

# Nominal Compound Compositionality: A Multilingual Lexicon and Predictive Model

Silvio Cordeiro<sup>1,2</sup>

[sr.cordeiro@inf.ufrgs.br](mailto:sr.cordeiro@inf.ufrgs.br)

Carlos Ramisch<sup>2</sup>

[carlos.ramisch@lif.univ-mrs.fr](mailto:carlos.ramisch@lif.univ-mrs.fr)

Aline Villavicencio<sup>1</sup>

[avillavicencio@inf.ufrgs.br](mailto:avillavicencio@inf.ufrgs.br)

<sup>1</sup> Federal University of Rio Grande do Sul (Brazil)

<sup>2</sup> Aix Marseille Université, CNRS, LIF UMR 7279 (France)

## Summary

- We collect human **compositionality judgments** for  $180 \times 3$  **nominal compounds** in English, French and Portuguese.
  - Example: insurance company > climate change > ... > milk tooth > ... > nut case > eager beaver > cloud nine.
- We consider a **predictive model** that uses **word embeddings** to predict the compositionality of these nominal compounds.
- We thoroughly **evaluate** several aspects of the model and overcome **state of the art** results in standard datasets.

## 1. Predictive model

### Hypothesis:

A compound  $w_1w_2$  is compositional  $\Leftrightarrow \overrightarrow{w_1w_2}$  is similar to  $\overrightarrow{w_1} + \overrightarrow{w_2}$

### Model:

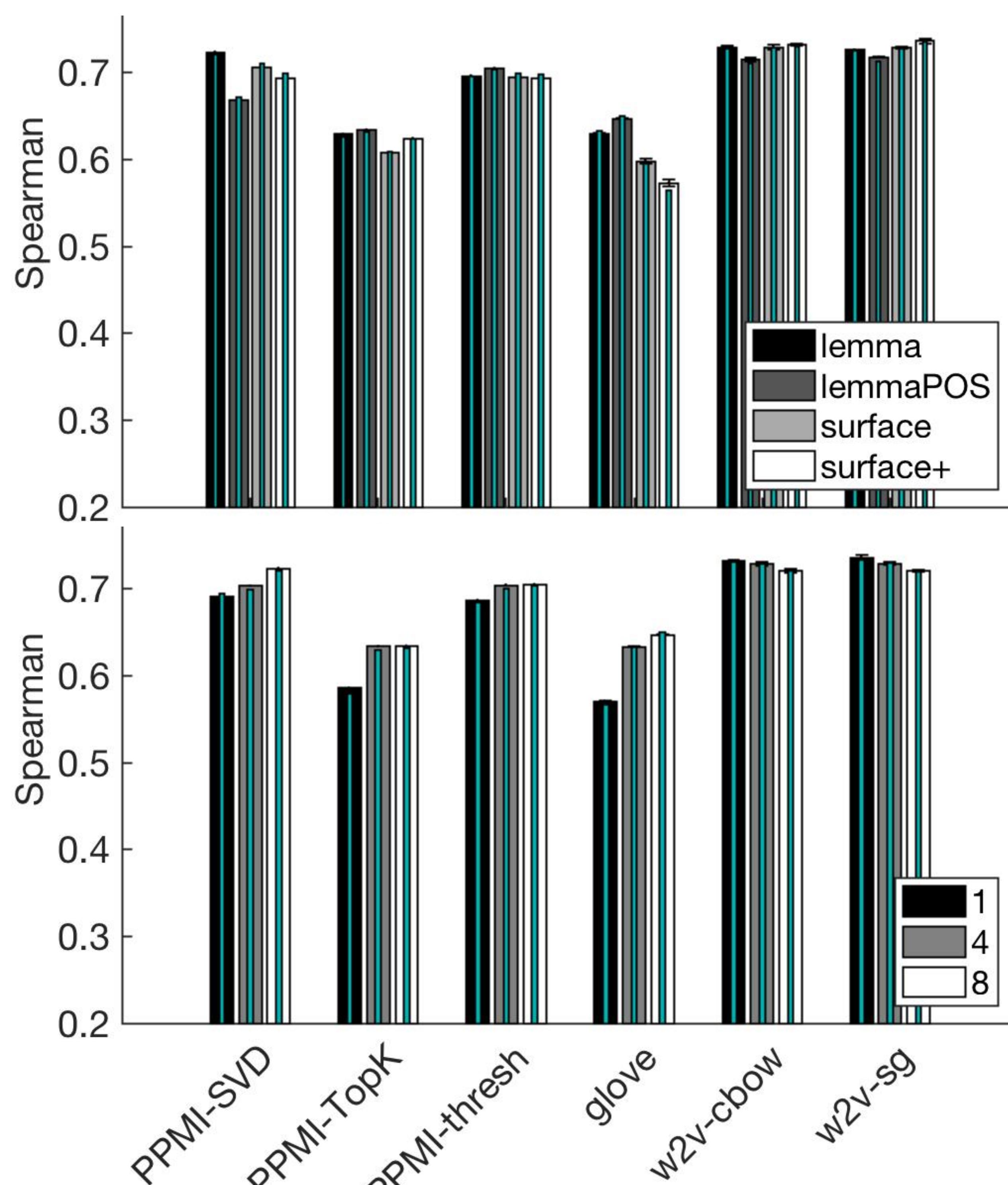
- Build word embeddings for  $w_1w_2$  and components  $w_1$  and  $w_2$
- Add vectors  $\overrightarrow{w_1} + \overrightarrow{w_2}$ , assuming the compound is compositional
- Predict compositionality score using cosine similarity
- Calculate Spearman rank correlation between prediction and human judgments

