



IKER
GAZTE
NAZIOARTEKO
IKERKETA EUSKARAZ

I. IKERGAZTE

NAZIOARTEKO IKERKETA EUSKARAZ

2015eko maiatzaren 13, 14 eta 15
Durango, Euskal Herria

ANTOLATZAILEA:
Udako Euskal Unibertsitatea (UEU)

INGENIARITZA ETA ARKITEKTURA

**Euskarazko izena+aditza
konbinazioak corpusetik
automatikoki erauztea eta
idiomatikotasunaren arabera
karakterizatzea**

A. Gurrutxaga, I. Alegria eta X. Artola

599-606 or.
<https://dx.doi.org/10.26876/ikergazte.i.82>

ANTOLATZAILEA:



BABESLEAK:



eman ta zabal zazu



LAGUNTZAILEAK:



Euskarazko izena+aditza konbinazioak corpusetik automatikoki erauztea eta idiomatikotasunaren arabera karakterizatzea

Gurrutxaga, A.¹

Alegria, I.² eta Artola, X.²

¹ Elhuyar Fundazioa. Zelai Haundi kalea, 3, Osinalde industrialdea, 20170 Usurbil

² Ixa taldea - UPV/EHU. Manuel Lardizabal 1, 48014 Donostia

Laburpena

Euskarazko izena+aditza egiturako unitate fraseologikoak (UFak) corpusetik automatikoki erauzi eta idiomatikotasun-mailaren arabera karakterizatze lan esperimentalak egin dugu. Corpusetik hautagaiak erauzteko sistema landu ondoren, idiomatikotasunaren lau ezaugarri edo propietateetako bakoitza neurtzeko teknikak garatu eta ebaluatu ditugu, hiru adituk eskuz sailkatutako erreferentzia erabiliz. Hiru kategoria bereizi dira: esapide idiomatikoa, kolokazioa eta konbinazio librea. Azkenik, ezaugarri bakunen neurketak ikasketa automatikoko sailkatze-ataza batean konbinatu dira. Ondorio nagusia da arlo honetan estandar diren agerkidetze-tekniken emaitzak modu esanguratsuan gainditu direla, batez ere teknika semantikoaren bidez, baina baita malgutasun morfosintaktikoaren neurketaren bidez ere.

Hitz gakoak: fraseologia konputazionala, idiomatikotasuna, esapide idiomatikoak, kolokazioak

Abstract

We present an experimental study on the automatic extraction of phraseological units of noun+verb structure in Basque, and their characterization according to the idiomaticity level. After automatically extracting candidates from corpora, we develop several techniques for quantifying the four basic properties of idiomaticity, using for evaluation a gold standard of candidates classified by three experts. We use three categories: idioms, collocations and free combinations. Finally, the results of those experiments have been combined using Machine Learning for classification. The results show that the standard cooccurrence techniques are significantly outperformed by semantic measures, and, to a lower extent, by measures of morphosyntactic flexibility.

Keywords: computational phraseology, idiomaticity, idioms, collocations

1 Sarrera eta motibazioa

Hitz anitzeko unitateek edo unitate fraseologikoen (UFek) egiteko giltzarria dute, hiztegigintzan eta hizkuntzen irakaskuntzan ez ezik, hizkuntzaren prozesamendu automatikoan ere (Sag *et al.*, 2010). Gaur egun aski onartua dago hizkuntzaren funtzionamendua ezin dela osagai bakunen konbinazio libreaz soilik azaldu, hiztunek erabiltzen dituzten elementu batzuk nolabaiteko unitate aurrez eratuak baitira (Fillmore, 1979, 92). Horrelakoak dira, esaterako, *adarra jo* eta *zarata atera*. Lehenaren esanahia ezin igarri osagaien esanahietatik; bigarrenean, *atera* aditzak adiera berezia du, gertuago baitago ‘egin’ edo ‘sortu’ adieratik, *ateraren* ohiko adieratik baino.

