

TOPIC EXPRESSIONS IN SPANISH: contrasting corpus and questionnaire data in the analysis of prepositional synonymy

1. BACKGROUND

Spanish prepositions

- 17 basic prepositions, hundreds of compound prepositions (*NGLE*)
- The compound prepositions have appeared (historically) in order to specify certain meanings and relieve the basic prepositions of excessive usage work load
 - The topic meaning (‘about’) is one case in point (Halliday 1967)
- Basic prepositions *de* and *sobre* have always been used with this meaning (DE CIVITATE DEI, SED HAC SUPER RE NIMIS DIXI) (Bassols de Climent 1967, I)
- *Acerca de* and *en torno a* appear later as alternatives
- Further expressions: *en cuanto a*, (*con*) *respecto a/de*, *en relación a/con*, *en lo tocante a*,...

The four prepositional topic expressions

- de* (‘of, from’) polysemous, topic only one out of dozens of senses
- sobre* (‘on, over’) polysemous, topic only one out of a dozen or so senses
- acerca de* (‘about’) monosemous, formally related to locative *cerca de* (‘near, close to’)
- en torno a* (‘around’) polysemous: locative, approximative vs. abstract (=topic) meaning

Aim of the study

- Address the differences in usage between four near-synonymous prepositions
- Are the expressions synonymous or not? To what degree?
- Compare different analyses (Arppe 2006, 2008; Liu 2013, Vanhatalo 2003)
 - What can corpus analysis tell us about synonymy
 - Compare the results of corpus analysis with questionnaire data
 - What other information might be used and/or needed in order to reach a “full” account of this issue?

2. METHODOLOGY AND CORPUS

Corpus analysis

- Annotation of examples
- Individual characterization of each expression
- Quantitative analysis: multinomial logistic regression (SPSS)

Questionnaire data

- Quantitative and qualitative data analysis
- Comparison of the results

Corpus del español (Davies 2002-)

- Large diachronic corpus of Spanish (100 million words)
 - Comparable to the COCA-corpus
- 4 x 100 examples of the four topic expressions (TE)
- Manual annotation of the 400 examples
- 15 syntactic and four semantic factors

Syntactic factors

HEAD/GOV word class (N, V, A, 0 / N, pron, V, conj)
HEAD/GOV complexity (NP(1-3) vs. N)
HEAD/GOV number (sg vs. pl)
HEAD/GOV definiteness (def vs. indef)
HEAD/GOV determinacy (det vs. indet)
HEAD/GOV modification (modif vs. unmodif)
HEAD/GOV attribute (attr vs. no attr)
HEAD before/after the TE

Semantic factors

Word class (of HEAD and GOV element) = Communication/cognition/action/general Animacy (of HEAD / GOV element) = Human/Non-human/collective/Unspecified Presence of other Topic expression (e.g. *de* and *sobre*) in the same clause Abstract/figurative vs. concrete reading

MULTINOMIAL LOGISTIC REGRESSION DATA

| Case Processing Summary | | | | Model Fitting Information | | | | | |
|-------------------------|---------------|-----------------|---------------------|---------------------------|-------------------------|------------------------|---------|----|------|
| Marker | de | N | Missing Percentages | Model | Fitting Criteria | Likelihood Ratio Tests | | | |
| sobre | Intercept | 89 | 25,3% | Model | -2 Log Likelihood df | Chi-Square df | Sig. | | |
| | sobre | 100 | 25,3% | | | | | | |
| | acerca de | 95 | 24,2% | | | | | | |
| | en torno a | 88 | 25,0% | | | | | | |
| Valid | Intercept | 392 | 100,0% | Intercept Only | 594,716 | 271,840 | 323,075 | 24 | ,000 |
| | Missing | 0 | | | | | | | |
| | Total | 400 | | | | | | | |
| | Subpopulation | 65 ^a | | | | | | | |

a. The dependent variable has only one value observed in 24,00% subpopulations.

