

Onto.PT: Automatic construction of a Lexical Ontology for Portuguese

Hugo Gonalo Oliveira¹, Paulo Gomes

{hroliv,pgomes}@dei.uc.pt

Cognitive & Media Systems Group
CISUC, University of Coimbra

Lisbon, August 16, 2010

¹ supported by FCT scholarship grant SFRH/BD/44955/2008



- 1 Introduction
 - Lexical ontologies
 - Goals
- 2 Approach
 - Information extraction from text
 - Synset discovery
 - Merging synset-based resources
 - Weighting triples
 - Assigning terms to synsets
 - Assigning terms to synsets
 - Knowledge organisation
- 3 Current results
 - Relation extraction
 - Synset discovery
 - Wordnet establishment
- 4 Concluding remarks



Today's applications

- Need to understand information conveyed by natural language



Today's applications

- Need to understand information conveyed by natural language
- Therefore, demand better access to knowledge on words and their meanings!



Today's applications

- Need to understand information conveyed by natural language
- Therefore, demand better access to knowledge on words and their meanings!
- Encoded in lexical ontologies



Lexical ontologies

- Such as Princeton WordNet [Fellbaum, 1998]



Lexical ontologies

- Such as Princeton WordNet [Fellbaum, 1998]
 - ▶ Ontology + lexicon [Hirst, 2004]



Lexical ontologies

- Such as Princeton WordNet [Fellbaum, 1998]
 - ▶ Ontology + lexicon [Hirst, 2004]
 - ▶ Knowledge structured on words and their meanings



Lexical ontologies

- Such as Princeton WordNet [Fellbaum, 1998]
 - ▶ Ontology + lexicon [Hirst, 2004]
 - ▶ Knowledge structured on words and their meanings
 - ▶ Cover the whole language
 - ▶ Not based on a specific domain



Lexical ontologies

- Such as Princeton WordNet [Fellbaum, 1998]
 - ▶ Ontology + lexicon [Hirst, 2004]
 - ▶ Knowledge structured on words and their meanings
 - ▶ Cover the whole language
 - ▶ Not based on a specific domain
- Typically handcrafted...
 - ▶ Construction and maintenance involve time-consuming human effort!



Onto.PT

- Automatic construction of a lexical ontology for Portuguese



Onto.PT

- Automatic construction of a lexical ontology for Portuguese
- Extracted from different sources



Onto.PT

- Automatic construction of a lexical ontology for Portuguese
- Extracted from different sources
 - ▶ Manually created thesauri
 - ▶ Language dictionaries/encyclopedias
 - ▶ Corpora



Onto.PT

- Automatic construction of a lexical ontology for Portuguese
- Extracted from different sources
 - ▶ Manually created thesauri
 - ▶ Language dictionaries/encyclopedias
 - ▶ Corpora
- Modelled after Princeton WordNet



Onto.PT

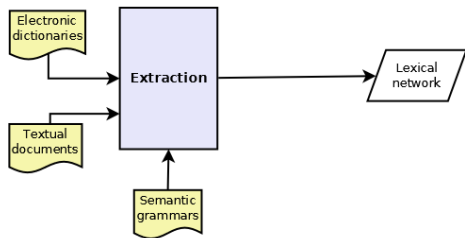
- Automatic construction of a lexical ontology for Portuguese
- Extracted from different sources
 - ▶ Manually created thesauri
 - ▶ Language dictionaries/encyclopedias
 - ▶ Corpora
- Modelled after Princeton WordNet
 - ▶ Synsets: groups of synonymous words
 - ▶ Synset-based relational triples



Onto.PT

- Automatic construction of a lexical ontology for Portuguese
- Extracted from different sources
 - ▶ Manually created thesauri
 - ▶ Language dictionaries/encyclopedias
 - ▶ Corpora
- Modelled after Princeton WordNet
 - ▶ Synsets: groups of synonymous words
 - ▶ Synset-based relational triples
- WSD based on the knowledge already extracted, not on the context





