

Donostia/San Sebastián 2016

“5.2.8 DSS2016 Behagunea - Hirikia kaia” web aplikazioa

Proiektuaren Memoria

Elhuyar Hizkuntza eta Teknologia

Proiektu arduraduna (Elhuyar): Iñaki San Vicente Roncal	DSS2016ko arduraduna (Hirikia kaia): Aitzol Astigarraga
--	--

2016eko maiatzaren 10a

Aurkibidea

1 Sarrera.....	4
2 Web aplikazioa – Soluzio teknologikoa.....	6
2.1 Sarrera.....	6
2.2 Sistemaren arkitektura.....	7
2.3 Erabiltzailearen interfazea.....	8
2.3.1 Pantaila nagusia, nabigazio orokorra.....	9
2.3.2 Estatistikak.....	11
2.3.3 Kudeatzailearen aukerak.....	13
2.4 Iritzi bilketa: Iturriak.....	15
2.5 Aipamen biltzailea: MSM.....	18
2.6 Iritzien analisia: prozesamendu kate eleaniztuna.....	20
Oinarrizko prozesamendu linguistikoa: Ixa-pipes.....	22
3 Software librea, lizentziak.....	24
4 Etorkizunerako Ondarea.....	27
5 Proiektuaren ebaluazioa.....	28

1 Sarrera

Donostia/San Sebastián 2016 Fundazioak (aurrerantzean, DSS2016EU) DSS2016ko ekitaldien inguruan herritarrek idatziko dituzten iritziak automatikoki aztertzen dituen web aplikazioa garatzeko Elhuyar Fundazioaren zerbitzuak kontratatu ditu.

Txosten honek garatutako plataforma digitalaren ezaugarriak eta egindako lanaren inguruko xehetasunak ematen ditu. Garatutako soluzioa ikusgai dago <http://behagunea.dss2016.eu/> helbidean. Orokorrean, “Behagunearen” helburua Donostia 2016 ekimenaren markoan antolaturiko ekitaldiek gizartean duten oihartzuna aztertzea da, ondorengo ezaugarri nagusiak dituelarik:

- Sistemak DSS2016 ekimenaren barruan ematen diren proiektuen inguruko iritzia jasotzen du, denbora errealean. Iritzi bat publikatu eta handik 10 minutura iritzi horiek linguistikoki prozesatuta eta polaritatea aztertuta “Behagunean” ikusgai daude.
- Aipamenak hainbat iturritatik eskuratzen ditu aplikazioak horretarako sortutako berriazko softwareari esker. Egunkariak, sare sozialak, eta blog-ak ustiatzeko gai da.
- Iritzien analisi aurreratua. Sistemak mezuen polaritatea detektatzen du, mezuetan aipatzen diren gaiak eta haien egileak identifikatzen ditu. DSS2016ren domeinura egokitutako iritziak aztertzeko teknologia garatu da, 4 hizkuntzatarako, Euskara barne.
- Datuen bisualizazio aurreratua: Sistemak DSS2016 ekimenaren barruan ematen diren proiektuen inguruko iritzia jasotzen du,

denboran zeharreko aldaketak erakusten ditu, egile aktiboenak identifikatzen ditu, eta jasotako datuen gaineko estatistikak sortzen ditu, datu guztiak modu erosoan kontsultatzeko bisualizazio interaktibo aurreratua eskaintzen dituelarik.

- Zaintza sistemaren kudeaketa aurreratua: Kudeatzaileak aukera du monitorizatu nahi diren gaiak zein iturriak konfiguratzeko, eta estatistika atal berezitua ere badu. Horretaz gain, jasotako iritzien eskuzko berrikuspena egiteko aukera ematen zaio.

Sistema irekien filosofiarekin guztiz bat datorren soluzio teknologikoa garatu da. Proiektuan zehar garatu diren teknologia guztiak software librea erabiliz egin dira, eta lizentzia libreekin eskuragarri daude sarean.

Txosten hau honela antolatzen da hemendik aurrera; hurrengo atalean Elhuyar Fundazioak garatutako soluzio teknologikoa deskribatzen da. 3. atalak 2016 ondorenerako proiektuan dituen jarraipen asmoak aurkezten dira. 4. atalak proiektuaren eta garatutako teknologiaren ebaluazioa egiten da.

2 Web aplikazioa – Soluzio teknologikoa

2.1 Sarrera

Aurreko atalean azaldu dugun bezala, proiektu honen helburua aplikazio bat garatzean datza, zeinaren bitartez automatikoki analizatuko diren Donostia 2016 ekimeneko ekitaldien inguruan herritarrek idazten dituzten mezuak, denbora errealean. Testu horiek linguistikoki analizatzen ditu, eta bertan azaltzen diren iritzien sailkapena burutzen du, zein gairi dagozkion identifikatuz eta iritzi horiek positibo, negatibo edo neutralak diren esanez.

Datu horiek webgunean publikoki ikusgai daude, modu erakargarrian, erabiltzaileei nabigazio eroso eta intuitiboa eskainiz eta datuen kontsulta aurreratua ahalbidetuz.

Teknologiari dagokionez, iritzi bat adierazten duen testua aztertzeko orduan, gai izan behar dugu, testuan adierazi diren elementu positiboak edo negatiboak identifikatzeko, eta iritzia zeri buruzkoa den detektatzeko, edo iritzia nork adierazi duen esateko. Horretarako, hizkuntzaren prozesamendurako teknologiak erabiltzea ezinbestekoa da, testuaren analisi egoki bat. Gure proposamenak Ixa-Pipes tresna linguistikoak erabiltzen ditu oinarrizko analisia burutzeko, EHUko IXA taldeak garatuak. Baina ez hori bakarrik, informazioaren erauzketa ere egin behar da, iritzien erauzketa. Horretarako erregela eta ikasketa automatikoan oinarritutako algoritmoak, artearen egoerakoen parekoak, erabiltzen ditu gure proposamenak.

Garatutako soluzioa, domeinura egokitutakoa da,, izan ere, emaitza domeinu kulturalera eta euskal herriko egoera soziolinguistikora

moldatutako sistema bat da, hots DSS2016ren beharretara egokitutako sistema.