| Parameter Estimates | | | | | | | 95% Confidence Interval for Exp. (B) | | |
|---------------------|----------------|------------|-------|----------|------|--------|--------------------------------------|-------------|---------|
| Marker ^a | B | Std. Error | Wald | df | Sig. | Exp(B) | Lower Bound | Upper Bound | |
| sobre | Intercept | -5,402 | 1,044 | 26,762 | 1 | ,000 | | | |
| | O_fig | 1,931 | ,660 | 7,687 | 1 | ,006 | 6,240 | 1,710 | 22,768 |
| | Ag_anim_hum | -1,272 | ,813 | 6,140 | 1 | ,013 | ,280 | ,102 | ,778 |
| | HEAD_1 | 2,496 | ,771 | 10,160 | 1 | ,001 | 11,860 | 2,576 | 52,969 |
| | Genre_News | 21,476 | ,414 | 2695,315 | 1 | ,000 | 21,2369 | 9,43869 | 47,4465 |
| | Genre_Fiction | 3,069 | ,669 | 21,039 | 1 | ,000 | 21,520 | 5,799 | 79,870 |
| | Genre_Academic | 21,013 | ,485 | 1878,687 | 1 | ,000 | 1,39659 | 5,16569 | 3,45568 |
| | O_det_det | 1,215 | ,445 | 4,352 | 1 | ,037 | 2,546 | 1,056 | 6,125 |
| | H_def_active | 2,380 | ,842 | 7,990 | 1 | ,005 | 10,810 | 2,075 | 56,233 |
| | Intercept | -5,439 | ,997 | 29,787 | 1 | ,000 | | | |
| acerca de | O_fig | 2,496 | ,686 | 13,223 | 1 | ,000 | 12,134 | 3,160 | 46,589 |
| | Ag_anim_hum | ,976 | ,898 | 3,055 | 1 | ,080 | ,377 | ,137 | ,109 |
| | HEAD_1 | 1,058 | ,665 | 2,526 | 1 | ,112 | 2,881 | ,782 | 10,160 |
| | Genre_News | 21,527 | ,445 | 2335,827 | 1 | ,000 | 2,23469 | 9,33269 | 5,74469 |
| | Genre_Fiction | 3,615 | ,646 | 31,347 | 1 | ,000 | 37,163 | 10,482 | 131,765 |
| | Genre_Academic | 21,243 | ,514 | 1707,846 | 1 | ,000 | 1,80269 | 6,14369 | 8,00029 |
| | H_def_det | 1,215 | ,457 | 7,059 | 1 | ,006 | 2,386 | 1,375 | 4,252 |
| | H_def_active | 1,559 | ,731 | 4,555 | 1 | ,033 | 4,756 | 1,136 | 19,912 |
| | Intercept | -5,268 | 1,121 | 21,369 | 1 | ,000 | | | |
| | O_fig | 2,238 | ,712 | 10,081 | 1 | ,001 | 10,261 | 2,541 | 41,431 |
| en torno a | Ag_anim_hum | -1,289 | ,827 | 12,206 | 1 | ,001 | ,112 | ,040 | ,351 |
| | HEAD_1 | 2,821 | ,811 | 12,123 | 1 | ,001 | 16,794 | 3,427 | 82,296 |
| | Genre_News | 21,132 | ,400 | 2711,001 | 1 | ,000 | 1,50689 | 1,06469 | 1,56869 |
| | Genre_Fiction | 2,400 | ,736 | 10,821 | 1 | ,001 | 11,025 | 2,603 | 44,692 |
| | Genre_Academic | 20,283 | ,500 | | 1 | ,000 | 6,43768 | 6,43768 | 6,43768 |
| | O_det_det | 1,826 | ,485 | 14,199 | 1 | ,000 | 6,207 | 2,401 | 16,043 |
| | H_def_active | 1,924 | ,899 | 9,081 | 1 | ,002 | 16,848 | 2,998 | 97,932 |

a. The reference category is: de.