Examples

- From dictionaries:
 - ▶ `tenreiro, n -- terneiro, novilho ou bezerro.`
 - `terneiro` SYNONYM_OF `tenreiro`
 - `novilho` SYNONYM_OF `tenreiro`
 - `bezerro` SYNONYM_OF `tenreiro`



Examples

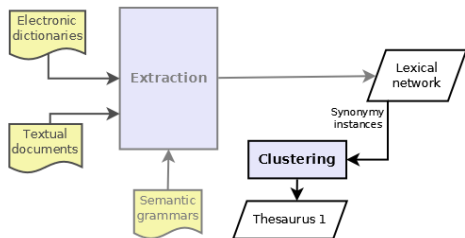
- From dictionaries:
 - ▶ *tenreiro*, n -- *terneiro*, *novilho* ou *bezerro*.
 - *terneiro* SYNONYM_OF *tenreiro*
 - *novilho* SYNONYM_OF *tenreiro*
 - *bezerro* SYNONYM_OF *tenreiro*
 - ▶ *ébola*, n -- virose que provoca febres e hemorragias
 - *ébola* CAUSATION_OF *febres*
 - *ébola* CAUSATION_OF *hemorragias*



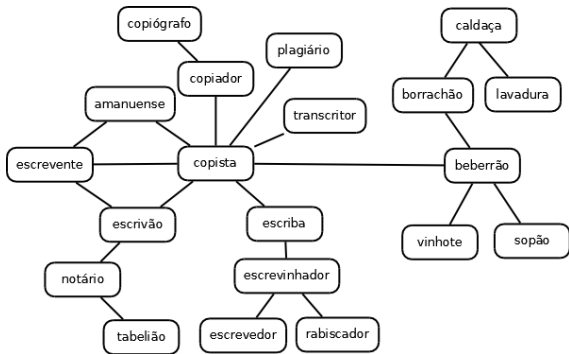
Examples

- From dictionaries:
 - ▶ *tenreiro*, n -- *terneiro*, *novilho* ou *bezerro*.
 - *terneiro* SYNONYM_OF *tenreiro*
 - *novilho* SYNONYM_OF *tenreiro*
 - *bezerro* SYNONYM_OF *tenreiro*
 - ▶ *ébola*, n -- *virose* que provoca *febres* e *hemorragias*
 - *ébola* CAUSATION_OF *febres*
 - *ébola* CAUSATION_OF *hemorragias*
- From textual corpora:
 - ▶ O *automobilismo* (também conhecido como *corridas de automóveis* ou *desporto motorizado*) é um *desporto*...
 - *automobilismo* SYNONYM_OF *corridas_de_automóveis*
 - *automobilismo* SYNONYM_OF *desporto_motorizado*
 - *desporto* HYPERNYM_OF *automobilismo*

