



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Dezvoltarea platformei de management de cunoaștere



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÂNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRI ȘI INOVĂRI



Creșterea Capacității Administrative
a Sistemului Public de CDI



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Dezvoltarea platformei de management de cunoaștere

Cuprins

I. Dezvoltarea componentei de prevenire timpurie (early warning) a platformei de management de cunoaștere	6
I.1. Colectarea textelor și analiza surselor	6
I.1.1. Analiza surselor	2
I.1.2. Repository de știri	12
I.2. Validarea umana (TAGy)	13
I.2.1. Consolidarea platformei electronice de identificare și clasificare a semnalelor slabe (TAGy)	13
I.2.2. Procesul de evaluare și validare umană a știrilor	14
II. Implementarea componentei de analiză semantică din cadrul platformei de management de cunoaștere	15
II.1. Dezvoltarea scenariului de analiza semantică	15
II.2. Analiza automată	15
Listă Anexe:	40
Listă Figuri:	40



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Platforma de management de cunoaștere își propune să deservească structurile MENCS responsabile în domeniul CDI precum și UEFISCDI, principala instituție cu rol în finanțarea cercetării din România. Mecanismele dezvoltate în proiect pot susține orientarea strategică a organizațiilor cu activitate de cercetare. De asemenea, modulele ale acestei platforme pot contribui la consultări multi-actor pe teme de CDI, precum și la elaborarea politicilor CDI bazate pe evidente. Premisa este că de înțelegerea tehnologiilor emergente – tehnologii cu potențial de adoptare pe scară largă și/sau impact major asupra unuia sau mai multor sectoare economice – depinde capacitatea actorilor din ecosistemul de cercetare, dezvoltare, inovare de a își construi strategii și planuri pe termen lung. Identificarea din timp a tehnologiilor relevante va crește competitivitatea actorilor amintiți și capacitatea lor strategică.

Platforma de management de cunoaștere se constituie ca un „radar” de identificare a tendințelor tehnologice - un sistem inovativ care combină analiști umani și inteligența artificială pentru selectarea dintr-un set larg de stiri a celor care reprezintă „semnale slabe”, respectiv care descriu fenomene emergente (cu posibilitate de amplificare), în special din domeniile tehnologice. Acest demers a presupus dezvoltarea unei baze cu stiri, preluate din peste 300 de surse online, care în acest moment înglobează peste 650.000 de stiri unice din perioada 2010-2015. Procesul de identificare a fenomenelor emergente are două componente: o componentă de validare umană și o componentă automată care folosește algoritmi de inteligență artificială și tehnici de procesare a limbajului natural. Rezultatele celor două componente majore, precum și restul stirilor din repozitoriu au fost integrate într-un modul software de vizualizare multicriterială (**R7.2 Software de vizualizare multicriterială și de tip rețea SVMR**, pe baza unui sistem de filtrare multiplă (<http://radarrepository.uefiscdi.ro>), care permite utilizatorilor accesul online pe baza unui cont cu credențiale unice.

Mecanismul de validare umană a presupus un proces continuu de lectură și încadrare a stirilor din baza de date în categoria „semnalelor slabe” sau a „non-semnalelor”. Acest efort a revenit unei echipe formate din 20 de studenți masteranzi cu specializări diverse, recrutați pe baza abilităților de înțelegere a textelor în limba engleză. Aceștia din urmă au fost organizați în echipe de câte două persoane, schimbându-și compoziția la fiecare flux de lucru săptămânal. Interacțiunea s-a realizat online pe o platformă de evaluare colaborativă a stirilor tehnologice, denumită sugestiv, TAGy (**R7.1 1 platformă operațională privind evaluarea colaborativă a stirilor tehnologice**, <http://tagy.uefiscdi.ro/Account/Login.aspx>), dezvoltată în cadrul proiectului.

Analiza automată presupune identificarea de *patternuri* cu ajutorul unor algoritmi specifici de clasificare de texte: LDA (Latent Dirichlet Analysis) - pentru identificarea de pattern-uri independente de analiza umană și SVM (Support Vector Machine) - pentru predicție de semnale slabe și non-semnale pe baza exemplelor validate manual.

În vederea prelucrării stirilor rezultate în urma evaluării umane, inclusiv cu scopul furnizării de materiale pregătitoare pentru analiza automată, a fost nevoie de construirea unei ontologii proprii pentru semnale slabe. Aceasta activitate a presupus crearea unor rețele semantice care să includă termenii specifici; cu alte



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNĂLȚĂMANTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRII ȘI INOVĂRII



Creșterea Capacității Administrative
a Sistemului Public de CDI



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

cuvinte a unor dicționare de termeni cu structura arborescentă (**R7.5 Un scenariu de analiza semantica**), în programul Tropes Zoom. Aceste dezvoltări permit încadrarea automată a stărilor cu acuratețe pe taxonomii multiple, precum și vizualizarea lor multicriterială.

Un alt modul dezvoltat în cadrul activității este acela de vizualizare a stărilor pe categorii semantice (**R7.6 Un modul de vizualizare a stărilor pe categorii semantice**, <http://greuceanu.uefiscdi.ro/Login.aspx>). Această componentă permite și ierarhizarea unor seturi de stări împărțite pe categorii semantice și a fost folosit cu succes în cadrul unor exerciții de orientare strategică a Institutelor Naționale de Cercetare Dezvoltare (INCD). (Vezi [anexa Platforma sondaje](#))

Pentru accesarea modulelor de software dezvoltate în cadrul proiectului, potențialii utilizatori pot solicita credențiale de acces la adresa info_radar@uefiscdi.ro. În cadrul anexelor acestui raport, pentru 2 dintre platforme (R7.1 și R7.2) au fost create conturi de acces (username și parolă) în vederea vizualizării datelor prezentate.

Platforma de management de cunoaștere la nivelul structurilor MECS responsabile în domeniul CDI și al UEFISCDI a presupus realizarea, testarea și operaționalizarea următoarelor module:

- I. **Dezvoltarea componentei de prevenire timpurie (*early warning*) a platformei de management de cunoaștere**
 - 1.1. Colectarea textelor și analiza surselor
 - Identificare unui set relevant de site-uri de stări preponderent tehnologice din întreaga lume (minim 50)
 - Crearea unei baze de stări (repozitoriu) din aceste surse de stări
 - Analiza surselor
 - 1.2. Validare umană
 - Dezvoltare unei platforme tip *gaming* pentru evaluarea umană a stărilor
 - Evaluarea umană a unui set de stări folosind platforma dezvoltată (minim 10.000 stări evaluate săptămânal)
- II. **Implementarea componentei de analiză semantică din cadrul platformei de management de cunoaștere**
 - 2.1. Dezvoltarea scenariului de analiză semantică
 - 2.2. Dezvoltarea componentei de analiză automată a stărilor

În prezentul raport sunt descrise activitățile desfășurate pe fiecare componentă în parte.



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

I. Dezvoltarea componentei de prevenire timpurie (early warning) a platformei de management de cunoaștere

I.1. Colectarea textelor și analiza surselor

Un număr mare de știri (preponderent tehnologice) sunt preluate constant ca exemple/ resurse de învățare și prelucrate sub diferite formate pentru analize (analiza manuală - TAGy, analize semantice, analiza automată, analiza surselor). Știrile sunt preluate de pe o serie de platforme online anterior identificate dar și actualizate pe parcursul proiectului.

De la începutul anului 2015, procesul de colectare a știrilor s-a automatizat prin colectarea lor din RSS feed-uri. În perioada de raportare au fost colectate prin această metodă 193.431 de știri. Pe tot parcursul proiectului s-au analizat informațiile și conținutul știrilor astfel încât, pe de o parte unele dintre platforme au fost eliminate datorită faptului că nu au generat suficiente resurse pentru identificarea tendințelor emergente iar pe de altă parte, au fost adăugate platforme noi, cu potențial. Lista finală a resurselor online (platforme de știri) este evidențiată în figura de mai jos, sau poate fi vizualizată accesând adresa: (<http://192.168.10.21/RSSfeeder/>) în Fig.1.1 Lista finală platforme RSS feeder.



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÂNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRII ȘI INOVĂRII



Creșterea Capacității Administrative
a Sistemului Public de CDI

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Fig.1.1. Lista finala platforme RSS feeder

Situatie RSS feeder	
Nr crt	Platforma
1	www.firstpost.com
2	www.sciencedaily.com
3	www.technology.org
4	www.forbes.com/technology
5	techcrunch.com
6	www.huffingtonpost.com
7	www.entrepreneur.com
8	www.zdnet.com
9	phys.org
10	www.businessinsider.com/sai
11	www.venturebeat.com
12	www.gizmag.com
13	www.dailymail.co.uk/sciencetech
14	www.cleantechnica.com
15	www.techradar.com
16	www.salon.com
17	www.theguardian.com
18	www.firstpost.com/tech
19	www.biospace.com
20	www.scientificamerican.com
21	www.nanowerk.com
22	www.livescience.com
23	www.theverge.com/tech
24	www.foodnavigator.com
25	news.discovery.com
26	www.popsci.com
27	www.cbc.ca/news/technology
28	www.newscientist.com
29	www.iflscience.com
30	www.usatoday.com/tech/
31	www.biosciencetechnology.com
32	www.reuters.com/news/technology
33	www.azonano.com
34	www.space.com
35	www.electronicweekly.com
36	www.extremetech.com
37	www.securityweek.com
38	www.naturalnews.com
39	www.science20.com
40	www.foxnews.com/tech
41	www.nanotech-now.com
42	blogs.sap.com/innovation
43	www.genengnews.com
44	www.independent.co.uk/life-style/gadgets-and-tech/
45	news.sciencemag.org
46	www.techinasia.com
47	www.environmentalleader.com
48	www.bbc.com/technology
49	newsoffice.mit.edu
50	campustechnology.com
51	www.greentechmedia.com
52	www.xconomy.com/channels
53	www.bbc.com/science_and_environment
54	www.3ders.org
55	www.foodnavigator-asia.com
56	www.npr.org/sections/technology
57	www.news.com.au/technology
58	spectrum.ieee.org
59	www.medgadget.com
60	timesofindia.indiatimes.com/home/science
61	www.3dprintingindustry.com
62	http://abcnews.go.com/Technology
63	www.perfscience.com



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

64	www.latimes.com/science	97	www.psfk.com
65	www.impactlab.net	98	www.inside3dprinting.com
66	www.sci-news.com	99	www.roadtovr.com
67	www.bbc.com/education	100	www.autoweek.com
68	www.sciencealert.com	101	www.army-technology.com
69	www.socialnewsdaily.com	102	www.eschoolnews.com
70	www.goodnewsnetwork.org	103	www.water-technology.net
71	www.edsurge.com	104	www.neurosciencenews.com/neuroscience-topics
72	www.aerospace-technology.com	105	edition.cnn.com/TECH
73	scitechdaily.com	106	www.theengineer.co.uk
74	www.washingtonpost.com/business/technology	107	www.theregister.co.uk/science
75	www.healthcareitnews.com	108	www.biopharma-reporter.com
76	uncovercalifornia.com/business/technology	109	www.business-standard.com/technology-news
77	www.technologyreview.com	110	www.npr.org/sections/research-news
78	www.biospectrumasia.com	111	www.photonics.com
79	robohub.org	112	www.cnn.com/id/10000876
80	smartcitiescouncil.com	113	www.japantoday.com/category/technology
81	www.azosensors.com	114	www.recyclinginternational.com
82	www.autonews.com	115	www.bbc.com/health
83	www.theverge.com/science	116	www.uncovermichigan.com/business/technology
84	www.aerotechnews.com	117	www.designboom.com/technology
85	www.the-scientist.com	118	www.bbc.com/future
86	www.recyclingtoday.com	119	http://www.business-standard.com/category/finance
87	searchenginewatch.com	120	www.insidehpc.com/category/news-analysis/
88	www.sfgate.com/technology	121	ec.europa.eu/information_society
89	www.technewsworld.com	122	www.optics.org
90	http://gadgets.ndtv.com/news	123	www.technode.com
91	profit.ndtv.com/news/banking-finance	124	www.biotech-now.org
92	www.abc.net.au	125	www.springwise.com
93	www.foodproductiondaily.com	126	www.mainenewsonline.com/business/technology
94	www.in-pharmatechnologist.com	127	utsandiego.com/headlines/business/technology
95	www.mnn.com	128	biopharma-asia.com
96	www.renewableenergyworld.com	129	www.todavonline.com/tech



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÎNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRI ȘI INOVĂRI



Creșterea Capacității Administrative
a Sistemului Public de CDI



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

130	www.augmentedrealitytrends.com	163	www.mininginnovationnews.com/category/new-tech
131	www.japantimes.co.jp/tech	164	www.moreinspiration.com
132	www.singularityhub.com	165	www.watertechonline.com/topics/6126-environment
133	www.nordiclifescience.org	166	www.watertechonline.com/topics/14556-Industrial
134	www.nanomagazine.co.uk/category&id=172&Itemid	167	www.alternative-energy-news.info
135	www.bionews.org.uk	168	www.neurogadget.com
136	www.wsj.com/news/technology	169	www.rand.org/topics/terrorism-and-homeland-security
137	defensetech.org	170	www.forumforthefuture.org
138	www.waste-management-world.com/recycling.html	171	www.collective-evolution.com/category/sci-tech
139	www.pbs.org/wgbh/nova/next/	172	www.frontlinedesk.com/science-and-fiction
140	www.watertechonline.com	173	www.itwire.com/itwire-rss/science-news
141	www.biologynews.net/	174	www.voicechronicle.com/category/tech-and-science
142	www.euronews.com/sci-tech/	175	www.telecomstechnews.com
143	www.enn.com	176	www.rand.org/topics/law-and-business
144	www.technabob.com	177	www.techthefuture.com
145	www.postscapes.com/internet-of-things-news	178	www.rand.org/topics/science-and-technology
146	defense-update.com	179	www.rand.org/topics/education-and-the-arts
147	www.chinatechnews.com	180	www.iotnewsnetwork.com
148	www.nsf.gov	181	www.rand.org/topics/health-and-health-care
149	www.shapeways.com/blog	182	www.clickgreen.org.uk/rss/research
150	www.nasa.gov	183	www.economist.com/topics/life-sciences
151	www.asmarterplanet.com	184	www.frontlinedesk.com/technology
152	www.itechfuture.com	185	phidgets.wordpress.com
153	www.european-biotechnology-news.com	186	www.rand.org/topics/cyber-warfare
154	www.economist.com/sections/science-technology	187	www.theguardian.com/society/social-trends
155	www.gigaom.com/channel/science-energy	188	www.augmented.org/blog
156	www.theaustralian.com.au/business/technology	189	www.news-medical.net/?tag=/Biosensor
157	www.economist.com/sections/international	190	www.rand.org/topics/children-and-families
158	www.foodtechnology.co.nz	191	www.usnews.com/topics/subjects/biology
159	www.watertechonline.com/topics/14555-municipal	192	www.techanalyst.co
160	www.wfs.org/category/user-interest-tags/scitech	193	www.greenwisebusiness.co.uk
161	www.watertechonline.com/topics/6122-Drinking-Water		
162	news.stanford.edu/news/socsci/		



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNĂLȚĂMĂNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRII ȘI INOVĂRII



Creșterea Capacității Administrative
a Sistemului Public de CDI



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

194	www.waste-management-world.com/recycling.html
195	www.eurekalert.org/bysubject/mathematics.php
196	www.eurekalert.org/bysubject/education.php
197	www.eurekalert.org/bysubject/oceanography.php
198	www.clickgreen.org.uk/rss/research
199	www.watertechonline.com/topics/6122-Drinking-Water
200	www.cnbc.com/id/44877279
201	www.foodtechnology.co.nz
202	www.gigaom.com/channel/science-energy
203	www.treehugger.com/business
204	www.watertechonline.com/topics/14556-Industrial
205	www.treehugger.com/technology
206	www.collective-evolution.com/category/sci-tech
207	www.nasa.gov
208	physicsworld.com
209	www.european-biotechnology-news.com
210	www.treehugger.com/energy
211	www.techthefuture.com
212	humanoids.io
213	www.iotnewsnetwork.com
214	www.mininginnovationnews.com/category/new-technology
215	www.nytimes.com/business/smallbusiness
216	www.itwire.com/itwire-rss/science-news
217	www.economist.com/topics/life-sciences

217	www.economist.com/topics/life-sciences
218	www.forbes.com/green-tech
219	www.alternative-energy-news.info
220	www.rand.org/topics/terrorism-and-homeland-security
221	www.theguardian.com/society/social-trends
222	www.news-medical.net/?tag=Biosensor
223	www.economist.com/topics/space-technology
224	www.rand.org/topics/law-and-business
225	www.voicechronicle.com/category/tech-and-science
226	www.rand.org/topics/science-and-technology
227	www.frontlinedesk.com/science-and-fiction
228	www.rand.org/topics/health-and-health-care
229	www.rand.org/topics/education-and-the-arts
230	www.usnews.com/topics/subjects/biology
231	www.rand.org/topics/cyber-warfare
232	www.augmented.org/blog
233	www.greenwisebusiness.co.uk
234	phidgets.wordpress.com
235	www.quantamagazine.org/category/biology-2
236	www.frontlinedesk.com/technology
237	www.rand.org/topics/children-and-families
238	www.quantamagazine.org/category/computer-physics-2
239	www.quantamagazine.org/category/computer-science-2
240	www.techanalyst.co
241	www.emdt.co.uk



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÂNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRII ȘI INOVĂRII

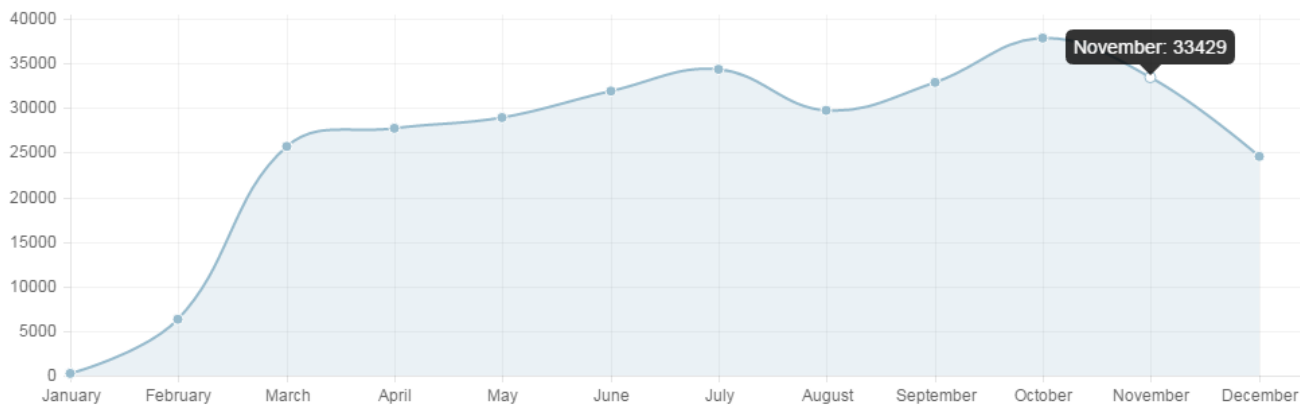


Creșterea Capacității Administrative
a Sistemului Public de CDI

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Fig. 1.2. Distribuția lunară a colectării de știri, 2015

Statistici pe luna



Știrile preluate prin RSS feeder au fost transferate în platforma Repository (această platformă cuprinde toate știrile colectate de-a lungul timpului, inclusiv prin RSS feeder; știrile sunt identificate prin denumirea care conține un cod txt, numele platformei de care aparține știrea, în format txt dar și Excel).