Euskararen kasuan, hitz anitzeko unitateen prozesamenduan egindako aurrerapausoak terminologiaren erauzketara bideratu dira (Alegria *et al.*, 2004, 2006; Saralegi *et al.*, 2008) eta corpusean automatikoki etiketatzen datu-base lexikaletan adierazita dauden UFak (Urizar, 2012).

Ikerkuntza honen bidez, ekarpen bat egin nahi izan dugu euskarazko fraseologia konputazionalaren arloan. Zehazki, euskarazko izena+aditza egiturako UFak corpusetik automatikoki erauzi eta idiomatikotasun-mailaren arabera karakterizatze lan esperimentalak egin dugu.

- Osagaien ordezkagarritasuna: konbinazioaren estatistikak osagaiak sinonimoez ordezkatzuz sortzen diren konbinazioen estatistikekin konparatzea.

Konparazio horien emaitza zenbat eta desberdinago, hautagaia hainbat eta idiomatikoago. Agerkide-tza da gehien erabili den teknika, eta hori neurtzeko elkartze-neurriak hasieratik erabili dira arlo honetako ikerkuntzan (Church eta Hanks, 1990; Smadja, 1993). Beste ezaugarriak kuantifikatzeko ikerlanak nabarmen ugartu dira azken hamarkadan (Baldwin *et al.*, 2003; Van de Cruys eta Moirón, 2007; Fazly eta Stevenson, 2007), dela hautagaien rankingak osatzeko, dela automatikoki sailkatzeko.

2.3 Ikerketaren helburuak

Ikerkuntza-lan honen helburu nagusia izan da corpusetatik **izena+aditza** osaerako UFAk automatikoki eskuratzeko eta haien idiomatikotasunaren arabera karakterizatzeko teknikak garatzea, eta ikergai hauen inguruko ezagutza eskuratzeko eta sortzea: a) idiomatikotasunaren eta haren propietate bakoitzaren neurketen artean dagoen korrelazioa; b) UFAren propietateen ebidentzia enpirikoak zerbateraino datozen bat teoria fraseologikoen aurreanekin; eta c) idiomatikotasuna konplexua izanik, propietateen neurketak konbinatuz emaitza hobeak lortzen diren ala ez jakitea.

3 Lan esperimentalak

3.1 Lan esperimentalaren diseinua

1 taulan eman dugu aztergaitzat izan ditugun **izena+aditza** osaerako konbinazio-moten errepertorioa.

1	izena _{subjektua} + aditza : <i>burua joan, eguzkia sartu; gogoak eman, loak hartu</i>
2	izena _{objektua} + aditza : <i>hanka sartu, zubiak eraiki, lan egin, itxurak egin, zarata atera, gola sartu</i>
3	izena _{subjektuaren pred.} + aditza : <i>beldur izan, falta izan, giro egon</i>
4	izena _{objektuaren pred.} + aditza : <i>atsegin ukan/*edun, damu ukan/*edun</i>
5	izena _{datiboa} + aditza : <i>edanari eman, bideari ekin, lanari lotu</i>
6	izena _{adjuntua} + aditza : <i>mendean hartu, adarretatik heldu, harira etorri, aurrez ikusi</i>

1 Taula: **izena+aditza** osaerako konbinazio-moten errepertorioa.

Bi karakterizazio-ataza bereizi dira: a) ranking-ataza, non propietate bakoitzaren neurketen emaitzak hautagaiak ordenatzeko erabili baitira; eta b) sailkatze-ataza, non esperimentu bakunak emaitzak konbinatu baititugu, ikasketa automatikoko teknikak erabiliz. Esperimentuetarako, 74 milioi hitzeko kazetaritza-corpus bat erabili dugu, Euskaldunon Egunkariako eta Berriako artikuluek osatua.

