- Pentru fiecare topic, sistemul este capabil să determine dacă este vorba despre un topic pozitiv sau negativ. Acest lucru este foarte important în cazul review-ului de produse, al politicii, etc
- De asemenea, pentru fiecare topic în parte, pe baza articolelor din ziare diferite care ating același subiect, să se creeze rezumatul

La volatilitatea știrilor din ziua de azi, crearea unui corpus de antrenare la nivel de categorie (politic, financiar, tabloid, etc) este un lucru trivial, însă nu se poate spune același lucru și despre clasificarea la nivel de topic.

O bună metodă de creare a unui corpus de antrenare la nivel de topic s-a dovedit a fi Similaritatea Jaccard, care cu toate că pare destul de simplă, s-a dovedit în testele mele, în special pe zona clasificării articolelor din presa online, a avea rezultate foarte bune.

Distanța Jaccard are rezultate în general foarte bune, dar în capitolul 5 vom arata diferite metode de îmbunătățire, prin adăugarea de semidistanțe Jaccard (numărul de apariții al fiecărui cuvânt, dimensiunea documentelor, etc), care ne vor duce gradul de acuratețe care în mod normal la distanța Jaccard clasică este undeva între 70-80 la sută. În cazul testelor mele, în gama 98-100 la sută. Însă și distanța Jaccard clasică poate fi folosită dacă dorim doar crearea unui corpus de antrenare.

## **Analiza Sentimentelor**

Așa cum am observat în capitolul 3, analiza sentimentelor bazate pe sentimentele asociate fiecărei componente sau subcomponente ale obiectului analizat. De asemenea, am observat că în momentul în care facem această analiză, este foarte important să ne uităm și la zona geografică, deoarece se pot observa diferențe de gusturi. Cel mai elocvent exemplu este vânzarea de telefoane iPhone 6 și iPhone 6

plus: dacă în Statele Unite și în Europa iPhone 6 a stat de aproape 6 ori mai bine decât iPhone 6 plus, însă nu același lucru se poate spune și despre Asia, unde iPhone 6 plus a fost un succes. S-a observat că ecranul mai mare este un avantaj pentru piața din Asia față de cea din Europa sau Statele Unite.

Dacă ar fi să facem această analiză la nivel mondial și să coroborăm rezultatele cu cele de vânzări, atunci am putea concluziona că iPhone 6 este net superior față de iPhone 6 plus, iar ca principal argument ar fi dimensiunea foarte incomodă a telefonului și faptul că este mai degrabă o tabletă decât un telefon și poate în unele situații inclusiv prețul poate fi un argument. Spre exemplu, dacă un telefon iPhone valorează în Statele Unite suma de X USD, atunci în România acesta va avea valoarea X EUR, iar în UK X GBP și având în vedere diferența de curs valutar, atunci diferența de preț este chiar una substanțială.

Putem concluziona că analiza sentimentelor la nivel mondial poate avea rezultate eronate sau poate chiar controversate. Ca să-l cităm pe Mark Twain: "There are three kinds of lies: lies, damned lies and statistics" ceea ce face foarte grea analiza sentimentelor la nivel mondial fără să ne uităm și la zona geografică.

## Sistem Automat de Monitorizare Media

În capitolul 5, putem observa o schimbare de paradigmă. De-a lungul zecilor de ani pretrecuți în studierea procesării limbajului natural, majoritatea tehnicilor dezvoltate au fost concepute pentru a compensa lipsa puterii de procesare de la momentul respectiv. Tehnici care păreau atunci evidente și foarte ușor de implementat nu puteau fi folosite cu adevărat în practică datorită complexității algoritmilor, a cantității uriașe de date necesare și a puterii de procesare. Însă în ziua de azi, având în vedere progresul uimitor în materie de hardware și al cloud computing-ului, putem observa cum tehnici foarte simple (poate chiar vechi) sau cel puțin suficient de intuitive pot fi folosite cu o rată mult mai mare de succes decât algoritmi foarte complecși.