Pentru o selecție mai eficientă a știrilor care sunt supuse validării umane prin platforma Tagy, s-a realizat o ancoră între cele două platforme (Tagy și Repository), astfel automatizându-se mai multe procese care stăteau la baza activităților asociate celor două instrumente (vezi descrierea acestora două în [Anexa 1 Repository Tagy](#)).

I.1.1. Analiza surselor

Analiza și clasificare surse

Pornind de la sursele amintite mai sus s-a realizat o analiză, atât din punct de vedere cantitativ (numărul de știri preluate din RSS feeder și introduse în platforma Tagy, data ultimei extrageri, numărul de intrări în Tagy) dar și calitativ (procentul de semnale slabe pentru fiecare platformă a platformelor de știri (241 de platforme). Rezultatele acestei analize sunt utile pentru clasificarea platformelor de știri în funcție de relevanța informațiilor furnizate (cu cât numărul de SS este mai mare raportat la numărul de știri votate, cu atât este mai relevantă platforma în sine). O previzualizare a clasificării platformelor se poate observa în *figura 2* de mai jos.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Fig. 2. Captura clasificare platforme

Platforma	Data ultimei extrageri din RSS	Nr. Stiri RSS (txt) 02.07	Nr intrari TAGy 15.07	Nr. SS(TAGy)	Nr. NS(TAGY)	% SS TAGy
http://www.3ders.org/	30-Jun-15	528	422	28	394	6.64%
http://www.news-medical.net/?tag=/Biosensor	24-Jun-15	10	7	3	4	42.86%
http://nanomagazine.co.uk/index.php?option=com_sectionex&view=category&id=172&Itemid=158	30-Jun-15	100	58	19	39	32.76%
http://www.moreinspiration.com/article/6144/control-windscreens-wipers-with-your-eyes	22-Jun-15	33	10	3	7	30.00%
http://www.alternative-energy-news.info/	13-Apr-2015	27	4	1	3	25.00%
http://3dprintingindustry.com/	30-Jun-15	464	266	11	255	4.14%
http://www.abc.net.au/science/news/	30-Jun-15	304	236	2	234	0.85%
http://www.aerotechnews.com/news/	30-Jun-15	318	283	4	279	1.41%
http://neurogadget.com/	24-Jun-15	26	21	4	17	19.05%
http://www.augmentedrealitytrends.com/	30-Jun-15	110	138	2	136	1.45%
http://www.azonano.com/	30-Jun-15	925	493	68	425	13.79%
http://www.azosensors.com/	30-Jun-15	338	424	34	390	8.02%

Pentru a scoate în evidență câteva rezultate ale analizei mai sus menționate se poate observa numărul de platforme de știri care au mai mult de 50/100/200/300/500/1000 de stiri precum și numărul de platforme care depășesc procentul de 2% semnale slabe conform procesului de vot din Tagy.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Stiri in Tagy	Numar platforme
>50	148
>100	120
>200	82
>300	58
>500	40
>1000	19

% SS	Numar platforme (cu peste 50 de stiri in Tagy)
>10%	20
>5%	44
>3%	59
>2%	74

% SS	Numar platforme(cu peste 100 de stiri in Tagy)
>10%	18
>5%	39
>3%	54
>2%	66

Analiza similitudinii

O parte din textele colectate sunt analizate cu software-ul de plagiat în vederea identificării de surse adiționale pentru fiecare știre în parte. Sunt astfel generate rapoarte ce prezintă procente ale similarității știrilor cu texte din alte surse de pe internet.

Pentru completarea analizei a fost dezvoltat intern un soft pentru preluarea datei știrilor din toate sursele identificate.

În perioada raportată a fost realizată o analiză a surselor unui eșantion de știri din Repository (din aproximativ 60 de platforme) cu ajutorul software-ului Plagiarism Detector și al software-ului dezvoltat intern de preluare automată a știrilor și transformarea rezultatelor în format excel. Astfel, a rezultat documentul [Anexa 2.Rezultate surse adiționale stiri similitudine](#) . Acesta prezintă aproximativ 700 de știri pentru care au fost identificate câte alte 10 surse pe care se repetă fiecare știre cu minim 50% procent de similaritate. De asemenea, data știrii din fiecare sursă a fost extrasă și este vizibilă în excelul cu rezultate. Din analiză au rezultat și o serie de site-uri de tip agregator, prezentate în figura de mai jos:

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Fig. 3. Agregatoare identificate în analiza surselor

Nr crt	Agregatoare
1	http://www.ideaslaboratory.com/ (AGG)
2	http://able2know.org/topic/226001-132 (AGG)
3	http://campustechnology.com/articles/2014/11/20/stanford-researchers-create-computer-vision-algorithm-for-describing-visual-scenes.aspx (AGG)
4	http://cordis.europa.eu/news/rcn/122134_en.pdf (AGG)
5	http://euroiphonenews.com/ (AGG)
6	http://jqj.umd.edu/news/entropy-nations (AGG)
7	http://lightspeedindia.wordpress.com/ (AGG)
8	http://nanocomputer.com/ (AGG)
9	http://neuronclub.org/2014/01/ (AGG)
10	http://newenergyandfuel.com/ (AGG)
11	http://porqueno.tumblr.com/ (AGG)
12	http://potentiamed.com/category/darpa/ (AGG)
13	http://programacion2z.wordpress.com/ (AGG)
14	http://redmondmag.com/home.aspx?Page=5 (AGG)
15	http://scienceofsingularity.com/ (AGG)
16	http://semimd.com/blog/tag/soitec/ (AGG)
17	http://sheerwind.com/media/releases (AGG)
18	http://southendaquarist.weebly.com/ (AGG)
19	http://vivoforums.com/news/ (AGG)
20	http://wiche.edu/state-highlights/idaho (AGG)
21	http://www.aacn.org/wd/Cetests/media/A142301.pdf (AGG)
22	http://www.ahpanet.com/?page=LatestNews (AGG)
23	http://www.anl.gov/news-room (AGG)
24	http://www.auburnsentinel.com/ (AGG)
25	http://www.bioquicknews.com/node/1505 (AGG)
26	http://www.brijj.com/vk-gupta (AGG)
27	http://www.ehexperts.us/ (AGG)
28	http://www.futuretimeline.net/blog.htm (AGG)
29	http://www.ox.ac.uk/news/science-blog (AGG)
30	http://www.uc.edu/news/NR.aspx?id=20860 (AGG)
31	http://www.umdrightnow.umd.edu/ (AGG)
32	http://www.veooz.com/news/xHWsVg5.html (AGG)
33	http://www.wiche.edu/sara (AGG)

Acest pas va ajuta la următoarele analize de suprapunere și comparare a surselor în vederea:

- Realizării unui top al surselor de știri;
- Identificării de noi radare;
- Identificării sursei primare a știrilor.

Tot în perioada raportată au fost realizate câteva analize cu software-ul de plagiat *Plagiarism Detector* pe un eșantion de știri validate ca semnale slabe în platforma TAGy. 1309 stiri au fost trecute prin software-ul de detectare a similitudinii, astfel au fost generate 1309 rapoarte. O parte din rapoartele rezultate pot fi consultate în [Anexa 3. Rapoarte PD SS](#) (125 rapoarte).

Din cele 1309 rapoarte doar pentru 781 de știri au fost identificate surse cu procent de minim 50%.

De asemenea au fost identificate și agregatoare (o parte dintre acestea pot fi vizualizate mai jos). Aceste surse nu au fost adăugate în RSS.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Fig.4 Agregatoare

1	http://andreswhy.blogspot.com/ (AGG)
2	http://arrgt13.rssing.com/chan-7597928/all_p5.html (AGG)
3	http://berkeley914.rssing.com/chan-28666368/all_p1.html (AGG)
4	http://biodips.com/a/00000000b2d63fb9 (AGG)
5	http://cambridge597.rssing.com/chan-11168856/all_p2.html (AGG)
6	http://cancerlive.net/ (AGG)
7	http://cars.psgghost.net/tag/solar/ (AGG)
8	http://cavendish81.rssing.com/chan-11168859/all_p3.html (AGG)
9	http://chayroot57.rssing.com/chan-27480809/all_p1.html (AGG)
10	http://chrysobalanaceae57.rssing.com/chan-27523191/all_p2.html (AGG)
11	http://corsy14.rssing.com/chan-8100371/all_p3.html (AGG)
12	http://coulterneb14.rssing.com/chan-8101289/all_p6.html (AGG)
13	http://daily-nano-news.blogspot.com/2015/05/microcombing-creates-stronger-more.html (AGG)
14	http://dawdlers4.rssing.com/chan-3600489/all_p6.html (AGG)
15	http://dcircuits.com/?paged=3&m=201505 (AGG)
16	http://discountbook.in/?paged=2 (AGG)
17	http://drumbledore55.rssing.com/chan-24107249/all_p1.html (AGG)
18	http://drunkschmuck.tumblr.com/page/70 (AGG)
19	http://emvco.us/?p=54 (AGG)
20	http://engineering.rice.edu/biomimetic/ (AGG)
21	http://forum.boinaslava.net/showthread.php?13662-%CA%EE%F1%EC%E8%F7%E5%F1%EA%E8-%ED%EE%E2%E8%ED%E8/page!
22	http://forums.spacebattles.com/threads/ibm-sets-new-tape-storage-record.335843/ (AGG)
23	http://genesisananotech.com/solar-cell/ (AGG)
24	http://geneticcenter.com/2013/10/ (AGG)
25	http://godwinaruana.me/tag/research/ (AGG)
26	http://helmet511.rssing.com/chan-33977129/all_p1.html (AGG)
27	http://hulotheseism.rssing.com/chan-1613869/all_p274.html (AGG)
28	http://ieeetesm.org/news/ (AGG)
29	http://industry1090.rssing.com/chan-4098869/all_p10.html (AGG)
30	http://jaunders55.rssing.com/chan-24149330/all_p1.html (AGG)
31	http://lidar82.rssing.com/chan-28258599/all_p1.html (AGG)
32	http://light783.rssing.com/chan-14382296/all_p2.html (AGG)
33	http://ma.webradar.me/portal/40012187 (AGG)



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

34	http://machines158.rssing.com/chan-23604288/all_p2.html (AGG)
35	http://mashable.com/category/dna/ (AGG)
36	http://misplead56.rssing.com/chan-24596183/all_p2.html (AGG)
37	http://mixedcrunch.blogspot.com/ (AGG)
38	http://nanobrainimplant.com/ (AGG)
39	http://naval240.rssing.com/chan-10893106/all_p2.html (AGG)
40	http://news.dxy.cn/bbs/thread/30820918 (AGG)
41	http://news.wisc.edu/releases/19049 (AGG)
42	http://noncustodial56.rssing.com/chan-24737094/all_p1.html (AGG)
43	http://otididae56.rssing.com/chan-24957105/all_p1.html (AGG)
44	http://outboasts56.rssing.com/chan-24937888/all_p2.html (AGG)
45	http://ow.ly/LBgjtj (AGG)
46	http://pests394.rssing.com/chan-28610491/all_p1.html (AGG)
47	http://plug-ugly3.rssing.com/chan-3371028/all_p82.html (AGG)
48	http://plug-ugly3.rssing.com/chan-3371028/all_p85.html (AGG)
49	http://prepayments12.rssing.com/chan-7399440/all_p3.html (AGG)
50	http://princeton675.rssing.com/chan-30037752/all_p3.html (AGG)
51	http://processing284.rssing.com/chan-19070398/latest.php (AGG)
52	http://projectavalon.net/forum4/showthread.php?62082-Technological-advances-that-will-directly-affect-you-in-the-next-2-ye
53	http://pseudofeminine56.rssing.com/chan-25453727/all_p2.html (AGG)
54	http://reawakes56.rssing.com/chan-25545354/all_p2.html (AGG)
55	http://reddylab.org/news/ (AGG)
56	http://research.uconn.edu/2014/09/ (AGG)
57	http://roscosmos1.rssing.com/chan-19077869/all_p1.html (AGG)
58	http://rustiness56.rssing.com/chan-25708429/all_p4.html (AGG)
59	http://safety1713.rssing.com/chan-8246759/all_p1.html (AGG)
60	http://sarahazablog.tumblr.com/ (AGG)
61	http://savings695.rssing.com/chan-7053140/all_p49.html (AGG)
62	http://scenographical53.rssing.com/chan-23736376/all_p1.html (AGG)
63	http://schen.ucsd.edu/lab/news/Kurzweil.pdf (AGG)
64	http://semidivisively56.rssing.com/chan-25857449/all_p1.html (AGG)
65	http://semimd.com/blog/tag/imec/ (AGG)
66	http://sewers71.rssing.com/chan-34134844/all_p1.html (AGG)



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÎNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRII ȘI INOVĂRII



Creșterea Capacității Administrative
a Sistemului Public de CDI



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

<http://www.bioquicknews.com/node/2391> (AGG)
<http://www.bioquicknews.com/node/2473> (AGG)
<http://www.bioquicknews.com/node/2530> (AGG)
<http://www.bioquicknews.com/node/2532> (AGG)
<http://www.bioquicknews.com/node/2553> (AGG)
<http://www.bluesoleil.com/Life/281.html> (AGG)
<http://www.caltech.edu/tags/energy> (AGG)
<http://www.cs.berkeley.edu/> (AGG)
http://www.eduloc.com/hot_news/773.html (AGG)
<http://www.elektrosmog.com/english-1/> (AGG)
<http://www.frogheart.ca/?tag=silicene> (AGG)
<http://www.hearingaidforums.com/showthread.php?t=16385-Cochlear-Implants-with-No-Exterior-Hardware-Its-Here> (AGG)
<http://www.helpcash.gq/microloans.html> (AGG)
<http://www.iaams.ca/aggregator/10.1038/magazine?page=45> (AGG)
<http://www.notey.com/text/blogs/typing> (AGG)
<http://www.taodocs.com/p-10117841.html> (AGG)
<http://www.therxforum.com/showthread.php?t=725297&p=11042860> (AGG)
<http://www.uc.edu/news/NR.aspx?id=20860> (AGG)
<http://www.waseda.jp/top/en-news/25025> (AGG)
<http://www.wfs.org/aggregator/sources/4?page=7> (AGG)
<http://www.xradia.com/press-releases/> (AGG)
<https://facebookaim.wordpress.com/> (AGG)
<https://hisaurav.wordpress.com/> (AGG)
<https://medicupdate.wordpress.com/> (AGG)
<https://plus.google.com/events/c1Inj82cdksanhko39e1g1cg1fk> (AGG)
<https://www.broadinstitute.org/files/news/media-kit/2015/CRISPR-FengZhang-presskit.pdf> (AGG)
<https://www.facebook.com/hot877/posts/10152910257361220:0> (AGG)
<https://www.facebook.com/pages/Inventers-Hub/282412521925675?fref=photo> (AGG)
<https://www.facebook.com/pages/Wonder-Of-The-Science/1511872849028922> (AGG)
https://www.facebook.com/permalink.php?id=472130596221947&story_fbid=628078567293815 (AGG)
<https://www.hzdr.de/db/Cms?pOid=44032> (AGG)
<https://www.pinterest.com/autonomousnow/autonomous-vehides-news/> (AGG)
<https://zedie.wordpress.com/tag/iphone/> (AGG)

Toate rapoartele au fost transpuse într-un excel care conține informația de care aveam nevoie pentru analiză.

O parte din documentul excel poate fi vizualizat în imaginile de mai jos, iar documentul integral poate fi consultat în [Anexa 2.Rezultate surse aditionale stiri similitudine](#).



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÂNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRI ȘI INOVAȚII



Creșterea Capacității Administrative
a Sistemului Public de CDI

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Astfel, utilizatori puteau acorda următoarele scoruri: 1-platformă bună, 2-platformă medie(se poate păstra) sau 3 – platformă propusă spre eliminare.

Pentru fiecare platformă s-a calculat media, iar cele ce aveau scor de minim 2.5 urmau a fi mai atent analizate în fluxurile următoare, pentru a decide eliminarea sau pastrarea lor.

