



FILTRAREA MESAJELOR ÎN SERVICII DE MICRO-BLOGGING

LUCRARE DE LICENȚĂ

Absolvent: **Adrian Cornel PETRI**

Coordonator
științific: **Șef lucrări ing Cosmina IVAN**

2014



FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE
DEPARTAMENTUL CALCULATOARE

DECAN,
Prof. dr. ing. Liviu MICLEA

DIRECTOR DEPARTAMENT,
Prof. dr. ing. Rodica POTOLEA

Absolvent: **Adrian Cornel PETRI**

**FILTRAREA MESAJELOR ÎN SERVICII
DE MICRO-BLOGGING**

1. **Enunțul temei:** *Implementarea unei aplicații care să filtreze mesajele aparținând unui serviciu de micro-blogging.*
2. **Conținutul lucrării:** *Introducere, Obiectivele proiectului, Studiu bibliografic, Analiză și fundamentare teoretică, Proiectare de detaliu și implementare, Testare și validare, Manual de instalare și utilizare, Concluzii, Bibliografie, Anexă*
3. **Locul documentării:** Universitatea Tehnică din Cluj-Napoca, Departamentul Calculatoare
4. **Consultanți:**
5. **Data emiterii temei:** 1 februarie 2014
6. **Data predării:** 11 septembrie 2014

Absolvent: _____

Coordonator științific: _____

Cuprins

Capitolul 1. Introducere	1
1.1. Contextul proiectului	1
1.2. Provocări.....	2
Capitolul 2. Obiectivele Proiectului	3
2.1. Specificația proiectului	3
2.2. Obiective generale	3
2.2.1. Cerințe funcționale	4
2.2.2. Cerințe non-funcționale	5
2.2.3. Cerințe de implementare.....	6
2.3. Conținutul lucrării.....	6
Capitolul 3. Studiu Bibliografic.....	8
3.1. Clasificarea și filtrarea textelor folosind algoritmi clasici.....	8
3.1.1 SIFT	8
3.1.2 Fast Data Finder	8
3.1.3 Info Scope	9
3.2 Clasificatori de tip machine learning	11
3.2.1 Clasificatori de tip arbori de decizie	11
3.2.2 Clasificatori de tip reguli de decizie	11
3.1.1. Rețele neuronale	12
3.1.2. Asamblorul de clasificatori.....	12
3.2 Micro-blogging	13
3.2.1. Microblogging-ul local	13
3.2.2. Microblogging-uri uzuale	14
Capitolul 4. Analiza și fundamentarea teoretică	17
4.1. Clasificare bazată pe conținutul mesajelor	17
4.1.1. Module de prelucrare pre-filtrare.....	17
4.2. Modulul de comparare a cuvintelor	21
4.3. Algoritmul de clasificare	22
4.4. Pseudocodul algoritmului	23
4.5. Protocolul HTTP.....	25
Capitolul 5. Proiectare de Detaliu si Implementare	27
5.1. Modul de definire al sistemului	27

5.1.1.	Cursul de execuție al aplicației	27
5.1.2.	Arhitectura conceptuală generală	28
5.2.	Modelul cazurilor de utilizare.....	30
5.2.1.	Extragerea url-ului din mesaj	30
5.2.2.	Traducerea mesajelor care nu sunt în engleză	30
5.2.3.	Eliminarea semnelor de punctuație și a cuvintelor uzuale	31
5.2.4.	Traducerea cuvintelor în mod corect chiar dacă acestea nu conțin aceleași forme ale cuvintelor.	31
5.3.	Diagrama cazurilor de utilizare	32
5.4.	Diagrama de pachete.....	32
5.5.	Diagrama de clase.....	33
5.5.1.	filter.xhtml	34
5.5.2.	SearchBean	34
5.5.3.	PostReceiver	34
5.5.4.	DatabasePost și DatabaseConnection	34
5.5.5.	Post	34
5.5.6.	PostTransformers	35
5.5.7.	InformationExtractor	35
5.5.8.	EnglishTranslator.....	36
5.5.9.	Common words și comparator.....	36
5.6.	Diagrama secvențială.....	37
5.7.	Tehnologii folosite în aplicație	39
5.7.1.	Tehnologiile de front-end	39
5.7.2.	Tehnologiile de back-end	40
5.7.3.	Alte tehnologii	44
5.7.4.	Integrarea cu alte sisteme	44
Capitolul 6. Testare și validare		47
6.1.	Validarea datelor.....	47
6.2.	Testarea aplicației	47
Capitolul 7. Manual de instalare și utilizare		50
7.1.	Instalarea aplicației	50
7.1.1.	Aplicația de micro-blogging	50
7.1.2.	Aplicația de dicționar.....	51
7.1.3.	Aplicația de filtrare de mesaje	51

7.2. Utilizarea aplicației	51
Capitolul 8. Concluzii	53
8.1. Contribuții personale	53
8.2. Dezvoltări ulterioare	54
Bibliografie	55

Capitolul 1. Introducere

1.1. Contextul proiectului

Internetul este o rețea vastă de calculatoare care leagă milioane de rețele mai mici din întreaga lume (inter = între, net = rețea). Internetul este o modalitate de comunicare a calculatoarelor. O rețea este un grup de calculatoare și echipamente de calcul conectate pentru a partaja informații și resurse. Calculatoarele dintr-un birou, sunt des interconectate pentru a putea utiliza aceleași fișiere și imprimante. Toate calculatoarele legate la Internet pot schimba informații între ele. Este la fel de ușor de comunicat cu un calculator din cealaltă parte a lumii ca și cu unul din aceeași cameră. Înainte de a putea schimba informații, două calculatoare de pe Internet trebuie să se "gasească" reciproc și să comunice într-un limbaj accesibil amândouă.

Odată cu dezvoltarea internetului a apărut nevoia omului de a discuta asupra anumitor teme cu alți utilizatori sau de a afla păreriile lor. De asemenea omul a simțit nevoie de a-și putea scrie gândurile online, atât sub anonim cât și sub numele lor adevărat cu scopul de a se descarca. În acest context a apărut termenul de blogging. Definiția pe care Wikipedia o dă blogului este: "Blogul este un website în care intrările sunt scrise în ordine cronologică și afișate în ordine cronologică inversă. Cuvântul blog poate avea și funcția gramaticală de verb, însemnând menținerea sau adăugarea de conținut într-un blog." Primul blog a apărut în anul 1997 și a prins un contur tot mai mare ajungând ca în zilele noastre cuvântul „blog” în limba engleză să fie atât substantiv cât și verb(to blog something)[1].

Deși blog-ul era foarte popular, postările aici erau de obicei lungi iar utilizatorii pierdeau mult timp căutând informații utile. În acest context apare micro-bloggingul. Micro-bloggingul este un mediu care există sub formă de blogging. Un microblog diferă față de un blog tradițional prin conținutul său care de obicei este mai mic atât în mărimea mesajului dar și a fișierului. Microblogul permite utilizatorilor să facă schimb de elemente mici de conținut, cum ar fi fraze scurte, imagini individuale, sau link-uri video. Aceste mesaje mici sunt numite micro-postări.

Ca și în cazul blogging-ului tradițional, micro-bloggerii pot posta mesaje despre teme care variază de la cele simple, cum ar fi "ceea ce fac eu acum," la cele tematice, cum ar fi "mașini sport". De asemenea există și micro-bloguri comerciale, pentru a promova site-uri web, servicii și produse, etc.

Utilizatorii aveau acum la dispoziție microblogging-ul și puteau căuta sau posta mesaje în interiorul serviciilor. Totuși căutarea presupunea ca utilizatorul să verifice toate mesajele și să aleagă din interiorul lor doar pe acelea care îl interesau. Deși în serviciile de micro-blogging postările erau împărțite pe conținut, căutarea în interiorul lor presupunea un timp îndelungat. În acest context apare filtrarea de mesaje după conținut. În acest tip de sortare fiecare utilizator se presupune a lucra independent unul față de celălalt. Ca și rezultat sortarea după conținutul mesajului poate exploata doar postările. După aplicarea acestui filtru de mesaje se vor afișa doar informațiile relevante față de mesajul nostru.

1.2. Provocări

Pentru a putea filtra mesajele din cadrul micro-bloggingului va trebui să identificăm cele mai potrivite soluții pentru realizarea unui set de provocări. În cele ce urmează vor fi enumerate câteva provocări identificate în cadrul proiectului de față:

Mesaje scurte. O provocare poate fi considerată limitarea lungimii textului. În micro-bloggingul folosit pentru filtrare limitarea este de 500 de caractere. În termenii clasificării textele scurte conțin date distribuite de aceea poate fi considerată o provocare clasificarea și filtrarea lor.

Identificarea cuvintelor cheie. Este foarte important să detectăm cuvintele cheie ale mesajului pentru ca ulterior să putem filtra mesajele ca relevante și irelevante față de topicul nostru. Mesajul nostru poate fi unul foarte general („Tehnologia Informației”), unul general („Baze de date”), sau chiar unul foarte specific („MySQL Database”).

Limbaje informale. O altă problemă ar putea fi considerată structura informală a serviciilor de micro-blogging. Aceasta este mai puțin structurată, poate conține abrevieri datorită limitării textului iar utilizatorii modifică anumite cuvinte sau le combină cu emoticoane. Așadar vocabularul este foarte mare iar surse externe precum Wikipedia nu vor fi suficiente pentru a aduce destule informații despre text. De asemenea este necesar să putem identifica care cuvinte sunt comune și care cuvinte cheie sunt utile pentru clasificarea textelor. Datorită limbajelor informale, numărul cuvintelor care trebuie blocat să intre în lista de mesaje trebuie să fie mai mare.

Schimbări ale vocabularului. Vocabularul din cadrul serviciului de micro-blogging se schimbă în mod constant datorită cuvintelor și a frazelor noi. De exemplu, filme și produse noi pot apărea ca și denumiri și pot deveni foarte populare. Astfel clasificarea de texte nu va putea fi una statică, va trebui să răspundă la schimbările din cadrul vocabularului.

Limbaje diferite. Micro-bloggingul este folosit de utilizatori de pe întreaga suprafață a planetei și de aceea mai multe limbi vor fi utilizate în postări.

Așadar sistemul de filtrare va trebui să fie unul dinamic, ar trebui să urmărească îndeaproape sistemul de micro-blogging și să răspundă la mesaje. Pentru clasificarea textelor ar trebui depășită provocarea cu referire la limita textului, iar cuvintele să poată fi extrase într-o modalitate informală și multilinguală. În final sistemul ar trebui să fie capabil să identifice cuvintele cheie și să filtreze mesajele irelevante.

Capitolul 2. Obiectivele Proiectului

2.1. Specificația proiectului

Proiectul de față propune implementarea unei aplicații care să se ocupe cu preluarea, prelucrarea și filtrarea mesajelor provenite dintr-un serviciu de micro-blogging. Aplicația trebuie să fie capabilă să dividă mesajele din micro-blogging în două mari categorii: relevante față de mesajul dat sau irelevante față de mesajul dat. De asemenea aplicația trebuie să fie capabilă să filtreze mesajele chiar dacă acestea sunt în altă limbă, au în componența sa un url sau conțin cuvinte doar asemănătoare cu cele folosite în mesajul original.

2.2. Obiective generale

Obiectivul principal este de a construi un filtru de mesaje care să fie capabil să sorteze mesajele provenite dintr-un serviciu de micro-blogging. Pentru aceasta ne vom concentra în principal asupra conținutului mesajelor aduse. Prin gruparea lor în mesaje relevante sau mesaje irelevante vom putea obține doar mesajele care vor avea același subiect de discuție sau idee principală cu mesajul dat. Studiile arată faptul că, două mesaje au idei aproximativ asemănătoare dacă cuvintele folosite în mesaje au aproximativ același câmp semantic.[2]

În primul rând ne vom concentra asupra extragerii cuvintelor principale din interiorul mesajului. Având în vedere faptul că, atunci când scriem un mesaj, folosim multe cuvinte de legătură, prepoziții, pronume, conjucții și alte părți de vorbire care nu sunt relevante, acestea vor fi excluse din algoritm. Pentru a face acest lucru, vom extrage din mesajele noastre cele mai utilizate 250 de cuvinte din limba engleză. Conform studiilor, masa de bază a unui vocabular este în principal formată din puține cuvinte. Dacă extrăgeam mai multe cuvinte riscam faptul că, mulți termeni relevanți pentru mesajul meu să nu fie incluși.

În ceea ce privește cuvintele relevante care vor fi folosite, acestora li se va da o pondere. Această pondere reprezintă în principal importanța pe care cuvântul o va avea în mesajul respectiv. Cu cât cuvântul apare mai des în mesaj, cu atât el va fi mai important. Ponderea cuvântului se va calcula doar în mesajul respectiv iar durata de execuție a calculului va fi una foarte rapidă.

Micro-bloggingul asupra căruia se va face filtrarea permite doar mesaje cu lungimea de maxim 500 de cuvinte. Astfel, pentru a se face mai bine înțeles utilizatorul folosește link-uri pe care le adaugă în mesaj. Prin aceste link-uri se va redirecționa cititorul către o altă pagină, pagină care va conține mai multe informații în legătură cu subiectul precizat. Algoritmul utilizat își propune să preia toată informația utilă din link-ul respectiv și să o adauge mesajului. Astfel limitarea de 500 de cuvinte va putea deveni mult mai mare. De notificat ar mai putea fi și faptul că mesajul se va adăuga mesajului inițial, informația din acesta rămânând utilă în continuare. Prin informația provenită din link-uri vom putea îmbunătăți algoritmul de clasificare.

O altă provocare a proiectului ar fi mesajele în altă limbă. Întodeauna comunicarea a fost limitată la doar o singură limbă. Evoluțiile actuale din domeniul tehnologic au reușit însă să rupă această barieră și să facă traduceri între texte dintr-o limbă în alta. În cazul în care mesajul din micro-blogging este într-o altă limbă decât engleza, se va face traducerea lui iar după traducere se va aplica algoritmul de filtrare.

De foarte puține ori, în momentul în care comparăm două mesaje, se întâmplă să găsim aceleași cuvinte în exact aceeași formă. Astfel, deși probabil rădăcina cuvântului va fi aceeași, cuvintele vor diferi. O altă problemă des întâlnită se referă la faptul că, pentru a spune același lucru, persoane diferite folosesc cuvinte diferite. Algoritmul de filtrare își propune să poată trece peste aceste impedimente, și indiferent dacă cuvântul respectiv este o altă formă a cuvântului de bază sau chiar și un sinonim, să funcționeze exact ca și cum cuvintele întâlnite sunt aceleași.

Pentru operația de filtrare clasificatorii sunt folosiți pentru a determina care mesaje vor fi relevante față de mesajul introdus de noi și care vor fi irelevante. Clasificatorii vor fi construiți în special în funcție de ponderile cuvintelor care apar în mesaje. Astfel fiecare mesaj va avea o pondere, pondere care nu va însemna nimic altceva decât ponderea de relevanță a textului. De notificat ar fi și faptul că fiecare din obiectivele descrise mai sus vor fi folosite pentru filtrarea mesajelor.

Sistemul de filtrare a fost construit ca fiind unul dinamic. Fiecare modificare a unui mesaj din baza de date aduce la rândul său o modificare ulterioară a mesajului din filtru însă numai la următoarea rulare a filtrului. O altă caracteristică care face ca sistemul nostru să fie unul dinamic se referă la faptul că, ponderea de calcul a mesajelor relevante să se calculeze dinamic, în funcție și de mărimea mesajului. Teoria din spatele acestei idei constă în faptul că, conform studiilor, cu cât mesajul este mai lung, oamenii tind să folosească mai puține cuvinte de specialitate și cât mai multe uzuale, care să îl facă pe receptor să se simtă familiar cu discuția. În figura 2.1 ne este prezentată schema generală a proiectului

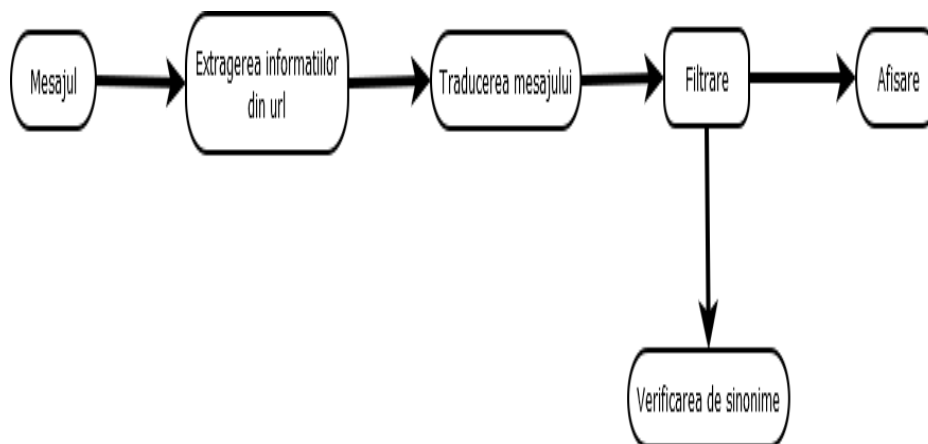


Figura 2.1 Schema generală a sistemului

2.2.1. Cerințe funcționale

Cerințele funcționale reprezintă acele funcționalități concrete pe care sistemul trebuie să le ofere astfel încât să răspundă cerințelor de business. Se bazează pe analiza și modelarea proceselor de lucru.

Principala cerință funcțională pe care proiectul trebuie să o furnizeze este aceea că trebuie să ofere rezultate optime la un număr mare de mesaje. Algoritmul de filtrare trebuie să fie independent de numărul de mesaje și să afișeze atâtea rezultate câte va fi nevoie.

Filtrarea trebuie să fie de asemenea independentă de mărimea mesajului. Deși micro-bloggingul este limitat la 500 de cuvinte, datorită folosirii url-urilor și datorită

preluării informației utile din interiorul url-ului, mesajul nostru poate ajunge până la 10 000 de cuvinte sau chiar depăși acest prag. Algoritmul este unul dinamic, de aceea numărul de cuvinte nu va putea influența sortarea.

Utilizatorul va putea introduce spre filtrare un mesaj care va trebui să aibă aceeași lungime ca și mesajele din micro-blogging. În cazul în care mesajul va avea o lungime mai mare decât cea din micro-blogging utilizatorul va fi avertizat cu privire la faptul că a depășit dimensiunea maximă a mesajului. De asemenea în cazul în care utilizatorul introduce un mesaj gol, acesta va fi avertizat asupra acestui fapt.

O altă caracteristică funcțională se poate referi la faptul că utilizatorul va putea vizualiza mesajele în funcție de relevanța lor. Astfel mesajul cel mai semnificativ, respectiv cel care are cea mai mare legătură cu mesajul dat de către noi, va fi afișat primul. Mesajul cel mai puțin semnificativ însă cu o semnificație mai mare decât valoarea de prag va fi afișat ultimul.

2.2.2. Cerințe non-funcționale

Cerințele non-funcționale reprezintă o constrângere de cost, calitate, construcție, livrare sau operarea unui produs, componentă a unui produs, sau funcția a unui produs. În domeniul software cerințele non-funcționale se referă la acele cerințe care nu implică realizarea unei funcționalități, dar care sunt necesare pentru ca funcționalitățile să poată fi utilizate. Includ cerințe de dimensionare, scalabilitate, securitate, performanță, ușurință de utilizare.

Dimensiunea unei aplicații reprezintă numărul de kilobytes pe care aplicația o ocupă pe server. **Dimensionarea** reprezintă caracteristica care ne arată cât de compact este programul nostru. Direcția de evoluție a programării este aceea în care se insistă pe aplicații distributive. Serviciile web, aplicațiile locale, api-urile reprezintă modalitățile prin care se vor îmbunătăți atât distributivitatea cât și dimensionarea. Cerința pentru proiectul actual este aceea de a folosi pe cât posibil aceste tehnologii.

