

Words and their secrets

Diana Santos & Maria José Bocorny Finatto

Technologies Group

ESSLLI 2010, Copenhagen

WATS: Words and their secrets, ESSLLI 2010 Diana Santos & Maria José Bocorny Finatto SINTEF Natural Language



References for ''Words and Their Secrets'' at ESSLLI 2010

Diana Santos & Maria José Bocorny Finatto

Works cited in the course

Baayen, R. Harald. Word frequency distributions. Kluwer Academic Publishers, 2001.

- Bacelar do Nascimento, Maria Fernanda, José Bettencourt Gonçalves, Lucília Chacoto, Paula Neto & Luísa Alice Santos Pereira. "Ambiguidade morfológica no Português Fundamental". In Actas do 1º Encontro de Processamento da Língua Portuguesa (escrita e falada) (EPLP'93) (Lisboa, 25-26 de Fevereiro de 1993), 1993, pp. 101-106.
- Bick, Eckhard. The Parsing System "Palavras": Automatic Grammatical Analysis of Portuguese in a Constraint Grammar Framework. Aarhus University Press, 2000.
- Bindi, Remo, Nicoletta Calzolari, Monica Monachini, Vito Pirrelli & Antonio Zampolli. "Corpora and Computational Lexica: Integration of Different Methodologies of Lexical Knowledge Acquisition", *Literary and Linguistic Computing* **9**, No. 1, 1994, pp. 29-46.
- Bortolini, U., C. Tagliavini & A. Zampolli. Lessico di frequenza della lingua italiana contemporanea. IBM Italia, 1981.
- Bowerman, Melissa. "The origins of children's spatial semantic categories: cognitive versus linguistic determinants", in John Gumperz & Stephen C. Levinson (eds.), *Rethinking linguistic relativity*. Cambridge University Press, Cambridge, pp. 145-176.
- Brill, Eric. "A simple rule-based part of speech tagger", *Proceedings of the Third Conference on Applied Natural Language Processing* (Trento, Italy), 1992, pp. 152-155.
- Nicoletta & Remo Bindi. "Acquisition of Lexical Information from a Large Textual Italian Corpus", in Hans Karlgren (ed.), *Proceedings of COLING'90* (Helsinki, August 1990), Vol 1, pp. 54-59.
- Carlson, Lauri. "Aspect and Quantification". In Philip Tedeschi & Annie Zaenen (eds.), Syntax and Semantics, Volume 14: Tense and Aspect, Academic Press, 1981, pp. 31-64.
- Catford, J.C. A Linguistic Theory of Translation: An Essay in Applied Linguistics, Oxford University Press, 1967.
- Cherry, Lorinda L. "PARTS A System for Assigning Word Classes to English Text", Computer Science Technical Report #81, Bell Lab., Murray Hill, N.J., 1978.
- Chesterman, Andrew. Contrastive functional analysis. Amsterdam: Benjamins, 1998.
- Church, Kenneth Ward. "A stochastic Parts Program and Noun Phrase Parser for Unrestricted Text", *Proceedings of the Second Conference on Applied Natural Language Processing* (ACL), 1988, pp. 136-143.
- Church, Kenneth & William Gale. "Inverse document frequency (IDF): a measure of deviations from Poisson". In David Yarowsky & Kenneth Church (eds.), *Proceedings of the Third Workshop on Very Large Corpora* (Cambridge, MA, EUA, 30 June 1995), 1995a, pp. 121-130.
- Church, Kenneth & William Gale. "Poisson mixtures", *Journal of Natural Language Engineering* 1, 2, 1995b, pp. 163-190.
- Church, Kenneth & Patrick Hanks. "Word Association Norms, Mutual Information and Lexicography", *Computational Linguistics* **16**, 1, 1991, pp. 22-29, 1991.
- Cruse, Alan. *Meaning in Language: An Introduction to Semantics and Pragmatics*. Oxford. Oxford University Press, 2004.
- DeRose, Stephen J. "Grammatical category disambiguation by statistical optimization", *Computational Linguistics* **14**, 1, Jan. 1988, pp. 31-39.

