



BERRIA

Euskarazko kalitate handiko corpusa sortu du EHUko Informatika fakultateko Ixa ikerketa taldeak: EusCrawl. Nahi duenaren esku jarri du, gainera. Euskarazko hainbat komunikabideren edukietan oinarritu du. Corpusak funtsezkoak dira hizkuntzan oinarritutako aplikazioak sortzeko.

# Euskara ona makinentzat

Jakes Goikoetxea Donostia

**G**aur egungo herritarrak hizkuntzari lotutako aplikazio ugari dituinguruan: sakelakoan edo ordenagailuan idazten ari dela hurrengo hitza izan daitekeena erakusten dutenak, hitzegaldetu eta erantzuten duten laguntzaile adimendunak (Alexa, Siri, Google Assistant...), bezeroen arretarako txatbotak, itzultzaile automatikoak...

Gertuko, ohiko eta berehalako bihurtu diren aplikazio horien muinean, baina, itzaleko lana eta entrenamendua dago. Zehazki, adimen artifiziala, *machine learning*-a eta *deep learning*-a. *Machine learning*-a, ikasketa auto-

matikoa, adimen artifizialaren adar bat da: algoritmoei datu pila bat ematen dizkiete, eta gai dira datu horietan patroiak, ereduak, identifikatzeko eta iragarpenak egiteko. *Deep learning*-a, ikasketa sakona, ikasketa automatiko mota bat da, sare neuronaletan oinarritua.

Funtsean, hizkuntzaren prozesamendua da. Ordenagailuek gizakien hizkuntza ulertzeko ahalgina, adimen artifiziala erabiliz. Teknologia horrek datu pila bat behar ditu, sistemak ikasteko eta entrenatzeko. Hizkuntzan oinarritutako teknologia denez, testuak eman behar zaizkio. Baita euskaraz ere. Ixa taldeak kalitate handiko euskarazko corpus bat osatu du, eta eskura jarri du, nahi duenak erabil dezan: EusCrawl ([ixa.ehu.es/euscrawl](http://ixa.ehu.es/euscrawl)). EHU

Euskal Herriko Unibertsitateko ikerketa talde bat da Ixa. Hizkuntzaren tratamendu automatikoan lan egiten du.

Euskararen kasuan, hizkuntza gutxitua denez, zaila du testu corpus erraldoiak biltzea. Euskarazko corpusak badaude. Batzuk ez daude erabili nahi dituenaren eskura. Beste batzuk bai: «Facebookek, Googlek eta horiek sortuak», azaldu du Aitor Soroak, Ixa taldeko eta Hitz ikerketa zentroko ikertzaile eta EHUko Informatika fakultateko irakasleak. «Baina corpus haiek erabat automatikoki sortuta zeuden: web osoa hartzen dute; webguneak zer hizkuntzatan dauden bereizteko programa automatikoak aplikatzen dituzte; eta edukiak hizkuntzaren arabera automatikoki sailkatzen eta biltzen dituzte.

Hizkuntza guztien corpusak dituzte».

Googlerenak eta Meta AI-renak –lehen Facebook– dira euskarazko testu masa handienak: Googlerena, mC4, mila milioi hitzeko; Meta AIrena, CCI00, 416 milioi hitzeko. Haien kalitatea, baina, zailtasun jarri da, euskarazko edukiak bereizteko programa automatikoek hainbat akats egiten dituztelako.

## Komunikabideak

Kalitate oneko euskarazko corpusa sortu du Ixak. Ez du, Interneteko erraldoien gisan, begi estuko sarea hartu eta Internet osoan arrantza egin, horrela euskarazkoak bai, baina zaborra eta beste hizkuntza batzuetan dauden edukiak ere harrapatzen direlako. «Guk kontrakoa egin dugu», ar-

gitu du Soroak. «Aurrena iturri on batzuk identifikatu genituen; gero programa batzuk sortu genituen, haietatik informazioa xurgatzeko; eta horrela sortu da corpusa».

Edukiak Creative Commons lizentzia librearekin banatzen dituzten Interneteko zenbait webgune aukeratu zituzten, komunikabideak batez ere: Tokikom (tokiko 76 komunikabide biltzen dituen elkarte), BERRIA, eskualdeetako *Hitza* egunkariak, euskarazko Wikipedia, *Argia* eta Bilbo Hiria irratia. Haiak sortutako edukiak xurgatu egin zituzten –*crawl*, ingelesez–. Emaitza: EusCrawl. 12,5 milioi dokumentu eta 423 milioi hitz.