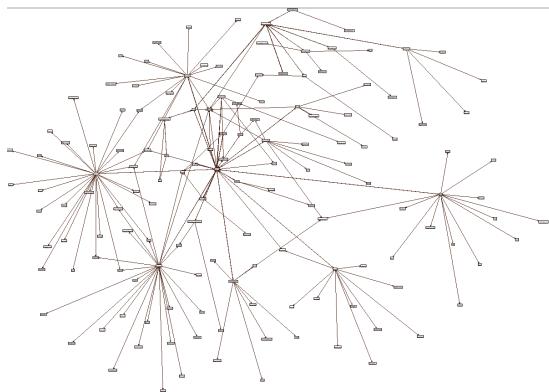




Synonymy lexical network – example

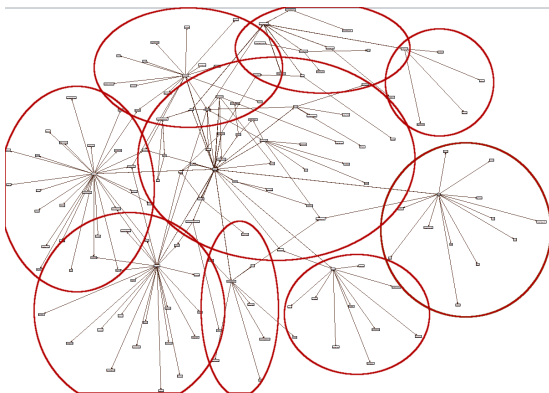


Synonymy networks tend to have a clustered structure



- Goal: Identify synsets taking advantage of clusters

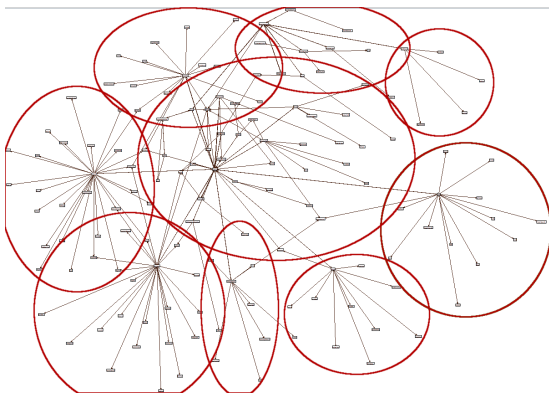
Synonymy networks tend to have a clustered structure



- Goal: Identify synsets taking advantage of clusters
- Approach: Clustering algorithm over the synonymy lexical network (see poster [Oliveira and Gomes, 2010])



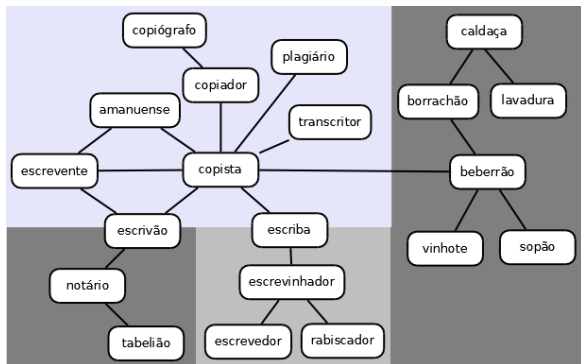
Synonymy networks tend to have a clustered structure

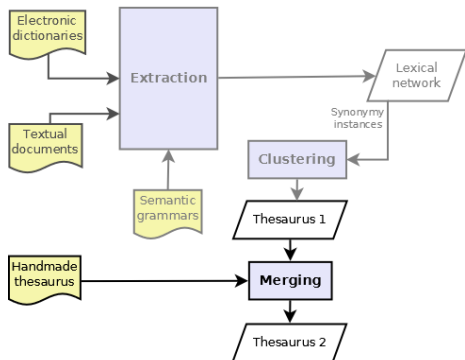


- Goal: Identify synsets taking advantage of clusters
- Approach: Clustering algorithm over the synonymy lexical network (see poster [Oliveira and Gomes, 2010])
- Keep ambiguity: clusters might be overlapping!



Clustering – example





Merging synsets from different thesauri

For each synset $T_i \in T$, select $B_j \in B$ with higher $c = |T_i \cap B_j| / |T_i \cup B_j|^2$

- $B_1 = (\textit{diva}, \textit{beldade}, \textit{beleza}, \textit{deidade}, \textit{deusa}, \textit{divindade})$
- $B_2 = (\textit{divindade}, \textit{deidade}, \textit{deus}, \textit{nume})$

²Jaccard coefficient



Merging synsets from different thesauri

For each synset $T_i \in T$, select $B_j \in B$ with higher $c = |T_i \cap B_j| / |T_i \cup B_j|^2$

- $B_1 = (\textit{diva}, \textit{beldade}, \textit{beleza}, \textit{deidade}, \textit{deusa}, \textit{divindade})$
- $B_2 = (\textit{divindade}, \textit{deidade}, \textit{deus}, \textit{nume})$
- $T_1 = (\textit{divindade}, \textit{diva}, \textit{deusa})$

²Jaccard coefficient

Merging synsets from different thesauri

For each synset $T_i \in T$, select $B_j \in B$ with higher $c = |T_i \cap B_j| / |T_i \cup B_j|^2$

