

Recunoașterea umorului în texte

Țifrea Oana
Facultatea de Informatică,
Universitatea Alexandru Ioan Cuza,
Iași,
coordonator: Dan Cristea,

29 iunie 2008

Cuprins

1	Introducere	5
1.1	O privire de ansamblu asupra umorului	6
1.2	Tipuri de umor	8
1.2.1	Gramatica umorului	10
1.3	Teoriile umorului	11
1.4	Criterii pentru determinarea succesului umorului	13
1.5	Umorul din punct de vedere social	14
2	Cercetările în domeniul umorului computațional	16
2.1	Generare de umor	16
2.1.1	LBJOJG	16
2.1.2	Jape	17
2.1.3	HAHAcronym	17
2.1.4	WISCRAIC	17
2.1.5	MSG	18
2.1.6	Jester	18
2.2	Recunoașterea umorului	18
2.2.1	Recunoașterea umorului fără înțelegerea sensului	18
2.2.2	Aplicarea clasificatorilor de texte și a unor attribute pentru recunoașterea umorului	21
3	Experimente de identificare a umorului prin implementarea algoritmiilor clasici de clasificare de texte	26
3.1	Clasificarea de texte	26
3.2	Metode de reprezentare a documentelor	28
3.2.1	N-gramele	29
3.3	Mașinile cu vector suport (Support Vector Machine SVM)	30
3.4	Clasificatorul Bayes Naiv	34

4	Experimente în detectarea automată a umorului și rezultatele obținute	36
4.1	Corpusurile	36
4.2	Experimentele	41
5	Interpretarea rezultatelor	46
5.1	Observații privind rezultatele obținute	46
5.2	Posibilități de îmbunătății programul de recunoaștere al umorului în texte	46
5.3	Concluzii	47

Listă de figuri

3.1	Maparea pe un alt spațiu	31
3.2	SVM-hiperplanul și vectorii suport	33
4.1	Numărul de cuvinte din componența glumelor	37
4.2	Numărul cuvintelor pe categorii din componența textelor non-umoristice Corpus 1	38
4.3	Distribuția numărului de cuvinte pentru datele nonumoristice după a doua filtrare	40
4.4	Experimentul I	42
4.5	Experimentul II	44

Listă de tabele

1.1	Scara nivelurilor de de atașament	14
2.1	Recunoașterea umorului prin caracteristici	21
2.2	Rezultatele obținute folosind euristicile stilistice	25
2.3	Rezultatele obținute folosind algoritmi de învățare automată .	25
4.1	Structura Datelor după prima filtrare	37
4.2	Media numărului de cuvinte pentru datele nonumoristice . . .	39
4.3	Numărul de componente al vectorilor	42
4.4	Numărul de componente al vectorilor	44

Capitolul 1

Introducere

”Raționamentul corect se naște din experiență. Experiența se naște din raționament eronat. Concluzie: Raționamentul corect se naște din raționament eronat.”(anonim)

La înțelegerea acestor cuvinte suntem contrariați de absurditatea logicii și izbucnim în râs.

Umorul are un caracter specific (literar, subtil sau fin) prin soluții neașteptate, caraghioase care pot produce ilaritate. Persoanele cu umor sunt acele persoane care prin comportare sau prin vorbe, în anumite contexte, declanșează râsul.

Simțul umorului este influențat de tradițiile, cultura, istoria unui popor, sau diferă după poziția pe scara ierarhiei sociale sau după etate. Nu numai că variază de la o persoană la alta, dar se poate ca aceeași persoană să găsească o glumă ca fiind amuzantă într-o zi și în altă zi nu, depinzând de starea de spirit a persoanei, de evenimentele recent petrecute în viața persoanei respective.

Umorul diferă de asemenea după anumite perioade istorice, multe glume din trecut nemaifiind actuale deoarece a dispărut contextul care permitea perceperea lor ca având un anumit haz. Umorul poate fi usturător prin satiră, ironie, batjocură; cinic, sau blând, binevoitor, plin de înțelegere, autocritic.

Lucrarea de față își propune să depisteze umorul fără a înțelege sensul mesajului. Se vor încerca algoritmi clasici de clasificare de texte și diferite atribute ale textelor umoristice și neumoristice determinate euristic.

Umorul computațional este un domeniu în care există unele abordări de găsimă a unui șablon universal pentru generarea și recunoașterea umorului în texte.

Experimentele testate în această lucrare cu privire la recunoașterea umorului în texte sunt comparabile cu cele existente în domeniu.

1.1 O privire de ansamblu asupra umorului

Comunicarea om-calculator nu mai constituie demult un deziderat iluzoriu al inteligenței artificiale. Pentru ca această comunicare să fie una cât mai apropiată de comunicarea interumană, calculatorul (va trebui nu numai să recunoască, ci și să folosească și umorul. Mai mult, umorul oferă profunzimi ale limbajului uman- referindu-ne la cel real, complex, un limbaj creativ nu doar o mulțime de propoziții standard. Reușind să modelăm înțelegerea și generarea umorului de către calculatoare, câștigăm o mai bună imagine asupra modului în care creierul uman funcționează nu doar în privința umorului ci a limbajului și cunoșterii în general. Sunt multe situații în interacțiunea interumană unde umorul joacă un rol important permițând continuitatea conversației, întărind relațiile interumane. Urmați de paradigma CASA (Computers are Social Actors)(Calculatoarele sunt actori sociali) ne putem aștepta ca umorul să joace un rol similar și în interacțiunea om-calculator.

În accepțiunea populară, calculatoarele nu vor putea să folosească și să aprecieze umorul. Calculatoarele ficționale și roboții au fost mereu imaginați ca fiind fără de umor chiar dacă pot utiliza limbajul natural. Dar și șahul a fost odată considerat ca fiind un domeniu al oamenilor iar acum calculatoarele joacă la nivel grandmaster.

Agenții sociali și inteligenți au devenit o paradigmă pentru rezolvarea și descrierea problemelor în stilul oamenilor. Cercetările în privința acestor agenți includ capacitatea de percepție a dorințelor, credințelor și a intențiilor. Dar precum spune și Roddz Cowie, *Dacă vom arăta emoții la acești agenți cu siguranță ne vom aștepta ca ei să aibă și un pic de simț al umorului* [Binsted et al., 2006].

Pentru a înțelege umorul trebuie să-l plasăm în mediul său natural, care este societatea, trebuie mai întâi să-i determinăm funcția utilă, care este o funcție socială. Dintre **influențele pozitive ale umorului** putem menționa:

- afectează atenția și memoria [Baym, 1995];
- facilitează interacțiunile sociale [Binsted et al., 2006];
- ameliorează problemele de comunicare [Bergson, 1992];

- ajută la armonizarea unei conversații [Bergson, 1992];
- poate stabili un punct comun între partenerii de discuție [Hewitt, 2002];
- face conversația plăcută [Nijholt, 2005];
- contribuie la motivare, la atenție, la înțelegerea și captarea informațiilor și dezvoltarea unui sentiment afectiv a mesajului [Nijholt, 2006] [Binsted et al., 2006];
- poate înlesni problemele de comunicare ce pot apărea în interacțiunea dintre agenți și om, pentru că umorul este un mecanism primar de stabilire a individualității, întărindu-se raporturile acestei comunicări (omul se simte apreciat atunci când alții îi recunosc glumele ceea ce îmbunătățește [Baym, 1995] [Black and Forro, 1999];
- umorul înlesnește comunicarea și învățarea individuală și ajută la sincronizarea învățării în grup atât la adulți cât și la tineri [Binsted et al., 2006];
- facilitează crearea unei înțelegeri comune, ajută la generarea solidarității și a identității de grup [Binsted et al., 2006];
- reduce stressul [Binsted et al., 2006];
- stimulează creativitatea și îmbunătățește comunicarea, morala și productivitatea [Stock and Strapparava, 2006];
- atragerea atenției auditoriului [Stock and Strapparava, 2006];
- ajută la memorare [Stock and Strapparava, 2006];

Umorele trebuie să răspundă unor exigențe ale vieții în comun, având o semnificație socială. Nerespectarea acestor exigențe de utilizare a umorului poate genera efecte contrare pentru care umorul a fost inițial creat. Pot apărea astfel tensiuni de ordin personal și social:

- câteodată este greu de înțeles glumele celorlalți deoarece indivizii, deși au același set de cunoștințe, au raționamente total diferite. Mediul social al umorului este vast și umorul poate fi văzut ca o parte din alte multe acțiuni sociale. Potrivit sociologilor impactul unei glume și a umorului este foarte mare în viața de zi cu zi;

- în cadrul grupului -umorul se bazează pe normele grupului, cunoștințe, practici și probleme, generând identitatea umană dar și pe cea de grup. Vorbim despre o identitate individuală, pentru că fiecare din cei implicați transmit părți de mesaje care apoi sunt filtrate de fiecare individ în parte. Deci se poate spune că umorul utilizat în grup definește grupul respectiv. Unii oameni se pot simți frustrați datorită neînțelegerii umorului celorlalți, și, deși umorul ar trebui să detensioneze, utilizat greșit poate înrăutăți situația;
Din influențele negative ale umorului putem aminti:
- poate jigni;
- poate inhiba comunicarea din cauza stilurilor diferite de a face glume;
- poate crea tensiuni din punct de vedere profesional

[Black and Forro, 1999]

1.2 Tipuri de umor

Există numeroase forme și tipuri de umor. Dintre acestea, câteva sunt adecvate pentru folosire în dezbateri, fiind deopotrivă și cele mai des întâlnite:

- **Anecdota:** orice întâmplare interesantă care ajută vorbitorul să clarifice o chestiune. De multe ori cei care dezbat sunt încurajați să folosească studii de caz sau să construiască scenarii. Unele dintre ele pot căpăta o nuanță comică, nu doar pentru a capta atenția, ci și pentru a accentua o anumită idee.

Acordorul: Bună dimineața, am venit să vă acordez pianul.

Domnul Bergovici: Dar eu nu am cheamat nici un acordor.

Acordorul: Știu, vecinul dumneavoastră m-a chemat [Solomivici, 2002].

- **Exagerarea (Hiperbola):** exagerarea unor trăsături, defecte sau a neadecvării unei acțiuni.

Simon și Dov stau de vorbă la cafenea

-În timpul călătoriei mele în Africa, povestește Simon, am întâlnit un negru atât de negru, încât a trebuit să aprind lumina în plină zi, ca să-l pot vedea.

-*Asta nu-i nimic, spune Dove, când am fost în Spania am văzut un om atât de slab că trebuia să între de două ori într-o încăpere ca să-ți poți da seama că e acolo* [Solomivici, 2002].