Raport voturi platforme pe liste de stiri

Lista de stiri: Lista 14.08_21.08.2015;
Generat de: Radar Admin

Sumar voturi:

Platforma	Numar voturi 1	Numar voturi 2	Numar voturi 3	Total voturi	Media
defensetech.org	0	1	0	1	2
defense-update.com	0	1	0	1	2
edition.cnn.com/TECH	0	4	2	6	2.33
http://abcnews.go.com/Technology	1	7	2	10	2.1
news.discovery.com	2	9	1	12	1.92
news.sciencemag.org	4	7	0	11	1.64
newsoffice.mit.edu	3	4	0	7	1.57
robohub.org	2	1	0	3	1.33
scitechdaily.com	4	6	0	10	1.6
searchenginewatch.com	0	5	3	8	2.38
smartcitiescouncil.com	0	5	2	7	2.29
spectrum.ieee.org	3	5	0	8	1.62
techcrunch.com	5	6	2	13	1.77
timesofindia.indiatimes.com/home/science	2	7	3	12	2.08
uncovercalifornia.com/business/technology	0	5	2	7	2.29
utsandiego.com/headlines/business/technology	1	2	1	4	2
www.3ders.org	2	7	1	10	1.9
www.3dprintingindustry.com	3	8	1	12	1.83
www.abc.net.au	0	3	4	7	2.57
www.aerospace-technology.com	0	7	1	8	2.12
www.aerotechnews.com	2	8	0	10	1.8
www.army-technology.com	0	3	2	5	2.4
www.asmarterplanet.com	0	2	0	2	2
www.augmented.org/blog	0	1	0	1	2
www.augmentedrealitytrends.com	0	2	3	5	2.6
www.azonano.com	4	4	1	9	1.67

Fig 6.1 Raport vot platforme



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Raport voturi platforme pe liste de stiri

Lista de stiri: Lista 09_16.10.2015;

Generat de: Radar Admin

Sumar voturi:

Platforma	Numar voturi 1	Numar voturi 2	Numar voturi 3	Total voturi	Media
biopharma-asia.com	0	5	2	7	2.29
blogs.sap.com/innovation	0	6	3	9	2.33
campustechnology.com	0	6	3	9	2.33
defensetech.org	0	1	1	2	2.5
edition.cnn.com/TECH	0	6	1	7	2.14
http://abcnews.go.com/Technology	0	6	2	8	2.25
http://gadgets.ndtv.com/news	1	7	3	11	2.18
news.discovery.com	0	9	1	10	2.1
news.sciencemag.org	0	9	0	9	2
news.stanford.edu/news/socsci/	0	1	0	1	2
newsoffice.mit.edu	0	7	2	9	2.22
profit.ndtv.com/news/banking-finance	0	2	8	10	2.8
robohub.org	0	6	1	7	2.14
scitechdaily.com	1	6	0	7	1.86
searchenginewatch.com	0	3	2	5	2.4
smartcitiescouncil.com	0	2	0	2	2
spectrum.ieee.org	0	7	0	7	2
techcrunch.com	0	8	3	11	2.27
timesofindia.indiatimes.com/home/science	1	7	1	9	2
uncovercalifornia.com/business/technology	0	5	1	6	2.17
www.3ders.org	0	7	3	10	2.3
www.3dprintingindustry.com	0	7	2	9	2.22
www.abc.net.au	0	3	3	6	2.5



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÂNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRI ȘI INOVAȚII



Creșterea Capacității Administrative
a Sistemului Public de CDI

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Raport voturi platforme pe liste de stiri

Lista de stiri: Lista 09_16.10.2015;

Generat de: Radar Admin

Sumar voturi:

Platforma	Numar voturi 1	Numar voturi 2	Numar voturi 3	Total voturi	Media
www.business-standard.com/technology-news	0	2	1	3	2.33
www.cbc.ca/news/technology	0	7	2	9	2.22
www.chinatechnews.com	0	1	2	3	2.67
www.cleantechnica.com	1	8	0	9	1.89
www.clickgreen.org.uk/rss/research	0	2	0	2	2
www.dailymail.co.uk/sciencetech	2	6	3	11	2.09
www.designboom.com/technology	1	0	0	1	1
www.economist.com/sections/science-technology	0	3	0	3	2
www.economist.com/topics/life-sciences	0	1	0	1	2
www.edsurge.com	0	5	2	7	2.29
www.electronicweeky.com	0	5	0	5	2
www.environmentalleader.com	0	6	3	9	2.33
www.eschoolnews.com	0	3	1	4	2.25
www.euronews.com/sci-tech/	0	3	0	3	2
www.extremetech.com	0	9	0	9	2
www.firstpost.com/tech	0	6	4	10	2.4
www.foodnavigator-asia.com	0	4	3	7	2.43
www.foodproductiondaily.com	0	3	3	6	2.5
www.forumforthefuture.org	0	1	0	1	2
www.foxnews.com/tech	0	6	4	10	2.4
www.futurity.org	1	9	0	10	1.9
www.genengnews.com	1	8	1	10	2
www.gizmag.com	2	7	1	10	1.9
www.goodnewsnetwork.org	0	3	5	8	2.62
www.greentechmedia.com	0	6	1	7	2.14

Fig 6.2 Raport vot platforme

I.1.2. Repository de știri

Această aplicație/platformă conține toate știrile existente în format excel și txt. *Repository* este disponibilă la adresa: <http://radarrepository.uefiscdi.ro/>.

În perioada raportată au fost realizate următoarele funcționalități:

- Accesul pentru vizualizare pe bază de cont personalizat;
- Corelarea bazei de date Repository cu platforma TAGy în vederea uniformizării informațiilor.



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Toate funcționalitățile acestei platforme în versiunea din perioada raportării și felul în care aceasta se corelează cu platforma TAGy sunt prezentate în [Anexa 1 Repository Tagy](#).

I.2. Validarea umana (TAGy)

1. Consolidarea platformei electronice de identificare și clasificare a semnalelor slabe, în urma feedback-ului constant de la utilizatori/administratori
2. Procesul de evaluare și validare umană a știrilor

I.2.1. Consolidarea platformei electronice de identificare și clasificare a semnalelor slabe (TAGy)

În perioada de raportare s-a finalizat dezvoltarea, actualizarea și operationalizarea platformei electronice de evaluare colaborativă a știrilor tehnologice, **Tagy (R7.1 tagy.uefiscdi.ro)**. Acesta a fost un proces continuu prin care echipa proiectului a integrat sugestiile venite din partea utilizatorilor și a rezolvat problemele tehnice aparute.

Astfel, în vederea eficientizării procesului de validare umană, în cadrul platformei au fost dezvoltate câteva noi funcționalități (vizualizarea noilor aplicații ale platformei este disponibilă în [Anexa 4 Platforma Tagy](#)):

- **Rapoarte**
 - **Raport general NS** – este la baza un raport general care ia în considerare doar non-semnalele, și este folosit în vederea evaluării activității de validare umană; [Anexa 5 Raport general NS](#)
 - **Raport vot platforme** – Raport care scoate în evidență o evaluare prin vot a celor mai bune platforme din punct de vedere al conținutului știrilor. Procesul de votare s-a introdus tocmai pentru a primi un feedback din partea validatorilor cu privire la calitatea știrilor intrate în procesul de evaluare și implicit asupra surselor de proveniență. În funcție de scorurile obținute, echipa proiectului ia în calcul, împreună cu alte criterii, eliminarea sau păstrarea sursei respective pentru fluxurile viitoare. [Anexa 6 Raport vot platforme](#)
- **Integrarea Tagy și Repository**
 - În perioada de raportare s-a finalizat interoperabilitatea dintre cele două platforme; Astfel în platforma Repository se regăsesc toate știrile colectate de-a lungul timpului prin diverse modalități, inclusiv RSS feed (identificate unic printr-un cod txt, numele platformei de unde provine știrea, în format txt dar și xcel); Integrarea presupune ca alocarea știrilor în platforma Tagy se face automat din Repository – astfel ca există o evidență a știrilor care au intrat în fiecare săptămână în procesul de votare. [Anexa 1 Repository Tagy](#)
 - La sfârșitul proiectului, Tagy îndeplinește toate condițiile și înglobează toate funcționalitățile necesare pentru a fi o platformă operațională privind evaluarea colaborativă a știrilor tehnologice



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

I.2.2. Procesul de evaluare și validare umană a știrilor

În perioada iulie 2015 – noiembrie 2015, grupul de cunoaștere tacită inițial (CT7 din echipa proiectului) și-a păstrat același rol de monitorizare și coordonare în procesul de validare umană pe platforma electronică Tagy.

Procesul de evaluare a știrilor a demarat cu alocarea către evaluatorii umani a unui număr de 60 de texte în calitate de extractor și 60 de texte ca checker, pe care fiecare dintre cei 20 de colaboratori au trebuit să le voteze într-un interval de o săptămână. Lista de știri alocate evaluatorilor însumează 1200 de texte selectate din mai multe surse (platforme de știri).

Ulterior, ritmul de lucru a crescut, listele alocate săptămânal conțineau 2400 de texte (120 texte/extractor, 120 texte/checker pentru fiecare persoană), ajungându-se la un maxim de 2500 de texte evaluate săptămânal.

În urma unei evaluări intermediare realizate asupra procesului de validare umană, s-a luat decizia suplimentării numărului de știri alocate evaluatorilor; într-o primă etapă s-a mărit fluxul la 200 știri/checker, 200 știri/extractor iar ulterior s-a ajuns la 300/checker/extractor ceea ce a făcut posibilă evaluarea a cel puțin 4000 știri săptămânal la un număr de 18 evaluatori umani. Până la sfârșitul perioadei de raportare au rămas 13 evaluatori, iar volumul de știri evaluate a scăzut proporțional.

Pe toată perioada derulării proiectului au fost încărcate/ alocate/validate uman – 47 de liste de texte – adică un total de 149 717 de știri evaluate.

Pe tot parcursul acestui flux de lucru, echipa proiectului a monitorizat evoluția voturilor și a modului în care colaboratorii externi au reușit să se integreze și să participe la acest proces.

Monitorizarea s-a concretizat în analize și statistici relevante care scot în evidență informații relevante, atât din punct de vedere cantitativ (numărul de texte evaluate, numărul de semnale slabe sau nonsemnale votate, procentajul acestora în totalul știrilor, procentajul SS pe diferite surse, procentajul în funcție de categoria aleasă a SS, etc) dar și calitativ (tipologia checker/extractor și modul cum aceștia sunt fie autoritari/persuasivi, etc). Aceste analize au fost realizate pe baza seturilor de date de-a lungul desfășurării activității în cadrul proiectului.

Pentru a scoate în evidență câteva rezultate ale acestor statistici vom puncta faptul că din cele **79997** de știri evaluate în perioada de raportare, **2347** au fost votate ca semnale slabe, adică un procent de **3%** în medie. Acest procentaj a scăzut în timp ca urmare a creșterii exigențelor cu care s-au tratat știrile, astfel că dacă în primele sesiuni procentul era undeva la 20%, apoi acesta a scăzut la 16%, iar în prezent pe ultimele seturi de date ajungându-se la 3%, și ca urmare la o medie generală sub 4%. Pentru informații mai detaliate cu privire la statistici și analiza seturilor de date se poate accesa [Anexa 7 Statistici și analize actualizare](#).



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

II. Implementarea componentei de analiză semantică din cadrul platformei de management de cunoaștere

II.1. Dezvoltarea scenariului de analiza semantică

Aceasta activitate a presupus crearea unor rețele semantice care să includă termenii specifici, cu alte cuvinte a unui dicționar cu structura arborescentă (**R7.5 Un scenariu de analiza semantică**), în programul Tropes Zoom. Aceste dezvoltări permit încadrarea automată a stirilor cu acuratețe pe aceste taxonomii multiple, precum și vizualizarea lor multicriterială.

Situația elaborării dicționarilor până la sfârșitul proiectului este detaliată în [Anexa 8. Evidența domeniilor și tendințe](#).

În această perioadă au fost finalizate domeniile care au avut dicționare dezvoltate anterior: *Ageing, Electrical machinery, Environmental challenges, Environmental technology, Health, Medical equipment, Oceanography, Pharmaceuticals, Shipbuilding, Space, Transport, Water industry*.

De asemenea au mai fost dezvoltate dicționare pentru tendințe mai noi și alte 7 domenii: 3d printer, additive manufacturing, artificial intelligence, artificial organ, big data, bioelectronics, biosensor, childcare, cloud computing, crowdsourcing, cybersecurity, driverless, energy storage, entertainment, fashion, fuel cell vehicles, gender, genomics, graphene, internet of things, migration, nuclear fusion, poverty, precision agriculture, quantum computing, recyclable thermoset plastic, renewable energy, retail, smart cities, sport, virtual reality, wearable technology, bionics, distributed manufacturing, precise genetic engineering.

Dicționarele pe domenii dar și cel pentru tendințe realizate în perioada iunie-noiembrie sunt integrate în versiunea Dicționarilor Smart și tendințe din anexele: *Smart* din [Anexa 9 SmartDictionary final](#) și [Anexa 10. Dictionary tendințe](#).

II.2. Analiza automată

În perioada raportată, s-a finalizat dezvoltarea *in-house* a software-urilor necesare dezvoltării componentei de analiză automată și semantică a platformei de management de cunoaștere.

1. Software Dezvoltare Scenariu Semantic (SDSS)

A fost implementat sistemul informatic care permite clasificarea semantică a unui corpus mare de știri într-un număr de clustere la alegerea analistului uman. Pentru fiecare cluster, software-ul furnizează un set de cuvinte cheie, care vor fi incluse în dicționarul semantic. Acest software utilizează metoda de procesare a limbajului natural LDA (Latent Dirichlet Allocation).



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Latent Dirichlet Analysis pentru clusterizarea nesupervizata a fluxului de stiri

Am implementat un proces generativ de clusterizare nesupervizata a fluxului de stiri pentru determinarea automata a similaritatilor computabile in corpus.

Intuitia principala a acestei tehnici este ca putem asocia in mod automat fiecarui cuvint o probabilitate de a semnala un tip de similaritate manifestata in corpus. Ceea ce se obtine este un vector de probabilitati pentru fiecare cuvint, dimensiunea vectorului fiind determinate de numarul de tipuri de similaritati considerate.

Inputul este consituit de doua variabile:

- Numarul de tipuri de similaritati dorite
- Corpusul de analizat

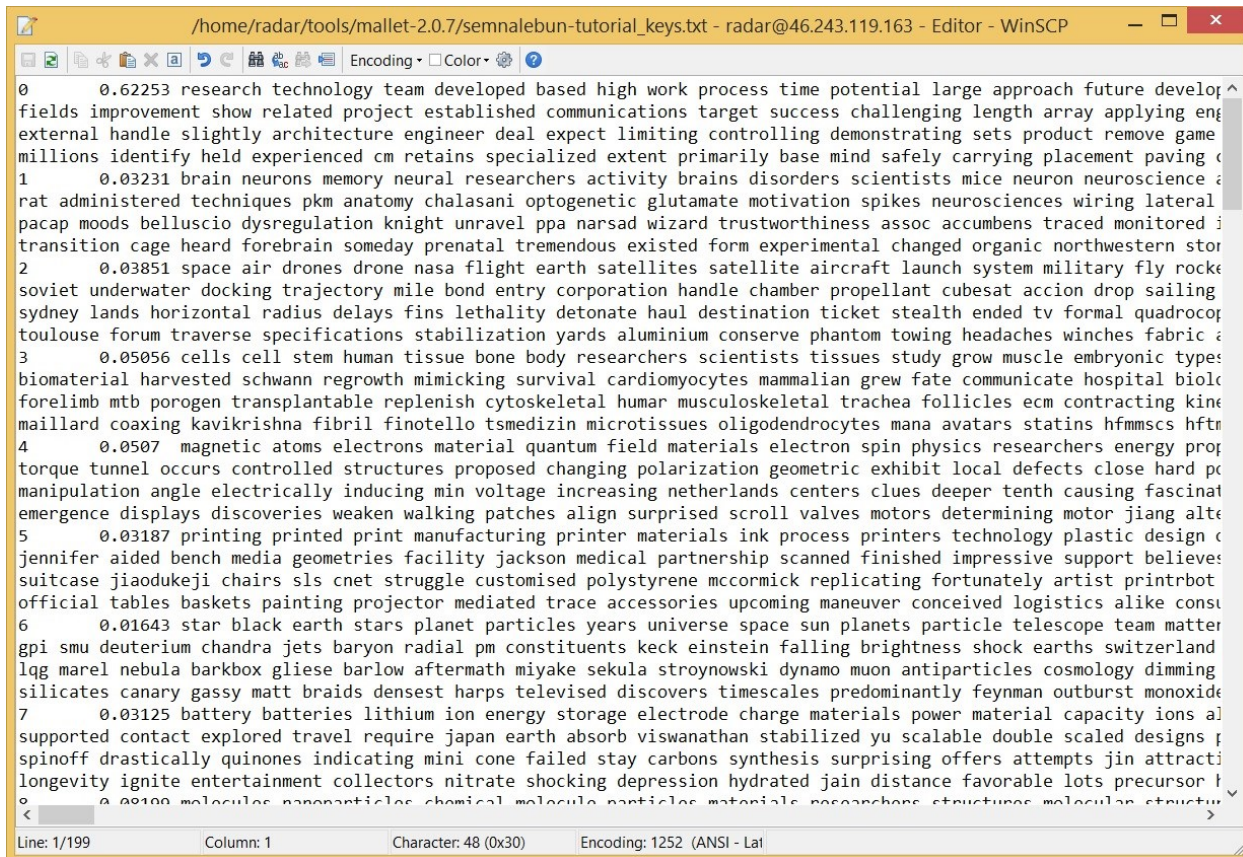
In perioada de raportare am folosit 50 topice(aceasta repartizare a fost facuta pentru imbunatatirea predictiei de incadrare a stirilor) si *recurrent neural network* pentru definirea similaritatii la nivel de cuvinte. Aceasta clusterizare paralela a cuvintelor, independent de topice, este cuplata la inputul LDA pentru clusterizarea corpusului intr-o maniera care reduce fenomenul *data sparseness*. Rezultatul este o acoperire a corpusului mult mai mare, pentru aceeasi valoare a parametrului alfa.

Se maximizeaza probabilitatea de ocurenta a unui cuvint tinind cont de contextul imediat inconjurator:

Corpusul nostru a fost format din peste 5000 de stiri text iar outputul a fost format din:

- Un set de cuvinte cheie pentru fiecare topic

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013



```
0 0.62253 research technology team developed based high work process time potential large approach future develop
fields improvement show related project established communications target success challenging length array applying eng
external handle slightly architecture engineer deal expect limiting controlling demonstrating sets product remove game
millions identify held experienced cm retains specialized extent primarily base mind safely carrying placement paving c
1 0.03231 brain neurons memory neural researchers activity brains disorders scientists mice neuron neuroscience a
rat administered techniques pkm anatomy chalasani optogenetic glutamate motivation spikes neurosciences wiring lateral
pacap moods belluscio dysregulation knight unravel ppa narsad wizard trustworthiness assoc accumbens traced monitored i
transition cage heard forebrain someday prenatal tremendous existed form experimental changed organic northwestern stor
2 0.03851 space air drones drone nasa flight earth satellites satellite aircraft launch system military fly rocke
soviet underwater docking trajectory mile bond entry corporation handle chamber propellant cubesat accion drop sailing
sydney lands horizontal radius delays fins lethality detonate haul destination ticket stealth ended tv formal quadrocop
toulouse forum traverse specifications stabilization yards aluminium conserve phantom towing headaches winches fabric a
3 0.05056 cells cell stem human tissue bone body researchers scientists tissues study grow muscle embryonic types
biomaterial harvested schwann regrowth mimicking survival cardiomyocytes mammalian grew fate communicate hospital biolo
forelimb mtb porogen transplantable replenish cytoskeletal humar musculoskeletal trachea follicles ecm contracting kine
maillard coaxing kavikrishna fibril finotello tsmedizin microtissues oligodendrocytes mana avatars statins hfmmcs hftr
4 0.0507 magnetic atoms electrons material quantum field materials electron spin physics researchers energy prop
torque tunnel occurs controlled structures proposed changing polarization geometric exhibit local defects close hard pe
manipulation angle electrically inducing min voltage increasing netherlands centers clues deeper tenth causing fascinat
emergence displays discoveries weaken walking patches align surprised scroll valves motors determining motor jiang alte
5 0.03187 printing printed print manufacturing printer materials ink process printers technology plastic design c
jennifer aided bench media geometries facility jackson medical partnership scanned finished impressive support believe
suitcase jiaodukeji chairs sls cnet struggle customised polystyrene mccormick replicating fortunately artist printrobot
official tables baskets painting projector mediated trace accessories upcoming maneuver conceived logistics alike consu
6 0.01643 star black earth stars planet particles years universe space sun planets particle telescope team matter
gpi smu deuterium chandra jets baryon radial pm constituents keck einstein falling brightness shock earths switzerland
lqg marel nebula barkbox gliese barlow aftermath miyake sekula stroynowski dynamo muon antiparticles cosmology dimming
silicates canary gassy matt braids densest harps televised discovers timescales predominantly feynman outburst monoxide
7 0.03125 battery batteries lithium ion energy storage electrode charge materials power material capacity ions al
supported contact explored travel require japan earth absorb viswanathan stabilized yu scalable double scaled designs p
spinoff drastically quinones indicating mini cone failed stay carbons synthesis surprising offers attempts jin attracti
longevity ignite entertainment collectors nitrate shocking depression hydrated jain distance favorable lots precursor f
< 0.00100 molecules nanoparticles chemical molecule particles materials researchers structures molecular structur
>
```

Fig. 7. Cuvintele cheie pentru fiecare cluser

Primul este numarul topicului (numerotarea incepe de la 0), iar al 2 numar reprezinta importanta acestuia.