- $B_1 = (\textit{diva}, \textit{beldade}, \textit{beleza}, \textit{deidade}, \textit{deusa}, \textit{divindade})$
- $B_2 = (\textit{divindade}, \textit{deidade}, \textit{deus}, \textit{nume})$
- $T_1 = (\textit{divindade}, \textit{diva}, \textit{deusa})$
 - ▶ $c(T_1, B_1) = \frac{1}{3}$
 - ▶ $c(T_1, B_2) = \frac{1}{6}$

²Jaccard coefficient



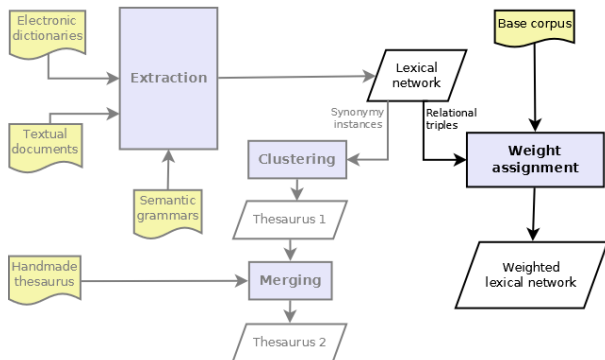
Merging synsets from different thesauri

For each synset $T_i \in T$, select $B_j \in B$ with higher $c = |T_i \cap B_j| / |T_i \cup B_j|^2$

- $B_1 = (\textit{diva}, \textit{beldade}, \textit{beleza}, \textit{deidade}, \textit{deusa}, \textit{divindade})$
- $B_2 = (\textit{divindade}, \textit{deidade}, \textit{deus}, \textit{nume})$
- $T_1 = (\textit{divindade}, \textit{diva}, \textit{deusa})$
 - ▶ $c(T_1, B_1) = \frac{1}{3}$
 - ▶ $c(T_1, B_2) = \frac{1}{6}$
- $N = B_1 \cup T_1 = (\textit{diva}, \textit{beldade}, \textit{beleza}, \textit{deidade}, \textit{deusa}, \textit{divindade})$

²Jaccard coefficient





Weighting triples

- Frequency of extraction



Weighting triples

- Frequency of extraction
- Corpus distributional similarity metrics (e.g. LSA [Deerwester et al., 1990], PMI [Turney, 2001])



Weighting triples

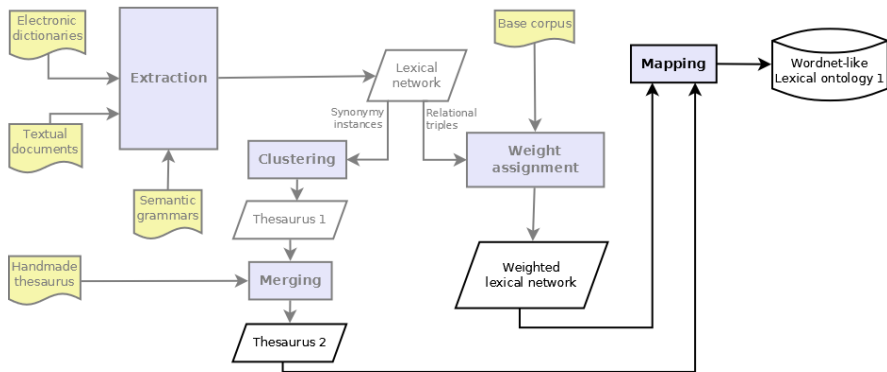
- Frequency of extraction
- Corpus distributional similarity metrics (e.g. LSA [Deerwester et al., 1990], PMI [Turney, 2001])
- Web distributional similarity metrics (e.g. WebJaccard, WebOverlap [Bollegala et al., 2007])



Weighting triples

- Frequency of extraction
- Corpus distributional similarity metrics (e.g. LSA [Deerwester et al., 1990], PMI [Turney, 2001])
- Web distributional similarity metrics (e.g. WebJaccard, WebOverlap [Bollegala et al., 2007])
- Define filters based on weights





Mapping methods

- Input:
 - ▶ Thesaurus T , containing synsets
 - ▶ Term-based semantic network, N , where each edge has a type R



Mapping methods

- Input:
 - ▶ Thesaurus T , containing synsets
 - ▶ Term-based semantic network, N , where each edge has a type R
- Goal: map $a R b \in N$ to $A R B, (A, B) \in T$