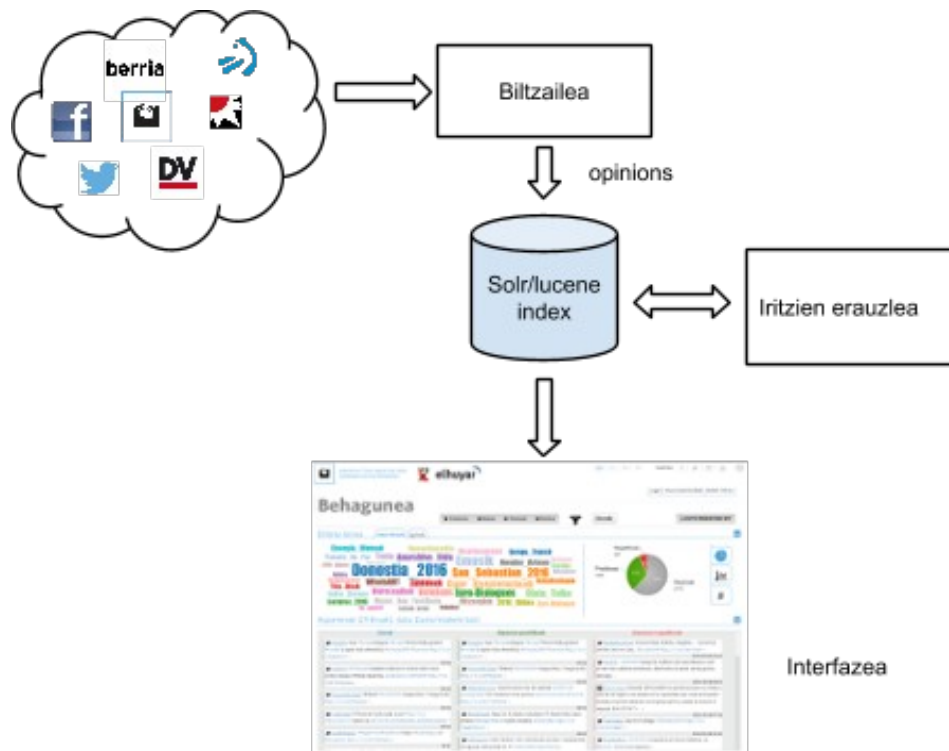
Ondorengo ataletan garatutako aplikazioaren ezaugarri zehatzak aurkezten dira, eta erabiliko den teknologiaren inguruko xehetasunak ematen dira.

2.2 Sistemaren arkitektura

Web aplikazioa MVC (Model–view–controller) arkitektura batean oinarrituta dago. Arkitektura honek datuak eta negozioaren logika erabiltzailearen interfaze eta gertaerak kudeatzeko moduluetatik bereizten ditu, etorkizuneko garapenak eta mantentze lanak erraztuz. Zehazki, Django WAFa (Web Application Framework) erabili da interfaze nagusia garatzeko eta datu-biltegiarekin komunikazioa inplementatzeko..

Sistema hiru modulu nagusiz osatua dago:

- Biltzailea (crawler-a): Iritzien bilketarako moduluak interneteko iturri ezberdinetatik jasoko ditu edukiak biltegi batean gordetzeko.
- Iritzien erauzlea: Prozesamendu kateak jasotako edukiak prozesatu eta bertako iritzien analisia egiten du. Erauzle honek prozesamendu linguistiko osoa burutzeko katea integratzen du, lematizazio, etiketatze linguistiko eta polaritate sailkatzailea barne hartzen dituelarik.
- Iritziak aztertzeke interfazea: Bildutako testuetan aurkitutako iritziak modu grafikoan kontsultatzeko funtzionalitateak eskaintzen ditu, bilketa moduluaren kudeaketa egitea ahalbidetzen du, eta sailkapen automatikoaren emaitzak berrikusteko aukera ematen du.



Irdia 2: Proposatutako sistemaren arkitektura, MVC patroian oinarritutakoa.

2.3 Erabiltzailearen interfazea

RWD (Responsive Web Design) filosofian oinarritutako web-interfazea garatu da, tresnaren funtzionalitatea modu intuitiboan erabiltzea ahalbidetzen duena, eta interfazea bera zenbait gailutara (tabletak, smartphone-ak, PCak...) egokitzen dela bermatzen duena. Interfazea eleaniztuna da 4 hizkuntzetan (Euskara, gaztelania, ingelesa eta frantsesa).

Ondoren erabiltzailearen interfazearen pantaila nagusiak eta eskaintzen dituen nabigazio aukerak azaltzen dira. Hurrengo ataletan agertuko diren bisualizazio taula eta horien gainean implementatutako elkarrekintza oro d3.js javascript liburutegiarekiko

2.3.1 Pantaila nagusia, nabigazio orokorra

Pantaila nagusiak DSS2016ko proiektuen uneko argazkia erakusten du, denbora errealean. 3 elementu nagusi daude pantaila honetan: goiburua filtro eta elkarrekintza aukerekin, proiektu zein egileen bisualizazioak, eta txioen erakusleihoa.

The screenshot shows the website interface for Donostia 2016. It is divided into three main sections:

- 1. Header and Navigation:** Includes the website logo, language options (EU, ES, EN, FR), a search bar, and filter buttons for 'Orokorra', 'Bakea', 'Ahotsak', 'Bizitza', and 'Denak'. There is also an 'AUTO-FRESKATZEA OFF' button.
- 2. Main Content Area:** Features a word cloud with terms like 'Donostia 2016', 'San Sebastian 2016', 'Emusik', and 'Hondar Artea'. To the right is a pie chart showing the distribution of tweets: Positiboak (1444, 32%), Negatiboak (287, 19%), and Neutroak (2778, 62%).
- 3. Tweet Feed:** A list of tweets under the heading 'Aipamenak ([Filtroak]: data: Duela hilabete bat)'. The feed is categorized into 'Denak', 'Aipamen positiboak', and 'Aipamen negatiboak'. Each tweet includes the user's name, profile picture, text, and a timestamp.