- Repartizarea documentelor pe topic:



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

#doc	name	topic	proportion	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
10	file:/home/radar/tools/semnale/R_www.sciencealert.com%202015%2000023.txt	6	0.533364	46	0.206657	33	0.089925	14	0.042637	40	0.039271	8	0.066195	1	0.050088	20	0.026278	45	0.018033	23	0.019779

Fig. 8. Repartizarea documentelor pe topic

Fisierul 10 are topicul nr 6 ca topic principal cu 53%, topicul 46 cu 20%

2. Software pentru Adnotarea Textelor (SAT)

A fost dezvoltat un sistem de adnotare automată a știrilor pe domenii descrise printr-un corpus de texte. Acest software utilizează tehnica SVM (Support Vector Machine) și permite ca, pornind de la un set de texte exemplu dintr-un domeniu (ex. foraj marin), să selecteze din corpusul de știri pe cele care se încadrează în acest domeniu (știrile sunt ordonate în funcție de probabilitatea de încadrare). Dezvoltarea software-ului SAT are la baza "Data mining" – procesul de analiza a unor cantități mari de date și de extragere a informațiilor relevante. Este un exercitiu interdisciplinar, în realizarea cărui, statistica, tehnologia bazelor de date, recunoașterea de tipare, inteligența artificială și vizualizarea își au rolul lor.

Pentru asta am folosit WEKA - o colecție de instrumente de vizualizare și algoritmi pentru analiza datelor și modelarea predictivă. Pentru a aplica acești algoritmi pe corpusul nostru de știri va trebui să procesăm datele.

Primul pas este să transformăm colecția de fișiere .txt într-un format pe care weka îl poate înțelege. Textul va fi transformat în atribute cu ajutorul clasei StringToWordVector, folosind algoritmi de tokenizing și stemming.





Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Se va face o filtrare a atributelor folosind clasa Attribute Selection, pentru a le clasa în funcție de evaluările lor individuale (Ranker)

Următorul pas este alegerea unui clasificator, pentru a crea modelul de clasificare automată a stărilor. Support Vector Machine este o tehnică de clasificare bazată pe teoria învățării statistice care a fost aplicată cu mult succes în multe probleme neliniare de clasificare și pentru mulți mulțimi foarte mari de date. Ideea algoritmului SVM este de a găsi un hiperplan care împarte optim setul de date de antrenament. Hiperplanul optim se poate distinge prin marginea de separare maximă dintre toate punctele de antrenare și hiperplan. Pentru o problemă într-un spațiu bidimensional algoritmul caută după o dreaptă care separă „cel mai bine” punctele din clasa pozitivă de punctele din clasa negativă. Hiperplanul este caracterizat printr-o funcție de decizie de forma:

$$f(x) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$$

unde w este vectorul pondere, perpendicular pe hiperplan, „ b ” este un scalar care reprezintă marginea hiperplanului, „ x ” este esantionul curent testat și \cdot reprezintă produsul scalar. Sgn este funcția semn care întoarce 1 dacă valoarea obținută este mai mare sau egală cu 0 și -1 altfel. Dacă w este de lungime 1, atunci $\langle w, x \rangle$ este de lungimea lui x de-a lungul direcției lui w . În general w va fi scalat prin $\|w\|$.

În partea de antrenare algoritmul trebuie să găsească vectorul normal „ w ” care conduce la cea mai mare margine „ b ” a hiperplanului. Problema pare foarte ușor de rezolvat dar trebuie să reamintim că linia optimă de clasificare trebuie să clasifice corect toate elementele generate cu aceeași distribuție. Există o multitudine de hiperplane care îndeplinesc cerințele clasificării, dar algoritmul încearcă să determine hiperplanul optim. Acest algoritm de învățare are loc într-un spațiu în care există produs scalar iar pentru datele liniar separabile construiește funcția f din date empirice. Se bazează pe două idei principale: dintre toate hiperplanele de separare există un hiperplan optim unic care se distinge prin marginea maximă de separare dintre orice punct de antrenare și hiperplan. A doua idee este: capacitatea hiperplanului de a separa clasele descrește o dată cu creșterea marginii. Pentru datele de antrenament care nu sunt liniar separabile printr-un hiperplan de separare în spațiul de intrare ideea algoritmului SVM este de a proiecta datele de intrare într-un spațiu de dimensiune mai mare prin intermediul unei funcții Φ , și să încerce să găsească acolo un hiperplan de separare cu marginea maximă. Aceasta conduce la o graniță de separare neliniară în spațiul inițial de intrare. Datorită faptului că toți vectorii apar în produse scalare și prin utilizarea funcției nucleu $\phi(w), \phi(x)$, este posibil să calculăm hiperplanul de separare fără a proiecta explicit datele în noul spațiu de trasaturi.

Pentru a găsi hiperplanul de separare optim care se distinge prin marginea maximă, trebuie să rezolvăm următoarea funcție obiectiv:

$$\begin{aligned} \text{minimize } \tau(\mathbf{w}) &= \frac{1}{2} \|\mathbf{w}\|^2 \\ \mathbf{w} \in H, b \in \mathbb{R} \\ \text{subject to } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) &\geq 1 \text{ for all } i = 1, \dots, m \end{aligned}$$

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Constrângerile asigură faptul că $f(x_i)$ va fi +1 pentru $y_i=+1$ și -1 pentru $y_i=-1$. Această problemă este atractivă din punct de vedere al calculului deoarece poate fi abordată prin rezolvarea unei probleme de programare patratică pentru care există algoritmi eficienți. Funcția τ se numește funcția obiectiv cu inegalități de constrângeri. Acestea formează așa numita problemă de optimizare primară. Pentru rezolvarea acestei probleme este mult mai convenabil să lucrăm cu problema duală prin introducerea multiplicatorilor Lagrange $\alpha_i \geq 0$ și a Lagrangianului care conduce la problema duală de optimizare:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1)$$

Cu multiplicatorii Lagrange $\alpha_i \geq 0$. Observăm că restricțiile sunt incluse în partea a doua a Lagrangianului și nu mai trebuie să fie aplicate separat. Lagrangianul L trebuie să fie maximizat în raport cu variabilele duale α_i , și minimizat în raport cu variabilele primare w și b . Aceasta conduce la:

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

Vectorul soluție este o extensie în termeni de exemple de antrenament. Observăm de altfel că soluția w este unică (datorită convexității stricte pentru problema primară de optimizare), coeficienții α_i , nu trebuie să fie unici. În acord cu teorema Karush-Kuhn-Tucker (KKT) doar multiplicatorii Lagrange α_i care sunt diferiți de zero sunt puncte de sprijin, corespunzătoare constrângerilor care se întâlnesc. Formal, pentru toți $i=1, \dots, m$ acesta poate fi scris:

$$\alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1] = 0 \text{ for all } i = 1, \dots, m$$

Exemplele x_i pentru care $\alpha_i > 0$ se numesc vectori suport. Această terminologie se referă la termeni corespondenți din teoria convexității. În acord cu condițiile KKT, aceștia vor sta exact pe margine. Toate celelalte exemple de antrenament rămase sunt nerelevante. Prin eliminarea variabilelor primare w și b din Lagrangian obținem așa numita problemă de optimizare duală, care este problema care se rezolvă de obicei în practică

$$\text{maximize}_{\alpha \in \mathbb{R}^m} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

care se mai numeste si functia tinta, cu urmatoarele restrictii:

$$\text{subject to } \alpha_i \geq 0 \text{ for all } i = 1, \dots, m \text{ and } \sum_{i=1}^m \alpha_i y_i = 0$$

Astfel hiperplanul poate fi scris, in problema de optimizare duala ca:

$$f(x) = \text{sgn} \left(\sum_{i=1}^m y_i \alpha_i \langle x_i, x \rangle + b \right)$$

unde b este calculat utilizând condițiile KKT. Structura problemei de optimizare este foarte similară cu cea aparută în formularea Lagrange mecanică. În rezolvarea problemei duală apare frecvent cazul în care doar o submulțime de restricții devine activă. De exemplu, dacă vrem să păstrăm o bilă într-o cutie atunci aceasta în mod uzual se va rostogoli într-un colț. Restricțiile corespund peretilor care nu sunt atinși de bilă și care sunt nerelevanți în acel context deci pot fi eliminați. Totul a fost formulat în produse scalare. La nivel practic aceasta oferă posibilitatea ca algoritmul să lucreze într-un spațiu de dimensiune mare. Astfel noile sabloane $\Phi(x_i)$ pot fi rezultatul mapei datelor de intrare originale x_i într-un spațiu de dimensiune mai mare utilizând funcția Φ . Maximizarea funcției țintă și evaluarea funcției de decizie implică calculul produsului scalar $\phi(x), \phi(x)$ într-un spațiu de dimensiune mai mare. Aceste calcule costisitoare sunt reduse semnificativ prin utilizarea unui nucleu pozitiv definit k , astfel ca $k(x, x') := \langle \Phi(x), \Phi(x') \rangle$. Această substituție, care este referită uneori ca trucul nucleu, este utilizată pentru a extinde clasificarea cu hiperplane la SVM-ul neliniar. Trucul nucleu poate fi aplicat de vreme ce toți vectorii de trasaturi apar doar în produse scalare. Vectorii de trasaturi devin o expresie în spațiul de trasaturi și asadar Φ va fi funcția prin care reprezentăm vectorii de intrare în noul spațiu. Astfel funcția de decizie va avea următoarea formă:

$$f(x) = \text{sgn} \left(\sum_{i=1}^m y_i \alpha_i k(x, x_i) + b \right)$$

Principalul avantaj al acestui algoritm este că nu necesită transpunerea tuturor datelor de intrare într-un spațiu de dimensiune mai mare. De aceea nu vor exista nici calcule costisitoare ca în cazul rețelelor neuronale. De asemenea obținem o micșorare a setului de date în faza de antrenare când vom considera doar vectorii suport care de obicei sunt în număr redus. Un alt avantaj al acestui algoritm este că permite utilizarea datelor de intrare cu un număr oricât de mare de trasaturi fără a crește exponențial timpul de antrenare. Această caracteristică nu este adevărată pentru rețelele neurale, de exemplu algoritmul backpropagation are probleme când trebuie să lucreze cu un număr mare de trasaturi. Singura problemă care apare la SVM este numărul de vectori suport rezultați. Când numărul acestora crește timpul de răspuns în faza de testare crește și el liniar.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Dupa alegerea corecta a parametrilor clasificatorului SVM, se va crea un model pornind de la un set de texte exemplu dintr-un domeniu. Pentru verificarea acuratetei modelului se va folosi cross-validation. Cross-validarea este o tehnica de validare pentru a evalua modul in care rezultatele unei analize statistice se va generaliza la un set de date independent . Acesta este utilizat in principal in cazul in care obiectivul este predictia , si se doreste estimarea cu exactitate modul in care un model predictiv va performa in practica.

Am folosit un corpus de training din 19 domenii. Dupa cross-validare s-a obtinut o acuratete de 93,70 % a modelului.

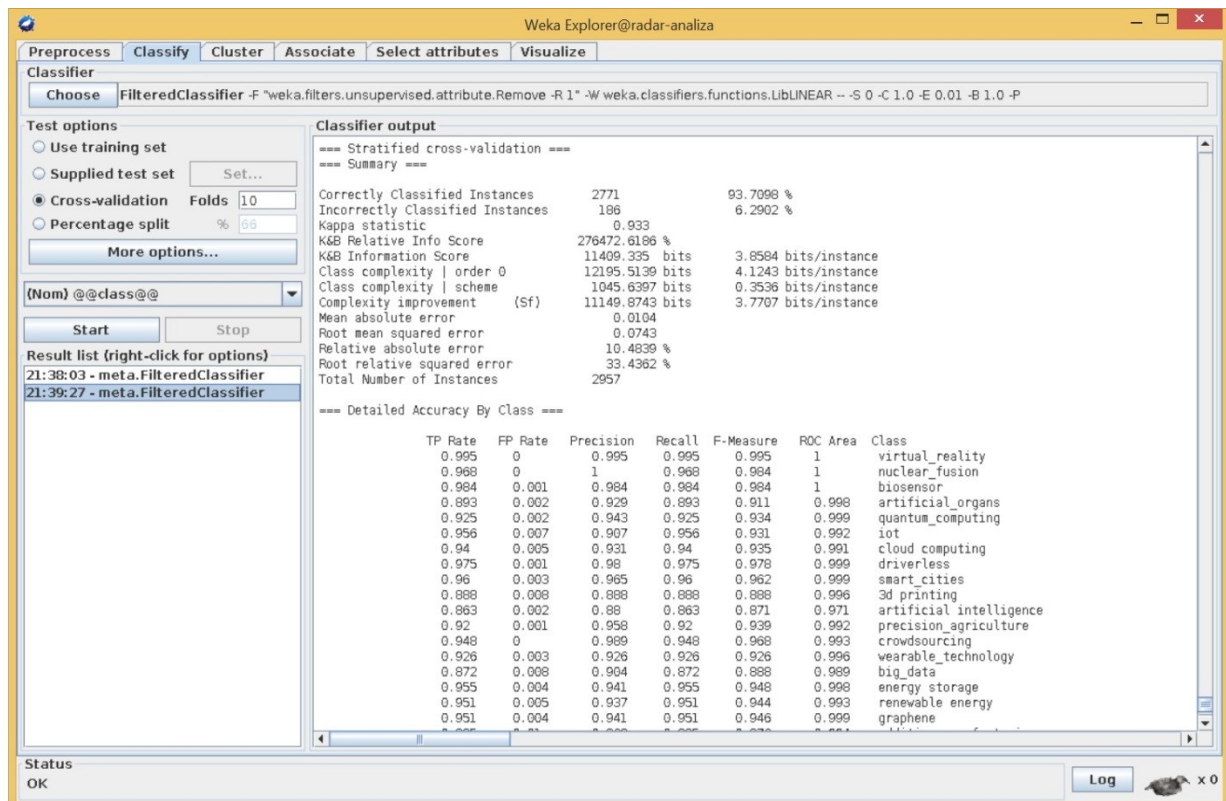


Fig. 9

Urmatorul pas este testarea unui corpus nou avand la baza modelul creat. Acesta este cel mai important pas deoarece scopul experimentului este ca modelul sa repartizeze corect stiri dintr-un corpus la prima vedere. Am folosit 400.000 de stiri iar outputul a fost :

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

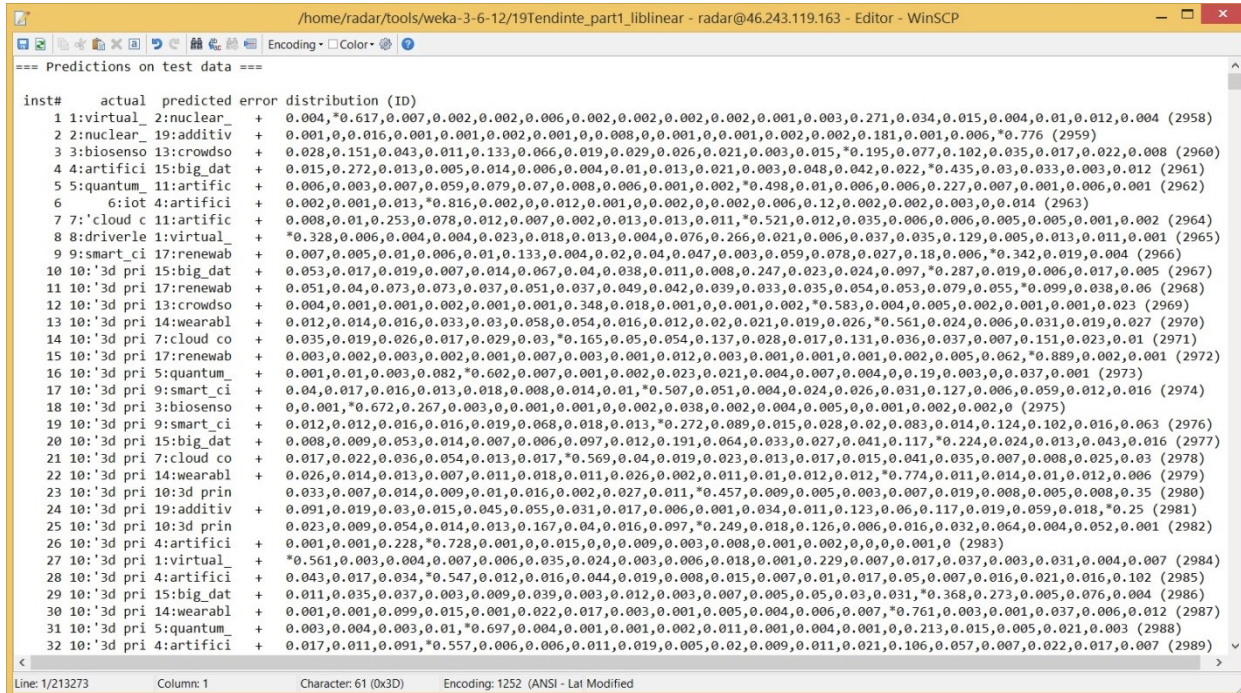


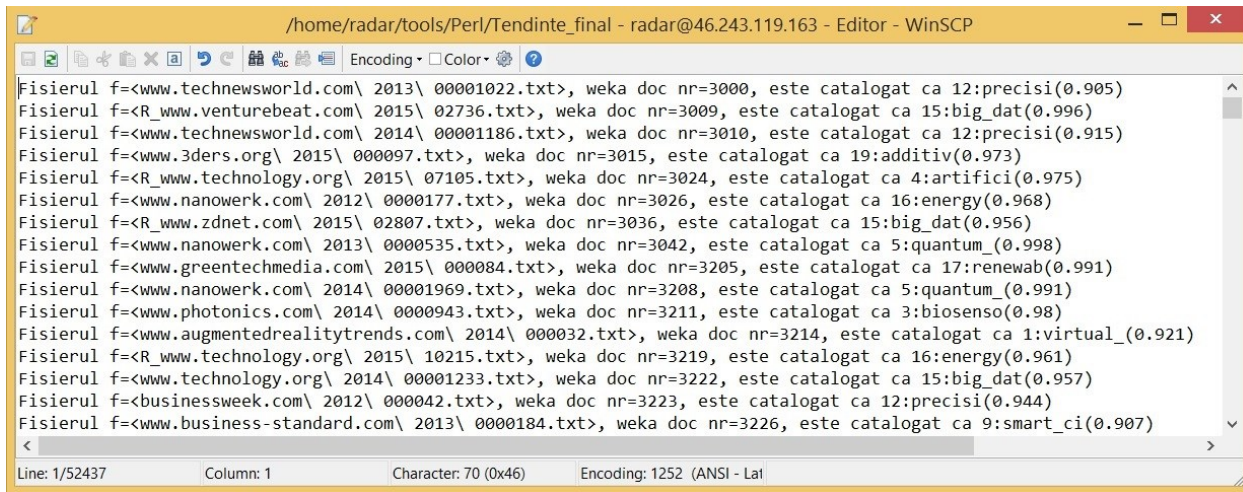
Fig. 10. Probabilitatea de incadrare a fiecărei stiri

Fiecare stare are un ID (numarul din paranteze de la finalul randului) pentru identificare. După aplicarea unor funcții logistice, L2 regularized logistic regression în cazul nostru, vom putea vedea probabilitatea de încadrare a fiecărei stiri.