- Dixon, R.M.W. "A method of semantic description", in Danny D. Steinberg & Leon A. Jakobovits (eds), *Semantics: An interdisciplinary reader in philosophy, linguistics and philosophy*, Cambridge: Cambridge University Press, 1971, pp. 436-471.
- Dorow, Beate. "A Graph Model for Words and their Meanings". PhD Thesis, IMS, Stuttgart University, 2006. <u>http://elib.uni-stuttgart.de/opus/volltexte/2007/2985/pdf/diss_27022007.pdf</u>.
- Dunning, Ted. "Accurate Methods for the Statistics of Surprise and Coincidence", *Computational Linguistics* **19**, Number 1, March 1993, pp. 61-74.
- Edmonds, Philip & Graeme Hirst. "Reconciling fine-grained lexical knowledge and coarse-grained ontologies in the representation of near-synonyms". In *Proceedings of the Workshop on Semantic Approximation, Granularity, and Vagueness (KR-2000)*, Breckenridge, Colorado, 2000.
- Ellegård, Alvar. The syntactic structure of English texts: a computer-based study of four kinds of text in the Brown University corpus. Gothenburg: Acta Universitatis Gothoburgensis, 1970.
- Ellis, John M. Language, Thought and Logic. Evanston, IL: Northwestern University Press, 1993.
- Fellbaum, Christiane (ed.). *WordNet: An Electronic Lexical Database*, with a preface by George Miller. The MIT Press, May 1998.
- Garside, Roger, Geoffrey Leech & Geoffrey Sampson. *The Computational Analysis of English: A Corpus-Based Approach*, Longman, 1987.
- Goodman, Nelson. "Seven strictures on similarity". In Nelson Goodman (ed.), Problems and projects. Indianapolis, IN: Bobbs-Merrill, 1972, pp. 437-447.
- Gonçalo Oliveira, Hugo, Diana Santos & Paulo Gomes. "Extracção de relações semânticas entre palavras a partir de um dicionário: o PAPEL e sua avaliação", *Linguamática* **2**, 1, April 2010, pp. 77-94.
- Green, T. R. G. "The Necessity of Syntax Markers: Two Experiments with Artificial Languages", *Journal* of Verbal Learning and Verbal Behavior 18, 4, Aug 1979, pp. 481-496.
- Greene, Barbara B. & Gerald M. Rubin. "Automated Grammatical Tagging of English". Providence, R.I.: Department of Linguistics, Brown University, 1971.
- Grefenstette, Gregory. *Explorations in Automatic Thesaurus Discovery*. Boston, MA: Kluwer Academic Press, 1994.
- Grefenstette, Gregory & Pasi Tapanainen. "What is a word, What is a sentence? Problems of Tokenization", *Proceedings of the 3rd International Conference on Computational Lexicography* (COMPLEX'94), 1994, pp. 79-87.
- Gruber, Thomas R. "A Translation Approach to Portable Ontology Specifications", *Knowledge* Acquisition 5(2), 1993, pp. 199-220.
- Halliday, M.A.K. *Computational and Quantitative Studies*, vol 7 in the Collected Works of MAK Halliday edited by Jonathan J. Webster, London, New York: Continuum, 2005.
- He, Ying & Mehmet Kayaalp. "A Comparison of 13 Tokenizers on MEDLINE", Technical Report LHNCBC-TR-2006-003. <u>http://lhncbc.nlm.nih.gov/lhc/docs/reports/2006/tr2006003.pdf</u>
- Heiden, Serge. "Interface hypertextuelle à un espace de cooccurrences: implémentation dans Weblex", in Gérard Purnelle, Cédrick Fairon & Anne Dister (eds.), 7^{ième} Journées internationales d'Analyse Statistique des Données Textuelles (JADT'04) "Le poids des mots" 10 - 12 Mars 2004, vol 1, Presses Universitaires de Louvain, Louvain-la-Neuve, Belgique, 2004, pp. 577-588.
- Hindle, Donald. "Acquiring Disambiguation Rules from Text", Proceedings of ACL 1989, pp. 118-125.
- Hindle, Donald & Mats Rooth. "Structural Ambiguity and Lexical Relations", *Computational Linguistics* **19**, 1, March 1993, pp. 103-120.
- Hirst, Graeme. "Ontology and the lexicon". In Steffen Staab & Rudi Studer (eds.). Handbook on ontologies, Springer, 2004, pp. 209-229.

- Jurafsky, Daniel & James Martin. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Prentice Hall, 2000.
- Justeson, John S. & Slava M. Katz. "Technical terminology: some linguistic properties and an algorithm for identification in text", *Natural Language Engineering* **1**, 1995, pp. 9-27.
- Katz, Slava M. "Distribution of content words and phrases in text and language modelling", *Natural Language Engineering* **2**, 1996, pp.15-59.
- Kennedy, Graeme. "Over once lightly". In Carol E. Percy, Charles F. Meyer & Ian Lancashire (eds.), Synchronic corpus linguistics: papers from the sixteenth International Conference on English Language Research on Computerized Corpora (ICAME 16), Rodopi, Amsterdam – Atlanta, GA, 1996, pp. 253-62.
- Kilgarriff, Adam. "Which words are particularly characteristic of a text? A survey of statistical approaches", *Proceedings of AISB Workshop on Language Engineering for Document Analysis and Recognition* (Sussex, April 1996), pp. 33-40.
- Kilgarriff, Adam. ""I don't believe in word senses"", *Computers and the Humanities* **31** (2), 1997, pp. 91-113.
- Kilgarriff, Adam. "Language is never ever random", *Corpus Linguistics and Linguistic Theory* **1**, 2, 2005, pp. 263-276.
- Kilkki, Kalevi. "A practical model for analyzing long tails", *First Monday* **12**, 5, May 2007, <u>http://www.firstmonday.org/issues/issue12_5/kilkki/</u>
- Klein, Sheldon & Robert F. Simmons. "A computational approach to grammatical coding of English words", *Journal of the Association for Computing Machinery* **10**: 334-347.
- Krenn, Brigitte & Christer Samuelsson. "The Linguist's Guide to Statistics: DON'T PANIC", 21 May 1997.
- Macklovitch, Elliott. "Where the Tagger Falters", Proceedings of the 4th International Conference on Theoretical and Methodological Issues in Machine Translation (Montréal, June 25-27, 1992), pp. 113-126.
- Manning, Chris & Hinrich Schütze. Foundations of Statistical Natural Language Processing. MIT Press, May 1999.
- Marshall, I. "Choice of grammatical word-class without global syntactic analysis: Tagging words in the LOB Corpus", *Computers in the Humanities* **17** (1983), pp. 139-150.
- Martins,JoãoPavão.KnowledgeRepresentation.IST.http://www.cse.buffalo.edu/~rapaport/663/S02/martinssneps.pdfIST.
- Medeiros, José Carlos. "Avaliação de Correctores Ortográficos", Actas do XI Encontro da Associação Portuguesa de Linguística, Lisbon: Colibri, 1996, pp. 73-91.
- Medeiros, José Carlos, Rui Marques & Diana Santos. "Português Quantitativo", Actas do 1.0 Encontro de Processamento de Língua Portuguesa (Escrita e Falada) - EPLP'93 (Lisbon, 25-26 February 1993), 1993, pp. 33-38.
- Miller, George A., Richard Beckwith, Christiane Fellbaum, Derek Gross & Katherine Miller. "Introduction to WordNet: An On-line Lexical Database" (revised August 1993), *Five papers on WordNet*, 1993, pp. 1-25.
- Monachini, Monica & Nicoletta Calzolari. "Standardization in the lexicon". In Hans van Halteren (ed.), *Syntactic Wordclass Tagging*, Dordrecht/Boston/London: Kluwer Academic Publishers, 1999, pp. 149-174.
- Mosteller, Frederick & David L. Wallace. Inference and Disputed Authorship. 1964.