3.2 UF hautagaiak erauztea

UF hautagaiak testetik automatikoki erauzteko eta forma kanonikoa esleitzeko prozesu automatikoa garatu da. Erauzketa-prozesuaren deskribapen zehatza Gurrutxaga eta Alegria (2011) lanean egin da. Bi urratsetan antolatua dago: bigrama-sorkuntza eta bigramen forma kanonikoa lortzeko normalizazioa. Lehen urratserako, Ngram Statistics Package erabili dugu¹. Erauzketaren aurretik, corpusa linguistikoki prozesatu da, Ixa taldearen Eustagger etiketatzailearen bidez lehenik (Alegria *et al.*, 2002), eta, ondoren, ikerkuntza honetan landutako prozesu batzuk inplementatu dira, corpusak erauzketarako behar diren ezaugarri batzuk izan dituzan.

¹<http://search.cpan.org/dist/Text-NSP>

3.3 Ebaluazioa: metodologia eta baliabideak

Ranking bidezko atazan, heinen korrelazio-koefizienteak erabili ohi dira, eta guk Kendall τ aukeratu dugu, berdinketak kudeatzeko modalitatean (Kendall τ_B). Batez besteko doitasuna (AP - *Average Precision*) ere erabili dugu, kategoria bakoitzeko hautagaiekiko ebaluazioa egiteko (AP_{UF} , AP_{id} eta AP_{col}). Sailkapen automatikoaren bidezko atazen ebaluazioan, Weka tresnak eskaintzen dituen zenbait neurri erabili ditugu: ondo sailkatutako instantzia-kopurua (CC), klase bakoitzaren F neurriak, F_{mikro} eta F_{makro} .

Ebaluaziorako *gold standard*tzat, erauzketatik ausaz ateratako eta aditu-talde batek eskuz sailkatutako 1 145 konbinazioko multzo bat erabili dugu, hiru kategoriatan banatua: esapide idiomatikoak (80), kolokazioak (268) eta konbinazio libreak (797). Anotatzaileen arteko adostasunerako, adostasun ertainaren mailako 0,55eko Fleiss κ lortu dugu (Landis eta Koch, 1977). Horrez gain, erauzketaren parametro batzuk doitzeko, beste erreferentzia bat ere eratu dugu, iturri hauetako izena+aditza unitatez osatua: a) Euskaltzaindiaren *Hiztegi Batua*; b) Ibon Sarasolaren *Euskal Hiztegia*; c) Elhuyar Fundazioaren *Euskara-Castellano/Castellano-Vasco Hiztegia*; d) *Intza proiektua*²; eta e) Ixa taldearen EDBL datu-base lexikala (Aldezabal *et al.*, 2001).

3.4 Idiomatikotasunaren ezaugarrien edo propietateen banakako neurketa

3.4.1 Idiosinkrasia estatistikoa, agerkidetza neurtuz

Bigramen agerkidetza-datuen analisi estatistikoa S. Everten UCS toolkit³ paketearen bidez egin dugu (Evert, 2005). Elkartze-neurri hauek kalkulatu ditugu: z neurria, t neurria, khi karratua (χ^2), egiantz-arrazoiaren logaritmoa, Fisherren test zehatza, elkarrekiko informazioa (MI), MI^3 eta f .

3.4.2 Konposizionaltasun semantikoa, antzekotasun distribuzionala neurtuz

Hauek dira UFen konposizionaltasuna karakterizatzeko gure metodologiaren oinarriak:

- UF hautagaien testuinguruak haren osagai bakunen testuinguruekin konparatzea.
- Osagai bakunen testuinguruetan konbinazioaren testuinguruak ez sartzea. Esaterako, *mahaia jaso* bigramaren agerpena atzematen denean, testuinguruko hitzek *mahaia jasoren* testuinguru-dokumentua elikatu dute, eta ez dira *mahai* eta *jasoren* testuingurutzat kontsideratu.