Conform glosarului de termeni științifici, **scalabilitatea** este însușirea unui sistem, echipament sau dispozitiv software de a permite schimbări importante ale dimensiunilor sau capacităților sale, cu costuri acceptabile, fără dificultăți și cu păstrarea caracteristicilor și a performanțelor inițiale. O bună diviziune a claselor și o modularizare cât mai bună a proiectului îi poate conferi acestuia o scalabilitate ridicată.

Securitatea reprezintă caracteristica prin care un client are acces doar la datele la care îi permite developer-ul și nu poate accesa alte date. În cazul aplicațiilor web, securitatea este cea mai importantă caracteristică. Aplicația mea trebuie să respecte atât securitatea datelor din micro-blogging cât și securitatea codului din filtru de mesaje.

În ingineria software, **performanța** se caracterizează ca fiind capacitatea sistemului de a răspunde repede și stabil cerințelor utilizatorului. Având în vedere faptul că proiectul are la bază un algoritm, timpul de răspuns și stabilitatea algoritmului este foarte importantă. Un timp de răspuns lent va face algoritmul inutil și ineficient pe când un timp de răspuns rapid însă greșit va anula întreaga funcționalitate a proiectului.

Ușurința de utilizare se referă la capacitatea unei aplicații de a fi înțeleasă și folosită cu ușurință de către un utilizator. Având în vedere faptul că funcționalitățile aplicației sunt puține, deoarece ideea de bază a proiectului constă în algoritmul folosit, dar și faptul că interfața grafică este una prietenoasă, ușurința de utilizare a clientului este una foarte mare.

2.2.3. Cerințe de implementare

Cerințele de implementare se referă la acele tehnologii și limbaje de programare care vor trebui folosite pentru dezvoltarea aplicației. Având în vedere faptul că aplicația mea este un site web, aceasta trebuie să conțină atât parte de front-end cât și parte de back-end.

Când discutăm despre front-end, vorbim despre acea parte a site-ului sau a aplicației web, pe care o putem vedea și cu care interacționează vizitatorii. Pentru afisarea paginii vom folosi pagini jsp, pentru stilizare se vor folosi css-uri iar pentru operațiile asincrone se va folosi javascript.

Rolul back-endului este cel de management de conținut. Ca și limbaj de bază al aplicație în care va fi scrisă întreaga logică a aplicației se va folosi limbajul Java. Pentru partea de bază de date se va folosi sistemul de gestiune de baze de date MySql.

Ca și platforme pentru dezvoltarea aplicației se va folosi NetBeans 7.4 iar pentru accesul la baza de date s-a instalat programul XAMPP împreună cu extensia sa de baze de date, PhpMyAdmin. Pentru rularea micro-blogging-ului se va folosi platforma Eclipse(programul a fost creat pe această platformă).

Alte cerințe de implementare suplimentare ar fi: documentarea și comentarea codului (pentru a fi mai ușor de înțeles și modificat de diferite persoane), testarea manuală și automată (pentru a ne asigura de faptul că programul funcționează bine și cu parametrii de performanță acceptați), realizarea unei aplicații modulare care să poată fi ușor integrată cu alte sisteme.

2.3. Conținutul lucrării

În cele ce urmează se vor dezvolta și se va explica în mod mai clar toate cele 8 secțiuni evidențiate, lucrarea fiind împărțită în 8 capitole:

- În capitolul I, Introducere – se va specifica contextul în care a fost creat proiectul nostru dar și diferitele provocări pe care le vom întâmpina în dezvoltarea programului.
- În capitolul II, Obiectivele proiectului s-au descris Specificațiile și obiectivele aplicației împreună cu cerințele funcționale, non-funcționale și de implementare ale aplicației.
- În capitolul III, Studiu bibliografic se va face un studiu bibliografic cu referire la aspectele teoretice folosite în lucrarea de față.
- În capitolul IV, Analiză și fundamentare teoretică, se descriu algoritmi și noțiunile teoretice utilizate în aplicație. Scopul acestui capitol este de a ne explica principiile funcționale ale aplicației implementate. La sfârșitul acestui capitol se prezintă soluția propusă.
- În capitolul V, Proiectare de detaliu și implementare, se descriu pe larg modulele sistemului de construire a ierarhiilor, funcționarea lor și modul în care au fost implementați diverși algoritmi prezentați în capitolul anterior.
- În capitolul VI, Testare și validare se prezintă metodele de testare folosite pentru validarea aplicației.
- În capitolul VII, Manual de instalare și utilizare, sunt detaliate resursele software și hardware necesare pentru instalarea și rularea aplicației, precum și

o descriere pas cu pas a procesului de instalare. Se descrie și modul de utilizare al aplicației pas cu pas, din punctul de vedere al utilizatorului.

- În capitolul VIII, Concluzii se prezintă un rezumat al contribuțiilor personale, o analiză critică a rezultatelor obținute și o descriere a posibilelor dezvoltări și îmbunătățiri ulterioare.

Capitolul 3. Studiu Bibliografic

Lucrarea noastră conține două direcții importante. Prima direcție se referă la clasificarea și filtrarea textelor. În această parte se va discuta în principal despre filtrarea textelor. În cea de-a doua parte se va discuta despre conceptul de micro-blogging și se va analiza modul în care micro-bloggingul realizează filtrarea de mesaje, despre modul cum micro-bloggingul local a fost creat și despre cum se aseamănă cu alte micro-blogginguri.

3.1. Clasificarea și filtrarea textelor folosind algoritmi clasici

În literatura de specialitate se găsesc două tipuri de sortări de mesaje: Sortarea după conținutul mesajului și filtrarea după componentele sociale.

3.1.1 SIFT

Sortarea după conținutul mesajului este cea mai bună modalitate de sortare. În acest tip de sortare fiecare utilizator se presupune a lucra independent unul față de celălalt. Ca și rezultat sortarea după conținutul mesajului poate exploata doar documentele.

Programatorul Yan a implementat o filtrare bazată pe conținut pentru prima dată pentru un program de știri, într-un sistem intitulat SIFT[3]. Profilele pentru SIFT au fost construite manual prin specificarea cuvintelor care sunt de preferat să fie adăugate sau ignorate și trebuia să fie updatat manual dacă utilizatorul dorea să le schimbe. Pentru fiecare format, 20 de articole erau disponibile în fiecare zi. Articolele puteau fi alese selectiv folosind World Wide Web.

SIFT oferea două facilități pentru a ajuta utilizatorul cu crearea profilelor. Yan a dezvoltat SIFT pentru a studia eficient algoritmi pentru filtrarea informațiilor. În SIFT colecții mari de informații erau comparate cu fiecare articol care sosea în Internet News printr-un server central. Eficiența era obținută prin gruparea profilelor în grupuri în așa fel încât filtrarea se poate face pe grupuri. SIFT nu făcea deosebire între locurile în care apăreau cuvintele așa că, cuvintele care apăreau în titlu, în corpul mesajului sau în semnătură aveau de obicei aceeași pondere.

3.1.2 Fast Data Finder

Câteva sisteme de filtrarea de texte comerciale se bazează pe construcția manuală a profilelor. Sistemul „Fast Data Finder”[3], un produs al Parcel Inc, folosește mii de unități de procesare care sunt optimizate pentru operații cum ar fi ponderea termenilor, constrângerile de proximitate etc. Fiecare unitate de procesare este programată pentru un singur task iar grupuri de unități de procesare sunt create pentru fiecare profil în parte. Căutări simultane pe multiple profile sunt suportate așadar un profil multilingual poate fi implementat ca și un set de profile unilinguale, unul pentru fiecare limbă. Tool-uri automate au fost create pentru a ajuta utilizatorul să-și creeze profiluri. Spre deosebire însă de SIFT, algoritmul nu este dinamic, așa că profilele de Fast Data Finder nu se adaptează automat.

3.1.3 Info Scope

Tehniciile de filtrare a textelor în mod adaptiv caută să facă utilizatorul să caracterizeze un profil din documentele anterioare și să folosească acel profil pentru a stabili o bază de reper pentru noile documente sosite. Stevens a dezvoltat un sistem intitulat „Info Scope” [3] care se baza pe această abordare. Info Scope implementa filtrare adaptivă de text, sugera reguli de funcție pentru observarea comportamentului utilizatorului și le oferea pentru acceptare sau uneori modificare. Aceste sugestii aveau la bază acțiuni simple, observabile cum ar fi timpul petrecut de utilizator pentru a citi stirile sau dacă un document citit a fost salvat pentru referințe ulterioare. Prin evitarea cerințelor pentru feedback-ul utilizatorilor despre articole individuale Info-Scope a fost creat pentru a minimiza încărcarea de informații pe server.

Dacă SIFT trata știrile ca fiind o colecție monolitică de articole, Info-Scope face filtrări clare între grupuri, subiecte și chiar și persoane individuale. Stevens a construit în acest scop nivele de abstractizare astfel încât acestea să fie importante pentru utilizator. InfoScope a implementat această idee pentru a dezintegra grupuri de știri și a forma grupuri virtuale individuale. Pentru a defini știri pentru grupuri virtuale cu posibilitate de a suprascrise sursele, utilizatorii pot reorganiza informația în acordanță cu modelul personal cognitiv al părților care participă la comunicare pe care ei doresc să le observe.

InfoScope nu a fost însă fără limitări. Sistemul experimental pe care Stevens l-a dezvoltat a fost capabil să proceseze informațiile doar din header-ul fiecărui articol. Astfel potențialul algoritmului de a explora texte mari era limitat așa că a fost nevoie de folosirea unei metode care să aleagă doar cuvintele exacte din headere.

Datorită costului lor redus, disponibilitatea unui volum mare de mesaje, și ușurința de a recunoaște informații noi, Internet News și poșta electronică au fost domenii populare pentru cercetarea informațiilor de filtrare. Din păcate, aceste domenii sunt slab potrivite pentru experimente formale, deoarece rezultatele reproductibile sunt greu de obținut. Din acest motiv, foarte puține se știe despre eficacitatea SIFT-ului sau InfoScope-ului. Stevens a raportat faptul că opt din zece cititori experimentați de Internet News au preferat InfoScope la software-ul lor principal într-un studiu inițial, și că toți cei cinci utilizatori din a doua evaluare au raportat că mai puține articole neinteresante au fost prezentate în prima săptămână și mai multe articole interesante au fost citite în a 10-a săptămână decât în prima. Deoarece SIFT a fost dezvoltat pentru a studia eficiența, mai degrabă decât probleme de eficacitate, cu atât mai puțin sunt disponibile informații cu privire la eficacitatea acestuia.

Învățând mai multe despre eficiența unei tehnici de filtrare a textului prevede ca tehnica să fie evaluate în condiții experimentale controlate. Și pentru că performanța tehnicilor de filtrare de text poate varia semnificativ atunci când are nevoie de informații diferite și sunt folosite colecții de documente, compararea rezultatelor între sisteme este facilitată atunci când acești factori sunt considerați constante. Evaluarea de rutare TREC a oferit o întâlnire fără precedent pentru acest tip de evaluare a performanțelor. Realizat anual, începând din 1992, cea mai recentă evaluare de rutare (TREC-5) a atras participarea de la 14 grupuri de cercetare.

3.1.4 TREC NIST

NIST [3] oferă fiecărui participant cincizeci de subiecte și un set mare (de obicei sute) de documente de formare și evaluări relevante pentru fiecare subiect. Participanții își dezvoltă sistemele lor de filtrare a textului, folosind aceste date ca și cum ar reprezenta un feedback explicit cu privire la utilitatea fiecărui document de formare dat de către un utilizator, iar apoi trebuie să își înregistreze profilurile lor cu NIST înainte de a primi documentele de evaluare. Profilele sunt apoi utilizate de către sistemele de filtrare a textului care le-au generat pentru a pondera un set nevăzut în prealabil de documente de evaluare, iar primele câteva mii de documente se vor depune la NIST pentru evaluare.

Pentru a obține rezultate reproductibile, este necesar să se facă unele ipoteze foarte puternice cu privire la natura sarcinii de filtrare de informații. În TREC se presupune că judecățile umane cu privire la satisfacerea unei informații de către un document sunt binar evaluate (de exemplu, un document este relevant pentru o informare sau nu este) și constant, (nu contează cine face această hotărâre sau atunci când se face). Relevanța, conceptul fundamental pe care se bazează această metodologie, de fapt, nu îndeplinește simultan aceste două ipoteze. Judecata de relevanță umană prezintă variabilitate semnificativă între mai mulți evaluatori, dar și variabilitate pentru același evaluator de-a lungul timpului. Mai mult decât atât, evaluatorii găsesc uneori că este dificil să facă o judecată de relevanță binară pe o combinație specifică a unui document și o nevoie de informare. Cu toate acestea, măsurătorile de performanță bazate pe un set comun de hotărâri oferă o bază de principiu pentru a compara performanța relativă a diferitelor tehnici de filtrare a textului pe bază de conținut.

Evaluarea Trec se bazează pe măsuri de eficacitate, care sunt frecvent utilizate pentru sistemele de filtrare a textului. Eficiența sistemelor de filtrare a textului este de obicei caracterizată prin trei statistici: "Precizie", "Reapelare" și "Eșec". Precizia este fracțiunea de detectat (și astfel relevante) din documentele care sunt de fapt relevante pentru utilizator, în timp ce reapelarea este fracțiunea din setul actual de documente relevante care sunt clasificate în mod corect ca fiind relevante de către sistemul de filtrare a textului. Atunci când sunt utilizați împreună, precizie și reapelarea definesc eficacitatea de detecție. Eșecul (fracțiunea de documente non-relevante care sunt clasificate de către sistemul de drept potențial relevant) este folosită pentru a măsura eficiența respingerii.

În TREC, aproape toate sistemele de filtrare a textului produc rezultate în funcție de pondere. În consecință, precizia și eșecul de la mai multe valori de rechemare sunt raportate, iar "precizia medie" (aria de sub curba de precizie-eșec) este raportată pentru utilizare atunci când este nevoie de o singură măsură de eficacitate. Precizia medie se calculează prin alegerea de seturi succesive de documente cu pondere mare și care au un grad mare de reapelare. Precizie medie se calculează prin alegerea de seturi succesiv mai mari de documente din partea de sus a listei pe locul care duce la o reapelare din ce în ce mai mare. Precizia este apoi calculată pentru fiecare set, o tehnică de interpolare este utilizată pentru a estima precizia între punctele de rechemare observate, iar aria de sub curba interpolată este raportată ca precizia medie pentru o informație individuală. Procesul se repetă pentru mai multe informații, iar media valorilor obținute în acest mod este raportat ca precizia medie pentru sistemul testat. Deoarece atât precizia cât și reapelarea variază între zero și unu, valorile mai mari de precizie medie sunt considerate a fi mai bune, valoarea ideală fiind în mod logic unu.

3.2 Clasificatori de tip machine learning

În ultimii ani, managementul documentelor în funcție de conținut a luat o amploare deosebită datorită creșterii disponibilității documentelor în formă digitală dar și nevoii de a le accesa în moduri flexibile. Categorisirea de text (sau clasificarea de text), activitatea de a asocia un document cu o categorie tematică provenită dintr-un anumit set este unul dintre aceste task-uri. Un algoritm de tip machine learning presupune faptul că, în momentul în care are loc o filtrare de mesaje, sistemul „învață” anumite stereotipuri astfel încât la următoarea filtrare algoritmul să fie mai bun.

3.2.1 Clasificatori de tip arbori de decizie

Un clasificator de tip arbori de decizie [4] este de fapt o reprezentare simbolică a unui copac în care nodurile interne sunt termenii, ramurile care pleacă de la nodurile interne sunt ponderile pe care acei termeni îi au în document iar frunzele sunt categoriile. Un astfel de clasificator categorizează un document de test prin testare recursivă a ponderiilor din nodurile interne și adăugarea lor într-un vector până în momentul în care o ajungem la o frunză.

Există un număr foarte mare de pachete standard pentru folosirea clasificatorilor de tip decision-tree. Printre cele mai populare putem aminti ID3 (folosit în Fuhr 1991), C4.5 (folosit în Cohen and Hirsh 1998) și C5 (folosit în Li și Jain 1998).

O metodă posibilă a algoritmului de învățare pentru un arbore de decizie și o categorie constă în folosirea strategiei „divide and conquer”. Strategia verifică dacă toate exemplele au aceeași categorie; dacă nu se va selecta un termen și se va diviza arborele nostru în funcție de termenul ales. Procesul se va repeta recursiv și pe sub-arbori până când fiecare frunză va fi asignată unei anume categorii. Pasul cheie al algoritmului este alegerea unui termen potrivit după care să se facă subdivizarea, această alegere făcându-se pe baza unui cumul de informații sau unui criteriu entropic. Majoritatea algoritmilor care au la bază arbori de decizie includ metode pentru creșterea mărimii arborilor (și reducerea complexității lor) dar și pentru tăierea ramurilor inutile (pruning algorithm).

Arborii de decizie au fost folosiți de-a lungul timpului fie ca algoritmi de clasificare primară, fie ca și clasificare secundară sau ca și membri ai unei clasificări.

3.2.2 Clasificatori de tip reguli de decizie

Un clasificator de tip reguli de decizie [4] reprezintă un clasificator creat prin metode din forma normală disjunctivă. Forma normală disjunctivă a unei mesaje reprezintă o modalitate de reprezentare a mesajului care folosește doar simboluri denumiți literari. Literarii din premise (mesajele inițiale) denotă fie prezența fie absența cuvântului de bază în document pe când clauzele de decizie denotă dacă documentul poate aparține categoriei sau nu. Reguliile DNF (forma normală disjunctă) sunt similare cu cele din arborii de decizie prin faptul că pot codifica orice funcție booleană. Pe de altă parte un avantaj al folosirii DNF l-ar putea reprezenta faptul că acesta din urmă tinde să genereze o clasificare mult mai compactă.

Metodele bazate pe reguli tind să aleagă din toate reguliile pe cele mai bune conform anumitor criterii de minimalitate. În timp ce arborii de decizie au fost construiți printr-o strategie de tip dop-down sau divide-and-conquer reguliile DNF au fost construite printr-o strategie bottom-up. Inițial fiecare exemplu este văzut ca și o serie de clauze. Acest set de clauze este deja asignat categoriei, însă are un scor redus. Utilizatorul va aplica o serie de procese de generalizare în care regula va fi simplificată printr-o serie de modificări (îndepărtarea premiselor din clauze sau unirea clauzelor), care vor maximiza compactivitatea și nu vor micșora acoperirea clasificatorului. La sfârșitul clasificării, un algoritm de tăiere a clauzelor neimportante este aplicat, în cele din urmă rămânând doar clauzele importante. În acest moment, în funcție de ponderea obținută, putem să ne dăm seama dacă algoritmul aparține metodei respective sau nu.

Reguliile DNF variază larg în ceea ce privește formele lor de apariție (metode, euristici și criterii). Printre reguliile DNF care au fost aplicate la filtrarea de mesaje amintim CHARADE (Moulinier and Ganascia 1996), DL-ESC (Li and Yamanishi 199), RIPPER (Cohen 1995), Scar (Moulinier 1996) și SWAP (Apte 1994).

Metodele de mai sus folosesc logica propozițională însă au fost folosiți algoritmi în care s-a folosit și tehnica First-Order-Logic. Cohen a testat comparativ DNF și FOL și a ajuns la concluzia că First-Order-Logic este mai limitat în eficacitate decât DNS.