Coloana a 2-a reprezintă domeniul corect al stirii, iar mai departe vedem probabilitatea de încadrare pentru fiecare domeniu.

Pentru o vizualizare mai bună, formatam aceste rezultate folosind un program perl și obținem:

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013



```

/home/radar/tools/Perl/Tendinte_final - radar@46.243.119.163 - Editor - WinSCP
Encoding: Color
Fisierul f=<www.technewsworld.com\ 2013\ 00001022.txt>, weka doc nr=3000, este catalogat ca 12:precisi(0.905)
Fisierul f=<R_www.venturebeat.com\ 2015\ 02736.txt>, weka doc nr=3009, este catalogat ca 15:big_dat(0.996)
Fisierul f=<www.technewsworld.com\ 2014\ 00001186.txt>, weka doc nr=3010, este catalogat ca 12:precisi(0.915)
Fisierul f=<www.3ders.org\ 2015\ 000097.txt>, weka doc nr=3015, este catalogat ca 19:additiv(0.973)
Fisierul f=<R_www.technology.org\ 2015\ 07105.txt>, weka doc nr=3024, este catalogat ca 4:artifici(0.975)
Fisierul f=<www.nanowerk.com\ 2012\ 0000177.txt>, weka doc nr=3026, este catalogat ca 16:energy(0.968)
Fisierul f=<R_www.zdnet.com\ 2015\ 02807.txt>, weka doc nr=3036, este catalogat ca 15:big_dat(0.956)
Fisierul f=<www.nanowerk.com\ 2013\ 0000535.txt>, weka doc nr=3042, este catalogat ca 5:quantum_(0.998)
Fisierul f=<www.greentechmedia.com\ 2015\ 000084.txt>, weka doc nr=3205, este catalogat ca 17:renewab(0.991)
Fisierul f=<www.nanowerk.com\ 2014\ 00001969.txt>, weka doc nr=3208, este catalogat ca 5:quantum_(0.991)
Fisierul f=<www.photonics.com\ 2014\ 0000943.txt>, weka doc nr=3211, este catalogat ca 3:biosenso(0.98)
Fisierul f=<www.augmentedrealitytrends.com\ 2014\ 000032.txt>, weka doc nr=3214, este catalogat ca 1:virtual_(0.921)
Fisierul f=<R_www.technology.org\ 2015\ 10215.txt>, weka doc nr=3219, este catalogat ca 16:energy(0.961)
Fisierul f=<www.technology.org\ 2014\ 00001233.txt>, weka doc nr=3222, este catalogat ca 15:big_dat(0.957)
Fisierul f=<businessweek.com\ 2012\ 000042.txt>, weka doc nr=3223, este catalogat ca 12:precisi(0.944)
Fisierul f=<www.business-standard.com\ 2013\ 0000184.txt>, weka doc nr=3226, este catalogat ca 9:smart_ci(0.907)
Line: 1/52437 Column: 1 Character: 70 (0x46) Encoding: 1252 (ANSI - Lat)

```

Fig. 11. Repartizarea stiriilor pe domenii

Pentru determinarea pragului de acuratete folosim un procedeu de maximizare a tolerantei. Prin rezolvarea ecuatiilor de mai jos obtin un pran optimal pentru a asigura o calitate ridicata a clusterizarii:

$$P\left(\int_{-\infty}^{\bar{x}+ks} f(y) dy \geq 1 - \beta\right) = 1 - \alpha \quad P\left(\int_{-\infty}^{(\bar{x}+ks-\mu)/\sigma} \varphi(y) dy \geq 1 - \beta\right) = 1 - \alpha$$

$$P\left((\bar{x} + ks - \mu) / \sigma \geq u_{1-\beta}\right) = 1 - \alpha$$

3. Folosirea Retelelor Neuronale pentru gasirea similaritatii

Clusterizarea automata a stiriilor, prezentata in sectiuna anterioara, este primul pas in gasirea similaritaiilor intre diferte parti constituent a unui flux de stiri. De la nivelul macro analiza, (document) se continua descendend catre nivelul de relational al cuvintelor. Focul este adus catre relatiile semantice dintre cuvinte.

In general o topica este descrisa de un set ce poate fi considerat din punct de vedere practic ca deschis. Cu alte cuvinte, un analizator trebuie sa aiba capacitatea sa recunoasca o topica chiar daca cuvinte noi sint folosite. Acest lucru se realizeaza prin determinarea similaritii semantice intre cuvinte. In acest scop se folosesc



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

modele neurale care sînt capabile să atingă un grad ridicat de acuratețe în precizarea similarității individuale (ca opusă similarității combinatoriale).

Determinarea similarității se face prin calcul vectorial. Contextul este reprezentat ca vectori iar similaritatea este calculată via proximitate vectorială dată de o metrică oarecare. Specific, în acest caz, folosim produsul punctual și metrica euclidiană.

Folosim procuul "sac-de-cuvinte". Peste corpusul de stiri, calculăm și maximizăm următoarea serie de probabilități:

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

Deci vom găsi setul de features care maximizează probabilitatea de mai sus.

Acest proces este unul iterativ. Optimizarea este obținută pas cu pas, repetind calcularea max loglikelihood de mai sus.

În figura următoare se prezintă arhitectura rețelei neurale recurente.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

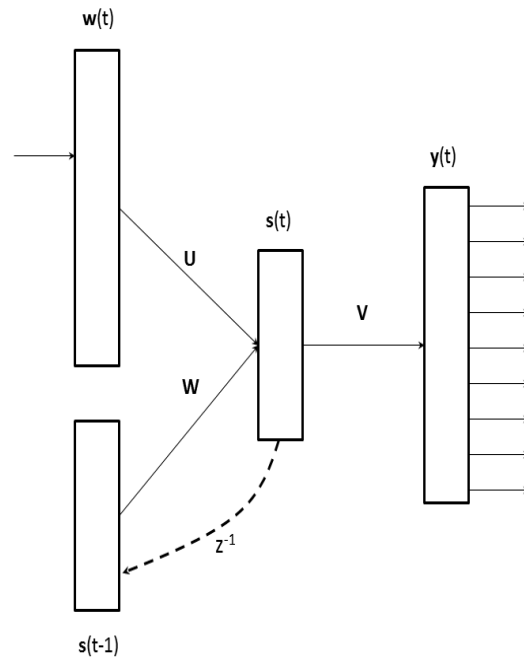


Fig. 12.

Formula de liniarizarea a inputului este:

$$s(t) = f(Uw(t) + Ws(t-1)) \quad (1)$$

$$y(t) = g(Vs(t)), \quad (2)$$

where

$$f(z) = \frac{1}{1 + e^{-z}}, \quad g(z_m) = \frac{e^{z_m}}{\sum_k e^{z_k}}. \quad (3)$$

În acest fel obținem un set de cuvinte similare pentru fiecare target de care avem nevoie. Prezentăm câteva exemple mai jos.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

```
brains  
frontal_lobe  
cerebellum  
cerebral_cortex  
temporal_lobe  
cortex  
brainstem  
frontal_lobes  
neural  
brain_frontal_lobe  
nerve_cells  
neurons  
brain_circuitry  
neurological  
temporal_lobes
```

Fig. 13. Multimea de Similaritate pentru brain

Un system care detecteaza semnale slabe vs. non-semnale pentru rapoarte asupra activitatii de cercetare ale creierului, va fi activat nu numai de cuvintul "brain", dar si de cuvinte ca "nerve_cells", sau "neurons".

Probabilitatile a priori sunt calculate de RNR (rețeaua neuronală recurentă). Aceste probabilitati sunt a priori pentru ca au fost calculate pentru cuvinte , intr-un mod individual. Ele reprezinta doar asocieri unul la unul intre cuvinte. Semantica unei topici se calculeaza via nivelul de discriminativitate intre topici , peste nivelul individual.

De exemplu , pentru similaritatile individuale obtinute in aput de RNR pentru "dentist", vezi figura ..., avem

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

```

orthodontist
Dentist
periodontist
doctor
pediatric_dentist
dentists
oral_surgeon
dental
dermatologist
plastic_surgeon
endodontist
cosmetic_dentist
chiropractor
dental_hygienist
dentistry
    
```

Fig. 14. Similarități individuale dentist

“plastic surgeon”, poate într-adevăr activa în mod corect în modul; care caută semnale slabe pentru “dentist”. Dar, această probabilitate trebuie ponderată via probabilității globale.

Cu o tehnică de redistribuire a probabilităților, rezolvind și problema probabilităților asociatelor cu cuvinte încă neprocesate, folosind o formulă, (Katz):

$$P_{bo}(w_i | w_{i-n+1} \dots w_{i-1}) = \begin{cases} d_{w_{i-n+1} \dots w_i} \frac{C(w_{i-n+1} \dots w_{i-1} w_i)}{C(w_{i-n+1} \dots w_{i-1})} & C(w_{i-n+1} \dots w_i) > k \\ \alpha_{w_{i-n+1} \dots w_{i-1}} P_{bo}(w_i | w_{i-n+2} \dots w_{i-1}) & \end{cases}$$

Cu C reprezentând numărul de ocurențe pentru o anumită topică.

4. Annotarea Elementelor Semnificative

În detectarea semnalelor slabe vs. tari, nu numai cuvintele cheie și cele similare lor sunt importante, sau domeniul și conectivitățile semantice găsite prin procedeele prezentate în secțiunile de mai sus.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

În această secțiune vom prezenta modulul de anotare a elementelor semnificative dintr-un document. După identificarea cuvintelor cheie și încadrarea într-un topic a unui document, acest modul detectează agenții și relațiile dintre cuvintele de interes și agenții respective.

La baza acestui modul stă un model al limbajului generat prin tehnica Named Entity Recognition (NER), sau Identificarea Entităților Semnificative (IES). Tehnica IES a fost dezvoltată pe baza algoritmilor HMM and SVM.

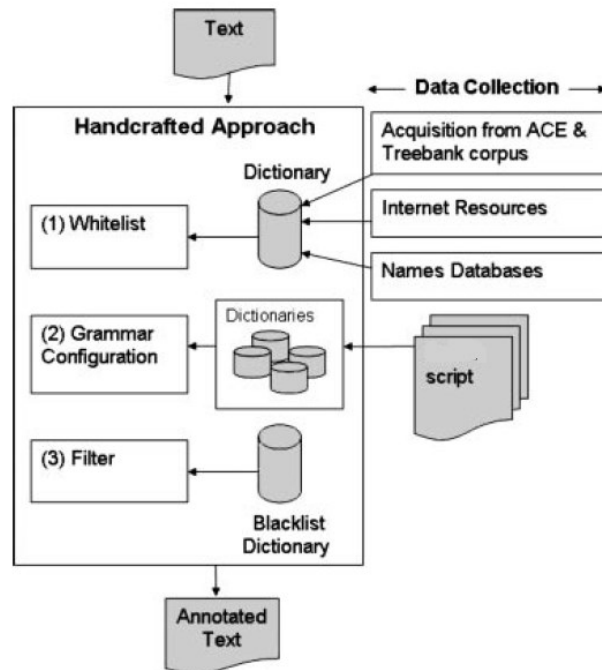


Fig. 15.

Modelul general este prezentat în *Figura de mai sus*. Detectarea agenților și tipului este un proces în care se creează un model, care este folosit pentru anotarea agenților în documentele de interes. Formula este dată via un proces Gibbs (lanțuri Markov)

$$P_A(\mathbf{s}^{(t)} | \mathbf{s}^{(t-1)}) = \frac{P_M(s_i^{(t)} | \mathbf{s}_{-i}^{(t-1)}, \mathbf{o})^{1/c_t}}{\sum_j P_M(s_j^{(t)} | \mathbf{s}_{-j}^{(t-1)}, \mathbf{o})^{1/c_t}}$$

Ceea ce se obține în output este:



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Researchers at UC San Diego School of Medicine conducted the first population-based study that characterizes the association and temporal relationship between gastrointestinal stromal tumors (GIST) and other cancers. The results indicate that one in 5.8 patients with GIST will develop additional malignancies before and after their diagnosis. Specifically, patients with GIST are more likely to develop other sarcomas, non-Hodgkin's lymphoma, carcinoid tumors, melanoma, colorectal, esophageal, pancreatic, hepatobiliary, non-small cell lung, prostate and renal cell cancers. "Only 5 percent of patients with gastrointestinal stromal tumors have a hereditary disorder that predisposes them to develop multiple benign and malignant tumors," said Jason K. Sicklick, MD, assistant professor of surgery and UC San Diego Moores Cancer Center surgical oncologist. "The research indicates that these patients may develop cancers outside of these syndromes, but the exact mechanisms are not yet known." The researchers said further studies are needed to understand the connection between GIST and other cancers, but the findings may have clinical implications. "Patients diagnosed with gastrointestinal stromal tumors may warrant consideration for additional screenings based on the other cancers that they are most susceptible to contract," said co-author James D. Murphy, MD, assistant professor of radiation oncology and UC San Diego Moores Cancer Center radiation oncologist. When compared to the United States population, the researchers found that people with GIST had a 44 percent increased prevalence of cancers occurring before a GIST diagnosis and a 66 percent higher risk of developing cancers after diagnosis. The most common tumors were those of the genitourinary tract, breast, respiratory and blood. Non-Hispanic patients had a higher incidence of other cancers before a GIST diagnosis. Patients whose tumors were smaller than 10 centimeters had a higher probability of a second cancer than patients whose growth was larger. People with tumors smaller than 2 cm had the greatest likelihood of developing additional malignancies, both before and after diagnosis.

Fig. 16

Putem astfel sa punem in relatie nu numai substantivele commune dar si substantivele proprii care definesc agentii care sint mentionati in document.

In exemplul de mai sus reusim sa extragem automat urmatoarele fapte:

GIST este un tip de tumoare (gastrointestinal)

San Diego este un centru de cercetare responsabil de o descoperire cu implicatii importante

D Murph este un cercetator Implicat in GIST

In viitor, adica in prelucrarea documentelor ulterioare, putem urmarii ceea ce se intampla cu agentii, stiind ca ei sint conectati la o anumita topica. Astfel se asigura un control al acuratetei procesarii semnalelor , mult mai ridicat.

Concluzii

Am dezvoltat un sistem de procesare a semnalelor pe baza a patru instrumente fundamentale:

- SVM
- LDA
- RNR
- NER

Rezultatele au confirmat puterea de procesare ridicata a sistemului implemetat.

Prezentarea algoritmilor utilizati pentru clasificarea si adnotarea automata a stirilor:

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

I. LDA:

1. Se instalează softul MALLET - <http://mallet.cs.umass.edu/>
2. Următorul pas este importarea datelor – acestea se pun într-un singur folder.



```
radar@radar-analiza: ~/tools/mallet-2.0.7
radar@radar-analiza:~/tools/mallet-2.0.7$ ./bin/mallet import-dir --input "calea
  catre folder" --output "nume fisier unic".mallet --keep-sequence --remove-stop
  words
```

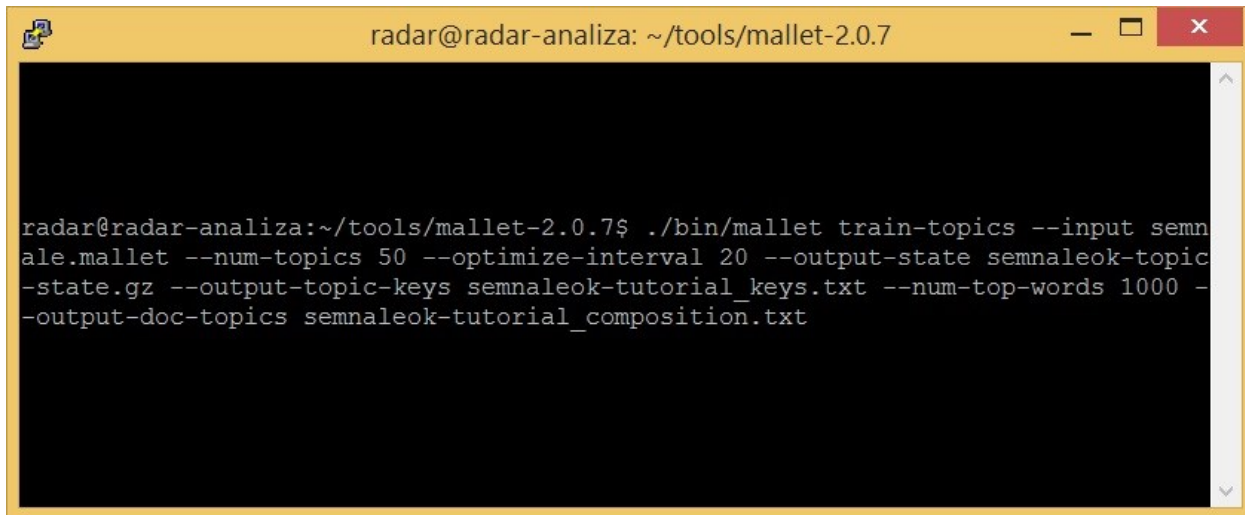
Fig. 16

Pentru a lucra cu întregul corpus va trebui să transformăm toate fișierele într-un singur fișier de tip .mallet. Pentru acest pas folosim comanda **import**

Comanda `--keep-sequence` păstrează ordinea inițială a fișierelor, iar `--remove-stopwords` șterge cuvintele care ar obstrucționa analiza corectă.

