

My love affair with the Web

...and why it's over!

Marco Baroni

Center for Mind/Brain Sciences
University of Trento

BOTWU
Forlì, 24 June 2013

[Change photo](#)

Marco Baroni [Edit](#)

Tenured Researcher, Center for Mind/Brain Sciences, University of Trento (Italy) [Edit](#)

[Computational Linguistics](#) - [Linguistics](#) - [Semantics](#) [Edit](#)

Verified email at unitn.it [Edit](#)

My profile is private [Edit](#) [Add homepage](#)

Citation indices

	All	Since 2008
Citations	2185	1735
h-index	24	21
i10-index	45	39

Citations to my articles



Select: [All](#), [None](#) [Actions](#)

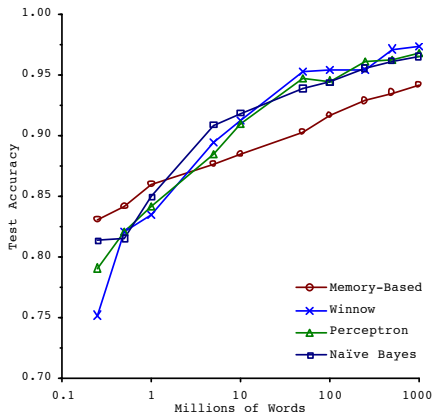
Show: [20](#) [1-20](#) [Next >](#)

Title / Author

Cited by Year

<input type="checkbox"/>	The WaCky wide web: a collection of very large linguistically processed web-crawled corpora M Baroni, S Bernardini, A Ferraresi, E Zanchetta Language Resources and Evaluation 43 (3), 209-226	227	2009
<input type="checkbox"/>	BootCaT: Bootstrapping corpora and terms from the web M Baroni, S Bernardini Proceedings of LREC 4, 1313-1316	217	2004
<input type="checkbox"/>	Large linguistically-processed web corpora for multiple languages M Baroni, A Kilgariff Proceedings of the Eleventh Conference of the European Chapter of the ...	113	2006
<input type="checkbox"/>	Introducing and evaluating ukwac, a very large web-derived corpus of english A Ferraresi, E Zanchetta, M Baroni, S Bernardini Proceedings of the 4th Web as Corpus Workshop (WAC-4) Can we beat Google, 47-54	95	2008

More data trump better algorithms



Michele Banko and Eric Brill, ACL 2001

Momentous events in the history of the Web as Corpus

- 1997 Radev and McKeown: Building a generation knowledge source using Internet-accessible newswire, ANLP
- 2002 Keller, Lapata and Ourioupina: Using the Web to overcome data sparseness, best paper award at EMNLP
- 2003 Special issue of Computational Linguistics on the Web as Corpus
- 2004 First version of BootCaT available! :-)
- 2006 WAC2, ACL Special Interest Group on Web-as-Corpus is born
- 2006 Google trillion-word Web 1T 5-Gram collection
- 2007 Googleology is bad science! (Adam Kilgarriff, Computational Linguistics)
- 2008 WaCky site up with billion-word Italian, German and English corpora freely available

A shocking revelation

BootCaT was invented in 2001 at CMU:

Rayid Ghani, Rosie Jones and Dunja Mladenić:
Mining the web to create minority language
corpora. **CIKM 2001**

Futurology

Web-as-Corpus 2014 as seen from 2004



- ▶ We've got the linguist search engine
- ▶ Data retrieval and pre-processing are huge topics in computational linguistics
- ▶ Cleaning Web data is a solved problem
- ▶ The Web is exactly the same as in 2004, just bigger

Flavio De Benedictis, Stefano Faralli and Roberto Navigli:
**GlossBoot: Bootstrapping Multilingual Domain Glossaries
from the Web**

We present GlossBoot, an effective minimally-supervised approach to acquiring wide-coverage domain glossaries for many languages. For each language of interest, given a small number of hypernymy relation seeds concerning a target domain, we bootstrap a glossary from the Web for that domain by means of iteratively acquired term/gloss extraction patterns. Our experiments show high performance in the acquisition of domain terminologies and glossaries for three different languages.

Cleaning the Web

The boilerplate challenge

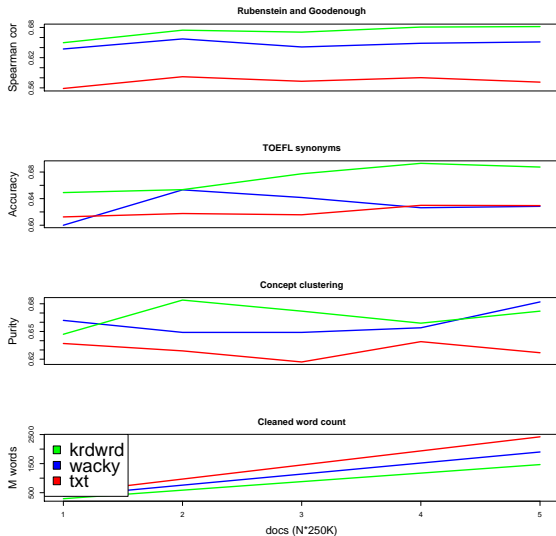


Cleaning, a solved problem?

	<i>features</i>	<i>precision</i>	<i>recall</i>	<i>F</i>
<i>wacky heuristic</i>	NA	80%	99%	88%
<i>text cues only</i>	21	92%	93%	92%
<i>DOM cues only</i>	13	89%	91%	90%
<i>visual cues only</i>	8	90%	93%	91%
<i>full krdwr</i>	42	93%	92%	92%

Does cleaning really help?

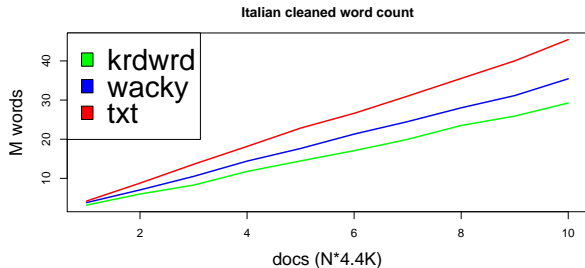
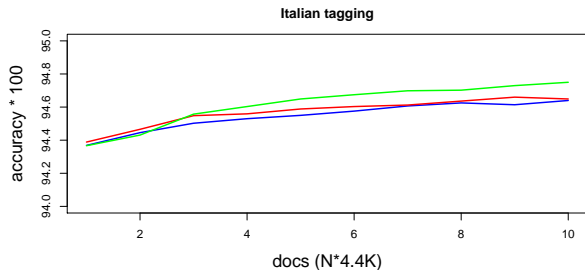
Various English semantic tasks



Collaboration with Egon Stemle

Does cleaning really help?

Italian Part-of-Speech tagging



The Web is no longer what it used to be!

[illegible]

And this is the end...

Thank you!