Mapping methods

- Input:
 - ▶ Thesaurus T , containing synsets
 - ▶ Term-based semantic network, N , where each edge has a type R
- Goal: map $a R b \in N$ to $A R B, (A, B) \in T$
- Output: semantic network W , whose nodes are synsets, which relate to other synsets by means of semantic relations (wordnet)



Mapping procedures

- Baseline
 - ▶ A and B are random synsets containing a and b respectively



Mapping procedures

- Baseline
 - ▶ A and B are random synsets containing a and b respectively
- Related proportion [Gonçalo Oliveira and Gomes, 2010]



Mapping procedures

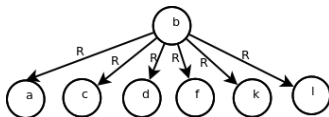
- Baseline
 - ▶ A and B are random synsets containing a and b respectively
- Related proportion [Gonçalo Oliveira and Gomes, 2010]
- Cosine similarity



Related proportion

Assignment of a (in $a R b$) to A :

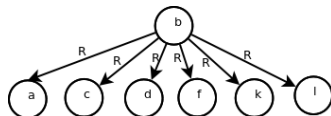
- 1 Fix b



Related proportion

Assignment of a (in $a R b$) to A :

1 Fix b

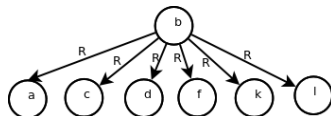


2 $S_a \subset T : S_{ai} \in S_a, a \in S_{ai}$

Related proportion

Assignment of a (in $a R b$) to A :

- 1 Fix b



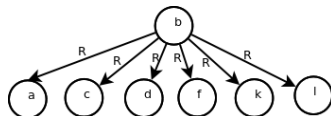
- 2 $S_a \subset T : S_{ai} \in S_a, a \in S_{ai}$

▶ a is not in T ? create synset $A = (a), a \rightarrow A$

Related proportion

Assignment of a (in $a R b$) to A :

- 1 Fix b

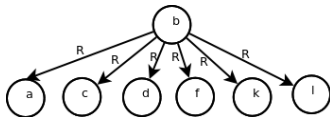


- 2 $S_a \subset T : S_{ai} \in S_a, a \in S_{ai}$
 - ▶ a is not in T ? create synset $A = (a), a \rightarrow A$
- 3 For each $S_{ai} \in S_a,$

Related proportion

Assignment of a (in $a R b$) to A :

- 1 Fix b

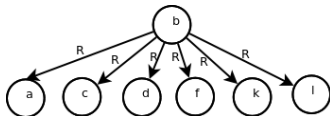


- 2 $S_a \subset T : S_{ai} \in S_a, a \in S_{ai}$
 - ▶ a is not in T ? create synset $A = (a), a \rightarrow A$
- 3 For each $S_{ai} \in S_a$,
 - ▶ $p_{ai} = \frac{n_{ai}}{|S_{ai}|}$, n_{ai} = number of terms $t_j \in S_{ai} : (t_j R b)$

Related proportion

Assignment of a (in $a R b$) to A :

- 1 Fix b



- 2 $S_a \subset T : S_{ai} \in S_a, a \in S_{ai}$

▶ a is not in T ? create synset $A = (a), a \rightarrow A$

- 3 For each $S_{ai} \in S_a$,

▶ $p_{ai} = \frac{n_{ai}}{|S_{ai}|}$, n_{ai} = number of terms $t_j \in S_{ai} : (t_j R b)$

★ $S_{a1} = (a, c, d, e), p_{a1} = \frac{3}{4}$

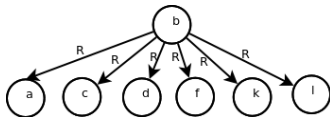
★ $S_{a2} = (a, f, g), p_{a2} = \frac{2}{3}$

★ $S_{a3} = (a, h, i, j), p_{a3} = \frac{1}{4}$

Related proportion

Assignment of a (in $a R b$) to A :