3. Apoi se construiește modelul folosind comanda **train topic**

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013



```
radar@radar-analiza: ~/tools/mallet-2.0.7

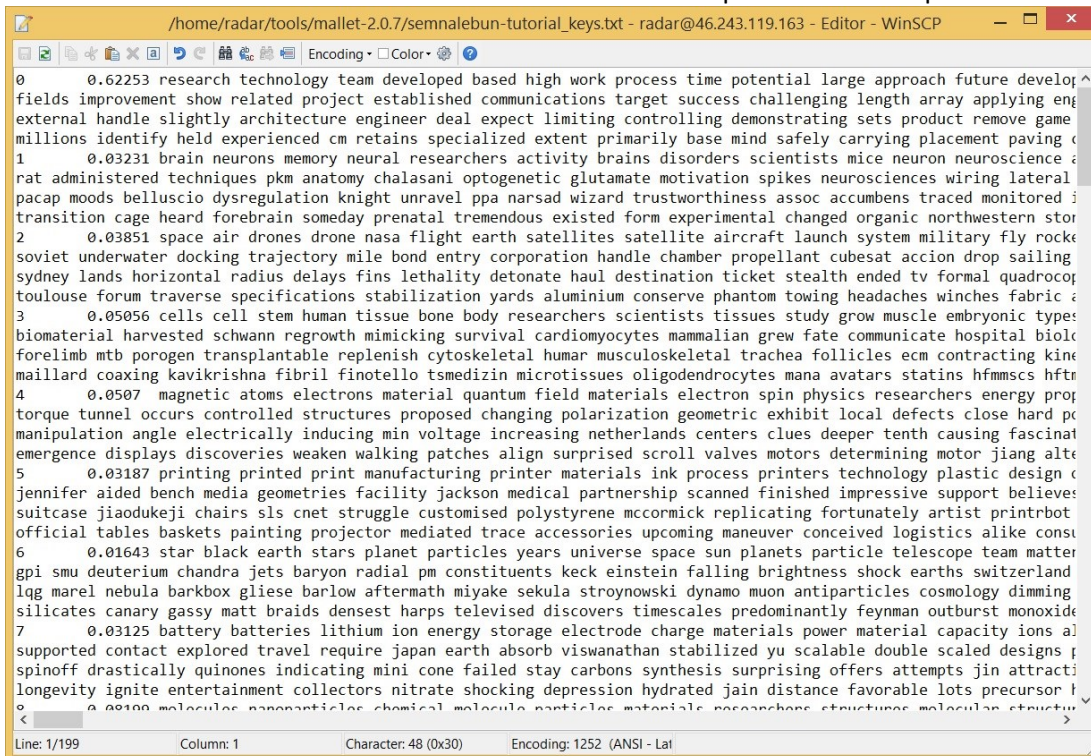
radar@radar-analiza:~/tools/mallet-2.0.7$ ./bin/mallet train-topics --input semn
ale.mallet --num-topics 50 --optimize-interval 20 --output-state semnleok-topi
c-state.gz --output-topic-keys semnleok-tutorial_keys.txt --num-top-words 1000 -
-output-doc-topics semnleok-tutorial_composition.txt
```

Fig.

Fig. 17

4. Aceasta comanda:

- Deschide fisierul *.mallet
- Mallet creeaza 50 de topic
- Creeaza un fisier txt ce contine toate cuvintele cheie pentru fiecare topic



```
/home/radar/tools/mallet-2.0.7/semnalebun-tutorial_keys.txt - radar@46.243.119.163 - Editor - WinSCP

0 0.62253 research technology team developed based high work process time potential large approach future develop
fields improvement show related project established communications target success challenging length array applying eng
external handle slightly architecture engineer deal expect limiting controlling demonstrating sets product remove game
millions identify held experienced cm retains specialized extent primarily base mind safely carrying placement paving c
1 0.03231 brain neurons memory neural researchers activity brains disorders scientists mice neuron neuroscience a
rat administered techniques pkm anatomy chalasani optogenetic glutamate motivation spikes neurosciences wiring lateral
pacap moods belluscio dysregulation knight unravel ppa narsad wizard trustworthiness assoc accumbens traced monitored i
transition cage heard forebrain someday prenatal tremendous existed form experimental changed organic northwestern stor
2 0.03851 space air drones drone nasa flight earth satellites satellite aircraft launch system military fly rocke
soviet underwater docking trajectory mile bond entry corporation handle chamber propellant cubesat accion drop sailing
sydney lands horizontal radius delays fins lethality detonate haul destination ticket stealth ended tv formal quadrocop
toulouse forum traverse specifications stabilization yards aluminium conserve phantom towing headaches winches fabric a
3 0.05056 cells cell stem human tissue bone body researchers scientists tissues study grow muscle embryonic type:
biomaterial harvested schwann regrowth mimicking survival cardiomyocytes mammalian grew fate communicate hospital biol
forelimb mtb porogen transplantable replenish cytoskeletal humar musculoskeletal trachea follicles ecm contracting kin
maillard coaxing kavikrishna fibril finotello tsmedizin microtissues oligodendrocytes mana avatars statins hfmmcs hftr
4 0.0507 magnetic atoms electrons material quantum field materials electron spin physics researchers energy prop
torque tunnel occurs controlled structures proposed changing polarization geometric exhibit local defects close hard po
manipulation angle electrically inducing min voltage increasing netherlands centers clues deeper tenth causing fascinat
emergence displays discoveries weaken walking patches align surprised scroll valves motors determining motor jiang alte
5 0.03187 printing printed print manufacturing printer materials ink process printers technology plastic design c
jennifer aided bench media geometries facility jackson medical partnership scanned finished impressive support believes
suitcase jiaodukeji chairs sls cnet struggle customised polystyrene mccormick replicating fortunately artist printrobot
official tables baskets painting projector mediated trace accessories upcoming maneuver conceived logistics alike consu
6 0.01643 star black earth stars planet particles years universe space sun planets particle telescope team matter
gpi smu deuterium chandra jets baryon radial pm constituents keck einstein falling brightness shock earths switzerland
lqg marel nebula barkbox gliese barlow aftermath miyake sekula stroynowski dynamo muon antiparticles cosmology dimming
silicates canary gassy matt braids densest harps televised discovers timescales predominantly feynman outburst monoxide
7 0.03125 battery batteries lithium ion energy storage electrode charge materials power material capacity ions al
supported contact explored travel require japan earth absorb viswanathan stabilized yu scalable double scaled designs p
spinoff drastically quinones indicating mini cone failed stay carbons synthesis surprising offers attempts jin attracti
longevity ignite entertainment collectors nitrate shocking depression hydrated jain distance favorable lots precursor h
8 0.02100 molecules nanoparticle chemical molecule particles materials researchers structures molecular structur
< >
```

Fig. 18

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

5. Deschidem fișierul ***-tutorial_composition.txt** într-un fișier nou Excel pentru a-l putea vizualiza corect. Acesta conține repartizarea fiecărui fișier pe topic

#doc	name	topic	proportion	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
1	file:/home/radar/tools/semnale/R_www.designboom.com_tech																				
2	0	file:/home/radar/tools/semnale/R_www.designboom.com_tech				39	0.277821	20	0.187166	40	0.164548	0	0.11204	46	0.08407	25	0.068541	7	0.045645	19	0.018472
3	1	file:/home/radar/tools/semnale/R_www.sciencemag.com%2020				22	0.214547	46	0.164823	35	0.150674	28	0.127508	30	0.09281	14	0.087013	40	0.076184	38	0.04649
4	2	file:/home/radar/tools/semnale/impactlab_2014%2000443.txt				30	0.448803	46	0.272684	20	0.195362	7	0.052564	9	0.00775	0	0.004654	33	0.020216	40	0.001212
5	3	file:/home/radar/tools/semnale/ScienceDaily_2014%2013987.txt				14	0.448082	46	0.140153	35	0.105713	33	0.060963	2	0.059229	40	0.03686	22	0.036725	31	0.029699
6	4	file:/home/radar/tools/semnale/www.sciencedaily.com%202015				14	0.359871	33	0.170119	46	0.152603	21	0.066664	0	0.059325	40	0.051886	23	0.033613	13	0.021317
7	5	file:/home/radar/tools/semnale/www.bbc.com_science_and_em				48	0.219896	4	0.207092	0	0.197864	41	0.194263	46	0.061226	33	0.053433	22	0.05212	40	0.001047
8	6	file:/home/radar/tools/semnale/www.technology.org%202015				28	0.296084	0	0.265733	22	0.085786	25	0.06376	46	0.062324	16	0.056411	5	0.053899	17	0.034367
9	7	file:/home/radar/tools/semnale/www.sciencedaily.com%202015				44	0.368963	27	0.206499	47	0.13203	46	0.10307	45	0.05274	33	0.031917	36	0.030874	0	0.029076
10	8	file:/home/radar/tools/semnale/Popsci_2014%2000023.txt				38	0.222533	46	0.13892	12	0.096507	10	0.081639	25	0.059544	0	0.057911	40	0.052498	24	0.052075
11	9	file:/home/radar/tools/semnale/futurity_medicine%200361.txt				47	0.378321	33	0.162615	0	0.127305	35	0.126122	8	0.066195	1	0.050088	20	0.026278	45	0.018033
12	10	file:/home/radar/tools/semnale/popsci_2013%2000167.txt				6	0.533364	46	0.206657	33	0.089925	14	0.042637	40	0.039271	0	0.037186	47	0.023465	23	0.019779
13	11	file:/home/radar/tools/semnale/ec.europa.eu%202015%2000000				25	0.276307	0	0.180237	5	0.108727	29	0.097413	46	0.087095	20	0.070473	16	0.065779	23	0.041013
14	12	file:/home/radar/tools/semnale/robohub%2000487.txt				29	0.30224	28	0.228548	46	0.150817	2	0.088671	33	0.075642	22	0.066678	39	0.029872	40	0.023291
15	13	file:/home/radar/tools/semnale/www.futurity.org_med%202015				0	0.192346	17	0.173918	18	0.148694	22	0.106699	46	0.093831	3	0.075822	49	0.056141	36	0.053346
16	14	file:/home/radar/tools/semnale/phys.org%2000496.txt				8	0.198204	0	0.156843	43	0.155436	22	0.091144	33	0.072959	39	0.072386	17	0.06434	34	0.056458
17	15	file:/home/radar/tools/semnale/phys.org%2000211.txt				0	0.312547	19	0.252106	39	0.166254	17	0.126157	33	0.05304	46	0.02561	40	0.012371	1	0.011299
18	16	file:/home/radar/tools/semnale/robohub%2000420.txt				2	0.436618	14	0.152237	38	0.152217	18	0.076733	39	0.057908	21	0.057808	0	0.011798	46	0.009016
19	17	file:/home/radar/tools/semnale/tech_review%2000540.txt				23	0.216371	30	0.174658	12	0.142512	46	0.130763	0	0.108477	11	0.101223	40	0.07881	20	0.02357
20	18	file:/home/radar/tools/semnale/www.theengineer.co.uk%20201				27	0.383377	13	0.154947	0	0.118458	46	0.097306	43	0.061148	4	0.060839	34	0.054462	41	0.040771
21	19	file:/home/radar/tools/semnale/R_www.sciencedaily.com%2020				9	0.526979	0	0.185876	18	0.090561	46	0.079422	33	0.078779	34	0.019023	40	0.00674	47	0.003342
22	20	file:/home/radar/tools/semnale/futurity_medicine%2000250.txt				36	0.413157	34	0.23818	33	0.120327	1	0.10767	46	0.064921	0	0.049012	25	0.042929	24	0.02388
23	21	file:/home/radar/tools/semnale/www.biopace.com%202015%2				3	0.470691	33	0.169381	34	0.085801	46	0.072429	27	0.066624	0	0.049012	25	0.042929	24	0.02388
24	22	file:/home/radar/tools/semnale/www.technology.org%202015%				9	0.275269	3	0.219834	0	0.161323	46	0.088479	34	0.074882	33	0.051773	1	0.038468	18	0.027837
25	23	file:/home/radar/tools/semnale/R_www.sciencedaily.com%2020				13	0.449895	0	0.088805	33	0.087236	32	0.065807	18	0.057842	37	0.048184	1	0.043371	42	0.043349
26	24	file:/home/radar/tools/semnale/futurity_sci_tech%2000756.txt				22	0.319901	43	0.256153	0	0.168415	21	0.069503	33	0.066002	46	0.023837	16	0.021618	37	0.021479
27	25	file:/home/radar/tools/semnale/R_www.sciencedaily.com%2020				34	0.573683	1	0.112063	0	0.08014	39	0.068039	46	0.066002	13	0.044275	33	0.02127	23	0.020756
28	26	file:/home/radar/tools/semnale/www.3ders.org%202015%20027				18	0.446022	46	0.174439	0	0.158234	26	0.092667	33	0.05581	5	0.033781	20	0.013176	13	0.008625
29	27	file:/home/radar/tools/semnale/R_www.3dprintingindustry.com				5	0.340622	2	0.316035	0	0.183055	46	0.047077	39	0.033046	40	0.023381	33	0.017515	20	0.008749
30	28	file:/home/radar/tools/semnale/www.azonano.com%202015%2				17	0.390229	1	0.220382	0	0.158248	3	0.098338	19	0.050587	46	0.046384	33	0.013986	44	0.013369

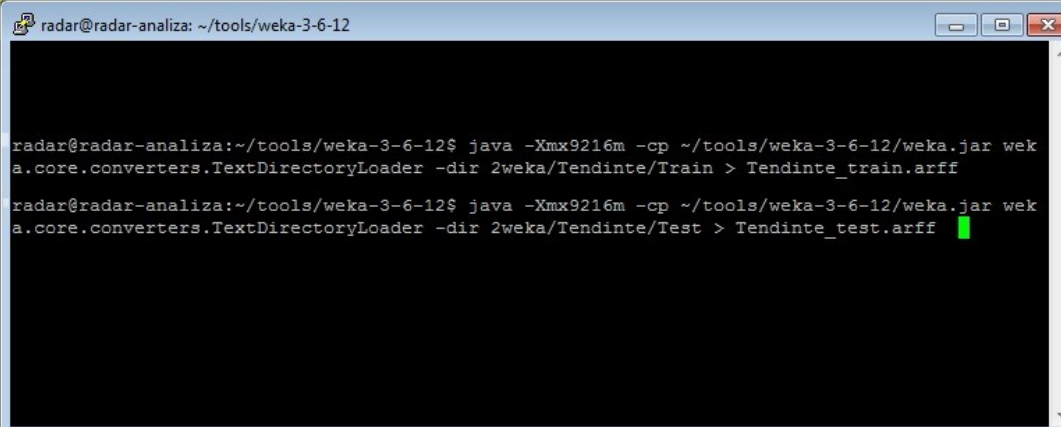
Fig. 19

Exemplu: Fișierul 10 are topicul nr 6 ca topic principal cu 53%, topicul 46 cu 20%.

II. SVM:

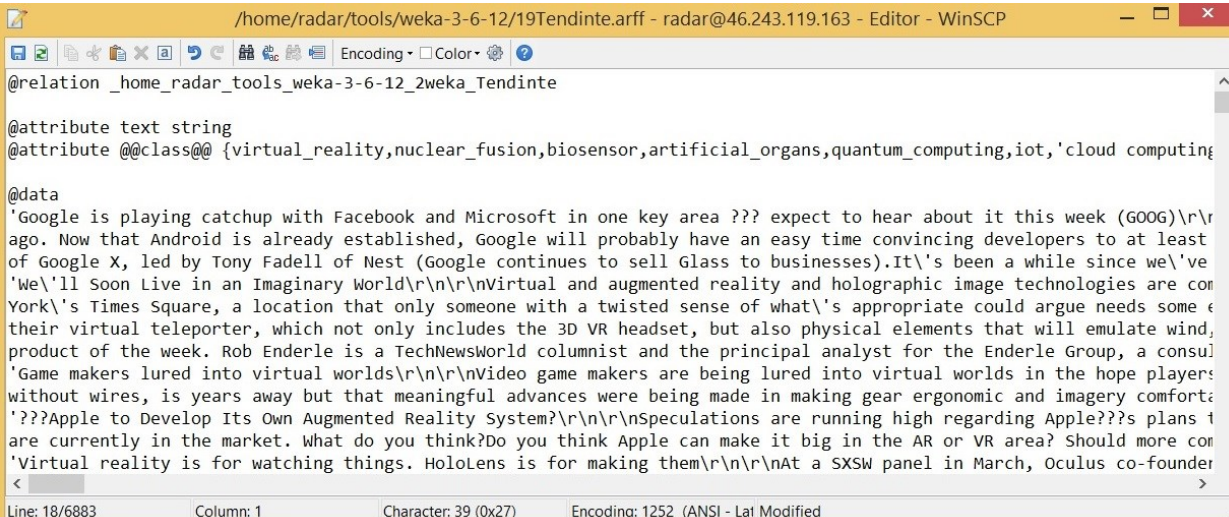
1. Se instalează WEKA, un software "open-source" de învățare automată și LibSVM, un clasificator SVM
2. Cele 2 corpuri de texte, cel de învățare și cel de test vor fi transformate în fișiere ARFF (**Attribute-Relation File Format**). Fiecare text va fi reprezentat ca un rând în noul fișier. Pentru asta vom folosi un program Java, numit TextDirectoristoARFF. Pentru a utiliza acest program, fiecare categorie din colecția de texte trebuie să aibă propriul său director.