În ziua de azi, posibilitățile tehnice sunt mult mai avansate, ceea ce permite utilizarea unor algoritmi atât intuitivi cât și foarte intenși din punct de vedere de procesare și în același timp să obținem un timp de execuție foarte bun. Tehnologiile din spatele Siri sunt în general euristice și se bazează pe cantități foarte mari de date și pe un data-center dedicat. Însă în același timp, rezultatele sunt statistic întotdeauna 77 la sută corecte iar timpul de procesare este insesizabil.

În platforma de pe site, se găsesc atât implementări adaptate pentru limba română a algoritmilor existenți, cât și implementări pentru algoritmi mai simplii, cu rezultate mai bune, însă care necesită resurse hardware mai complexe.

De asemenea, o parte din o mare parte din algoritmi folosiți în platforma de sumarizare oferă rezultate atât de bune deoarece infrastructura pe care rulează este suficientă pentru cerințele proiectului.

## Contribuții personale

În capitolul 5 am construit un sistem automat de monitorizare media specializat pentru limba română, care implementează următoarele funcții: topic clustering,

sentiment analysis și sumarizare. Toate aceste funcții sunt oferite gratuit utilizatorilor, pentru un număr nelimitat de conturi și un număr nelimitat de cuvinte cheie. Majoritatea sistemelor de monitorizare media sunt excesiv de scumpe și acest lucru este nejustificat pentru anul 2014. Cu toate că acest lucru va produce o multitudine de discuții în această zonă, trebuie luată o inițiativă.

De asemenea, pentru partea de clasificare semantică, există foarte multe probleme ce trebuie rezolvate. În articolele scrise despre clasificarea documentelor am rezolvat probleme concrete, în special din zona de antispam, care au dus mai departe la 4 patente și nenumărate articole despre acest subiect, în special în zona de clasificare documente și împărțirea pe categorii.

În zona de procesare a limbajului natural, în 2010, am creat doua start-up-uri, în care puteam oferi căutarea de restaurante în limbaj natural "aș vrea o ciorbă de burtă ieftină în drumul taberei" sau "vreau o ciorbă de fasole în apropiere" sau oferirea unui algoritm de măsurare a influenței blogerilor în funcție de cum se propagă informația între bloguri.

## Direcții Viitoare de Cercetare

Având în vedere resursele care vor fi puse la dispoziție pe site, se poate continua în orice direcție de cercetare în zona de NLP pentru limba română. În sistemul de monitorizare media dezvoltat, pe partea de clasificare semantică am obținut rezultate foarte bune folosind un corpus de antrenare generat cu distanța Jaccard a.î. să putem avea o învățare supervizată.

Ca și resurse lingvistice disponibile pe site, se pot găsi următoarele elemente:

- expresii de concluzionare - acestea se pot dovedi utile atât pentru algoritmi de sumarizare, cât și pentru algoritmi dedicați determinării subiectelor dezbătute în articolele analizate. Multe lucrări de specialitate se folosesc de aceste expresii de concluzionare pentru efectuarea diferitelor task-uri, însă la momentul actual există foarte puține astfel de resurse pentru limba română. De aici și inițiativa de a pregăti un set de resurse pentru oricine dorește să continue sau să îmbunătățească cercetarea în acest domeniu, în special pentru limba română. Exemple de expresii de concluzionare: "în concluzie", "pe scurt", "această lucrare", etc.
- acronime - în analiza unui text, este foarte important să se poată determina că este vorba despre același lucru atunci când se folosește ori titlul complet al unei expresii și acronimul său. De asemenea, sunt foarte puține resurse ce conțin acronime dedicate pentru România
- abrevieri - în cazul știrilor publicate în presă sau în articolele științifice nu este cazul, dar se poate întâmpla ca într-unul din textele analizate să fie nevoie de înțelegerea abrevierilor pentru a se putea lucra mai ușor.

De asemenea, așa cum am precizat încă de la început, unul din scopurile principale ale lucrării a fost oferirea unui set de resurse sintactice și semantice pentru limba română, atât cod sursă și librării cât și fișiere de date.

