



# EVALITA 2009

## The Lexical Substitution Task

A. Toral  
Istituto di Linguistica Computazionale (CNR)





# Outline

- Introduction to the Lexical Substitution Task
- Evaluation
  - Dataset
  - Metrics
  - Participants
  - Results
  - Discussion
- Conclusion and future work



# Lexical Substitution

- Given a word in a context → find substitutes that fit (synonyms in context)
  - [...] nel tribunale di Pescara, riprendeva il **processo** contro [...] (*[...] in the court of Pescara, the **trial** against [...] resumed*)
    - Giudizio
  - [...] deve avere convinto i vertici della necessità di accelerare il **processo** di ristrutturazione (*[...] should have convinced the top brass of the need to accelerate the restructuring **process***)
    - Sviluppo, procedura



# Lexical Substitution

- Alternative to evaluate WSD
- Advantages
  - Absence of pre-defined inventory
    - Allows to evaluate unsupervised systems
    - WSD inventories too fine-grain
- First time for IT
  - EN in Semeval 2007



# Dataset

- Words to be substituted chosen from LRs (IWN, PSC) following diverse criteria to guarantee
  - High level of polysemy
  - Representativeness of the language
- Contexts: sentences from ISST where these words appear
- 210 words x 10 contexts (= 2,010) annotated with substitutes by 3 annotators
  - 300 trial
  - 1,710 test



# Annotation

- Tool developed for Semeval'07

logged as guest

## *LexSub* An interface for Lexical Substitution

Please replace the word in bold with a substitute which preserves the meaning of the sentence:

### Sentence #1926:

Alle nove , nel tribunale di Pescara , a poche centinaia di metri dal carcere , riprendeva il **processo** contro la banda Battestini e tutte le forze dell' ordine della citta' erano mobilitate attorno al palazzo di giustizia per garantirne \*NE la sicurezza .

Substitute:

OK

nil  extra responses  name  used a dictionary

Target word is part of  
phrase:

Comments:

**Reminder:** "You are free to consult a dictionary or thesaurus if it helps, but not another person. Please tick the dictionary box if you did consult a dictionary for any of the items for this word"

---

[< previous](#) | [next >](#) | [summaries](#) | [instructions](#) | [logout](#)

processo.n 1925 :: sviluppo 1;sistema 1;procedura 1;

processo.n 1926 :: giudizio 1;



# Metrics

- Two scoring types:
  - Best → best substitute
  - out-of-ten (oot) → best 10 substitutes
- Metrics
  - Precision (P)
  - Recall (R)
  - F-measure
  - ModeP, ModeR → like P, R but consider only substitute by majority of annotators



# Baseline

- No discrimination between contexts
- Semantic relations from LRs
  - Fetch all senses in LRs that correspond to word and extract synonyms and hyperonyms
  - Synonyms weight 3 times hyperonyms
  - If a word extracted more than once → sum weights
  - Lists of substitutes output ordered by weight
  - 3 runs:
    - PSC
    - IWN
    - PSC+IWN





# Participants

System	Techniques, tools	Lexicons	Corpora
Basile_a	WSD	ItalWordNet	
Basile_b	n-grams IR	ItalWordNet De Mauro	ItWaC
De Cao et al.	LSA SVD	ItalWordNet MultiWordNet	Repubblica
Baroni et al.	WSM PoS-tagger Lemmatiser	None	Repubblica ItWac Italian Wikipedia



# Results: best

Run	Precision	Recall	F	mode P	mode R
uniba2	<b>8.16</b>	<b>7.18</b>	<b>7.64</b>	10.58	10.58
baroniCutugnoLenciPucci	6.26	6.01	6.13	<b>11.28</b>	<b>10.84</b>
uniba1	6.80	5.53	6.10	8.90	8.90
uniba3	6.28	5.46	5.84	8.13	8.13
decao3	3.95	3.21	3.54	6.58	6.58
decao2	3.90	3.17	3.50	6.71	6.71
decao1	3.16	3.16	3.16	6.97	6.97
decao4	3.52	2.80	3.12	5.03	5.03
baseline_psc	<b>10.86</b>	<b>9.06</b>	<b>9.88</b>	<b>13.94</b>	<b>13.94</b>
baseline_iwn_psc	9.71	8.19	8.89	13.16	13.16
baseline_iwn	2.72	1.78	2.15	2.19	2.19



# Results: oot

Run	Precision	Recall	F	mode P	mode R
uniba2	<b>41.46</b>	<b>36.50</b>	<b>38.82</b>	<b>47.23</b>	<b>47.23</b>
uniba1	37.74	30.69	33.85	34.84	34.84
uniba3	28.54	24.79	26.53	34.58	34.58
decao3	23.48	19.11	21.07	26.58	26.58
decao2	23.00	18.72	20.64	26.32	26.32
decao1	20.09	20.09	20.09	27.74	27.74
decao4	18.62	14.78	16.48	20.52	20.52
baroniCutugnoLenciPucci	16.65	16.00	16.32	24.97	24.00
baseline_iwn_psc	<b>27.52</b>	<b>23.23</b>	<b>25.19</b>	<b>37.24</b>	<b>32.39</b>
baseline_psc	23.00	19.20	20.93	26.97	26.97
baseline_iwn	14.51	9.51	11.49	12.77	12.77



# Results: IT vs EN

---

## Scoring type Measure Italian English Difference

---

<b>best</b>	<b>P</b>	10.86	12.90	18.78
	<b>R</b>	9.06	12.90	42.38
	<b>mode P</b>	13.94	20.73	48.70
	<b>mode R</b>	13.94	20.73	48.70
<b>oot</b>	<b>P</b>	41.46	69.03	66.50
	<b>R</b>	36.50	68.90	88.80
	<b>mode P</b>	47.23	66.26	40.30
	<b>mode R</b>	47.23	66.26	40.30

---



# Conclusions, future

- Higher participation than WSD → higher interest?
- Different approaches presented
  - Complementary? Combine?
- Baselines: PSC much better than IWN
  - Some systems exploited IWN, replace by SIMPLE?
- Results far behind english → room for improvement



# End

## Thanks! Questions?

**EVALITA 2009**  
The Lexical Substitution Task

A. Toral  
Istituto di Linguistica Computazionale (CNR)