- **Ironia:** folosirea cuvintelor pentru a exprima altceva decât în mod normal. De obicei vorbitorul spune opusul la ce se gândește sau a ce așteaptă publicul de la el.
Rabinul nu este deloc mulțumit de atitudinea plină de umilință a lui Gelb.
-Măi Gelb, nu ești tu atât de mare pe cât de mic vrei să pari [Solomivici, 2002].
- **Revenirea:** vorbitorul comite o eroare în mod intenționat, numai pentru a reveni și a corecta repede acea eroare (în cazul unei dezbateri aceasta metodă poate fi folosită pentru a da dreptate inițial oponenților, și a reveni apoi pentru a arăta cât de greșit ar fi fost acest lucru).
-Am auzit că băiatul tău s-a căsătorit. Din dragoste sau pentru bani?
-Din dragoste... pentru bani [Solomivici, 2002].
- **Satira:** o formă de sarcasm care scoate în evidență lipsurile unei idei/acțiuni/persoane.
Știi bancul cu statuia? Statu-ia tot...
- **Subestimarea:** transformarea a ceva mare sau important în ceva normal sau chiar mult mai mic/mai puțin important decât este în realitate
Doi ardeleni tăiau o bombă cu fierăstrăul. Vine al treilea și îi întreabă:
- Mă, dar dacă explodează?
La care ceilalți:
-Apăi nu-i bai, că mai avem una!
- **Umorul de situație:** umor care vine din experiența proprie. Puteți fi siguri că publicul nu are de unde să cunoască situația descrisă.
Aflat în delegație constată că și-a uitat papucii acasă. Se așează la o masă și îi scrie o scrisoare soției:
Dragă Lea,
Trimite-mi numai decît papucii tăi. Am scris "papucii tăi" fiindcă dacă aș fi scris "papucii mei" ai fi citit papucii mei și mi-ai fi trimis papucii tăi. Ce să fac eu cu papucii tăi? Așa că îți scriu foarte clar ca să înțelegi-"papucii tăi" și nu "papucii mei", ca să citești papucii tăi și să îmi trimiți papucii mei...[Solomivici, 2002]

1.2.1 Gramatica umorului

Varietatea elementelor și tipurilor care contribuie la crearea diferitelor forme de umor este nemăsurată, totuși aceste elemente trebuie să conțină un ingredient de bază și indispensabil: un impuls sau o urmă de agresivitate sau chiar răutate. Asta este tocmai ce filozoful francez Henri Bergson numea o anestezie măreață a inimii [Bergson, 1992]. Mulți teoreticieni ai umorului au încercat să determine logica sau gramatica umorului. Să luăm de exemplu următoarele glume:

1. *Un masochist e o persoană căreia, deși îi place să facă un duș fierbinte dimineața, face unul cald.*
2. *Domnule aș vrea să vă cer mâna fiicei dumneavoastră. De ce nu? Deja ai avut restul.*

Analizând cele două glume, dorința noastră este de a identifica un șablon care subliniază cele două situații. Gluma masochistului care se pedepsește și nu își satisface plăcerea zilnică este guvernată de regula care prezintă o **inversare a situației logice**. Mâna fetei este percepută în primul caz în sensul metaforic dar apoi în sensul de bază. Altfel ascultătorul percepe situația în **două cadre incompatibile**. Minteia lui trebuie să opereze simultan în două lungimi de undă. Evenimentul nu este numai asociat unui singur cadru (cum este în cazul normal), ci împărțit în două.

În umor, atât creația cât și percepția unei glume implică plăcerea de a face salt de la un sens la altul. Șablonul comun care subliniază cele două glume este percepția unei situații în două cadre incompatibile, care își găsesc un punct comun, ceea ce implică inteligență. Totuși umorul vine de cele mai multe ori din percepția unei relații între două contexte consistente, dar incompatibile și acest tip de activitate mentală pare să fie încântătoare pentru ființele umane. Tot ce este **contrar** nu numai **obișnuinței** dar și **simțurilor** sau **rațiunii** sau ceva care **lasă loc de interpretare** este amuzant. Sensul logic comun este neviabil pentru uzul practic. Nu poate fi reperat, de aceea trebuie să încercăm să evităm funcționările cele mai comune. Umorul joacă un rol important în învățarea și comunicarea acestui aspect [Minski, 1981].

Umorul în mod sigur se bazează pe complicație. În primul rând sunt două scenarii, cel natural și cel neașteptat și de obicei acestea două sunt în opoziție. Un cuvânt cheie sau un declanșator îndreaptă mintea ascultătorului spre o altă viziune.

1.3 Teoriile umorului

Teoriile umorului sunt o încercare de a găsi un șablon comun pe baza căruia amuzamentul este construit. Marile școli de filozofie au încercat să găsească răspunsuri la scopul și originea umorului, precum și modul de manifestare. E.M. Blistein identifică trei mari categorii: [Blistein, 1964]

- Teoria de superioritate

Potrivit teoriei de superioritate a umorului chiar și cel mai subtil umor își are originea într-o anumită cruzime și plăcere, altfel umorul vine dintr-un sentiment de superioritate față de cei de care râdem. Pasiunea pentru amuzament nu este altceva decât un sentiment de măreție crescută în noi dintr-o concepție de superioritate în comparație cu inferioritatea celorlalți. Adesea absurditatea, măreția, ciudățenia, infirmitatea provoacă râsul.

Această teorie conține trei părți:

- orice situație amuzantă are un câștigător; Pe parcursul istoriei oamenii au folosit umorul pentru a concura cu alte persoane, făcându-i ținta glumelor lor. Câștigătorul este cel care reușește să îl ia în râs pe celălalt. Potrivit teoriei de superioritate oamenii care nu respectă un standard al societății tind să fie ridicoli.
- nepotrivirea este mereu prezentă în umor;
- umorul necesită un element surpriză;

- Teoria de nepotrivire

Când glumele sunt examinate din lumina teoriei de nepotrivire pot fi observate două obiecte prezentate sub un singur concept sau cadru. Conceptul poate fi aplicat ambelor obiecte și obiectele devin similare. Pe măsură ce gluma progresează, discrepanța dintre aceste obiecte devine tot mai evidentă. Teoria de nepotrivire se axează pe elementul surpriză. Umorul este creat dintr-un conflict între ceea ce se așteaptă și ceea ce se întâmplă de fapt într-un text scris. Aceasta are în vedere cel mai evidente forme de umor: ambiguitate, dublul înțeles care în mod deliberat induce în eroare auditoriul, urmată de ultima frază a unei glume. În acest fel umorul încalcă o regulă importantă a utilizării limbii: Comunicarea trebuie să fie cât mai clară.

- Teoria de eliberare

Teoria de eliberare este de natură fiziologică și psihofiziologică [Rutter, 1997]. În afară că exprimă o bătălie între două persoane sau un grup din societate, umorul poate exprima și un conflict cu noi înșine, ne eliberează de grija de a fi noi înșine. Seriozitatea este un stress obișnuit asupra minții în fața ordinii evenimentelor urmate cu o regularitate, pe când umorul este o pierdere neașteptată sau o domolire a acestui stress [Hazlitt, 1963].

Oamenii care sunt tensionați izbucnesc în râs de îndată ce tensiunea dispare. Potrivit teoriei de eliberare, elementul central în umor nu este nici sentimentul de superioritate nici noutatea ci sentimentul de eliberare care vine din eliberarea inhibițiilor.

Să luăm ca exemplu următoarea glumă:

Un om intră într-o cofetărie și comandă o prăjitură, dar cum i se aduce prăjitura cere un pahar cu lichior în locul ei. A băut paharul și voia să plece fără să fi plătit pentru el. Proprietarul l-a oprit:

-Nu ați plătit lichiorul.

-Dar v-am dat prăjitura în locul lui.

-Dar nu ați plătit nici pentru aceasta.

-Dar nici nu am mâncat-o [Freud, 1957].

Rezultatul unei asemenea glume este că ne simțim contrariați apoi în încurcătură; după aceea râdem și găsim absurditatea logicii amuzantă. Personificarea perfectă pentru această teorie poate fi găsită în teoria lui Sigmund Freud. Freud privește umorul ca o metodă de a se elibera de cenzură - termen folosit pentru inhibițiile interne care ne previn să ne eliberăm de toate impulsurile noastre naturale. Prin cenzură nu se referă doar la impulsurile de natură sexuală, ci și la cele răutăcioase. Teoria de superioritate a lui Freud se bazează pe răutate. El distinge între glumele nevinovate și cele tendențioase și deși afirmă că există și glume inocente, acestea foarte rar reușesc să provoace o izbucnire în râs ca cele răutăcioase. El explică de ce glumele tind să fie compacte și condensate, cu dublu înțeles. Este pentru a-i păcăli pe cenzorii care văd doar partea nevinovată a înțelesului și nu reușesc să perceapă dorințele interzise [Freud, 1957]. În cele din urmă nu explică umorul doar după sensurile acestor dispozitive comice; ele ar fi fără

sens dacă omul nu ar putea, sub acoperire, să dea frâu liber dorințelor reprimate. Glumele sunt comparate cu visele, pentru că datorită lor putem afla mai multe despre activitatea normală a psihicului fiecăruia. Scopul umorului este de a readuce în lumea adulților plăcerea care caracteriza copilăria noastră, când nu aveam nevoie de glume pentru a ne simți fericiți.

1.4 Criterii pentru determinarea succesului umorului

Criteriul prin care se determină dacă umorul este bun sau rău depinde în mare măsură de gustul și preferințele personale, și, în particular, depinde de stil și tehnica umoristică. Totuși, se pare că sunt trei criterii de care depinde succesul umorului: **originalitatea, accentul și economia**. Meritul originalității este de la sine înțeles; oferă elementul esențial: surpriza, care depășește așteptările noastre. Accentul se obține prin exagerare sau simplificare, două criterii care conduc la efectele umorului. Economia în umor nu înseamnă curaj, ci mici aluzii în locul atacului frontal, indicii implicite în locul declarațiilor explicite. Fiecare din aceste trei teorii pot explica câteva dintre tipurile de umor, dar nu pot constitui explicațiile pentru toate tipurile de umor. Dar umorul este în mod sigur cu mai multe fețe. Poate fi agresiv sau poate batjocoritor, poate fi inteligibil, poate exprima eliberarea din tensiune sau libertatea de griji, poate fi jucăuș, sau inteligent, chiar poate deveni serios, dar nu poate fi fals. Umorul, chiar și la cel mai exagerat nivel de răutate, **nu poate abandona adevărul**.

Sunt trei condiții pentru existența umorului [Veatch, 1998]:

- V Violarea unui anumit angajament al ascultătorului despre cum lucrurile trebuie să fie.
- N Receptorul are sentimentul dominant că o anumită situație este normală.
- S Simultaneitate: cele două condiții anterioare N și V trebuie să fie prezente în mintea receptorului în același timp.

În altă ordine de idei, o situație este amuzantă când situația pare normală dar în același timp ceva apare în neregulă.

Veatch descrie condiția V ca violare a moralității subiective, o violare a lucrului de care este atașat receptorul. Poate fi descrisă cel mai bine în acest fel, pentru că depinde de ceea ce receptorul experimentează ca violare. Depinde de angajamentul sau atașamentul său față de o situație dată. Veatch continuă explicând că gradul de atașament al receptorului este important deoarece o situație i se poate părea amuzantă, jignitoare sau neremarcantă. Descrie acest fenomen cu trei niveluri de angajament și consecințe, relația atașament-umor fiind descrisă mai bine în Tabelul 1.

Nivel	Logica	Atașament	Receptorul		
			O înțelege	îl jignește	Observă umorul
Nivelul 1	non-V	nimic	nu	nu	nu
Nivelul 2	V și N	slab	da	nu	da
Nivelul 3	V și non-V	puternic	da	da	nu

Tabela 1.1: Scara nivelurilor de de atașament

Factorul de normalitate (totalitatea informațiilor normale dintr-o glumă) împreună cu nivelul de atașament, au o profundă importanță în felul cum receptorul percepe experiența violării.

1.5 Umore din punct de vedere social

Simțul umorului poate fi:

- individual (doar unor indivizi li se pot părea amuzante unele situații);
- universal (de ex: clovnii pentru majoritatea oamenilor sunt stimulatori de umor);
- cultural (de ex: vestimentația africanilor pot provoca ”șocuri ” culturale europenilor);

Ceea ce poate pentru un individ este foarte amuzant pentru altcineva poate să nu însemne nimic. Recunoașterea umorului nu este doar subiectivă, ci este și personală. Umorel caracterizează interacțiunea persoanelor în societate și răspunsul nostru la această interacțiune trebuie înțeles în plin context [Mulder and Anton, 2001]. în construcția unui mesaj, vorbitorii și scriitorii

pornesc cu niște cunoștințe pe care ei le presupun a fi cunoscute de mai mulți destinatari și la care noua informație (cea pe care o vor să o transmită) poate fi atașată. Presupunerea este bazată pe cunoștințe acumulate ale lumii sau pe experiențele cu membrii comunității. Când comunicarea are succes, așteptările vorbitorului sau ale scriitorului sunt în concordanță cu percepția destinatarului și noua informație se poate lega cu celelalte.