Antzekotasun distribuzionaleko tekniken bidez egin dugu testuinguruen konparazioa (Gurrutxaga eta Alegria, 2012). Lehenik, bektore-espazioaren ereduan (VSM – *Vector Space Model*) ohikoak diren neurri batzuk esperimatu ditugu [Jaccard koefizientea, kosinua eta Jensen-Shannon dibergentzia; Berry-Roggheren (1974) R balioa eta Wulffek (2008) egindako bi hedapenak]. VSMren implementazio berezia den ezkutuko semantikaren analisia (LSA – *Latent Semantic Analysis*) izeneko teknikaren aplikazioa ere esploratu dugu. Infomap softwarea erabili dugu horretarako⁴. Azkenik, dokumentu arteko antza kalkulatzeko IR arloko indize batzuk ere aplikatu ditugu. Lemur Toolkit aukeratu dugu testuinguru-dokumentuen arteko antza kalkulatzeko⁵.

3.4.3 Malgutasun morfosintaktikoa, erreferentzia-portaerarekiko distantzia neurtuz

Metodologiaren oinarria da aztergai dugun bigrama bakoitzaren portaera morfosintaktikoa erreferentzia-portaera batekin konparatzea. Horretarako, bi konparazio-prozedura hauek esperimatu ditugu:

- Portaera orokorrarekiko konparazioa (Barkema, 1994; Fazly eta Stevenson, 2007; Wulff, 2008): kategoria-osaera bereko (gurean, izena+aditza) konbinazioen batez besteko portaerarekikoa.
- Osagaien portaerarekiko konparazioa (Bannard, 2007): konbinazioaren osagai batek beste osagaiaren kategoriako edozein hitzekin osatutako konbinazioen batez besteko portaerarekikoa.

²<http://intza.armiarma.com>

³<http://www.collocations.de/software.html>

⁴<http://infomap-nlp.sourceforge.net/>

⁵<http://www.lemurproject.org>

Kontuan hartutako aldakuntzak. Malgutasun morfosintaktikoa neurtzeko aztertzen diren aldakuntzen multzoa hizkuntzaren ezaugarrien araberakoa da. Arlo honetako ikerlan esanguratsu batzuk (Martínez, 1996; Oyharçabal, 2006; Odriozola, 2010) eta EDBLko aditz-lokuzioen gauzatzeko eskemak (Urizar, 2012) kontuan izanik, kasu hauek hautatu ditugu malgutasun morfosintaktikoa karakterizatzen:

- Izenaren ezker- eta eskuin-hedapenak: determinatzailea (*liburu **bat** irakurri dut*); izenondoa (*liburu **interesgarria** irakurri nuen*); izenlaguna (***gustuko** liburuak irakurtzea*); eta erlatiboa (***irakurri dudan** liburua; anaiak **irakurritako** liburu batzuk*).
- Hedapen bat baino gehiago konbina daitezke aldakuntza berean: *liburu interesgarri bat irakurri dut*.
- Mugatasuna: izenaren edo haren hedapenen mugatasun-aldakuntzak. Mugatasun-informazioa sintagmaren azken osagaiak darama. Kasu sinple batzuk: *liburua/liburu**ak**/zenbait liburu**Ø**/liburu**ok** irakurri*; *egunkari**an**/egunkari**etan**/hiru egunkari**tan**/egunkari**otan** irakurri*.
- Osagaien ordena (IZE ADI / ADI IZE): *liburua irakurri dut/irakurri dut liburua*.

Bigramaren portaera ezagutzeko, aldakuntza bakoitzean duen maiztasuna kontatzen dugu. Bestetik, malgutasuna kalkulatzeko erabili ditugun bi konparazio-prozeduretako portaera ezagutzeko behar diren kontaktak ere egiten dira. Neurri hauek erabili ditugu portaeren arteko distantzia kalkulatzeko: Kullback-Leibler dibergentzia (Fazly eta Stevenson, 2007); Wulffen (2008) SSD neurria (*sum of squared deviations*) eta entropia erlatiboa; eta Bannarden (2007) CPMI (*conditional pointwise mutual information*).