Irudia 3: Interfazearen pantaila nagusia

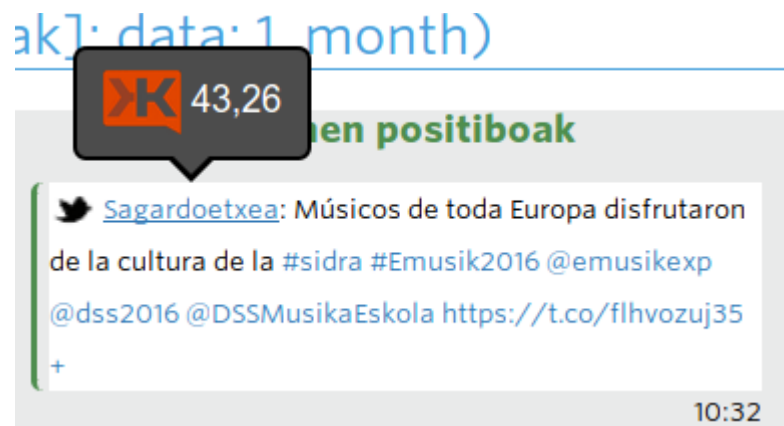
Elementu nagusia bisualizazioen atalak osatzen du. Bertan, gaien eta iritzi-emaien lainoak dira elementu garrantzitsuena. Lainoak gai esanguratsuenak eta iritzi-emaiak aktiboenak begiratu bakarrean ikusteko aukera ematen du. Modu lehenetsiak gai edo proiektuen lainoa ikusten da, egileen lainoa ikusteko dagokion fitxa aukeratu behar dugularik. Lainoen eskuinaldean, iritzien laburpen atala dugu ikusgai. Bertan, iritzi 3 informazio ditugu eskuragarri: (i) iritzien positibo eta negatiboen oreka

sektore-diagrama baten bitartez, (ii) iritzien eboluzioa denboran zehar erakusten duen grafikoa eta (iii) jasotako iritzi kopuruen inguruko estatistikak.



Irudia 4: Bisualizazioak

Aipamenen erakusleihoan txioak 3 zutabetan agertzen dira. Lehen zutabeaz aipamenak erakusten dira, eta hurrengoetan positibo ala negatibo sailkatutako aipamenak. Iritzi bakoitzaren iturria, iturri horren “ospea”, data eta egilea ere erakutsiko dira (ikus irudia 5).



Irudia 5: Aipamen baten ezaugarriak

Elkarrekintza bideratzeko, hainbat filtro eskaintzen dira, goiburuko menuan eskuragarri. Filtroek ondorengo aukerak eskaintzen dituzte:

- Epe ezberdinen arteko iritziak bakarrik ikusi.
- Hizkuntza jakin bateko iritziak ikusi.

- “Ospe” maila batetik gorako egileen iritziak bakarrik ikusi.
- Iturri mota jakin bateko iritziak bakarrik ikusi (prentsa, sare-sozialak).
- Itsasargi zehatz bateko proiektuen inguruko aipamenak ikusi.

Filtro hauetako bat edo gehiago aukeratutakoan pantailako elementu guztien informazioa filtroen aukeraketarekin arabera eguneratzen da, koherentzia mantenduz. Filtro horietaz gain, bai bisualizazioek eta baita aipamenek ere elkarrekintza eskaintzen dute:

- Zerrendako iritzi bakoitza interaktiboa izango da, iritziaren egilea klikagarria izanik. Egilearen gainean klikatuz gero pertsona horrek emandako iritzi guztiak bistaratuko dira. Aipamen bat bere jatorrizko helbidean ikusteko aukera ere badugu “+” ikurra sakatuz.
- Egile eta gaien lainoetako elementuak ere interaktiboak dira, eta filtro modura funtzionatzen dute, hots, elementu bat klikatzerakoan, interfazea freskatu egiten da, elementu horrekin lotutako informazioa erakutsiz.

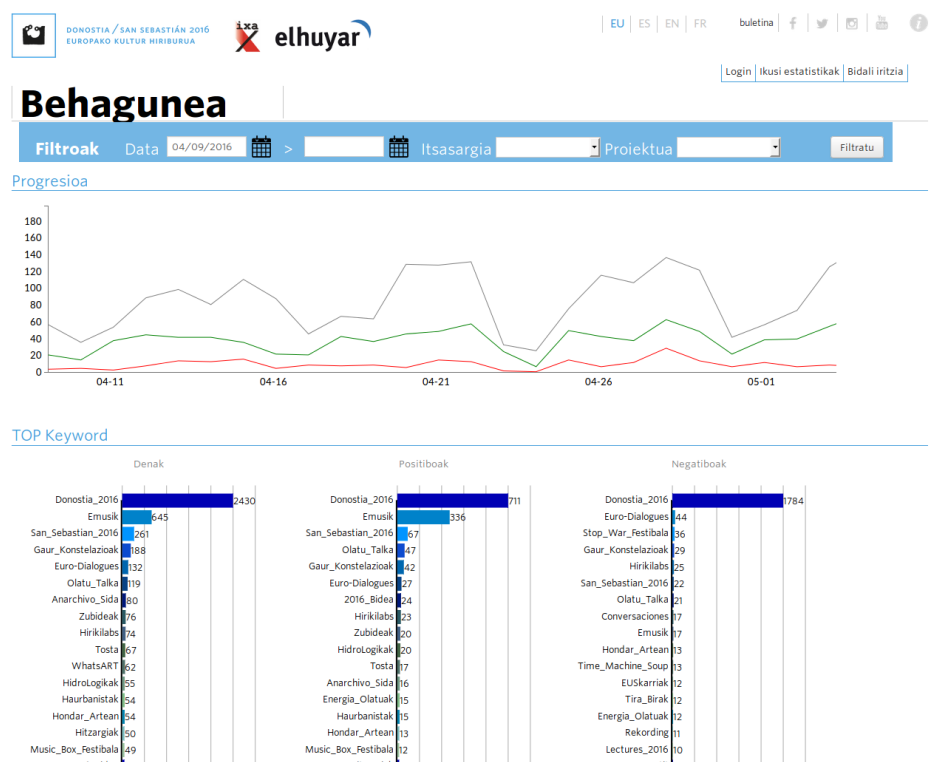
Azkenik, beste bi elementu aipatuko ditugu. Batetik freskatze automatikoa ahalbidetzeko botoia. Honekin pantailako informazioa 15 minutuero automatikoki eguneratzen da. Horrela “Behagunea” pantaila nagusi batean eduki liteke eta informazio eguneratua ikusi uneoro. Bestetik “Iritzia bidali” aukeraren bidez edozein erabiltzailek bere iritzia emateko aukera du, zuzenean plataformaren bitartez.