Umorul și comunicarea în sine depind de acumularea de înțelegeri împărțite, de stilul oratoric al oamenilor aparținând unei anumite culturi. Comunicarea se rupe când nivelul de cunoștințe anterior acumulate de către vorbitor/scriitor și ascultător/cititor nu sunt la fel. Aceasta este adevărat pentru orice tip de comunicare dar ruperea este în mod special evidentă în cazul umorului, a cărui percepție depinde direct de concurența între fapte și impresii disponibile atât vorbitorului/scriitorului cât și ascultătorului/cititorului.

Sunt multe glume care semantic pot însemna același lucru, dar în termeni de pragmatism și cultură pot apărea neînțelegeri care pot împiedica înțelegerea glumei. Acestea sunt, ceea ce numim "glume culturale" sau "glume etnice" și ar fi dificil să spunem că există glume universale. Poate o glumă universală este o glumă biculturală. Multe glume etnice sunt interschimbabile, depinzând totuși, de audiență și de cel ce spune gluma. Bancuri despre olteni pot deveni despre moldoveni pentru oamenii din diferite regiuni ale țării.

Capitolul 2

Cercetările în domeniul umorului computațional

Studierea prin mijloace computaționale a umorului este un domeniu ale cărui baze au început să se pună abia în ultimii ani, neexistând o teorie general acceptată (ca de exemplu Centering Theory în Teoria Discursului). În continuare sunt trecute în revistă principalele cercetări dedicate acestui domeniu pe plan mondial.

2.1 Generare de umor

2.1.1 LBJOJG

LBJOJG(Light Bulb Joke Generator) a fost dezvoltat de Attardo și Raskin în 1993 și după cum îi spune și numele, este un generator de glume de tipul ”*De câți (substantiv) este nevoie pentru a (verb)*”, însă era foarte limitat deoarece nu ansamblează sau analizează atribute ale glumelor [Attardo and Raskin, 1994].

De exemplu șablonul:

*How many [substantiv] does it take to change a light bulb?
[număr1].[număr1-număr2] to [activitate1] and [număr2] to [activitate2].*

cu intrarea:

(poles (activitate1 hold the light bulb) (număr1 five) (activity2 turn the table he's standing on)(număr2 four))

va genera următoarea glumă:

How many Poles does it take to change a light bulb? Five. One to hold the light bulb and four to turn the table he's standing on.

2.1.2 Jape

Cercetările lui Binsted și Ritchie au condus la model semantic și sintactic de reguli ce a produs generatorul de glume numit JAPE. JAPE folosește substituția cuvintelor, a silabelor pentru a crea ambiguitate fonologică. Este alcătuit dintr-un vocabular format din 59 de cuvinte omofone, 14 șabloane și un validator de glume [Binsted and Ritchie, 1997].

Exemplu: *What do you call a quirky quantifier? An odd number.*

2.1.3 HAHAcronym

Un alt proiect a fost HAHAcronym dezvoltat de Stock și Strapparava pentru un sistem care generează versiuni amuzante ale acronimelor deja existente. Efectul comic s-a obținut mai ales prin exploatarea teoriei de nepotrivire. Algoritmul urmărește prima dată parsarea acronimelor. Se va păstra o parte din cuvintele care definesc acronimul iar înlocuirea celorlalte cuvinte se face prin:

- utilizarea unui câmp semantic opus
- păstrarea literei inițiale, a ritmului și rimei

Exemplu: ACM(Association for Computing Machinery) devine Association for Confusing Machinery [Stock and Strapparava, 2003]

2.1.4 WISCRAIC

WISCRAIC este un generator de glume utilizat pentru cei ce învață limba engleză. Un exemplu de glumă generată de acesta este:

What bird is lowest in spirits? A bluebird

(ambiguitate deoarece *blue* poate fi culoarea sau poate însemna depresie) [Binsted-McKay, 2000]

2.1.5 MSG

MSG este un program care convertește cuvintele alfanumerice din parole în propoziții amuzante. Programul ia ca argument 8 caractere și va trebui să genereze propoziții ușor de memorat. De exemplu, folosind șablonul:

$$(cuv1=Nume\ Personă)+(cuv2=Verb\ Pozitiv)+(cuv3=Nume\ Personă+''s'')+(cuv4=Substantiv\ Comun)+,while+(cuv5=Nume\ Personă)+(cuv6=Verb\ Negativ)+(cuv7=Nume\ Personă+''s'')+(cuv8=Substantiv\ comun)$$

cu șirul alfabetic *AjQA3Jtv* se obține *Arafat joined Quayle's Ant, while Tarar Jeopardized thurmond's vase*. [McDonough, 2001]

2.1.6 Jester

Jester este un sistem online de recomandare de glume. În funcție de alegerile efectuate se determină statistic printr-un algoritm de tipul cel mai apropiat vecin, care sunt gusturile utilizatorului în privința glumelor. [Goldberg et al., 2000]

2.2 Recunoașterea umorului

2.2.1 Recunoașterea umorului fără înțelegerea sensului

Această secțiune se bazează pe articolul publicat de [Sjobergh and Araki, 2007]. Pentru a recunoaște umorul fără a face toate conexiunile dintre cuvintele unei propoziții se calculează valori ale unor atribute ce pot caracteriza un text. În funcție de valorile acestor caracteristici se încearcă o clasificare a textelor. Se folosesc combinații de atribute pentru a vedea dacă acestea sunt suficiente pentru recunoașterea umorului sau nu.

Caracteristicile

- **Cea mai apropiată glumă:** Din datele de antrenament se caută gluma cea mai apropiată de textul pe care îl testăm. Gradul de apropiere

dintre 2 texte se calculează după numărul de cuvinte comune celor 2 texte.

- **Cea mai apropiată non-glumă:** Pe același principiu se determină cea mai apropiată non-glumă.
- **Cele mai apropiate 5 vecine:** Se determină cele mai apropiate 5 propoziții în datele de antrenament, fiecare cu o pondere egală cu numărul de cuvinte care se suprapun. Glumele au semn pozitiv și non-glumele au semn negativ.
- **Cuvinte amuzante:** S-a observat că unele cuvinte sunt comune dar unele sunt specifice doar glumelor. Pentru a surprinde acest aspect, cuvintele care apar măcar de 5 ori în datele de antrenament și dacă apar de 5 ori mai des în glume decât în non-glume sunt păstrate într-o listă. Fiecărui cuvânt îi este asignat o pondere care este frecvența relativă cuvântului printre glume (numărul de glume în care apare cuvântul/numărul total de glume) împărțită la frecvența relativă a cuvântului printre non-glume.
- **Ambiguitatea cuvintelor:** se calculează uitându-ne într-un dicționar online și numărând sensurile cuvintelor.
 - **Ambiguitatea medie:** numărul mediu de ambiguități într-o propoziție (media numărului de sensuri pentru fiecare cuvânt);
 - **Ambiguitatea maximă:** cea mai mare valoare a numărului de sensuri pentru un cuvânt dintr-o propoziție.
- **Cuvinte murdare:** numărul de cuvinte murdare prezente în propoziții. O listă cu 2500 de cuvinte murdare downloadată de pe Internet a fost folosită pentru a se decide dacă un cuvânt este murdar sau nu.[Sjobergh and Araki, 2007]
- **Numărul de cuvinte prezente din lista:** *you, your, I, me, my, man, woman, he, she, his, her, guy, și girl*; și negațiile, numărul de apariții ale lui *not sau nt*.
- **Pronunția:** folosind Dicționarul CMU¹ de Pronunție pentru a găsi pronunția cuvintelor au fost calculate:

¹<ftp://ftp.cs.cmu.edu/afs/cs.cmu.edu/data/anonftp/project/fgdata/dict/>

- **Rime:** numărul de cuvinte perechi care au cel puțin 4 litere, cel puțin una din ele este vocală, pronunțată la fel la sfârșitul cuvântului;
 - **Similarități** la fel ca și rimele dar folosind începutul cuvintelor în locul sfârșitului;
 - **Cuvintele noi:** cuvintele care nu apar în Dicționarul de traduceri CMU²
- **Cuvintele care se repetă:** numărul de cuvinte care apar de mai mult decât odată în propoziție și care au mai mult de 5 litere;
 - **Subșiruri care se repetă** cel mai lung subșir din propoziție care este prezent mai mult decât odată. Media cuvintelor care se repetă și a subșirurilor este aceeași, dar divizată după lungimea propozițiilor;
 - **Antonimia:** caută antonimia unui cuvânt în dictionary.com. Dacă oricare din aceste antonime listate este prezent în propoziție, un scor de 1 împărțit la numărul de antonime posibile îi este atribuit.
 - **Scorul maxim de antonimie:** cea mai mare valoare (perechea de antonime cu cele mai puține alte valori de antonime);
 - **Antonimia:** suma tuturor valorilor antonimelor calculate;

Pentru fiecare caracteristică este calculată o funcție de prag astfel încât media entropiei³ acestora să fie cât mai mică posibil. Pentru a clasifica o dată de test, și a vedea cărui grup aparține, este verificată pentru fiecare caracteristică proporția dintre exemplele pozitive și cele negative. Proporția exemplilor pozitive pentru fiecare grup la care aparține este înmulțită și comparată apoi cu produsul proporțiilor caracteristicilor negative. Dacă produsul exemplilor pozitive este mai mare atunci exemplul este pozitiv. Metoda are

²<ftp://ftp.cs.cmu.edu/afs/cs.cmu.edu/data/anonftp/project/fgdata/dict/>

³Măsură statistică care calculează gradul de împrăștiere a datelor. Dacă S este un exemplu de date de antrenament p_{\oplus} este proporția datelor pozitive, iar p_{\ominus} este cea a datelor negative atunci

$$Entropia(S) = -p_{\ominus} \log_2 p_{\ominus} - p_{\oplus} \log_2 p_{\oplus}$$

	Cu	fără
Toate atributele	85.4%	50.0%
Similaritate	75.7%	83.8%
Cuvinte prezente în glume	84.1%	76.8%
Ambiguitate	62.5%	84.8%
Stil	59.1%	85.4%
Idiomuri	63.5%	85.0%

Tabela 2.1: Recunoașterea umorului prin caracteristici

avantajul de a fi foarte rapidă.

Deoarece sunt foarte multe caracteristici care reprezintă aproape aceeași informație, autorii încercă să elimine caracteristicile redundante sau nefolositoare pentru a mări performanța. (Tabelul 2.1) Eliminarea unor caracteristici este importantă deoarece așa putem afla ce informație este utilă în detectarea umorului. Procesul se desfășoară astfel: se elimină caracteristicile pe rând. Atributul, care atunci când nu este prezent dă cel mai bun rezultat, este eliminat. Când toate caracteristicile sunt eliminate, cel mai bun rezultat (împreună cu structura caracteristicilor prezente) este păstrat și după acest rezultat se poate vedea care caracteristici sunt mai importante.

Corpusul

Este format din 6800 de glume colectate de pe Internet și 6800 de propoziții colectate din BNC (British National Corpus). Toate propozițiile au lungimea cuprinsă între și 80 de cuvinte, cu o medie de 15 cuvinte.

2.2.2 Aplicarea clasificatorilor de texte și a unor atribute pentru recunoașterea umorului

Se poate încerca recunoașterea umorului folosind învățarea automată, utilizând clasificatorul Bayes naiv și Mașini cu vector suport.

