

## Idiomatikotasunaren karakterizazio automatikoa: izen+aditza konbinazioak

**Doktoregaia:** Antton Gurrutxaga Hernaiz

**Zuzendariak:** Iñaki Alegria Loinaz, Xabier Artola Zubillaga

Lengoaia eta Sistema Informatikoak Saila  
Euskal Herriko Unibertsitatea / Universidad del País Vasco

2014ko uztailearen 21a

## urratsak egiten, pausoak ematen...



*eu urrats = en step*

## urratsak egiten, pausoak ematen...



eu *urrats* = en *step*  
⇒  
en to **do/make/give?** *steps*

urratsak egiten, pausoak ematen...



Martin Luther King

to pull sb's leg = adarra jo



I'm pulling your leg!

≠



Adarra jotzen ari naiz!

to pull sb's leg = adarra jo



≠



I'm pulling your leg!

Adarra jotzen ari naiz!

erabilera "idiomatikoa"

to put sb's leg  
(to tease, to joke)

⇓

=

norbaiti adarra jo  
(burla egin, iseka egin)

## Unitate fraseologikoez ari gara!

## Unitate fraseologikoez ari gara!

hiztunok erabiltzen ditugun hitz-konbinazio  
"preferentzialak" edo unitate "aurrefabrikatuak"

konbinazio idiosinkratikoak,  
sistemaren gramatika-arauen zein  
semantikaren arabera soilik  
ezin aurreikus edo azal  
daitezkeenak

ezaugarri semantiko, morfosintaktiko  
eta lexikal bereziak dituztenak



## Testuingurua

- ▶ Euskaraz, HAULen errepresentazioa/identifikazioa, termino-erauzketa eta entitateen arloa landu dira, batez ere
  - **Euskarazko UFak erauzteko eta karakterizatzeko teknologia garatu beharra**
- ▶ Testuinguru orokorra: fraseologia konputazionalan, UFak karakterizatzeko ohiko metodoa osagaien agerkidetzan oinarritua da
  - **Azkenaldiko teknika berriak esperimentatzeko eta garatzeko premia**

## Helburuak

### Helburu nagusia

Corpusetatik *izena+aditza* osaerako unitate fraseologikoak (UFak) automatikoki eskuratzeko eta haien idiomatikotasunaren arabera karakterizatzeko teknikak ikertzea, garatzea eta konbinatzea

## Helburuak

### Helburu zehatzak

- 1 UFen **idiomatikotasunaren** definizio operatiboa: haren osagai diren propietate neurgarriak zehaztea eta UFen sailkapen-eredua lantzea
- 2 UFen erauzketa eta karakterizazio automatikoaren **ebaluazio-metodologia** zehaztea, eta behar diren baliabideak eratzea
- 3 Ikergetzat ditugun euskarazko **izena+aditza osaerako unitateen** ezaugarriak deskribatzea
- 4 Idiomatikotasunaren propietate bakoitza neurtzeko **lan esperimentalak**, eta horien emaitzak konbinatzea
- 5 Emaitzak analizatzea eta **ondorioak** ateratzea, batez ere, euskarazko **baliabide lexikal konputazionalak** eratzeari eta **hiztegitantzari** begira

## Helburuak

### Galderak

- ▶ Zenbaterainoko korrelazioa dago idiomatikotasun-mailaren eta haren propietate bakoitzaren neurketen artean?
- ▶ Zenbateraino datoz bat UFen propietateen ebidentzia enpirikoak teoria fraseologikoak UFetarako oro har zein UF-kategoria bakoitzerako aurretan duenarekin?
- ▶ Hobetu daiteke UFen karakterizazioa idiomatikotasunaren propietate bakunen kuantifikazioaren emaitzak konbinatuz?

## Aurkezpenaren eskema

- 1 UFen idiomatikotasunaren eta karakterizazioaren marko teorikoa
- 2 UFen erauzketa eta karakterizazio automatikorako teknikak
- 3 Lan esperimentalak eta emaitzen analisia
- 4 Ondorioak, ekarpenak eta etorkizuneko lanak

## Aurkezpenaren eskema

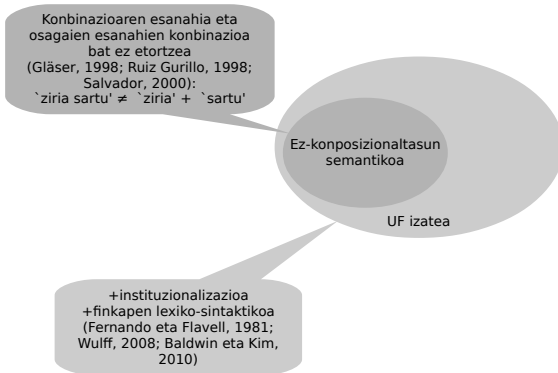
- 1 UFen idiomatikotasunaren eta karakterizazioaren marko teorikoa
- 2 UFen erauzketa eta karakterizazio automatikorako teknikak
- 3 Lan esperimentalak eta emaitzen analisia
- 4 Ondorioak, ekarpenak eta etorkizuneko lanak

## Idiomatikotasunaren definizioak

Konbinazioaren esanahia eta osagaien esanahien konbinazioa bat ez etortzea (Gläser, 1998; Ruiz Gurillo, 1998; Salvador, 2000):  
'ziria sartu' ≠ 'ziria' + 'sartu'

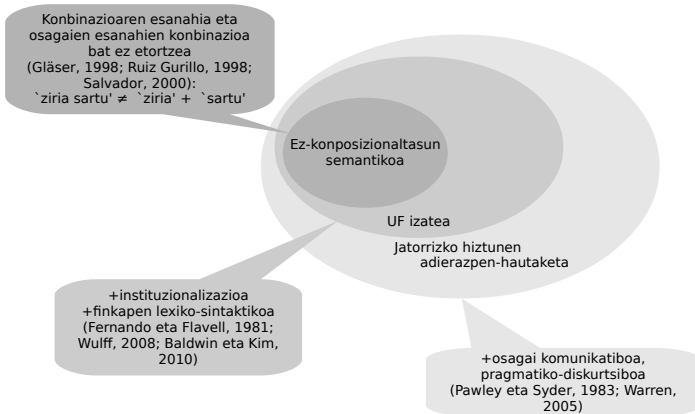
Ez-konposizionaltasun semantikoa

# Idiomatikotasunaren definizioak

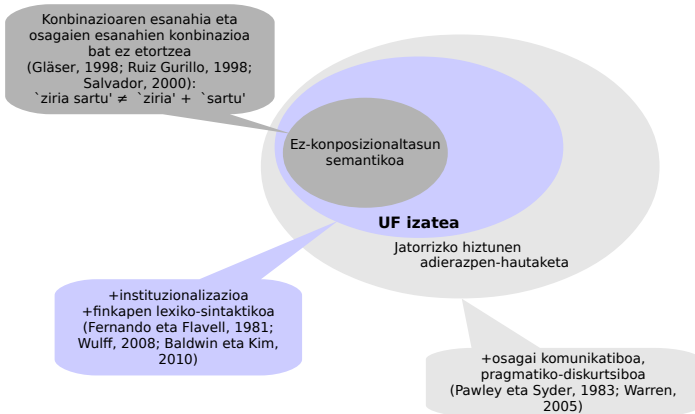




## Idiomatikotasunaren definizioak



## Idiomatikotasunaren definizioak



## Idiomatikotasunaren definizioa

### Gure hautua

- ▶ Idiomatikotasuna konbinazio bat UF izatea determinatzen duen propietatea da, muineko ezaugarritzat **idiosinkrasia** duena
- ▶ Idiomatikotasuna konplexua eta graduala da, eta bere barnean zenbait propietate hartzen ditu: **instituzionalizazioa**, **ez-konposizionaltasun semantikoa** (osoa edo partziala), eta **finkapena** (morfosintaktikoa zein lexikala)

## Idiomatikotasunaren osagaiak

- ▶ **Instituzionalizazioa:** adiera hertsian, idiosinkrasia estatistikoa (konbinazioa ausazkoa ez izatea)
  - ▶ behatutako maiztasuna > osagaiak ausaz konbinatuz itxaron daitekeen maiztasuna
- ▶ Ez-konposizionaltasuna: konbinazioaren esanahia  $\neq$  osagaien esanahien konbinazioa. Konposizionaltasun-mailak:
  - ▶ 'adarra jo'  $\neq$  'adarra' + 'jo'
  - ▶ 'zubiak eraiki'  $\leftarrow$  'zubiak' + 'eraiki'
  - ▶ 'atentzioa eman' = 'atentzioa' + 'eman<sub>atentzio</sub>'
- ▶ Finkapen morfosintaktikoa: aldakuntza batzuk izateko murriztapenak
  - ▶ \**adarra ipini / adarrak ipini; \*adar ederra jo nion / adarra ederki jo nion*
- ▶ Finkapen lexikala: osagaiak sinonimoez ordezkatzeko murriztapenak
  - ▶ \**meza aditu / meza entzun; hanka sartu  $\neq$  zangoa sartu*

## Idiomatikotasunaren osagaiak

- ▶ **Instituzionalizazioa:** adiera hertsian, idiosinkrasia estatistikoa (konbinazioa ausazkoa ez izatea)
  - ▶ behatutako maiztasuna > osagaiak ausaz konbinatuz itxaron daitekeen maiztasuna
- ▶ **Ez-konposizionaltasuna:** konbinazioaren esanahia  $\neq$  osagaien esanahien konbinazioa. Konposizionaltasun-mailak:
  - ▶ 'adarra jo'  $\neq$  'adarra' + 'jo'
  - ▶ 'zubiak eraiki'  $\leftarrow$  'zubiak' + 'eraiki'
  - ▶ 'atentzioa eman' = 'atentzioa' + 'eman<sub>atentzio</sub>'
- ▶ **Finkapen morfosintaktikoa:** aldakuntza batzuk izateko murriztapenak
  - ▶ \**adarra ipini / adarrak ipini; \*adar ederra jo nion / adarra ederki jo nion*
- ▶ **Finkapen lexikala:** osagaiak sinonimoez ordezkatzeko murriztapenak
  - ▶ \**meza aditu / meza entzun; hanka sartu  $\neq$  zangoa sartu*

## Idiomatikotasunaren osagaiak

- ▶ **Instituzionalizazioa:** adiera hertsian, idiosinkrasia estatistikoa (konbinazioa ausazkoa ez izatea)
  - ▶ behatutako maiztasuna > osagaiak ausaz konbinatuz itxaron daitekeen maiztasuna
- ▶ **Ez-konposizionaltasuna:** konbinazioaren esanahia  $\neq$  osagaien esanahien konbinazioa. Konposizionaltasun-mailak:
  - ▶ 'adarra jo'  $\neq$  'adarra' + 'jo'
  - ▶ 'zubiak eraiki'  $\leftarrow$  'zubiak' + 'eraiki'
  - ▶ 'atentzioa eman' = 'atentzioa' + 'eman<sub>atentzio</sub>'
- ▶ **Finkapen morfosintaktikoa:** aldakuntza batzuk izateko murriztapenak
  - ▶ \**adarra ipini / adarrak ipini; \*adar ederra jo nion / adarra ederki jo nion*
- ▶ **Finkapen lexikala:** osagaiak sinonimoez ordezkatzeko murriztapenak
  - ▶ \**meza aditu / meza entzun; hanka sartu  $\neq$  zangoa sartu*

## Idiomatikotasunaren osagaiak

- ▶ **Instituzionalizazioa:** adiera hertsian, idiosinkrasia estatistikoa (konbinazioa ausazkoa ez izatea)
  - ▶ behatutako maiztasuna > osagaiak ausaz konbinatuz itxaron daitekeen maiztasuna
- ▶ **Ez-konposizionaltasuna:** konbinazioaren esanahia  $\neq$  osagaien esanahien konbinazioa. Konposizionaltasun-mailak:
  - ▶ 'adarra jo'  $\neq$  'adarra' + 'jo'
  - ▶ 'zubiak eraiki'  $\leftarrow$  'zubiak' + 'eraiki'
  - ▶ 'atentzioa eman' = 'atentzioa' + 'eman<sub>atentzio</sub>'
- ▶ **Finkapen morfosintaktikoa:** aldakuntza batzuk izateko murriztapenak
  - ▶ \**adarra ipini / adarrak ipini; \*adar ederra jo nion / adarra ederki jo nion*
- ▶ **Finkapen lexikala:** osagaiak sinonimoez ordezkatzeko murriztapenak
  - ▶ \**meza aditu / meza entzun; hanka sartu  $\neq$  zangoa sartu*

## Idiomatikotasun-continuuma

- ▶ Idiomatikotasunaren propietateak hein desberdinean konbinatzen dira, eta UFak sailkatzeko irizpideak ezartzearen emaitza graduazio moduko continuum bat da (Sinclair, 1996; Ruiz Gurillo, 1998; Bannard et al., 2003; Katz eta Giesbrecht, 2006; Wulff, 2008)



## Idiomatikotasun-continuumak

- ▶ Idiomatikotasunaren propietateak hein desberdinean konbinatzen dira, eta UFAk sailkatzeko irizpideak ezartzearen emaitza graduazio moduko continuum bat da (Sinclair, 1996; Ruiz Gurillo, 1998; Bannard et al., 2003; Katz eta Giesbrecht, 2006; Wulff, 2008)

BAINA

- ▶ Continuum horretan eremu edo zona desberdinak proposatu dira:  
*‘Collocations’ and ‘idioms’ represent two large and amorphous subgroups of FEIs on a continuum.* (Moon, 1998)

## Sintagma-erako UFak sailkatzeko proposamen batzuk

Author	Opaque, invariable unit	Partially motivated unit	Phraseologically bound unit
Vinogradov (1947)	Phraseological function	Phraseological unity	Phraseological communication
Amosova (1963)	Idiom	Idiom (nor differentiated)	Phraseme, or phraseoloid
Cowie (1981)	Pure idiom	Figurative idiom	Restricted collocation
Mel'čuk (1988)	Idiom	Idiom (nor differentiated)	Collocation
Gläser (1988)	Idiom	Idiom (nor differentiated)	Restricted collocation
Howarth (1996)	Pure idiom	Figurative idiom	Restricted collocation
Corpas Pastor (1996)	Locución	Locución	Colocación
Burger (1998)	Idiom	Partial idiom	Collocation