- 1 Fix b



- 2 $S_a \subset T : S_{ai} \in S_a, a \in S_{ai}$

▶ a is not in T ? create synset $A = (a), a \rightarrow A$

- 3 For each $S_{ai} \in S_a$,

▶ $p_{ai} = \frac{n_{ai}}{|S_{ai}|}$, n_{ai} = number of terms $t_j \in S_{ai} : (t_j R b)$

★ $S_{a1} = (a, c, d, e), p_{a1} = \frac{3}{4}$

★ $S_{a2} = (a, f, g), p_{a2} = \frac{2}{3}$

★ $S_{a3} = (a, h, i, j), p_{a3} = \frac{1}{4}$

▶ $a \rightarrow S_{a1}$



Cosine similarity

- 1 $M = \textit{term-term}$ matrix based on the adjacencies of the lexical network



Cosine similarity

- 1 $M = \textit{term-term}$ matrix based on the adjacencies of the lexical network
- 2 Collect all the synsets with a , $S_a \subset T$, and all synsets with b , $S_b \subset T$



Cosine similarity

- 1 $M = \textit{term-term}$ matrix based on the adjacencies of the lexical network
- 2 Collect all the synsets with a , $S_a \subset T$, and all synsets with b , $S_b \subset T$
- 3 For each $A \in S_a$ and $B \in S_b$, with terms $A_i \in A$ and $B_j \in B$:

$$\textit{sim}(A, B) = \frac{\sum_{i=1}^{|A|} \sum_{j=1}^{|B|} \cos(A_i, B_j)}{|A||B|}$$



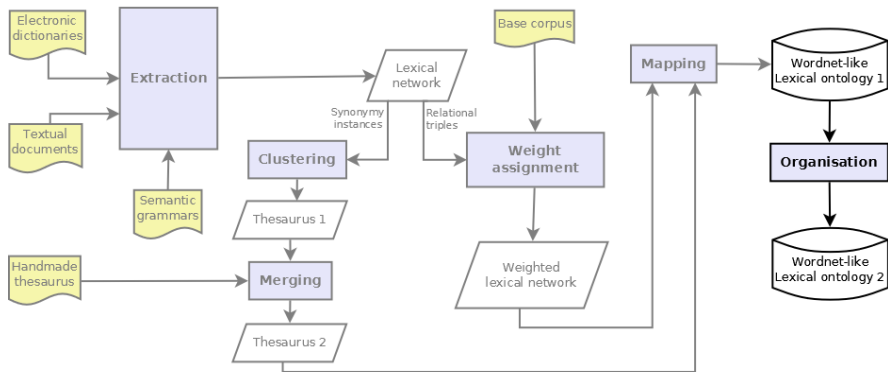
Cosine similarity

- 1 $M = \textit{term-term}$ matrix based on the adjacencies of the lexical network
- 2 Collect all the synsets with a , $S_a \subset T$, and all synsets with b , $S_b \subset T$
- 3 For each $A \in S_a$ and $B \in S_b$, with terms $A_i \in A$ and $B_j \in B$:

$$\textit{sim}(A, B) = \frac{\sum_{i=1}^{|A|} \sum_{j=1}^{|B|} \cos(A_i, B_j)}{|A||B|}$$

- 4 Select the pair of synsets with the highest similarity





Knowledge organisation

- Transitivity

- ▶ if R is transitive (e.g. SYNONYMY, HYPERNYMY, ...):
 $(A R B) \wedge (B R C) \rightarrow (A R C)$



Knowledge organisation

- Transitivity

- ▶ if R is transitive (e.g. SYNONYMY, HYPERNYMY, ...):
 $(A R B) \wedge (B R C) \rightarrow (A R C)$

- Inheritance

- ▶ if R is not a HYPERNYMY or HYPONYMY relation:
 $(A \text{ HYPERNYM_OF } B) \wedge (A R C) \rightarrow (B R C)$



Triples extracted from dictionaries

- Dicionário da Língua Portuguesa (PAPEL 2.0)
- Dicionário Aberto (DA)