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013



```
radar@radar-analiza: ~/tools/weka-3-6-12
radar@radar-analiza:~/tools/weka-3-6-12$ java -Xmx9216m -cp ~/tools/weka-3-6-12/weka.jar weka
a.core.converters.TextDirectoryLoader -dir 2weka/Tendinte/Train > Tendinte_train.arff
radar@radar-analiza:~/tools/weka-3-6-12$ java -Xmx9216m -cp ~/tools/weka-3-6-12/weka.jar weka
a.core.converters.TextDirectoryLoader -dir 2weka/Tendinte/Test > Tendinte_test.arff
```

Fig. 21



```
/home/radar/tools/weka-3-6-12/19Tendinte.arff - radar@46.243.119.163 - Editor - WinSCP
@relation _home_radar_tools_weka-3-6-12_2weka_Tendinte
@attribute text string
@attribute @@class@@ {virtual_reality,nuclear_fusion,biosensor,artificial_organs,quantum_computing,iot,'cloud computing
@data
'Google is playing catchup with Facebook and Microsoft in one key area ??? expect to hear about it this week (GOOG)\r\r
ago. Now that Android is already established, Google will probably have an easy time convincing developers to at least
of Google X, led by Tony Fadell of Nest (Google continues to sell Glass to businesses).It\'s been a while since we\'ve
'Ve\'ll Soon Live in an Imaginary World\r\r\r\nVirtual and augmented reality and holographic image technologies are cor
York\'s Times Square, a location that only someone with a twisted sense of what\'s appropriate could argue needs some €
their virtual teleporter, which not only includes the 3D VR headset, but also physical elements that will emulate wind,
product of the week. Rob Enderle is a TechNewsWorld columnist and the principal analyst for the Enderle Group, a consul
'Game makers lured into virtual worlds\r\r\r\nVideo game makers are being lured into virtual worlds in the hope players
without wires, is years away but that meaningful advances were being made in making gear ergonomic and imagery comfort:
'???Apple to Develop Its Own Augmented Reality System?\r\r\r\nSpeculations are running high regarding Apple???s plans t
are currently in the market. What do you think?Do you think Apple can make it big in the AR or VR area? Should more cor
'Virtual reality is for watching things. HoloLens is for making them\r\r\r\nAt a SXSW panel in March, Oculus co-founder
```

Fig. 22

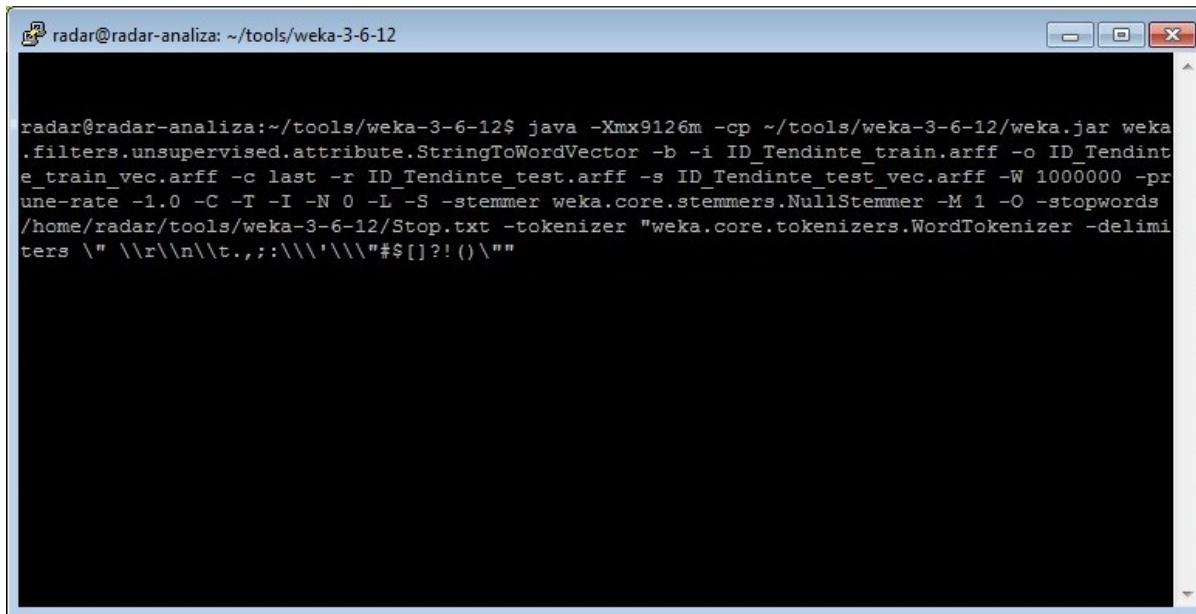
Iar outputul va arata astfel

3. Convertirea atributelor de tip "string" in attribute numerice

Fisierul de la pasul anterior are doar 2 attribute: "class" – care reprezinta categoriile corpusului si "text" – continutul textelor. Trebuie să transformam acest fișier într-un format de unde putem extrage caracteristici (attribute) și sa avem valori numerice pentru ele.

Pentru asta vom folosi tot un program java numit **StringtoWordVector**

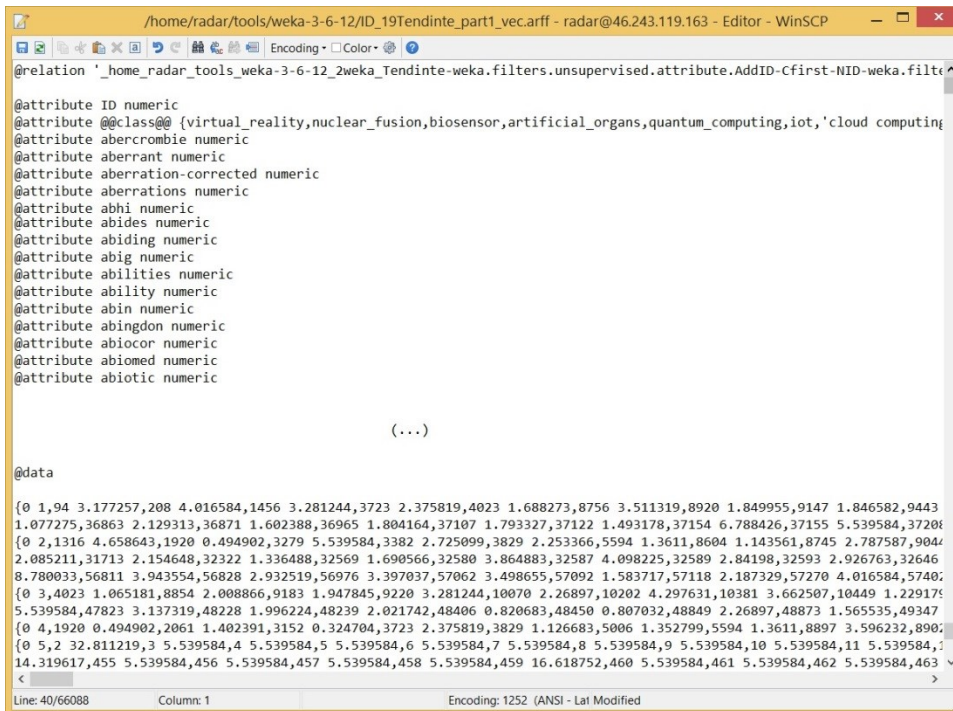
Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013



```
radar@radar-analiza: ~/tools/weka-3-6-12
radar@radar-analiza:~/tools/weka-3-6-12$ java -Xmx9126m -cp ~/tools/weka-3-6-12/weka.jar weka
.filters.unsupervised.attribute.StringToWordVector -b -i ID Tendinte_train.arff -o ID Tendint
e_train_vec.arff -c last -r ID Tendinte_test.arff -s ID Tendinte_test_vec.arff -W 1000000 -pr
une-rate -1.0 -C -T -I -N 0 -L -S -stemmer weka.core.stemmers.NullStemmer -M 1 -O -stopwords
/home/radar/tools/weka-3-6-12/Stop.txt -tokenizer "weka.core.tokenizers.WordTokenizer -delimi
ters \" \\r\\n\\t.,:;\\'\\/\\\"#\$%!?!()\""
```

Fig. 23

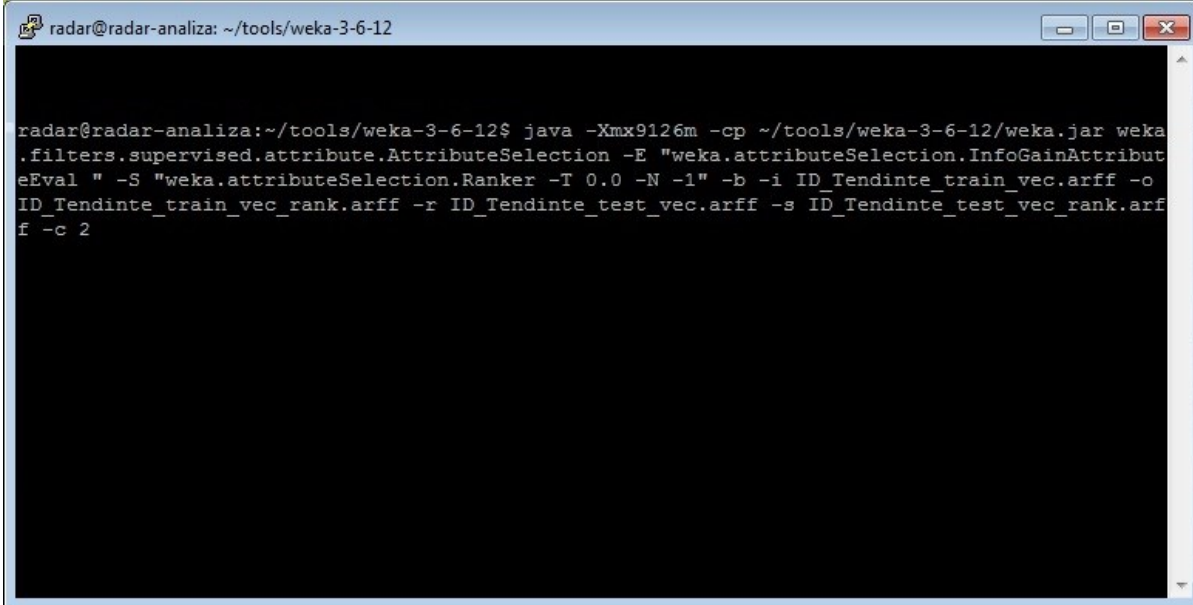
Rezultand fisierul de mai jos



```
/home/radar/tools/weka-3-6-12/ID_19Tendinte_part1_vec.arff - radar@46.243.119.163 - Editor - WinSCP
@relation '_home_radar_tools_weka-3-6-12_2weka_Tendinte-weka.filters.unsupervised.attribute.AddID-Cfirst-NID-weka.filt
@attribute ID numeric
@attribute @@class@@ {virtual_reality,nuclear_fusion,biosensor,artificial_organs,quantum_computing,iot,'cloud computing
@attribute abercrombie numeric
@attribute aberrant numeric
@attribute aberration-corrected numeric
@attribute aberrations numeric
@attribute abhi numeric
@attribute abides numeric
@attribute abiding numeric
@attribute abig numeric
@attribute abilities numeric
@attribute ability numeric
@attribute abin numeric
@attribute abingdon numeric
@attribute abiocor numeric
@attribute abiomed numeric
@attribute abiotic numeric
(...)
@data
{0 1,94 3.177257,208 4.016584,1456 3.281244,3723 2.375819,4023 1.688273,8756 3.511319,8920 1.849955,9147 1.846582,9443
1.077275,36863 2.129313,36871 1.602388,36965 1.804164,37107 1.793327,37122 1.493178,37154 6.788426,37155 5.539584,37208
{0 2,1316 4.658643,1920 0.494902,3279 5.539584,3382 2.725099,3829 2.253366,5594 1.3611,8604 1.143561,8745 2.787587,904
2.085211,31713 2.154648,32322 1.336488,32569 1.690566,32580 3.864883,32587 4.089225,32589 2.84198,32593 2.926763,32646
8.780033,56811 3.943554,56828 2.932519,56976 3.397037,57062 3.498655,57092 1.583717,57118 2.187329,57270 4.016584,5740;
{0 3,4023 1.065181,8854 2.008866,9183 1.947845,9220 3.281244,10070 2.26897,10202 4.297631,10381 3.662507,10449 1.229175
5.539584,47823 3.137319,48228 1.996224,48239 2.021742,48406 0.820683,48450 0.807032,48849 2.26897,48873 1.565535,49347
{0 4,1920 0.494902,2061 1.402391,3152 0.324704,3723 2.375819,3829 1.126683,5006 1.352799,5594 1.3611,8897 3.596232,890;
{0 5,2 32.811219,3 5.539584,4 5.539584,5 5.539584,6 5.539584,7 5.539584,8 5.539584,9 5.539584,10 5.539584,11 5.539584,1
14.319617,455 5.539584,456 5.539584,457 5.539584,458 5.539584,459 16.618752,460 5.539584,461 5.539584,462 5.539584,463
Line: 40/66088 Column: 1 Encoding: 1252 (ANSI - Lat Modified)
```

Fig. 24

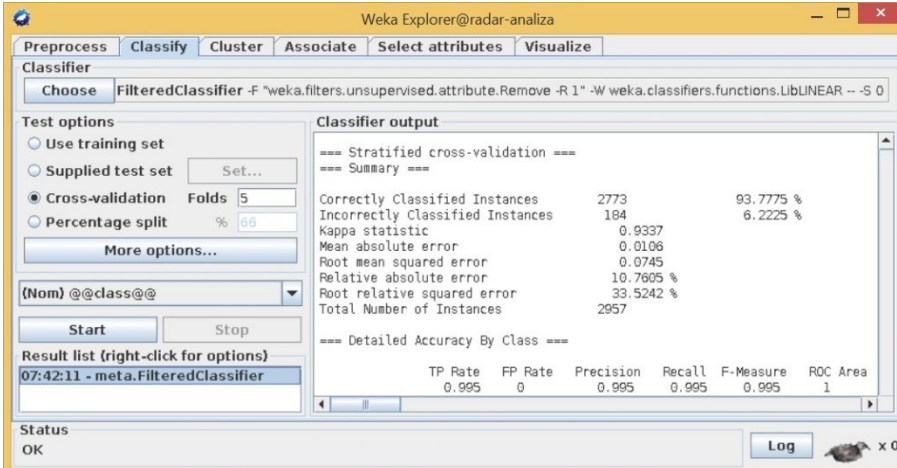
Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013



```
radar@radar-analiza: ~/tools/weka-3-6-12
radar@radar-analiza:~/tools/weka-3-6-12$ java -Xmx9126m -cp ~/tools/weka-3-6-12/weka.jar weka
.filters.supervised.attribute.AttributeSelection -E "weka.attributeSelection.InfoGainAttribut
eEval " -S "weka.attributeSelection.Ranker -T 0.0 -N -1" -b -i ID_Tendinte_train_vec.arff -o
ID_Tendinte_train_vec_rank.arff -r ID_Tendinte_test_vec.arff -s ID_Tendinte_test_vec_rank.arf
f -c 2
```

Fig. 25

4. Urmatorul pas este sa folosim functia **weka.filters.supervised.attribute.AttributeSelection**, care selecteaza cele mai importante atribute. Este o etapa care dureaza destul de mult, dar ofera rezultate foarte bune.
5. Dupa alegerea corecta a parametrilor clasificatorului SVM, se va crea un model pe baza textelor folosite pentru invatare. Pentru verificarea acuratetei modelului se va folosi cross-validation. Cross-validarea este o tehnica de validare pentru a evalua modul in care rezultatele unei analize statistice se va generaliza la un set de date independent.
In acest caz acuratetea modelului este de 93,77%



Weka Explorer@radar-analiza

Classifier: **FilteredClassifier** -F "weka.filters.unsupervised.attribute.Remove -R1" -W weka.classifiers.functions.LibLINEAR -- -S 0

Test options: Use training set, Supplied test set, Cross-validation (Folds: 5), Percentage split (%: 66)

Classifier output:

```

=== Stratified cross-validation ===
=== Summary ===
Correctly Classified Instances      2773           93.7775 %
Incorrectly Classified Instances    184            6.2225 %
Kappa statistic                     0.9337
Mean absolute error                  0.0106
Root mean squared error              0.0745
Relative absolute error              10.7605 %
Root relative squared error          33.5242 %
Total Number of Instances           2957

=== Detailed Accuracy By Class ===

```

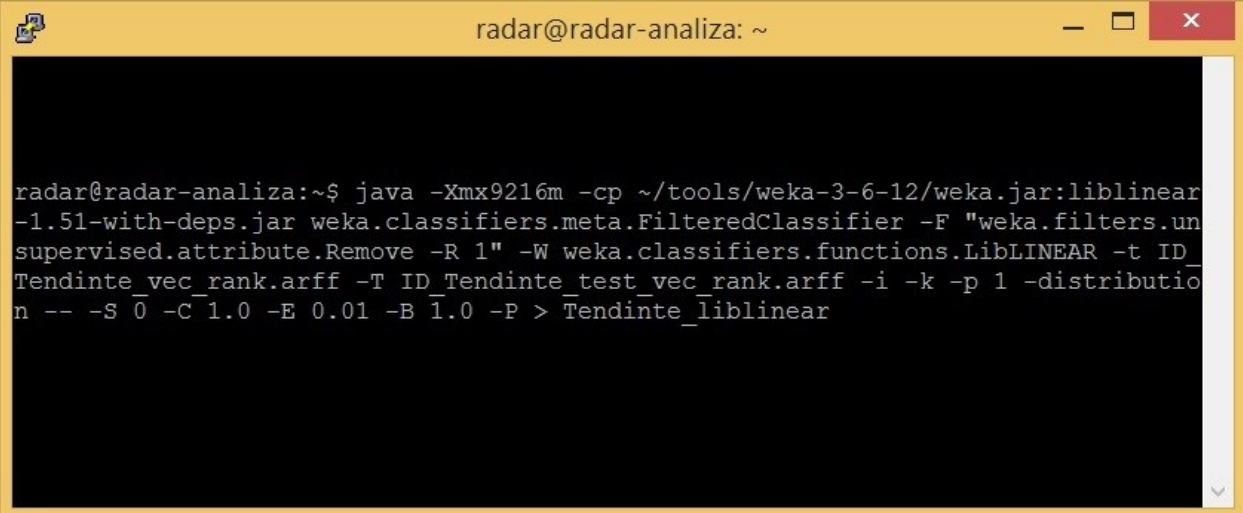
	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
	0.995	0	0.995	0.995	0.995	1

Status: OK

Fig. 26

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

6. După găsirea celor mai buni parametri și construirea modelului final aplicăm același clasificator pentru corpul de test, în cazul nostru – Liblinear



```
radar@radar-analiza: ~  
radar@radar-analiza:~$ java -Xmx9216m -cp ~/tools/weka-3-6-12/weka.jar:liblinear-1.51-with-deps.jar weka.classifiers.meta.FilteredClassifier -F "weka.filters.unsupervised.attribute.Remove -R 1" -W weka.classifiers.functions.LibLINEAR -t ID_Tendinte_vec_rank.arff -T ID_Tendinte_test_vec_rank.arff -i -k -p 1 -distribution -- -S 0 -C 1.0 -E 0.01 -B 1.0 -P > Tendinte_liblinear
```