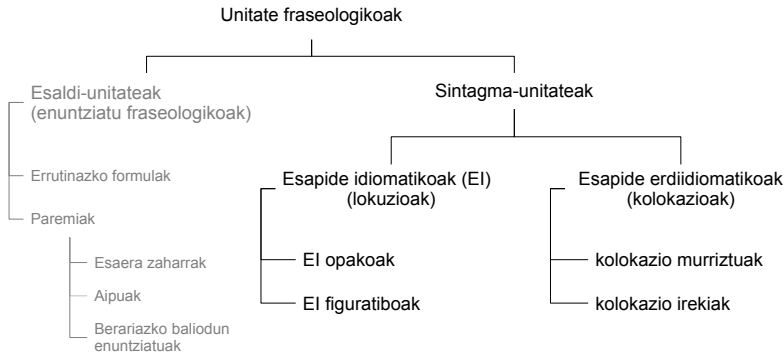
## Sintagma-erako UFak sailkatzeko proposamen batzuk

Author	Opaque, invariable unit	Partially motivated unit	Phraseologically bound unit
Vinogradov (1947)	Phraseological function	Phraseological unity	Phraseological communication
Amosova (1963)	Idiom	Idiom (nor differentiated)	Phraseme, or phraseoloid
Cowie (1981)	Pure idiom	<b>Figurative idiom</b>	Restricted collocation
Mel'čuk (1988)	Idiom	Idiom (nor differentiated)	Collocation
Gläser (1988)	Idiom	Idiom (nor differentiated)	Restricted collocation
Howarth (1996)	Pure idiom	<b>Figurative idiom</b>	Restricted collocation
Corpas Pastor (1996)	Locución	Locución	Colocación
Burger (1998)	Idiom	<b>Partial idiom</b>	Collocation

## Sintagma-erako UFak sailkatzeko proposamen batzuk

Author	Opaque, invariable unit	Partially motivated unit	Phraseologically bound unit
Vinogradov (1947)	Phraseological function	Phraseological unity	Phraseological communication
Amosova (1963)	Idiom	Idiom (nor differentiated)	Phraseme, or phraseoloid
Cowie (1981)	Pure idiom	Figurative idiom	<b>Restricted collocation</b>
Mel'čuk (1988)	Idiom	Idiom (nor differentiated)	Collocation
Gläser (1988)	Idiom	Idiom (nor differentiated)	<b>Restricted collocation</b>
Howarth (1996)	Pure idiom	Figurative idiom	<b>Restricted collocation</b>
Corpas Pastor (1996)	Locución	Locución	Colocación
Burger (1998)	Idiom	Partial idiom	Collocation

## UFen sailkapen-eredua



Non da esapide idiomatikoaren eta kolokazioaren arteko muga, edo zertan da bien arteko alde zehatza?

## Esapide idiomatikoak (lokuzioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, **maiztasun txikiagokoak kolokazioak baino** (Moon, 1998)
- ▶ Ez-konposizionaltasuna: erabatekoa edo handia
  - ▶ Opakoak (deskodetzeko arazoak): *adarra jo, ziria sartu, erreka jo, gorriak ikusi*
  - ▶ Figuratiboak (deskodetzeko aukera; deskonposagarriak): *burua hautsi, garunak urtu, zubiak eraiki*
- ▶ Finkapen morfosintaktikoa: aski handia; batzuk aldaezinak (*alde egin*); beste batzuk, gehienak, malgutasun gutxikoak (*adarra jo, gerrikoa estutu*)
- ▶ Finkapen lexikala: aski handia, osagaiak ordezkatzeko murriztapenak; batez ere, opakoetan (*hanka sartu  $\neq$  zangoa sartu*); hala ere: *burua makurtu/apaldu; hautsak harrotu/eraiki*

## Esapide idiomatikoak (lokuzioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, maiztasun txikiagokoak kolokazioak baino (Moon, 1998)
- ▶ Ez-konposizionaltasuna: **erabatekoa edo handia**
  - ▶ Opakoak (deskodetzeko arazoak): *adarra jo, ziria sartu, erreka jo, gorriak ikusi*
  - ▶ Figuratiboak (deskodetzeko aukera; deskonposagarriak): *burua hautsi, garunak urtu, zubiak eraiki*
- ▶ Finkapen morfosintaktikoa: aski handia; batzuk aldaezinak (*alde egin*); beste batzuk, gehienak, malgutasun gutxikoak (*adarra jo, gerrikoa estutu*)
- ▶ Finkapen lexikala: aski handia, osagaiak ordezkatzeko murriztapenak; batez ere, opakoetan (*hanka sartu  $\neq$  zangoa sartu*); hala ere: *burua makurtu/apaldu; hautsak harrotu/eraiki*

## Esapide idiomatikoak (lokuzioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, maiztasun txikiagokoak kolokazioak baino (Moon, 1998)
- ▶ Ez-konposizionaltasuna: erabatekoa edo handia
  - ▶ Opakoak (deskodetzeko arazoak): *adarra jo, ziria sartu, erreka jo, gorriak ikusi*
  - ▶ Figuratiboak (deskodetzeko aukera; deskonposagarriak): *burua hautsi, garunak urtu, zubiak eraiki*
- ▶ Finkapen morfosintaktikoa: **aski handia**; batzuk aldaezinak (*alde egin*); beste batzuk, gehienak, malgutasun gutxikoak (*adarra jo, gerrikoa estutu*)
- ▶ Finkapen lexikala: aski handia, osagaiak ordezkatzeko murriztapenak; batez ere, opakoetan (*hanka sartu  $\neq$  zangoa sartu*); hala ere: *burua makurtu/apaldu; hautsak harrotu/eraiki*



## Esapide idiomatikoak (lokuzioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, maiztasun txikiagokoak kolokazioak baino (Moon, 1998)
- ▶ Ez-konposizionaltasuna: erabatekoa edo handia
  - ▶ Opakoak (deskodetzeko arazoak): *adarra jo, ziria sartu, erreka jo, gorriak ikusi*
  - ▶ Figuratiboak (deskodetzeko aukera; deskonposagarriak): *burua hautsi, garunak urtu, zubiak eraiki*
- ▶ Finkapen morfosintaktikoa: aski handia; batzuk aldaezinak (*alde egin*); beste batzuk, gehienak, malgutasun gutxikoak (*adarra jo, gerrikoa estutu*)
- ▶ Finkapen lexikala: **aski handia**, osagaiak ordezkatzeko murriztapenak; batez ere, opakoetan (*hanka sartu  $\neq$  zangoa sartu*); hala ere: *burua makurtu/apaldu; hautsak harrotu/eraiki*

## Esapide erdiidiomatikoak (kolokazioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, esapide idiomatikoak baino **maiztasun handi**agoak; eskola britainiarrean, ezaugarri zentrala da
- ▶ Ez-konposizionaltasuna: oro har, erdikonposizionalak. Polartasuna (Hausmann, 1989):
  - ▶ Oinarriak “ohiko” edo oinarritzko adiera atxikitzen du
  - ▶ Kolokatiboa: ahula da (euskarri-aditzak: *etzula egin, lotsa eman*), edo kolokazioari espezifikoa zaion adiera du (*gola sartu, zarata atera*)
- ▶ Finkapen morfosintaktikoa: apala, ez da ezaugarri nabarmena (Heid, 1994)
- ▶ Finkapen lexikala: ezaugarri bereizgarria (Melčuk, 1994); izaera polarrarekin erlazionatua
  - ▶ Kolokazio murriztuetan, kolokatiboa ezin da ordezkatu, edo oso aukera gutxi daude
  - ▶ Kolokazio irekiaren kategoria eztabaidatua da (Cowie, 1998: *pay one's respects/a compliment/court to sb*); hautatze-murrizketak (Fontenelle, 1998)

## Esapide erdiidiomatikoak (kolokazioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, esapide idiomatikoak baino maiztasun handiagokoak; eskola britainiarrean, ezaugarri zentrala da
- ▶ Ez-konposizionaltasuna: oro har, **erdikonposizionalak**. Polartasuna (Hausmann, 1989):
  - ▶ Oinarriak “ohiko” edo oinarritzko adiera atxikitzen du
  - ▶ Kolokatiboa: ahula da (euskarri-aditzak: *etzula egin, lotsa eman*), edo kolokazioari espezifikoa zaion adiera du (*gola sartu, zarata atera*)
- ▶ Finkapen morfosintaktikoa: apala, ez da ezaugarri nabarmena (Heid, 1994)
- ▶ Finkapen lexikala: ezaugarri bereizgarria (Melčuk, 1994); izaera polarrarekin erlazionatua
  - ▶ Kolokazio murriztuetan, kolokatiboa ezin da ordezkatu, edo oso aukera gutxi daude
  - ▶ Kolokazio irekiaren kategoria eztabaidatua da (Cowie, 1998: *pay one's respects/a compliment/court to sb*); hautatze-murrizketak (Fontenelle, 1998)

## Esapide erdiidiomatikoak (kolokazioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, esapide idiomatikoak baino maiztasun handiagokoak; eskola britainiarrean, ezaugarri zentrala da
- ▶ Ez-konposizionaltasuna: oro har, erdikonposizionalak. Polartasuna (Hausmann, 1989):
  - ▶ Oinarriak “ohiko” edo oinarritzko adiera atxikitzen du
  - ▶ Kolokatiboa: ahula da (euskarri-aditzak: *eztula egin, lotsa eman*), edo kolokazioari espezifikoa zaion adiera du (*gola sartu, zarata atera*)
- ▶ Finkapen morfosintaktikoa: **apala**, ez da ezaugarri nabarmena (Heid, 1994)
- ▶ Finkapen lexikala: ezaugarri bereizgarria (Melčuk, 1994); izaera polarrarekin erlazionatua
  - ▶ Kolokazio murriztuetan, kolokatiboa ezin da ordezkatu, edo oso aukera gutxi daude
  - ▶ Kolokazio irekiaren kategoria eztabaidatua da (Cowie, 1998: *pay one's respects/a compliment/court to sb*); hautatze-murrizketak (Fontenelle, 1998)

## Esapide erdiidiomatikoak (kolokazioak)

### Ezaugarri nagusiak

- ▶ Idiosinkrasia estatistikoa: oro har, esapide idiomatikoak baino maiztasun handiagokoak; eskola britainiarrean, ezaugarri zentrala da
- ▶ Ez-konposizionaltasuna: oro har, erdikonposizionalak. Polartasuna (Hausmann, 1989):
  - ▶ Oinarriak “ohiko” edo oinarritzko adiera atxikitzen du
  - ▶ Kolokatiboa: ahula da (euskarri-aditzak: *etzula egin, lotsa eman*), edo kolokazioari espezifikoa zaion adiera du (*gola sartu, zarata atera*)
- ▶ Finkapen morfosintaktikoa: apala, ez da ezaugarri nabarmena (Heid, 1994)
- ▶ Finkapen lexikala: **ezaugarri bereizgarria** (Melčuk, 1994); izaera polarrarekin erlazionatua
  - ▶ Kolokazio murriztuetan, kolokatiboa ezin da ordezkatu, edo oso aukera gutxi daude
  - ▶ Kolokazio irekiaren kategoria eztabaidatua da (Cowie, 1998: *pay one's respects/a compliment/court to sb*); hautatze-murrizketak (Fontenelle, 1998)

## Aurkezpenaren eskema

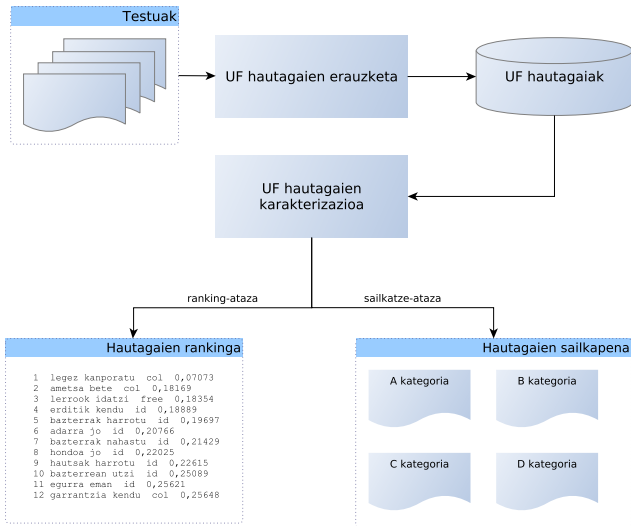
- 1 UFen idiomatikotasunaren eta karakterizazioaren marko teorikoa
- 2 UFen erauzketa eta karakterizazio automatikorako teknikak
- 3 Lan esperimentalak eta emaitzen analisia
- 4 Ondorioak, ekarpenak eta etorkizuneko lanak

## UFen erauzketa, fraseologia konputazionalaren egitekoetako bat

### Atazak (Urizar, 2012)

- ▶ **Erauzketa:** testuetatik UFak eskuratzea, eskuarki lexikoiak edo hiztegiak elikatzeko
- ▶ Identifikazioa: UFen testuetako banakako agerpenak atzematea
- ▶ Interpretazioa: UFen barne-sintaxia eta semantika desanbiguatzea
- ▶ Deskribapena: UF bakoitzak izan ditzakeen bariazio morfosintaktikoak zehaztea

# Urratsak





## Hautagaien erauzketa

### Konbinazioak sortzeko zenbait irizpide

- ▶ Konbinazioaren osaera
  - ▶ Lemak, hitz-formak, bestelako informazio morfosintaktikoa (kasua, mugatasuna...)
- ▶ Konbinazioak osatzeko metodoak
  - ▶ Distantzia: hitz baten inguruko “leiho-zabalera” (*window span*) baten barneko hitzekin osatzen dira konbinazioak
  - ▶ Erlazio morfosintaktiko jakin bat duten hitzen arteko konbinazioak