Ikasketa sakonaren barruan hizkuntza ereduak deitutako teknika edo teknologia dago, sare



© JONUBER/FOKU



neuroaletan oinarritua. Hizkuntza eredu horiek testuarekin entrenatzen dituzte: «Testua irakurtzen dute, hizkuntzaren patroiak ikusten dituzte, eta, milioika eta milioika hitz irakurri, ikasi egiten dute. Horrela, hizkuntzari buruzko eredu matematiko erraldoi bat sortzen duzu. Hori bai, ondo idatzitako testuak eman behar zaizkio». Ez dute ikasi bakarrik egiten. Testu berriak sortzeko gai ere badira. Hizkuntza ereduak hizkuntzaren egitura egitura matematiko bihurtzen dute, nolabait esateko. «Hizkuntzari buruzko aplikazioak egiteko gaur egungo tresna onenak dira», nabarmendu du du Soroak.

Ikasi egiten dute, «baina ez dugu oso ondo zer ikasten duten», onartu du Soroak. Kutxa beltzean parekatu ditu. «Ikerkuntza arazotik dago: jakitea zer dagoen hori barruan eta zer egin hori hobeto kontrolatzeko. Ataza bat ematen diozu, eta ikasi egiten du. Input bat ematen diozu, eta output bat ematen dizu».

Ez dira kontrolatzeko errazak, sistema errealak baitira. Hainbat parametro dituzte. Bi adibide ezagun: GPT-3 hizkuntza ereduak, esaterako, 175.000 milioi parametro ditu; Bert-large-k, berriaz, 350 milioi.

### Hizkuntza ereduak

Ixa taldeak, EusCrawl osatu eta gero, euskararen bi hizkuntza eredu sortzeko eta entrenatzeko erabili zuen. Hizkuntza eredu horietako bat gaur egun euskararako dagoen eredu handiena da, 355 milioi parametrokoa. Hizkuntza eredu berri horiek euskarazko beste corpus batzuekin ere entre-

### Euskarazko ahalik eta corpus handiena biltzen saiatu behar dugu, euskarazko aplikazio hobek nahi baditugu»

**Aitor Soroa**

EHUko Informatika fakultateko irakasle eta Ixako eta Hitz-eko ikertzailea

natu zituzten: Googleren mC4rekin eta Meta Alren CC100ekin, beraz, EusCrawl baino kalitate txarragoak.

Probetan ikusi zuten baitetz, EusCrawlen testuen kalitatea besteena baino hobea zela, baina ezusteko ondorio bat ere ateratu zuten: hizkuntza eredu guztiei hizkuntzaren prozesamendurako zenbait eginkizun jarri zizkieten, eta emaitzak berdintsuak izan ziren. Alegia, ez zegoen hainbeste alderik EusCrawlekin entrenatutako eta Googleren eta Meta Alren corpusekin entrenatutako hizkuntza eredu artean.

«Horrek erakusten digu muntro hauek entrenatzeko garrantzitsuagoa dela testuen kantitatea, kalitatea baino», ondorioztatu du Ixak. «Beraz, euskarazko ahalik eta corpus handiena biltzen saiatu behar dugu, euskarazko tresna eta aplikazio hobek nahi baditugu».

Euskarazko komunikabide askoren testuak xurgatu dituzte. Handien artean, EITBren edukiak falta dira. Argialetxeen ere bai. Eta beste hainbat: sare sozialak... Sare sozialei dagokien, hizkuntza ereduak hizkuntza horren zenbat eta erregistro gehiago izan, orduan eta eredu aberatsagoa da.

«Hala ere», ohartarazi du Soroak, «euskararen corpus ezagun guztiak bilduta ere, hizkuntza na-

gusien tamainatik oso urrun gertatu ginatete, eta horrek euskarazko hizkuntza ereduak goi borte bat ezartzen die». Arriskua: euskararentzat sor daitezkeen tresnen kalitatea ezizatea ingelesarentzat sortzen diren parekoa, esaterako.

Egoera horri aurre egiteko, Ixak bi helburu estrategiko ezarri ditu: alde batetik, corpus handiagoak biltzea, euskarazko testu ekoizle guztien edo gehien testuak erabili ahal izatea; bestetik, testu guxtiagorekin ikasteko gai izango diren hizkuntza ereduak ikerketa bultzatzea, hizkuntza gutxitua izatearen mugei aurre egin ahal izateko. Beste herrialde batzuetan ikertzaileei lizentzia librerik gabeko testuak erabiltzen uzten die, corpusak osatzeko, teknologia horien garapena lehenetsia baita.

### Estrategiaren premia

Ixak sortu du EusCrawl, Ixak erabili du hizkuntza ereduak sortzeko, eta Ixak planteatu ditu helburu estrategikoak. EHUko ikerketa talde bat da Ixa. Ez al dute hizkuntzaren prozesamendurako teknologiek herri estrategia bat behar? Ez al da erakundeen lana halako estrategiak bultzatzea eta gartzea? «Apustu bat behar da, baina ez Ixarena bakarrik, erakundeen apustua behar da», Soroaren iriztia.

Ikertzaileak Espainia aipatu du. Espainiako Gobernuak badu Hizkuntzaren Teknologiak Bultzatzeko Plana. Eusko Jaurlaritza ere antzeko zerbait sortzeaz hitz egiten ari dela aipatu du. «Euskararentzat behar-beharrezkoa da halako estrategia bat, bai ikerkuntzaren aldetik, bai aplikazioen aldetik. Ez badugu euskaraz egiten, jendeak beste hizkuntza batean egingo du». Euskal Herriari goi mailako ikerketa taldeak daude hizkuntzaren prozesamenduan.