- în primul rând se pot regăsi implementări pentru toți algoritmi descriși în lucrare
- diferite resurse tehnice pentru pre-procesarea textelor scrise în limba română (e.g. transformarea unui text scris cu diacritice într-un text ușor digerabil de algoritmi dezvoltați în această lucrare)
- întreaga arhitectură (inclusiv codul sursă) al aplicației de monitorizare media

Așadar, se poate merge în mai multe direcții de cercetare, însă eu personal o voi continua în crearea unui asistent bazat pe comunicarea audio.

## Referințe

1. T. K. Landauer, P. W. Foltz and D. Laham, Discourse Processes, Pages 259-284, Introduction to Latent Semantic Analysis, 1998
2. G. Golub and W. Kahan, J. Siam, Calculating the Singular Values and Pseudo-Inverse of a Matrix, 1965
3. Scott Deerwester, Susan T. Dumais, George W. Furnas and Thomas K. Landauer, Richard Harshman, Indexing by Latent Semantic Analysis, 1990
4. Maciej Ceglowski, Aaron Coburn, and John Cuadrado, National Institute for Technology and Liberal Education Middlebury College, Semantic Search of Unstructured Data using Contextual Network Graphs, 2003
5. Gibson W, Ace Book, ISBN: 0-441-56959-5, Neuromancer, 1984
6. Tufis Dan, Filip Florin Gh, Limba Română în Societatea Informațională - Societatea Cunoașterii, 2002
7. Hristea Florentina, Editura Universității din București, Introducere în procesarea limbajului natural cu aplicații în Prolog, 2000
8. Chomsky Noam, Mouton and Co, Language Syntax, 1957
9. Claudiu Mihaila, Sumarizare Automată focalizată Temporală, 2008
10. Regina Barzilay and Michael Elhadad, Mathematics and Computer Science Dept, Ben-Gurion University, Using Lexical Chains for Text Summarization, 2001
11. Xiaojun Wan, Jianguo Xiao, Institute of Computer Science and Technology, Peking University, Beijing, Single Document Keyphrase Extraction Using Neighborhood Knowledge, 2008
12. Bruce Krulwich and Chad Burkey, Center for Strategic Technology Research Andersen Consulting LLP, Learning user information interests through the extraction of semantically significant phrases, 1996
13. Ken Barker and Nadia Cornacchia, School of Information and Technology Engineering, University of Ottawa, Noun Phrase Heads to Extract Document Keyphrases, 2000
14. Alberto Munoz, Compound Key Word Generation from Document Databases Using A Hierarchical Clustering ART Model, Intel, 1997
15. Rada Mihalcea and Paul Tarau, Department of Computer Science, University of North Texas, TextRank: Bringing Order into Texts, 2004
16. Andrew H. Morris, George M. Kasper, Dennis A. Adams, Advances in Automated Text Summarization, The effects and Limitations of Automated Text Condensing on Reading Comprehension Performance, 1989