- Nicolaeva, T. M. "Soviet Developments in Machine Translation: Russian Sentence Analysis", *Mechanical Translation* **5**, 2, November 1958, pp. 51-59.
- Pinker, Steven. The stuff of thought: Language as a Window into Human Nature. Allen Lane, 2007.
- Pym, Anthony. Translation and Text Transfer. An Essay on the Principles of Intercultural Communication. Frankfurt am Main, Berlin, Bern, New York, Paris, Vienna: Peter Lang, 1992. Revised online version, Tarragona: Intercultural Studies Group, 2010, http://www.tinet.cat/~apym/publications/TTT_2010.pdf
- Richardson, Stephen. "Determining Similarity and Inferring Relations in a Lexical Knowledge Base", Ph.D. thesis, The City University of New York, 1997, Microsoft Research Report MSR-TR-97-02, <u>ftp://ftp.research.microsoft.com/pub/tr/tr-97-02.doc</u>.
- Rivenc, Paul. "Vocabulário frequente e vocabulário disponível". In Bacelar do Nascimento, Maria Fernanda, Paul Rivenc & Maria Luísa Segura da Cruz. *Português Fundamental*, Volume II, *Métodos e Documentos*, tomo 2, *Inquérito de Disponibilidade*, Lisbon: Centro de Linguística de Universidade de Lisboa, 1987, pp. 3-26.
- Ruppenhofer, Josef, Michael Ellsworth, Miriam R. L. Petruck, Christopher R. Johnson & Jan Scheffczyk. *FrameNet II: Extended Theory and Practice*. Printed June 15, 2010. <u>http://framenet.icsi.berkeley.edu/index.php?option=com_wrapper&Itemid=126</u>
- Sampson, Geoffrey. *English for the Computer: The SUSANNE Corpus and Analytic Scheme*. Clarendon Press, Oxford, 1995.
- Sampson, G.R. "Review of Christiane Fellbaum (ed.), Wordnet: An Electronic Lexical Database, 1998", International Journal of Lexicography 13, 2000, pp. 54-59.
- Sampson, Geoffrey. *The 'Language Instinct' Debate*. March 2005. London & New York: Continuum International. Enlarged and revised edition of *Educating Eve*, Cassell, 1997.
- Santos, Diana. Translation-based corpus studies: Contrasting English and Portuguese tense and aspect systems. Amsterdam/New York, NY: Rodopi, 2004.
- Santos, Diana. "What is natural language? Differences compared to artificial languages, and consequences for natural language processing". Invited lecture, SBLP2006 and PROPOR'2006, Itatiaia, RJ, Brazil, 15 May 2006. <u>http://www.linguateca.pt/Diana/download/SantosPalestraSBLPPropor2006.pdf</u>
- Santos, Diana & Caroline Gasperin. "Evaluation of parsed corpora: experiments in user-transparent and user-visible evaluation". In Manuel González Rodríguez & Carmen Paz Suárez Araujo (eds.), *Proceedings of LREC 2002, the Third International Conference on Language Resources and Evaluation* (Las Palmas de Gran Canaria, Spain, 29-31 May 2002), ELRA, 2002, pp. 597-604.
- Santos, Diana, Luís Costa & Paulo Rocha. "Cooperatively evaluating Portuguese morphology". In Nuno J. Mamede, Jorge Baptista, Isabel Trancoso & Maria das Graças Volpe Nunes (eds.), Computational Processing of the Portuguese Language, 6th International Workshop, PROPOR 2003, Faro, 26-27 June 2003, Proceedings, Springer Verlag, 2003, pp. 259-266.
- Saussure, Ferdinand de. *Cours de Linguistique Générale*, publié par Charles Bally & Albert Sechehaye avec la collaboration de Albert Riedlinger. Payot, Paris, 1972. First edition: 1915.
- Sinclair, John. "Corpus Evidence in Language Description". In Wichmann, Anne, Steven Fligelstone, Tony McEnery & Gerry Knowles (eds.), *Teaching and language corpora*. London & New York, Longman, 1997, pp. 27-39.
- Snell-Hornby, Mary. Verb-descriptivity in German and English: A contrastive study in semantic fields, Carl Winter Universitätsverlag, Heidelberg, 1983.
- Sovran, Tamar. "Between similarity and sameness", Journal of Pragmatics 18, 4, 1992, pp. 329-344.
- Sparck Jones, Karen. "What's new about the Semantic Web?: some questions", *ACM SIGIR Forum* **38**, 2, December 2004, COLUMN: Invited talks, pp. 18-23.