2.3.2 Estatistikak

Goiburuko menuan “Ikusi estatistikak” klikatuta, estatistika orokorren atala zabalduko dugu. Bertan iritzien eboluzioaren grafikoa aurkituko

dugu lehenengo gauza, eta ondoren, gai aipatuenen, eta, prentsako iturrien eta Twitterreko egileen rankingak ikus daitezke. Bloke bakoitzean 3 ranking erakusten dira. Lehenengoan elementu aipatuena eta bakoitzak jasotako aipamen kopurua agertzen da. Erdiko rankingak aipamen positiboak dagokien informazioa ematen du, berriz ere elementu aipatuena eta jaso dituzten aipamen positibo kopurua, eta azken zutabea aipamen negatiboak dagozkienak. Egileei dagozkien rankingen kasuan, egile bakoitzak botatutako aipamen kopuruak erakusten dituzte rankingak.

Pantaila nagusiaren antzera, informazioa filtratzeko aukerak ere eskaintzen dira. Kasu honetan, data zehatzen arteko bilaketak egin ahal izango ditugu, proiektuaren edo itsasargien arabera.



Irudia 6: Estatistika orokorrak, iritzien denboran zeharreko eboluzioaren grafikoa eta gai aipatuena rankingak ikusten dira.

2.3.3 Kudeatzailearen aukerak

Behaguneak kudeatzaile rol bat ere implementatzen du. Kudeatzaile horrek funtzionalitate aurreratuak ditu eskura, batetik, datu-bilketaren ezarpenak aldatu ahal izateko, eta, bestetik, jasotako datuen gaineko kontrola eramateko. Ondoren kasu bakoitzean ematen diren aukerak azaltzen ditugu.

- **Gakoen kudeaketa:** aipamenak aurretik definitutako gako-hitz zerrenda baten arabera biltzen dira. Kudeatzaileak zerrenda hori ikusi eta editatzeko aukera du.

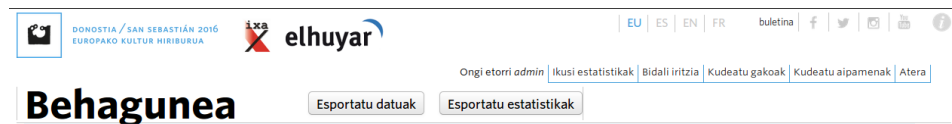
The screenshot shows the 'Behagunea' administration interface. At the top, there are navigation links for EU, ES, EN, FR, and a 'buletina' button. Below that, there are links for 'Ongi etorri admin', 'Ikusi estatistikak', 'Bidali iritzia', 'Kudeatu gakoak', 'Kudeatu aipamenak', and 'Atera'. The main heading is 'Behagunea' with a 'Sortu gako berria' button. Below the heading, there is a search bar and a 'Show 125 entries' dropdown. The main content is a table with the following columns: Mota, Hizkuntza, Kategoria, Azpi-kategoria, Terminoa, and Pantaila-etiketa. The table lists various entries with their respective languages, categories, sub-categories, terms, and labels. At the bottom, there is a pagination control showing 'Showing 1 to 25 of 695 entries' and a 'Previous' button with a list of page numbers (1, 2, 3, 4, 5, ..., 28, Next).

Mota	Hizkuntza	Kategoria	Azpi-kategoria	Terminoa	Pantaila-etiketa	
Press	eu	general	main	Donostia_2016	Donostia_2016	Editatu 13
Press	eu	general	main	D2016	Donostia_2016	Editatu 14
Press	eu	general	main	D552016	Donostia_2016	Editatu 15
Press	es	general	main	Donostia_2016	Donostia_2016	Editatu 16
Press	es	general	main	D2016	Donostia_2016	Editatu 17
Press	es	general	main	D552016	Donostia_2016	Editatu 18
Press	es	general	main	San_Sebastian_2016	Donostia_2016	Editatu 19
Press	en	general	main	San_Sebastian_2016	San_Sebastian_2016	Editatu 20
Press	en	general	main	Donostia_2016	San_Sebastian_2016	Editatu 21
Press	en	general	main	D2016	San_Sebastian_2016	Editatu 22
Press	en	general	main	D552016	San_Sebastian_2016	Editatu 23
Press	fr	general	main	Saint_Sebastien_2016	Saint_Sebastien_2016	Editatu 24
Press	fr	general	main	San_Sebastian_2016	Saint_Sebastien_2016	Editatu 25
Press	fr	general	main	Donostia_2016	Saint_Sebastien_2016	Editatu 26
Press	fr	general	main	D2016	Saint_Sebastien_2016	Editatu 27
Press	fr	general	main	D552016	Saint_Sebastien_2016	Editatu 28
Press	eu	general	AGENDA	2016_AGENDA	2016_AGENDA	Editatu 33
Press	es	general	AGENDA	2016_AGENDA	2016_AGENDA	Editatu 34
Press	en	general	AGENDA	2016_CALENDAR	2016_CALENDAR	Editatu 35
Press	fr	general	AGENDA	Calendrier_2016	Calendrier_2016	Editatu 36
Press	eu	general	ESPACIO 2016	2016_Gunea	2016_Gunea	Editatu 42
Press	es	general	ESPACIO 2016	Espacio_2016	Espacio_2016	Editatu 43
Press	es	general	ESPACIO 2016	space_2016	Espacio_2016	Editatu 44
Press	en	general	ESPACIO 2016	2016_Space	2016_Space	Editatu 45
Press	fr	general	ESPACIO 2016	Espace_2016	Espace_2016	Editatu 46

Irudia 7: Gakoak kudeatzeko interfazea

Gako bakoitza kategoria eta proiektu batekin lotuta dago. Horrela, posible da proiektu edo ekimen berriak jarraitzen hasia, dagoeneko jarraitzen ari garen ekimenei gako-hitz berriak gehitzea edo daudenak kentzea.