## Hautagaien erauzketa

### Konbinazioak sortzeko zenbait irizpide

- ▶ Konbinazioaren osaera
  - ▶ Lemak, hitz-formak, bestelako informazio morfosintaktikoa (kasua, mugatasuna...)
- ▶ Konbinazioak osatzeko metodoa
  - ▶ Distantzia: hitz baten inguruko “leiho-zabalera” (*window span*) baten barneko hitzekin osatzen dira konbinazioak
  - ▶ Erlazio morfosintaktiko jakin bat duten hitzen arteko konbinazioak



## Prozesamendu linguistikoa

### Prozesamendu-mailak

- ▶ Prozesamendu estandar minimoa: lematizazioa eta etiketatze morfosintaktikoa
  - ▶ Lema eta kasu bereko forma flexionatuen ondoriozko aldakuntzak forma kanoniko berarekin lotzeko

***Adarra jotzeko** gogoz, ala? Ez **adarrrik jo** niri.  
Atzoko partidaren hiru **gol sartu** zituen. **Gola sartzen** ahalegindu da. Horrela jarraituz gero, ez du **golik sartuko**. Nork **sartu** ditu gaurko bi **golak**?*
  - ▶ Kategorio-osaera jakineko konbinazioak sortzeko: izena+aditza, izena+izena, izena+izenondua...
- ▶ Prozesamendu aurreratua: azaleko edo sakoneko analisi sintaktikoa
  - ▶ Abantailak: doitasun handiagoko metodoa
  - ▶ Arazoak: analizatzaile sintaktikoetan hizkuntza bakoitzean erdietsi den kalitatearen mendeotasuna

## Idiomatikotasuna karakterizatzeko teknikak

Oinarrizko karakterizazio-teknika: idiosinkrasia estatistikoa agerkidetzaren bidez neurtzea

- ▶ Osagaien agerkidetzaren informazioa prozesatzea, elkartze-neurriak erabiliz (AM - *association measures*; Evert, 2005)
- ▶ Behatutako maiztasunak eredu desberdinen arabera estimatutako ausazko maiztasunekin konparatzea
- ▶ Ohiko AM batzuk:  $z$  neurria,  $t$  neurria,  $\chi^2$  (khi karratua), LLR (*log-likelihood ratio*, egiantz-arrazoiaren logaritmoa), Fisherren test zehatza, MI (*mutual information*, elkarrekiko informazioa),  $MI^3$ ,  $f$  (maiztasuna)

## Idiomatikotasuna karakterizatzeko teknikak

### Karakterizazio-teknika aurreratuak

- ▶ Konposizionaltasuna: konbinazioaren eta osagaien arteko **antzekotasun distribuzionala**
  - ▶ WSM (word space model, *hitz-espazioaren ereduak*): Fazly eta Stevenson (2007); Wulff (2008)
  - ▶ LSA (latent semantic analysis, *ezkutuko semantikaren analisiak*): Katz eta Giesbrecht (2006); Krčmár et al. (2013)
  - ▶ Adiera-indukzioa: Korkontzelos eta Manandhar (2009)
- ▶ Malgutasun morfosintaktikoa: hautagaien portaera morfosintaktikoa batez besteko erreferentzia-portaera batekin konparatzea (banaketen arteko distantzia)
  - ▶ Portaera orokorra: *izena+aditza* osaerako konbinazioen batez besteko portaera (Barkema, 1994; Fazly eta Stevenson, 2007; Wulff, 2008)
  - ▶ Osagaien portaera: konbinazioaren osagai batek beste osagaiaren kategoriako edozein hitzekin osatutako konbinazioen batez besteko portaera (Bannard, 2007)
- ▶ Malgutasun lexikala: osagaien ordezkagarritasuna
  - ▶ Lin (1999); Fazly eta Stevenson (2007); Van de Cruys eta Moirón (2007)

## Idiomatikotasuna karakterizatzeko teknikak

### Karakterizazio-teknika aurreratuak

- ▶ Konposizionaltasuna: konbinazioaren eta osagaien arteko antzekotasun distribuzionala
  - ▶ WSM (word space model, *hitz-espazioaren eredu*): Fazly eta Stevenson (2007); Wulff (2008)
  - ▶ LSA (latent semantic analysis, *ezkutuko semantikaren analisia*): Katz eta Giesbrecht (2006); Krčmár et al. (2013)
  - ▶ Adiera-indukzioa: Korkontzelos eta Manandhar (2009)
- ▶ Malgutasun morfosintaktikoa: hautagaien portaera morfosintaktikoa batez besteko **erreferentzia-portaera batekin konparatzea** (banaketen arteko distantzia)
  - ▶ Portaera orokorra: *izena+aditza* osaerako konbinazioen batez besteko portaera (Barkema, 1994; Fazly eta Stevenson, 2007; Wulff, 2008)
  - ▶ Osagaien portaera: konbinazioaren osagai batek beste osagaiaren kategoriako edozein hitzekin osatutako konbinazioen batez besteko portaera (Bannard, 2007)
- ▶ Malgutasun lexikala: osagaien ordezkagarritasuna
  - ▶ Lin (1999); Fazly eta Stevenson (2007); Van de Cruys eta Moirón (2007)

## Idiomatikotasuna karakterizatzeko teknikak

### Karakterizazio-teknika aurreratuak

- ▶ Konposizionaltasuna: konbinazioaren eta osagaien arteko antzekotasun distribuzionala
  - ▶ WSM (word space model, *hitz-espazioaren eredu*): Fazly eta Stevenson (2007); Wulff (2008)
  - ▶ LSA (latent semantic analysis, *ezkutuko semantikaren analisia*): Katz eta Giesbrecht (2006); Krčmár et al. (2013)
  - ▶ Adiera-indukzioa: Korkontzelos eta Manandhar (2009)
- ▶ Malgutasun morfosintaktikoa: hautagaien portaera morfosintaktikoa batez besteko erreferentzia-portaera batekin konparatzea (banaketen arteko distantzia)
  - ▶ Portaera orokorra: *izena+aditza* osaerako konbinazioen batez besteko portaera (Barkema, 1994; Fazly eta Stevenson, 2007; Wulff, 2008)
  - ▶ Osagaien portaera: konbinazioaren osagai batek beste osagaiaren kategoriako edozein hitzekin osatutako konbinazioen batez besteko portaera (Bannard, 2007)
- ▶ Malgutasun lexikala: **osagaien ordezkagarritasuna**
  - ▶ Lin (1999); Fazly eta Stevenson (2007); Van de Cruys eta Moirón (2007)

## Sailkapen automatikoa

- ▶ Elkartze-neurrien emaitzak konbinatzea. Pecina eta Schlesinger (2006); Lin et al. (2008):  $f$ ,  $t$  neurria, LLR, MI eta  $\chi^2$  (ezaugarri guztien konbinazioa da onena, eta hurrena  $t$  neurria).
- ▶ Agerkidetzaren eta konposizionaltasunaren neurketak konbinatzea. Venkatapathy eta Joshi (2005): ekarpen handiena kolokazioaren izenari dagokion aditzaren eta kolokazioaren arteko antzekotasunak egiten du.
- ▶ Lau propietateen neurketak konbinatzea. Fazly eta Stevenson (2007): konbinazioaren aditza ere erabiltzen dute; ekarpen handiena malgutasun sintaktikoarekin eta lexikalarekin lotutako ezaugarriek egiten dute.



## Aurkezpenaren eskema

- 1 UFen idiomatikotasunaren eta karakterizazioaren marko teorikoa
- 2 UFen erauzketa eta karakterizazio automatikorako teknikak
- 3 Lan esperimentalak eta emaitzen analisia
- 4 Ondorioak, ekarpenak eta etorkizuneko lanak

## Esperimentuen diseinuaren elementuak

- 1 Idiomatikotasunaren propietateekin lotutako **behagaiak** aukeratzea, eta horietako bakoitza neurtzeko estrategia zehaztea
- 2 Esperimentuetako karakterizazio-atazak deskribatzea
- 3 Ikergaia zehaztea: *izena+aditza* konbinazioak definitzea
- 4 Corpusetik *izena+aditza* konbinazio hautagaiak erauzteko prozesua
- 5 Ebaluazio-metodologia antolatzea, eta behar diren erreferentzia-baliabideak garatzea

## Esperimentuen diseinuaren elementuak

- 1 Idiomatikotasunaren propietateekin lotutako behagaiak aukeratzea, eta horietako bakoitza neurtzeko estrategia zehaztea
- 2 Esperimentuetako **karakterizazio-atatak** deskribatzea
- 3 Ikergaia zehaztea: *izena+aditza* konbinazioak definitzea
- 4 Corpusetik *izena+aditza* konbinazio hautagaiak erauzteko prozesua
- 5 Ebaluazio-metodologia antolatzea, eta behar diren erreferentzia-baliabideak garatzea

## Esperimentuen diseinuaren elementuak

- 1 Idiomatikotasunaren propietateekin lotutako behagaiak aukeratzea, eta horietako bakoitza neurtzeko estrategia zehaztea
- 2 Esperimentuetako karakterizazio-atazak deskribatzea
- 3 Ikergaia zehaztea: **izena+aditza konbinazioak** definitzea
- 4 Corpusetik *izena+aditza* konbinazio hautagaiak erazteko prozesua
- 5 Ebaluazio-metodologia antolatzea, eta behar diren erreferentzia-baliabideak garatzea

## Esperimentuen diseinuaren elementuak

- 1 Idiomatikotasunaren propietateekin lotutako behagaiak aukeratzea, eta horietako bakoitza neurtzeko estrategia zehaztea
- 2 Esperimentuetako karakterizazio-atazak deskribatzea
- 3 Ikergaia zehaztea: *izena+aditza* konbinazioak definitzea
- 4 Corpusetik *izena+aditza* konbinazio **hautagaiak erazteko prozesua**
- 5 Ebaluazio-metodologia antolatzea, eta behar diren erreferentzia-baliabideak garatzea

## Esperimentuen diseinuaren elementuak

- 1 Idiomatikotasunaren propietateekin lotutako behagaiak aukeratzea, eta horietako bakoitza neurtzeko estrategia zehaztea
- 2 Esperimentuetako karakterizazio-atazak deskribatzea
- 3 Ikergaia zehaztea: *izena+aditza* konbinazioak definitzea
- 4 Corpusetik *izena+aditza* konbinazio hautagaiak erauzteko prozesua
- 5 **Ebaluazio-metodologia** antolatzea, eta behar diren erreferentzia-baliabideak garatzea

## Behagaiak

Hipotesia da behagai hauen bidez neur daitezkeela idiomatikotasunaren lau propietateak edo osagaiak:

<b>Propietatea</b>	<b>Behagaia</b>
Instituzionalizazioa (idiosinkrasia estatistikoa)	Agerkidetza
Konposizionaltasun-maila	Antzekotasun distribuzionala
Malgutasun morfosintaktikoa	Erreferentzia-portaera batekiko distantzia
Malgutasun lexikala	Osagaien ordezkagarritasuna

## Karakterizazio-atazak

### Ranking-ataza

- ▶ Continuumaren ideian oinarrituta, UF hautagaiak idiomatikotasun-mailaren arabera ordenatzea, rankingak antolatzea
- Idiomatikotasunaren propietateen banakako neurketa

### Sailkatze-ataza

- ▶ UF hautagaiak sailkapen-ereduan bereizi ditugun kategorietan banatzea
- Propietateen banakako neurketen konbinazioa, ikasketa automatikoan



## izena+aditza osaerako konbinazioen errepertorioa

### Zabala (2004), Altzibar (2005) eta Urizarren (2012) lanetan oinarritua

1 izena<sub>subjektua</sub> + aditza

1.1 Absolutibodun sintagma: *burua joan, eguzkia sartu, ardoa garrastu, esnea galdu*

1.2 Ergatibodun sintagma: *gogoak eman, loak hartu, suak hartu, ilunak jo*

2 izena<sub>objektua</sub> + aditza: *hanka sartu, zubiak eraiki, lan egin, min hartu, bizarra egin, itxurak egin, aukera eman, erabakia hartu, zarata atera, eskerrak eman, urratsak egin, adostasuna lortu, gola sartu, elkartasuna adierazi*

3 izena<sub>subjektuaren pred.</sub> + aditza: *beldur izan, falta izan, giro egon*

4 izena<sub>objektuaren pred.</sub> + aditza: *atsegin ukan, damu ukan*

5 izena<sub>datiboa</sub> + aditza: *edanari eman, bideari ekin, lanari lotu*

6 izena<sub>adjuntua</sub> + aditza: *mendean hartu, borrokan sartu, martxan jarri, adarretatik heldu, larrutik ordaindu, burutik kendu, harira etorri, gogora ekarri, aurrera eraman, zerbitzura egon, sareetara bidali, muturreraino eraman, aurrez ikusi, arduraz jokatu, oinarritzat hartu*

## UF hautagaiak erauzteak

### Corpus-baliabideak eta prozesamendua

- ▶ Kazetaritza-corpus bat (75 milioi hitz), bi iturritatik eratua
  - ▶ *Euskaldunon Egunkaria*: 2001-2002 (28 milioi hitz)
  - ▶ *Berría*: 2006-2010 (47 milioi hitz)
- ▶ Etiketatzeko linguistikoa: lematizazioa eta analisi morfologikoa
  - ▶ UPV/EHUko Ixa taldearen Eustagger etiketzailea
- ▶ Eustaggerren irteeraren tratamendua
  - ▶ Izenen informazioa: forma lema kategoria azpikategoria kasua mugatasuna
  - ▶ Aditzen informazioa: printzipioz, lema kategoria; baina PART analisisa dutenetan, analisi-kateko zenbait informazio atxikitzen ditugu:
    - ▶ *herrialde aurreratuak, gobernuaren aliatuak* modukoak ez erauzteko
    - ▶ Partizipio jokatu gabearren gaineko erlatiboak bereiz itatzeko (*hartutako erabakien ondorioak*)