Relation	Arguments	PAPEL 2.0	DA	Examples
Synonymy	noun,noun	37,452	20,910	<i>auxílio, contributo</i>
	verb,verb	21,465	8,715	<i>tributar, colectar</i>
	adj,adj	19,073	7,353	<i>flexível, moldável</i>
	adv,adv	1,171	605	<i>após, seguidamente</i>
Hypernymy	noun,noun	62,591	59,887	<i>planta, salva</i>
Part-of	noun,noun	2,805	1,795	<i>cauda, cometa</i>
	noun,adj	3,721	4,902	<i>tampa, coberto</i>
Member-of	noun,noun	5,929	1,564	<i>ervilha, Leguminosas</i>
	adj,noun	883	59	<i>celular, célula</i>
Causation	noun,noun	1,013	264	<i>fricção, assadura</i>
	adj,noun	498	166	<i>reactivo, reacção</i>
	verb,noun	6,399	5,714	<i>limpar, purgação</i>
Purpose	noun,noun	2,886	1,760	<i>defesa, armadura</i>
	verb,noun	5,192	3,383	<i>fazer_rir, comédia</i>
	verb,adj	260	186	<i>corrigir, correccional</i>

Table: Examples of triples



Relations extracted from Wikipedia abstracts

Relation	Quantity	Example	Sample	Correct	Agreement
Synonymy	11,862	<i>estupro, violação</i>	286	86,1%	91,2%
Hypernymy	29,563	<i>estilo_de_música, folk</i>	322	59,1% ³	93,1%
Part-of	1,287	<i>jejuno, intestino</i>	268	52,6%	78,4%
Causation	520	<i>parasita, doença</i>	244	49,6%	79,5%
Purpose	743	<i>construção, terracota</i>	264	57,0%	82,2%

Table: Examples and validation of relations

³In 30%, grammars/tagger could not identify complete scientific names as in:
O Iriatherina weneri é uma espécie de peixe de aquário
 → *peixe_de_aquário* HYPERNYM_OF *weneri*



Thesaurus

- TeP⁴ thesaurus
- OpenThesaurus.PT (OT)⁵
- Clustered PAPEL (CLIP)
- TeP merged with OT, merged with CLIP (TOP)

		TeP	OT	CLIP	TOP
Words	Quantity	17,158	5,819	23,741	30,554
	Ambiguous	5,867	442	12,196	13,294
	Most ambiguous	20	4	47	21
Synsets	Quantity	8,254	1,872	7,468	9,960
	Avg. size	3.51	3.37	12.57	6.6
	Biggest	21	14	103	277

Table: (Noun) thesauruses in numbers.

⁴<http://www.nilc.icmc.usp.br/tep2/index.htm>

⁵<http://openthesaurus.caixamagica.pt/>

Manual validation

	Sample	Correct	Incorrect	N/A	Agreement
CLIP	519 sets	65.8%	31.7%	2.5%	76.1%
CLIP'	310 sets	81.1%	16.9%	2.0%	84.2%
TOP	480 sets	83.2%	15.8%	1.0%	82.3%
TOP'	448 sets	86.8%	12.3%	0.9%	83.0%

Table: Results of manual synset validation.

- CLIP' and TOP' only consider synsets with 10 or less words.
 - ▶ The quality is higher for smaller synsets.



Resulting WordNet – related proportion

		Hypernym_of	Part_of	Member_of
Term-based triples		62,591	2,805	5,929
1st	Mapped	27,750	1,460	3,962
	Same synset	233	5	12
	Already present	3,970	40	167
Semi-mapped triples		7,952	262	357
2nd	Mapped	88	1	0
	Could be inferred	50	0	0
	Already present	13	0	0
Synset-based triples		23,572	1,416	3,783

Table: Results of triples mapping



Automatic validation

For each triple, $A R B$

- 1 Compile a set of textual patterns denoting R , e.g.:
 - ▶ (hypo) é um|uma (tipo|forma|variedade|...)* de (hyper)
 - ▶ (whole/group) é um (grupo|conjunto|...) de (part/member)



