

Baliabide lexiko-semantikoak hizkuntz industriarako

Jose Maria Arriola / Arantza Diaz de Ilarraza / Kepa Sarasola

hizkuntzaren inguruan industria berria sortzen ari da, konputagailuaz baliaturik hizkuntza tratatzea helburu duena. Arlo honek aurrera egin dezan, ezinbestekoak dira hitzen esanahia emango duten baliabide lexikalak. Europako Batasuneko *Language Engineering* programaren irizpideetan baliabide lexikalek duten oinarritzko papera azpimarratuta azaltzen da.

Mintzagai dugun Hiztegia 2002 proiektuak ere Europako Elkartearen laguntza du (FEDER, 2FD97-2000-2001) eta baliabide lexikalak sortzea jomuga izan duten honako proiektu hauekin dago lotuta: Wordnet², EuroWordnet eta ITEM³. Proiektu honen bidez, IXA⁴ taldeak ondorengo baliabide lexikalak erdietsi nahi ditu:

- Euskal Hiztegiaren bertsio egituratua; horretarako, TEIko (*Text Encoding Initiative*) gida-lerroei jarraiki. SGML (*Standard Generalized Markup Language*) lengoia estandarra erabiliko da.
- Euskal Hiztegiaren ezagutza-base lexikala: bertatik ateratako erlazio semantikoez osatua.



(aleman, espainiera, estoniera, frantses, ingeles, italiara, nederlandera eta txekierara) zabaldu da.

EBL gehienak ingeleserako sortu direnez, gainontzeko hizkuntzak teknologia berriekiko egoera ahulean daude. Egoera horri aurre egiteko bi irtenbide osagarri ikusten ditugu:

- Euskal Wordnet: EuroWordnet egokitzea, ingelesezko kontzeptuei euskarazkoak lotuz.

Baliabide horiek sortzean, besteak beste, ondorengo produktu komertzialak garatzea dugun helburu:

- Euskal Hiztegiaren bertsio elektronikoko egituratua (CD-ROMean, Interneten edota testu-prozesadoreetan integratua)
- Euskararako *thesaurus* bat testu-prozesadoreetan integratua: sinonimia, hiperonimia, hiponimia eta bestelako kontzeptuen erlazioak kontsultatu ahal izateko.

1. Hizkuntza bakoitzarentzat dauden corpus eta hiztegiatik abiatuta EBLak sortzea. Gure kasuan, iturri lexikal gisa. Euskal Hiztegia baliatu dugu. Lehenbiziko lana Euskal Hiztegia SGML-TEI estandarrei jarraituz egituratzea izan da. Horrela bada, euskara aztergai edo lan-tresna duen edonorentzat baliagarri izango da. Bertsio egituratu honen definizioak aztertuz, hainbat erlazio lexiko-semantiko aterako ditugu: sinonimia, hiperonimia (*klase-azpiklase* erlazioa; adibidez: animalia-intsektu), meronimia (*osoa-partea* erlazioa; adibidez: txori-moko), etab.

2. Beste hizkuntzetarako EBLak sortzeko, ingeleserako egin diren EBLez baliatzea. Gure kasuan, EuroWordnet abiapuntuz hartuz euskararako Wordnet-a egin nahi dugu, ingelesezko kontzeptuei euskarazkoak lotuz. Euskararako Wordnet hori egiteko, metodo erdiautomatikoak erabiliko ditugu, baina gero emaitzak eskuz orraztuko ditugu.

Irudian ikus daiteke EuroWordnet-i euskarazko kontzeptuak lotzeko erabiltzen dugun interfazea⁵.

¹ <http://www.let.uva.nl/~ewn/>
² <http://www.cogsci.princeton.edu/~wn>
³ <http://sensei.ieec.uned.es/item/>
⁴ <http://ixa.si.ehu.es/>
⁵ EuroWordnet kontsultatu nahi duenak jo beza ondorengo helbidera: <http://ixa.si.ehu.es/wei3.html>

Proiektuaren izenburua:

HIZTEGIA 2002: Hizkuntz teknologiko aplikazioetarako ezagutza-base lexiko-semantikoa

Helburua:

Euskararen ontologia eraikitzea, hau da, euskarazko hitzen esanahiak modu hierarkikoan definitzea eta beren arteko hainbat erlazio lexiko-semantikorekin osatzea (sinonimia, hiperonimia, meronimia...). Horretarako, Ibon Sarasolaren Euskal Hiztegia eta EuroWordnet¹ ontologia izango ditugu oinarri.

Bi azpigelburu ditugu: alde batetik, Euskal Hiztegi adieren definizioetatik erlazio lexiko-semantikoak ateratzea eta, bestetik, EuroWordnet-eko ingelesezko kontzeptuei euskarazko kontzeptuak lotzea.

Zuzendaria:

Arantza Diaz de Ilarraza.

Ikerketa-taldea:

Lan-taldea: Agirre E., Ansa O., Arregi X., Arriola J.M., Artola X., Diaz de Ilarraza A., Garcia E., Irurtzun A., Lersundi M., Martinez D., Pociello E., Sarasola K., Sorro A., Urkia M., Zubezu E.

Partaideak: UPV-EHUko IXA taldea, UZEI eta Ibon Sarasola

Saila:

Lengoia eta Sistema Informatikoak

Fakultatea:

UPV-EHUko Informatika (Donostia)