## UF hautagaiak erauzteak

### izena+aditza osaerako konbinazio hautagaiak lortzea

- ▶ **Bigrama-sorkuntza**
  - ▶ Ngram Statistics Package (NSP)<sup>a</sup> (Pedersen eta Banerjee, 2010)
  - ▶ Erauzte-parametroak
    - ▶ Maiztasun-atariak: 3, 10, 30, 50
    - ▶ Leiho-zabalera:  $w = \pm 1$ ,  $w = \pm 5$
- ▶ **Forma kanonikoa esleitzea: bigramen normalizazioa**
  - ▶ Aditz jakin batekin konbinaturik agertzen diren kasu bereko izen-formak normalizatu egin ditugu, maiztasun handieneko forma eta mugatasuna esleituz

---

<sup>a</sup><http://www.d.umn.edu/~tpederse/nsp.html>

## UF hautagaiak eraztea

*egunkari\_irakurri* lema-bikotearen kasu eta mugatasun desberdineko konbinazio batzuk

```
egunkarietan_egunkari_IZE_ARR_INE_NUMP<>irakurri_ADI<>31 74 1956  
egunkarian_egunkari_IZE_ARR_INE_NUMS<>irakurri_ADI<>43 76 1956  
egunkariak_egunkari_IZE_ARR_ABS_NUMP<>irakurri_ADI<>67 486 1956  
egunkaria_egunkari_IZE_ARR_ABS_NUMS<>irakurri_ADI<>126 1003 1956
```

Normalizazioaren emaitza

```
egunkarian_egunkari_IZE_ARR_INE_NUMS<>irakurri_ADI<>74 150 1956  
egunkaria_egunkari_IZE_ARR_ABS_NUMS<>irakurri_ADI<>193 1489 1956
```

## Ebaluazio-prozedura

- 1 Irizpideak
- 2 Hiztegi-erreferentzia
- 3 Ebaluazio-erreferentzia sailkatua
- 4 Metrikak

## Ebaluazioa: irizpideak

### Kontuan hartu beharreko alderdi nagusiak

- ▶ UF kategoriak bereizi nahi ditugu → **UF bai/ez sailkapena aski ez**
- ▶ Euskaraz UF-kategoriak bereizten dituen iturri lexikografikorik ez dago → adituen eskulanaren bidezko ebaluazio-erreferentzia bat eratu behar
- ▶ Antzekotasun distribuzionalaren eta malgutasun morfosintaktikoaren neurketek hautagaien testuinguru-informazio aberatsa behar dute → hautagaien agerpen-kopurua txikiegia ez izatea (20-50)
- ▶ Ikasketa automatikoa: sistemak ikasteko adinako kopuruan izan behar ditu UF-kategoriaren instantziak → moduren bat behar dugu ebaluazio-erreferentzian UFak gutxiegia ez izateko

## Ebaluazioa: irizpideak

### Kontuan hartu beharreko alderdi nagusiak

- ▶ UF kategoriak bereizi nahi ditugu → UF bai/ez sailkapena aski ez
- ▶ Euskaraz UF-kategoriak bereizten dituen iturri lexikografikorik ez dago → **adituen eskulanaren bidezko ebaluazio-erreferentzia** bat eratu behar
- ▶ Antzekotasun distribuzionalaren eta malgutasun morfosintaktikoaren neurketek hautagaien testuinguru-informazio aberatsa behar dute → hautagaien agerpen-kopurua txikiegia ez izatea (20-50)
- ▶ Ikasketa automatikoa: sistemak ikasteko adinako kopuruan izan behar ditu UF-kategorien instantziak → moduren bat behar dugu ebaluazio-erreferentzian UFac gutxiegia ez izateko

## Ebaluazioa: irizpideak

### Kontuan hartu beharreko alderdi nagusiak

- ▶ UF kategoriak bereizi nahi ditugu → UF bai/ez sailkapena aski ez
- ▶ Euskaraz UF-kategoriak bereizten dituen iturri lexikografikorik ez dago → adituen eskulanaren bidezko ebaluazio-erreferentzia bat eratu behar
- ▶ Antzekotasun distribuzionalaren eta malgutasun morfosintaktikoaren neurketek hautagaien testuinguru-informazio aberatsa behar dute → **hautagaien agerpen-kopurua txikiegia ez izatea** (20-50)
- ▶ Ikasketa automatikoa: sistemak ikasteko adinako kopuruan izan behar ditu UF-kategoriaren instantziak → moduren bat behar dugu ebaluazio-erreferentzian UFAk gutxiegia ez izateko



## Ebaluazioa: irizpideak

### Kontuan hartu beharreko alderdi nagusiak

- ▶ UF kategoriak bereizi nahi ditugu → UF bai/ez sailkapena aski ez
- ▶ Euskaraz UF-kategoriak bereizten dituen iturri lexikografikorik ez dago → adituen eskulanaren bidezko ebaluazio-erreferentzia bat eratu behar
- ▶ Antzekotasun distribuzionalaren eta malgutasun morfosintaktikoaren neurketek hautagaien testuinguru-informazio aberatsa behar dute → hautagaien agerpen-kopurua txikiegia ez izatea (20-50)
- ▶ Ikasketa automatikoa: sistemak ikasteko adinako kopuruan izan behar ditu UF-kategorien instantziak → moduren bat behar dugu **ebaluazio-erreferentzian UFak gutxiegia ez** izateko

## Hiztegi-erreferentzia

### Iturriak

- ▶ HB: Euskaltzaindiaren *Hiztegi Batua*<sup>a</sup>
- ▶ EH: Ibon Sarasolaren *Euskal Hiztegia*
- ▶ ELH: Elhuyar Fundazioaren *Euskara-Castellano/CastellanoVasco Hiztegia*<sup>b</sup>
- ▶ Intza: *Intza proiektua*<sup>c</sup>
- ▶ EDBL: Ixa taldearen EDBL datu-base lexikala<sup>d</sup>

<sup>a</sup><http://www.euskaltzaindia.net/hiztegibatua>; [http://www.euskaltzaindia.net/eaeb\\_gunetik\\_deskargatua](http://www.euskaltzaindia.net/eaeb_gunetik_deskargatua) (2010-06-04ko bertsioa).

<sup>b</sup>2010-06-02ko bertsioa.

<sup>c</sup><http://intza.armiarma.com> (2010-06-02ko bertsioa).

<sup>d</sup><http://ixa2.si.ehu.es/edbl> (2010-06-22ko bertsioa).

**Emaitza:** izena+aditza osaerako 3 720 UF

## Ebaluazio-erreferentzia

### Laginaren hautaketa

- ▶ Hautatutako erauzketa:  $f \geq 30$ ,  $w = \pm 1$ , bigramak normalizatuta
- ▶ UFen dentsitatea handitzeko: AMen araberako ranking bakoitzetik lehen  $n$  hautagaien multzoak hartu, eta horien bildura egin;  $n = 2\,000$  hautatu dugu. Bigrama-kopurua: 4 334
- ▶ 1 200 bigramako ausazko azpimultzoa aukeratu dugu, aditu-talde batek eskuz sailka dezan

## Sailkatze-lana

- ▶ Hiru hizkuntzalariz osatutako aditu-talde batek 5 kategoriatan sailkatu du (`id_op`, `id_fig`, `col_res`, `col_open`, `free`)
- ▶ Sailkatze-prozesurako irizpideak Krennen (2004) lanean inspiratuta eman ditugu
- ▶ `izena+aditza` ez diren konbinazioak (55) kenduta: 1 145eko lagina
- ▶ 5 kategoriako ereduari uko egin eta 3 kategoriakoa hobetsi dugu: `id (id_op, id_fig)`, `col (col_res, col_open)`, `free`. Arrazoiak:
  - ▶ Esapide idiomatiko opakoen kategoria oso instantzia gutxiri esleitu diete adituek
  - ▶ Kolokazio irekien kategoria oso labaina da, eta zalantza ugari sortu ditu sailkatzaileen artean
- ▶ ITA (etiketatzailen arteko adostasuna): Fleiss  $\kappa = 0,58$  (adostasun ertaina)
- ▶ Erreferentzia sailkatua: `id 80`; `col 268`; `free 797`

## Sailkatze-lana

- ▶ Hiru hizkuntzalariz osatutako aditu-talde batek 5 kategoriatan sailkatu du (`id_op`, `id_fig`, `col_res`, `col_open`, `free`)
- ▶ Sailkatze-prozesurako irizpideak Krennen (2004) lanean inspiratuta eman ditugu
- ▶ `izena+aditza` ez diren konbinazioak (55) kenduta: 1 145eko lagina
- ▶ 5 kategoriako ereduari uko egin eta 3 kategoriakoa hobetsi dugu: **id** (`id_op`, `id_fig`), **col** (`col_res`, `col_open`), **free**. Arrazoiak:
  - ▶ Esapide idiomatiko opakoen kategoria oso instantzia gutxiri esleitu diete adituek
  - ▶ Kolokazio irekien kategoria oso labaina da, eta zalantza ugari sortu ditu sailkatzaileen artean
- ▶ ITA (etiketatzailen arteko adostasuna): Fleiss  $\kappa = 0,58$  (adostasun ertaina)
- ▶ Erreferentzia sailkatua: `id 80`; `col 268`; `free 797`

## Erreferentzia sailkatuaren osaera

kategoria	Hiztegi-erreferentzia		Totala
	bai	ez	
id	23	57	80
col	41	227	268
free	10	787	797
Totala	74	1 071	1 145

Eskuz sailkatutako ebaluazio-erreferentziako bigramen eta hiztegi-erreferentziaren arteko konparazioa.

### Nabarmentzekoak

- ▶ Hiztegia aberasteko aukera handia (batez ere, col: % 15,3 daude hiztegietan)
- ▶ Predikatu konplexu batzuk (*jaramon egin*, *zor izan*) esapide erdiidiomatikoetan sailkatu dira, eta hori ez dator bat euskarazko fraseologian horrelakoak aditz-lokuziotzat jotzeko tradizioarekin

→ komenigarria da, kontrasterako, konbinazio horiek (17) esapide idiomatikotzat sailkatuta esperimentuen emaitzetan igarriko litzatekeen eragina aztertzea

## Ebaluazio-metriak

### Ranking-ataza

- ▶ Ideia da nolabait konparatzea idealki ordenatutako rankinga neurri bakoitzak sortutako rankingarekin
  - ▶ Kendall  $\tau_B$  heinen korrelazio-koefizientea, berdinketak kudeatzeko
- ▶ Batez besteko doitasunak ere kalkulatu dira (*AP* - *average precision*):  $AP_{UF}$ ,  $AP_{id}$  eta  $AP_{col}$ 
  - ▶ Neurri batek Ufekiko oro har eta UF-kategoria bakoitzarekiko duen portaera ikusteko

### Sailkatze-ataza

- ▶ Wekak eskaintzen dituen neurri hauek erabili ditugu:
  - ▶ Zuzen sailkatutako instantzia-kopurua (*Correctly Classified Instances* - CCI) Zehaztasunaren baliokidea da (*accuracy*)
  - ▶ Kategoría bakoitzeko  $F$  neurriak:  $F_{id}$ ,  $F_{col}$ ,  $F_{free}$
  - ▶ Batez besteko  $F$  neurri haztatua (*Weighted Average F-measure*) Mikrobatezbestekoa da ( $F_{mikro}$ )

Mikrobatezbestekoaz gain, makrobatezbestekoa ere ( $F_{makro}$ ) kalkulatu dugu

## Idiosinkrasia estatistikoaren neurketa

### Agerkidetza-datuen prozesamendua: AMen kalkulua

- ▶ NSPren irteera UCS toolkitaren bidez prozesatzea<sup>a</sup> (Evert, 2005)
- ▶ Erabilitako AMak:  $z$  neurria,  $t$  neurria,  $\chi^2$  (khi karratua), LLR (*log-likelihood ratio*, egiantz-arrazoiaren logaritmoa), Fisherren test zehatza, MI (*mutual information*, elkarrekiko informazioa),  $MI^3$ ,  $f$  (maiztasuna)

<sup>a</sup><http://www.collocations.de/software.html>



## Idiosinkrasia estatistikoaren neurketa

### Agerkidetza-datuen prozesamendua: AMen kalkulua

- ▶ NSPren irteera UCS toolkitaren bidez prozesatzea<sup>a</sup> (Evert, 2005)
- ▶ Erabilitako AMak:  $z$  neurria,  $t$  neurria,  $\chi^2$  (khi karratua), LLR (*log-likelihood ratio*, egiantz-arrazoiaren logaritmoa), Fisherren test zehatza, MI (*mutual information*, elkarrekiko informazioa),  $MI^3$ ,  $f$  (maiztasuna)

<sup>a</sup><http://www.collocations.de/software.html>

	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
ausazko rankinga	0,000	0,308	0,070	0,234
$z$ neurria	(-0,038)	0,297	0,109	0,204
$t$ neurria	0,197	0,455	0,084	0,383
$\chi^2$	(-0,037)	0,302	0,119	0,206
LLR	0,156	0,427	0,100	0,335
Fisher	0,156	0,426	0,100	0,335
MI	(-0,121)	0,257	0,086	0,182
$MI^3$	0,103	0,389	0,107	0,291
$f$	0,189	0,436	0,074	0,379

## Idiosinkrasia estatistikoaren neurketa

### Agerkidetza-datuaren prozesamendua: AMen kalkulua

- ▶ NSPren irteera UCS toolkitaren bidez prozesatzea<sup>a</sup> (Evert, 2005)
- ▶ Erabilitako AMak:  $z$  neurria,  $t$  neurria,  $\chi^2$  (khi karratua), LLR (*log-likelihood ratio*, egiantz-arrazoiaren logaritmoa), Fisherren test zehatza, MI (*mutual information*, elkarrekiko informazioa),  $MI^3$ ,  $f$  (maiztasuna)