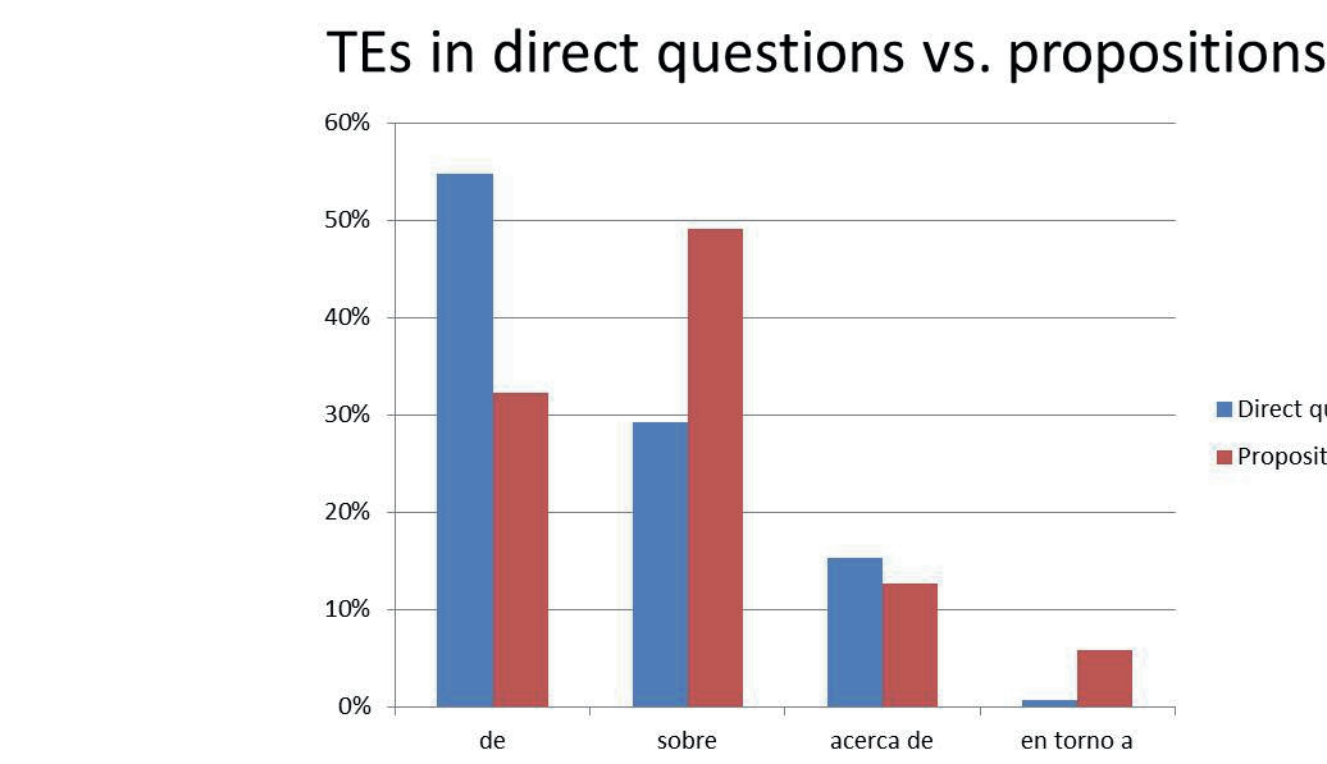
| Likelihood Ratio Tests | | | | |
|------------------------|------------------------------------|------------|----|------|
| Effect | -2 Log Likelihood of Reduced Model | Chi-Square | df | Sig. |
| Intercept | 342,333 | 70,892 | 3 | ,000 |
| O_fig | 290,544 | 18,903 | 3 | ,000 |
| Ag_anim_hum | 294,938 | 23,297 | 3 | ,000 |
| HEAD_1 | 294,173 | 22,532 | 3 | ,000 |
| Genre_News | 412,033 | 140,393 | 3 | ,000 |
| Genre_Fiction | 320,882 | 48,241 | 3 | ,000 |
| Genre_Academic | 338,035 | 67,395 | 3 | ,000 |
| O_det_det | 288,290 | 16,650 | 3 | ,001 |
| H_def_active | 286,015 | 14,375 | 3 | ,002 |

The chi-square statistic is the difference in -2 log-likelihoods between the final model and a reduced model. The reduced model is formed by omitting an effect from the final model. The null hypothesis is that all parameters of that effect are 0.

Data from questionnaires A + B

| HEAD | de | sobre | acerca de | en torno a | Mean |
|------|-------------|-------------|-------------|--------------|---------------|
| N | 59 % / 56 % | 48 % / 64 % | 54 % / 48 % | 56 % / 0 % | 54 % |
| V | 41 % / 44 % | 52 % / 36 % | 46 % / 52 % | 44 % / 100 % | 46 % |
| sum | 100 % | 100 % | 100 % | 100 % | 100 % |
| Com | 43 % / 51 % | 43 % / 46 % | 46 % / 53 % | 24 % / 0 % | 39 % |
| Cog | 34 % / 29 % | 31 % / 30 % | 25 % / 26 % | 27 % / 5 % | 29 % |
| Gral | 23 % / 20 % | 26 % / 24 % | 28 % / 21 % | 49 % / 95 % | 32 % |
| sum | 100 % | 100 % | 100 % | 100 % | 100 % |
| N | 423 / 452 | 282 / 349 | 153 / 85 | 79 / 21 | 937 / 907 |
| | 30 % / 38 % | 45 % / 50 % | 16 % / 9 % | 8 % / 2 % | 100 % / 100 % |

- High correlation ($r = 0,98$) between the two questionnaires
- A: Topic expressions given, 49 TEs to fill in ($N = 937$); B: no topic expressions given, 53 TEs to fill in ($N = 1040$)
- First two rows (N, V): A) $p \leq 0$, $\chi^2 = 12,8$ (3 df); B) $p \leq 0$, $\chi^2 = 18,5$ (df = 3)
- Following three rows: A) (Com, Cog, Gral): $p = 0,045$, $\chi^2 = 12,8$ (6 df); B) $p \leq 0$, $\chi^2 = 30,5$ (df = 6)



4. RESULTS AND DISCUSSION

What do the different approaches tell us about the synonymy relations between the four TEs?