## 2. Evaluated parameters

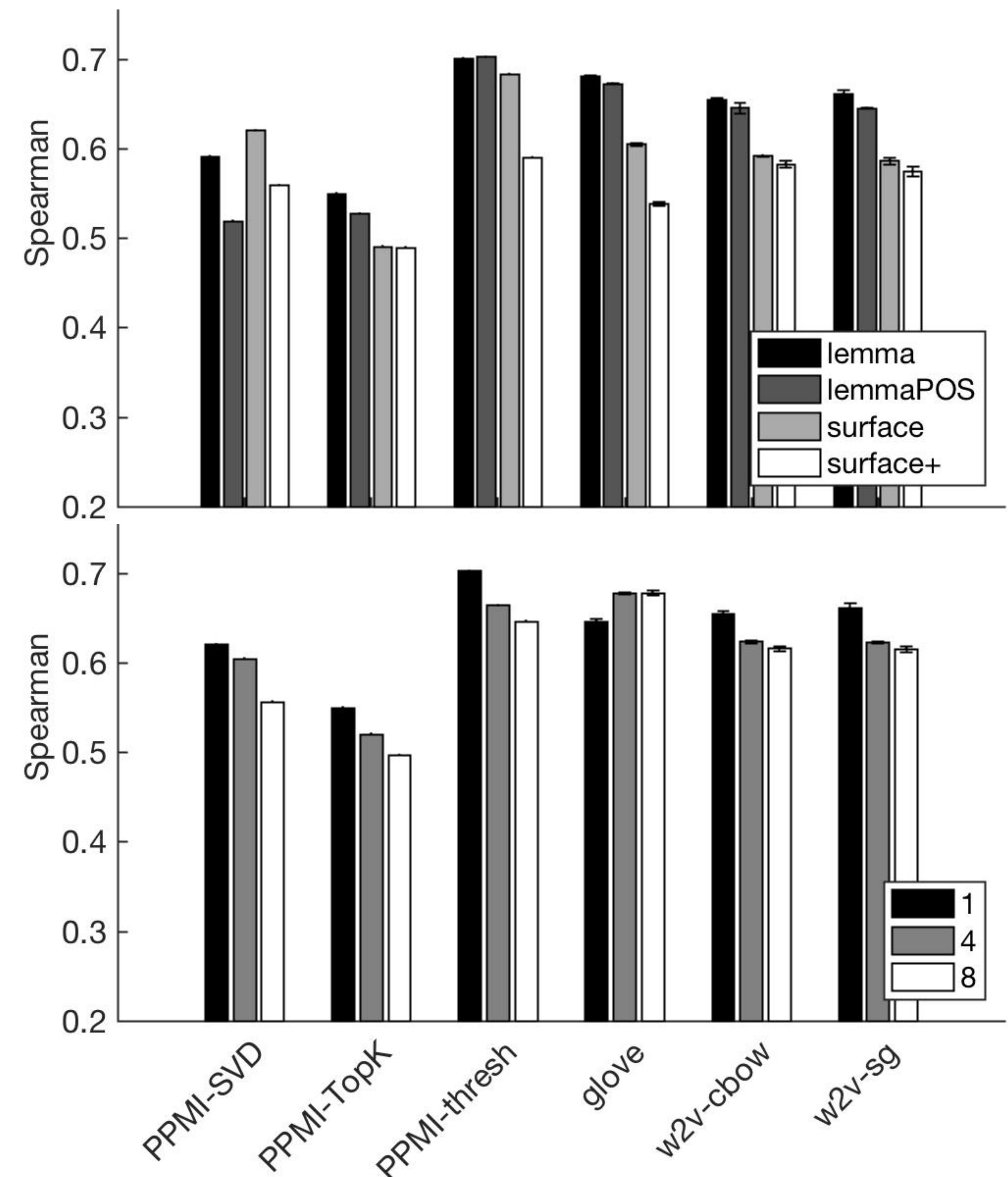
- Languages: English, French
- Models: PPMI matrix, GloVe, word2vec cbow & skipgram
- Preprocessing: surface+<sup>stopw</sup>, surface, lemmaPOS, lemma
- Context window sizes: 1, 4, 8 words to the left/right
- Dimensions: 250, 500, 750

## 3. Best correlations for English

- strict** evaluation (remove missing compounds)
- loose** evaluation (use fallback = average of predictions)



## 4. Best correlations for French



## 5. State of the art (English only)

| Dataset                                   | Model & Parameters | Spearman $\rho$ |
|---|--------------------|-----------------|
| Reddy et al [2011]                        |                    | .71             |
| Salehi et al [2015]                       |                    | .80             |
| Best w2v (sg, WF=surface, D=750, W=1)     |                    | .82 / .80       |
| Best PPMI(thresh, WF=surface, D=750, W=8) |                    | .80 / .80       |
| Dataset from Farahmand et al [2015]       |                    |                 |
| Yazdani et al [2015]                      |                    | .49             |
| Best w2v (sg, WF=lemma, D=500, W=1)       |                    | .51 / .47       |
| Best PPMI(svd, WF=lemma, D=750, W=4)      |                    | .52 / .45       |

## 6. Conclusions

- Lemmas better than surface forms for French
- Small windows better for French, not relevant for English
- More dimensions is consistently better
- Classical PPMI models are comparable to word2vec
- Two papers in ACL 2016: Ramisch et al & Cordeiro et al