<sup>a</sup><http://www.collocations.de/software.html>

	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
ausazko rankinga	0,000	0,308	0,070	0,234
$z$ neurria	(-0,038)	0,297	0,109	0,204
$t$ neurria	<b>0,197</b>	<b>0,455</b>	0,084	<b>0,383</b>
$\chi^2$	(-0,037)	0,302	0,119	0,206
LLR	0,156	0,427	0,100	0,335
Fisher	0,156	0,426	0,100	0,335
MI	(-0,121)	0,257	0,086	0,182
$MI^3$	0,103	0,389	0,107	0,291
$f$	0,189	0,436	0,074	0,379

## Idiosinkrasia estatistikoaren neurketa

### Agerkidetza-datuen prozesamendua: AMen kalkulua

- ▶ NSPren irteera UCS toolkitaren bidez prozesatzea<sup>a</sup> (Evert, 2005)
- ▶ Erabilitako AMak:  $z$  neurria,  $t$  neurria,  $\chi^2$  (khi karratua), LLR (*log-likelihood ratio*, egiantz-arrazoiaren logaritmoa), Fisherren test zehatza, MI (*mutual information*, elkarrekiko informazioa),  $MI^3$ ,  $f$  (maiztasuna)

<sup>a</sup><http://www.collocations.de/software.html>

	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
ausazko rankinga	0,000	0,308	0,070	0,234
$z$ neurria	(-0,038)	0,297	0,109	0,204
$t$ neurria	<b>0,197</b>	<b>0,455</b>	0,084	<b>0,383</b>
$\chi^2$	(-0,037)	0,302	<b>0,119</b>	0,206
LLR	0,156	0,427	0,100	0,335
Fisher	0,156	0,426	0,100	0,335
MI	(-0,121)	0,257	0,086	0,182
$MI^3$	0,103	0,389	0,107	0,291
$f$	0,189	0,436	0,074	0,379

## Konposizionaltasun semantikoaren neurketa

Antzekotasun distribuzionalaren hipotesia:

- ▶ *Difference of meaning correlates with difference of distribution.* (Harris, 1954)

UFen aplikatua

- ▶ *The underlying hypothesis is that semantically idiomatic MWEs will occur in markedly different lexical contexts to their component words.* (Baldwin eta Kim, 2010)

UFentzako funtsezko prozedura

- ▶ UF izateko hautagai den hitz-konbinazioaren testuinguruak haren osagai bakunen testuinguruekin konparatzea
- ▶ Konparazio hori osagai bakoitzarekiko egin dugu, bakoitzaren eragina bereiz aztertzeko (Wulff, 2008); horrez gain, baterako neurria ere kalkulatu da, bien batezbestekoaren bidez

## Konposizionaltasun semantikoaren neurketa

### Testuinguru-sorkuntza

- ▶ Ondoz ondoko agerkidetzak hartu dira bigramatzat; arrazoa da erauzketan ezarritako irizpide bera erabiltzea ( $w = \pm 1$ )
- ▶ Osagai bakunen testuinguruetan konbinazioaren testuinguruak ez sartzea
  - ▶ *mahaia jaso* bigramaren agerpena atzematen denean, testuinguruko hitzek *mahaia jasoren* testuinguru-dokumentua elikatu dute, eta ez dira *mahai* eta *jasoren* testuingurutzat kontsideratu
- ▶ Testuinguru-dokumentuak elikatzeke: eduki-hitzak bakarrik (izen, aditz eta adjektiboak)
- ▶ Esaldi osoa hartu dugu testuingurutzat (Reddy et al., 2011)

## Konposizionaltasun semantikoaren neurketa

### Testuinguruen prozesamendua - WSM (*word space model*)

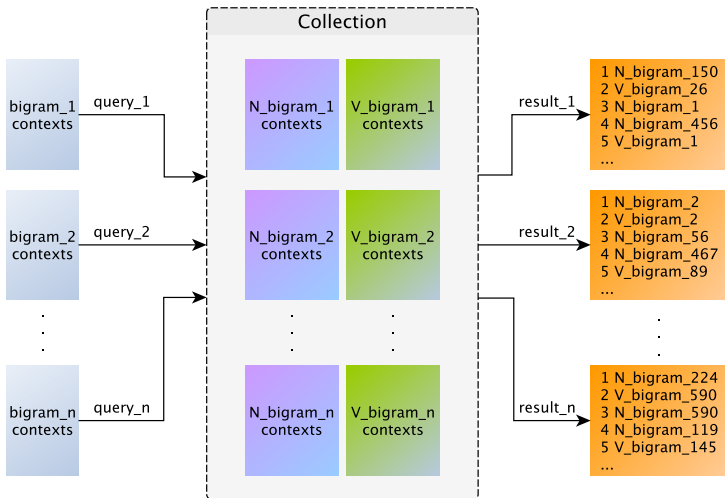
- ▶ Testuinguruak bektore gisa errepresentatuta
  - ▶ Bektoreetako balioak edo pisuak:  $f$ ,  $t$  neurria, LLR, MI eta Fisherren test zehatza
  - ▶ Dimentsioak gutxitzeko:
    - ▶ Ehuneko jakin bat: rankingaren % 75 eta % 50
    - ▶ Kopuru mugatu bat: rankingeko lehen 3 000, 2 000 eta 1 000 elementuak (Reddy et al., 2011; Mitchell eta Lapata, 2008)
    - ▶ Maiztasun-atari bat ( $f > 3$ )
  - ▶ Corpuseko maiztasun handieneko 100 hitzak ez dira kontuan hartu (*stop word* gisa prozesatu dira)
- ▶ Antzekotasun-neurriak
  - ▶ Jaccard koefizientea, kosinua eta Jensen-Shannon dibergentzia
  - ▶ Berry-Roggheren (1974)  $R$  balioa ( $R_{BR}$ ) eta Wulffen (2008) bi hedapenak ( $R_{W1}$  eta  $R_{W2}$ )

## Konposizionaltasun semantikoaren neurketa

### Testuinguruen prozesamendua - IR metodoak (Lemur)

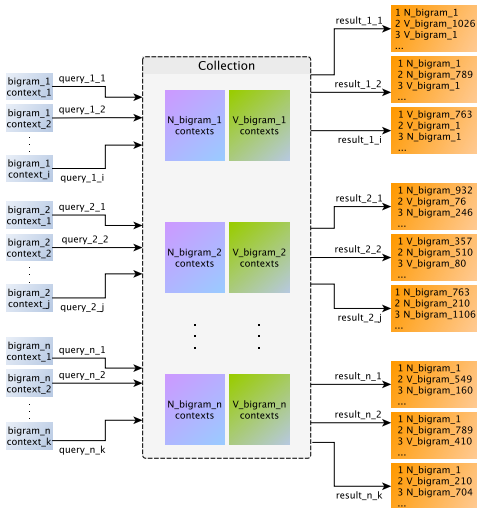
- ▶ Bigramen testuinguruen dokumentuak *query* edo kontsulta gisa erabiltzea bigramen osagaien testuinguruen dokumentuek osatzen duten bildumaren kontra
- ▶ Bi modalitate:
  - ▶ L1: bigrama baten testuinguruak dokumentu bakarrean bildu dira, eta orobat osagaien testuinguruak
  - ▶ L2: bigramen testuinguru-esaldiak dokumentu banatan sartu dira; beraz, bigrama bakoitzerako agerkidetza adina kontsulta-dokumentu eratu dira
- ▶ Indizeak: Indri, KL, tf-idf, Okapi, cos (idf)

## Lemur - L1 modalitatea





# Lemur - L2 modalitatea



## Malgutasun morfosintaktikoa

### Prozeduraren oinarriak

- ▶ Bigramaren portaera bi erreferentzia-portaera mota hauekin konparatzea
  - ▶ Portaera orokorra
  - ▶ Osagaien portaera
- ▶ Konbinazio baten portaera = aldakuntza morfosintaktikoen banaketa
- ▶ Aldakuntzen hautaketa
  - ▶ EDBLko gauzatze-eskemak (Urizar, 2012) aztertu ondoren *izena* osagaiaren aldakuntzetan jarri dugu arreta; aditzaren aldakuntzak ez dira idiomatikotasunerako esanguratsuak
- ▶ Murriztapen-gramatika bat garatu dugu hautatutako aldakuntzak detektatzeko

## Malgutasun morfosintaktikoa - Aztertutako aldakuntzak I

### Izenaren hedapenak

- ▶ Determinatzailea: *liburu **bat** irakurri dut; **zenbat** liburu irakurri dituzu?*
- ▶ Izenondoa: *liburu **interesgarria** irakurri nuen*
- ▶ Izenlaguna: ***gustuko** liburuak irakurtzea; **italierazko** liburua irakurri*

Hedapen bat baino gehiago konbina daitezke aldakuntza berean: *liburu **interesgarri bat** irakurri dut; **lau** liburu **hauek** irakurri ditut; anaiak irakurritako **frantsesezko** liburu **eder batzuk**.*

### Erlatiboa

- ▶ ***irakurri dudan** liburua; anaiak **irakurritako** frantsesezko **liburu** eder batzuk*

## Malgutasun morfosintaktikoa - Aztertutako aldakuntzak II

Mugatasuna: NUMS, NUMP, MG, PH

- ▶ *liburu interesgarria* (NUMS) / *interesgarriak irakurri nuen/nituen* (NUMP)
- ▶ *liburu bat irakurri dut* (MG)
- ▶ *ez dut liburu hori irakurri* (NUMS)
- ▶ *hiru liburu hauek irakurriko ditut* (NUMP)

Ordena: IZE ADI / ADI IZE

- ▶ *liburu interesgarri bat irakurri dut / irakurri dut gomendatu didazun liburua*

## Malgutasun morfosintaktikoa - Malgutasun-neurriak

### Banaketen arteko distantziaren neurriak

- ▶ **RFR** (*relative frequency ratio*): aldakuntza bakoitzerako behatutako maiztasun erlatiboaren eta itxarondakoaren arteko zatiduren batura, Barkematik (1994) egokitu duguna
- ▶ **NSSD** (*normalized sum of squared deviations*): aurreko maiztasunen kenduren karratuen batura normalizatua (Wulff, 2008)
- ▶  $H_{rel}$ : behatutako banaketaren entropia erlatiboa, sistemaren entropia maximoarekiko (Wulff, 2008)
- ▶ **KL**: Kullback-Leibler dibergentzia (Fazly and Stevenson, 2007)
- ▶ **CPMI** (*conditional pointwise mutual information*): elkarrekiko informazio puntual baldintzazkoa (Bannard, 2007)

## Malgutasun lexikala - Ordezagarritasuna

### Prozeduraren oinarria

- ▶ Konbinazioaren osagaietako bakoitzaren ordezagarritasuna konputatzea, ordezkatzat osagaiaren sinonimoak, kuasisinonimoak edo semantikoki erlazionatutako hitzak erabiliz

### Ordezkoak eskuratzeko baliabideak

- ▶ ELH\_SK: *Sinonimoen Kutxa* (Elhuyar Fundazioa)
- ▶ EusWN 3.0: WordNeten euskarazko bertsioa<sup>a</sup> (Ixa taldea)
- ▶ Baliabide konbinatuak eta hedatuak:
  - ▶ ELHWN: aurreko bien bildura
  - ▶ WNhedap: bigramaren osagai bakoitzak EusWN 3.0-n dituen ahaideak (*siblings*, edo hiperonimo bereko hitzak) gehituz osatutako baliabide hedatua
  - ▶ ELHWNhedap: ELHWN eta WNhedap baliabideen bildura
- ▶ EB\_antzdistr: bigramen osagaien thesaurus distribuzional bat eratu dugu corpusetik, semantikoki antzekoenak diren hitzez osatua

<sup>a</sup><http://ixa2.si.ehu.es/cgi-bin/mcr/public/wei.consult.perl>

## Malgutasun lexikala - Ordezgarritasuna

Ordezgarritasuna neurtzeko baliabideen ezaugarriak: ordezkoen kopuruak

<b>kategoria</b>	<b>ELH_SK</b>	<b>EusWN 3.0</b>	<b>ELHWN</b>	<b>WNhedap</b>	<b>ELHWNhedap</b>
IZE	29 407	29 603	52 434	548 454	570 404
ADI	9 333	7 898	15 061	87 567	94 371
<b>Totala</b>	<b>38 740</b>	<b>37 669</b>	<b>67 495</b>	<b>636 021</b>	<b>664 775</b>

Adiera-bereizketa kontuan hartu gabe eraturako sinonimo-bikoteen bost bildumak

Ordezko konbinaziorik ez duten edo corpusean ordezkoaren agerpenik ez duten bigramen kopuruak

<b>Baliabidea</b>	<b>ordezkorik ez</b>
ELH_SK	674
EusWN 3.0	552
ELHWN	475
WNhedap	299
ELHWNhedap	238
EBantzdistr_20	207

## Malgutasun lexikala - Ordezgarritasuna

### Ordezgarritasunaren neurriak

- ▶ Van de Cruys eta Moirónen (2007)  $R$  indizea, bigramaren eta ordezkoko aldaeren arteko KL dibergentzia oinarritua
  - ▶  $R_{nv}$  eta  $R_{vn}$ , hurrenez hurren, aditz batek izen jakin batekiko duen joera eta izen batek aditz batekiko duen joera dira, bere ordezkoez izen edo aditz horiekiko duten joerarekin konparatuta
- ▶ Fazly eta Stevenson (2007)  $\text{Fixedness}_{\text{lex}}$  neurria, bigramaren eta haren osagai bakoitza ordezkatzeko sortutako aldaeren MI balioen arteko  $z$  neurria dena
  - ▶ Egileen baterako neurketaz gain ( $z_{\text{PMI}_{\text{NV}}}$ ),  $z_{\text{PMI}_{\text{V}}}$  eta  $z_{\text{PMI}_{\text{N}}}$  ordezkagarritasunak bereiz kalkulatutako ditugu



## Esperimentu bakunen emaitza hautatuak

	neurria	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
	ausaz	0,000	0,309	0,070	0,234
AM	$t$ neurria	<u>0,197</u>	<u>0,455</u>	0,084	<u>0,383</u>
	$\chi^2$	(-0,037)	0,302	<u>0,119</u>	<u>0,206</u>
DSim	L2_Indri_hit_erl_NV	<u>0,322</u>	<u>0,566</u>	0,294	0,354
	L2_KL_hit_erl_NV	0,309	0,551	<u>0,320</u>	0,332
	L2_Indri_rankV_hazt	0,320	0,552	0,130	<u>0,431</u>
MSFlex	CPMI_osag_mugat	<u>0,154</u>	0,432	0,113	0,329
	$H_{rel\_mugat}$	0,119	<u>0,434</u>	0,134	0,313
	CPMI_orok_ADJ	(0,045)	0,387	<u>0,202</u>	0,237
	CPMI_izena_mugat	0,132	0,425	<u>0,101</u>	<u>0,331</u>
LFlex	$z\_PMI\_V\_WNhedap$	<u>0,110</u>	0,365	0,120	0,259
	$R_{nv\_EBantzdistr\_20}$	0,066	<u>0,381</u>	0,066	<u>0,323</u>
	$z\_PMI\_V\_ELHWNhedap$	0,108	0,357	<u>0,122</u>	0,250

Idiomatikotasunaren lau osagaiak neurtzeko egindako esperimentu bakunen emaitza onenen laburpena. Kendall  $\tau_B$  koefizientea (idiomatikotasun-ranking ideal batekiko hein-korrelazioa); eta batez besteko doitasunak ( $AP$ ), UFen rankingerako ( $AP_{UF}$ ), oro har, eta  $AP$  espezifikoak esapide idiomatikoetarako eta kolokazioetarako.