3.4.4 Malgutasun lexikala, osagaien ordezkagarritasuna neurtuz

Bigramen osagaien ordezkagarritasuna neurtzeko, osagaien ordezkoak behar ditugu. Horretarako baliabide hauek erabili ditugu: Elhuyarren *Sinonimoen Kutxa*⁶; eta Ixa taldearen Euskal WordNet (Pociello *et al.*, 2011). Halaber, baliabideon estaldura handitzeko, konbinatu egin ditugu, baita hedatu ere, bigramen osagai bakoitzaren Euskal WordNeteko ahaideak (*siblings*) gehituz.

Ordezkarritasuna neurtzeko, Van de Cruys eta Moirónen (2007) R_{nv} eta R_{vn} indizeak erabili ditugu, eta Fazly eta Stevenson (2007) Fixedness_{lex} neurria.

3.4.5 Esperimentu bakunen emaitzen analisia

2 taulan laburbildu ditugu teknika bakoitzarekin lortutako emaitza onenak. 2 irudian, horien P/R (doitasuna/estaldura) kurbak.

2 Taula: Esperimentu bakunetan teknika bakoitzarekin lortutako emaitza onenak.

teknika	τ_B	AP_{UF}	AP_{Id}	AP_{Co1}
ausaz	0,000	0,309	0,070	0,234
AM	0,197	0,455	0,119	0,383
DSim	0,322	0,566	0,320	0,431
MSFlex	0,154	0,434	0,202	0,331
LFlex	0,110	0,381	0,122	0,323

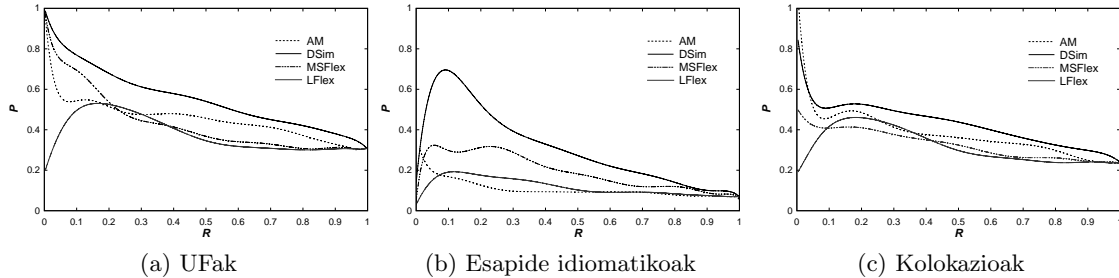
Emaitza horien alderdi esanguratsuenak:

- Antzekotasun distribuzionaleko neurriak (DSim) dira idiomatikotasun-mailarekin korrelazio onena dutenak⁷. Bestetik, esapide idiomatikoen erauzketan da are nabariagoa horien nagusitasuna.
- Kolokazio-erauzketan, AM neurri onena (t neurria) DSim neurrien lehiakide handia da, baina AP_{Co1} emaitza onena lortu duena aditzaren semantika neurtzen duen Indri indize bat izan da. Emaitza hori oso garrantzitsua da, kolokazioak erdikonposizionalak direlako ikuspegiarekin bat baitator.

⁶ *Sinonimoen Kutxa*. 2010. 3. ed. Elhuyar Fundazioa. Usurbil.

⁷ Zehazki, testuinguru-dokumentuen arteko antza neurtzen duten Indri indizea eta KL dibergentzia dira neurri onenak, L2 esperimentu-modalitatean.

- MSFlex emaitza onenak esapide idiomatikoekin lortu dira (AP_{id} ; zehazki, mugatasun-aldakuntzen neurketetan). Gainerakoan, ez dituzte agerkidetzatza-esperimentuen emaitzak gainditzen.
- LFlex neurketetan, emaitzak espero baino txarragoak dira.