3.1.1. Rețele neuronale

Un sistem bazat pe rețele neuronale [4] este un sistem care este compus în principal din unități, fiecare unitate de intrare reprezentând un termen, fiecare unitate de ieșire reprezentând categoria de interes iar legăturile reprezintă ponderiile. Pentru a clasifica un document, se vor afișa pe intrări termenii, se va rula algoritmul, iar pe ieșiri se vor putea vizualiza categoriile. O metodă tipică în rețelele neuronale este propagarea înspre înapoi. Astfel, termenii vor fi introduși pe intrare și dacă o eroare de clasificare are loc, rețelele care conțin această eroare se vor reinițializa. Prin această modalitate, vom avea costuri de recuperare mici (nu se vor înlocui decât rețelele afectate) și securitate sporită (datorită apariției metodei de recuperare din eroare).

Cel mai simplu tip de clasificare bazat pe rețele neuronale este „preceptron”. Acesta este un algoritm de clasificare liniar. Alte tipuri de clasificatori neuronali implementează o formă de logistică regresivă și sunt doar în stadiul de propunere.

Un algoritm non-linear bazat pe rețele neuronale este în schimb o rețea cu unul sau mai multe straturi de unități pe care algoritmul de clasificare îl reprezintă ca și interacțiuni de tip înalt pe care rețeaua trebuie să fie capabilă să-l învețe. Totuși comparând algoritmul non-linear cu cel linear s-a ajuns la concluzia că foarte puține îmbunătățiri au fost aduse.

3.1.2. Asamblorul de clasificatori

Asamblorul de clasificatori [4] este bazat pe ideea că, în momentul în care avem un task care necesită expertiză de specialitate pentru a putea fi rezolvat, mai mulți experți ar putea fi mai utili în rezolvarea lui decât unul singur. Traducând această idee în filtrarea de mesaje rezultă faptul că, aplicând o serie de clasificatori asupra aceluiași text,

posibilitatea ca textul sa fie clasificat corect este mai mare. Astfel asamblorul de clasificatori este caracterizat de o serie de clasificatori și un algoritm de decizie.

Problema principală a acestui algoritm este aceea că, pentru a asigura eficacitate bună, clasificatorii ar trebui să fie independenți unul de celălalt. Aceștia ar putea să difere atât ca și conținut cât și ca modalitate. Astfel, folosind mai mulți algoritmi de filtrare timpul de execuție va crește.

3.2 Micro-blogging

Termenul de blog este o prescurtare a termenului “webblog” și se referă la un jurnal creat online unde articolele sunt afișate în ordine invers cronologică, ultimul scris fiind primul care apare. În prezent în lume există peste 3 milioane de blog-uri având diferite teme.

3.2.1. Microblogging-ul local

Microblogging-ul local este unul de tip open-source. Acesta a fost simplificat deoarece tema mea își propune să pună un mai mare accent pe filtrarea de mesaje și nu pe design-ul și implementarea microblogging-ului. Totuși acesta se încadrează perfect în categoria microblogging-ului și de aceea am să-l prezint ca și un microblog în comparație cu celelalte microblogging-uri. Principalul obiectiv a microblogging-ului personal este acela de a permite persoanelor atât postarea cât și citirea anumitor mesaje.

Prima pagină a micro-blogging-ului este pagina de logare. Această pagină este una asemanătoare celorlalte pagini de logare obișnuită de pe alte aplicații. Pe lângă logare, utilizatorul va mai putea să-și creeze cont pe aplicație. Pentru a putea posta mesaje sau pentru a putea citii mesajele celorlalți utilizatori, utilizatorul nostru trebuie în primul rând să se logheze. Vizualizarea mesajelor nu este permisă fără logare.

Pentru înregistrarea contului, utilizatorul are nevoie doar de aceleași date ca și pentru logare adică username și parolă. Aplicația a fost importată doar pentru postarea de mesaje și de aceea complexitatea sa nu este foarte importantă.

Odată logați în sistemul nostru avem 4 posibilități: Crearea mesajului, găsirea unui utilizator, vizualizarea postărilor unui utilizator (follow) oricând acesta le va modifica și vizualizarea pe moment a postărilor. Pentru crearea unei postări avem la dispoziție 500 de caractere. Restricția a fost adăugată deoarece aplicația face parte din categoria micro-bloggingului iar micro-bloggingul nu permite decât mesaje scurte. De asemenea în această pagină se pot vizualiza și mesajele tuturor utilizatorilor care sunt „urmăriți” de către noi.

În pagina care se ocupă cu găsirea unui utilizator avem un câmp care va trebui completat cu numele utilizatorului. Dacă utilizatorul va fi găsit în sistemul de micro-blogging, utilizatorul nou găsit va apărea într-un tabel iar sistemul îi va da posibilitatea utilizatorului curent să urmărească postările introduse de către acel utilizator.

În pagina în care putem vedea postările utilizatorului, avem, la fel ca și în pagina anterioară un input în care putem introduce un nume. După ce utilizatorul a fost găsit ne vor apărea în pagină postările sale. Spre deosebire de prima pagină, în care ne apar toate postările utilizatorilor urmăriți, în această pagină ne apar doar postările utilizatorului

respectiv. Utilizatorul va putea fi schimbat oricând în câmpul de căutare în tabel urmând să ne apară postările celui alt utilizator.

Din punct de vedere tehnic microblogging-ul este o aplicație aparent simplă însă cu o complexitate și un număr al claselor destul de mare, clasele fiind special create pentru a îndeplini funcția de securitate. A fost creat în limbajul Java și are la bază tehnologia Spring. Aplicația se mapează la 4 tabele în baza de date: `authorities`, `follower`, `user` și `post`.

Tabela `post` este cea mai importantă pentru sistemul nostru de filtrare deoarece din această tabelă se vor extrage mesajele dorite. Aceasta conține câmpurile de `username`, `password` și `enabled`. De notificat ar fi aici și faptul că parola este salvată în baza de date sub formă criptată. În continuare se vor prezenta câteva servicii de micro-blogging și se vor compara cu blogul local.

3.2.2. *Microblogging-uri uzuale*

Termenul de microblogging nu se referă la altceva decât un blog în care conținutul este mai mic atât în ceea ce privește mesajele cât și în ceea ce privește dimensiunea fișierelor trimise. Pe scurt, în interiorul unui microblog, utilizatorii trimit unul altuia mesaje scurte, poze individuale sau link-uri către alte site-uri web. În ultimii ani, micro-bloggingul a avut o creștere considerabilă comparativ cu părintele sau mai mare blog-ul, în principal datorită simplității sale. Printre microblogging-urile mai importante putem aminti: Twitter, Binify, Status Net, Jaiku, Qaiku.

3.2.2.1. *Twitter*

Twitter[10] este cel mai important serviciu de micro-blogging. Acesta permite utilizatorilor să înregistreze mesaje care au dimensiunea de maxim 140 de caractere. Microblogging-ul personal pe de altă parte permite utilizarea a 500 de caractere. Twitter este un microblogging multilingval, adică opțiunile și conținutul acestuia poate apărea în mai multe limbi. Postările însă nu vor fi modificate ele putând apărea în limba care au fost create și sub modalitatea în care au fost create. Micro-bloggingul personal folosește doar engleza ca și limbă de bază.

O altă caracteristică a Twitter-ului care poate fi utilizată în compararea acestuia cu micro-bloggingul personal se referă la conceptul de “tagging”. Un tag este un cuvânt cheie sau un termen asignat unei informații care ne poate duce către pagina de unde aceea informație provine. În Twitter putem da tag unei postări de Twitter sau oricărei alte pagini Web. Microblogging-ul personal, fiind o creație mult mai simplă nu permite tagging-ul sub nici o formă.

Micro-bloggingurile au de obicei două tipuri de conversații. Un tip de conversație este acela public, în care oricine poate vedea mesajele care au fost schimbate între oameni. Un alt tip de conversație poate fi cel privat în care doar persoanele care participă la conversație pot vedea mesajele. În ceea ce privește Twitter-ul acesta permite atât conversații private cât și conversații publice. Microblogging-ul personal permite doar acele tipuri de conversații care sunt publice, mesajele postate de cineva putând fi văzute de oricine urmărește aceea persoană.

Aplicațiile web pot avea servicii de tip open source. Prin intermediul serviciilor de tip open-source putem accesa date din baza de date a acelor aplicații. Microblogging-urile fiind aplicații web au posibilitatea de a acorda anumite servicii de tip open source. Atât Twitter-ul cât și microblogging-ul personal oferă utilizatorilor acces la toate mesajele postate.

Micro-bloggingurile pot avea o serie de funcționalități suplimentare, însă principala lor funcționalitate ar fi aceea de a posta și de a citi mesaje. Postarea și citirea mesajelor presupune operații diferite și pot fi făcute prin mai multe modalități. În cazul twitter-ului, putem posta mesaje prin web, sms sau api. Postarea mesajelor prin web presupune accesarea paginii web unde vom scrie și posta mesajul. Postarea prin sms a fost creată pentru a spori flexibilitatea microblogging-ului și constă în scrierea unui mesaj pe telefon și respectarea anumitor indicații astfel încât mesajul dorit să poată să ajungă în locul dorit de noi. Postarea prin API este foarte asemănătoare cu postarea prin web, însă în loc să deschidem o pagină web vom accesa o aplicație desktop special creată de Twitter care ne va permite să scriem mesajele pe pagina dorită. Citirea mesajelor poate fi făcută prin web, sms, api și feed. Un web feed sau un news feed este un format de date utilizat pentru a oferi utilizatorilor conținutul actualizat frecvent. În cazul microblogging-ului personal citirea și postarea mesajelor se poate face doar prin web.

3.2.2.2. *StatusNet*

StatusNet[11] este cel mai complex microblogging existent pe piață. Deși este foarte asemănător cu Twitter și are aproximativ aceleași funcționalități, StatusNet asigură utilizarea unor dezvoltări noi.

O primă caracteristică care ar putea fi folosită în comparația cu microblogging-ul personal ar fi faptul că StatusNet asociază utilizatorii în diferite grupuri. Grupurile pot fi formate după diferite criterii, utilizatorii putând aparține mai multor grupuri iar un grup poate cuprinde mai mulți utilizatori. Microblogging-ul personal, spre deosebire de StatusNet nu permite aceasta dezvoltare, fiecare utilizator fiind considerat ca și persoană diferită.

RSS este una din dezvoltările noi adoptate de către StatusNet. RSS provine de la Rich Site Summary și este un protocol prin care se pot trimite informații complexe (imagini, audio) și care are ca și scop principal transformarea informației într-una publică, informația putând fi accesată din oricare altă pagină web. Microblogging-ul personal nu folosește protocolul rss însă oferă accesibilitate la informație.

FOAF sau Friend of Friend este o ontologie folosită în special în căutarea utilizatorilor. Astfel algoritmul începe de la utilizatorul actual și caută recursiv prin toți “prietenii” acestuia până când utilizatorul dorit va fi găsit. Microblogging-ul personal nu permite acest tip de căutare însă ar putea fi implementată ca și o dezvoltare ulterioară.

Bookmarklet-ul este un bookmark păstrat în browser care conține comenzi Javascript pentru a extinde funcționalitatea browser-ului. StatusNet folosește acest tip de dezvoltare nouă în special pentru a crește distributivitatea aplicației. Microblogging-ul local, are secvențe de Javascript integrate direct în cod și nu folosește bookmarklet-ul.

În ceea ce privește postarea și scrierea de mesaje, StatusNet este superior atât microblogging-ului local cât și Twitter-ului. Postarea mesajelor se poate face prin web, xmpp, sms, mail sau api. Postarea prin xmpp este o dezvoltare nouă a StatusNet. Xmpp

este un protocol special folosit in comunicare care salvează mesajele intr-un format asemănător XML-urilor. Microblogging-ul local folosește în comunicare tehnica clasică prin care mesajele se salvează ca și șiruri. Postarea mesajelor prin E-mail este posibilă printr-o mapare dintre email-ul destinație și pagina pe care se va face postarea mesajelor. Citirea se face la fel ca și scrierea prin toate modalitățile posibile. În tabelul de mai jos se vor prezenta modalitățile prin care se pot posta și citi mesaje in cele mai populare servicii de micro-blogging. Primele 4 coloane se referă la postarea mesajelor iar ultimele 4 coloane se referă la citirea mesajelor.

Micro-blogging	WEB	XMPP	EMAIL	SMS	WEB	XMPP	EMAIL	SMS
StatusNet	DA	DA	DA	DA	DA	DA	DA	DA
Jaiku	DA	DA	NU	NU	DA	DA	NU	NU
Qaiku	DA	DA	DA	NU	DA	DA	NU	NU
Twitter	DA	NU	NU	DA	DA	NU	NU	DA

Tabela 3.1 Tipurile de comunicații acceptate de micro-blogginguri[12]

În concluzie, studierea caracteristicilor fiecărui sistem de micro-blogging în parte a ajutat la înțelegerea scopului micro-bloggingului. Astfel, în acest moment se știe motivul pentru care micro-bloggingul a fost creat, și modalitățile prin care mesajele circulă prin interiorul serviciilor de micro-blogging.

Capitolul 4. Analiza și fundamentarea teoretică

Filtrarea postărilor irelevante în legătură cu un subiect dat poate fi realizată prin clasificarea fiecărei postări ca fiind relevantă sau irelevantă și afișare pe ecran. Pentru aceasta s-a propus o metodă pentru a identifica o metodă de clasificare care să filtreze mesajele.

Clasificarea are la bază o serie de module de transformare și decizie. Pentru fiecare postare vom returna valorile care au relevanța cu mesajul scris de noi. Extragerea cuvintelor importante din interiorul mesajului nostru este foarte important pentru algoritm. Din fericire tehnologiile actuale fac acest lucru mult mai ușor decât în trecut.

Vom folosi o serie de module pre-filtrare care au diferite roluri, rolul principal fiind acela că se va obține un mesaj simplificat mult mai ușor de manipulat. Urmează apoi etapa de filtrare unde se vor filtra mesajele bazate pe conținut, urmând ca mai apoi, în etapa finală să se afișeze acele informații care sunt relevante cu mesajul dat de noi.

4.1. Clasificare bazată pe conținutul mesajelor

Scopul clasificării este acela de a filtra mesajele care au texte irelevante față de un mesaj dat de noi, elementele textuale sunt principalele caracteristici din sistemul nostru. Caracteristicile textuale ar trebui să reprezinte cât de relevant este subiectul textului pentru mesajul nostru. Fiecare mesaj are un subiect diferit (sau subiecte), de aceea avem nevoie de o metodă pentru a extrage în mod automat temele-cheie din mesajul nostru și măsura similaritatea textului cu acesta.

4.1.1. Module de prelucrare pre-filtrare

Pentru a putea filtra textul avem în primul rând nevoie ca acesta să fie modificat. Textele trebuie să ajungă la algoritmul de filtrare într-o formă minimă astfel încât, procesări ulterioare să nu mai trebuiască făcute. În funcție de momentele de apariție vor fi prezentate modulele de preprocesare.

4.1.1.1. Modulul de extragere a conținutului HTML

Micro-bloggingul este limitat la 500 de caractere. Acest lucru poate fi considerat ca un impediment deoarece explicațiile și detaliile mesajelor sunt restrânse. Astfel, mulți utilizatori apelează la opțiunea de a introduce un url care să conțină mai multe informații cu legătură la subiect. Modulul de extragere a conținutului HTML își propune să extragă informația utilă din interiorul url-ului și să o aprobeze textului.

Pentru a putea extrage datele utile dintr-un html, trebuie mai întâi să verificăm dacă cuvântul respectiv este sau nu este un url. Clienții nu au în micro-blogging-ul local nici un loc special în care să poată insera url-uri de aceea sunt nevoiți să le introducă în interiorul mesajului. Pentru a verifica faptul că avem un url în mesaj vom parcurge mesajul și vom aplica un algoritm de verificare de url pe fiecare cuvânt.

Modalitatea de verificare dacă un cuvânt este un url sau nu o reprezintă încercarea de conexiune la acel url. Dacă cuvântul nostru este url , rezultatul browser-ului va avea codul 200. Dacă cuvântul nu este url, rezultatul browser-ului v-a avea alt cod. Această metodă este una mai lentă însă este cea mai eficientă. O altă modalitate ar fi putut fi folosirea expresiilor regulate. În această metodă se lua cuvântul și se verificau anumite tipare(de exemplu dacă cuvântul începe cu „http” sau dacă conține unul din domeniile de internet(ro,com,de etc).Metoda cu expresii regulate ar fi fost mult mai rapidă însă ar fi fost și mult mai ineficientă.Odată ce știm că cuvântul este un url putem parsa informațiile din el.

Pentru a parsa url-ul și pentru a scoate informațiile utile vom folosi un html parser predefinit de java. În informatică parsarea este parcurgerea și analizarea unui text, cu identificarea atomiilor care îi corespund în raport cu o gramatică formală.În cazul nostru, documentul de parsat este întreaga pagină html.

Algoritmul de parsare a unui text HTML nu este foarte complicat.Acesta ia ,într-o primă fază, întregul text html și îl consideră ca fiind un text normal.Fiecare algoritm de parsare conține un limbaj lexical. Limbajul lexical nu este altceva decât totalitatea cuvintelor cheie ale limbajului respectiv(în cazul nostru este vorba despre HTML).

Ce-a de-a doua fază a algoritmului constă în parcurgerea mesajului caracter cu caracter și verificarea cu limbajul lexical. Astfel dacă în parcurgerea mesajului vom găsi caracterul „<” acesta nu indică altceva decât începerea unui tag html. În momentul în care se găsește un caracter special se va memora poziția lui. Caracterul „/” ar putea însemna sfârșitul unui tag html doar dacă este precedat de caracterul special „<”.

Ce-a de-a treia fază a algoritmului reprezintă parcurgerea mesajului cuvânt cu cuvânt. În momentul în care a fost găsit un cuvânt cheie cum ar fi „body” este păstrată poziția sa în mesaj.Dacă de exemplu avem același element de mai multe ori se va păstra fiecare locație a sa

În ultima fază a algoritmului, html-ul nostru este practic tradus într-un vector în care fiecare element are un anumit simbol și o anumită poziție. Nu ne rămâne de făcut decât să supunem acel vector unor serii de reguli pentru a scoate cuvintele utile din html. Se consideră text util, textul dintre ultimul tag de „>” și primul tag de „/” care este precedat de untag de „>”. De exemplu dacă avem următoarea sintaxă html „<tr><td><p>Acesta este un mesaj</p></td></tr>” textul „acesta este un mesaj” va fi considerat ca și text datorită regulilor descrise. Bineînțeles acesta este un exemplu în care a fost descrisă o regulă simplă de aflare a conținutului. În figura 1.1 avem schema modului de extragere de text dintr-o pagină html.

