

# Why So Down? The Role of Negative (and Positive) Pointwise Mutual Information in Distributional Semantics

Alexandre Salle<sup>1</sup> Aline Villavicencio<sup>1,2</sup>

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

<sup>2</sup>School of Computer Science and Electronic Engineering, University of Essex (UK)

alex@alexsalle.com avillavicencio@inf.ufrgs.br

## Abstract

In distributional semantics, the pointwise mutual information (*PMI*) weighting of the cooccurrence matrix performs far better than raw counts. There is, however, an issue with unobserved pair cooccurrences as *PMI* goes to negative infinity. This problem is aggravated by unreliable statistics from finite corpora which lead to a large number of such pairs. A common practice is to clip negative *PMI* ( $-PMI$ ) at 0, also known as Positive *PMI* (*PPMI*). In this paper, we investigate alternative ways of dealing with  $-PMI$  and, more importantly, study the role that negative information plays in the performance of a low-rank, weighted factorization of different *PMI* matrices. Using various semantic and syntactic tasks as probes into models which use either negative or positive *PMI* (or both), we find that most of the encoded semantics and syntax come from positive *PMI*, in contrast to  $-PMI$  which contributes almost exclusively syntactic information. Our findings deepen our understanding of distributional semantics, while also introducing novel *PMI* variants and grounding the popular *PPMI* measure.

## 1 Introduction

Dense word vectors (or embeddings) are a key component in modern NLP architectures for tasks such as sentiment analysis, parsing, and machine translation. These vectors can be learned by exploiting the distributional hypothesis (Harris, 1954), paraphrased by Firth (1957) as “*a word is characterized by the company that it keeps*”, usually by constructing a cooccurrence matrix over a training corpus, re-weighting it using Pointwise Mutual Information (*PMI*) (Church and Hanks, 1990), and performing a low-rank factorization to obtain dense vectors.

Unfortunately,  $PMI(w, c)$  goes to negative infinity when the word-context pair  $(w, c)$  does not

appear in the training corpus. Due to unreliable statistics, this happens very frequently in finite corpora. Many models work around this issue by clipping negative *PMI* values at 0, a measure known as Positive *PMI* (*PPMI*), which works very well in practice. An unanswered question is: “*What is lost/gained by collapsing the negative *PMI* spectrum to 0?*”. Understanding which type of information is captured by  $-PMI$  can help in tailoring models for optimal performance.

In this work, we attempt to answer this question by studying the kind of information contained in the negative and positive spectrums of *PMI* ( $-PMI$  and  $+PMI$ ). We evaluate weighted factorization of different matrices which use either  $-PMI$ ,  $+PMI$ , or both on various semantic and syntactic tasks. Results show that  $+PMI$  alone performs quite well on most tasks, capturing both semantics and syntax, in contrast to  $-PMI$ , which performs poorly on nearly all tasks, except those that test for syntax. Our main contribution is deepening our understanding of distributional semantics by extending Firth (1957)’s paraphrase of the distributional hypothesis to “*a word is not only characterized by the company that it keeps, but also by the company it rejects*”. Our secondary contributions are the proposal of two *PMI* variants that account for the spectrum of  $-PMI$ , and the justification of the popular *PPMI* measure.

In this paper, we first look at related work (§2), then study  $-PMI$  and ways of accounting for it (§3), describe experiments (§4), analyze results (§5), and close with ideas for future work (§6).

## 2 Related Work

There is a long history of studying weightings (also known as association measures) of general (not only word-context) cooccurrence matrices; see Manning et al. (1999); Jurafsky (2000) for an

overview and Curran and Moens (2002) for comparison of different weightings. Bullinaria and Levy (2007) show that word vectors derived from *PPMI* matrices perform better than alternative weightings for word-context cooccurrence. In the field of collocation extraction, Bouma (2009) address the negative infinity issue with *PMI* by introducing the normalized *PMI* metric. Levy and Goldberg (2014) show theoretically that the popular Skip-gram model (Mikolov et al., 2013) performs implicit factorization of shifted *PMI*.

Recently, work in explicit low-rank matrix factorization of *PMI* variants has achieved state of the art results in word embedding. GloVe (Pennington et al., 2014) performs weighted factorization of the log cooccurrence matrix with added bias terms, but does not account for zero cells. Shazeer et al. (2016) point out that GloVe’s bias terms correlate strongly with unigram log counts, suggesting that GloVe is factorizing a variant of *PMI*. Their SwiVel model modifies the GloVe objective to use Laplace smoothing and hinge loss for zero counts of the cooccurrence matrix, directly factorizing the *PMI* matrix, sidestepping the negative infinity issue. An alternative is to use *PPMI* and variants as in Kiela and Clark (2014); Polajnar and Clark (2014); Milajevs et al. (2016); Salle et al. (2016); Xin et al. (2018). However, it is not clear what is lost by clipping the negative spectrum of *PMI*, which makes the use of *PPMI*, though it works well in practice, seem unprincipled.

In the study of language acquisition, Regier and Gahl (2004) argue that indirect negative evidence might play an important role in human acquisition of grammar, but do not link this idea to distributional semantics.