2 Irudia: Idiomatikotasun-rankingen P/R (doitasuna/estaldura) kurbak.

3.5 Propietateen integrazioa: ikasketa automatikoa

Ikasketa automatikoko esperimentuak egiteko, Weka paketea erabili dugu⁸. Sei algoritmorekin probak egin ondoren (Gurrutxaga eta Alegria, 2013), bi hauekin lortu ditugu emaitza onenak⁹: Logistic Regression eta SMO (*Sequential Minimal Optimization*)¹⁰.

Atributu gisa, aurreko ataleko neurketen emaitzak erabili ditugu, eta, Fazly eta Stevensoni (2007) jarraituz, konbinazioaren aditza ere bai. Ebaluazioan, balidazio gurutzatua erabili dugu, ebaluazio-erreferentzia ez delako oso handia (1 145 instantzia). Atributuak hautatzeko iragazkiak ikaste-multzoan dauden instantziak soilik azter ditzan, AttributeSelectedClassifier metasailkatzailea hautatu dugu.

3 taulan, lau datu-multzorekin lortutako emaitzak bistaratu ditugu: DSim (antzekotasun distribuzio-naleko neurketak); 4 osag. (lau propietateen neurketak); 4 osag.+ad. (aurreko atributuak eta aditza); eta CS-BF (AttributeSelectedClassifier metasailkatzailean, CfsSubsetEval1 ebaluatzaileaz, eta BestFirst bilaketa-metodoaz osatutako iragazkia).

3 Taula: Ikasketa automatikoko esperimentuen emaitzak.

Atrib.	Metod.	CC	F_{id}	F_{col}	F_{free}	F_{mikro}	F_{makro}
Oinarri-lerroa		69,607	0,000	0,000	0,821	0,571	0,274
DSim	LR	74,061	0,270	0,468	0,842	0,714	0,527
4 osag.	LR	73,362	0,355	0,487	0,837	0,722	0,560
	SMO	76,070	0,300	0,479	0,858	0,731	0,546
4 osag.+ad.	SMO	76,856	0,418	0,544	0,858	0,754	0,607
CS-BF	LR	75,721	0,339	0,487	0,854	0,732	0,560

Emaitza onenak SMO algoritmoarekin lortu dira, 4 osag.+ad. datu-multzoa erabiliz. Beraz, ezau-garri guztiak egiten dute beren ekarpena. Hala ere, datu-multzo hori osatzeko prozesamendua handia da, eta LR metodoak atributu-hautaketa automatikoarekin lortutako emaitzek atea irekitzen diote bideragarritasunaren eta emaitzen kalitatearen arteko erlazio on bat lortzeko aukerari.

4 Ondorioak

Ondorio nagusia da ataza honetan estandar diren agerkidetzatza-tekniken emaitzak modu esanguratsuan gaudituz direla, batez ere teknika semantikoaren bidez, baina baita malgutasun morfosintaktikoaren neur-

⁸<http://www.cs.waikato.ac.nz/ml/weka/>

⁹Gainerako lauak: Naive Bayes, C4.5 decision tree (j48), RandomForest eta PART.

¹⁰Support Vector Machine edo sostengu-bektoreen makinaren modalitatearen inplementazio bat.

ketaren bidez ere. Aldiz, finkapen lexikalaren neurketak ez du espero zen mailako emaitza izan.

Bestetik, fraseologiaren aurrean batzuen ebidentzia experimentalak lortu ditugu: idiomatikotasunaren konplexutasuna, konposizionaltasunaren UF-kategorien zeharreko graduazioa, eta kolokazioen erdikonposizionaltasuna eta malgutasun morfosintaktikoa.

Ikerketa honen ekarpenak baliagarriak dira etorkizuneko hiztegitintzak automatizaziorantz izango duen bilakabidean, eta hizkuntzaren prozesamenduko arloko zenbait atazatan, hala nola datu-base lexikalen elikatzean, corpusen etiketatzean eta, testuinguru eleaniztunean aplikatuta, itzulpen automatikoan.