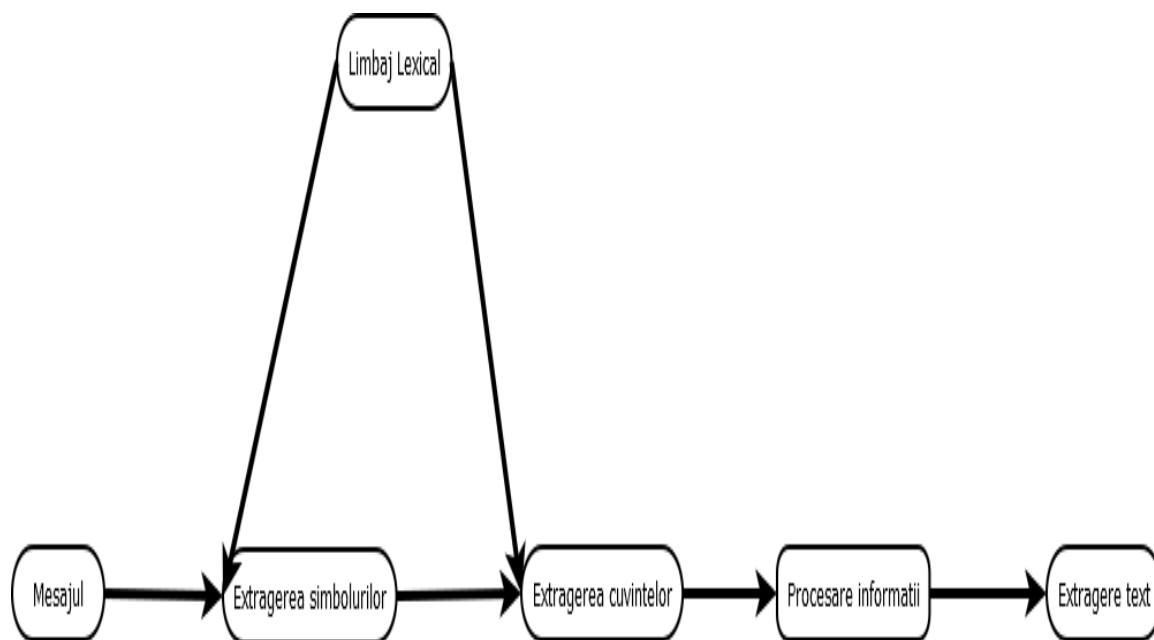


Figura 4.1 Schema modului de extragere de text html

4.1.1.2. Modulul de traducere al textului

Textele scrise în limbi diferite au fost pentru mult timp o restricție. Multe dintre informații se pierdeau doar pentru că nu exista acces la ele. În zilele noastre însă, în principal datorită creșterii distributivității aplicațiilor, putem folosi texte din alte limbi ca și cum acestea au fost scrise în limbile proprii. Ca și limbă comună tuturor textelor am ales engleza pentru că aceasta este documentată mai bine, sporind eficiența algoritmului.

Ca și în cazul anterior, înainte de a face o traducere avem nevoie de o detecție a limbii utilizate. În cazul în care am fi avut mesajul nostru în engleză și noi am fi încercat să traducem mesajul în engleză, toate atributele de calitate ale proiectului ar fi scăzut considerabil. Pentru traducerea de mesaje am folosit un detector de limbaj. Acesta conține cele mai utilizate cuvinte dintr-o anumită limbă și face căutarea cuvintelor în mod secvențial. Astfel dacă propoziția mea este „Vreau să iau licența” acesta va căuta în baza lui de date cuvintele în funcție de ordinea de apariție în propoziție. Dacă la cuvântul „vreau”, găsește ca și limbă, limba română se oprește din căutat și returnează ca și limbă descoperită limba română. Dacă găsește mai multe limbi, acesta va traduce următoarele cuvinte până când se va ajunge la o detecție de limbă. Dacă după terminarea tuturor cuvintelor din propoziție nu se va detecta nici o limbă, algoritmul va selecta una din ele după anumite criterii de relevanță.

Algoritmul de traducere de texte este unul foarte complex care se face de către mai multe computere. Aceste computere utilizează un proces numit ‘traducere automată statistică’, acesta fiind de fapt un mod elevat de a spune că traducerile sunt generate de computerele noastre pe baza modelelor găsite în cantități mari. Dar să revenim puțin asupra lucrurilor. Dacă ați dori să învățați pe cineva o limbă nouă, ați începe prin a-l învăța cuvinte din vocabular și reguli gramaticale pentru a explica modul de construire al propozițiilor. Un computer poate învăța o limbă străină în același mod – utilizând un vocabular și un set de reguli. Dar limbile sunt complicate și, după cum vă poate spune

orice persoană care învață o limbă străină, există excepții de la aproape orice regulă. Când se încearcă includerea tuturor acestor excepții și a excepțiilor de la respectivele excepții într-un program de computer, calitatea traducerii începe să se deterioreze.

Algoritmul de traducere are o abordare diferită. În loc să încercăm să învățăm computerele toate regulile unei limbi, lăsăm computerele să descopere singure regulile. Computerele realizează acest lucru analizând milioane și milioane de documente care au fost deja traduse de traducători reali. Aceste traduceri provin din cărți, de la organizații ca ONU și de pe site-uri web din întreaga lume. Computerele scanează aceste texte căutând modele importante din punct de vedere statistic, și modele create între traducere și textul original cu probabilitate redusă să fi apărut în mod întâmplător. După ce computerul găsește un model, îl poate utiliza pentru a traduce texte similare pe viitor. După ce acest proces se va fi repetat de miliarde de ori, se vor obține miliarde de modele și un program de computer foarte inteligent. Totuși, pentru unele limbi avem la dispoziție mai puține documente traduse și prin urmare mai puține modele detectate de aplicația noastră software. Din acest motiv, calitatea traducerii variază în funcție de limbă și de perechea de limbi. Suntem conștienți că traducerile nu sunt întotdeauna perfecte, dar furnizând constant noi texte traduse, putem crea programe mai inteligente și traduceri mai bune.

4.1.1.3. Modulul de eliminare a semnelor de punctuație și a cuvintelor uzuale

În momentul în care se scrie un text, majoritatea oamenilor nu sunt atenți asupra modului în care spațiile dintre cuvinte și semnele de punctuație sunt puse. De aceea, semnul de punctuație poate apărea înaintea cuvântului sau după. În acest sens un algoritm de eliminare a semnelor de punctuație este necesar. Algoritmul verifică dacă cuvântul începe sau se termină cu un semn de punctuație. În oricare din aceste situații semnul de punctuație se va înlătura.

Dacă dorim să facem o filtrare de mesaje este știut faptul că, cuvintele de legătură au o prezență mai mare în fiecare propoziție. Prezența lor mare va determina ca ele să aibă conform algoritmului o pondere foarte mare și să încurce foarte mult filtrarea.

Conform studiilor, masa de bază a unui vocabular este în principal formată din puține cuvinte. Pentru a face acest lucru, vom extrage din mesajele noastre cele mai utilizate 250 de cuvinte din limba engleză. Dacă extrăgeam mai multe cuvinte riscam faptul că, mulți termeni relevanți pentru mesajul meu să nu fie incluși. Din algoritm au fost excluse părți de vorbire precum prepoziții, conjuncții, pronume, interjecții, verbe uzuale (am, has) etc. În tabelul de mai jos ne sunt date toate cuvintele care fac parte din masa vocabularului englez și au fost excluse din algoritmul meu.

the	be	to	of	and	a	in	that	have	I
it	for	not	on	with	as	you	do	at	this
but	his	by	from	they	we	say	her	she	or
an	will	my	one	all	would	there	their	what	so
up	out	if	about	who	get	which	go	me	when
make	can	like	time	no	just	him	know	take	people
into	year	your	good	some	could	them	see	other	than
then	now	look	only	come	back	after	use	two	how

Tabela 4.2 Cele mai întâlnite 80 de cuvinte din vocabularul englez[5]

4.2. Modulul de comparare a cuvintelor

Modulul de comparare a cuvintelor nu este un modul pre-filtrare ci aparține algoritmului de filtrare de aceea a fost descris separat. Acest modul se ocupă cu compararea a două cuvinte și verifică dacă acestea reprezintă același lucru.

În filtrarea de mesaje se întâmplă foarte rar să se găsească cuvinte identice. De exemplu cuvintele “filtering” și “filter” sunt două forme ale aceluiași verb “to filter”, însă deoarece nu sunt identice nu s-ar considera egale în comparare. După cum putem observa, cuvântul de bază din cele două verbe este filter. Această ar fi și ideea care ar putea combate această problemă, și anume, folosirea unui algoritm de extragere a bazei cuvântului.

Pentru aceasta s-a instalat pe calculatorul propriu un dicționar al limbii engleze cu care s-a făcut legătura. Acest dicționar conține în baza lui de date toate mapările dintre cuvinte și cuvintele lor de bază. Ideea de a folosi această metodă de comparare consider că este foarte bună deoarece cuvântul de bază poate conține și ideea de bază a semnificației lui. De exemplu substantivul “filter” și verbul “filtering” au aceeași idee de baza(filtrarea).

Oamenii comunică deseori aceeași idee însă în moduri diferite. În acest caz extragerea ideii de bază nu este suficientă ci este nevoie și de un analizator de sinonime. De exemplu “people” și “humans” sunt sinonime însă în momentul comparării lor acestea nu vor fi pozitiv comparate niciodată.

Dicționarul ne-a fost util și în momentul acesta deoarece acesta conține în baza lui de date legăturile dintre toate sinonimele. Algoritmul de comparare verifică într-o primă etapă dacă cuvintele sunt egale („filtering cu filtering”). În ce-a de-a doua etapă se va verifica dacă cuvintele de bază din cele două cuvinte sunt egale („filter” cu „filtering”) urmând ca în a treia etapă să se verifice dacă sinonime ale cuvintelor de bază sunt identice.

Deși pare riscant algoritmul nostru are ca scop găsirea oricăror două cuvinte care au ceva comun între ele. Șansele ca mai multe cuvinte dintr-o propoziție să aibă legătură cu mai multe cuvinte dintr-o altă propoziție fără ca logica propozițiilor să aibă ceva comun sunt infime.

De notificat ar mai fi aici și faptul că, în momentul în care se elimină cuvintele de legătură, se folosește algoritmul de comparare între cuvântul respectiv și cuvântul de eliminat.

4.3. Algoritm de clasificare

În abordarea tradițională, frecvența termenilor sau inversul frecvenței termenilor sunt măsurători folosite în clasificarea documentelor. Aceste măsurători indică cât de important este un cuvânt pentru un document într-o anumită sintaxă. În frecvența inversă a termenilor, importanța unui cuvânt pentru un mesaj este cu atât mai mare cu cât cuvântul este folosit de mai multe ori și este cu atât mai mică când cuvântul este folosit de mai puține ori. Similaritățile celor două documente pot fi calculate cu ajutorul similarităților ponderilor cuvintelor din fiecare document.

Inversul frecvenței termenilor după pondere este o modalitate populară de sortare și este folosită atât în motoare de căutare cât și în diferiți algoritmi de clasificare de text. Cu ajutorul algoritmului care are la bază inversul frecvenței termenilor, cuvintele care se întâmplă foarte des în propoziție, de exemplu „the” sau „end” vor avea ponderi mai mici. Totuși în cazul în care textul conține foarte multe cuvinte de legătură și foarte puține cuvinte, acesta ar putea apărea ca fiind relevant. De aceea s-a ales ca aceste cuvinte să se elimine complet din algoritm și să ia ponderea 0.

Având în vedere faptul că aceste metode nu sunt potrivite complet cu clasificarea noastră vom propune o nouă metodă care va calcula proximitatea unui text față de postările din micro-blogging. Principala provocare aici ar putea fi care cuvinte importante sunt cu adevărat importante în contextul de față și care vor fi doar cuvinte de legătură. Modulul de eliminare a cuvintelor comune prezentat anterior s-a ocupat de această problemă, în mesajele de clasificat rămânând doar cu cuvinte cheie.

Spre exemplu să presupunem că avem o serie de postări în micro-blogging și că avem un mesaj care va sorta mesajele noastre ca fiind relevante sau irelevante cu topicul de discuție. Să presupunem de asemenea faptul că cuvântul MySQL va apărea des în propoziția noastră deoarece postarea noastră se referă la baze de date. În lista de postări extrase din microblogging vor putea fi sute de mesaje iar dintre acestea doar puține se vor referi la baze de date. Așadar frecvența apariției cuvântului MySQL este mai mică deoarece ponderea cuvintelor în care a apărut cuvântul „MySQL” este mai mică. Dacă însă vom extrage un cuvânt comun precum „children” sau „tomorrow”, este de așteptat ca aceste cuvinte să aibă o frecvență de apariție mai mare și să apară mai des decât cuvântul MySQL. Algoritmul nostru ar trebui să identifice faptul că „tomorrow” este un cuvânt comun și să-i atribuie o pondere cât mai mică. Cuvintele precum „and” și alte cuvinte de legătură nu vor fi luate în considerare.

Prin folosirea ideii de comparare după frecvența termenilor, vom putea calcula cât de relevant este cuvântul respectiv pentru mesajul nostru. Vom considera mesajul nostru ca fiind unul local, iar mesajul provenit din server-ul de micro-blogging unul global. În figura de mai jos ne este prezentat modul în care se va calcula ponderea unui cuvânt în mesajul local:

$$\text{ponderea locală} = \frac{\text{numărul de apariții al cuvântului}}{\text{numărul total de cuvinte}}$$

Figura 4.3 Modul de calcul al mesajului local

unde deîmpărțitul va reprezenta numărul de apariții a cuvântului în mesajul local iar împărțitorul va reprezenta suma aparițiilor tuturor cuvintelor în mesajul local.

Relevanța unui text față de o postare nu va fi așadar altceva decât suma ponderilor din aceea postare raportată la numărul de cuvinte din postarea respectivă.

$$\text{relevanța} = \frac{\text{totalul ponderilor}}{\text{numărul total de cuvinte}}$$

Figura 4.4 Relevanța unei postări

Dimensiunea textului pentru fiecare postare este limitată la 500 de caractere în microblogging-ul local ceea ce face algoritmul de clasificare mult mai greu. Pentru a trece peste această problemă se folosesc link-uri care conțin text. În zilele noastre este un lucru comun să se folosească url-uri pentru o mai bună documentare a informației prezentate. Conținutul url-ului va fi extras cu ajutorul unui serviciu de extragere de conținut HTML și adăugat conținutului fiecărei postări.

Prin algoritmul nostru, vom putea să calculăm relevanța unui text sau a unei postări. Scorul are limita inferioară 0 iar limita superioară 1. Pentru fiecare mesaj adăugat ponderea mesajelor se schimbă, algoritmul trebuind reluat din nou. Mesajele ambigue pot crea cele mai multe erori de clasificare pe când mesajele clare, cu cuvinte de topic specifice, pot face ca sistemul să dea cele mai bune rezultate. Vom folosi tipul mesajului ca și un indicator din care sistemul trebuie să învețe, în funcție de acesta dându-se valoarea pragului de relevanță.

4.4. Pseudocodul algoritmului

În figura 4.5 este prezentată partea de atribuire de ponderi. Prezentarea este făcută în pseudocod și are rolul de a prezenta mai bine algoritmul folosit.

Inițial vom avea textul local, venit din interfața grafică și textele globale, textele aparținând bazei de date a micro-bloggingului. Textul local va fi inițial despărțit în cuvinte. După despărțirea în cuvinte se va verifica fiecare cuvânt în parte. Dacă cuvântul nostru este cuvânt comun se va trece mai departe. Dacă cuvântul nostru a mai fost găsit în mesaj, atunci se va lua ponderea sa și se va incrementa cu 1. Dacă cuvântul nostru nu este cuvânt comun și este folosit pentru prima dată în mesajul nostru, atunci se va insera în lista de ponderi cuvântul nostru cu ponderea 1.

După ce toate cuvintele au fost parcurse, se va lua lista de ponderi și se va împărți fiecare cuvânt la numărul total de cuvinte. Vom face acest lucru pentru ca ponderea fiecărui cuvânt să fie după formula descrisă în figura 4.3. În figura 4.4 va fi prezentată o monștră a algoritmului folosit.

```

var textLocal=ui.getTestLocal();
var lista<texte>=database.getTexteFromMicroblogging();

vector cuvinteLocale=textLocal.split();
listaPonderi=null;

for(i in cuvinteLocale)
    if(nuEsteCuvantComun(i))
        if(listaPonderi.get(i)==0)
            listaPonderi.add(i,1);
        else
            pondere=listaPonderi.get(i);
            pondere=pondere+1;
            listaPonderi.add(i,pondere);
        endif
    endif
endfor

for(i in listaPonderi)
    pondereFinala=listaPonderi[i]/nrCuvinte
endfor

```

Figura 4.5 Pseudocod aplicare ponderi

După ce sunt atribuite ponderile pentru fiecare cuvânt va urma ca toate mesajele din baza de date să fi parcurse. Fiecare mesaj va fi despărțit la rândul său în cuvinte. După aceea urmează ca fiecare cuvânt din lista de ponderi să se compare cu fiecare cuvânt din mesaj. Dacă cuvintele sunt identice se va incrementa o variabilă cu valoarea ponderii cuvântului. După calcularea tuturor ponderilor, se va împărți ponderea rezultată la numărul de cuvinte al textului global. Dacă această valoare este mai mare decât un anumit prag atunci mesajul este considerat relevant. În figura 4.6 este prezentat algoritmul de atribuire a ponderilor mesajelor globale și verificarea relevanței acestora.

De notat ar mai putea fi și faptul că, modulele de prefiltrare nu au fost introduse în pseudocod. Astfel, se consideră că textele globale nu conțin url (conținutul url-ului fiind deja extras) și sunt în limba engleză. De asemenea egalitatea din pseudocodul nostru se va face cu ajutorul utilitarului Wordnet creat de universitatea din Princeton. Mesajele care vor trece pragul de relevanță se vor trece într-o listă separată însă pentru înțelegerea algoritmului s-a adăugat doar un comentariu sugestiv.

În concluzie folosirea pseudocodului a avut rolul de a clarifica modalitatea de funcționare a algoritmului și punerea în practică a formulelor utilizate.

```

var textLocal=ui.getTestLocal();
var lista<texte>=database.getTexteFromMicroblogging();

for(i in lista<texte>)
    pondereTextGlobal=0;
    cuvinte textGlobal=i.split();
    for(j in listaPonderi)
        for(k in textGlobal)
            if(j.cuvant==k)
                pondereTextGlobal=pondereTextGlobal+j.pondere;
            endif
        endfor
    endfor
    pondereTextGlobal=pondereTextGlobal/nrCuvinteTextGlobal;
    if(pondereTextGlobal)>valoare de prag
        //mesajul este relevant
    endif
endfor

```

Figura 4.6 Pseudocod algoritm filtrare

4.5. Protocolul HTTP

Hypertext Transfer Protocol sau cunoscut sub numele de HTTP este un protocol la nivel de aplicație ce constituie fundamentul pentru World Wide Web (www). Acesta are la bază protocolul de comunicare TCP/IP folosit pentru a trimite date pe Internet. Aceste date pot fi fișiere HTML, fișiere audio sau video, poze sau pur și simplu text. Specificația HTTP exemplifică modul în care sunt construite și trimise datele clienților către server și modul în care serverele răspund acestor cereri.

Există trei particularități care fac protocolul HTTP să fie simplu, dar și puternic în același timp:

- HTTP este “connectionless”: se referă la faptul că browser-ul clientului inițiază o cerere server-ului după care se deconectează de la server și așteaptă răspunsul. Server-ul procesează cererea, se reconectează cu clientul și îi trimite răspunsul acestuia.
- HTTP este “media independent”: se referă la faptul că prin intermediul protocolului se pot trimite orice fel de formate media atâta timp cât server-ul și clientul știu să proceseze conținutul acestora.
- HTTP este “stateless”: din prima caracteristică menționată mai sus rezultă faptul că atât server-ul cât și clientul au cunoștință unul de celălalt doar pe perioada unei cereri, după furnizarea răspunsului acestia nu mai au nicio legătură.

În Figura 4.7. de mai jos este prezentată arhitectura de bază a tuturor aplicațiilor care se bazează pe protocolul HTTP.