Rada Mihalcea și Carlo Strapparava au ales să își restricționeze studiul la *one-linere* [Mihalcea, May 2006], această secțiune bazându-se pe studiile celor doi.

Un *one-liner* este o propoziție cu efect comic și cu o structură lingvistică interesantă: sintaxă simplă, folosirea deliberată a unor instrumente retorice

(e.g. aliterații, rime) și utilizarea frecventă a unor structuri lingvistice menite să atragă atenția cititorilor. În timp ce glumele mai lungi tind să producă umor printr-o structură narativă mai complexă, onelinerele produc efectul comic dintr-o lovitură, cu foarte puține cuvinte. Acest lucru face ca acest tip de text să fie folosit pentru recunoașterea automată a umorului, deoarece efectul comic se produce în prima și singura propoziție.

Exemplu: *Everyone has a photographic memory. Not everyone has film.*

Datele umoristice au fost alcătuite din onelinere colectate de pe Internet folosind procesul de bootstrapping. Datele nonumoristice au fost selectate și structurate astfel încât să fie structural și stilistic similare onelinerelor. [Mihalcea and Pulman, 2007]

Corpusul

- **Date negative:**

Au testat pe 3 tipuri diferite de exemple negative, o propoziție având în medie 10-15 cuvinte. Colecția cuprinde:

1. **Titluri Reuters** extrase din (8.20.1996-8.19.1997). Titlurile sunt formate din propoziții scurte menite să atragă atenția la fel ca cele din onelinere.
2. **Proverbe** extrase dintr-o colecție de proverbe. Proverbele sunt texte care transmit de obicei într-o propoziție, fapte importante sau experiențe care sunt considerate adevărate de oameni. Proprietatea lor de a fi condensate și de transmite un mesaj într-o singură propoziție le face foarte asemănătoare cu onelinerele. Defapt unele glume încearcă să reproducă proverbele, cu un efect comic, ca în exemplul următor:
Beauty is in the eye of the beer holder, preluată din *Beauty is in the eye of the beholder*.
3. **British National Corpus (BNC)** propoziții extrase din BNC un corpus balansat care acoperă stiluri, genuri și domenii diferite. Propozițiile au fost colectate astfel încât să fie similare ca și conținut cu onelinerele: s-a folosit un sistem de colectare a informațiilor implementând un model vectorial pentru a identifica propozițiile cele mai asemănătoare cu fiecare din cele 16000 de onelinere.

- **Datele pozitive:**

Pentru a colecta foarte multe date este destul de dificil deoarece majoritatea site-urilor Web sau liste de mail fac public nu mai mult de 50-100 de glume. Pentru a depăși această problemă, s-a implementat o tehnică de colectare automată a glumelor pornind cu câteva glume manual identificate. Algoritmul identifică apoi în mod automat o listă de pagini web în care se găsește gluma respectivă. Paginilor astfel găsite li se aplică 2 constrângeri.

1. **constrângerea implementată:** este un set de cuvinte cheie care au legătură cu tema;

Setul de cuvinte cheie folosit în implementare este alcătuit din 6 cuvinte cheie care au legătură cu tematica căutată: *onliner, onliner, one-liner, humor, humour, joke, funny*.

De exemplu: <http://www.berro.com/Jokes>

<http://www.mutedfaith.com/funny/life.htm>

2. **constrângerea modelată:** exploatează structura HTML-urilor paginilor web;

Aceasta se bazează pe ipoteza că paginile web tind să folosească enumerații atunci când au o colecție de date de același tip. De exemplu dacă una din online.rele colectate manual este prezentă într-o pagină web precedată de un tag HTML de tip \ulcorner , atunci printre alte linii care au același tag se consideră că ar fi onelinere.

[Mihalcea and Strapparava, 2006]

După două iterații ale algoritmului de căutare, pornind de fiecare dată de la un set mic de câte 10 onelinere s-a ajuns la colectarea a 24000 posibile *onelinere*. După ce s-au eliminat duplicatele prin algoritmul de cea mai lungă secvență comună, s-a ajuns la un set de 16000 de onelinere, care pot fi folosite în experimente.

Prin procesul automat s-au identificat prin selecția aleatorie 200 de glume, cu o probabilitate de 9% de eroare, ceea ce nu are un impact prea mare în procesul de învățare [Mihalcea, May 2006].

Experimentele

Au avut în vedere tehnici automate de clasificare folosind euristici bazate pe attribute stilistice specifice umorului (aliterații, antonimii, cuvinte cu sens conotativ), cu un framework de învățare formulat ca un clasificator tipic de text. S-a încercat identificarea unui set de attribute care să fie atât semnificative cât și fezabil de implementat utilizând algoritmi existenți.

- **Aliterația** : proprietățile structurale și fonetice ale glumelor sunt cel puțin la fel de importante ca și conținutul. *Onelinerele* se bazează pe producerea de efecte umoristice prin atragerea atenției cititorului prin aliterații, repetarea unor cuvinte, rime care au un efect comic.

Următoarele onelinere sunt un exemple de glume care includ una sau mai multe lanțuri de aliterații.

Veni, Vidi, Visa: I came, I saw, I did a little shopping.

Infants dont enjoy infancy like adults do adultery.

Pentru a extrage aceste caracteristici se identifică și se numără aliterațiile sau lanțurile de rime din fiecare set de date. Lanțurile sunt extrase automat utilizând un index creat cu ajutorul dicționarului de pronunții CMU.

- **Antonimia**: umorul se bazează de obicei pe nepotriviri, opuneri și alte forme de contradicții.
Spre exemplu, efectul comic produs de următoarele onelinere este datorat prezenței antonimelor. *A clean desk is a sign of a cluttered desk drawer. Always try to be modest and be proud of it!* Sursa lexicală folosită pentru identificarea antonimelor este WordNet⁴.
- **Construcții cu contotații sexuale**: de exemplu următoarele onelinere cuprind astfel de expresii. *The sex was so good that even the neighbors had a cigarette. Artificial Insemination: procreation without recreation.* Pentru a forma un lexicon pentru identificarea acestor attribute, s-au extras din Wordnet toate synset-urile marcate ca făcând parte din domeniul sexualității. Lista este apoi procesată pentru a elimina cele cu o polisemie mai mare de 4. Unele onelinere conțin toate cele 3 attribute:

⁴<http://wordnet.princeton.edu/>

Euristica	One-linere Reuters	One-linere BNC	One-liners proverbe
Aliterații	74.31%	59.34%	53.30%
Antonimii	55.65%	51.40%	50.51%
Toate	76.73%	60.63%	53.71%

Tabela 2.2: Rezultatele obținute folosind euristicile stilistice

Clasificatorul	Reuters	BNC	Proverbe
Bayes Naiv	96.67%	73.22%	84.81%
SVM	96.09%	77.51%	84.48%

Tabela 2.3: Rezultatele obținute folosind algoritmi de învățare automată

Behind every great man is a great woman, and behind every great woman is some guy staring at her Behinds.

Pe lângă aceste caracteristici stilistice experimentele au avut în vedere și caracteristici de conținut, în care recunoașterea umorului este formulată ca o problemă de clasificare.

Rezultatele

Primul set de experimente (Tabela 2.2) a evaluat acuratețea de clasificare folosind atributele de mai sus: aliterații, antonimii, cuvinte cu conotație sexuală. Acestea sunt atribute numerice care se comportă ca euristici, și singurul parametru este determinarea unui prag indicând valoarea minimă pentru a spune dacă un text poate fi clasificat ca fiind umoristic sau nu. Acest prag poate fi obținut utilizând arbori de decizie aplicați pe un set de date mic de exemplu 1000, iar pe restul de 15000 să fie testat. Luând în considerare că titlurile de articole reprezintă indicatori stilistici, indicatorii din titlurile Reuters sunt cele mai diferite de onelinere. Pentru toate seturile de date atributul de aliterație pare să fie cel mai important.

Al doilea set de experimente (tabela 2.3) se ocupă cu evaluarea atributelor determinate de conținutul lor. încurajați de rezultatele obținute în cele două experimente s-a construit un al treilea experiment care încearcă să exploateze atributele de stil și de conținut. Toate evaluările folosesc "10-fold cross-validation" (validare încrucișată prin împărțirea în 10 părți).

Capitolul 3

Experimente de identificare a umorului prin implementarea algoritmiilor clasici de clasificare de texte

3.1 Clasificarea de texte

Scopul clasificării de texte este de a asigura itemi uneia sau mai multor categorii predefinite pe baza conținutului contextual.

Funcțiile optimale de categorizare pot fi învățate din datele de antrenament. [Dumais, 1998]

Definiție 1 *Un program învață dintr-o experiență E dată de o clasă de activități T și o măsură a performanței P , dacă performanța activității T , măsurată cu P , se îmbunătățește odată cu experiența E .*

[Mitchell, 1997]

În cazul nostru un program spunem că învață să recunoască umorul din texte dacă își îmbunătățește performanțele măsurate prin abilitatea lui de a recunoaște umorul ceea ce implică analiza textelor, obținută prin filtrarea diferitelor tipuri de texte.

O problemă de detectare a umorului se poate formula astfel:

- **Obiectivul T:** recunoașterea și clasificarea umorului în texte;

- **Măsurarea performanței P:** procentajul de situații umoristice corect clasificate;
- **Experiența E:** o bază de date cu texte corect clasificate;

Atributele care influențează alegerea antrenării corecte:

1. Tipul de experiență din care va învăța sistemul

Tipul de experiență din care învață programul are un impact foarte mare în succesul sau eșecul mașinii de învățare. Un atribut cheie ar fi dacă antrenarea va oferi feedback:

- **direct:** Sistemul va învăța din exemple de antrenare directe constituite din cuvinte sau expresii și probabilitatea lor de a produce situații amuzante.
- **indirect:** Alternativ, poate avea doar informație indirectă alcătuită din propoziții și clasificarea în amuzante sau nu. Aici cel care învață se confruntă cu o problemă în plus de credibilitate: determinarea gradului în care fiecare cuvânt are importanță foarte mare în stabilirea deciziei finale. Credibilitatea poate fi o problemă extrem de dificilă deoarece un cuvânt poate schimba sensul unei propoziții.

2. Gradul în care cel care învață controlează secvența de exemple de antrenare

De exemplu cel care învață se poate baza pe un *profesor* care să selecteze un grup de propoziții pentru filtrare. Dacă avem de a face cu un text pe tema naturii și vrem să știm dacă este amuzant sau nu o euristică destul de intuitivă ar fi să căutăm în datele de antrenare acele propoziții care sunt și ele pe tema naturii. În mod alternativ, cel care învață poate selecta și propune un număr de grupuri de cuvinte pe care le consideră dificile și poate solicita răspunsul corect de la profesor.

3. Cât de bine reprezentată este distribuția exemplilor din care va învăța algoritmul nostru

În general învățarea este mult mai credibilă când datele de antrenament urmează o distribuție similară datelor de test. În exemplu nostru de clasificare de texte este foarte probabil să nu știm răspunsul corect.

Teoriile curente din învățarea automată se bazează pe presupunerea că distribuția datelor de antrenament este identică cu distribuția exemplilor de test. Chiar dacă trebuie să facem această presupunere, în practică această regulă este încălcată de cele mai multe ori.

Definiție 2 *Un clasificator de texte este o funcție care primește la input un document și îl încadrează într-o categorie y dintr-o mulțime predefinită de clase $y_1 \dots y_k$.*

Secțiunile 3.3 și 3.4 prezintă detaliat doi dintre cei mai folosiți algoritmi.