- Different approaches yields different kinds of results, i.e. the TEs can be distinguished on different levels, depending on what one looks at and for
- *De* and *sobre* are the default TEs, *acerca de* and *en torno de* are alternatives
- *De* is lexically more restricted than *sobre*, but syntactically more independent
- **Multinomial regression analysis** brings forward very detailed information
 - *De* and *en torno a* are best described by the model
 - No significant effect of the semantic classes (com, cog, gral)
- The **questionnaire data** –analyzed on a more holistic level and not annotated– confirm some of the observations made on the corpus data
 - *sobre* is the default TE, followed by *de* and, to a lesser extent, *tent*, *acerca de*
- But the questionnaires also highlights other aspects
 - Highly marginal use of *en torno a* when not explicitly mentioned
 - *En torno a* is preferredly used with verbs (vs. nouns in log. regr.)
 - *Acerca de* is preferred in more “formal” contexts
- What I haven’t done:
 - Detailed collocational/collostructional analysis
 - Annotate the questionnaire data and compare it with the corpus data
 - Thorough analysis of the answers to the open questions included in the questionnaires

References

- Arppe, Antti (2006): “Complex phenomena deserve complex explanations – choosing how to think in Finnish”, paper presented at the QUITL2 Conference, Osnabrück, Germany, 2 June 2006.
- Arppe, Antti (2008): *Univariate, bivariate, and multivariate methods in corpus-based lexicography – a study of synonymy*. Publications of the Department of General Linguistics, University of Helsinki.
- Bassols de Climent, Mariano (1967): *Sintaxis latina*. 2 vols. Madrid: CSIC.
- Davies, Mark (2002-): *Corpus del español. 100 million words. 1200s-1900s*. Available on line at <http://www.corpusdelespanol.org>.
- Halliday, Michael A.K. (1967): “Notes on transitivity and theme in English. Part 2”, *Journal of Linguistics*, 3, 199-244.
- Liu, Dilin (2013): “Salience and construal in the use of synonymy: A study of two sets of near-synonymous nouns”, *Cognitive Linguistics*, 24-1, 67-113.
- NGLE = Real Academia Española y AALE (2009): *Nueva gramática de la Lengua Española*. 2 vols. Madrid: Espasa.
- Vanhatalo, Ulla (2003): “Kyselytestit vs. korpuslingvistiikka lähisynonyymien semanttisten sisältöjen arvioinnissa – Mitä vielä keskeistä ja tärkeistä?”, *Virtutiäjä*, 107:3, pp. 351-369.