- Steiner, George. *After Babel: aspects of language and translation*. Oxford: Oxford University Press, 1992 (1st edition 1975).
- Stolz, Walter S., Percy H. Tannenbaum & Frederick V. Carstensen. "A stochastic approach to the grammatical coding of English", *Communications of the ACM* **8**, 6, June 1965, pp. 399-405.
- Talmy, Leonard. "How language structures space". In H. Pick & L. Acredolo (eds.), *Spatial orientation: theory, research, and application*. New York, Plenum Press, 1983, pp. 225-282.
- Tversky, Amos. "Features of similarity", Psychological Review 84, pp. 327-352.
- Underwood, Nancy, Patrizia Paggio & Gurli Rohde. "A methodology for evaluating Spelling Checker functionality: Developing test suites for Danish", in Kimmo Koskenniemi (ed.), Short papers presented at the Tenth Scandinavian Conference on Computational Linguistics (Helsinki, 29-30th May 1995), 1995, pp. 76-85.
- Veale, Tony. "Enriched Lexical Ontologies: Adding new knowledge and new scope to old linguistic resources", ESSLLI 2007, Dublin, <u>http://afflatus.ucd.ie/papers/Essilli_EnrichedLexiOnto.pdf</u>
- Wilks, Yorick. "Is Word Sense Disambiguation Just One More NLP Task?", *Computers and the Humanities* **34**, 1-2, April 2000, pp. 235-243.
- Wilks, Yorick & John Tait. "A Retrospective view of Synonymy and Semantic Classification". In John I. Tait (ed.), *Charting a New Course: Natural Language Processing and Information Retrieval: Essays in Honour of Karen Spärck Jones*, Springer, 2005, pp. 255-282.
- Wilks, Yorick & Christopher Brewster. "Natural Language Processing as a Foundation of the Semantic Web", *Foundations and Trends in Web Science* **1**, 3, March 2009, pp. 199-327.

Image credits

Contar carneiros: <u>http://www.oasrs.org/conteudo/agenda/noticias-detalhe.asp?noticia=1549</u> (obtained 4 June 2010).

Topic maps: <u>http://en.wikipedia.org/wiki/File:TopicMapKeyConcepts2.PNG</u> (obtained 16 June 2010)

Expectations: Joakim Krøvel, Scanpix, www.scanpix.no

Personal acknowledgements

Diana Santos thanks Danilo Giampiccolo for help in tracing Bortolini et al. (1981) right from Italy, Eric Atwell for permission to use an interesting email exchange, back in February 2010, Anton Landmark for the Humpty Dumpty quotation and for encouragement and criticism, Nuno Cardoso, Tormod Håvaldsrud and all other colleagues at SINTEF who attended a preliminary version of this course, Doris Lund from the SINTEF library for granting access to a vast number of books and papers, João Pavão Martins for the SNePs figure, Maarten Marx for permission to use his site and photographs therein, and last but not least, acknowledges financial support in the scope of the Linguateca project, co-financed by the Portuguese government, the European Union (FEDER and FSE) under POSC/339/1.3/C/NAC contract, UMIC and FCCN.



























- Direct reference
- Ostensive denotation
- Harwired in the brain (mental images)
- More enlightened views/aproaches
- - Words (and the other mechanisms of language) represent classes of different objects which are considered, for the purpose of conceptualization, as similar
 - Words are a prerequisite for thought and communication
 - Words (and the rest of language, including communication patterns) are learned through interaction with the language community (and especially the mother)







Institute of Language and Linguistics, Federal University of Rio Grande do Sul, Brazil

Technologies Group













Spell checking Issues in spell checking What to encode in the dictionary? Identify / detect incorrectly spelt words Rare words may correspond to errors Suggest corrections Some of the most frequent errors (exchange between false friends) can Automatically correct only be detected in context: Words are defined as sequences of "word-proper" characters, ■ its / it's (en) separated by word separators ■ å/og (no) "Incorrectly" spelt means ■ à/a (pt) two / to (en) not belong in the dictionary not being accepted by a set of (language-specific) given rules What if the error is absence or addition of word-separator? callback/ call back (this problem is compounded in languages with compounds) not numbers or simple letters Fee dback Avoid correction of (some) proper names WATS:: Basic technologies: spell checking and POS tagging WATS:: Basic technologies: spell checking and POS tagging



Fundação para a Competação Científica Nacional

Institute of Language and Linguistics, Federal University of Rio Grande do Sul, Brazil

Natural Language

Technologies Group

Spelling correction

- How to evaluate/rank the best suggestions?
- How to provide measures to compare different spellcheckers?
- Number of correct corrections/Number of (first) corrections suggested
- Number of correct corrections/Number of errors
- I don't likke cracrashes
- How to count the number of errors? Words with errors?
- And how to count the number of correct suggestions if the number of words can be different after correction?
- There are incorrect corrections which are nevertheless useful!
- WATS:: Basic technologies: spell checking and POS tagging

Further examples

Dirigi lhe	dirigi-lhe	2 or 1 words / 2 or 1 errors
Senti-la-hia	senti-la-ia	3 or 1 words / 1 error
Ta, to		
'Tás, 'tamos	estás, estam	os errors?
diversidade.Nesse	diversidade.	Nesse 1,2 or 3 words?
dêmo	dê-mo	2 or 1 words / 1 error
auto-denominado-se	auto-denom	inando-se
PhG	PhD	
rock'n'roll, Toys'R'Us,	90's, M'Gladb	ach, 2000-2010,
R&D, A4, UB40,		
WATS:: Basic technologies: spell c	hecking and POS tag	ging

















PoS tagging

- Apparently the easiest and best defined task...Manual or manually validated vs. automatic
- For each **word**, assign the correct part-of-speech
- Word?