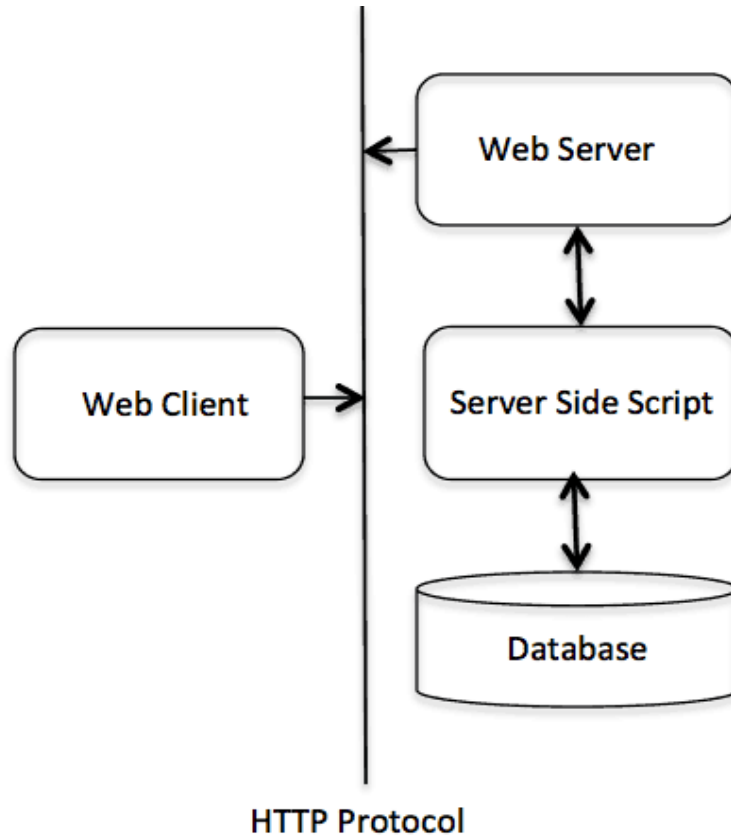


Figura 4.7 Arhitectura aplicațiilor bazate pe HTTP

Clientul trimite o cerere server-ului prin intermediul protocolului HTTP cu toate informațiile de care acesta are nevoie. Server-ul procesează datele, interacționează eventual cu o bază de date după care îi răspunde clientului cu un cod de succes sau de eroare și cu informațiile cerute în caz de succes.

Capitolul 5. Proiectare de Detaliu si Implementare

În acest capitol sunt descrise specificațiile aplicației, arhitectura generală și diagramele implicate în realizarea proiectului.

5.1. Modul de definire al sistemului

În figura 5.1 este prezentat un curs complet al execuției principale aplicației. O primă acțiune în momentul pornirii este aceea de a insera un text care va fi apoi considerat ca și text de pornire pentru filtrarea de mesaje. De notificat este faptul ca acest fișier poate fi scris doar în engleză, trebuie sa aibă o lungime de minim 15 cuvinte și nu poate conține url-uri.

5.1.1. Cursul de execuție al aplicației

Aplicația conține 2 fâșii de execuție care sunt independente până în momentul clasificării, comunicarea dintre ele în acel moment devenind crucială.

După ce utilizatorul apasă butonul „Search” mesajul de căutat va avea de suferit o serie de transformări. Prima dintre ele este aceea că îi vor fi eliminate semnele de punctuație din interiorul mesajului. Mesajul trebuie să conțină doar cuvinte. Urmează mai apoi partea de traducere. Aici se va verifica dacă mesajul de tradus este în limba engleză. În cazul în care limba detectată este alta decât engleza se va traduce mesajul în limba engleză. Următorul modul este acela de extragere de cuvinte uzuale. După ieșirea din acest modul, mesajul nostru nu va conține decât cuvintele relevante cuvintele de legătură fiind eliminate. În acest moment se va putea aplica algoritmul de atribuire a ponderilor mesajului nostru.

De cealaltă parte vom avea partea de modulele care se vor ocupa cu mesajele provenite din server-ul de micro-blogging. Într-o primă etapă se vor extrage mesajele din baza de date. Pentru ca aceasta să poate fi executată va fi nevoie de o conexiune la baza de date a micro-bloggingului respectiv. Odată ce mesajele sunt salvate din baza de date și preluate într-o listă, se vor apela aproximativ aceeași pași de pre-clasificare ai mesajului ca și în cazul anterior. Singura diferență notabilă ar fi aici modulul de extragere de informație din url. În acest modul se verifică dacă textul conține un url, iar dacă acesta va conține se va appenda la textul nostru mesajul din interiorul url-ului.

Odată ajunse în etapa de clasificare, mesajele vor fi parcurse pentru compararea cuvintelor și atribuirea ponderilor. După această etapă se va verifica dacă suma ponderiilor din fiecare text este mai mare decât un anumit prag iar mesajele relevante se vor afișa pe ecran. În acest moment toți algoritmi au luat sfârșit, operația putând fi executată din nou, fie pentru același mesaj, fie pentru mesaje diferite.

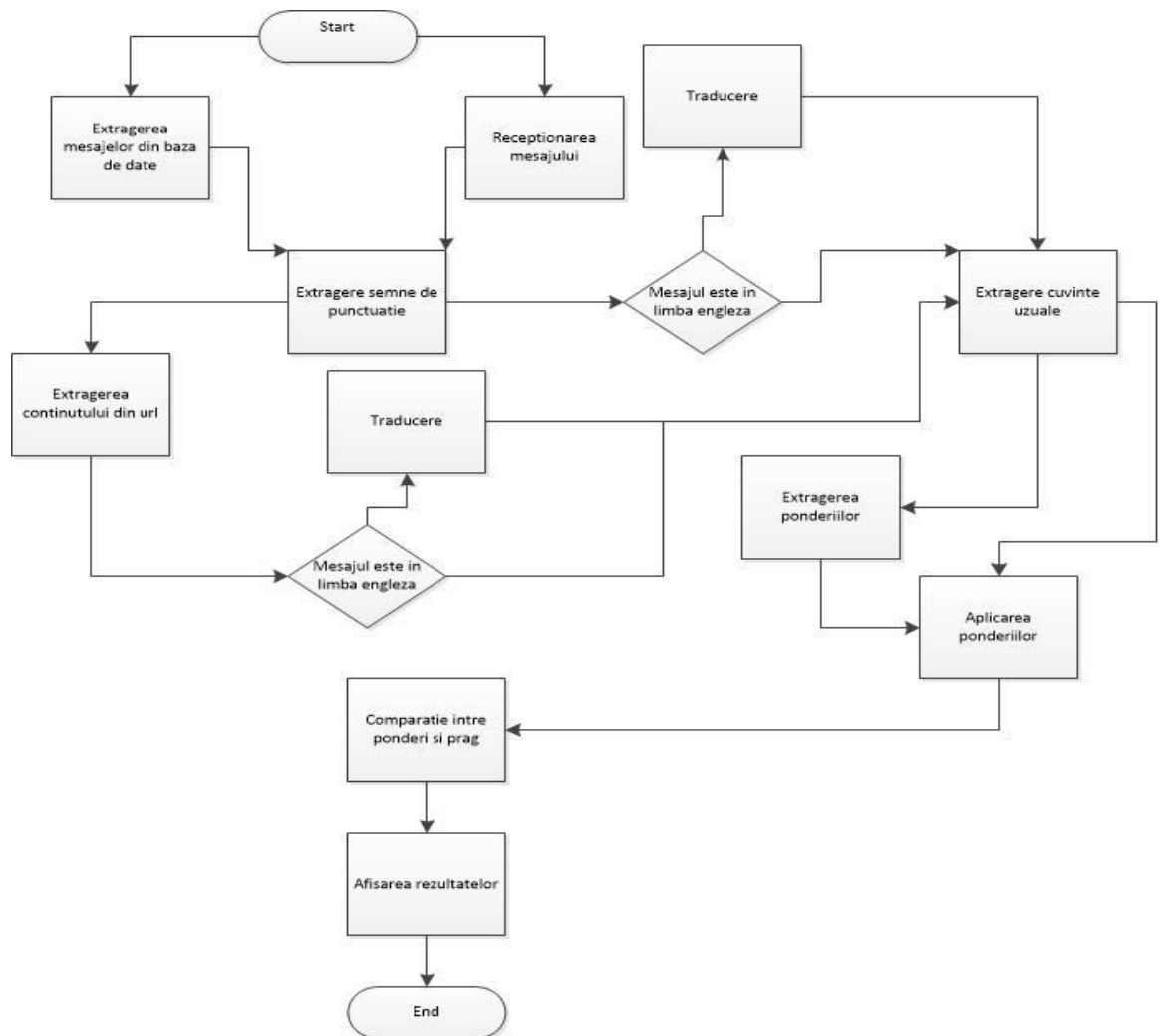


Figura 5.1 Flow-ul de execuție al aplicației

5.1.2. Arhitectura conceptuală generală

Aplicația noastră are ca scop principal filtrarea mesajelor provenite dintr-un serviciu de micro-blogging și de aceea am considerat ca în diagrama conceptuală să apară și aplicația de micro-blogging. De asemenea, filtrul nostru mai folosește și o bază de date de dicționar de aceea aplicația de dicționar de la wordnet va mai apărea și ea în diagramă. Nu în ultimul rând aplicația mai folosește și un serviciu de translate folosit de microsoft care va apărea și el în diagramă.

Utilizatorul va putea interacționa cu toate cele 3 aplicații, aplicația de dicționar fiind disponibilă doar în varianta desktop. De notificat ar fi și faptul că linia dublă din desen semnifică interacțiune pe când linia simplă reprezintă apel.

Aplicația de dicționar este una simplă și este compusă din partea de „User Interface” și de clasele care se ocupă cu maparea la baza de date. Dacă se va dori faptul

ca un nou cuvânt sau o nouă bază a unui cuvânt va trebui să se interacționeze cu această aplicație.

Aplicația de micro-blogging este cea pe care se va face filtrarea de mesaje. Ca și în cazul anterior, pentru a putea crea un mesaj va trebui ca mesajul să fie inserat în aplicația de micro-blogging. Aplicația de micro-blogging are la bază tehnologia Spring. În figură au fost prezentate pachetele din aplicație într-o formă simplistă pentru ușurința înțelegerii.

Aplicația de filtrare este aplicația centrală a diagramei. Motivul pentru care baza de date este apelată din pachetul de servicii este acela că, baza de date este apelată prin intermediul unei librării. Legătura cu aplicația de micro-blogging se va face doar la nivelul bazei de date și nu va exista cod comun între cele două aplicații. O ultimă componentă folosită este api-ul oferit de microsoft pentru traducerea mesajelor. Acesta nu se află pe calculatorul local de aceea este doar apelat în aplicația noastră el rulând în cloud.

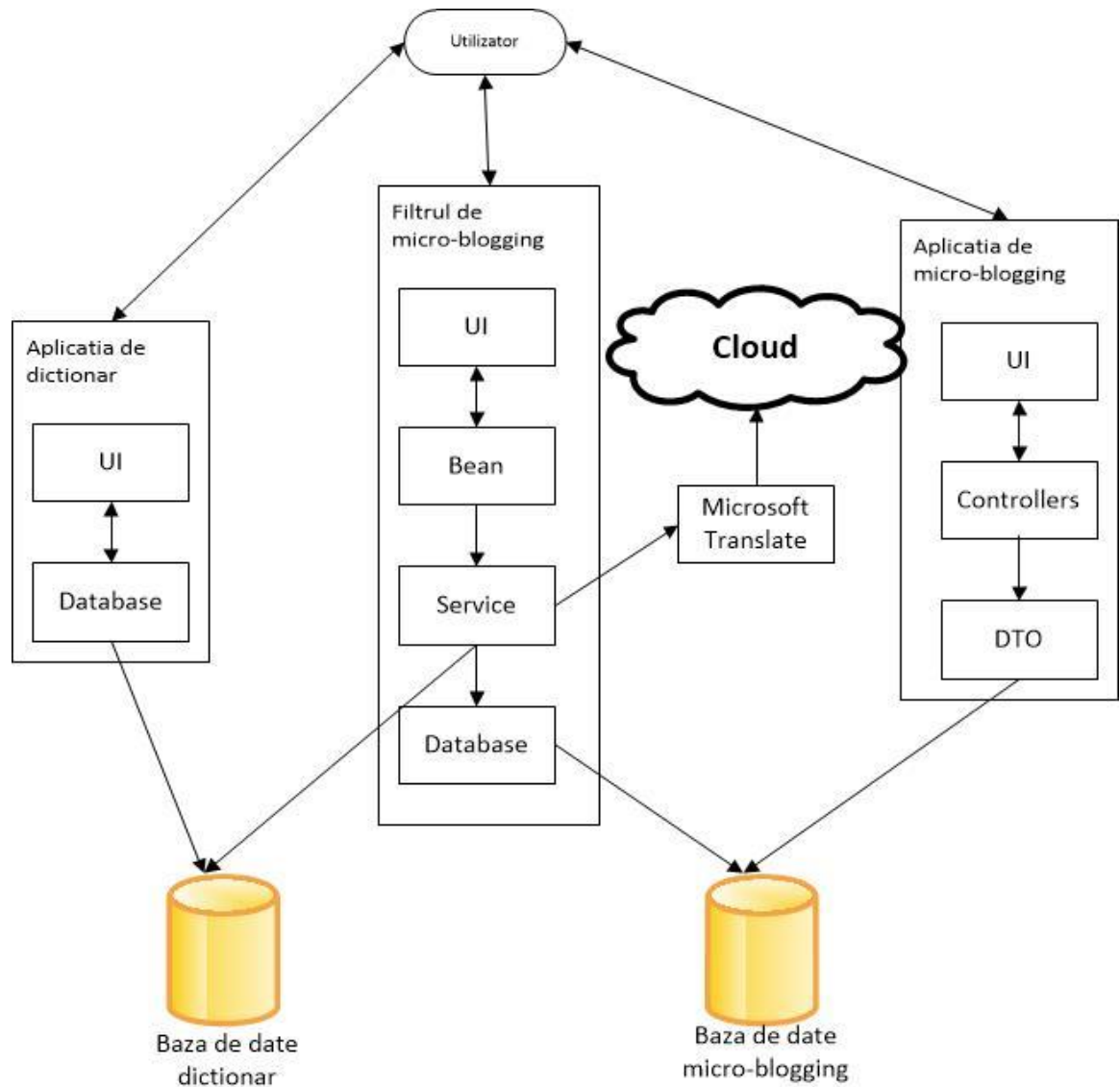


Figura 5.2 Diagrama conceptuală a aplicației

5.2. Modelul cazurilor de utilizare

Sistemul propus și dezvoltat este caracterizat de următoarele cazuri de utilizare:

- extragerea url-ului din mesaj
- traducerea mesajelor care nu sunt în engleză
- eliminarea semnelor de punctuație și cuvintelor uzuale
- traducerea textelor în mod corect chiar dacă acestea nu conțin aceleași forme ale cuvintelor

5.2.1. Extragerea url-ului din mesaj

Actorul primar:Sistemul

Participanți și interese: Utilizatorul setează un mesaj de filtrare

Precondiții: Utilizatorul cunoaște aplicația și scopul ei

Postcondiții: Afișarea mesajului împreună cu url-ul.

Scenariul principal de succes:

- utilizatorul pornește sistemul de filtrare
- sistemul găsește mesajul care conține url-ul
- sistemul extrage conținutul url-ului și îl apendează mesajului
- mesajul este afișat împreună cu url-ul.

Fluxuri alternative(de eroare)

- în caz că sistemul nu poate identifica în mod corect url-ul se va continua algoritmul de filtrare iar url-ul va fi considerat ca și un cuvânt obișnuit.

5.2.2. Traducerea mesajelor care nu sunt în engleză

Actorul primar:Sistemul

Participanți și interese: Utilizatorul dorește să vizualizeze și informațiile care nu sunt în engleză

Precondiții: Utilizatorul cunoaște aplicația și scopul ei

Postcondiții:Afișarea mesajelor din alte limbi decât engleza.

Scenariul principal de succes:

- utilizatorul pornește sistemul de filtrare
- sistemul găsește un mesaj care nu este în engleză
- sistemul traduce mesajul și îl pune în locul mesajului existent
- mesajul este afișat în limba engleză

Fluxuri alternative(de eroare)

- în caz că sistemul nu cunoaște cuvântul de tradus, se va traduce cu cea mai apropiată formă a sa.

5.2.3. *Eliminarea semnelor de punctuație si a cuvintelor uzuale*

Actorul primar: Utilizatorul

Participanți și interese: Utilizatorul setează un mesaj de filtrare

Precondiții: Utilizatorul cunoaște aplicația și scopul ei

Postcondiții: Sortarea mesajelor în mod corect

Scenariul principal de succes:

- utilizatorul pornește sistemul de filtrare
- sistemul a întâlnit un semn de punctuație sau un cuvânt uzual
- sistemul elimină din cuvânt semnul de punctuație respectiv elimină cuvântul uzual din mesaj
- sortarea mesajelor a fost făcută în mod corect

Fluxuri alternative(de eroare):

- în cazul în care punctuația face parte din cuvânt, cuvântul nu va fi supus în mod corect filtrării
- în cazul în care cuvântul uzual este un cuvânt relevant pentru aplicație(extrem de rar), ponderea cuvântului se va pierde.

5.2.4. *Traducerea cuvintelor în mod corect chiar dacă acestea nu conțin aceleași forme ale cuvintelor.*

Actorul primar: Utilizatorul

Participanți și interese: Utilizatorul setează un mesaj de filtrare

Precondiții: Utilizatorul cunoaște aplicația și scopul ei

Postcondiții: Sortarea mesajelor în mod corect

Scenariul principal de succes:

- utilizatorul pornește sistemul de filtrare
- sistemul desparte mesajele din baza de date și le compară cu cuvintele relevante din mesajul de filtrat
- sistemul verifică la comparare dacă cuvintele de comparat au același cuvânt de bază sau sunt sinonime

Fluxuri alternative(de eroare)

- cuvântul de bază lipsește din dicționarul local
- cuvântul de comparat nu apare ca și sinonim în dicționarul local

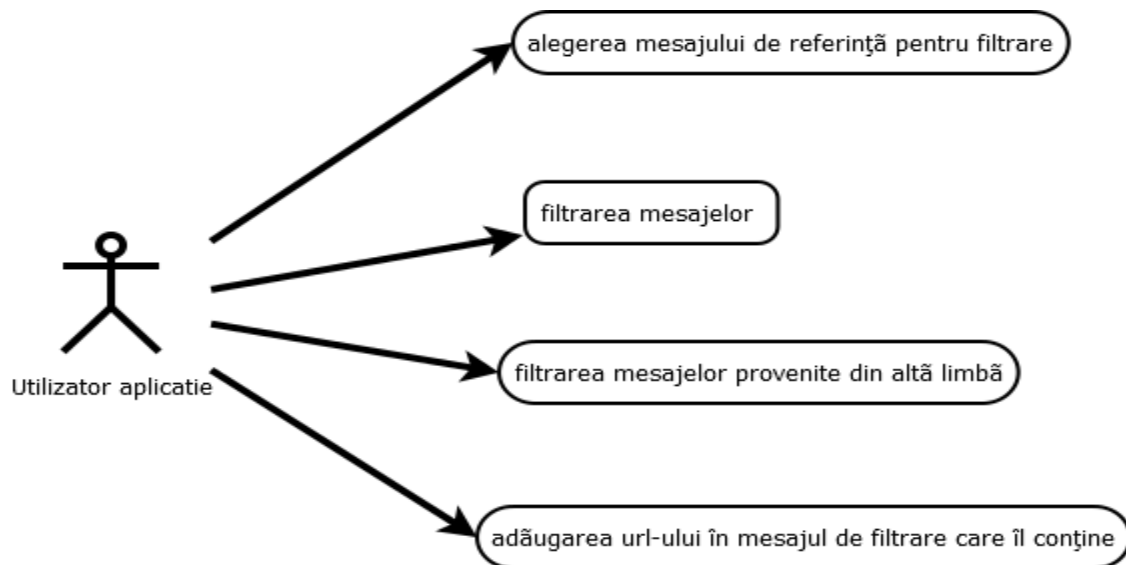


Figura 5.3 Diagrama use-case a aplicației

5.3. Diagrama cazurilor de utilizare

În figura 5.2 este prezentată diagramă cazurilor de utilizare. Actorul este utilizatorul care poate efectua una din următoarele operațiuni:

- alegerea mesajului de referință pentru filtrare
- filtrarea mesajelor
- filtrarea mesajelor provenite din altă limbă
- adăugarea url-ului în mesajul de filtrare care îl conține

5.4. Diagrama de pachete

Lucrarea de față a fost creată folosind tehnologia JSF ca și tehnologie de bază iar ca și tehnologie de back-end clasele au fost grupate în pachete astfel încât să se respecte regulile unei aplicații pe layere. Primul pachet este cel de view în care avem paginiile JSF ale aplicației. Cel de-al doilea pachet este pachetul de styles în care vom păstra fișierele css ale aplicației. De notat ar fi aici și faptul că aceste două clase sunt cele care sunt responsabile cu interfața grafică a utilizatorului.