17. Gerard Salton, Amit Singhal, Mandar Mitra Chirs Buckley, *Advances in Automated Text Summarization, Automatic Text Structuring and Summarization*, 1989
18. Andrew Goldberg, CS838-1 - *Advanced NLP, Automatic Summarization*, 2007
19. Eduard Hovy and Chin-Yew Lin, *Advances in Text Summarization, Automatic Text Summarization in Summarist*, 1999
20. George Petre, *SpamConference 2010, Facebook - Another Breach in the Wall*, 2010
21. Edgar GONZÁLEZ, Maria FUENTES, TALP Research Center, Spain, *A New Lexical Chain Algorithm Used for Automatic Summarization*, 1999
22. Catalin Cosoi, Madalin Vlad, Valentin Sgarciu, *DAAAM 2008, proceedings*, pg.0309-0310, indexat ISI Web of Knowledge, *Introducing Syntactic Features into a Bayesian Classifier*, 2008
23. Kyoung-Min Kim, Jin-Hyuk Hong, Sung-Bae Cho, *Information Processing and Management* 43 (2007) 225–236, *A semantic Bayesian network approach to retrieving information with intelligent conversational agents*, 2007
24. Sylvain Auroux, *Histoire Epistemologie Langage, La notion de linguistique generale*, 1988
25. Noam Chomsky, *Lectures on Government and Binding*, 1981
26. W. Gerrod Parrott, *Emotions in Social Psychology: Key Readings*, 2000
27. Nitin Indurkha, Fred J. Damerau, *Handbook of Natural Language Processing*, 2000
28. Bing Liu, *Sentiment Analysis and Opinion Mining, Sentiment Analysis and Opinion Mining*, 2012
29. L. Ku, Y. Liang and H. Chen, *Proceedings of AAAI-2006, Opinion extraction, summarization and tracking in news and blog corpora*, 2006
30. Li Zhuang, Feng Jing, Xiao-Yan Zhu, *Proceedings of the 15th ACM international conference on Information and knowledge management, Movie review mining and summarization*, 2006
31. Hatzivassiloglou, V. and K. McKeown., *Proceedings of Annual Meeting of the Association for Computational Linguistics, Predicting the semantic orientation of adjectives*, 2007
32. Wilson et al, *Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis*, 2007
33. Kanayama and Nasukawa, *Proc. of EMNLP'06*, pp. 355-363, *Fully Automatic Lexicon Expansion for Domain-oriented Sentiment Analysis*, 2006
34. Qiu et al, *Stud Health Technol Inform.*, *Understanding the psychological motives behind microblogging*, 2010
35. Inderjeet Mani, Mark T. Maybury, Cambridge, MA: The MIT Press., *Advances in Automatic Text Summarization*, 1999
36. Bo Pang, Lillian Lee, *Proceedings of ACL '04 Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics, A sentimental education: sentiment analysis using subjectivity summarization based on minimum cuts*, 2004
37. Gerard Salton, Amit Singha, *Proceedings of ACL '04 Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics, Information Processing & Management*, 1997
38. GChinatsu Aonet Mary Ellen Okurowskit James Gorklinskyt Bjornar Larsent, *SRA Internatlon, A Scalable Summarization System Using Robust NLP*, 2003

39. E Nistor, E Roman, *Revue Roumaine de Linguistique*, Attempts at Automatic Detection of Some Semantic Deviations, 1980
40. Luhn H. P., *IBM Journal of Research*, Automatic Creation of Literature Abstracts, 1959
41. Cremmins, Edward T., ISI Press, *The Art of Abstracting*, 1982
42. Inderjeet Mani, Eric Bloedorn, *AAAI '98/IAAI '98 Proceedings of the fifteenth national/tenth conference on Artificial intelligence/Innovative applications of artificial intelligence* Pages 820-826 , Machine learning of generic and user-focused summarization, 1982
43. Barker, K. and Cornacchia, N, *Proceedings of the 13th Biennial Conference of the Canadian Society on Computational Studies of Intelligence*, pages 40–52, Using noun phrase heads to extract document keyphrases, 2000
44. Daniel Marcu, *The MIT Press - ISBN 0-262-13372-5*, *The Theory and Practice of Discourse Parsing and Summarization*, 2000
45. Mark Twain, *The Chapters from the North American Review*, *Mark Twain's Own Autobiography*, 1898
46. AC Cosoi, MS Vlad, V Sgarciu, *US Patent 8,131,655*, 2012, Spam filtering using feature relevance assignment in neural networks, 2007
47. AC Cosoi, MS Vlad, V Sgarciu, *MIT Spam Conference*, *Asigning Relevancies to Individual Features for Large Patterns in Artmap Networks*, 2007
48. AC Cosoi, *Proc. MIT Spam Conf*, *A False Positive Safe Neural Network; The Followers of the Anatrium Waves*, 2008
49. AC Cosoi, MS Vlad, V Sgarciu, *AQTR 2008*, vol.III, *Proceedings*, pg.230-233, *indexat? ISI Web of Knowledge*, *On neural networks and the future of spam*, 2008
50. AC COSOI CM COSOI V SGARCIU B DUMITRU, *Proceedings of DAAAM 2009*, Viena, Austria, 25-28 nov, pp.103-104, ISSN 1726-9679, ISBN 978-3-901509-70-4, *indexat? ISI Web of Knowledge*, *The Romanian Social Fractal*, 2009
51. AC Cosoi, *Virus Bulletin*, *The medium or the message? Dealing with image spam*, 2006
52. Maria Corduneanu Carmen Maria Cosoi Catalin Alexandru Cosoi Madalin Vlad Valentin Sgarciu, *INFORMACIJSKA DRU?BA ? IS 2009*, *A HYBRID NEURAL NETWORK MODEL FOR SPAM DETECTION*, 2009
53. C.Cosoi; M.St.Vlad; V.Sgarciu, *Annals of DAAAM for 2007 & Proceedings of the 18th International DAAAM Symposium*, ISBN 3-901509-58-5, vol. 18, pp 617-618, october 2007, *indexat? ISI Web of Knowledge*, *Using a Neural Network to Determine a Cat's Emotional State by its Meow*, 2007
54. Catalin Cosoi, *USPTO 8695100*, *Systems and methods for electronic fraud prevention*, 2014
55. Catalin Cosoi, *USPTO 8572184*, *Systems and methods for dynamically integrating heterogeneous anti-spam filters*, 2013
56. Catalin Cosoi, *USPTO 8335383*, *Image filtering systems and methods*, 2012
57. Catalin Cosoi, Madalin Vlad, Valentin Sgarciu, *USPTO 8131655*, *Spam filtering using feature relevance assignment in neural networks*, 2012
58. Musat Claudiu, Catalin Cosoi, *USPTO 8010614*, *Systems and methods for generating signatures for electronic communication classification*, 2011
59. Catalin Cosoi, *USPTO 7751620*, *Image spam filtering systems and methods*, 2010