- **Estatistika aurreratuak:** estatistiken atalean kudeatzaileak bi aukera gehigarri ditu eskuragarri, erabiltzaile arruntekin alderatuta: “Esportatu datuak” eta “Esportatu estatistikak”. “Esportatu datuak” aukerak



Irudia 9: Kudeatzailearen aukera osagarriak estatistiken atalean

- **Zuzenketa:** Iritziaren azterketa automatikoa izanik, demagun erabiltzaile gisa oker sailkatutako iritzi bat topatzen dugula. Kudeatzaileak aukera du jasotako iritziak edozein unetan errepasatu, eta hala beharrezko ikusten badu horiek polaritatea zuzentzeko klik soil bat eginez. Aipamenak ezabatzeko aukera ere badu kudeatzaileak, adibidez, gaikoa ez den aipamen bat antzematen bada, edo inolako iritzirik adierazten ez duten spam mezuak detektatzen dituenean. Honek iritzi erauzketa hobetzen lagunduko digu epe ertainera, sistemak bere erroreetatik ikasi ahal izango baitu.

ID	Testua	Polaritatea	Data	Iturria	Hizk
103010	Aurtengo edizioa Impact Hub kolektiboak eta Donostia 2016k elkarlanean antolatatu dute. [Donostia_2016] ESTEKA	<input checked="" type="checkbox"/> P <input type="checkbox"/> N <input type="checkbox"/> Neu [P]	2016-05-07 15:59		eu
103011	Easo plazako antzinako baina publikoetatik irten da Jane-en ibilaldia ren bosgarren edizioa. Amaren historia beste modu batera ezguzteko aukera izan dute bertaratuak. Izan ere, Jane Jacobs pentsalari eta urbanista sozialean oinarrituta, boluntarioek gidatutako ibilaldia egin dute. Hain zuzen ere, Amarako eta Impact Hub-eko arkitekto Victor Aspe izan da ibilaldian azalpenak eman dituen. Halere, bertaratu diren auzotarrek ere euren istorioak kontatu dituzte. Aurtun, Impact Hub kolektiboaz gain, Donostia 2016k ere antolatatu du jarduera. daren ardua eta kudeaketa Impact Hub kolektiboarena eta Donostia 2016rena da. [Donostia_2016] ESTEKA	<input checked="" type="checkbox"/> P <input type="checkbox"/> N <input type="checkbox"/> Neu [P]	2016-05-07 15:59		eu
103012	Aurtengo ibilaldia Donostia 2016rekin elkarlanean egin dutenez, paraleleki beste jarduera batzuk egongo dira. [Donostia_2016] ESTEKA	<input type="checkbox"/> P <input type="checkbox"/> N <input checked="" type="checkbox"/> Neu [N]	2016-05-07 15:59		eu
103013	Erihorabuen a la @EMNLAlun, Habéis estado a la altura de la Capitalidad @DSS2016. Sois #orgullosedadad. https://t.co/Tenw35540 [DSS2016] ESTEKA	<input checked="" type="checkbox"/> P <input type="checkbox"/> N <input type="checkbox"/> Neu [P]	2016-05-07 16:10		es
103015	Imágenes del día de hoy en el #emusik2016 en @vitoragasteiz https://t.co/9ar4caXfM [Emusik2016] ESTEKA	<input type="checkbox"/> P <input type="checkbox"/> N <input checked="" type="checkbox"/> Neu [NEU]	2016-05-07 16:12		es
103016	Boteretean denboraren jabetza funtzioa da. 3 alderdiak orekatu behar dira : politika, familia eta pertsonala' E. Epelde #AnarchivSida [AnarchivSida] ESTEKA	<input checked="" type="checkbox"/> P <input type="checkbox"/> N <input type="checkbox"/> Neu [N]	2016-05-07 16:13		eu
103022	#Donostia2016 Cucurucho de bambú con picos y jamón de bellota! https://t.co/4UDtpYgkN [donostia2016] ESTEKA	<input checked="" type="checkbox"/> P <input type="checkbox"/> N <input type="checkbox"/> Neu [P]	2016-05-07 16:26		es
103024	Palladio tocado por la Jove Orquesta de Banyoles. Me encanta esta canción #emusik2016 #orquesta #escuelamusica https://t.co/HwG3M7N [Emusik2016] ESTEKA	<input checked="" type="checkbox"/> P <input type="checkbox"/> N <input type="checkbox"/> Neu [P]	2016-05-07 16:27		es
103025	Today enjoy the concerts of music students from all around Europe! #Emusik2016 [Emusik2016] ESTEKA	<input checked="" type="checkbox"/> P <input type="checkbox"/> N <input type="checkbox"/> Neu [P]	2016-05-07 16:31		en
103028	Hoy, sábado 7 de mayo, a las 20:00hrs clausura del Festival #Emusik2016 en Sagües: Actuación de una Big Band... https://t.co/TawU35Nhw [Emusik2016] ESTEKA	<input type="checkbox"/> P <input type="checkbox"/> N <input checked="" type="checkbox"/> Neu [NEU]	2016-05-07 16:40		es

Irudia 9: aipamenak kudeatzeko atala. Aipamenak bere polaritatearen arabeko koloreaz erakusten dira. Aipamen baten polaritatea zuzentzeko polaritate zuzena markatu besterik ez dago. Zuzendutakoak urdinez nabarmentzen dira.

2.4 Iritzi bilteta: Iturriak

Behatokian iritziak biltzeko hautatutako edukiak DSS2016Eurekin batera landu dira. Ondoren jarraitzen diren iturrien zerren da ematen da:

- Sare sozialak:
 - Twitter
 - Instagram eta Facebook biltzea ere aztertu da, baina lehena baztertu egin da testuzko eduki minimoa izanik, iritziak jasotzeko zailtasunak planteatzen dituelako. Facebook-i dagokionez, erabiltzaileek mezu positiboak botatzeko joera nabarmena dute, zenbait ikerketek erakutsi duten bezala, mezu positiboek negatiboek baino

eragin handiagoa dutelako, jarraitzaile gehiago lortu araziz. Ondorioz Facebook ez erabiltzea erabaki da.