| anim | uniqu | uniqu+ | H_fig | H_wc | H_wc+ | H_lex | HEAD | H_num | H_def | H_det | H_modif | H_attr | H_compl | H_pos | Nr. | Nr. | Source | Genre | Marker | EXAMPLE | G_lex | GOV | G_num | G_def | G_det | G_modif | G_attr | G_compl | G_fig | G_wc | | |
|------|-------|--------|-------|------|-------|-------|--------------|---------------|-------|-------|---------|--------|---------|-------|-----|-----|--------|-----------------------|-----------------------|--------------------------|----------------------------|------------------------|-------|-------|-------|---------|--------|---------|-------|------|----------|---------|
| col. | 1 | 1 | 1 | fig | gral | m_com | informe | N | sg | 0 | 1 | 0 | 0 | 1 | pre | 1 | 2 | Bolivia:ERBOL:N | | sobre | la Comisión de acciones | N | sg | 1 | 1 | 1 | 1 | 1 | 3 | 1 | act | |
| hum | 1 | 1 | 1 | fig | com | com | conferencias | N | pl | 1 | 1 | 0 | 1 | 2 | pre | 2 | 5 | Enc: Crítica litet:AC | | sobre | Ralph Waldo El principio | N | sg | 1 | 1 | 1 | 1 | 0 | 2 | 1 | abs_obj | |
| hum | 1 | 1 | 1 | conc | com | com | platicar | V | sg | act | inf | | 0 | 1 | 1 | pre | 12 | 40 | Entrevista (PAI-OR) | sobre | y en las platicas tema | N | sg | 1 | 1 | 1 | 0 | 0 | 1 | 1 | abs_obj | |
| | | | | | | | O | | | | | | | | | | 128 | 88 | Enc: José Gao: AC | en torno a | (El pensamiento) filosofía | N | sg | 1 | 1 | 1 | 1 | 0 | 2 | 1 | abs_obj | |
| | 0 | 0 | 0 | de | fig | gral | act | gírar | V | sg | act | inf | | 0 | 0 | 0 | pre | 139 | 111 | Enc: Enseñanza:AC | en torno a | como tal, lo cie idea | N | sg | 1 | 1 | 1 | 0 | 1 | 2 | 1 | cog |
| hum | 1 | 1 | 1 | conc | com | com | questionarie | N | pl | 0 | 0 | 0 | 1 | 1 | 2 | pre | 185 | 240 | Entrevista (PAI-OR) | acerca de | consenso total lo que | pron | sg | 1 | 1 | 1 | 0 | 1 | 2 | 1 | pred_cog | |
| unsp | 1 | 1 | 1 | fig | cog | cog | dudas | N | pl | 0 | 0 | 0 | 0 | 0 | pre | 216 | 87 | Arg:Prensa:21_N | acerca de | El ministro de (que con) | sg | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | abs_obj | |
| hum | 0 | 0 | 0 | de | conc | com | com | hablar pestes | V | sg | act | ind | | 1 | 0 | 1 | pre | 273 | 117 | Táhtalo en el tr | acerca de | el cuñado de A modales | N | pl | 1 | 1 | 1 | 0 | 1 | 2 | 1 | abs_obj |
| hum | 1 | 1 | 1 | fig | gral | act | hacer | V | pl | act | inf | | 1 | 0 | 1 | pre | 274 | 83 | Entrevista (Zed-OR) | acerca de | muchos años y esto | pron | sg | 0 | 0 | 0 | 0 | 0 | 0 | 1 | abs_obj | |
| hum | 1 | 1 | 1 | fig | cog | cog | contento | A | sg | 1 | 0 | 0 | 0 | 0 | 0 | pre | 301 | 574 | Entrevista (ABC19-OR) | de | , tiene en la caltener | V | sg | act | inf | | | | | 0 | abs_obj | |
| unsp | 1 | 1 | 1 | fig | cog | cog | eso | pron | sg | 0 | 1 | 0 | 0 | 0 | 1 | pre | 326 | 1450 | Habla Culta: M19-OR | de | No, no, no, no, conservar | V | sg | act | inf | | | | | 0 | act | |
| hum | 1 | 1 | 1 | conc | com | com | hablar | V | sg | act | ind | | 0 | 0 | 0 | pre | 365 | 154 | Habla Culta: La19-OR | de | cristiano. Esto iesto misi | pron | sg | 0 | 0 | 1 | 0 | 1 | 1 | 1 | abs_obj | |

3. DATA

Corpus del Español – overview

| PREP | de | sobre | acerca de | en torno a | Total |
|------------------------|-----------|--------|-----------|------------|-----------|
| topic uses | 100 | 100 | 100 | 100 | 400 |
| total | 1334 | 255 | 100 | 290 | 1978 |
| percentage | 7,5 % | 39,2 % | 100 % | 34,5 % | 20,2 % |
| Total uses | 1 163 904 | 33 572 | 782 | 1006 | 1 199 264 |
| Expected nr of TE uses | 87 315 | 13 165 | 782 | 347 | 242 521 |

• Characteristics of Head element (HEAD)

| HEAD | de | sobre | acerca de | en torno a | mean |
|------|-------|-------|-----------|------------|-------|
| N | 37 % | 73 % | 51 % | 78 % | 60 % |
| V | 54 % | 27 % | 42 % | 20 % | 36 % |
| A | 9 % | 0 % | 2 % | 0 % | 3 % |
| O | 0 % | 0 % | 5 % | 2 % | 2 % |
| Sum | 100 % | 100 % | 100 % | 100 % | 100 % |

• Characteristics of Governed element (GOV)

| GOV | de | sobre | acerca de | en torno a | mean |
|------|-------|-------|-----------|------------|-------|
| N | 61 % | 91 % | 80 % | 87 % | 79 % |
| pron | 24 % | 9 % | 15 % | 10 % | 15 % |
| V | 10 % | 0 % | 1 % | 0 % | 3 % |
| conj | 5 % | 0 % | 4 % | 3 % | 3 % |
| sum | 100 % | 100 % | 100 % | 100 % | 100 % |

p = 0,000, $\chi^2 = 67,426$ (9 df); $p \leq 0$, $\chi^2 = 47,466$ (9 df)
N = noun, V = verb, A = adjective, O = no head; N = noun, pron = pronoun,
V = verb, conj = (subordinate clause introduced by a) conjunction