

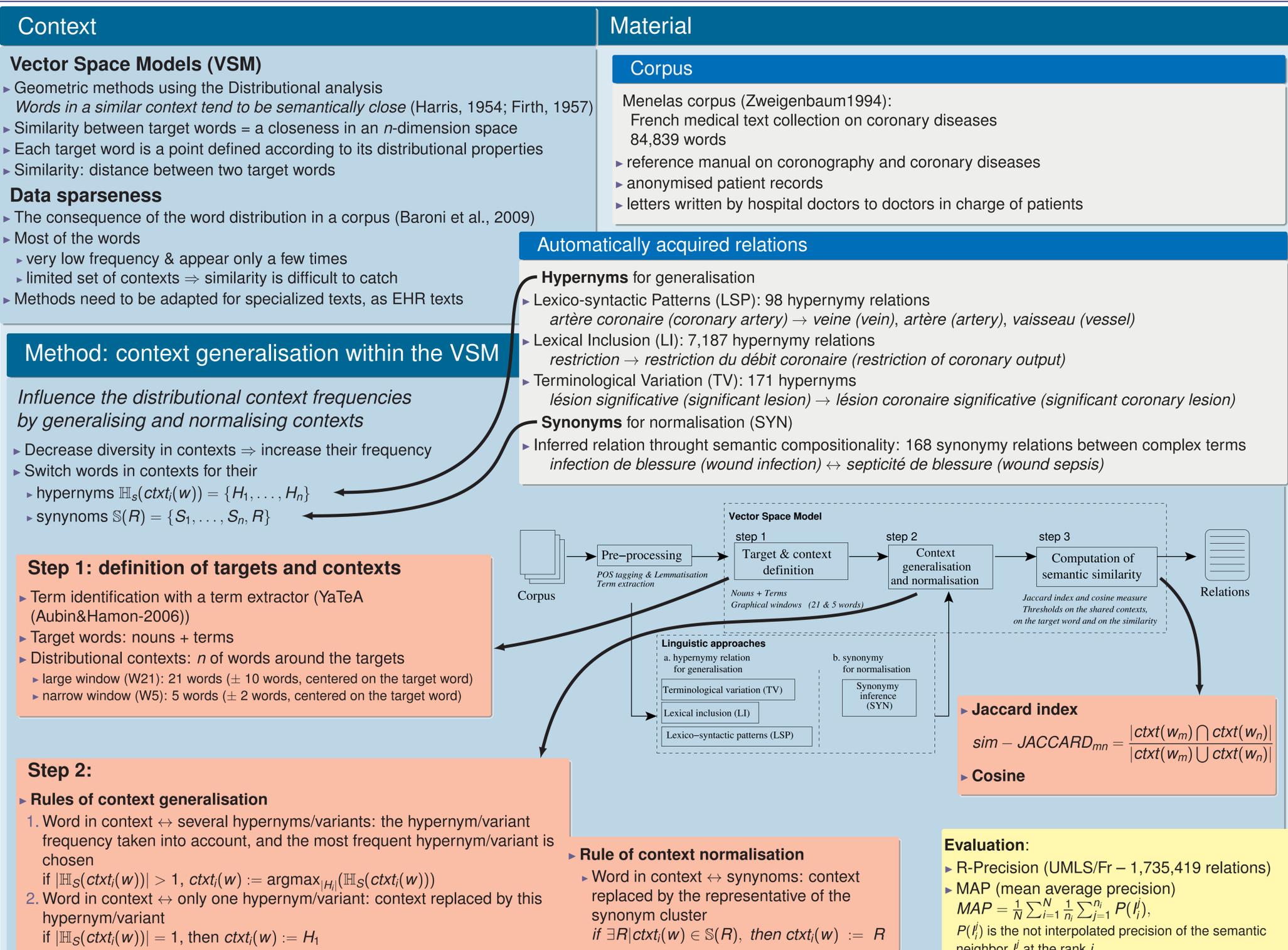
# Distributional context generalisation and normalisation as a mean to reduce data sparsity: evaluation of medical corpora