3.2 Metode de reprezentare a documentelor

Din documentele în forma lor originală nu se poate învăța. Ele trebuie să fie transformate pentru a se potrivi cu formatul algoritmilor de învățare. Deoarece cei mai mulți din algoritmi de învățare folosesc reprezentarea atribut-valoare, acest lucru înseamnă transformarea textului într-un vector. În primul rând toate documentele trebuie să fie **pre-procesate**. Acest lucru înseamnă de obicei eliminarea cuvintelor care nu prezintă importanță, aducerea cuvintelor la forma de bază, transformarea în litere mici.

După acest proces are loc **transformarea**. Fiecare cuvânt va corespunde unei dimensiuni (cuvintele identice aparțin aceleiași dimensiuni). Notăție: cuvântul w_i corespunzând dimensiunii i a spațiului vectorului.

Cea mai comună metodă este cea numită TF-IDF (Term Frequency Inverse Document Frequency). TFIDF(i, j) este a i -a coordonată a documentului j :

$$TFIDF(i, j) = TF(i, j)IDF(i) \quad (3.1)$$

$$IDF(i) = \log \frac{N}{DF(i)} \quad (3.2)$$

TF(i, j) reprezintă de câte ori al i -lea cuvânt se găsește în documentul j .

N numărul de documente

DF(i) numărul de documente care conțin cuvântul i măcar odată.

Documentele transformate formează împreună matricea termenilor documentului. Este de dorit ca documentele de lungime diferite să aibă aceeași lungime, care se realizează prin așa numita normalizare a documentelor.

Dimensiunea unui spațiu de vectori este foarte mare, ceea ce reprezintă un dezavantaj în învățarea automată, de aceea frecvent se apelează la metode

de reducere a dimensionalității. Sunt două posibilități: fie selectarea unui subset din atributele inițiale, fie integrarea mai multor atribute în unul.

Această metodă de reprezentare are proprietatea că pe măsură ce un cuvânt este mai frecvent în toate documentele cu atât este mai puțin valoros (nu oferă informații utile pentru o clasificare). [Pilasz, 2005]

3.2.1 N-gramele

Pentru a putea recunoaște sau genera o glumă, programul trebuie să fie capabil să proceseze secvențe de cuvinte. O metodă pentru această activitate ar putea fi N-gramele.

Definiție 3 *Un N-gram este un model care folosește probabilitatea condiționată de a prezice al N-lea cuvânt pe baza celor N-1 cuvinte anterioare.*

Se construiesc din statistici obținute dintr-un corpus mare de text folosind co-ocurența cuvintelor din corpus pentru a determina secvența de probabilități. Probabilitatea într-un model statistic precum cel al N-gramelor este dependentă de corpusul din care se face antrenarea. Dacă acest corpus este prea specific domeniului sau activității, programul nu va fi capabil să generalizeze.

Un bigram este un N-gram în care $N = 2$, iar pentru un trigram n este 3. Un bigram va folosi cuvântul anterior pentru a determina următorul cuvânt, iar un trigram va folosi 2 cuvinte anterioare. Probabilitatea bigramelor este una condiționată, formula pentru probabilitățile bigramelor fiind:

$$p(A|B) = \frac{p(A \wedge B)}{p(B)} \quad (3.3)$$

Pentru a calcula $p(B)$, următoarea formulă poate fi utilizată:

$$p(B) = \frac{\text{numărul de apariții ale lui B în text}}{\text{numărul de cuvinte în text}} \quad (3.4)$$

în mod similar,

$$p(A \wedge B) = \frac{\text{numărul de apariții ale lui B în text}}{\text{numărul de cuvinte în text}} \quad (3.5)$$

Aceasta înseamnă că $p(A \text{ — } B)$ este:

$$p(A|B) = \frac{\text{numărul de apariții ale lui A și B în text}}{\text{numărul de apariții ale lui B în text}} \quad (3.6)$$

De exemplu: *își dorește să asiste la ședințele cenaclului. Își dorește să asiste nu numai ca simplu spectator, ci să se implice, să se desfășoare, să se simtă util.* Pentru a afla care cuvânt este cel mai probabil să urmeze după *să* putem folosi bigrame. Avem perechile *să asiste* de 2 ori și *să se* de 3 ori. Deci $P(\text{se}/\text{să}) = 3/5$ iar $P(\text{asiste}/\text{să})=2/5$. Astfel folosind acest corpus spunem că după *să* va urma *se*. Dacă de exemplu în corpus ar fi fost doar prima propoziție ar fi fost alte rezultate.

Deci se poate observa necesitatea selectării unui corpus cât mai general, care să acopere datele de test.

3.3 Mașinile cu vector suport (Support Vector Machine SVM)

Intuitiv, într-o problemă de clasificare ar fi ideal să folosim cât mai multe caracteristici cu puțință ale datelor pentru a îmbunătăți rezultatul clasificării. În acest caz, cele mai multe clasificări suferă de așa numitul "small sample size effect". Adică, există un anumit număr optim de caracteristici de la care, dacă ne abatem, utilizând mai multe caracteristici în clasificare, performanța ar avea foarte mult de suferit.

Metoda bazată pe vectori de suport este o tehnică concepută pentru eficientizarea aproximării funcțiilor multidimensionale. Ideea de bază a SVM-urilor este de a determina un clasificator care minimizează riscul empiric (eroarea setului de antrenare sau acuratețea acestuia) și intervalul de încredere (erorile setului de test).

În 1965, Vapnik a propus o metodă de a găsi niște hiperplane care să "despartă" optim două clase, și care să nu depindă de estimarea unei probabilități. Acesta a fost baza teoriei mașinilor care învață bazându-se pe vectori de suport.

SVM-urile se bazează pe conceptul de plane de decizie (plane-hiperplane de separare) care definesc anumite "granițe". Un plan de decizie este un plan care separă un set de obiecte ce aparțin unor clase diferite.

Figura 3.1 prezintă ideea care stă la baza SVM-urilor. În figură observăm obiectele originale (din partea stângă a desenului) mapate (rearanjate), folosind un set de funcții matematice numite nuclee (kernels). Se poate vedea că obiectele mapate (din partea dreaptă a imaginii) sunt liniar separabile și, astfel, în loc să construim o curbă ca în figura 3.1, pentru a separa

obiectele, putem să construim o linie "optimă" care să separe obiectele albastre de cele roșii.¹

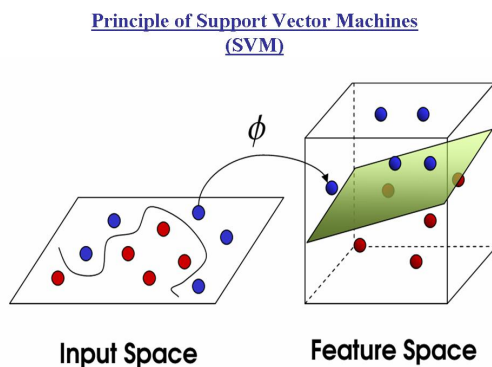


Figura 3.1: Maparea pe un alt spațiu

Problema clasificării poate fi restricționată, fără a restrânge generalitatea, la problema clasificării a doar două clase. În această problemă, obiectivul este să separăm cele două clase folosind o funcție indusă de exemplele pe care le avem la dispoziție. Scopul este de a obține o clasificare care funcționează bine și pe exemplele necunoscute încă (adică generalizează bine).

Există mai mulți clasificatori liniari care pot separa datele, dar numai unul dintre ei maximizează marginile (distanță între linie și cel mai apropiat punct din fiecare clasă). Acest clasificator liniar se numește **hiperplanul optim de separare**. Intuitiv, ne așteptăm ca această separare optimă să ne ajute cel mai mult la problema generalizării. Punctele care se află în contact cu zona de separare (marginea) se numesc **vectori de suport**.

Hiperplanul de separare

Să considerăm problema separării unui set de vectori de antrenare, care fac parte din două clase (unde x_i sunt datele de intrare iar y_i sunt clasele)

$$\mathbf{D} = \{(x^1, y^1), \dots, (x^l, y^l)\}, x \in \mathbf{R}, y \in \{-1, 1\} \quad (3.7)$$

separate cu un hiperplan.

$$\langle w, x \rangle + b = 0 \quad (3.8)$$

¹<http://www.imtech.res.in/raghava/rbpred/svm.jpg>

Se spune că setul de vectori este **separat optim de hiperplan** dacă este separat fără eroare și distanța între vectorul cel mai apropiat și hiperplan este maximă. Există o oarecare redundanță în ecuația 3.8, și fără a restrânge din generalitate este mai indicat să considerăm un hiperplan canonic [Mukherjee and Vapnik, 1999], unde parametrii w , b satisfac relația:

$$\min_i |\langle w, x_i \rangle + b| = 1 \quad (3.9)$$

Cu alte cuvinte norma vectorului ar trebui să fie egală cu inversa distanței între cel mai apropiat punct (obiect) din setul de date și hiperplan. Un hiperplan de separare în forma canonică trebuie să satisfacă următoarele relații,

$$y^i [\langle w, x^i \rangle + b] \geq 1, i = 1 \dots l. \quad (3.10)$$

Distanța $d(w, b; x)$ unui punct x față de hiperplanul (w, b) este,

$$d(w, b; x) = \frac{|\langle w, x^i \rangle + b|}{\|w\|}. \quad (3.11)$$

Hiperplanul optim este obținut prin maximizarea marginii care suferă constrângerea dată de ecuația 3.9. Hiperplanul care separă datele în mod optim este cel care minimizează relația,

$$\phi(w) = \frac{1}{2} \|w\| \quad (3.12)$$

Relația este independentă de b , pentru că dacă ecuația este adevărată (este un hiperplan de separare), modificarea lui b va produce "mișcarea" hiperplanului în direcția normală spre el însuși. Marginea rămâne neschimbată, dar hiperplanul nu va mai realiza o separare optimă, în sensul că va fi mai apropiat de una din clase. În continuare presupunem că următoarea inegalitate este satisfăcută,

$$\|w\| < A. \quad (3.13)$$

atunci, din 3.10 și 3.11,

$$d(w, b; x) \geq \frac{1}{A} \quad (3.14)$$

și deci hiperplanele nu pot fi la o distanță mai mică de $1/A$ față de oricare dintre puncte.

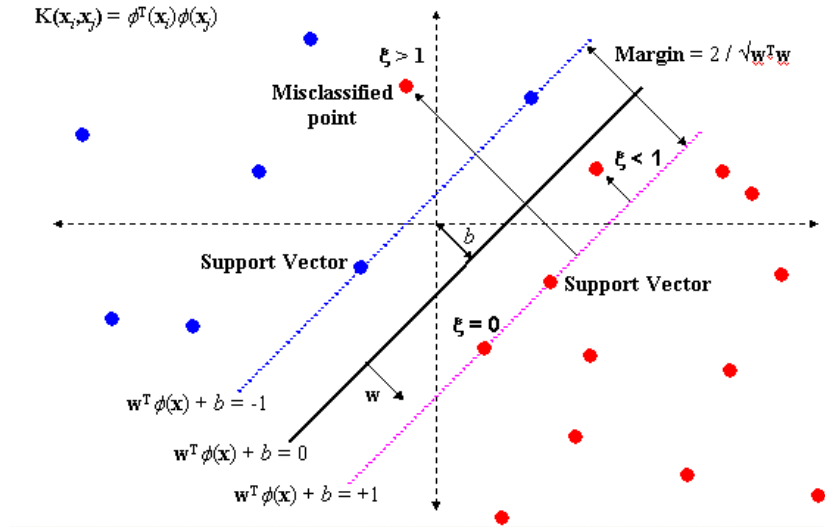


Figura 3.2: SVM-hiperplanul și vectorii suport

Lagrangianul trebuie minimizat după w , b și maximizat după $\alpha \geq 0$. (α_i sunt multiplicatori Lagrange):

$$\phi(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^l \alpha_i (y^i [\langle w, x^i \rangle + b] - 1) \quad (3.15)$$

Se poate arata ca ecuatia 3.15 se poate scrie si astfel:

$$\max_{\alpha} W(\alpha) = \max_{\alpha} - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{k=1}^l \alpha_k \quad (3.16)$$

Maparea liniară și funcțiile Kernel

După cum am mai spus, când nu se poate găsi o separare liniară satisfăcătoare SVM-urile pot mapa vectorul de intrare, x , pe un spațiu asociat, z , cu mai multe dimensiuni decât x . Astfel vom face maparea pe spațiul z și în acest spațiu asociat vom realiza separarea liniară.