Cel de-al treilea pachet este cel de bean. Acesta este cel care face legătura între evenimentele care se produc în interfața grafică și logica programului. Pachetul service este pachetul responsabil cu logica din spatele programului urmând iar pachetul database să facă legătura la baza de date.

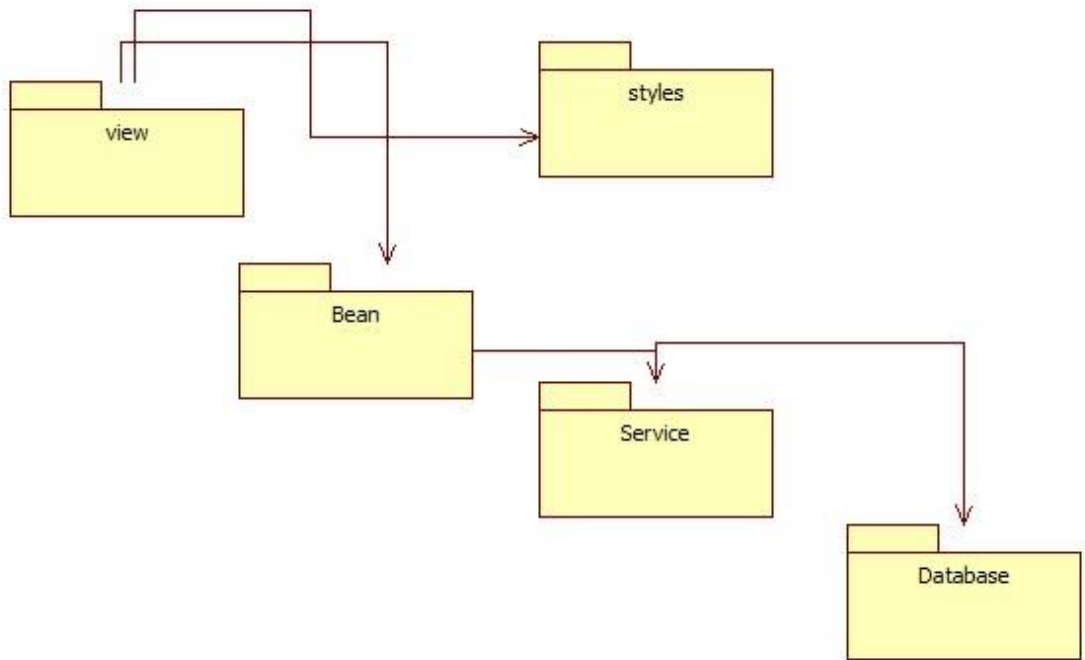


Figura 5.4 Diagrama de pachete

5.5. Diagrama de clase

În figura 5.5 se evidențiază diagrama de clase. Se poate observa faptul că este respectată o arhitectură pe layer. În cele ce urmează se vor explica clasele împreună cu metodele sale cele mai importante

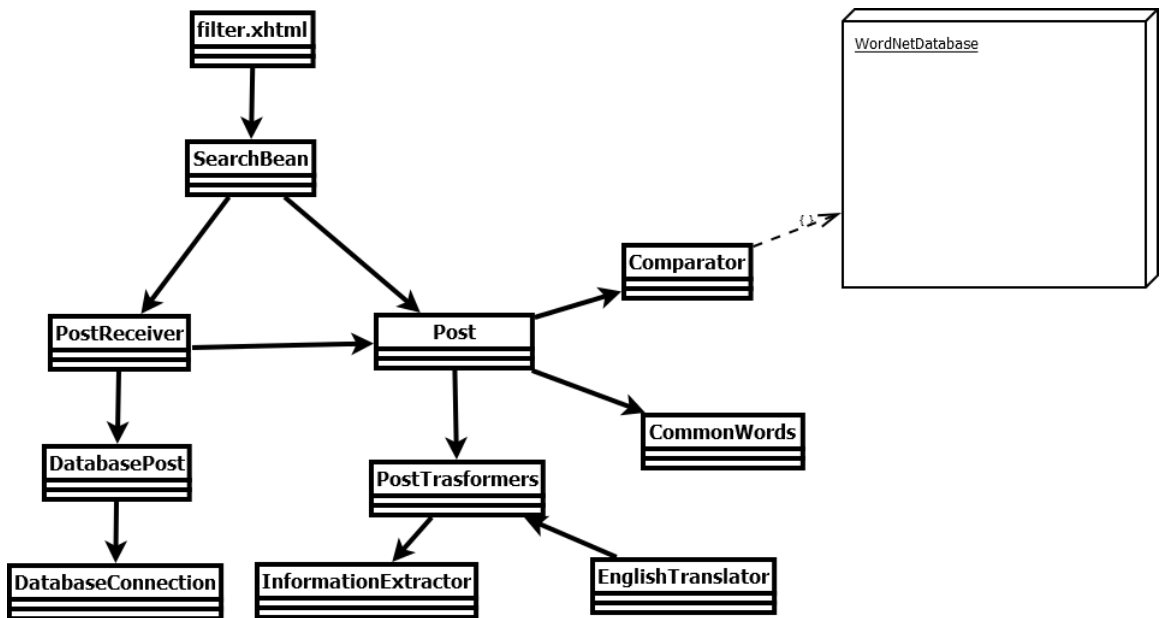


Figura 5.5 Diagrama de clase

5.5.1. *filter.xhtml*

Acest fișier este un fișier jsf specific în care este creată interfața grafică a aplicației. Conține practic o intrare de tip input în care utilizatorul poate insera textul dorit spre filtrare și un button care să trimită textul mai departe. De notificat ar fi faptul că interfața grafică nu este punctul forte al aplicației, punctul forte constând în aplicarea algoritmului.

5.5.2. *SearchBean*

Framework-ul de JSF presupune existența a cel puțin două entități. O pagină specifică JSF și un Java Managed Bean. În cazul nostru, SearchBean nu este altceva decât un Managed Bean. Acesta are un atribut de instanță searchText care se va update cu valoare din pagina jsf. De asemenea în interiorul acestui bean mai există și metoda de listener a butonului.

În momentul în care butonul este apăsat, execută două operații importante. Prima dintre ele este aceea că se apelează clasa PostReceiver care va aduce postările din baza de date. A doua operație importantă este aceea că, se va instanția o clasa Post cu mesajul din interfața grafică și se va porni algoritmul de filtrare de mesaje.

5.5.3. *PostReceiver*

Este clasa care se ocupă cu transformarea obiectelor provenite de la baza de date, în instanțe ale clasei Post. În această clasă, toate datele din baza de date vor fi încapsulate într-o listă de obiecte de tip „Post”. De notificat aici ar mai putea fi și faptul că această clasă s-a adăugat și din nevoia de respectare a arhitecturii pe layere, din managed bean nefiind îngăduită apelarea unei clase din pachetul de baze de date.

5.5.4. *DatabasePost și DatabaseConnection*

Aceste clase fac parte din pachetul „database” și sunt practic clasele care aduc datele din baza de date. Clasa „DatabaseConnection”, după cum sugestiv îi spune numele se ocupă cu conectarea la baza de date a micro-bloggingului iar în clasa „DatabasePost” avem metoda getAllPosts() care aduce toate datele din baza de date. De notificat aici ar mai putea fi și faptul că, la crearea clasei DatabaseConnection s-a folosit design-patternul Singleton

5.5.5. *Post*

Aceasta este practic clasa de bază a aplicației deoarece conține atât metode de instanțiere a variabilelor cât și algoritmul de filtrare propriu-zisă a mesajului. Din SearchFilter se va instanția mai întâi clasa Post cu mesajul din interfața grafică. În momentul în care s-a setat mesajul algoritmul de atribuirea a ponderilor începe.

În metoda splitMessage nu numai că se va despărți mesajul nostru în cuvinte însă le va fi atribuită și o pondere cuvintelor respective. Astfel, după despărțirea mesajului în cuvinte se va apela metoda removePunctuation(). Această metodă este una simplă, care

verifică dacă cuvântul începe sau se termină cu un semn de punctuație și îl elimină. După îndepărtarea semnelor de punctuație se va apela clasa `CommonWords` care verifică dacă cuvântul este unul uzual sau nu. Dacă nu este uzual și este primul cuvânt din propoziție cuvântul va fi adăugat unui vector și va avea ponderea 1. Dacă cuvântul este unul uzual atunci algoritmul îl va ignora complet și va merge mai departe. De asemenea dacă cuvântul nu este primul din propoziție, atunci cuvântul se va compara, folosind clasa `Comparator`, cu toate cuvintele din vector. Dacă compararea nu este una pozitivă, cuvântul va fi și el adăugat în vector cu ponderea 1. În momentul în care compararea este una pozitivă, cuvântul din vector cu care s-a făcut compararea va fi incrementat cu 1. Prin această modalitate, cuvintele care apar de cele mai multe ori vor avea o pondere mai mare decât cele care apar de mai puține ori. La final, vectorul se va mai parcurge odată, de această dată ponderile fiind împărțite la numărul de cuvinte utile (care nu sunt comune).

În metoda `seeRelevant()` se va face practic algoritmul de sortare în funcție de ponderile calculate în metoda de mai sus. În această metoda se va lua fiecare postare în parte și se va extrage mesajul. Prima operație care se va executa este aceea de transformare a mesajului (mesajul nu este pe engleză sau conține un url). După transformarea mesajului, se vor compara cuvintele din mesajul local cu cuvintele din mesajul din micro-blogging. Compararea se va face și de această dată cu ajutorul clasei `Comparator`. Dacă compararea va returna un rezultat pozitiv se va reține ponderea locală. La sfârșitul parcurgerii se va face o sumă pe toate ponderile reținute și se va împărți ponderea rezultată la numărul total de cuvinte.

5.5.6. *PostTransformers*

Această clasă are ca scop principal crearea unei legături dintre algoritmul de filtrare și feature-urile de extragere de conținut url și de traducere de text. Mai întâi se va apela extragerea de conținut url deoarece există șansa ca informația utilă care va fi extrasă de pe url să nu fie în limba engleză fiind nevoie de o traducere suplimentară pentru aceasta.

5.5.7. *InformationExtractor*

Clasa `InformationExtractor` este clasa care se preocupă cu extragerea informațiilor utile din url. Această clasă conține trei metode: `verifyIfWebsite`, `extractInformation` și `getInformationFromUrl`.

În metoda `verifyIfWebsite` se va verifica dacă stringul introdus este un url. Programul va încerca să creeze o conexiune pentru fiecare cuvânt. Dacă codul returnat de server este altul decât 200 atunci cuvântul introdus de mine nu este website. Acest algoritm este considerat a fi lent unul lent însă este foarte eficient. O altă abordare ar putea fi folosirea de expresii regulate. Deși pare mult mai rapid eficiența algoritmului va scădea considerabil.

În metoda `getInformationFromUrl` se va face extragerea propriu-zisă a url-ului. Pentru a extrage datele utile s-a folosit librăria Java `Jsoup`. Motivul pentru care s-a folosit această bibliotecă este ușurința în folosire (nu trebuie decât introduși parametrii în locul potrivit) dar și eficiența sa.

Metoda `extractInformation` este metoda de legătură dintre cele două metode. Astfel, în această metodă se vor evalua cuvintele url-ului unul câte unul și se va verifica dacă cuvântul respectiv este sau nu url. Dacă este url atunci se va extrage informația. De notificat ar mai putea fi și faptul că dacă mesajul conține mai multe url-uri atunci se vor extrage toate url-urile și se vor concatena mesajului.

5.5.8. *EnglishTranslator*

Clasa `EnglishTranslator` este clasa care verifică dacă mesajul este scris într-o altă limbă diferită de limba engleză, în acest caz traducându-l. Spre deosebire de clasa de mai sus această clasă conține o singură metodă intitulată sugestiv `translate`. În această metodă se va verifica mai întâi dacă mesajul scris este într-o limbă străină. Verificarea se va face folosind o librărie de detecție de cod. Această librărie nu conține decât o mare parte din vocabularul limbii respective iar pentru ca aceasta să poată detecta corect este nevoie de minimum 15 cuvinte.

Traducerea textului în limba engleză s-a făcut prin apelarea api-ului „bing translate” oferit de firma Microsoft. Acest api este de tip open-source și ne oferă până la două milioane de caractere de tradus. Pentru a-l putea folosi însă este nevoie de introducerea unui id și unei parole secrete.

5.5.9. *Common words și comparator*

Este clasa care conține într-o variabilă statică cele mai frecvente 250 de cuvinte din limba engleză. Clasa a fost făcută statică pentru a putea suporta eventuale dezvoltări ulterioare (adăugarea unui cuvânt care apare mult prea des în contextul proiectului deși nu face parte din cele mai întâlnite cuvinte). De notificat ar mai fi faptul că, clasa a fost creată în special pentru o mai bună observare a logicii algoritmului.

Clasa `comparator` este o altă clasă mai aparte. Această clasă are o logică proprie însă și face legătura cu o altă aplicație de dicționar. Aplicația „Dicționar” este o aplicație de tip open-source. Pentru folosirea ei vom avea nevoie de baza de date a aplicației (și implicit de aplicația propriu-zisă) și de o librărie oferită tot de aplicație. Librăria va crea un obiect de tip `WordNetDatabase` care va fi de fapt un obiect care mapează toate relațiile bazei de date. Ca și parametru de creare a metodei va trebui introdusă locația locală a bazei de date.

Astfel, metoda `compare` va instanția un obiect de tip `WordNetDatabase`. Metoda noastră are două cuvinte ca și parametrii. Într-o primă etapă metoda compară cele două cuvinte pentru a vedea dacă sunt egale. Acest lucru, în caz de egalitate reduce foarte mult timpul de execuție și sporește eficiența algoritmului. Dacă cuvintele nu sunt egale, cu ajutorul obiectului `WordNetDatabase` se vor extrage formele de bază ale cuvintelor. Apoi tot cu ajutorul obiectului de bază de date se vor extrage într-o listă toate sinonimele cuvântului de bază. Într-o ultimă etapă se vor parcurge listele și se vor compara termenii. Dacă termenii coincid la un moment dat algoritmul se va opri și va returna un răspuns pozitiv. Dacă niciunul din termenii nu coincid, algoritmul va întoarce un răspuns negativ. Deși pare un algoritm primejdios acesta va detecta dacă este vreo legătură între cuvintele comparate.

5.6. Diagrama secvențială

În figura 5.6 este prezentată digrama secvențială a aplicației. Motivul pentru care metodele și declanșatorii au fost scriși cu cifre este acela de a explica diagrama în detaliu pentru a putea fi înțeleasă de către toată lumea. În cele ce urmează voi scrie o mini-legendă a metodelor folosite mai jos

- 1.Utilizatorul introduce textul de convertit și apasă butonul de „Caută”
- 2.Datorită folosirii framework-ului de JSF, evenimentul din cadrul metodei caută() este activat
- 3.Se apelează clasa PostReceiver pentru aducerea mesajelor din baza de date
- 4.Se apelează clasa PostTransformers care mapează clasa Post la baza de date
- 5.Se apelează clasa DatabaseConnection pentru a obține o conexiune la baza de date
- 6.Se returnează conexiunea la baza de date
- 7.Se returnează o listă cu toate mesajele din baza de date.
- 8.Se returnează o listă care conține obiecte de tip Post mapate exact pe baza de date
- 9.Se apelează clasa Post pentru calcularea ponderiilor.
- 10.Se verifică dacă cuvântul de filtrat este cuvânt comun
- 11.Se returnează adevărat dacă cuvântul este cuvânt comun sau fals în caz contrar
- 12.Se compară cuvântul cu alte cuvinte care sunt conținute în mesajul meu și nu sunt comune.
13. Se returnează adevărat dacă cuvintele au aceeași bază sau sunt sinonime ori fals în caz contrar
- 14.Se apelează clasa PostTransformers care se ocupă cu algoritmi de pre-clasificare
- 15.Se apelează clasa InformationExtractor care are ca și rol extragerea url-ului.
- 16.Se returnează mesajul concatenat cu informațiile utile din url. Dacă mesajul nu conține url, se va returna mesajul gol.
- 17.Se apelează clasa EnglishTranslator care are ca și rol traducerea textelor în engleză.
- 18.Se returnează textul tradus în engleză sau același text, în cazul în care textul este deja în engleză.
- 19.Se returnează o variantă transformată a mesajului.
- 20.Se verifică dacă postarea este relevantă.Pentru aceasta fiecare cuvânt relevant din postarea locală va fi comparat cu fiecare cuvânt din postarea provenită din baza de date.
21. Se returnează adevărat dacă cuvintele au aceeași bază sau sunt sinonime ori fals în caz contrar
- 22.Se returnează o listă cu postările relevante
- 23.Datorită frame-workului de JSF, postările relevante vor apărea pe interfața grafică.
- 24.Utilizatorul va poate vizualiza mesajele dorite.

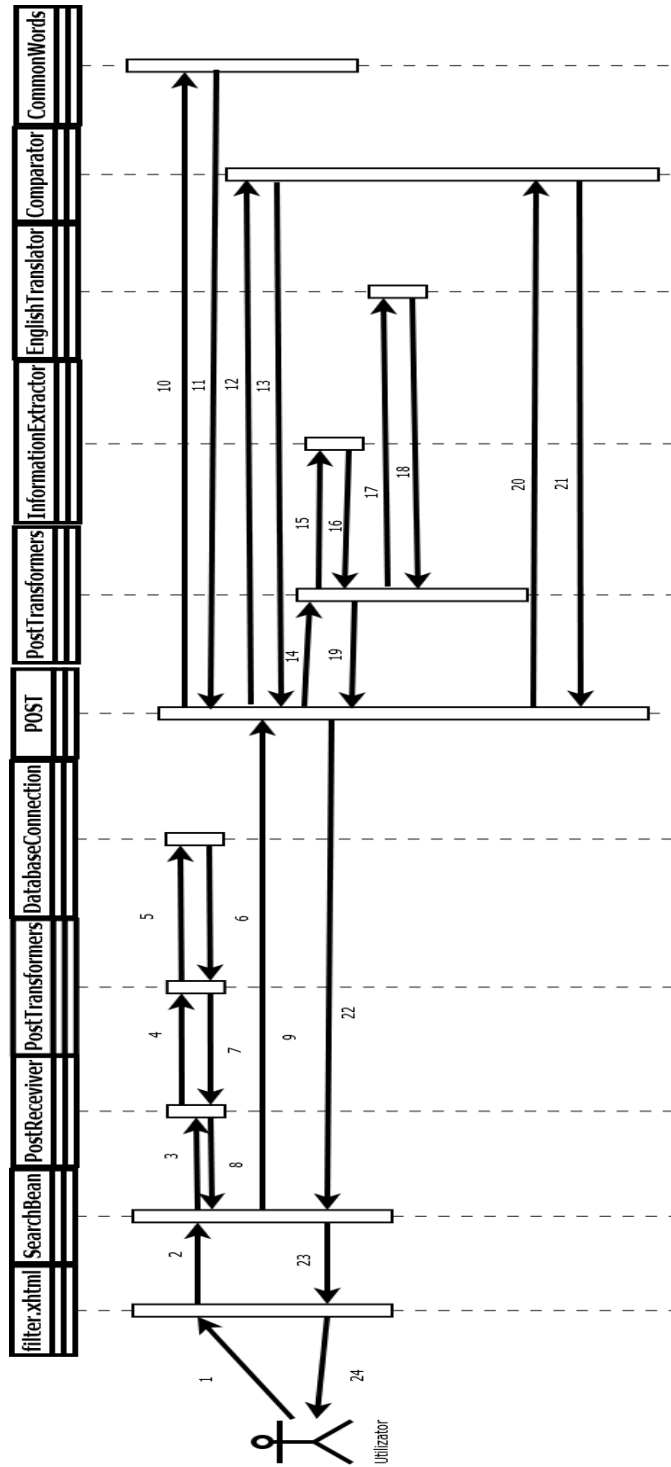


Figura 5.6 Diagrama secvențială

5.7. Tehnologii folosite în aplicație

Tehnologiile folosite în aplicație pot fi grupate în două mari categorii: tehnologii de front-end și tehnologii de back-end.