Natural Language

Technologies Group

- And multiwords? And named entities? And non-words?
- For each word, assign the correct part-of-speech
 Correct? Depends on the theory of grammar
 - Correct? Depends on the theor
 Only one tag?
 - Evaluation of PoS tagging... what is correct?
- For each word, assign the correct part-of-speech (PoS)
 Only PoS? Or morphology as well? Or subcategorization? Or everything?

WATS:: Basic technologies: spell checking and POS tagging

Some history of POS tagging

- Apparently the first machine disambiguation of natural language text was done by Russian researchers working on MT (Nicolaeva, 1958)
- Klein & Simmons (1962) develop a first component in a syntactic analysis program, which is part of a larger QA system
- Stolz et al. (1965) apply statistical methods: decisions ... based on conditional probabilities of various form classes in given syntactic environments -- Cherry (1978) assigns part of speech by rule
- Green & Rubin (1971) create the first annotated corpus, the Brown corpus, human revised; and Ellegård (1970) the first human annotation
- Marshall (1983) improving LOB based on Brown
- DeRose (1988), Church (1988), Garside et al. (1987), Hindle (1989): POS tagging as pre-processing

WATS:: Basic technologies: spell checking and POS tagging

Macklovitch (1992): First linguistic analysis? Generally speaking, a given tag set may be more or less suitable for certain applications after, before, until can be either IN or CS andogous to suppressing the distinction between verbs that subcategorize for an NP or for a sentence... -ed forms can be either VBD or VBN nouns or adjectives: JJ or NN Global dependencies (instead of "long-distance"): whether a verb is in imperative, present or subjunctive can depend on the whole sentence. Why bother? Evaluation relevance Automatic error detection – and maybe even correction

Brill tagger (1992) learns from its weaknesses

- "A simple Rule-Based PoS Tagger": robust and rules automatically acquired
- Currently called a hybrid method, because it uses machine learning, but requiring human annotated data
- First it assigns the most frequent tag to the already existing words in the training material; then uses the word endings out of the dictionary
- Comparing the output to human annotated material, it creates error triples: <old category, new category, frequency>
- Eight different patches are tried out, and the one which provides higher global error diminishing is added to the patch list
- 71 patches, 5% error in 5% of the Brown corpus

WATS:: Basic technologies: spell checking and POS tagging



Institute of Language and Linguistics, Federal University of Rio Grande do Sul, Brazil

Fundação para a Computação Científica Nacional

Measuring Portuguese POS ambiguity

- Medeiros et al. (1993): potential word classes in a corpus
 n/a vnp, v. adv. pf. cl
 - n/a, vpp, v, adv, pf, cl
 1.02494 classifications per form; 1.1398 class/form if only the three first are considered
- Bacelar do Nascimento et al. (1993): real word classes in a corpus From a corpus of transcribed oral speech, 700,000 words (25,107 types), reduced to the forms corresponding to lemmas with frequency > 40 (1553 lemmas): 65,000 forms, where there were potentially 834 ambiguous lemmas, corresponding to 1371 POS-ambiguous form (types), whose occurrences were then analysed in context
 N-ADJ: 143 types: 123 Noun, 121 Adj
 - N-ADJ-V: 66 types: 44 Noun, 57 Adj., 35 Verb

WATS:: Basic technologies: spell checking and POS tagging

Kennedy about the value of POS tagging

- When claims are made about the impressive accuracy with which grammatical tags can be assigned by machine, it is often not made clear to consumers that the high success rates [] are based on an averaging process. [] certain very frequent words or word classes can be tagged with virtually total accuracy, while for other items, accuracy rates of 80-85% are more typical. (Kennedy, 1996:253)
- 100 most frequent word types in LOB -> 49% of the tokens
- Ca. 2/3 (65, types, ca. 335,000 tokens) belong to one class only!

WATS:: Basic technologies: spell checking and POS tagging



Concluding remarks

- Beware of "easy" tasks, light hearted procedures
- Even for the least intelectually challenging task... Criteria for "wordness" have to be thought and decided upon. In linguistic textbooks tokenization is quickly dispatched as a relatively uninteresting pre-processing step performed before linguistic analysis is undertaken. In reality. tokenization is a non-
- trivial problem (Grefenstette & Tapanainen, 1994)
 In the next days this will be shown in other fields on natural language processing as well ...