## Esperimentu bakunen emaitza hautatuak

	neurria	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
	ausaz	0,000	0,309	0,070	0,234
AM	$t$ neurria	<u>0,197</u>	<u>0,455</u>	0,084	<u>0,383</u>
	$\chi^2$	(-0,037)	0,302	<u>0,119</u>	<u>0,206</u>
DSim	L2_Indri_hit_erl_NV	<b>0,322</b>	<u>0,566</u>	0,294	0,354
	L2_KL_hit_erl_NV	0,309	<u>0,551</u>	<u>0,320</u>	0,332
	L2_Indri_rankV_hazt	0,320	0,552	<u>0,130</u>	<u>0,431</u>
MSFlex	CPMI_osag_mugat	<u>0,154</u>	0,432	0,113	0,329
	$H_{rel}$ _mugat	0,119	<u>0,434</u>	0,134	0,313
	CPMI_orok_ADJ	(0,045)	0,387	<u>0,202</u>	0,237
	CPMI_izena_mugat	0,132	0,425	<u>0,101</u>	<u>0,331</u>
LFlex	$z_{PMI\_V\_WN}$ hedap	<u>0,110</u>	0,365	0,120	0,259
	$R_{nv\_EB}$ antzdistr_20	0,066	<u>0,381</u>	0,066	<u>0,323</u>
	$z_{PMI\_V\_ELHWN}$ hedap	0,108	0,357	<u>0,122</u>	0,250

Idiomatikotasunaren lau osagaiak neurtzeko egindako esperimentu bakunen emaitza onenen laburpena. Kendall  $\tau_B$  koefizientea (idiomatikotasun-ranking ideal batekiko hein-korrelazioa); eta batez besteko doitasunak ( $AP$ ), UFen rankingerako ( $AP_{UF}$ ), oro har, eta  $AP$  espezifikoak esapide idiomatikoetarako eta kolokazioetarako.

## Esperimentu bakunen emaitza hautatuak

	neurria	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
	ausaz	0,000	0,309	0,070	0,234
AM	$t$ neurria	<u>0,197</u>	<u>0,455</u>	0,084	<u>0,383</u>
	$\chi^2$	(-0,037)	0,302	<u>0,119</u>	<u>0,206</u>
DSim	L2_Indri_hit_erl_NV	<b><u>0,322</u></b>	<b><u>0,566</u></b>	0,294	0,354
	L2_KL_hit_erl_NV	0,309	0,551	<u>0,320</u>	0,332
	L2_Indri_rankV_hazt	0,320	0,552	0,130	<u>0,431</u>
MSFlex	CPMI_osag_mugat	<u>0,154</u>	0,432	0,113	0,329
	$H_{rel\_mugat}$	0,119	<u>0,434</u>	0,134	0,313
	CPMI_orok_ADJ	(0,045)	0,387	<u>0,202</u>	0,237
	CPMI_izena_mugat	0,132	0,425	<u>0,101</u>	<u>0,331</u>
LFlex	$z\_PMI\_V\_WNhedap$	<u>0,110</u>	0,365	0,120	0,259
	$R_{nv\_EBantzdistr\_20}$	0,066	<u>0,381</u>	0,066	<u>0,323</u>
	$z\_PMI\_V\_ELHWNhedap$	0,108	0,357	<u>0,122</u>	0,250

Idiomatikotasunaren lau osagaiak neurtzeko egindako esperimentu bakunen emaitza onenen laburpena. Kendall  $\tau_B$  koefizientea (idiomatikotasun-ranking ideal batekiko hein-korrelazioa); eta batez besteko doitasunak ( $AP$ ), UFen rankingerako ( $AP_{UF}$ ), oro har, eta  $AP$  espezifikoak esapide idiomatikoetarako eta kolokazioetarako.

## Esperimentu bakunen emaitza hautatuak

	neurria	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
	ausaz	0,000	0,309	0,070	0,234
AM	$t$ neurria	<u>0,197</u>	<u>0,455</u>	0,084	<u>0,383</u>
	$\chi^2$	(-0,037)	0,302	<u>0,119</u>	<u>0,206</u>
DSim	L2_Indri_hit_erl_NV	<b><u>0,322</u></b>	<b><u>0,566</u></b>	0,294	0,354
	L2_KL_hit_erl_NV	0,309	0,551	<b><u>0,320</u></b>	0,332
	L2_Indri_rankV_hazt	0,320	0,552	<u>0,130</u>	<u>0,431</u>
MSFlex	CPMI_osag_mugat	<u>0,154</u>	0,432	0,113	0,329
	$H_{rel\_mugat}$	0,119	<u>0,434</u>	0,134	0,313
	CPMI_orok_ADJ	(0,045)	0,387	<u>0,202</u>	0,237
	CPMI_izena_mugat	0,132	0,425	<u>0,101</u>	<u>0,331</u>
LFlex	$z\_PMI\_V\_WNhedap$	<u>0,110</u>	0,365	0,120	0,259
	$R_{nv\_EBantzdistr\_20}$	0,066	<u>0,381</u>	0,066	<u>0,323</u>
	$z\_PMI\_V\_ELHWNhedap$	0,108	0,357	<u>0,122</u>	0,250

Idiomatikotasunaren lau osagaiak neurtzeko egindako esperimentu bakunen emaitza onenen laburpena. Kendall  $\tau_B$  koefizientea (idiomatikotasun-ranking ideal batekiko hein-korrelazioa); eta batez besteko doitasunak ( $AP$ ), UFen rankingerako ( $AP_{UF}$ ), oro har, eta  $AP$  espezifikoak esapide idiomatikoetarako eta kolokazioetarako.

## Esperimentu bakunen emaitza hautatuak

	neurria	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
	ausaz	0,000	0,309	0,070	0,234
AM	$t$ neurria	<u>0,197</u>	<u>0,455</u>	0,084	<u>0,383</u>
	$\chi^2$	(-0,037)	0,302	<u>0,119</u>	<u>0,206</u>
DSim	L2_Indri_hit_erl_NV	<b><u>0,322</u></b>	<b><u>0,566</u></b>	0,294	0,354
	L2_KL_hit_erl_NV	0,309	0,551	<b><u>0,320</u></b>	0,332
	L2_Indri_rankV_hazt	0,320	0,552	<u>0,130</u>	<b><u>0,431</u></b>
MSFlex	CPMI_osag_mugat	<u>0,154</u>	0,432	0,113	0,329
	$H_{rel\_mugat}$	0,119	<u>0,434</u>	0,134	0,313
	CPMI_orok_ADJ	(0,045)	0,387	<u>0,202</u>	0,237
	CPMI_izena_mugat	0,132	0,425	<u>0,101</u>	<u>0,331</u>
LFlex	$z\_PMI\_V\_WNhedap$	<u>0,110</u>	0,365	0,120	0,259
	$R_{nv\_EBantzdistr\_20}$	0,066	<u>0,381</u>	0,066	<u>0,323</u>
	$z\_PMI\_V\_ELHWNhedap$	0,108	0,357	<u>0,122</u>	0,250

Idiomatikotasunaren lau osagaiak neurtzeko egindako esperimentu bakunen emaitza onenen laburpena. Kendall  $\tau_B$  koefizientea (idiomatikotasun-ranking ideal batekiko hein-korrelazioa); eta batez besteko doitasunak ( $AP$ ), UFen rankingerako ( $AP_{UF}$ ), oro har, eta  $AP$  espezifikoak esapide idiomatikoetarako eta kolokazioetarako.

## Esperimentu bakunen emaitza hautatuak

	neurria	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
	ausaz	0,000	0,309	0,070	0,234
AM	$t$ neurria	<u>0,197</u>	<u>0,455</u>	0,084	<u>0,383</u>
	$\chi^2$	(-0,037)	0,302	<u>0,119</u>	<u>0,206</u>
DSim	L2_Indri_hit_erl_NV	<b><u>0,322</u></b>	<b><u>0,566</u></b>	0,294	0,354
	L2_KL_hit_erl_NV	0,309	0,551	<b><u>0,320</u></b>	0,332
	L2_Indri_rankV_hazt	0,320	0,552	<u>0,130</u>	<b><u>0,431</u></b>
MSFlex	CPMI_osag_mugat	<u>0,154</u>	0,432	0,113	0,329
	$H_{rel\_mugat}$	0,119	<u>0,434</u>	0,134	0,313
	CPMI_orok_ADJ	(0,045)	0,387	<b><u>0,202</u></b>	0,237
	CPMI_izena_mugat	0,132	0,425	<u>0,101</u>	<u>0,331</u>
LFlex	$z\_PMI\_V\_WNhedap$	<u>0,110</u>	0,365	0,120	0,259
	$R_{nv\_EBantzdistr\_20}$	0,066	<u>0,381</u>	0,066	<u>0,323</u>
	$z\_PMI\_V\_ELHWNhedap$	0,108	0,357	<u>0,122</u>	0,250

Idiomatikotasunaren lau osagaiak neurtzeko egindako esperimentu bakunen emaitza onenen laburpena. Kendall  $\tau_B$  koefizientea (idiomatikotasun-ranking ideal batekiko hein-korrelazioa); eta batez besteko doitasunak ( $AP$ ), UFen rankingerako ( $AP_{UF}$ ), oro har, eta  $AP$  espezifikoak esapide idiomatikoetarako eta kolokazioetarako.

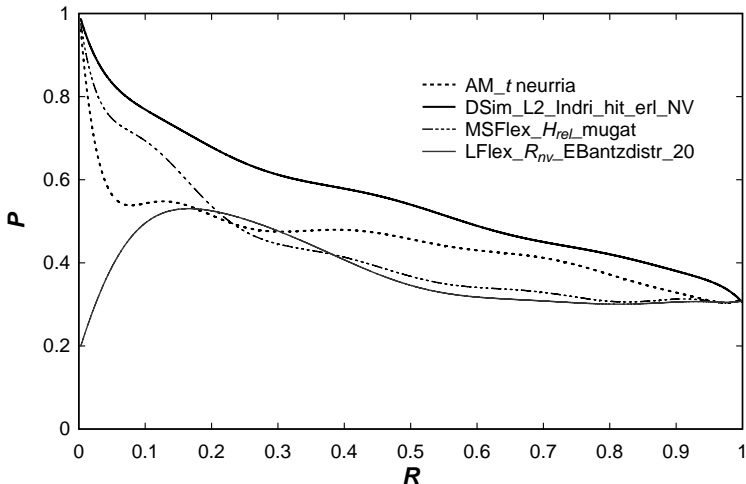
## Esperimentu bakunen emaitza hautatuak

	neurria	$\tau_B$	$AP_{UF}$	$AP_{id}$	$AP_{col}$
	ausaz	0,000	0,309	0,070	0,234
AM	$t$ neurria	<u>0,197</u>	<u>0,455</u>	0,084	<b><u>0,383</u></b>
	$\chi^2$	(-0,037)	0,302	<u>0,119</u>	0,206
DSim	L2_Indri_hit_erl_NV	<b><u>0,322</u></b>	<b><u>0,566</u></b>	0,294	0,354
	L2_KL_hit_erl_NV	0,309	0,551	<b><u>0,320</u></b>	0,332
	L2_Indri_rankV_hazt	0,320	0,552	<u>0,130</u>	<b><u>0,431</u></b>
MSFlex	CPMI_osag_mugat	<u>0,154</u>	0,432	0,113	0,329
	$H_{rel}$ _mugat	0,119	<u>0,434</u>	0,134	0,313
	CPMI_orok_ADJ	(0,045)	0,387	<b><u>0,202</u></b>	0,237
	CPMI_izena_mugat	0,132	0,425	<u>0,101</u>	<u>0,331</u>
LFlex	$z$ _PMI_V_WNhedap	<u>0,110</u>	0,365	0,120	0,259
	$R_{nv}$ _EBantzdistr_20	0,066	<u>0,381</u>	0,066	<u>0,323</u>
	$z$ _PMI_V_ELHWNhedap	0,108	0,357	<u>0,122</u>	0,250

Idiomatikotasunaren lau osagaiak neurtzeko egindako esperimentu bakunen emaitza onenen laburpena. Kendall  $\tau_B$  koefizientea (idiomatikotasun-ranking ideal batekiko hein-korrelazioa); eta batez besteko doitasunak ( $AP$ ), UFen rankingerako ( $AP_{UF}$ ), oro har, eta  $AP$  espezifikoak esapide idiomatikoetarako eta kolokazioetarako.