5.7.1. Tehnologiile de front-end

Când discutăm despre *front-end*, vorbim despre acea parte a site-ului sau a aplicației web, pe care o putem vedea și cu care interacționează vizitatorii. Front-End-ul are două părți: design-ul (partea creativă) și dezvoltarea interfeței (partea de cod sau implementare HTML, CSS).

5.7.1.1. Javascript

JavaScript (JS) este un limbaj de programare orientat obiect bazat pe conceptul prototipurilor. Este folosit mai ales pentru introducerea unor funcționalități în paginile web, codul Javascript din aceste pagini fiind rulat de către browser. Limbajul este binecunoscut pentru folosirea sa în construirea siturilor web, dar este folosit și pentru acesul la obiecte încastate (embedded objects) în alte aplicații.

Cea mai des întâlnită utilizare a JavaScript este în scriptarea paginilor web. Programatorii web pot îngloba în paginile HTML script-uri pentru diverse activități cum ar fi verificarea datelor introduse de utilizatori sau crearea de meniuri și alte efecte animate.

Browserserele rețin în memorie o reprezentare a unei pagini web sub forma unui arbore de obiecte și pun la dispoziție aceste obiecte script-urilor JavaScript, care le pot citi și manipula. Arborele de obiecte poartă numele de Document Object Model sau DOM. Există un standard W3C pentru DOM-ul pe care trebuie să îl pună la dispoziție un browser, ceea ce oferă premiza scrierii de script-uri portabile, care să funcționeze pe toate browserele. În practică, însă, standardul W3C pentru DOM este incomplet implementat. Deși tendința browserelor este de a se alinia standardului W3C, unele din acestea încă prezintă incompatibilități majore, cum este cazul Internet Explorer. [6]

În proiectul nostru, Javascript are rolul principal de a nu lăsa utilizatorul să caute un mesaj gol. Dacă utilizatorul va încerca să introducă un mesaj gol acesta va fi avertizat asupra faptului că mesajul este gol și că nu poate fi procesată o filtrare asupra unui mesaj gol.

5.7.1.2. CSS

CSS (Cascading Style Sheets) este un limbaj care permite programatorilor să definească aspectul unei pagini web separând conținutul paginii de design-ul acesteia. Efectul de cascadă, după cum îi spune și numele, este ceea ce face ca acest limbaj să fie unic. Acest efect reprezintă parcurgerea fișierelor de la început până la sfârșit extrăgând informațiile necesare. De exemplu dacă pentru un anumit titlu este definită o culoare de două sau mai multe ori, ultima definiție va fi luată în considerare.

Includerea unor stiluri definite în limbajul CSS într-o pagină JSF este una foarte simplă și se poate realiza în 3 moduri:

- Se creează un fișier separat cu stilurile definite după care acesta se include în fișierul JSF.
- Se scriu blocuri de stiluri direct în fișierul JSF.
- Se scriu stiluri pentru fiecare element în parte acestea fiind aplicate doar acelu element JSF.

Deși au apărut unele deficiențe de compatibilitate între browsere, majoritatea proprietăților CSS3 au fost implementate cu succes în variantele browserelor noi.

Elementul HTML `div` este definit de următoarele proprietăți CSS: dimensiunea în lungime este redată de valoarea în pixeli a proprietății `width`, folosește o bordură de 2 pixeli, o bordură solidă de culoare gri-închis definită de culoarea HEX `#333333`. Culoarea de fundal este gri deschis definită de HEX `#dddddd`. Bordura rotunjită este de 25 pixeli pentru toate cele 4 colțuri.

În proiectul nostru, `css`-urile au fost folosite pentru a oferi interfeței grafice un aspect cât mai plăcut. Imaginile și stilurile folosite în aplicație sunt de tip `open source`. Având în vedere faptul că aplicația este una care se ocupă cu filtrarea mesajelor s-a încercat adoptarea unui stil care să simuleze faptul că se folosește un stil didactic. Astfel, inputul a fost amplasat astfel încât să dea impresia utilizatorului că se scrie pe o tablă. De asemenea, tot în acest scop, culoarea font-ului a fost schimbată la alb pentru a simula scrierea cu cretă pe tablă. Butonului i-a fost adăugat un stil aparte care îl face mai interesant și mai ușor de integrat în aplicația noastră. Secțiunea de `div` a fost introdusă deoarece `css`-urile din aplicația mea folosesc într-un procent foarte mare `div`-uri pentru a importa stiluri aplicației.

5.7.2. Tehnologiile de back-end

Back-End-ul de obicei constă în trei părți: un server, o aplicație de interfață și o bază de date. Rolul principal al acestuia este cel de management de conținut.

5.7.2.1. Java

Java [7] este un limbaj de programare de nivel înalt, dezvoltat de JavaSoft, companie în cadrul firmei Sun Microsystems.

Limbajul împrumută o mare parte din sintaxă de la C și C++, dar are un model al obiectelor mai simplu și prezintă mai puține facilități de nivel jos. Un program Java compilat, corect scris, poate fi rulat fără modificări pe orice platformă care e instalată o mașină virtuală Java (engleză Java Virtual Machine, prescurtat JVM). Acest nivel de portabilitate (inexistent pentru limbaje mai vechi cum ar fi C) este posibil deoarece sursele Java sunt compilate într-un format standard numit cod de octeți (engleză `byte-code`) care este intermediar între codul mașină (dependent de tipul calculatorului) și codul sursă.

Dintre caracteristicile principale care m-au făcut să aleg limbajul Java ca și limbaj de bază în folosirea aplicației amintim:

- **simplitate**, elimină supraîncărcarea operatorilor, moștenirea multiplă și toate "facilitățile" ce pot provoca scrierea unui cod confuz.

- **robustețe**, elimină sursele frecvente de erori ce apar în programare prin eliminarea pointerilor, administrarea automată a memoriei și eliminarea fisurilor de memorie printr-o procedură de colectare a 'gunoiului' care rulează în fundal
- **complet orientat pe obiecte** - elimină complet stilul de programare procedural
- **usurința în ceea ce privește programarea în rețea**
- **securitate**, este cel mai sigur limbaj de programare disponibil în acest moment, asigurând mecanisme stricte de securitate a programelor concretizate prin: verificarea dinamică a codului pentru detectarea secvențelor periculoase, impunerea unor reguli stricte pentru rularea programelor lansate pe calculatoare aflate la distanță, etc

Java este un limbaj modern, din el putând deriva alte tehnologii și metode folosite în aplicație. Framework-urile și tehnologiile care au la bază Java și care au fost folosite în cadrul aplicației noastre vor fi descrise în secțiunile următoare..

5.7.2.1.1. JSF

Java Server Faces[13] este un framework folosit pentru dezvoltarea de aplicații web. Acesta permite crearea de UI folosind componente standard, reutilizabile. Sintaxa lui se aseamănă mult cu sintaxa HTML însă oferă și anumite facilități în plus. O altă caracteristică importantă a paginilor JSF este aceea că asigură persistența stării la nivel de WEB (permite crearea de UI folosind componente, salvează starea componentelor UI când clientul face o solicitare pentru o nouă pagină și o restaurează când solicitarea este returnată)

Motivul pentru care acest framework este util pentru aplicația mea este acela că, paginile JSF au o logică simplă, sunt ușor de folosit și întreținut. Pentru un algoritm simplu de filtrare acestea se mulează exact pe ce avem nevoie. De asemenea, framework-ul JSF oferă o organizare proprie, care face codul mai ușor de înțeles și urmărit.

5.7.2.1.2. Arhitectura Layers

Stilul Layers [14] (pe niveluri): adecvat aplicațiilor care pot fi descompuse în grupuri de subtaskuri, fiecare grup reprezentând un anumit nivel de abstractizare. Fiecare nivel utilizează servicii furnizate de nivelul(urile) de sub el. Arhitectura folosește un stil simplist, robust, sigur și de aceea se potrivește cerințelor aplicației mele

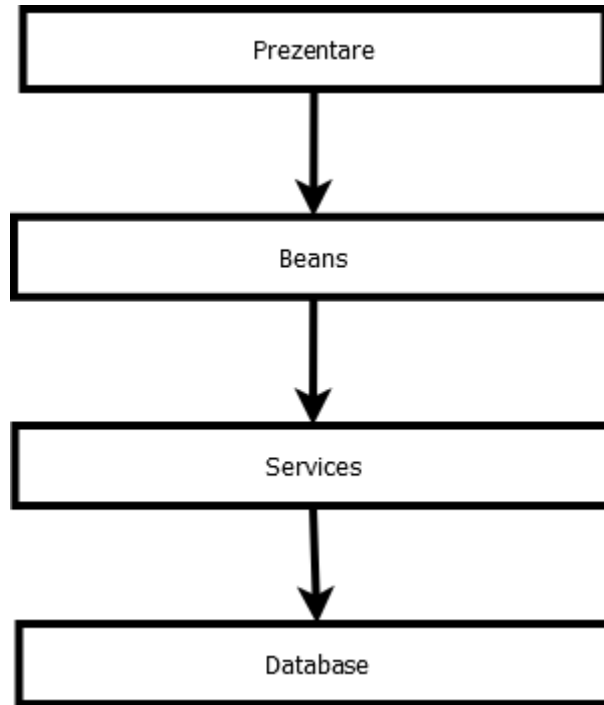


Figura 5.7 Modul de mapare al arhitecturii layers pe aplicație

5.7.2.1.3. Singleton

Design patterns (sau proiectarea modelelor) sunt un set de soluții ce pot fi aplicate la diverse probleme ce se întâlnesc frecvent în cadrul programării orientate pe obiect.

Modelul Singleton [15] este probabil cel mai simplu model. El este făcut pentru a permite accesul la o singură resursă care nu va fi niciodată duplicată, dar care trebuie făcută disponibilă în orice moment al execuției aplicației. Implementarea unui Singleton trebuie să satisfacă cele două condiții de acces global și instanțiere singulară a resursei respective. Are nevoie de un mecanism de acces la clasa Singleton fără a instanția un obiect și un mecanism de a păstra informația asupra resursei ce urmează a fi folosită. Astfel acest model poate fi cel mai ușor implementat prin crearea unei clase ce va conține o metodă care va crea o nouă instanță a clasei, dacă aceasta nu există deja, iar dacă există o va returna pe aceasta. Pentru a fi siguri că acea clasă nu va fi instanțiată în alt fel, constructorul este declarat protejat. Singleton este folosit în aplicația noastră pentru a instanția clasa DatabaseConnection (nu dorim ca să facem mai multe conexiuni la baza de date în cadrul aceleiași aplicații).

5.7.2.1.4. Bibliotecia JSOUP

Bibliotecia JSOUP [16] este o bibliotecă Java care se ocupă în principal cu parsarea de String-uri. Motivul pentru care s-a folosit bibliotecia JSOUP în aplicația noastră este pentru a parsa un site web și a scoate informațiile utile din interiorul său.

Pentru ca metoda să funcționeze se presupune că avem un url sub forma unui String și că dorim să extragem informațiile utile din acel url. String-ul poate proveni

dintr-un fișier de intrare sau dintr-un alt site web. Bibliotecă va oferi suport pentru toate aceste tipuri de proveniență.

Metoda care parsează string-ul și returnează conținutul url-ului este `Jsoup.parse(String html)` sau `Jsoup.parse(String html, String baseUrl)` dacă pagina provine de pe un alt site și se dorește să se ia url-ul absolut

Metoda `parse(String html, String baseUrl)` parsează url-ul într-un document nou. Argumentul `baseUrl` este folosit pentru a transforma url-ul relativ într-unul absolut. Dacă se cunoaște cu siguranță faptul că url-ul este o cale absolută atunci se va putea folosi metoda `parse(String html)`.

Atâta timp cât se va parsea un string care nu este null, este garantat faptul că vom avea o parsare efectuată cu succes și sensibilă (cu un document care conține cel puțin un tag de `<head>` și un tag de `<body>`).

Odată ce avem documentul putem accesa data din interiorul documentului folosind metodele definite în clasa `document` cu argumente specifice xml-urilor cum ar fi `Element` și `Node`.

5.7.2.2. *MySQL*

MySQL[17] este un sistem de gestiune a bazelor de date relaționale, produs de compania suedeză MySQL AB și distribuit sub Licența Publică Generală GNU. Este cel mai popular SGBD open-source la ora actuală, fiind o componentă cheie a stivei LAMP (Linux, Apache, MySQL, PHP).

Deși este folosit foarte des împreună cu limbajul de programare PHP, cu MySQL se pot construi aplicații în orice limbaj major. Există multe scheme API disponibile pentru MySQL ce permit scrierea aplicațiilor în numeroase limbaje de programare pentru accesarea bazelor de date MySQL, cum ar fi: C, C++, C#, Java, Perl, PHP, Python, FreeBasic, etc., fiecare dintre acestea folosind un tip specific API. O interfață de tip ODBC denumită `MyODBC` permite altor limbaje de programare ce folosesc această interfață, să interacționeze cu bazele de date MySQL cum ar fi ASP sau Visual Basic. În sprijinul acestor limbaje de programare, unele companii produc componente de tip COM/COM+ sau .NET (pentru Windows) prin intermediul cărora respectivele limbaje să poată folosi acest SGBD mult mai ușor decât prin intermediul sistemului ODBC. Aceste componente pot fi gratuite (ca de exemplu `MyVBQL`) sau comerciale.

Licența GNU GPL nu permite încorporarea MySQL în softuri comerciale; cei care doresc să facă acest lucru pot achiziționa, contra cost, o licență comercială de la compania producătoare, MySQL AB.

MySQL este componentă integrată a platformelor LAMP sau WAMP (Linux/Windows-Apache-MySQL-PHP/Perl/Python). Popularitatea sa ca aplicație web este strâns legată de cea a PHP-ului care este adesea combinat cu MySQL și denumit Duo-ul Dinamic. În multe cărți de specialitate este precizat faptul că MySQL este mult mai ușor de învățat și folosit decât multe din aplicațiile de gestiune a bazelor de date, ca exemplu comanda de ieșire fiind una simplă și evidentă: „exit” sau „quit”.

Pentru a administra bazele de date MySQL se poate folosi modul linie de comandă sau, prin descărcare de pe internet, o interfață grafică: `MySQL Administrator` și `MySQL Query Browser`. Un alt instrument de management al acestor baze de date este aplicația gratuită, scrisă în PHP, `phpMyAdmin`.

MySQL poate fi rulat pe multe dintre platformele software existente: AIX, FreeBSD, GNU/Linux, Mac OS X, NetBSD, Solaris, SunOS, Windows 9x/NT/2000/XP/Vista.

În aplicație, MySQL este folosit pentru a reține datele din aplicația de micro-blogging. Datele din baza de date vor trebui să poată fi extrase atât din interiorul filtrului cât și din interiorul aplicației de aceea sistemul de gestiune al bazei de date trebuie să fie flexibil. O altă caracteristică importantă a MySQL este securitatea. Doar aplicația de micro-blogging și aplicația de filtrare de mesaje pot avea acces la baza de date restul aplicațiilor având acces interzis. Acest lucru se va face folosind un username și o parolă secretă bazei de date.

5.7.3. Alte tehnologii

Deși aceste tehnologii face parte din back-endul aplicației nu pot fi considerate tehnologii specifice de back-end și de aceea vor fi tratate separat. Tehnologiile descrise în acest capitol sunt tehnologiile cele mai moderne pe care aplicația le deține

5.7.3.1. Aplicația WordNet

WordNet[8] este o bază de date lexicală din limba engleză. Substantive, verbe, adjective și adverbe sunt grupate în seturi de sinonime cognitive (synset-uri), fiecare exprimând un concept distinct. Synset-urile sunt interconectate prin intermediul unor relații conceptuale-semantice și lexicale. Rețeaua rezultată de cuvinte și concepte semnificativ legate poate fi navigată cu browser-ul. WordNet este, de asemenea, în mod liber și la dispoziția publicului pentru descărcare.

WordNet este o aplicație care instalează pe calculatorul local o aplicație de sine stătătoare care simulează funcționarea unui dicționar explicativ în limba engleză. Aplicația însă pune la dispoziția utilizatorilor o serie de librării prin care poate fi accesată baza de date a aplicației.

Pentru a accesa baza de date însă avem nevoie să setăm calea până la baza de date. După ce setăm calea, putem folosi metodele din librăriile de date pentru a obține acces la formele de bază și sinonimele cuvintelor.

Prima metodă care va fi apelată este metoda `getSynsets`. Această metodă primește ca și parametru un cuvânt și returnează totalitatea sinonimelor din baza de date ale cuvântului. Această metodă se va aplica pentru ambele cuvinte.

A doua metodă care se va aplica este metoda `getWordForms`. Această metodă interoghează baza de date a dicționarului și returnează forma de bază a cuvântului respectiv. Baza de date a aplicației WordNet este o bază de date relațională, bază de date a cărei calitate importantă este aceea că face query-urile să fie mult mai rapide.

5.7.4. Integrarea cu alte sisteme

Application Programming Interface (prescurtat API) reprezintă denumirea în limba engleză a unei interfețe pentru programarea de aplicații. De obicei este vorba despre interfața dintre programele de aplicație și sistemul de operare, care stabilește în

amănunt modul în care programele de aplicație pot accesa (apela) serviciile sistemului de operare sub care rulează. Pe scurt un API nu este nimic altceva decât o resursă de internet făcută disponibilă și care poate fi accesată din codul tău local.

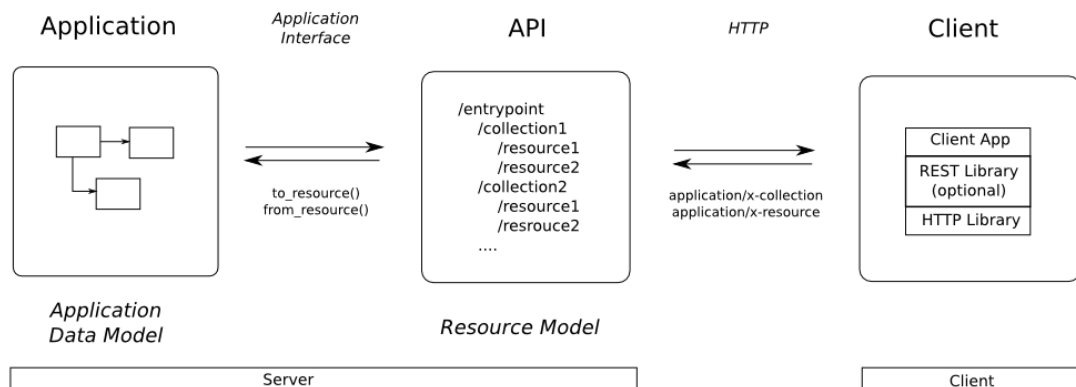


Figura 5.8 Desigul principal al unui API

5.7.4.1. Bing Translator API

Microsoft Translator oferă un set puternic de API-uri de servicii web pe care dezvoltatorii le pot utiliza pentru a profita de tehnologia sa de traducere în propria lor aplicații, servicii sau site-uri web. Acest API poate fi apelat într-un număr de moduri inclusiv un HTTP REST, un serviciu AJAX apelabile și un serviciu SOAP Web.