WATS:: Basic technologies: spell checking and POS tagging







Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds			
		Diana Santos	
VATS: Wo	ords and their s & Maria José Boc	secrets, ESSLLI 2010 orny Finatto	
iana Santos			





WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

Technologies Group

Fundamental vocabularies

- Frequency is not enough: what about availability?
- If *knife* is frequent, would not *fork* qualify as available and therefore required to be included as well?
- If one knows how to use *bachelor*, one knows the meaning of *married*

Rivenc (1987)	corpus	voc.	themes
Français fondamental	312135	806	15
Português fundamental	700000	1179	27+3
Español fundamental	800000	949	25
A list of themes/intere	est centers: elicit	ating words after	a theme (human

body, games, village, school, politics, ...). Threshold frequency: F₁ frequency of the highest ranked word, N the number of words requested, D is dispersion, K is a adjusted parameter

 $F_{L}=N^{*}K^{*}F_{l}/D$ WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds





Institute of Language and Linguistics, Federal University of Rio Grande do Sul, Brazil

Fundação para a Computação Científica Nacional



Natural Language

Technologies Group







Fundação para a Computação Científica Nacional



What is an ontology, part 2

- Even if it is not explicit, it always includes relations
- Does it include reasoning rules?
- Does it also include elements that can be obtained by reasoning?
- In the informatics community, Gruber's (1993) definition is accepted: An ontology is a formal explicit specification of a shared conceptualization for a domain of interest.
- In the linguistic community, I propose Veale's (2007) definition of lexical ontology: An ontology of lexical(-ized) concepts, used in NLP, serving as a lexical semantics (ESSLLI 2007, Enriched Lexical Ontologies)
- WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds



Same or similar revisited

- Similarity is relative, variable, culture dependant (Goodman, 1972)
- Circumstances alter similarities (Goodman, 1972)
- The similarity of objects is modified by the manner in which they are classified (Tversky, 1977)
 "similarity" is a sim that is attributed to a set of antitize, attributed by
- "similarity" is a sign that is attributed to a set of entities, attributed by someone and also interpreted by someone (Chesterman, 1998)
 similarity-as-trigger
 - similarity-as-attribution
- the greater the extension of the set of items assessed as being similar, the less the pertinent degree of similarity
- Tension between "oneness" and "separate individuation" (Sovran, 1992)

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

Example from Tversky (1977)

Question: To which country is Austria more similar to?

- Sweden, Poland, Hungary
- Sweden, Norway, Hungary

Let us try again

meanings

links)

- Germany, Denmark, The Netherlands
- Germany, Switzerland, The Netherlands

PhD thesis by Beate Dorow, IMS, 2006

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

A Graph Model for Words and their Meanings

consisting of words (the nodes) and relationships between them (the

Graph-theoretic approach to the automatic acquisition of word

[...] represent the nouns in a text in form of a semantic graph

Links in the graphs are based on cooccurrence of words in lists

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

Differences between languages

(1)"I want a different apple." "Why? They are all the same."

- (2) They wore the same dress.
- (3) I'll have the same as her (said to a waiter).
- (4) These two pens look similar, but one is more expensive than the other

English same is ambiguous between type and token identity

- Finnish: not the same item in (1) nor (2), but in (3).
- Portuguese: not the same item in (1): *são todas iguais*
- Portuguese: parecem iguais in (4)

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

SINTEF Natural Language Technologies Group





Sweden (49%)

Hungary (60%)









v of Rio Grande do Sul, Brazil

The \$\$ of plurals vs. singulars (Pinker 2007)

- Google sells index terms
 - In order for appropriate adverts to appear together with the results
 - "photo cameras" is more expensive than "photo camera"
 ... because it shows that people are undecided about which one to choose
- All conflation is of course reductive
 - squashes (En.) are ONLY vegetables, while squash is ambiguous (Dorow)
 pais (Pt.) can mean parents as well as fathers (plural of pai)
 - pais (1.2) can mean parents as well as failed's (plural of par)
 Bindi et al. (1994) describe the need for observation of word forms
 - contait (1/): ... Only three words out of twelve really apply to the lemma contatt The other nine either co-occur with the singular or with the plural

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

WordNet

- WordNet started as psycholexicologist's model of word meaning
- (psycholexicology = research concerned with the lexical component of language)
- An On-line Lexical Databae
- The initial idea was to "provide an aid to use in searching dictionaries conceptually (...) to be used in close conjunction with an on-line dictionary of the conventional type"
- Miller et al. (1993): a dictionary based on psycholinguistic principles
 - expose (psycholinguistic) hypotheses to the full range of vocabulary
 - organize lexical information in terms of word meanings, rather than word forms

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

WordNet ... and wordnets

- One the most well-know and used lexical resources for English
- An example/model for several other languages
- A lot of wordnets and wordnet-alignment word, Global WordNet conferences all around the world
- Free for use, abundant computational support
- Several new developments/augmentations:
 definitions, domains, addition of other sources, etc.
- But: are all uses warranted or appropriate? Is the underlying WordNet linguistic/semantic theory sound? Or applicable in every application?

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds

Sampson's (2000) critical remarks

- it seems surprising that a database constructed manually by academics with no access to a dictionary-publisher's archive could be a seriouis contender as the leading tool in this domain
- ... network of hyponymy relationships between nouns apparently requires some nodes which correspond to no single item of English
- The system is so naive that it (...) recognizes no distinction between the species/genus relationship, as in horse/animal, and the individual/universal relationship, as in Shakespeare/author, treating both indifferently as cases of "hyponymy"

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds





person

bad person

R

libertine









- When we say that the word *bake is* polysemous, we mean that the lemma *bake.v* (which has the word-forms *bake, bakes, baked,* and *baking*) is linked to three different frames:
 Apply heat: Michelle baked the potatoes for 45 minutes.
 Cooking creation: Michelle baked her mother a cake for her birthday.
 Absorb heat: The potatoes have to bake for more than 30 minutes.
- These constitute three different Lus [lexical units], with different definitions.
- Multiword expressions such as given name and hyphenated words like shut-eye can also be LUs.