60. Nasui D.V; Cosoi A.C; Sgarciu V; Rancea I, Proceedings of DAAAM 2009, Viena, Austria, 25-28 nov, pp.1153-1154, ISSN 1726-9679, ISBN 978-3-901509-70-4, indexat? ISI Web of Knowledge, Using Wireless Monitoring for Market Research Interviewers, 2009
61. Gams M; Cosoi A.C; Corduneanu M; Kolbe M; Vlad M.S; Sgarciu V, Proceedings of DAAAM 2009, Viena, Austria, 25-28 nov, pp.1415-1416, ISSN 1726-9679, ISBN 978-3-901509-70-4, indexat? ISI Web of Knowledge, A New Method for Fingerprint Clasification, 2009
62. C.Cosoi, V.Sgârciu, M.St.Vlad, Annals of DAAAM for 2007 & Proceedings of the 18th International DAAAM Symposium, ISBN 3-901509-58-5, vol. 18, pp 615-616, october 2007, indexat ISI Web of Knowledge, Asigning Relevancies to Individual Features for Large Patterns in Artmap Networks, 2007
63. Cosoi A.C.; Sgarciu V.; Vlad M.S, DAAAM 2008, proceedings, pg.0307-0308, indexat? ISI Web of Knowledge, Build Your Antiphishing Technology in Just 5 Minutes, 2008
64. Cosoi AC, Vlad MS, Sgarciu V, DAAAM 2008, proceedings, pg.0311-0312, indexat? ISI Web of Knowledge, Heuristic Image Recognition, 2008
65. Cosoi C.A; Cosoi C.M, Sgarciu V; Dumitru B, Vlad M.S, Proceedings of DAAAM 2009, Viena, Austria, 25-28 nov, pp.105-106, ISSN 1726-9679, ISBN 978-3-901509-70-4, indexat ISI Web of Knowledge., Spam / Twitter, 2009
66. Cosoi, A. C.; Vlad M. S. & Sgarciu, V, DAAAM International Scientific Book 2009, VOL.8, pp. 377-384 CHAPTER 39, ISBN 978-3-901509-69-8, ISSN 1726-9687; aparitie 2010, EBSCO bibliographic database, Heuristic Image Recognition, 2009
67. C.Cosoi; M.S.Vlad; V.Sgarciu , preprints, vol.3, CSCS 16, Bucharest, 22-25 May 2007, pg.33-37, ISBN 978-973-718-741-3 sau ISBN 978-973-718-744-4, Image Retrieval System Inspired from Antispam Filters, 2007
68. A.C.Cosoi; B.Dumitru; V.Sgarciu; M.St.Vlad; M.Corduneanu, IAFA 2009, vol.3, pp 36-40, ISSN: 2066-4451, The Romanian Blogosphere as a Social Fractal, 2009