Fig. 27

Obținând clasificarea pe categorii a fiecărei știri în formatul următor

7. Pentru a da o formă finală a rezultatelor vom folosi un program “perl” făcut special pentru acest tip de output.



Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

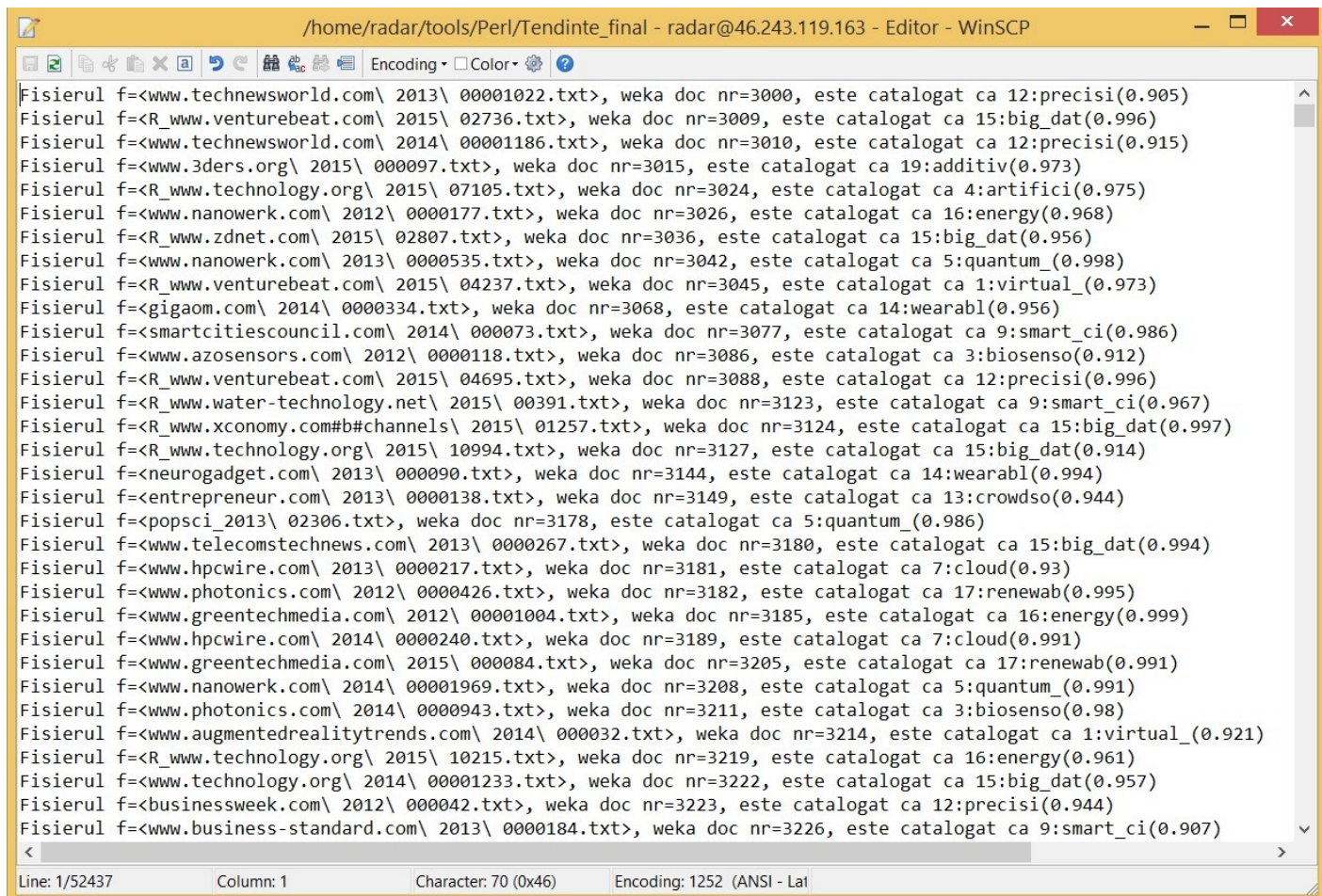
```

/home/radar/tools/weka-3-6-12/19Tendinte_part1_liblinear - radar@46.243.119.163 - Editor - WinSCP
Encoding - Color
=== Predictions on test data ===
inst# actual predicted error distribution (ID)
1 1:virtual_2:nuclear_ + 0.004,*0.617,0.007,0.002,0.002,0.006,0.002,0.002,0.002,0.002,0.001,0.003,0.271,0.034,0.015,0.004,0.01,0.012,0.004 (2958)
2 2:nuclear_19:aditiv + 0.001,0.016,0.001,0.001,0.002,0.001,0.008,0.001,0.001,0.002,0.002,0.181,0.001,0.006,*0.776 (2959)
3 3:biosenso13:crowdso + 0.028,0.151,0.043,0.011,0.133,0.066,0.019,0.029,0.026,0.021,0.003,0.015,*0.195,0.077,0.102,0.035,0.017,0.022,0.008 (2960)
4 4:artifici15:big_dat + 0.015,0.272,0.013,0.005,0.014,0.006,0.004,0.01,0.013,0.021,0.003,0.048,0.042,0.022,*0.435,0.03,0.033,0.003,0.012 (2961)
5 5:quantum_11:artific + 0.006,0.003,0.007,0.059,0.079,0.07,0.008,0.006,0.001,0.002,*0.498,0.01,0.006,0.006,0.227,0.007,0.001,0.006,0.001 (2962)
6 6:iot4:artifici + 0.002,0.001,0.013,*0.816,0.002,0.012,0.001,0.002,0.002,0.006,0.12,0.002,0.002,0.003,0.014 (2963)
7 7:'cloudc11:artific + 0.008,0.01,0.253,0.078,0.012,0.007,0.002,0.013,0.013,0.011,*0.521,0.012,0.035,0.006,0.006,0.005,0.005,0.001,0.002 (2964)
8 8:driverle1:virtual_ + *0.328,0.006,0.004,0.004,0.023,0.018,0.013,0.004,0.076,0.266,0.021,0.006,0.037,0.035,0.129,0.005,0.013,0.011,0.001 (2965)
9 9:smart_ci17:renewab + 0.007,0.005,0.01,0.006,0.01,0.133,0.004,0.02,0.04,0.047,0.003,0.059,0.078,0.027,0.18,0.006,*0.342,0.019,0.004 (2966)
10 10:'3d pri15:big_dat + 0.053,0.017,0.019,0.007,0.014,0.067,0.04,0.038,0.011,0.008,0.247,0.023,0.024,0.097,*0.287,0.019,0.006,0.017,0.005 (2967)
11 10:'3d pri17:renewab + 0.051,0.04,0.073,0.073,0.037,0.051,0.037,0.049,0.042,0.039,0.033,0.035,0.054,0.053,0.079,0.055,*0.099,0.038,0.06 (2968)
12 10:'3d pri13:crowdso + 0.004,0.001,0.001,0.002,0.001,0.001,0.348,0.018,0.001,0.001,0.002,*0.583,0.004,0.005,0.002,0.001,0.001,0.023 (2969)
13 10:'3d pri14:wearabl + 0.012,0.014,0.016,0.033,0.03,0.058,0.054,0.016,0.012,0.02,0.021,0.019,0.026,*0.561,0.024,0.006,0.031,0.019,0.027 (2970)
14 10:'3d pri7:cloudco + 0.035,0.019,0.026,0.017,0.029,0.03,*0.165,0.05,0.054,0.137,0.028,0.017,0.131,0.036,0.037,0.007,0.151,0.023,0.01 (2971)
15 10:'3d pri17:renewab + 0.003,0.002,0.003,0.002,0.001,0.007,0.003,0.001,0.012,0.003,0.001,0.001,0.001,0.002,0.005,0.062,*0.889,0.002,0.001 (2972)
16 10:'3d pri5:quantum_ + 0.001,0.01,0.003,0.082,*0.602,0.007,0.001,0.002,0.023,0.021,0.004,0.007,0.004,0.019,0.003,0.037,0.001 (2973)
17 10:'3d pri9:smart_ci + 0.04,0.017,0.016,0.013,0.018,0.008,0.014,0.01,*0.507,0.051,0.004,0.024,0.026,0.031,0.127,0.006,0.059,0.012,0.016 (2974)
18 10:'3d pri3:biosenso + 0,0.001,*0.672,0.267,0.003,0.001,0.001,0.001,0.002,0.038,0.002,0.004,0.005,0,0.001,0.002,0.002,0 (2975)
19 10:'3d pri9:smart_ci + 0.012,0.012,0.016,0.016,0.019,0.068,0.018,0.013,*0.272,0.089,0.015,0.028,0.02,0.083,0.014,0.124,0.102,0.016,0.063 (2976)
20 10:'3d pri15:big_dat + 0.008,0.009,0.053,0.014,0.007,0.006,0.097,0.012,0.191,0.064,0.033,0.027,0.041,0.117,*0.224,0.024,0.013,0.043,0.016 (2977)
21 10:'3d pri7:cloudco + 0.017,0.022,0.036,0.054,0.013,0.017,*0.569,0.04,0.019,0.023,0.013,0.017,0.015,0.041,0.035,0.007,0.008,0.025,0.03 (2978)
22 10:'3d pri14:wearabl + 0.026,0.014,0.013,0.007,0.011,0.018,0.011,0.026,0.002,0.011,0.01,0.012,0.012,*0.774,0.011,0.014,0.01,0.012,0.006 (2979)
23 10:'3d pri10:3dprin + 0.033,0.007,0.014,0.009,0.01,0.016,0.002,0.027,0.011,*0.457,0.009,0.005,0.003,0.007,0.019,0.008,0.005,0.008,0.35 (2980)
24 10:'3d pri19:aditiv + 0.091,0.019,0.03,0.015,0.045,0.055,0.031,0.017,0.006,0.001,0.034,0.011,0.123,0.06,0.117,0.019,0.059,0.018,*0.25 (2981)
25 10:'3d pri10:3dprin + 0.023,0.009,0.054,0.014,0.013,0.167,0.04,0.016,0.097,*0.249,0.018,0.126,0.006,0.016,0.032,0.064,0.004,0.052,0.001 (2982)
26 10:'3d pri4:artifici + 0.001,0.001,0.228,*0.728,0.001,0,0.015,0,0,0.009,0.003,0.008,0.001,0.002,0,0,0,0.001,0 (2983)
27 10:'3d pri1:virtual_ + *0.561,0.003,0.004,0.007,0.006,0.035,0.024,0.003,0.006,0.018,0.001,0.229,0.007,0.017,0.037,0.003,0.031,0.004,0.007 (2984)
28 10:'3d pri4:artifici + 0.043,0.017,0.034,*0.547,0.012,0.016,0.044,0.019,0.008,0.015,0.007,0.01,0.017,0.05,0.007,0.016,0.021,0.016,0.102 (2985)
29 10:'3d pri15:big_dat + 0.011,0.035,0.037,0.003,0.009,0.039,0.003,0.012,0.003,0.007,0.005,0.05,0.03,0.031,*0.368,0.273,0.005,0.076,0.004 (2986)
30 10:'3d pri14:wearabl + 0.001,0.001,0.099,0.015,0.001,0.022,0.017,0.003,0.001,0.005,0.004,0.006,0.007,*0.761,0.003,0.001,0.037,0.006,0.012 (2987)
31 10:'3d pri5:quantum_ + 0.003,0.004,0.003,0.01,*0.697,0.004,0.001,0.001,0.002,0.011,0.001,0.004,0.001,0,0.213,0.015,0.005,0.021,0.003 (2988)
32 10:'3d pri4:artifici + 0.017,0.011,0.091,*0.557,0.006,0.006,0.011,0.019,0.005,0.02,0.009,0.011,0.021,0.106,0.057,0.007,0.022,0.017,0.007 (2989)

```

Fig. 28

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013



```
/home/radar/tools/Perl/Tendinte_final - radar@46.243.119.163 - Editor - WinSCP
Encoding - Color
Fisierul f=<www.technewsworld.com\ 2013\ 00001022.txt>, weka doc nr=3000, este catalogat ca 12:precisi(0.905)
Fisierul f=<R_www.venturebeat.com\ 2015\ 02736.txt>, weka doc nr=3009, este catalogat ca 15:big_dat(0.996)
Fisierul f=<www.technewsworld.com\ 2014\ 00001186.txt>, weka doc nr=3010, este catalogat ca 12:precisi(0.915)
Fisierul f=<www.3ders.org\ 2015\ 000097.txt>, weka doc nr=3015, este catalogat ca 19:additiv(0.973)
Fisierul f=<R_www.technology.org\ 2015\ 07105.txt>, weka doc nr=3024, este catalogat ca 4:artifici(0.975)
Fisierul f=<www.nanowerk.com\ 2012\ 0000177.txt>, weka doc nr=3026, este catalogat ca 16:energy(0.968)
Fisierul f=<R_www.zdnet.com\ 2015\ 02807.txt>, weka doc nr=3036, este catalogat ca 15:big_dat(0.956)
Fisierul f=<www.nanowerk.com\ 2013\ 0000535.txt>, weka doc nr=3042, este catalogat ca 5:quantum_(0.998)
Fisierul f=<R_www.venturebeat.com\ 2015\ 04237.txt>, weka doc nr=3045, este catalogat ca 1:virtual_(0.973)
Fisierul f=<gigaom.com\ 2014\ 0000334.txt>, weka doc nr=3068, este catalogat ca 14:wearabl(0.956)
Fisierul f=<smartcitiescouncil.com\ 2014\ 000073.txt>, weka doc nr=3077, este catalogat ca 9:smart_ci(0.986)
Fisierul f=<www.azosensors.com\ 2012\ 0000118.txt>, weka doc nr=3086, este catalogat ca 3:biosenso(0.912)
Fisierul f=<R_www.venturebeat.com\ 2015\ 04695.txt>, weka doc nr=3088, este catalogat ca 12:precisi(0.996)
Fisierul f=<R_www.water-technology.net\ 2015\ 00391.txt>, weka doc nr=3123, este catalogat ca 9:smart_ci(0.967)
Fisierul f=<R_www.xonomy.com#b#channels\ 2015\ 01257.txt>, weka doc nr=3124, este catalogat ca 15:big_dat(0.997)
Fisierul f=<R_www.technology.org\ 2015\ 10994.txt>, weka doc nr=3127, este catalogat ca 15:big_dat(0.914)
Fisierul f=<neurogadget.com\ 2013\ 000090.txt>, weka doc nr=3144, este catalogat ca 14:wearabl(0.994)
Fisierul f=<entrepreneur.com\ 2013\ 0000138.txt>, weka doc nr=3149, este catalogat ca 13:crowdso(0.944)
Fisierul f=<popsci_2013\ 02306.txt>, weka doc nr=3178, este catalogat ca 5:quantum_(0.986)
Fisierul f=<www.telecomstechnews.com\ 2013\ 0000267.txt>, weka doc nr=3180, este catalogat ca 15:big_dat(0.994)
Fisierul f=<www.hpcwire.com\ 2013\ 0000217.txt>, weka doc nr=3181, este catalogat ca 7:cloud(0.93)
Fisierul f=<www.photonics.com\ 2012\ 0000426.txt>, weka doc nr=3182, este catalogat ca 17:renewab(0.995)
Fisierul f=<www.greentechmedia.com\ 2012\ 00001004.txt>, weka doc nr=3185, este catalogat ca 16:energy(0.999)
Fisierul f=<www.hpcwire.com\ 2014\ 0000240.txt>, weka doc nr=3189, este catalogat ca 7:cloud(0.991)
Fisierul f=<www.greentechmedia.com\ 2015\ 000084.txt>, weka doc nr=3205, este catalogat ca 17:renewab(0.991)
Fisierul f=<www.nanowerk.com\ 2014\ 00001969.txt>, weka doc nr=3208, este catalogat ca 5:quantum_(0.991)
Fisierul f=<www.photonics.com\ 2014\ 0000943.txt>, weka doc nr=3211, este catalogat ca 3:biosenso(0.98)
Fisierul f=<www.augmentedrealitytrends.com\ 2014\ 000032.txt>, weka doc nr=3214, este catalogat ca 1:virtual_(0.921)
Fisierul f=<R_www.technology.org\ 2015\ 10215.txt>, weka doc nr=3219, este catalogat ca 16:energy(0.961)
Fisierul f=<www.technology.org\ 2014\ 00001233.txt>, weka doc nr=3222, este catalogat ca 15:big_dat(0.957)
Fisierul f=<businessweek.com\ 2012\ 000042.txt>, weka doc nr=3223, este catalogat ca 12:precisi(0.944)
Fisierul f=<www.business-standard.com\ 2013\ 0000184.txt>, weka doc nr=3226, este catalogat ca 9:smart_ci(0.907)
Line: 1/52437 Column: 1 Character: 70 (0x46) Encoding: 1252 (ANSI - Lai)
```

Fig. 29

Fisierul **www.technewsworld.com 2013 00001022.txt** este catalogat ca 12:precision – prescurtare pentru categoria **precision_agriculture** cu o precizie de 90,5%



UNIUNEA EUROPEANĂ
Fondul Social European



GUVERNUL ROMÂNIEI
Ministerul Dezvoltării Regionale
și Administrației Publice



INOVAȚIE ÎN ADMINISTRAȚIE



Instrumente Structurale
2007-2013

Proiect cofinanțat din Fondul Social European, prin Programul Operațional "Dezvoltarea Capacității Administrative", în perioada 2007-2013

Listă Anexe:

Lista Anexe:

- Anexa 1_Repository_Tagy
- Anexa 2_Rezultate surse adiționale știri_similitudine
- Anexa 3_Rapoarte_PD_SS_iunie 2015.
- Anexa 4_Platforma Tagy
- Anexa 5_Raport general NS
- Anexa 6_Raport vot Platforme
- Anexa 7_Statistici si analize.actualizare
- Anexa 8 Domenii descriptori corpus_final
- Anexa 9 SmartDictionary_final
- Anexa 10.Dictionary tendinte

Listă Figuri:

- Fig.1.1 Lista finala platforme RSS feeder
- Fig. 1.2 Distributia lunara colectarii de stiri 2015
- Fig. 2. Captura clasificare platforme
- Fig.3. Agregatoare identificate în analiza surselor
- Fig.4 Agregatoare
- Fig.5 Vizualizare identificare surse stiri
- Fig 6.1 Raport vot platforme
- Fig 6.2 Raport vot platforme



UNITATEA EXECUTIVĂ PENTRU
FINANȚAREA ÎNVĂȚĂMÂNTULUI
SUPERIOR, A CERCETĂRII
DEZVOLTĂRII ȘI INOVĂRII



Creșterea Capacității Administrative
a Sistemului Public de CDI