Ruppenhofer et al. (2010)

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds



Concluding remarks

- There is a huge activity nowadays in (automatically or not) creating (lexical or not)) ontologies and merging or integrating them
- Unfortunately, many of the work is still based on ungrounded or naive assumptions
 - What is similarity
 - What is the purpose of the O
 - What are its units
- There are a lot of fancy tools and systems to deal with and visualize complex objects created from heaps of data

but their use is only as good as the underlying objects...

WATS:: Dictionaries, lexical networks, lexical ontologies, wordnets and wordclouds









Statistics is the branch of mathematics...









Institute of Language and Linguistics, Federal University of Rio Grande do Sul, Brazil

Fundação para a Computação Científica Nacional

Why connect the two disciplines at all?

- Because speech processing itself inspired by probability theory has been influential as a model for empirical language processing
 It remains to be assessed how successful statistical methods for speech have actually been
 - It remains to be assessed how close speech processing in fact is to machine translation (for example)
- Because information retrieval itself making heavy use of probabilistic models – has also been influential as a testbed for empirical language processing

WATS:: Lexical statistics



The dispersion index (Calzolari & Bindi, 1990)

- Measures the degree of fixity of the second word position with respect to the keyword, a measure of how frequency is distributed over the different positions of the window...
- Different slices of the multidimensional pie (the semantic hyperspace) carry with themselves a different bunch of word senses for the same word entry (Bindi et al., 1994)
- Italian quasi-synonyms: picollo, corto, breve, ristretto, esiguo, scarso, ridotto are the "units"
- Space: "word mates" (unambiguous words with whom the units keep company with, computed by mutual information)

WATS:: Lexical statistics

Natural Language

Technologies Group







Institute of Language and Linguistics, Federal University of Rio Grande do Sul, Brazil

Fundação para a Computação Científica Nacional

Predicting the occurrence of words

- A good keyword is one that behaves very differently from the null
- hypothesis (that the word is distributed according to a Poisson distr.)Variance and IDF correlate positively with good keywordness,
- entropy negatively
 Katz K-mixture has two parameters (α: fraction of relevant and
- irrelevant documents, and β : the average Poisson parameter) and corresponds to a convolution of Poisson distributions • $\beta = f/D^{*2DP} - 1$ Church & Gale (1995a)
 - $\alpha = f/D\beta$
- The main idea is that each Poisson distr. can model hidden variables such as what the documents are about, who wrote them, when they were written, what was going on in the world then

WATS:: Lexical statistics



Green (1979) and syntax markers

is required in order for a language to be learnable

problems about words (what are the units?)

-ing

WATS:: Lexical statistics

The marker hypothesis states roughly that "a small number of elements that signal the presence of particular syntactic constructions"

Markers in English: prepositions/closed words, suffixes such as -ly or

This is interesting food for thought also for natural language, although

the issue of what markers are is obviously subject to the same kind of



Katz (1996) continued

- Linguistically motivated approach (...) arriving at a coherent view of the word occurrence phenomenon without commitments to any particular, *a priori* assumed, stochastic mechanism
 - The probabilities of repeat occurrences do not depend on the relative frequency
 The continual presence of repeat occurrences in discourse is a general and widesread phenomenon () A principle distinction is identified between two
- The community presence on repeat occurrings in atomics in a general and widespread phenomenon (...) A principle distinction is identified between two probabilities of repeats (entering; and stayin in a document-level burst)
 Possion mixtures as two-stage stochastic mechanism for generating
- content works is incompatible with empirical data
- (discrete Poisson mixtures) limited in their capacity to provide satisfactory fit to the data because of their faulty functional form...

WATS:: Lexical statistics























Indexing

- This is the realm of information retrieval...
- Or the use of good "descriptors": what best than words themselves?
- Sparck Jones (2004) "lessons from information retrieval"
- away from lexical normalisation and towards relational simplification
 decreasing ontological expressiveness, epistemological commitment, and inferential power
- Shallow text operations (...) are right for information access. Information is
 primarily conveyed by natural language and this has to be shown to the user for
 them to assess
- and Wilks & Brewster (2009) state: The Semantic Web is nothing else other than scaling up natural language processing...

WATS:: Lexical statistics





Vagueness, ambiguity and multilingual issues		
	Diana Santo	DS
WATS: V	Nords and their secrets ESSLI	. J 2010
	for ab and men beereto, EbbEr	512010
Diana Sant	os & Maria José Bocorny Finatto	







Vagueness at all levels of description

- POS: the infamous case of past participles
- The case of *near*: adjective or preposition? (Manning & Schütze, 1999)
 The most famous case is, however, PP attachment. After discarding
- non V NP PP structures, Hindle and Rooth state:
 Disambiguating the test sample turned out to be a surprisingly difficult task. [...] more than 10% of the sentences seemed problematic to at least one author (Hindle & Rooth, 1993:112)

WATS:: Vagueness, ambiguity, and multilingual issues

Attempts to deal with vagueness

- In annotation, leave room for more than one category: HAREM and COMPARA
 do not force a choice when it is not required
- Identify contrastively vague categories in tense and aspect

not only coercion

 also aspectual classes or grammatical operators that can simultaneously mean more than one thing

The translation network

- linking two systems with different vague categories
- explaining and formalizing concrete translation issues

Slide 49 from Santos (2006)

WATS:: Vagueness, ambiguity, and multilingual issues

Again: How does a language choose its units?