Într-adevăr, există niște restricții în maparea neliniară care se poate aplica, dar, surprinzător, cele mai folosite funcții sunt acceptabile. Între acestea amintim pe cele polinomiale, așa numitele funcții RBF și câteva funcții

sinusoidale. Problema optimizării devine:

$$\max_{\alpha} W(\alpha) = \max_{\alpha} - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) + \sum_{k=1}^l \alpha_k \quad (3.17)$$

unde $K(x,y)$ este o funcție kernel (nucleu) care realizează maparea neliniară în spațiul asociat (feature space) și constrângerile rămân neschimbate.

3.4 Clasificatorul Bayes Naiv

Ideea de bază a algoritmului de învățare naivă este de a estima probabilitatea ca un document dat să aparțină unei anumite categorii. Clasificatorul Bayes naiv presupune independența cuvintelor, dar în ciuda acestei simplificări el are performanțe destul de bune.

În general instanțele care trebuie clasificate sunt caracterizate printr-un vector de atribute $\vec{a} = (a_1, a_2, \dots, a_n)$. Clasificatorul Bayes naiv asignează unei instanțe cea mai probabilă clasificare, sau o clasificare a posteriori dintr-un set finit de clase C .

$$c_{MAP} \equiv \operatorname{argmax} P(c|\vec{a}), \quad c \in C \quad (3.18)$$

care după ce este aplicată teorema lui Bayes devine

$$c_{MAP} \equiv \operatorname{argmax} P(c)P(\vec{a}|c), \quad c \in C \quad (3.19)$$

Probabilitatea posteriori $P(\vec{a}|c) = P(a_1, a_2, \dots, a_n|c)$ poate fi estimată direct din datele de antrenament dar nu sunt în general fezabil de interpretat dacă datele nu sunt vaste. Totuși presupunerea Bayes naivă că atributele ar fi condițional independente unele de altele - va da următorul rezultat:

$$P(a_1, a_2, \dots, a_n|c) = \prod_i P(a_i|c) \quad (3.20)$$

Astfel că această presupunere devine:

$$c_{NB} \equiv \operatorname{argmax} P(c) \prod_i P(a_i|c), \quad c \in C \quad (3.21)$$

[Mitchell, 1997]

În clasificarea de text se poate alege ca atribut poziția fiecărui cuvânt într-un document. Acest lucru înseamnă determinarea probabilității ca un anumit cuvânt w_k să apară la poziția j , dată fiind clasificarea țintă c_j , anume $P(a_j = w_k|c_j)$. Deoarece datele de antrenament sunt răzlețe introducem o nouă presupunere, anume: probabilitatea ca un anumit cuvânt w_k aflat la poziția j este identică cu probabilitatea ca același cuvânt să fie la o altă poziție m .

$$P(a_i = w_k|c_j) = P(a_m = w_k|c_j) \text{ pentru toți } i, j, k, m \quad (3.22)$$

Astfel estimăm probabilitatea $P(a_i = w_k|c_j)$ cu $P(w_k|c_j)$. Pentru a evita probabilitatea 0 se folosește aproximarea lui Laplace a probabilității.

$$P(w_k|c_j) = \frac{n_k + 1}{n + |Vocabulary|} \quad (3.23)$$

unde n_k reprezintă numărul de apariții al cuvântului w_k în toate documentele de clasa c_j ;

n_j reprezintă numărul total de poziții în documentul de clasă c_j ;

$|Vocabulary|$ reprezintă numărul de cuvinte distincte în toate documentele;

Cuvintele care nu se regăsesc în dicționar sunt ignorate. Eliminarea cuvintelor care apar frecvent sau a celor care sunt foarte rare pare un motiv destul de întemeiat. Cuvintele care apar foarte rar în documente este posibil să aibă un efect semnificativ la lungimea atributelor și predicțiile nu ar trebui să se bazeze pe niște observații rare. Eliminarea celor mai frecvente cuvinte este motivată de faptul că unele cuvinte, cum ar fi prepozițiile, s-ar putea să nu ofere informații utile. (Hovold 2005)

Capitolul 4

Experimente în detectarea automată a umorului și rezultatele obținute

4.1 Corpusurile

Pentru a observa evoluția algorimilor de învățare automată am folosit mai multe corpusuri.

1. Texte scurte

- **Datele pozitive** Am folosit corpusul de 6805 de glume culese manual de către Jonas Sjoberg [Sjobergh and Araki, 2007]. Lungimea textelor acestui corpus variază între 1 și 80, cu o medie de 13.
- **Datele negative** Au fost colectate urmărind să aibă aceeași structură ca cele umoristice (numărul de cuvinte și media trebuie să fie aproximativ aceleași atât pentru datele pozitive cât și pentru cele negative).

Pentru non-glume am folosit date din **American National Corpus**
<http://americannationalcorpus.org/OANC/>

Textele sunt structurate pe 6 categorii: *journal* (jurnal), *technical* (tehnic), *travel guides* (ghiduri turistice), *non fiction* (nonficțiune), *fiction* (ficțiune), *letter* (scrisoare).

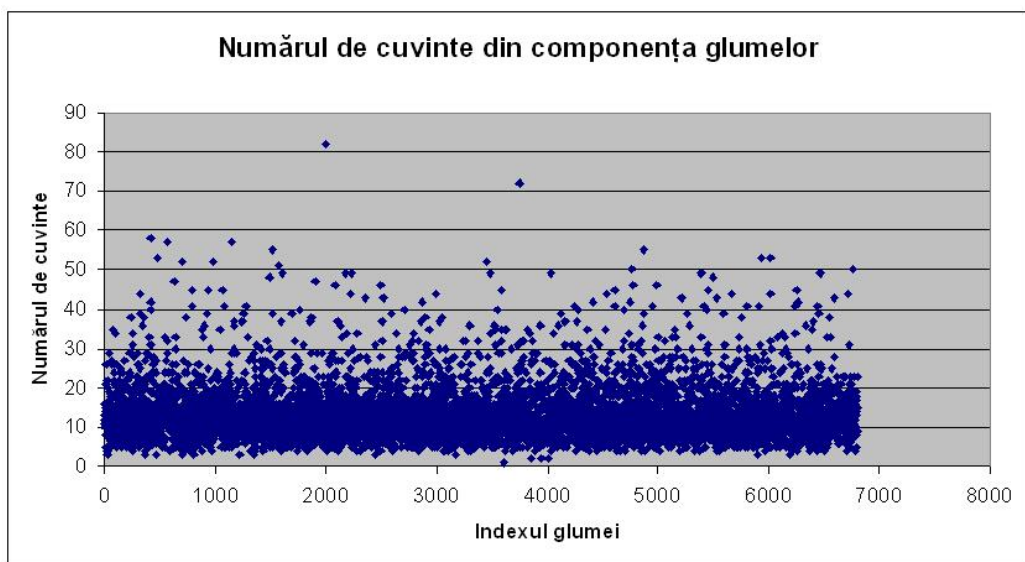


Figura 4.1: Numărul de cuvinte din componența glumelor

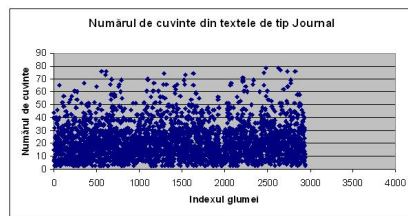
Datele au fost împărțite în propoziții și filtrate: au fost păstrate doar propozițiile între 1 și 80 de cuvinte.

Din figura se poate observa că majoritatea datelor pozitive sunt în intervalul 6-25. Astfel propozițiile din acest interval ar trebui să aibă o probabilitate mai mare de a fi alese pentru datele negative.

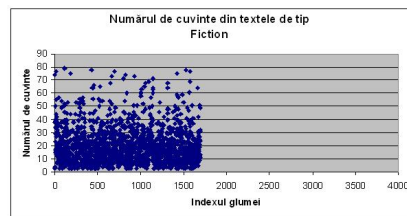
Dupa cum se poate observa din figura (media după numărul de cuvinte este 21, iar la datele umoristice este 13). Este nevoie de o nouă filtrare a datelor pentru a selecta mai multe date din intervalul 6-25. Dacă impunem

Clasificare	Media numărului de cuvinte	Numărul de texte
Journal	21.17453311	2945
Technical	15.25832609	3453
Travel Guides	20.37071886	2657
Non fiction	27.6681191	2754
Fiction	20.09540636	1698
Letter	18.62918248	2899

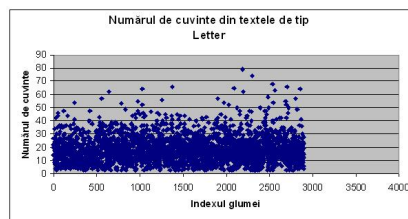
Tabela 4.1: Structura Datelor după prima filltrare



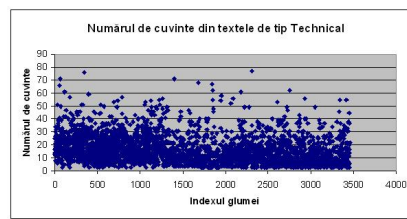
(a) Journal



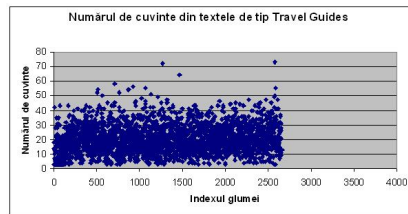
(b) Fiction



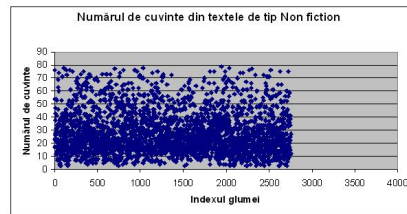
(c) Letter



(d) Technical



(e) Travel Guides



(f) NonFiction

Figura 4.2: Numărul cuvintelor pe categorii din componența textelor non-umoristice Corpus 1

Tip	Media
Journal	16.95258
Technical	14.61752
Travel Guides	16.74491
Non fiction	21.30893
Fiction	16.2572
Letter	16.61021

Tabela 4.2: Media numărului de cuvinte pentru datele nonumoristice

probabilitatea ca un text cu mai mult de 30 de cuvinte să fie 0.2, atunci obținem rezultatele din tabelul 4.2.