- **Prensa digitala:** Ondorengo taulan zehazten dira biltzen diren hedabide digitalak:

BERTAKOAK	
Hedabidea	Hizkuntza
DEIA	es,eu
Gara	eu,es,fr
Hitza	
<i>Busturialdeko Hitza</i>	eu
<i>Goierriko Hitza</i>	eu
<i>Irutxuloko Hitza</i>	eu
<i>Lea-Artibai eta Mutrikuko Hitza</i>	eu
<i>Oarso Bidasoko Hitza</i>	eu
<i>Bidasoko Hitza</i>	eu
<i>Oarsoaldeko Hitza</i>	eu
<i>Tolosaldeko Ataria</i>	eu
<i>Urola Kostako Hitza</i>	eu
Argia	eu
EL CORREO	es,eu
Info 7	es,eu,fr
Diario Vasco	es,eu
Radio Vitoria	es
Artezblai	es
Cuatro	es
Onda Vasca	es
Berria	eu
Goiena	eu
Noticias de Gipuzkoa	es,eu

AGENCIA EFE	es
EL PAIS (CULTURA)	es
EL PAIS SEMANAL	es
TV3	es,eu
Radio Donosti	es
Europapress	es
Eitb	es,eu
EUSKADI IRRATIA	eu
EL MUNDO	es
Teledonosti	es,eu
Hamaika tv	es,eu
Gaztezulo	eu
Naiz	es,eu,fr
RTVE	es
ABC	es
DIARIO DE NAVARRA	es
Radio France	fr
France Bleu	fr
El Periódico de Catalunya	es
La Vanguardia	es
NAZIOARTEKOAK	
Ouest-France	fr
Reuters Breaking Views	fr
Le Monde	fr
Liberation	fr
France 24	fr
Le soir	fr
Financial Times	en
Le Figaro	eu
The Times	en
L'Express	fr
The Economist	en
CNN	en
Arte tv	fr

2.5 Aipamen biltzailea: MSM

Iritzien behatoki bat martxan jartzeko lehen pausua sarean sakabanatuta daude iritziak atzeman eta horiek biltzea da. Baina ez hori bakarrik, gure kasuan prozesu hori iritziak argitaratu eta ahalik epe laburrenean burutzea da helburua, segundo edo minututako atzerapenaz.

Azkartasun hori ez da iturri guztien kasuan berdina. Sare sozialen kasuan edozein unetan sor daitezke mezu berriak, eta beraz etengabe “entzuten” egotea eskatzen dute, eduki berriak noiz agertuko adi. Aldiz prentsa digitalean iritzi artikulua sortzeko maiztasuna askoz ere baxuagoa da, eta ondorioz nahikoa izan daiteke egunean behin eduki berririk dagoen aztertzearekin. Gauza bera gertatuko litzateke blogekin, izan ere ezohikoa da astean 2-3 artikulua baino gehiago idaztea.

Aurre egin beharreko beste arazo bat edukia atzitzeko modua da. Sare sozialak, Twitter edo Facebook kasu, enpresa pribatuak dira, bere edukia API bidez jartzen dute eskuragarri, eta beraz, sare bakoitzarentzat ad-hoc prestatutako bilketa tresnak garatu behar dira, haiek jarritako baldintzak betez. Blog eta prentsaren kasuan, edukiak RSS bidez banatu izan dira. Azken aldian, zenbait komunikabide, Twitter bidez zabaltzen hasi dira haien eduki guztiak.

Aurrekoak kontutan izanik, “DSS2016 Behatokian” bi estrategia baliatzen dituen sistema garatu da: Media Source Monitor (MSM). MSM Javaz inplementatu da, eta hainbat funtzionalitate eskaintzen ditu, aipameneren bilketa, egileen eragin mailaren kalkulua, eta datu-biltegiarekin komunikazio barne:

- **Twitter bezeroa:** Public Streaming API delakoarekiko lotura sortzen du, APIra konektatzerakoan informazio geografikoaren, gako-hitzen, erabiltzaileen eta hizkuntzaren arabeko filtroak ezartzeko aukera ematen duelarik. Konexioa irekita mantentzen du etengabe, edukia sortu ahala jasotzen delarik. Konexioa eteten bada konexioa berrabiarazten du automatikoki.
- **RSS biltzailea:** Prentsa digitaleko eta blogetako edukien bilketa burutzeko RSS jarioetan oinarritutako biltzailea inplementatzen du MSM-k. RSSak aldizka irakurri behar dira eduki berriak iritsi ote den aztertzeko. Erabiltzailearen esku uzten da maiztasuna erabakitzea.
- **Egileen ospea:** Egile zerrenda bat emanda, MSM-k egile bakoitzaren “eragin maila” edo “ospea” kalkulatu du. Egile motaren arabera indize ezberdinak erabiltzen dira:
 - Klout: Twitter erabiltzaileen ospea neurtzeko Klout¹ API-a erabiltzen da.
 - Ahrefs: Web domeinuen ospea neurtzeko Ahrefs² indizea erabiltzen da, DomainStats.io API³ aren bitartez.

MSM-k konfigurazio fitxategi baten bitartez funtzionatu dezake, eta emaitzak irteera estandarrera bideratu, baina Behaguneako datu-biltegiarekin komunikazio zuzena ere badu. Horrela, parametro guztiak datu-biltegitik irakurtzeko gai da, eta jasotako edukiak datu-biltegitara bideratzeko ahalmena du.

¹ <https://klout.com/s/developers/v2>

² <https://ahrefs.com/>

³ <http://api.domainstats.io/>

2.6 Iritzien analisia: prozesamendu kate eleaniztuna

Jasotako mezuetan dauden iritzien analisia egiteko, Elhuyar Fundazioak garatutako iritzien analisirako softwarea erabili da, EliXa (San Vicente et al., 2015) izenekoa. Softwarea hau Javaz inplementatuta dago eta egungo bertsioa GPL lizentziarekin banatuko da. Software hau OpeNER proiektu Europarrean (FP7) eta Berztek proiektuan (Eortek IE12-333) garatutako teknologiaz baliatuko da.

Dataset	Aipamenak	P	N	NEU
eu	2937	931	408	1598
es	4754	1487	1303	1964
en	12,273	4,654	1,837	5,782
fr	11,071	3,459	2,618	4,994

Taula 1: Osatutako bildumen inguruko estatistikak.

EliXak ikasketa automatikoa erabiltzen du polaritatea sailkatzeko, polaritate lexikoen ezagutzarekin konbinatuta. Hurbilpen gainbegiratu dugu hau. Izan ere, sistema aurretik sailkatuta dauden adibideetatik eredu bat ikasten du, eta ondoren testu berria datozenean eredu horren arabera sailkatzen ditu. Hurbilpen hau erabiltzeko, beharrezkoa da ereduak ikastek sailkatuta dauden adibideak edukitzea. EliXaren arkitekturak ereduak ikasteko orduan hainbat ezaugarri gehitu edo kentzeko aukera ematen du, eta baita ikasitako ereduak ebaluatzeko ere. Horrela, uneko domeinura optimizatutako ereduak sortzea erraza eta eraginkorra bihurtzen da.