Emaitzak -  $P/R$  kurbak - UFak

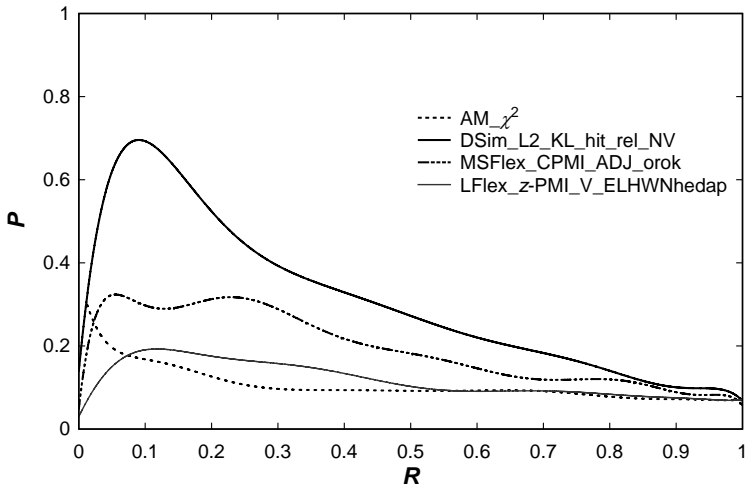
## Unitate fraseologikoak (esapide idiomatikoak+kolokazioak)





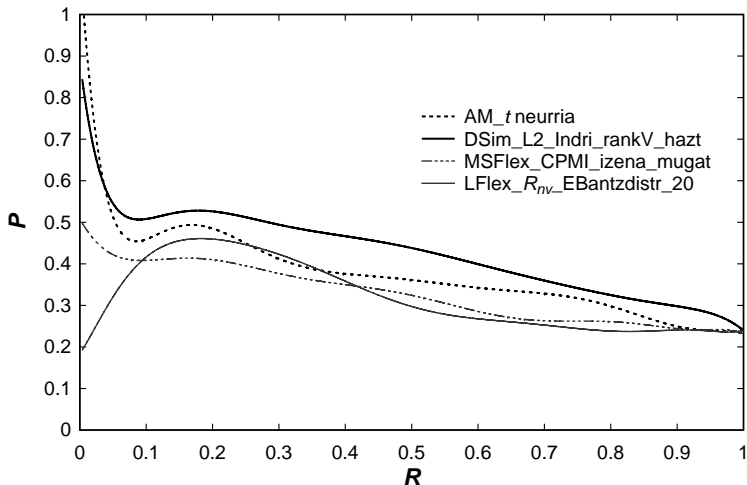
## Emaitzak - $P/R$ kurbak - esapide idiomatikoak

### Esapide idiomatikoak



Emaitzak -  $P/R$  kurbak - kolokazioak

## Kolokazioak



## Esperimentu bakunen emaitzen analisia

### Emaitzen alderdi nagusiak

- ▶ Antzekotasun distribuzionaleko neurriak (DSim) dira idiomatikotasun-mailarekin  $\tau_B$  korrelazio onena dutenak. Emaitza onenak IR tresna batekin lortu dira, L2 modalitateko esperimentuetan
- ▶ Esapide idiomatikoaren erauzketa ( $AP_{id}$ ) nabari da batez ere DSim neurrien gailentasuna
- ▶ Kolokazio-erauzketa, emaitza onena aditzaren semantika neurtzen duen Indri indize batek eman du (L2\_Indri\_rankV\_hazt), kolokazioen ezaugarri aipatuenetakoa den idiosinkrasia estatistikoa neurtzen duten AMen gainetik
- ▶ AMak: emaitza hobeak kolokazioekin esapide idiomatikoekin baino
- ▶ MSFlex neurketak, oro har, bigarren onenak dira  $AP_{id}$ -ren emaitzetan, baina  $AP_{col}$ -en emaitzak balio apalenen artean daude
- ▶ LFlex: emaitza eskasak, espero baino txarragoak

## Ikasketa automatikoa

### Esperimentuen diseinua

- ▶ **Sei metodo** (Weka toolkit): Naive Bayes, j48 (C4.5 erabaki-zuhaitza), Random Forest (RF), PART, SMO (SVM edo sostengu-bektoreen makinaren inplementazio bat) eta Logistic Regression (LR)
- ▶ Ereduaren balidazioa: ebaluazio-erreferentzia ez da oso handia (1 145 instantzia) → balidazio gurutzatua (*cross-validation*)

## Ikasketa automatikoa

### Esperimentuen diseinua

- ▶ **Sei metodo** (Weka toolkit): Naive Bayes, j48 (C4.5 erabaki-zuhaitza), **Random Forest** (RF), PART, **SMO** (SVM edo sostengu-bektoreen makinaren implementazio bat) eta **Logistic Regression** (LR)
- ▶ Ereduaren balidazioa: ebaluazio-erreferentzia ez da oso handia (1 145 instantzia) → balidazio gurutzatua (*cross-validation*)

## Ikasketa automatikoa

### Esperimentuen diseinua

- ▶ Sei metodo (Weka toolkit): Naive Bayes, j48 (C4.5 erabaki-zuhaitza), Random Forest (RF), PART, SMO (SVM edo sostengu-bektoreen makinaren inplementazio bat) eta Logistic Regression (LR)
- ▶ Ereduaren balidazioa: ebaluazio-erreferentzia ez da oso handia (1 145 instantzia) → **balidazio gurutzatua** (*cross-validation*)

## Datu-multzoak

- ▶ AM, DSim, MSFlex eta LFlex: propietate bakoitzaren neurketak  
Atributu-kopuruak: AM 7; DSim 40; MSFlex 34; LFlex 20
- ▶ 4 osag.: propietate guztien neurketak, batera
- ▶ 4 osag.+ad.: propietate guztien neurketak, eta konbinazioaren aditza, Fazlyri (2007) jarraituz, eta korrelazioa egiaztatu ondoren
- ▶ Atributu-hautaketa automatikoa
  - ▶ CS-BF: *Wekaren AttributeSelectedClassifier* metasailkatzailean, osagai hauek dituen iragazkia: `CfsSubsetEval1` ebaluatzailetzat, eta `BestFirst` bilaketa-metodotzat. Guztira, 37 atributu hautatu ditu iragazki horrek: AM 3; DSim 16; MSFlex 7; LFlex 2; 9 aditz.

## Ikasketa automatikoko esperimentuen emaitzak

Atrib.	Metod.	CCI	$F_{id}$	$F_{col}$	$F_{free}$	$F_{mikro}$	$F_{makro}$
Oinarri-lerroa		69,607	0,000	0,000	0,821	0,571	0,274
DSim	LR	<u>74,061</u>	0,270	<u>0,468</u>	<u>0,842</u>	<u>0,714</u>	<u>0,527</u>
	SMO	69,607	0,046	0,063	0,820	0,589	0,310
	RF	71,441	<u>0,279</u>	0,438	0,822	0,694	0,513
4 osag.	LR	73,362	<u>0,355</u>	<u>0,487</u>	0,837	0,722	<u>0,560</u>
	SMO	<u>76,070</u>	0,300	<u>0,479</u>	<u>0,858</u>	<u>0,731</u>	<u>0,546</u>
	RF	73,712	0,336	0,475	0,840	0,719	0,550
4 osag.+ad.	LR	62,795	0,274	0,490	0,751	0,656	0,505
	SMO	<u>76,856</u>	<u>0,418</u>	<u>0,544</u>	<u>0,858</u>	<u>0,754</u>	<u>0,607</u>
	RF	73,974	0,304	0,447	0,843	0,713	0,531
CS-BF	LR	<u>75,721</u>	<u>0,339</u>	<u>0,487</u>	<u>0,854</u>	<u>0,732</u>	<u>0,560</u>
	SMO	73,450	0,149	0,390	0,838	0,685	0,459
	RF	72,838	0,364	0,435	0,836	0,709	0,545



## Ikasketa automatikoko esperimentuen emaitzak

Atrib.	Metod.	CCI	$F_{id}$	$F_{col}$	$F_{free}$	$F_{mikro}$	$F_{makro}$
Oinarri-lerroa		69,607	0,000	0,000	0,821	0,571	0,274
DSim	LR	<u>74,061</u>	0,270	<u>0,468</u>	<u>0,842</u>	<u>0,714</u>	<u>0,527</u>
	SMO	69,607	0,046	0,063	0,820	0,589	0,310
	RF	71,441	<u>0,279</u>	0,438	0,822	0,694	0,513
4 osag.	LR	73,362	<u>0,355</u>	<u>0,487</u>	0,837	0,722	<u>0,560</u>
	SMO	<u>76,070</u>	0,300	0,479	<u>0,858</u>	<u>0,731</u>	0,546
	RF	73,712	0,336	0,475	0,840	0,719	0,550
4 osag.+ad.	LR	62,795	0,274	0,490	0,751	0,656	0,505
	SMO	<b><u>76,856</u></b>	<b><u>0,418</u></b>	<b><u>0,544</u></b>	<b><u>0,858</u></b>	<b><u>0,754</u></b>	<b><u>0,607</u></b>
	RF	73,974	0,304	0,447	0,843	0,713	0,531
CS-BF	LR	<u>75,721</u>	<u>0,339</u>	<u>0,487</u>	<u>0,854</u>	<u>0,732</u>	<u>0,560</u>
	SMO	73,450	0,149	0,390	0,838	0,685	0,459
	RF	72,838	0,364	0,435	0,836	0,709	0,545

## Ikasketa automatikoko esperimentuen emaitzak

Atrib.	Metod.	CCI	$F_{id}$	$F_{col}$	$F_{free}$	$F_{mikro}$	$F_{makro}$
Oinarri-lerroa		69,607	0,000	0,000	0,821	0,571	0,274
DSim	LR	<u>74,061</u>	0,270	<u>0,468</u>	<u>0,842</u>	<u>0,714</u>	<u>0,527</u>
	SMO	69,607	0,046	0,063	0,820	0,589	0,310
	RF	71,441	<u>0,279</u>	0,438	0,822	0,694	0,513
4 osag.	LR	73,362	<u>0,355</u>	<u>0,487</u>	0,837	0,722	<u>0,560</u>
	SMO	<u>76,070</u>	0,300	0,479	<u>0,858</u>	<u>0,731</u>	0,546
	RF	73,712	0,336	0,475	0,840	0,719	0,550
4 osag.+ad.	LR	62,795	0,274	0,490	0,751	0,656	0,505
	SMO	<b><u>76,856</u></b>	<b><u>0,418</u></b>	<b><u>0,544</u></b>	<b><u>0,858</u></b>	<b><u>0,754</u></b>	<b><u>0,607</u></b>
	RF	73,974	0,304	0,447	0,843	0,713	0,531
CS-BF	LR	<b><u>75,721</u></b>	<b><u>0,339</u></b>	<b><u>0,487</u></b>	<b><u>0,854</u></b>	<b><u>0,732</u></b>	<b><u>0,560</u></b>
	SMO	73,450	0,149	0,390	0,838	0,685	0,459
	RF	72,838	0,364	0,435	0,836	0,709	0,545

## Sailkapen automatikoaren emaitzen analisia

### Emaitzen alderdi nagusiak

- ▶ Ezagutza-iturri bakarreko lau datu-multzoetatik,  $D_{Sim}$  da emaitza onenak dituen, sailkatzaile guztietan (onena, LR)
- ▶ Oro har, algoritmo guztiek emaitza hobekak dituzte idiomatikotasunaren ezagutza-iturriak konbinatuz (4 osag. eta 4 osag.+ad. datu-multzoak), iturri bakarra erabiliz baino (onena: SVM familiako SMO algoritmoa)
- ▶  $CS-BF$  metasailkatzailearen LR esperimentuetako  $F_{mikro}$  eta  $F_{makro}$  (0,732 / 0,560) dira gehien hurbiltzen direnak SMO metodoak 4 osag.+ad. datu-multzoarekin lortutako emaitza onenetara (0,754 / 0,607)

## Predikatu konplexu batzuk birsailkatzearen eragina – Planteamendua

- ▶ Ebaluazio-erreferentzian badira konbinazio batzuk (esaterako, *jaramon egin, zor izan*), esapide erdiidiomatikotzat gabe, esapide idiomatikotzat jotzekoak liritekeenak. 17 dira. Horiek birsailkatuta: `id 97`; `col 251`; `free 797`
- ▶ Aurreikusitako ezaugarri nagusiak:
  - ▶ Maiztasun handikoak (ohiko `id` kategoriakoak baino handiagokoak)
  - ▶ Semantikoki erdikonposizionalak (izenak bere oinarrizko esanahia atxikitzen du)
  - ▶ Idiosinkrasia morfosintaktikoa (adib., forma kanonikoan, izena mugatzailerik gabe: *jaramon egin, zor izan*)
- ▶ Esperimentu bakun eta konbinatu guztiak egin ditugu sailkapen berriarekin

## Predikatu konplexu batzuk birsailkatzearen eragina - Emaitzak

### Esperimentu bakunak

- ▶  $\tau_B$ : igoera apalak (unitate gutxi birsailkatu dira)
- ▶ AM neurrien  $AP_{id}$ -ren igoera (batez ere,  $f$ -rekiko korrelazioa duten AMenak) → birsailkatzean,  $id$  kategoriakoen idiosinkrasia estatistikoa handitu egin da
- ▶ DSim:  $AP_{id}$ -ren balio hobexeagoak, aditzaren semantikarekiko neurketetan → birsailkatuek profil erdikonposizionala dute
- ▶ MSFlex neurketetan nabari da gehien  $AP_{id}$ -ren balioen igoera → birsailkatuak morfosintaktikoki idiosinkratikoak dira

### Ikasketa automatikoa

- ▶ Emaitzak, oro har, zertxobait apalagoak. Baina  $F_{id}$  hobeak, instantzia gehiago baitaude ikasteko
- ▶ CS-BF iragazkiak, proportzionalki, MSFlex atributu gehiago aukeratu ditu aurreko sailkapenarekin baino

→ Horrelako predikatu konplexuen profil berezia

## Aurkezpenaren eskema

- 1 UFen idiomatikotasunaren eta karakterizazioaren marko teorikoa
- 2 UFen erauzketa eta karakterizazio automatikorako teknikak
- 3 Lan esperimentalak eta emaitzen analisia
- 4 Ondorioak, ekarpenak eta etorkizuneko lanak