Amandine Périnet<sup>1</sup> and Thierry Hamon<sup>2,3</sup>

<sup>1</sup>LIMICS, INSERM, U1142, Université Paris 13, UPMC Univ Paris 06, Sorbonne Paris Cité, Villetaneuse, France

<sup>2</sup>LIMSI-CNRS, Orsay, France <sup>3</sup>Université Paris 13, Sorbonne Paris Cité, Villetaneuse, France

<sup>1</sup>amandine.perinet@edu.univ-paris13.fr, <sup>2</sup>hamon@limsi.fr



## Results

### Results obtained with a large window (21 words)

	Acquired Rel.		Rel. in UMLS		MAP		R-precision	
	JACC	COS	JACC	COS	JACC	COS	JACC	COS
<b>BASILINE</b>	406	428	4	44	0.406	0.188	0.250	0.118
<b>Generalisation</b>								
Variants	472	424	8	40	0.280	0.188	0.143	0.118
Lexical inclusion	328	223	4	26	0.454	0.110	0.250	0.133
Patterns	398	381	6	36	0.219	0.206	0.000	0.000
Variants + lex. inclusion	328	223	4	26	0.454	0.110	0.250	0.000
Patterns + lex. inclusion	338	220	4	26	0.454	0.101	0.250	0.000
3 generalisation methods	336	243	4	26	0.414	0.108	0.250	0.000
<b>Normalisation</b>								
Synonyms (Syn)	474	424	8	40	0.280	0.188	0.143	0.118
<b>Normalisation + generalisation</b>								
Syn/Variants	474	419	8	22	0.280	0.189	0.143	0.118
Syn/Lexical inclusion	366	279	6	14	0.440	0.105	0.333	0.000
Syn/Patterns	394	377	6	20	0.219	0.206	0.000	0.133
Syn/Variants + lex. inclusion	324	223	4	14	0.454	0.110	0.250	0.000
Syn/Patterns + lex. inclusion	394	373	6	20	0.219	0.206	0.000	0.133
Syn/3 generalisation methods	370	280	6	14	0.454	0.105	0.333	0.000

- Best results for a large window (21 words)
- Jaccard index:
  - Best MAP and R-precision
  - but few acquired relations found in UMLS
- Cosine: lower results
  - but more acquired relations found in UMLS

### General observations

- Context generalisation**
  - Quality improvement
  - Reduction of the number of relations acquired
  - Best MAP with
    - Jaccard and generalising with lexical inclusion
    - Cosine with pattern generalisation
  - Sequential generalisation: no influence of the order of the methods
  - Generalisation with several linguistic approaches does not improve the results
- Context normalisation**
  - Decrease of the results (Jaccard) or no effect (Cosine)
- Context normalisation and generalisation**
  - Jaccard:
    - Increase of the MAP and the R-precision when using all the hypernyms or variants and lexical inclusion
  - Cosine:
    - Improvement of the results when using lexical inclusion for generalising contexts

## Experiments

- Baseline:** VSM without any context generalisation (VSMOnly)
- Normalising distributional contexts (VSM/SYN)**
- Generalising distributional contexts**
  - Use of automatically acquired relations
  - 1 linguistic approach:** VSM/LSP, VSM/LI, VSM/TV
  - 2 approaches,** sequentially. eg: TV then LI (VSM/TV+LI)
  - 3 approaches** sequentially. eg: LSP then LI then TV (VSM/LSP+LI+TV) and all together (VSM/ALL3)
- Normalising then generalising distributional contexts (VSM/SYN/LSP, VSM/SYN/TV+LI, VSM/SYN/ALL3, ...)**

### Impact of the generalisation and the normalisation

- Positive impact:
  - Jaccard: when generalising the contexts with hypernyms issued from lexical inclusion
  - Cosine: when normalising and generalising with hypernyms provided by lexical inclusion
- No influence of context normalisation in most of the cases
- Reduction of the number of relations acquired

### Perspectives

- Manual analysis of the relations acquired and of the impact of the generalisation and normalisation process
- Comparison of our method with other dimension reduction methods as Random Indexing and LSA