EliXak Euskarri bektoredun makina (ingelesez Support Vector Machine, edo SVM) motako algoritmoa erabiltzen du sailkapena burutzeko eta polaritate lexikoetatik erauzitako informazioa ezaugarri linguistikoekin konbinatzen du. Erabiltzen dituen ezaugarri nagusiak hauek dira:

- Polaritate lexikoaren informazioa
- Informazio linguistikoa (adibidez, kategoria gramatikalak, puntuazio markak, ...)
- Egitura sintaktikoak (e.g., “egokia da” vs. “ez da egokia”)
- Informazio estatistikoa (adibidez, testu luzera, hitz positibo/negatibo kopurua, ...)
- Sare sozialekin lotutako informazioa (adibidez, emotikonoak, letra larrien kopurua, ...)
- Hitzen lema n-gramak
- PoS ngram-ak

Behagunea proiektuan, EliXaren hainbat egokitzapen burutu dira. Batetik, euskararako eta frantseserako polaritate sailkatzaileak garatu dira, horrek barne hartzen dituen lanekin: adibide bildumak sortu eta polaritateaz anotatu eta polaritate lexikoak sortu hizkuntza horietan. Bestetik, domeinura egokitutako ereduak sortu dira, horiek sortzeko ezaugarri egokienak zein diren ikertu dugularik.

DSS2016ren aplikazio eremura egokitutako adibideak edo corpusak bildu dira (2015eko urritik 2016ko martxoa bitarte), eta horien polaritatea eskuz etiketatu da, ikasketa automatikoko sistemak entrenatu ahal izateko. Taula 1-ek sortutako datu-bildumen inguruko informazioa ematen du. Ingelesezko (en) eta frantsesezko (fr) bildumen kasuan DSS2016ko informazio nahikorik ez zegoela eta domeinu orokorreko bildumak ere erabili dira adibidea osagarriak lortzeko.

Hizkuntza	doitasuna	fpos	fneg	fneu
eu	72.66	0.646	0.636	0.789
es	70.78	0.669	0.725	0.724
en	68.0946	0.685	0.525	0.722
fr	64.33	0.588	0.605	0.698

Taula 2: Entrenatutako ereduen doitasuna. Azken hiru zutabeek kategoria bakoitzeko lortzen den f1-score balioa ematen dute.

Sortutako sailkatzaileek hizkuntzaren arabera doitasun ezberdinak dituzte. Taula 2 ematen dira sailkatzaile ezberdinen dituzten doitasun mailak.

Oinarrizko prozesamendu linguistikoa: Ixa-pipes

EliXak polaritatearen azterketa egiteko testuak aurretik linguistikoki etiketatu behar ditu, hitzen lema, kategoria gramatikalak, entitateak eta abar identifikatuta eduki. EliXa aurretik etiketatutako testuen iritzia erazteko gai da, baina oinarrizko prozesamendu linguistikoa ere integratua du bere baitan, landu gabeko testu gordina aztertu ahal izateko.

Horretarako EHUko IXA taldean garatutako IXA-pipes⁴ (Agerri et al., 2014) tresna erabiltzen da. Ixa-pipes erabiltzeko prest datorren hizkuntza naturalaren prozesamendurako tresna multzoa da.

Behagunea proiektuaren aurretik, tokenizazioa (ixa-pipe-tok), lematizazioa eta analizatzaile morfologikoa (ixa-pip-pos), analizatzaile sintaktikoa (ixa-pipe-parse) eta entitate ezagutzailea (ixa-pipe-nerc) eskaintzen zituen gaztelania eta ingeleserako, eta partzialki galiziera ere bai. Gainera, Ixa-pipesen garapena OpeNER⁵ proiektuari lotuta dago, eta bertan garatutako prozesamendu katearekin bateragarria da. Horri esker,

⁴ <http://ixa2.si.ehu.es/ixa-pipes/>

⁵ <http://www.opener-project.eu/>

frantsesa, nederlandera, italiera eta partzialki alemana ere prozesatzeko aukera ematen du. “DSS2016 Behagunea” proiektuaren baitan, euskararako eta frantseserako tokenizataile eta analizataile linguistikoak garatu dira.

3 Software librea, lizentziak

Gure proposamenak software librearen filosofiarekin erabat bat egiten du. 2. atalean aurkeztutako software oro lizentzia libreekin banatutako teknologien gainean garatu da. 3 Taulak garapenean erabilitako osagai nagusien lizentziak eta proiektu honetan duten erabilera laburtzen du.

Softwarea	Lizentzia	“D2016 behatokia” Proiektuan duen eginkizuna	Eskuragarri non
Django Python Web framework	BSD	Interfazearen oinarria	https://www.djangoproject.com/
IXA-pipes	Apache 2.0	Testuen prozesamendu linguistikoa egiteko liburutegiak	http://ixaz.si.ehu.es/ixa-pipes/
Weka	GNU GPL	Ikasketa automatikoa burutzeko liburutegia	http://www.cs.waikato.ac.nz/ml/weka/
EliXa	GNU GPL	Polaritatearen azterketarako modulua.	https://github.com/Elhuyar/EliXa
Apache Solr	Apache 2.0	Datuen indexatzea, bilaketa motorra.	http://lucene.apache.org/solr/
D3.js	BSD	Bisualizaziorako erabilitako Javascript liburutegia	http://d3js.org/

Taula 3: Garapenean erabilitako software nagusiak eta beren lizentziak.

Garatutako softwarearen lizentziari dagokionez, 2 atalean, sistemaren arkitekturan deskribatzean aipatutako hiru moduluak modu independentean banatzen dira, bakoitza bere lizentziarekin, baina, beti ere, DSS2016EUren asmoekin bat eginez, lizentzia horiek libreak izanik. Horrenbestez ondorengo lizentziak proposatzen ditugu:

- Interfazea, Django-ren izpirituari jarraiki, BSD lizentziarekin banatuko da.
- Datuen bilketarako modula (MSM) GNU GPL (V3) lizentziarekin banatzen da.
- Azkenik, EliXa polaritatearen analisirako moduluak GNU GPL (V3) lizentzia du. Izan ere EliXaren egungo bertsioak GPL lizentziadun liburutegiak erabiltzen ditu, eta honek garatutako softwarea ere lizentzia horrekin banatzera behartzen du.