Pentru a nu avea texte doar din cele 6 categorii al corpusului ANC am mai introdus un corpus de 1000 de propoziții scurte cu media numărului de cuvinte egală cu 10. În concluzie avem pentru datele nonumoristice 8774 de texte cu o medie a numărului de cuvinte de 15.4.3

Datele pozitive și negative pot fi folosite pentru a învăța în mod automat modele computaționale de recunoaștere a umorului și pentru a evalua performanțele acestor modele. Se știe că prezența unui număr mare de date de antrenament au potențialul de îmbunătăți acuratețea procesului de învățare și, în același timp, oferă informații despre modul în care datele de diumenisiune mai mare pot afecta precizia clasificării. [Witten, 2000]

2. Computers Humor/Computers NonHumor

Pentru datele pozitive am folosit date de pe site-ul ¹ (un set de întâmplări amuzante ce implică prezența calculatoarelor), iar datele negative de pe site-ul² - ce conține articole despre calculatoare. În total am colectat 1566 de exemple pozitive și 1566 de exmple negative. Un exemplu de dată pozitivă

I was browsing the Internet when my friend came over and said he made a website. He told me to go a particular URL. When I went there, though, the browser said it was invalid. So I went to Google to search for it, and when I got to Google, he said, "Oh yeah, that's my web site."

iar unul de dată negativă

¹<http://linkenlim.vox.com/library/post/>

²<http://www.cs.cmu.edu/afs/cs/project/theo-11/www/>

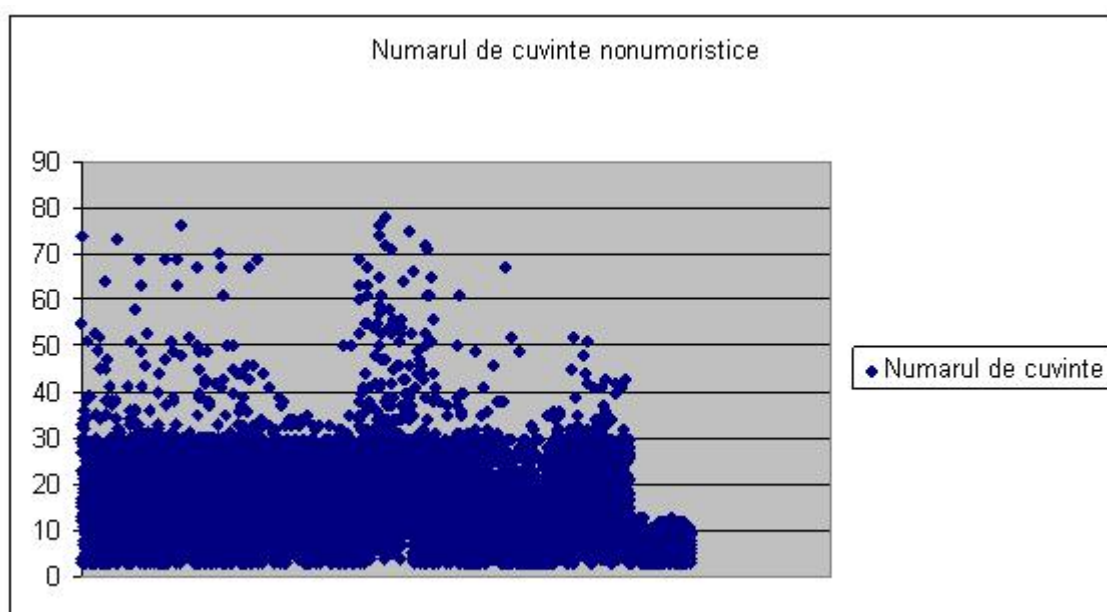


Figura 4.3: Distribuția numărului de cuvinte pentru datele nonnumoristice după a doua filtrare

Has anybody bought a Pentium motherboard? If so or you where I can buy it, please send me a E-mail. Thank you in advance. Pentium processors / motherboards are not available to the general public as of yet. Intel has released them to companies such as Gateway and Dell to do testing, etc. It'll be a while...

3. Jewish humor

Pentru a vedea dacă umorul poate fi clasificat am creat o bază de date cu glume evreiești împărțite în 14 categorii ³. *Bar mitzvah* (31 de documente), *Bris and Mohel* (18 documente), *Divorce* (19 documente), *Donations* (7 documente), *Drinks and alcohol* (15 documente), *Food* (43 de documente), *Golf* (14 documente), *Israel* (65 de documente), *Getting married* (11 documente), *Pesach* (20 de documente), *Rabbis* (152 de documente), *Seventieth birthday* (15 documente), *Shadchen* (7 documente), *Shmuck* (12 documente), *Wedding and anniversaries* (14 documente). Exemplu: *Rabbi Morris has just resigned and Issy, the shul president, goes to visit him. "Rabbi," Issy says, "I've just heard the news. I'm really sorry that you've decided to leave us." "Don't worry," says Rabbi Morris, "you'll have nothing to worry about. I'm going to recommend a successor whom I believe will be better than me." "But that's exactly what's worrying me," says Issy, "your predecessor told me exactly the same thing."*

4. Dirty humor/Non dirty humor

Din baza de date a lui Jonas Sjoberg [Sjobergh and Araki, 2007] s-au selectat manual 1000 de glume pentru adulți pentru a vedea dacă se poate face clasificarea umorului.

4.2 Experimentele

1. Fiecare cuvânt este un atribut sau rădăcina unui cuvânt este un atribut

Textele au fost parsate, s-au păstrat doar cuvintele care au frecvența de apariție între 3 și 200. Pe toate cele 4 tipuri de corpusuri s-au încercat două tipuri de abordări: în prima abordare fiecare cuvânt este un atribut iar în

³<http://www.awordinyoureye.com/category%20jokes%20home%20page.html>

Categoria	Rădăcina mai multor cuvinte un atribut	Orice cuvânt un atribut
Computers	5536	8252
Jewish	1561	1815
Dirty	988	920
Oneliners	4658	6113

Tabela 4.3: Numărul de componente al vectorilor

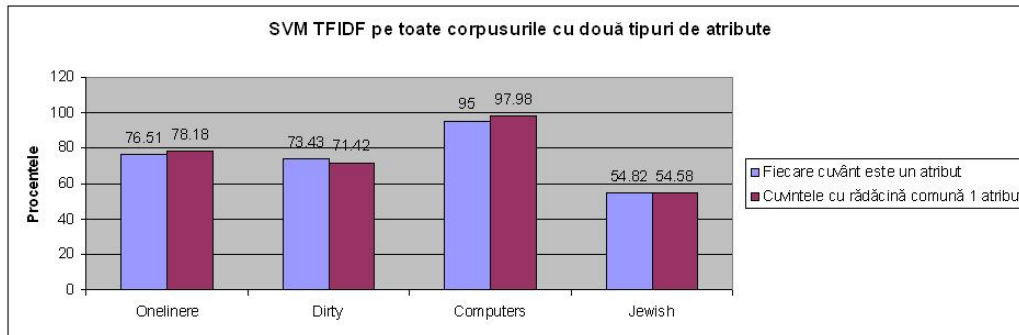


Figura 4.4: Experimentul I

cea de-a doua două cuvinte care provin din aceeași rădăcină sunt considerate cu un singur atribut (un exemplu de rădăcină ar fi *educ* care poate cuprinde printre altele *educate*, *education*, etc.). Pentru calculul valorii atributelor s-a folosit TFIDF.

S-a aplicat SVM cu factor de crossvalidare de 5%. Graficul 4.4 arată rezultatele obținute (procentele reprezintă acuratețea⁴ experimentului).

1. Adăugarea de noi attribute euristice

Vrem să vedem dacă acuratețea clasificării crește prin adăugarea de noi caracteristici la vectorul text. Pentru fiecare text se calculează următoarele attribute:

- **(a) Ambiguitatea medie:** Majoritatea glumelor tind să fie construcții ambigue prin cuvintele constituite. Utilizând Wordnetul⁵ am calculat media ambiguității medii a unui text (se adună numărul de synseturi pentru fiecare cuvânt/ numărul de cuvinte).

⁴raportul dintre numărul de texte corect clasificate/numărul total de texte*100

⁵<http://www.rednoise.org/rita/wordnet/documentation/index.htm>

- **(b)Antonimia medie:** Conform teoriei de nepotrivire unele construcții vor să inducă opusul mesajului citit. Utilizand Wordnetul ⁶ am calculat media antonimiei medii a unui text (se adună numărul de antonime pentruu fiecare cuvânt/ numărul de cuvinte).
- **(c)Cuvinte murdare:** Folosind o bază de cuvinte murdare, ⁷ am numărat cuvintele ”‘murdare’” din text.
- **(d)Ambiguitatea maximă:** Utilizand Wordnetul ⁸ am calculat numărul de sysnseturi pentru fiecare cuvânt iar apoi am determinat maximul valorilor obținute.
- **(e)Scorul maxim de antonimie:** Utilizând Wordnetul ⁹ am calculat valoarea maximă a numărului de antonime pe care le admit cuvintele din text.
- **(f)Lungimea maximă a proverbelor:** Folosind o bază de date de proverbe colectate de pe Internet ¹⁰ am determinat numărul maxim de cuvinte care se găsesc în proverbe și în textul căruia îi calculăm atributul.
- **(g)Numărul de semne de punctuație:** Pentru a detrmina influența semnelor de punctuație pentru un text în clasificare s-a calculat numărul de semne de punctuație al textului.
- **(h)Numărul mic de proverbe :** Am calculat numărul de proverbe cu proprietatea că 3 cuvinte valide (cu lungimea mai mare decât 3) din acele proverbe se întâlnesc și în textul analizat.
- **(i)Numărul mare de cuvinte din proverbe:** Am calculat numărul de proverbe cu proprietatea că 6 cuvinte valide (cu lungimea mai mare decât 3) din acele proverbe se întâlnesc și în textul analizat.
- **(j)Numărul de repetiții::** numărul de repetiții de cuvinte din textul studiat.

⁶<http://www.rednoise.org/rita/wordnet/documentation/index.htm>

⁷<http://www.rimboy.com/words/>

⁸<http://www.rednoise.org/rita/wordnet/documentation/index.htm>

⁹<http://www.rednoise.org/rita/wordnet/documentation/index.htm>

¹⁰<http://en.wikiquote.org/wiki/>

Atributul folosit	Rezultatul obținut
Unigrame	76.51%
Bigrame	78.36%
Trigrame	79.01%

Tabela 4.4: Numărul de componente al vectorilor

Rezultate: Am folosit toate atributele anterior calculate în care fiecare rădăcină este un atribut. Fiecare din atributele (a)...(j) au fost inserate câte unul. Rezultatele obținute sunt cele din figura 4.5 (folosind SVM-uri): Se poate observa că toate atributele cresc acuratețea clasificării, cel mai

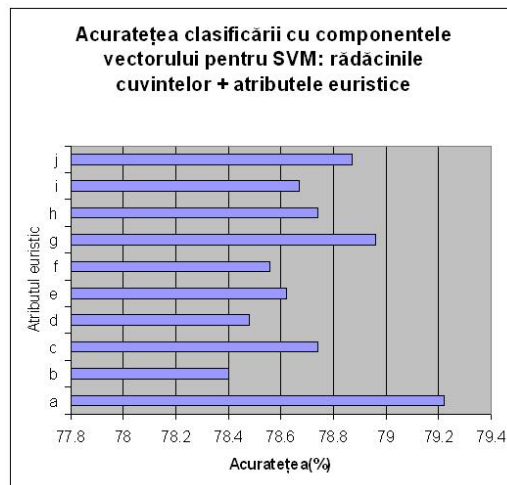


Figura 4.5: Experimentul II

important fiind atributul mediu de ambiguitate. Am creat vectori în care componente sunt rădăcinile cuvintelor textului și toate atributele euristice. Rezultatele sunt optimiste, cu un procent de acuratețe de 80.08% și rulând doar atributele euristice ca un vector obținem o acuratețe de 66.05% rulat pe SVM-uri.

Influența n-gramelor în clasificarea de texte

Pentru primul corpus am păstrat doar cuvintele care au o frecvență de apariție între 3 și 200. Rulând pe SVM-uri cu o acuratețe de 5% obținem: Se

poate observa că acuratețea crește odată cu creșterea numărului de grupuri de cuvinte ale n-gramelor.

Bayes Naiv

Am rulat algoritmul Bayes naiv iar rezultatele au arătat un procent de 50% de acuratețe pentru primul corpus, și de 45% pentru corpusul de glume evreiești.

Arbori de decizie

Atributele stilistice au fost rulate și cu arbori de decizie de unde s-a obținut o acuratețe de 51% demonstrându-se astfel că SVM-urile sunt o soluție mai bună pentru acest tip de clasificare.