- Talmy's (1983:277ff) suggestion:
- The majority of semantic domains in language are n-dimensional, with n a very large number. For example, no fewer that [] twenty parameters are relevant to the domain of spatial configuration as expressed by closed-class elements such as English prepositions and deictics. [List]
- With so many parameters, full domain coverage by fairly specific references would require thousands of distinct vocabulary items, [...]
- Rather that a contiguous array of specific references, languages instead exhibit a smaller number of such references in a scattered distribution over a semantic domain. That is, a fairly specific reference generally does not have any immediate neighbors of equal specificity.

WATS:: Vagueness, ambiguity, and multilingual issues

Natural Language

Technologies Group

Cont.

- General terms are necesssary for referring to insterstitial conceptual material, between the references of specific terms
- Their locations must nevertheless be to a great extent arbitrary, constrained primarily by the requirement of being "representative" of the lay of the semantic landscape, as evidenced by the enormous extent of non-correspondence between specific morphemes of different languages, even where these are spoken by the peoples of similar cultures.

WATS:: Vagueness, ambiguity, and multilingual issues







Institute of Language and Linguistics, Federal University of Rio Grande do Sul, Brazil

Fundação para a Computação Científica Nacional

Contrastive studies (according to Santos 1996)

- Universalism
- assume that differences are noise, and that they can be parametrized and done away at a deep enough level)
- Typology classify all languages on a number of axes, on the search of universal or frequent traits
- Relativism

take all languages as equals: the only unbiased way

WATS:: Vagueness, ambiguity, and multilingual issues

Contrastive studies (according to Pinker 2007)

Theories of language, in Pinker's (2007) words

- Extreme Nativism: born with 50,000 concepts (Fodor)
- Radical pragmatics: people can use a word to mean almost anything (Sperber and Wilson)
- Linguistic determinism: words determine thoughts (Sapir and Whorf) Pinker's moderate position ©: meanings of words are formulas in an abstract language of thought

WATS:: Vagueness, ambiguity, and multilingual issues

Contrastive studies (according to Chesterman 1998)

- Overview of the concept of equivalence in Translation Theory (pp.16-27) The equative view
- Signs represent meanings; meanings are absolute, unchanging, they are manifestations of the ideal, they are Platonic Ideas
- identity of meaning across translation The taxonomic view Different types of equivalence are argued to be appropriate in the translation of different kinds of texts
- Nida's formal equivalence vs dynamic equivalence The relativist view
- WATS:: Vagueness, ambiguity, and multilingual issues

Three ways of arriving at the relativist view

- From rational thinking: Logical rejection of sameness, replacing it by similarity, matching or family resemblance, or economical considerations
- equivalence depends only on what is offered, negotiated and accepted in the exchange situation (Pym, 1992/2010:46)
- From cognition: the interpretation of an utterance is a function of the utterance itself and the cognitive state of the interpreter: we interpret things in the light of what we already know ..
- From comparative literature and translation: TS is an empirical science whose aim is to determine the general laws of translation behaviour. Translations have many purposes and are of many kinds

WATS:: Vagueness, ambiguity, and multilingual issues





NATURE

descriptive verbs

descriptive verbs comprehend an activity nucleus (ANu) and a modificant (Mod) that can be expressed or rephrased by adjectives or manner adverbs, and often carries speaker's evaluation on some of the agents or on the action itself (Snell-Hornby, 1983)







Steiner's mystère supreme of anthropology

Why does *homo sapiens* whose digestive tract has evolved and functions in precisely the same complicated ways the world over, ... -why does this unified, though individually unique mammalian species does not use one common language? (Steiner, 1992 [1975] :52)

It is hardly to be found ONE distinction that is common across all

The comparison of languages is arguably the best mirror into language

translatable words (and not only) are a wonderful mirror to differences

Languages tend to evolve and age and innovate continuously

... and the comparison itself is best done through translation data

Words carve different domains in different languages, words are

different in different languages, the differences between inter-

in systematic organization of the languages (systematicy includes

WATS:: Vagueness, ambiguity, and multilingual issues

Concluding remarks

WATS:: Vagueness, ambiguity, and multilingual issues

natural languages

creativity)

To be or not to be: that's the question?

- It is remarkable how the verb to be is a complex problem for linguitsic description, and for translation, whose interpretation of this famous quote is difficult, to say the least
- The interpretation of be is an interesting chapter of natural language semantics. For the present purpose, it is enough to say that it ambiguously represents the operations of identity, membership and class inclusion. (Carlson, 1981: 156)
- the ambiguous noun time (Carlson, 1981:60) is translationally vindicated in Portuguese as follows: as a count noun, time is translated in Portuguese by vez ("turn"); as a mass noun, it represents the temporal domain (tempo). Cf. no. gang ("going"), fr. fois, it. volta ...

WATS:: Vagueness, ambiguity, and multilingual issues

SINTEF Natural Language Technologies Group



