5 Etorkizuneko lanak

Lehenik, informazio sintaktiko aberatsagoa (adib., osagaien mendekotasunak) erabiltzeak erauzketa hobetu dezake. Bestetik, interesgarria da interpretazio literala eta idiomatikoa izan ditzaketen konbinazioen agerpenak bereiztea. Azkenik, komeni da garatutako metodologia corpus paraleloei aplikatzea, bi hizkuntzako UF bikote baliokideak erauzteko eta karakterizatzeko.

Erreferentziak

- ALDEZABAL, I., O. ANSA, B. ARRIETA, X. ARTOLA, A. EZEIZA, G. HERNÁNDEZ, eta M. LERSUNDI. 2001. EDBL: A general lexical basis for the automatic processing of Basque. In *IRCS Workshop on linguistic databases*, 1–10, Philadelphia, Pennsylvania.
- ALEGRIA, I., MAXUX ARANZABE, AITZOL EZEIZA, NEREA EZEIZA, eta RUBEN URIZAR. 2002. Robustness and customisation in an analyser/lemmatiser for Basque. In *Proceedings of Workshop on “Customizing knowledge in NLP applications”*. *Third International Conference on Language Resources and Evaluation*, 1–6, Las Palmas de Gran Canaria.
- , A. GURRUTXAGA, P. LIZASO, X. SARALEGI, S. UGARTETXEA, eta R. URIZAR. 2004. An xml-based term extraction tool for Basque. In *LREC2004: 4th International Conference On Language Resources And Evaluation*, 1733–1736, Lisboa.
- , A. GURRUTXAGA, X. SARALEGI, eta S. UGARTETXEA. 2006. Elexbi, a basic tool for bilingual term extraction from Spanish-Basque parallel corpora. In *12th EURALEX International Congress*, 159–165, Turin.
- BALDWIN, T., C. BANNARD, T. TANAKA, eta D. WIDDOWS. 2003. An empirical model of multiword expression decomposability. In *Proceedings of the ACL 2003 Workshop on Multiword Expressions: analysis, acquisition and treatment*, volume 18, 89–96, Sapporo.
- , eta S.N. KIM. 2010. Multiword expressions. In *Handbook of Natural Language Processing, second edition*, ed. by N. Indurkha eta F. J. Damerau, 267–292. Boca Raton, AEB: CRC Press, Taylor and Francis Group.
- BANNARD, C. 2007. A measure of syntactic flexibility for automatically identifying multiword expressions in corpora. In *Proceedings of the Workshop on a Broader Perspective on Multiword Expressions*, 1–8, Praga. ACL.
- BARHEMA, HENK. 1994. Determining the syntactic flexibility of idioms. *Realising and Using English Language Corpora* 39–52.
- BERRY-ROGGHE, G.L.M. 1974. Automatic identification of phrasal verbs. *Computers in the Humanities* 16–26.
- CHURCH, K.W., eta P. HANKS. 1990. Word association norms, mutual information, and lexicography. *Computational linguistics* 16. 22–29.
- COWIE, A.P. 1998. *Phraseology: Theory, Analysis, and Applications*. Oxford University Press, USA.
- EVERT, S., 2005. *The statistics of word cooccurrences: Word pairs and collocations*. University of Stuttgart tesia.
- FAZLY, A., eta S. STEVENSON. 2007. Distinguishing subtypes of multiword expressions using linguistically-motivated statistical measures. In *Proceedings of the Workshop on A Broader Perspective on Multiword Expressions*, 9–16, Praga. ACL.