EusCrawl, euskarazko kalitate handiko corpus librea, ez da Euskal Herrira begira soilik egindako ekarpena. Munduko edozein ikertzailek erabili ahal izango du. Erabiltzen ari dira jada, BigScience proiektuan: hizkuntza eredu eleantatu eta librea sortu nahi dute –normalean enpresa handiek baino ez dituzte egiten, garestiak baitira–. Proiektu irekia eta kolektiboa da. Euskararen txertatu dute, EusCrawlen bidez. Beraz, sortutako hizkuntza ereduak euskaraz ere jakingo du.

EusCrawlekin badu beste erabilgarritasun bat: hizkuntzalaritza, hizkuntzaren azterketa. Corpusa bali dezakete erabilera ikusteko, esaterako.

**ARGI ALDIAN**  
Ana Galarraga Aiestaran

Elhuyar Zientzia

## Arren itxuraz, gizonaz jantzita

**E**rdialdeko eta Hego Amerikan, bada *Florisuga mellivora* izeneko kolibri espezie bat. Helberrak, hegazti eta beste animalia askotan gertatzen den bezala, arrek eta emeek itxura desberdina dute. Hegazti gehienetan, gazteek eme helduen tankera dute, eta helduek, arrek beste lumaje bat hartzen dute, normalean, koloretuagoa eta deigarriagoa. *Florisuga mellivora* kolibrietan, baina, alderantzizkoa gertatzen da: gaztetan ar helduen antza dute lumajearen, eta, helduek, emeek hartzen dute bestelako lumajea. Ez guztiak, ordea: emeen % 20k arren itxurari eusten diote.

Ikertzaileek ez dakite horren atzean arrazoi genetikoak ala inguruneak ote dauden, baina ondorio bat, behintzat, ikusi dute: ar-itxura duten emeek jazarpen txikiagoa jasaten dute, besteekin alderatuta. Hain zuzen ere, arrek dira jazarleak, eta, beraz, eme arrunten aurka erabiltzen dute indarkeria, jana eskuratzeko. Ar-itxura duten emeek, berriz, ez dute erasorik jasaten, eta ez dute eragozpenik arrek adina jateko.

Ugalketaren ikuspegitik, ar-itxura izatea ez da abantaila bat. Aitzitik, ohiko lumajea eta arra-

rena duten bi emeen artean, lehena aukeratzen dute arrek gorrea egiteko. Baina ar-itxurakoak ere gorrotean dituzte; beraz, itxura aldatzeko dakarren irabaziala galera baino handiagoa da.

Kolibri horiek burura ekartzeko dute gizonen aukerak eskuratzeko agatik gizon-itxura hartu duten emakumeen kasua. Historian, asko izan dira horrela jokatu dutenak, tartean, zientzialariak. Ezagunenetako bat Sophie Germain da. 1776an jaio zen, eta txikitatik zen zientziagilea. 18 urte zituela, Eskola Politeknikoa ireki zuten Parisen, baina gizonak baino ez zituzten onartzen. Alabaina, Germainek aurkitu zuen han ikasteko modua: postaz jasotzen zituen lezioak eta bidaltzen lanak, Monsieur Antoine-August Le Blanc izena hartuta.

Joseph Louis Lagrange irakaslea Le Blancen lan bikainaz jabetu zen, eta ezagutu nahi izan zuen. Germainek ezin izan zuen emakumea zela ezkutatua, baina Lagrangek ez zuen baztertu; alderantziz, Germaini matematikaren munduan aurrera egiten lagundu zion. Azkenean, matematikari handia izatera iritsi zen, eta Parisko Zientzien Akademiaren onartua izatea lortu zuen.

Garai bertsuan, Jeanne Baret botanikaria gizonaz jantzita ontziratua zen munduari bira eman zion lehen frantziar espedizioan. Izan ere, o ezustekoal, emakumeek debekatua zuten espedizioetan parte hartzea. Munduari itzulia eman zion lehen emakumea bilakatu zen, jada emakumez jantzita, eta botanikaren arloan sekulako ekarpena eginda.

Eta XXI. mendean? Bada, emakumeek sinatutako artikulu zientifikoek aukera gutxiago dituzte aldizkari ospetsuetan argitaratuak izateko, eta aipamen gutxiago jasoko dituzte. Eta jantziek ere elementu garrantzitsuak izaten jarraitzen dute, onartuak ala baztertuak izateko, adibidez, kongresuetan. Hautespene naturala deituko diote gero.

**Eta XXI. mendean? Bada, emakumeek sinatutako artikulu zientifikoek aukera gutxiago dituzte aldizkari ospetsuetan argitaratuak izateko, eta aipamen gutxiago jasoko dituzte**