Traducătorul API Microsoft [9], este un serviciu on-line, și este disponibil direct de la Windows Azure Marketplace. Microsoft Translator API este vândut ca un abonament lunar în funcție de numărul de caractere de text apelate din API. API-ul este disponibil gratuit pentru utilizare de până la 2 milioane de caractere pe lună.

Metodele pe care API-ul le oferă dezvoltatorului sunt enumerate mai jos:

- *Translate* – traduce textul dintr-o limbă în alta
- *TranslateArray* – traduce un șir de texte dintr-o limbă în alta
- *Speak* – Returnează o monștră audio a textului vorbit în limba dorită
- *Detect* – Detectează limba mesajului
- *DetectArray* - Detectează limbiile dintr-o listă de string-uri
- *GetTranslations* - Returnează un vector cu traduceri alternative ale textului
- *GetTranslationsArray* - Returnează un vector cu traduceri alternative la un vector de texte
- *GetLanguageNames*, *GetLanguagesForSpeak*, *GetLanguagesForTranslate* - returnează o listă cu limbile suportate.
- *AddTranslation*, *AddTranslationArray* - permite developerilor să adauge translatările proprii.

În plus față la utilizarea API-ului, developerii pot utiliza Microsoft Translator Widget pentru a obține traduceri instantanee de pagini web fără a fi nevoie de a scrie orice cod. Widget-ul este o mică bucată de cod HTML, care poate fi plasată oriunde pe site, și oferă utilizatorului final, o opțiune de a traduce textul de pe site în oricare dintre limbile acceptate. Widget-ul include caracteristici ale traducerii care vă vor permite să colaborați, să furnizați traduceri alternative, vă permit să editați traducerile automate și votați cea mai bună traducere pentru a îmbunătăți cadrul pentru comunitate.

Capitolul 6. Testare și validare

În cadrul acestui capitol se va prezenta modul de testare a aplicației precum și principalele clase de validare a datelor.

Pe parcursul dezvoltării unui sistem software, testarea este foarte importantă deoarece ajută la prevenirea erorilor și la rezolvarea acestora. De fiecare dată când o funcționalitate nouă este integrată în sistem este necesar să testăm atât acea funcționalitate în particular cât și întreg sistemul după integrarea acesteia pentru a observa modul în care aplicația se comportă. Pe parcursul dezvoltării proiectului au fost testate funcționalitățile și prevenite eventualele erori.

6.1. Validarea datelor

Validarea datelor se face în aplicația noastră cu ajutorul părții de front-end. Astfel, în momentul în care utilizatorul dorește să aleagă ca și mesaj de referință un mesaj gol va fi avertizat asupra faptului că mesajul ales este gol. De asemenea dacă dimensiunea mesajului va depăși limita de 500 de caractere, utilizatorul va fi avertizat din nou asupra faptului că mesajul este prea lung.

Principalul avantaj al validării datelor pe partea de front-end îl consideră faptul că întregul proces nu va funcționa dacă nu are date. Astfel, dacă mesajul nu va respecta cerințele utilizatorul va fi avertizat asupra acestui lucru iar întregul sistem va fi pus în așteptare. Sistemul va ieși din așteptare atunci când utilizatorul va introduce date potrivite. Un alt avantaj al folosirii validării pe partea de front-end ar putea fi considerat și costul redus dar și performanța. Cu doar câteva linii de cod (cost mic) utilizatorul va putea obține o aplicație sigură.

6.2. Testarea aplicației

Testarea aplicației face referire la verificarea componentelor cu scopul ca acestea să îndeplinească cerințele funcționale propuse. În procesul de testare am verificat periodic modul în care se comportă aplicația în diferitele navigatoare pe care le-ar putea folosi utilizatorii (Mozilla, Chrome, Internet Explorer sau Safari).

Cel mai important de testat în cadrul aplicației noastre va fi algoritmul de testare. Acesta, în cazul în care nu va funcționa sau va funcționa greșit, va produce rezultate greșite. Pentru a putea testa cât de bun este algoritmul nostru se vor introduce 9 mesaje care conțin informații despre subiecte diferite. Pentru testarea algoritmului s-a folosit și feed-backul dat de anumiți utilizator. Utilizatorul clasifică mesajele ca fiind relevante sau irelevante în funcție de condițiile sale proprii.

Pentru a putea să ne dăm seama dacă algoritmul nostru este unul bun sau ne trebuie să definim anumite metrici. Aceste metrici sunt:

$$\text{Acuratețe} = \frac{\text{număr relevant corect} + \text{număr irelevant corect}}{\text{numărul total de mesaje}}$$

Figura 6.1 Definirea acurateții

Acuratețea reprezintă raportul dintre numărul de clasificări corecte și numărul total de clasificări. În figura 6.1 am notat cu număr relevant corect ca fiind numărul de mesaje pe care algoritmul nostru l-a clasificat corect iar cu număr irelevant corect am notat numărul de mesaje irelevante pe care algoritmul l-a clasificat corect.

$$\mathbf{Precizie} = \frac{\mathbf{num\bar{a}r\ irelevant\ corect}}{\mathbf{num\bar{a}rul\ total\ de\ mesaje}}$$

Figura 6.2 Definirea preciziei

Precizia măsoară raportul dintre numărul de mesaje relevante clasificate corect și numărul total de mesaje. De notificat este faptul că numărul irelevant de mesaje va fi ignorat.

$$\mathbf{Reapelare} = \frac{\mathbf{num\bar{a}r\ irelevant\ incorect}}{\mathbf{num\bar{a}rul\ total\ de\ mesaje}}$$

Figura 6.3 Definirea reapelării

Reapelarea este opusul preciziei și reprezintă raportul dintre numărul de mesaje irelevante clasificate corect și numărul total de mesaje.

$$\mathbf{FScore} = \frac{\mathbf{2 * Precizie * Reapelare}}{\mathbf{Precizie + Reapelare}}$$

Figura 6.4 Definirea FScore

Scorul F sau FScore reprezintă cea mai complexă metrică de măsură a eficienței algoritmului și reprezintă media armonică dintre precizie și reapelare.

În tabelul 6.5 se pot vedea subiectele de discuție și valorile care au fost luate pentru metricile discutate.

Subiectul Mesajului	Acuratețe	Precizie	Reapelare	FScore
Sarcină-părinți	.855	.895	.833	.899
Film	.792	.89	.709	.789
Ecologie	.792	.787	.947	.86
Market	.878	.896	.977	.935
Programatori	.695	.722	.585	.646
Știri fructe	.791	.844	.882	.863
Actori comedianți	.462	.498	.572	.531
MySql	.694	.684	.57	.621
Bază de date	.882	.923	.873	.897

Tabela 6.5 Eficiența sistemului măsurată în metrici

Capitolul 7. Manual de instalare și utilizare

În acest subcapitol sunt prezentate cerințele pentru instalare și utilizare a aplicație precum și pașii necesari pentru aplicație.

7.1. Instalarea aplicației

După cum se știe deja, pentru ca filtrul să poată funcționa este nevoie de instalarea altor module suplimentare respectiv aplicația de micro-blogging și aplicația de dicționar.

7.1.1. Aplicația de micro-blogging

Micro-bloggingul local este o aplicație de tip open-source. Aceasta se găsește pe site-ul github și poate fi descărcată gratis. Pentru a instala aplicația avem nevoie de mediul virtual Eclipse și sistemul de gestiune a bazelor de date MySQL. De asemenea pentru a putea rula programul avem nevoie de Java Development Kit 1.6 sau o versiune mai nouă. Pentru a se putea crea baza de date a micro-bloggingului trebuie modificat fișierul context.xml cu username-ul și parola bazei de date locale. În acest moment aplicația poate fi rulată cu succes.

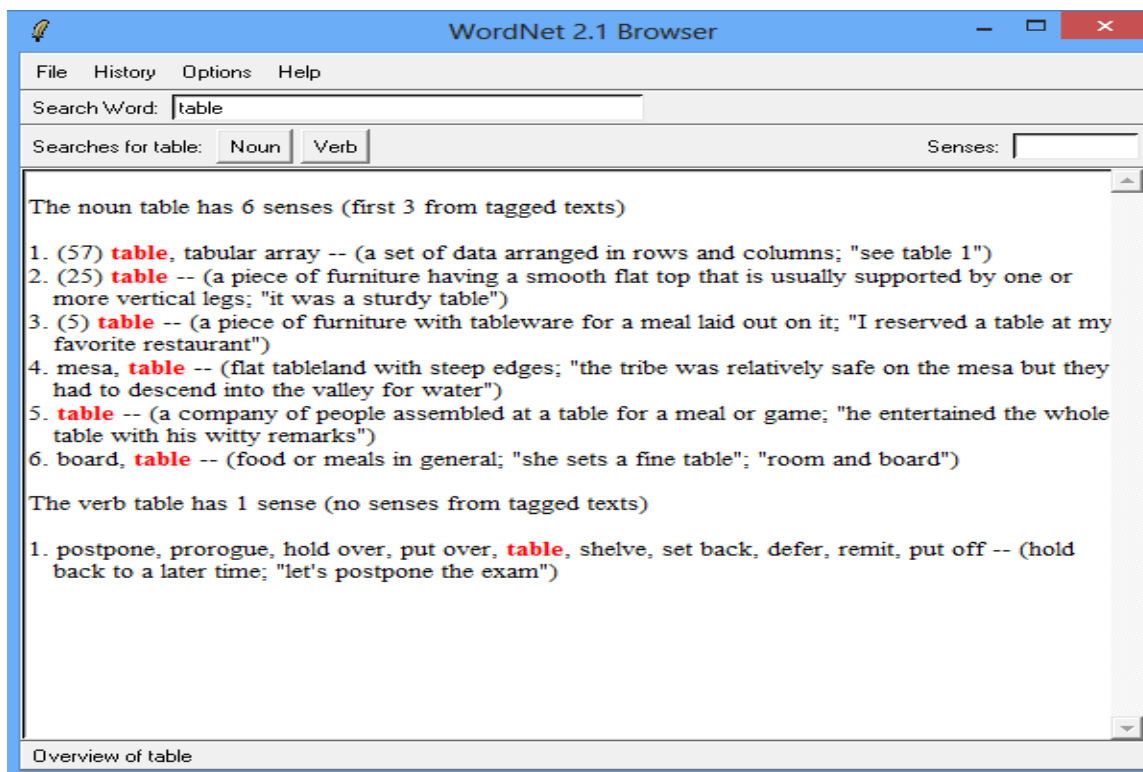


Figura 7.1 Aplicația WordNet

7.1.2. Aplicația de dicționar

Dicționarul nostru local este de fapt aplicația oferită de WordNet pentru dicționar. Aceasta se găsește ca fiind open source și poate fi descărcată de pe site-ul universității din Princeton. Pentru ca aplicația de micro-blogging să poată funcționa trebuie ca dicționarul să fie instalat precum în fișierul de instalare pe care îl conține („C:\Program Files \WordNet”). În figura 7.1 ne este prezentată interfața grafică a dicționarului.

7.1.3. Aplicația de filtrare de mesaje

Pentru ca aplicația de micro-blogging să poată funcționa avem nevoie ca și în cazul aplicației de micro-blogging de Java Development Kit 1.6 sau o versiune mai nouă. De asemenea aplicația a fost creată în mediul de dezvoltare NetBeans. Nu în ultimul rând, pentru a putea accesa și api-ul de traducere oferit de Microsoft vom avea nevoie și de o conexiune la Internet.

7.2. Utilizarea aplicației

Aplicația noastră se ocupă cu filtrarea mesajelor. Pentru a putea filtra mesajele avem nevoie ca mai întâi să avem anumite mesaje în aplicația de micro-blogging. Inserarea mesajelor nu se va putea face decât prin intermediul aplicației de micro-blogging. Pentru a putea insera un mesaj nu este nevoie decât să ne logăm (în cazul în care nu avem un cont, acesta va trebui creat) și să postăm mesajul. În figura 7.2 ne este prezentat modul în care se poate introduce un mesaj în aplicație.



Figura 7.2 Aplicația de micro-blogging

După ce avem mesaje în interiorul micro-bloggingului vom putea să aplicăm și procedura de filtrare de mesaje.



Figura 7.3 Filtrul de micro-blogging

După introducerea în căsuța a textului dorit spre filtrare și apăsarea butonului de search în tabel ne vor apărea mesajele relevante. Primul mesaj este considerat a fi relevant deoarece, după cum se poate observa apar foarte multe cuvinte care sunt asemănătoare ca formă cu textul nostru. În cel de-al doilea mesaj, deși este unul scurt, ponderea cuvântului filtering este una mare și de aceea mesajul nostru va fi considerat ca fiind unul relevant. După cum se poate observa, cel de-al treilea mesaj din figura 7.3 este traducerea în engleză a primului mesaj din figura 7.2 Astfel, sistemul nostru a tradus mesajul, a aplicat algoritmul și a considerat că informația conținută este una relevantă. Ca și concluzie se poate observa faptul că toate traducerile au fost făcute în mod corect.

Capitolul 8. Concluzii

În această secțiune se prezintă rezultatele în concordanță cu obiectivele propuse. Toate obiectivele propuse inițial au fost îndeplinite, astfel:

- Aplicația are abilitatea de a filtra textele provenite de la un serviciu de micro-blogging în texte relevante și texte irelevante
- Aplicația poate distinge în momentul sortării cuvintele uzuale cum ar fi prepozițiile, conjuncțiile de cuvintele folosite pentru sortare
- Algoritm funcționează în mod corect, ponderile fiind atribuite în funcție de frecvența cu care acestea apar în mesaj
- Aplicația permite extragerea conținutului provenit din link-uri, adăugarea acestuia în mesaj și sortarea lui
- Aplicația este capabilă să traducă mesaje provenite din altă limbă și să le adauge la algoritmul de filtrare
- Algoritm de sortare trebuie să fie capabil ca, în momentul în care avem de comparat 2 sinonime compararea sa aibă loc ca și când am avea de comparat 2 cuvinte identice
- Cerințele funcționale, nefuncționale și de implementare au fost respectate cu succes.

8.1. Contribuții personale

Ca și o contribuție am creat un algoritm de clasificare care se bazează pe atribuirea ponderilor cuvintelor din interiorul mesajelor, atribuire bazată pe frecvența mesajelor care apar în text. Am folosit aceste ponderi pentru a sorta mesajele în mesaje relevante și mesaje irelevante. În plus am folosit și sortarea mesajelor provenite din url-uri dar și mesajele provenite din alte limbi. Experimentele au arătat faptul că aplicația noastră este una destul de eficientă.

Am arătat în aplicație și faptul că link-urile incluse în mesaj sunt o sursă bună pentru informarea asupra subiectului respectiv. Am tratat mesajele cu link-uri prin extragerea informației utile din url și aplicarea filtrului de mesaje asupra sa. Experimentele au arătat faptul că eficiența filtrării are o pondere foarte mare.

Cu ajutorul experimentelor din filtrarea de mesaje am ajuns la valori ale acurateții între 75% și 85%. Mai departe am observat că fiecare tip de mesaj din cadrul experimentelor are un anumit topic și un anumit grad de acuratețe (în funcție de cât de specific este fiecare tip de mesaj).

În final, există o versiune online pentru filtrarea de mesaje. Acesta reține ponderea cuvintelor din mesajul filtrat și calculează pentru fiecare mesaj introdus valoarea sa de relevanță.

8.2. Dezvoltări ulterioare

Aplicația de față dispune de mai multe direcții de dezvoltări ulterioare cum ar fi:

- Posibilitatea de a alege din mai multe aplicații de micro-blogging. Utilizatorul va putea să-și aleagă din interfața grafică micro-bloggingul pe care se dorește să se facă filtrarea de mesaje. Dacă nici un micro-blogging nu va fi ales utilizatorul se va deduce faptul că utilizatorul dorește informații de pe toate serviciile de micro-blogging.
- Posibilitatea de a insera în câmpul de căutare un url. Astfel, se vor afișa mesajele relevante cu conținutul acelu url.
- Posibilitatea introducerii în interiorul câmpului a unui mesaj într-o altă limbă diferită de engleză. Utilizatorul ar trebui să primească informația în limba în care a fost scris și textul pentru căutarea informației.
- Adoptarea unor algoritmi și tehnologii care să mărească viteza de execuție a aplicației.
- Introducerea unor modele de tipar(“template-uri”) pentru tipurile uzuale de mesaje folosite

Bibliografie

- [1] „Blog” [Online], <http://ro.wikipedia.org/wiki/Blog>
- [2] „Proceedings of the 16th international conference on World Wide Web” New York, SUA , 2007.
- [3] Douglas W. Ouard, *The state of art in text filtering* , University of Maryland,2009
- [4] Fabrizio Sebastiani,*Machine learning in automated text categorization* , Consiglio Nazionale delle Ricerche,Italia
- [5] „Most Common Words in English”[Online], http://en.wikipedia.org/wiki/Most_common_words_in_English
- [6] „Javascript”, [Online], <http://ro.wikipedia.org/wiki/JavaScript>
- [7] „Java”, [Online], <http://profs.info.uaic.ro/~acf/java/curs/1/introducere.html>
- [8] „Wordnet”, [Online], <http://wordnet.princeton.edu/>
- [9] „Bing Translator”, [Online], <http://www.microsoft.com/en-us/translator/developers.aspx>
- [10] „Twitter”, [Online], <http://ro.wikipedia.org/wiki/Twitter>
- [11] „StatusNet”,[Online], <http://en.wikipedia.org/wiki/StatusNet>
- [12] „Comparasion of micro-blogging services”,[Online], http://en.wikipedia.org/wiki/Comparison_of_microblogging_services
- [13] „JSF”,[Online], http://en.wikipedia.org/wiki/JavaServer_Faces
- [14] Elliot Rusty Harold , *Java Network Programming*,O’Reilly
- [15] Alan Shalloway, James R. Trott, *Design Patterns Explained*, Pearson Education,2009
- [16] „JSOUP”, [Online], <http://jsoup.org/>
- [17] „MySql”, [Online], <http://ro.wikipedia.org/wiki/MySQL>

Anexa 1- Lista de figuri

Figura 2.1 Schema generală a sistemului.....	9
Figura 4.1 Schema modului de extragere de text html.....	19
Figura 4.3 Modul de calcul al mesajului local.....	22
Figura 4.4 Relevanța unei postări.....	23
Figura 4.5 Pseudocod aplicare ponderi.....	24
Figura 4.6 Pseudocod algoritm filtrare.....	25
Figura 4.7 Arhitectura aplicațiilor bazate pe HTTP.....	26
Figura 5.1 Flow-ul de execuție al aplicației.....	28
Figura 5.2 Diagrama conceptuală a aplicației.....	29
Figura 5.3 Diagrama use-case a aplicației.....	32
Figura 5.4 Diagrama de pachete.....	54
Figura 5.5 Diagrama de clase.....	33
Figura 5.6 Diagrama secvențială.....	38
Figura 5.7 Modul de mapare al arhitecturii layers pe aplicație.....	42
Figura 5.8 Desigul principal al unui API.....	45
Figura 6.1 Definirea acurateții.....	47
Figura 6.2 Definirea preciziei.....	48
Figura 6.3 Definirea reapelării.....	48
Figura 6.4 Definirea FScore.....	48
Figura 7.1 Aplicația WordNet.....	50
Figura 7.2 Aplicația de micro-blogging.....	51
Figura 7.3 Filtrul de micro-blogging.....	52

Anexa 2- Lista de tabele

Tabela 3.1 Tipurile de comunicații acceptate de micro-blogginguri.....16
Tabela 4.2 Cele mai întâlnite 40 de cuvinte din vocabularul englez.....21
Tabela 6.5 Eficiența sistemului măsurată în metrici.....49