Automatic validation

For each triple, $A R B$

- 1 Compile a set of textual patterns denoting R , e.g.:
 - ▶ (hypo) é um|uma (tipo|forma|variedade|...)* de (hyper)
 - ▶ (whole/group) é um (grupo|conjunto|...) de (part/member)
- 2 Score the triple with the help of Google:

$$\text{score} = \frac{\sum_{i=1}^{|A|} \sum_{j=1}^{|B|} \text{found}(A_i, B_j, R)}{|A| * |B|}$$



Automatic validation

For each triple, $A R B$

- 1 Compile a set of textual patterns denoting R , e.g.:
 - ▶ (hypo) é um|uma (tipo|forma|variedade|...)* de (hyper)
 - ▶ (whole/group) é um (grupo|conjunto|...) de (part/member)
- 2 Score the triple with the help of Google:

$$\text{score} = \frac{\sum_{i=1}^{|A|} \sum_{j=1}^{|B|} \text{found}(A_i, B_j, R)}{|A| * |B|}$$

Relation	Sample size	Validation
Hypernymy_of	419 synsets	44,1%
Member_of	379 synsets	24,3%
Part_of	290 synsets	24,8%

Table: Automatic validation of triples



Concluding remarks

- Answer to the growing demand on semantically aware applications



Concluding remarks

- Answer to the growing demand on semantically aware applications
- Lack of public domain lexico-semantic resources for Portuguese



Concluding remarks

- Answer to the growing demand on semantically aware applications
- Lack of public domain lexico-semantic resources for Portuguese
- Export resources to different data formats



Concluding remarks

- Answer to the growing demand on semantically aware applications
- Lack of public domain lexico-semantic resources for Portuguese
- Export resources to different data formats
- WSD without a context:
 - ▶ Clustering for establishing synsets



Concluding remarks

- Answer to the growing demand on semantically aware applications
- Lack of public domain lexico-semantic resources for Portuguese
- Export resources to different data formats
- WSD without a context:
 - ▶ Clustering for establishing synsets
 - ▶ Disambiguation of terms based on the extracted knowledge



Concluding remarks

- Answer to the growing demand on semantically aware applications
- Lack of public domain lexico-semantic resources for Portuguese
- Export resources to different data formats
- WSD without a context:
 - ▶ Clustering for establishing synsets
 - ▶ Disambiguation of terms based on the extracted knowledge
 - ▶ Further organisation



Concluding remarks

- Answer to the growing demand on semantically aware applications
- Lack of public domain lexico-semantic resources for Portuguese
- Export resources to different data formats
- WSD without a context:
 - ▶ Clustering for establishing synsets
 - ▶ Disambiguation of terms based on the extracted knowledge
 - ▶ Further organisation
- Check <http://ontopt.dei.uc.pt> for updates and available resources



References



Bollegala, D., Matsuo, Y., and Ishizuka, M. (2007).

Measuring semantic similarity between words using web search engines.

In *Proc. 16th International conference on World Wide Web (WWW'07)*, pages 757–766, New York, NY, USA. ACM.



Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990).

Indexing by latent semantic analysis.

Journal of the American Society for Information Science, 41:391–407.



Fellbaum, C., editor (1998).

WordNet: An Electronic Lexical Database (Language, Speech, and Communication).

The MIT Press.



Gonçalo Oliveira, H. and Gomes, P. (2010).

Towards the automatic creation of a wordnet from a term-based lexical network.

In *Proceedings of the ACL Workshop TextGraphs-5: Graph-based Methods for Natural Language Processing*.



Hirst, G. (2004).

Ontology and the lexicon.

In Staab, S. and Studer, R., editors, *Handbook on Ontologies*, International Handbooks on Information Systems, pages 209–230. Springer.



Oliveira, H. G. and Gomes, P. (2010).

Automatic creation of a conceptual base for portuguese using clustering techniques.

In *Proc. 19th European Conference on Artificial Intelligence (ECAI 2010)*.



Turney, P. D. (2001).

Mining the web for synonyms: PMI-IR versus LSA on TOEFL.

In *Proc. 12th European Conference on Machine Learning (ECML-2001)*, volume 2167, pages 491–502. Springer.



Thank you!