- FILLMORE, CHARLES J. 1979. On fluency. In *Individual Differences in Language Ability and Language Behavior*, ed. by D. Kempler eta W. S. Y. Wang, 85–101. New York: Academic Press.
- GRANGER, S., eta M. PAQUOT. 2008. Disentangling the phraseological web. In *Phraseology: An Interdisciplinary Perspective*, ed. by S. Granger eta F. Meunier, 27–50. John Benjamins Publishing.
- GURRUTXAGA, A., eta I. ALEGRIA. 2011. Automatic extraction of NV expressions in Basque: basic issues on cooccurrence techniques. In *Proceedings of the Workshop on Multiword Expressions: from Parsing and Generation to the Real World. ACL HLT 2011*, 2–7, Portland, Oregon.
- , eta —. 2012. Measuring the compositionality of NV expressions in Basque by means of distributional similarity techniques. In *Proceedings of the 8th international Conference on Language Resources and Evaluation - LREC 2012*, 2389–2394, Istanbul.
- , eta —. 2013. Combining different features of idiomaticity for the automatic classification of noun+ verb expressions in Basque. In *Proceedings of the 9th Workshop on Multiword Expressions (MWE 2013) NAACL HLT 2013*, volume 13, 116–125, Atlanta, Georgia.
- HEID, U. 2008. Computational Phraseology. An overview. *Phraseology: an interdisciplinary perspective* 337–360.
- LANDIS, J RICHARD, eta GARY G KOCH. 1977. The measurement of observer agreement for categorical data. *Biometrics* 159–174.
- MANNING, C. D., eta H. SCHÜTZE. 1999. *Foundations of Statistical Natural Language Processing*. Cambridge, Massachusetts: The MIT Press.
- MARTÍNEZ, A, 1996. Syntactic evidence in favour of degrees of incorporation in [*n+egin*] constructions. Eskuizkribua.
- ODRIOZOLA, J.C., 2010. Euskararen aditz-unitate fraseologikoen deskribapena. Unibertsitateko katedra plazarako lehiaketa. Zientzia eta Teknologia fakultatea.
- OYHARÇABAL, BERNARD. 2006. Basque light verb constructions. *Anuario del Seminario de Filología Vasca Julio de Urquijo. R. L. Trasken oroitzenetan ikerketak euskalaritzaz eta hizkuntzalaritza historikoaz* 40. 787–806.
- POCIELLO, E., E. AGIRRE, eta I. ALDEZABAL. 2011. Methodology and construction of the basque wordnet. In *Language resources and evaluation*, volume 45, 121–142. Springer.
- SAG, I., T. BALDWIN, F. BOND, A. COPESTAKE, eta D. FLICKINGER. 2010. Multiword expressions: A pain in the neck for NLP. *Computational Linguistics and Intelligent Text Processing* 189–206.
- SARALEGI, X., I. SAN VICENTE, eta A. GURRUTXAGA. 2008. Automatic extraction of bilingual terms from comparable corpora in a popular science domain. In *6th International Conference on Language Resources and Evaluations (LREC) - Building and using Comparable Corpora workshop*, 27–32.
- SINCLAIR, J. 1996. The search for units of meaning. *Textus* 9. 75–106.
- SMADJA, F. 1993. Retrieving collocations from text: Xtract. *Computational linguistics* 19. 143–177.
- URIZAR, RUBEN, 2012. *Euskal lokuzioen tratamendu konputazionala*. Donostia: Informatika Fakultatea, UPV/EHU tesia.
- VAN DE CRUYS, T., eta B.V. MOIRÓN. 2007. Semantics-based multiword expression extraction. In *Proceedings of the Workshop on A Broader Perspective on Multiword Expressions*, 25–32, Praga. ACL.
- WULFF, S. 2008. *Rethinking Idiomaticity*. Corpus and Discourse. New York: Continuum International Publishing Group Ltd.

Aipamenak

Tesi-lan hau Elhuyar Fundazioaren KONBITZ eta KONBITZ2 proiektuen testuinguruan egin da (EJren Ekonomia Garapena eta Lehiakortasuna sailaren SAIOTEK 2011 eta SAIOTEK 2012 programak). Ikerkuntzaren emaitzak MWE 2011 (ACL), LREC 2012 eta MWE 2013 (NAACL) nazioarteko biltzarretan aurkeztu dira ([Gurrutxaga eta Alegria, 2011](#), [2012](#), [2013](#)).

Elhuyarko I+Gko eta Ixa taldeko hainbat adituk lagundu digute ikerketan honetan. Eskerrik asko denoi.