Garatutako software oro biltegi publikoetan eskuragarri dago jada, ondorengo helbideetan:

- Interfazea: <https://github.com/Elhuyar/BehaguneUI>
- Datuen bilketa egiteko modula: <https://github.com/Elhuyar/MSM>
- EliXa Polaritate modula: <https://github.com/Elhuyar/EliXa>

Softwarea publiko egitearen helburua ez da soilik edonoren esku uztea, baizik eta etorkizunera begira, software hori hobetu eta erabiliko duen komunitate bat sortzea.

Software libreak zenbait abantaila sozial dakartza berekin, batez ez, bere erabilera eta banaketari dagokionean. Horrela, software libreak hurrengo lau askatasunak ematen dizkio erabiltzaileari:

- Programa nahi duen eran exekutatzeko askatasuna.

- Programaren iturburu kodea aztertu eta nahi duen eran aldatzeko askatasuna.
- Programaren kopia zehatzak banatzeko eskubidea.
- Aldatutako bertsioen kopiak banatzeko askatasuna.

Lau askatasun horiek oinarritzat hartuz erabiltzaileak norberaren konputazioaren kontrola bereganatu dezakete. Programaren banaketa eta erabilera ahalbidetzen duen sistema soziala bidezkoa da, eta, softwarea, banaketa eta erabilerari dagokionean, etikoa da.

4 Etorkizunerako Ondarea

“D2016 Behagunea” proiektuan garatutako software oro lizentzia libre baten bidez eskuragarri dagoela aipatu dugu. Proiektuaren baitan sortutako den aplikazioa edonork jarri dezake martxan, dela enpresa, erakunde, edo gizabanakoa.

Behatokiaren aplikazioaren oinarriko osagaiak NLP moduluak dira, testuaren analisi linguistiko automatikoak gauzatzen dutenak, alegia. Modulu horiek egun eskuragarri daude lizentzia librean, eta beren garapena etengabe ari da gauzatzen. Behagunean tresnak D2016 domeinura egokitzeko egindako lan guztia komunitatearen esku dago jada.

Elhuyarren asmoa da, proiektu honetan garatzen den teknologia beste eszenatoki batzuetara zabaltzea. Zentzu horretan, kulturarekin lotutako ekintzak dira aplikazio eremu nagusietako bat, adibidez, Zinemaldia, Behobia-Donostia lasterketa, edo Tabakalera proiektuaren inguruko jarrera. Orokorrean, herritarren iritziak bideratzeko tresna izango da Behatokia, erakundeek ere herritarrengana hurbiltzeko baliatu dezaketena.

Hizkuntza-teknologiaren ikuspuntutik, Euskararen garapen digital eta teknologikorako urrats berri bat eman da, eta euskararako analizatzaile linguistikoak jarri ditu euskal gizartearen eskura.

5 Proiektuaren ebaluazioa

Memoria hau idazteko garaian, “DSS2016 Behagunea” ez da oraindik denbora nahikoa igaro garatutako soluzioa publiko denetik, eta beraz hemen aurkezten den ebaluazioa mugatua da. Gure ebaluazioa 3 ardatz nagusiren inguruan burutuko da: garatutako softwarea, Donostia 2016 proiektuaren martxa neurtzeko izan duen erabilgarritasuna, eta herritarren parte-hartze dinamikak bultzatzeko garaian izan duen eragina.

Garatutako softwareari dagokionez, sailkatzaileen ebaluazioa txosten honen bigarren atalean ematen da. Garatutako tresnen inguruko dokumentazioa amaitzen ari gara memoria hau idazteko unean, baina etiketatzaile linguistikoaren inguruko dokumentazioa bakoitzari dagokion jatorrizko helbideetan aurki daiteke. Garatutako aplikazioa instalatzeko eta tresna guztien integrazioarako dokumentazioa memoria hau osatzearekin batera argitaratuko da.

Bestalde, goiz da software garatzaileen komunitatean garatutako tresnek izan duten eragina neurtzeko. Ixa-pipes tresna multzoak hainbat jarraitzaile dituela konfirmatu dugu eta mundu osoko proiektuetan erabiltzen ari dela ere bai. Zentzu horretan, uztailan egingo den Codefest udako laborategian “Behagunea”-k parte-hartuko du proiektu gisa, eta dagoeneko badakigu interesa sortu duela proiektuak. Proiektua gizarteratzeko eta garatutako tresnak berrerabiltzeko bidean urratsa garrantzitsua dugu hau.

Azkenik Behaguneak DSS2016 proiektuaren ibilbidea aztertzeko tresna gisa ere diseinatu zen. Zentzu horretan, martxan daraman 4 hilabeteetan informazio baliotsua jaso du, proiektu ezberdinek une ezberdinetan izan duten oihartzuna neurtzen ari delarik. Behagunea publiko egin ez bada

ere, mundu digitalean sortzen diren iritziak monitorizatzeko duen ahalmena dagoeneko datu baliagarriak erazteko balio duela erakusten ari da. Kuantitatiboki, 4 hilabeteetan 17.740 aipamen jaso ditu Behaguneak, 3.891 erabiltzaile ezberdinengandik (sare sozialetako atsegite eta bertxioak kontuan izan gabe). Epe honetan aipamen ezberdin gehien izan dituzten proiektuak azaltzen dira 4 taulan.

Proiektua	Aipamenak
Donostia_2016	11804
Euro-Dialogues	753
Emusik	735
Stop_War_Festibala	485
Gaur_Konstelazioak	441
Hirikilabs	379
Olatu_Talka	304
Elkarrizketak	267
Irekiera	190
2016_Bidea	157
Energia_Olatuak	125
Europa_Transit	115
Time_Machine_Soup	115
Labore	111
Rekording	107
Tosta	101
Hibrilaldiak	98
Tira_Birak	98
WhatsART	97

Taula 4. Oihartzun handiena lortu duten proiektuak, aipamen kopurua aipamen ezberdinei dagokie, sare sozialetako atsegite eta bertxioak kontuan izan gabe.

Hizkuntzei dagokienez, jasotako aipamenen %50 gaztelaniaz da, %37a euskaraz, %10 ingelesez eta gainontzeko %3a frantsesez.