## Ondorioak

- ▶ Idiomatikotasunaren izaera konplexua
- ▶ Idiosinkrasia semantikoaren nabarmentasuna, eta konposizionaltasunaren gradualtasuna
- ▶ Kolokazioen erdikonposizionaltasuna
- ▶ Kolokazioen malgutasun morfosintaktiko handia
- ▶ UF-kategoriaren eta konbinazioaren aditzaren arteko korrelazioa erabilgarria da ikasketa automatikoan
- ▶ Predikatu konplexu batzuen (*jaramon egin, zor izan* modukoak) profil berezia
- ▶ Teoria fraseologikoaren aurreanak eta lan esperimental batzuen emaitzak kontuan izanik, LFlex esperimentuen emaitzak ez dira espero genuen mailakoak

## LFlex esperimentuen emaitzak esplikatzen diren hipotesi batzuk

- ▶ Ordezkarritasuna egiaztatzen diren euskarazko baliabideak estaldura gutxiago izatea
- ▶ Euskara, batez ere kazetaritza-erregistroan, gutxiago finkatua izan liteke, eta, beraz, kolokatiboaren aukera lexikala ez litzateke hain argi finkatua
- ▶ LFlex emaitza konparatiboki onak lortu diren ikerlanetan kosinua erabili da DSim neurritzat, eta neurri horrek aski emaitza apalak izan ditu gure ikerketan ere



## Ekarpenak

- ▶ Corpus handi bat prozesatu da (74 milioi hitz), eta bigrama-bilduma baliagarria erauzi
- ▶ Ebaluaziorako erreferentziak: hiztegi-erreferentzia eta eskuz landutako erreferentzia sailkatua
- ▶ `izena+aditza` osaerako konbinazioak automatikoki erauzteko prozedura
- ▶ Idiomatikotasunaren propietateak neurtzeko teknikak euskararako egokitu eta garatu ditugu
- ▶ Idiosinkrasia semantikoa neurtzeko, IR arloko teknikak erabili ditugu, eta horien bidez lortu ere emaitza onenak
- ▶ Sailkatzaile automatikoak trebatu eta ebaluatu dira; ekarpen handiena egiten duten propietateen informazioa eskuratu dugu
- ▶ Garatutako metodoa erraz egokitu daiteke beste egitura morfosintaktiko batzuetako bigramak erauzteko (`izena+izena`, `izena+izenondoa`, `izenlaguna+izena`...); Aplikazioa: Elhuyar Fundazioaren *Web-corpusen Atariko* "Hitz-konbinazioak" atala<sup>1</sup>

---

<sup>1</sup><http://webcorpusak.elhuyar.org/cgi-bin/kolokatuak.py>

## Etorkizuneko lanak

- ▶ Informazio sintaktiko aberatsagoa erabiltzea: Maltixa mendekotasun-analizatzailea (Bengoetxea eta Gojenola, 2010)
- ▶ Interpretazio literalak eta idiomatikoa izan ditzaketen konbinazioen agerpenak bereiztea
- ▶ Beste corpus-mota batzuk erabiltzea (literarioak, espezializate-arlokoak, web-corpusak), arloak eta erregistroak LFlex neurketetan duten eragina aztertzeko
- ▶ Corpus paraleloak erabiltzea, bi hizkuntzako UF bikote baliokideak erauzteko eta karakterizatzeko
- ▶ UFak erauzi eta karakterizatzeko garatuko diren tresnetan aplikatzea. Tresna horiek hiztegi gintzan, itzulpen gintzan eta hizkuntzaren prozesamendu automatikoko hainbat atazatan izango dute aplikazioa

## Idiomatikotasunaren karakterizazio automatikoa: izen+aditza konbinazioak

**Doktoregaia:** Antton Gurrutxaga Hernaiz

**Zuzendariak:** Iñaki Alegria Loinaz, Xabier Artola Zubillaga

Lengoaia eta Sistema Informatikoak Saila  
Euskal Herriko Unibertsitatea / Universidad del País Vasco

2014ko uztailearen 21a

- ▶ Gurrutxaga, A. eta Alegria, I. (2011). Automatic extraction of NV expressions in Basque: basic issues on cooccurrence techniques. *Proceedings of the Workshop on Multiword Expressions: from Parsing and Generation to the Real World, 2-7* or. Portland, Oregon: Association for Computational Linguistics.
- ▶ Gurrutxaga, A. eta Alegria, I. (2012). Measuring the compositionality of NV expressions in Basque by means of distributional similarity techniques. *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC 2012)*, 2389-2394 or. Istanbul: ELRA.
- ▶ Gurrutxaga, A. eta Alegria, I. (2013). Combining different features of idiomaticity for the automatic classification of noun+ verb expressions in Basque. *Proceedings of the 9th Workshop on Multiword Expressions (MWE 2013) NAACL-HLT 2013*, 116-125 or. Atlanta, Georgia: Association for Computational Linguistics.

- ▶ Altzibar, X. (2005). Kolokazioak euskaraz. Zer axola duten kazetaritzan. *Euskarazko kazetaritzaren I. kongresua. Kazetaritza euskaraz: oraina eta geroa*, 383–395. UPV/EHU.
- ▶ Amosova, N.Ñ. (1963). *Osnovui anglijskoy frazeologii*. Leningrad: University Press.
- ▶ Baldwin, T. eta Kim, S. (2010). Multiword expressions. Indurkha, N. eta Damerau, F. J.(ed.), *Handbook of Natural Language Processing, second edition*, 267–292. CRC Press, Taylor and Francis Group, Boca Raton, AEB.
- ▶ Bannard, C., Baldwin, T. eta Lascarides, A. (2003). A statistical approach to the semantics of verb-particles. *Proceedings of the ACL 2003 Workshop on Multiword Expressions: analysis, acquisition and treatment* 18. lib., 65–72. Association for Computational Linguistics.
- ▶ Bannard, C. (2007). A measure of syntactic flexibility for automatically identifying multiword expressions in corpora. *Proceedings of the Workshop on a Broader Perspective on Multiword Expressions*, 1–8. Association for Computational Linguistics.
- ▶ Barkema, H. (1994a). Determining the syntactic flexibility of idioms. *Realising and Using English Language Corpora*, 39–52.
- ▶ Bengoetxea, K. eta Gojenola, K. (2010). Application of different techniques to dependency parsing of Basque. *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, 31–39. Association for Computational Linguistics.
- ▶ Berry-Rogghe, G. (1974). Automatic identification of phrasal verbs. *Computers in the Humanities*, 16–26.

- ▶ Burger, H. (1998). *Phraseologie: eine Einführung am Beispiel des Deutschen*. Berlin: Erich Schmidt.
- ▶ Corpas Pastor, G. (1996). *Manual de Fraseología Española*. Gredos, Madrid.
- ▶ Cowie, A. P. (1981). The treatment of collocations and idioms in learners' dictionaries. *Applied Linguistics*, 2(3):223–235.
- ▶ Cowie, A. P. (1998). Phraseological dictionaries: some east-west comparisons. Cowie, A., (ed.), *Phraseology: Theory, Analysis, and Applications*, 209–228. Oxford University Press, USA.
- ▶ Evert, S. (2005). *The Statistics of Word Cooccurrences: Word Pairs and Collocations*. Doktorego-tesia, University of Stuttgart.
- ▶ Fazly, A. eta Stevenson, S. (2007). Distinguishing subtypes of multiword expressions using linguistically-motivated statistical measures. *Proceedings of the Workshop on A Broader Perspective on Multiword Expressions*, 9–16. Association for Computational Linguistics.
- ▶ Fernando, C. eta Flavell, R. (1981). *On Idiom: Critical Views and Perspectives*. University of Exeter.
- ▶ Fontenelle, T. (1998). Discovering significant lexical functions in dictionary entries. Cowie, A., (ed.), *Phraseology: Theory, Analysis, and Applications*, 189–207. Oxford University Press, USA.
- ▶ Gläser, R. (1988). The grading of idiomaticity as a presupposition for a taxonomy of idioms. *Understanding the lexicon: Meaning, sense and world knowledge in lexical semantics*, 264–279.

- ▶ Gläser, R. (1998). The stylistic potential of phraseological units in the light of genre analysis. Cowie, A., (ed.), *Phraseology: Theory, Analysis, and Applications*, 125-143. Oxford University Press, USA.
- ▶ Harris, Z. (1954). Distributional structure. *Word*, 10(23):146-162.
- ▶ Hausmann, F. J. (1989). Le dictionnaire de collocations. *Wörterbücher, Dictionaries*, 1:1010-1019.
- ▶ Heid, U. (1994). On ways words work together - Research topics in lexical combinatorics. *Proceedings of the 6th International Congress of Lexicography - EURALEX '94*, 226-257.
- ▶ Howarth, P. A. (1996). *Phraseology in English academic writing: Some implications for language learning and dictionary making* 75. lib. Walter de Gruyter.
- ▶ Katz, G. eta Giesbrecht, E. (2006). Automatic identification of non-compositional multi-word expressions using Latent Semantic Analysis. *Proceedings of the Workshop on Multiword Expressions: Identifying and Exploiting Underlying Properties*, 12-19. Association for Computational Linguistics.
- ▶ Korkontzelos, I. eta Manandhar, S. (2009). Detecting compositionality in multi-word expressions. *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, 65-68. Association for Computational Linguistics.
- ▶ Krčmár, L., Jezek, K. eta Pecina, P. (2013). Determining compositionality of word expressions using word space models. *Proceedings of the 9th Workshop on Multiword Expressions (MWE 2013) NAACL HLT 2013*, 13:42-50.
- ▶ Krenn, B. (2004). Manual zur Identifikation von Funktionsverbgefügen und figurativen Ausdrücken in PP-Verb-Listen. *Austrian Research Institute for Artificial Intelligence*.

- ▶ Lin, D. (1998). Automatic retrieval and clustering of similar words. *Proceedings of the 17th International Conference on Computational Linguistics 2*. lib., 768–774. Association for Computational Linguistics.
- ▶ Lin, J., Li, S. eta Cai, Y. (2008). A new collocation extraction method combining multiple association measures. *Proceedings of International Conference on Machine Learning and Cybernetics 1*. lib., 12–17. IEEE.
- ▶ Mel'čuk, I. (1988). Semantic description of lexical units in an explanatory combinatorial dictionary: Basic principles and heuristic criteria. *International Journal of Lexicography*, 1(3):165–188.
- ▶ Mel'čuk, I. eta Wanner, L. (1994). Towards an efficient representation of restricted lexical cooccurrence. *Proceedings of the 6th International Congress of Lexicography - EURALEX '94* 94. lib.
- ▶ Mitchell, J. eta Lapata, M. (2008). Vector-based models of semantic composition. *Proceedings of ACL-08: HLT*, 236–244, Columbus, Ohio, AEB. Association for Computational Linguistics.
- ▶ Moon, R. (1998). *Fixed Expressions and Idioms in English: A Corpus-based Approach*. Clarendon Press Oxford.
- ▶ Pawley, A. eta Syder, F. (1983). Two puzzles for linguistic theory: Nativelike selection and nativelike fluency. *Language and Communication*, 191:225.
- ▶ Pecina, P. eta Schlesinger, P. (2006). Combining association measures for collocation extraction. *Proceedings of the COLING/ACL on Main Conference Poster Sessions*, 651–658. Association for Computational Linguistics.



- ▶ Pedersen, T., Banerjee, S., McInnes, B. T., Kohli, S., Joshi, M. eta Liu, Y. (2011). The Ngram Statistics Package (text::NSP): A flexible tool for identifying Ngrams, collocations, and word associations. *Proceedings of the Workshop on Multiword Expressions: from Parsing and Generation to the Real World*, 131–133. Association for Computational Linguistics.
- ▶ Reddy, S., McCarthy, D., Manandhar, S., eta Gella, S. (2011). Exemplar-based word-space model for compositionality detection: Shared task system description. *Proceedings of the Workshop on Distributional Semantics and Compositionality*, 54–60. Association for Computational Linguistics.
- ▶ Ruiz Gurillo, L. (1998). Una clasificación no discreta de las unidades fraseológicas del español. Wotjak, G., (ed.), *Estudios de fraseología y fraseografía del español actual*, 13–37. Lingüística Iberoamericana.
- ▶ Salvador, V. (2000). Idiomaticitat i discurs prefabricat. Salvador, V. eta Piquer, A.(ed.), *El discurs prefabricat*, 19–31. Universitat Jaume I.
- ▶ Sinclair, J. (1996). The search for units of meaning. *Textus*, 9(1):75–106.
- ▶ Urizar, R. (2012). *Euskal lokuzioen tratamendu konputazionala*. Doktorego-tesia, Informatika Fakultatea, UPV/EHU, Donostia.
- ▶ Van de Cruys, T. eta Moirón, B. (2007). Semantics-based multiword expression extraction. *Proceedings of the Workshop on A Broader Perspective on Multiword Expressions*, 25–32. Association for Computational Linguistics.
- ▶ Venkatapathy, S. eta Joshi, A. K. (2005). Measuring the relative compositionality of verb-noun (VN) collocations by integrating features. *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing*, 899–906. Association for Computational Linguistics.

- ▶ Vinogradov, V. (1947). Ob osnovnuikh tipakh frazeologicheskikh edinits v russkom yazuike. A .A. Shakhmatov, 1864-1920. *Sbornik statey i materialov*, 339-364. Mosku: Nauka.
- ▶ Warren, B. (2005). A model of idiomaticity. *Nordic Journal of English Studies*, 4(1):35-54.
- ▶ Wulff, S. (2008). *Rethinking Idiomaticity*. Corpus and Discourse. Continuum International Publishing Group Ltd, New York.
- ▶ Zabala, I. (2004). Los predicados complejos en vasco. Zabala, I., Pérez Gaztelu, E. eta García, L.(ed.), *Las fronteras de la composición en lenguas románicas y en vasco*, 445-567. Universidad de Deusto, Servicio de Publicaciones, Donostia.