Capitolul 5

Interpretarea rezultatelor

5.1 Observații privind rezultatele obținute

Rezultatele obținute sunt comparabile cu cele în domeniu chiar cu unele îmbunătățiri(66.05%) față de 60% folosind atribute stilistice[Mihalcea and Strapparava, 2006]. Una din diferențe ar putea fi faptul că experimentele noastre folosesc SVM-uri în locul arborilor de decizie și de asemenea folosim mai multe tipuri de atribute.

Clasificarea de texte este posibilă atâta timp cât avem suficient de multe date din fiecare categorie care să fie specifice categoriei respective. După cum s-a putut observa pentru datele din corpusul cu glume evreiești nu am reușit să obținem niște rezultate prea bune deoarece cele 14 categorii aveau prea puține date de antrenament. Deși am încercat un proces de clusterizare(K-means) pentru primul din 12000 de documente, câte 6000 din fiecare categorie, doar 200 au fost puse într-un cluster, iar restul în altul, rezultatele nu sunt prea optimiste (componente ale vectorilor sunt atributele structurale).

5.2 Posibilități de îmbunătății programul de recunoaștere al umorului în texte

- Crearea de corpusuri pentru alte limbi decât limba engleză.
- Crearea de corpusuri pe grupuri de utilizatori: clasificare de umor după vârstă, nivel de cultură, localizare.

- Testarea de noi atribute: de exemplu ar fi numărul de synseturi ale ultimului cuvânt din text, sau numărul maxim de cuvinte din synseturi diferite.
- Majoritatea glumelor se bazează pe existența unui cuvânt care face trecerea de la momente amuzante la cele neamuzante. Dacă s-ar găsi poziția aceluși cuvânt am putea exploata gluma mai bine. O modalitate ar fi să fie adnotate câteva glume și să se încerce căutarea de atribute care ar putea contribui la determinarea locului cuvântului respectiv (câteva din atribute ar putea fi: poziția verbelor, numărul de substantive).
- Încercarea de clusterizare a glumei: conform teoriei de nepotrivire există două cadre care sunt prezente într-o glumă de aceea s-ar putea încerca împărțirea cuvintelor existente pe clustere în funcție de distanță dintre synseturi. Dacă un cuvânt apare în mai multe astfel de clustere se consideră ca este posibil ca textul respectiv să fie amuzant.

5.3 Concluzii

”Umorul se construiește pe aspirațiile și limitările omului. Nu există mai multă logică ca în umor. Pentru că umorul este adevăr.” (Victor Borge) Odată cu apariția calculatoarelor a apărut și dorința de a face calculatorul să înțeleagă limbă vorbită și scrisă. Primele cercetări au condus la rezultate promițătoare și era de așteptat ca în anii cei mai apropiați calculatorul va fi în stare să prelucreze atât informația numerică, cât și cea textuală. Totodată, pe parcursul cercetărilor, s-a conturat complexitatea limbajelor naturale, în special ambiguitatea elementelor lor, specificul utilizării în societate, etc. Umorul este una din cele mai dificile caracteristici de recunoscut nu doar pentru un calculator ci și pentru oameni în general. Am încercat să arătăm că umorul poate fi recunoscut în texte, dar până să recunoască diferite tipuri de umor ale oamenilor este un drum lung.

Proiectul prezentat demonstrează că folosind învățarea automată, mașina este capabilă să recunoască umorul individual.

Există anumite elemente care alcătuiesc registrul umoristic al individului care va fi valabil și pentru:

- Nerecunoașterea ca persoană cu umor. Dacă în experiențe anterioare o persoană a avut succese sau eșecuri în producerea umorului, acest

lucru îi va afecta motivația actuală în dezvoltarea abilității de generare a umorului. La nivelul individului dacă acesta va simți că îi sunt recunoscute glumele va încerca să mai glumească.

- Recunoașterea umorului altora. De asemenea, o persoană poate fi motivată să-și dezvolte abilitatea de a înțelege gluma numai dacă a primit feed-back-uri pozitive în experiența anterioară. Iar dacă calculatorul va primi feedback despre modul cum a dat răspunsul la o întrebare va putea să se reantreneze cu noile cunoștințe.
- Aprecierea umorului se referă la atitudini. S-a observat că atitudinile către oamenii cu umor sunt legate de atitudinile față de umor, astfel că, dacă cineva afirmă că ”oamenii care glumesc chiar încearcă să mă manipuleze”, acest lucru reflectă concepția persoanei atât asupra utilizării umorului, cât și asupra umorului însuși. Astfel sunt oameni care resping umorul considerându-l o modalitate de manipulare.
- Râsul reprezintă un răspuns comportamental care poate fi legat sau nu de simțul umorului. Întrucât sunt persoane care râd și fără să se amuze și se amuză fără să râdă, râsul poate fi una din reacțiile posibile legate de simțul umorului. Cercetări suplimentare au relevat că umorul și râsul pot să aibă nu numai efecte pozitive asupra sănătății psihologice, dar chiar pot avea efecte negative [Kelly, 2003].
- Concepția persoanei, de asemenea, poate constitui un element al simțului umorului, în special atunci când acea concepție cuprinde o apreciere a absurdităților vieții.
- Utilizarea umorului ca un mecanism adaptativ reprezintă o componentă pozitivă a simțului umorului. A fi capabil de detașare față de probleme și de a le putea trata cu umor reprezintă un mijloc de protecție împotriva evenimentelor nefavorabile ale vieții. Umorul de autodeprecieri, dar cu autoacceptare reprezintă unul din cele mai mature mecanisme adaptative, însă umorul pe seama altora (ostil, agresiv) reprezintă un mecanism nevrotic de adaptare cu efect negativ în plan interpersonal.

O modalitate prin care s-ar putea realiza acest lucru ar fi prin analiza-manifestării umorului unui individ față de persoane diferite. Un exemplu ar fi arhiva de mesaje (de exemplu de tip Messenger) ale individului. O analiză

mai profundă ar putea releva ce se întâmplă cu evoluția umorului individual. Noutatea lucrării constă în primul rând în încercarea, pentru prima dată, a unei clasificări a umorului obținând o acuratețe de 73% pentru categoria dirty/non dirty humor și 54.8% pentru tipuri de bancuri evreiești. De asemenea utilizarea doar de atribute stilistice a dus până la o acuratețe de 66.7%. Experimentele au demonstrat deasemenea că n-gramele cresc acuratețea clasificării de la 76.5 pentru unigrame până la 79.8 pentru trigrame. Cel mai bun rezultat este obținut combinând atributele structurale cu cele stilistice, obținând o acuratețe de 80.08% prin aplicarea SVM-urilor cu factor de cross-validare de 5%. Conform informațiilor noastre este primă lucrare în limba română care abordează umorul computațional, iar prin posibilitățile de extindere mai sus menționate acuratețea algoritmului de clasificare va crește.

Bibliografie

- Savatore Attardo and Victor Raskin. Non-literality and non-bona-fide in language: An approach to formal and computational treatments of humor. *Pragmatics and Cognition*, 1994.
- Nancy Baym. The performance of humor in computer-mediated communication. *Journal of Computer-Mediated Communication, Volume 1 Issue 2*, 1995.
- Henri Bergson. *Teoria râsului*. Institutul European, 1992.
- Kim Binsted and Graeme Ritchie. Computational rules for generating punning riddles. *HUMOR, the International Journal of Humor Research, Mouton de Gruyter, volume 10-1*, 1997.
- Kim Binsted, Benjamin Bergen, Seana Coulson, Anton Nijholt, Oliviero Stock, Carlo Strapparava, Graeme Ritchie, Ruli Manurung, Helen Pain, Annalu Waller, and Dave O'Mara. Computational humor. *IEEE Intelligent Systems*, 21(2):59–69, 2006. ISSN 1541-1672.
- Justin McKay Binsted-McKay. Generation of idiom-based witticisms to aid second language learning. *Proceedings of the April Fools Day Workshop on Computational Humour, Trento, Italy*, 2000.
- Leah Black and Denise Forro. Humor in the academic library : You must be joking! or, how many academic librarians does it take to change a light-bulb? *Michigan State University Librarians at Michigan State University*, 1999.
- Elmer Blistein. Theories of humour. *Encyclopedia americana*, 1964.
- Susan Dumais. Using SVM for text categorisation. *Decision Theory and Adaptive Systems Group*, 1998.

- Sigmund Freud. *Jokes and their Relation to the Unconscious*. 1957.
- Ken Goldberg, T. Roeder, D. Gupta, and C. Perkins. Eigentaste: A constant time collaborative filtering algorithm. Technical Report UCB/ERL M00/41, EECS Department, University of California, Berkeley, 2000. URL <http://www.eecs.berkeley.edu/Pubs/TechRpts/2000/3869.html>.
- William Hazlitt. *Lectures on the English Comic Writers*. Wiley and Putnam, 1963.
- John Hewitt. *The Architecture of Thought: A New Look at Human Evolution*. Holmhurst House, Beds., 2002.
- Martin Rod.Anton Puhlik-Doris Patricia Larsen Gwen Gray Jeanette Weir Kelly. Individual differences in uses of humor and their relation to psychological wellbeing. *Journal of Research in Personality*, 2003.
- Craig McDonough. Mnemonic string generator:software to aid memory of random password. Technical report, CERIAS Technical Report. West Lafayette, IN: Purdue University., 2001.
- Rada Mihalcea. Learning to laugh (automatically):computational models for humor recognition. *Computational Intelligence*, 22:126–142(17), May 2006.
- Rada Mihalcea and Stephen G. Pulman. Characterizing humour: An exploration of features in humorous texts. In *CICLing*, pages 337–347, 2007.
- Rada Mihalcea and Carlo Strapparava. Technologies that make you smile: Adding humor to text-based applications. *IEEE Intelligent Systems*, 21(5):33–39, 2006. ISSN 1541-1672. doi: <http://dx.doi.org/10.1109/MIS.2006.104>.
- Marvin Minski. *Jokes and the relation of the cognitive unconscious*. 1981.
- Tom M. Mitchell. *Machine Learning*. McGraw-Hill, New York, 1997.
- Sayan Mukherjee and Vladimir Vapnik. Multivariate density estimation: an SVM approach. Technical Report AIM-1653, AI Memos (1959 - 2004), 1999.

- M.P. Mulder and Nijholt Anton. Humour research: State of the art. Technical report, Deliverable IST Programme on Future and Emerging Technologies. University of Twente, 2001.
- Anton Nijholt. Conversational agents a little humour too. *IEE Intelligent Systems*, pages 22–26, 2006.
- Anton Nijholt. Conversational agents, humorous act construction, and social intelligence. *University of Hertfordshire*, pages 1–8, 2005. ISBN=1-902956-45-X.
- Itvan Pilasyz. Text categorization and suport vector machine. *Departmenet of Measurment and Information System*, 2005.
- Jason Rutter. Stand-up as interaction: Performance and audience in comedy venues. *University of Salford*, 1997.
- Jonas Sjobergh and Kenji Araki. Recognizing humor without recognizing meaning. 2007.
- Tesu Solomivici. *5000 de ani de umor evreiesc - o antologie subiectivă*. Editura Tesu, București, 2002.
- Stock and Carlo Strapparava. Computational humor and prospects for advertising. *Rob Milne: A Tribute to a Pioneering AI Scientist, Entrepreneur and Mountaineer*. IOS Press., 2006.
- Oliviero Stock and Carlo Strapparava. Getting serious about the development of computational humor. *IJCAI*, pages 59–64, 2003.
- Thomas C. Veatch. A theory of humor. *International Journal of Humor Research*, 1998.
- Ian H. Witten. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations (The Morgan Kaufmann Series in Data Management Systems)*. 2000.