### 3 PMI & Matrix Factorization

**PMI:** A cooccurrence matrix  $M$  is constructed by sliding a symmetric window over the subsampled (Mikolov et al., 2013) training corpus and for each center word  $w$  and context word  $c$  within the window, incrementing  $M_{wc}$ . *PMI* is then equal to:

$$PMI(w, c) = \log \frac{M_{wc} M_{**}}{M_{w*} M_{*c}} \quad (1)$$

where  $*$  denotes summation over the corresponding index. To deal with negative values, we propose clipped *PMI*,

$$CPMI_z(w, c) = \max(z, PMI(w, c)) \quad (2)$$

which is equivalent to *PPMI* when  $z = 0$ .

**Matrix factorization:** LexVec (Salle et al., 2016) performs the factorization  $M' = WC^\top$ , where  $M'$  is any transformation of  $M$  (such as *PPMI*), and  $W, C$  are the word and context embeddings respectively. By sliding a symmetric window over the training corpus (window sampling), LexVec performs one Stochastic Gradient Descent (SGD) step every time a  $(w, c)$  pair is observed, minimizing

$$L_{wc} = \frac{1}{2} (W_w C_c^\top - M'_{wc})^2$$

Additionally, for every center word  $w$ ,  $k$  negative words (Mikolov et al., 2013) are drawn from the unigram context distribution  $P_n$  (negative sampling) and SGD steps taken to minimize:

$$L_w = \frac{1}{2} \sum_{i=1}^k \mathbf{E}_{c_i \sim P_n(c)} (W_w C_{c_i}^\top - M'_{wc_i})^2$$

Thus the loss function prioritizes the correct approximation of frequently cooccurring pairs and of pairs where either word occurs with high frequency; these are pairs for which we have more reliable statistics.

In our experiments, we use LexVec over Singular Value Decomposition (SVD) because a) Empirical results shows it outperforms SVD (Salle et al., 2016). b) The weighting of reconstruction errors by statistical confidence is particularly important for  $-PMI$ , where negative cooccurrence between a pair of frequent words is more significant and should be better approximated than that between a pair of rare words. GloVe’s matrix factorization is even more unsuitable for our experiments as its loss weighting — a monotonically increasing function of  $M_{wc}$  — ignores reconstruction errors of non-cooccurring pairs.

**Spectrum of PMI:** To better understand the distribution of *CPMI* values, we plot a histogram of  $10^5$  pairs randomly sampled by window sampling and negative sampling in fig. 1, setting  $z = -5$ . We can clearly see the spectrum of  $-PMI$  that is collapsed when we use *PPMI* ( $z = 0$ ). In practice we find that  $z = -2$  captures most of the negative spectrum and consistently gives better results than smaller values so we use this value for the rest of this paper. We suspect this is due to the large number of non-cooccurring pairs (41.7% in this sample) which end up dominating the loss function when  $z$  is too small.

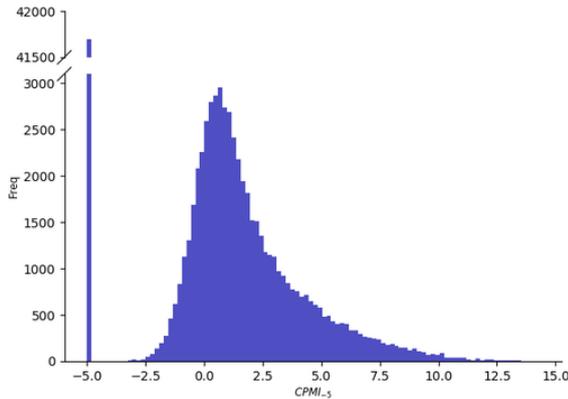


Figure 1:  $CPMI_{-5}$  histogram (bucket width equal to .2) of  $10^5$  sampled pairs using window sampling and negative sampling. Number of samples in interval:  $[-5, -5] = 41695$ ,  $(-5, 0] = 11001$ ,  $[-2, 0] = 10759$ ,  $(0, \infty) = 47304$

**Normalization:** We also experiment with normalized PMI ( $NPMI$ ) (Bouma, 2009):

$$NPMI(w, c) = PMI(w, c) / -\log(M_{wc}/M_{**})$$

such that  $NPMI(w, c) = -1$  when  $(w, c)$  never cooccur,  $NPMI(w, c) = 0$  when they are independent, and  $NPMI(w, c) = 1$  when they always cooccur together. This effectively captures the entire negative spectrum, but has the downside of normalization which discards scale information. In practice we find this works poorly if done symmetrically, so we introduce a variant called  $NNEGPMI$  which only normalizes  $-PMI$ :

$$NNEGPMI(w, c) = \begin{cases} NPMI(w, c) & \text{if } PMI(w, c) < 0 \\ PMI(w, c) & \text{otherwise} \end{cases}$$

We also experimented with Laplace smoothing as in Turney and Littman (2003) for various pseudocounts but found it to work consistently worse than both  $CPMI_z$  and  $NNEGPMI$  so we omit further discussion in this paper.

## 4 Materials

In order to identify the role that  $-PMI$  and  $+PMI$  play in distributional semantics, we train LexVec models that skip SGD steps when target cell values are  $> 0$  or  $\leq 0$ , respectively. For example,  $-CPMI_{-2}$  skips steps when  $CPMI_{-2}(w, c) > 0$ . Similarly, the  $+PPMI$  model skips SGD steps when  $PPMI(w, c) \leq 0$ . We compare these to

models that include both negative and positive information to see how the two interact.

We use the default LexVec configuration for all  $PMI$  variants: fixed window of size 2, embedding dimension of 300, 5 negative samples, positional contexts<sup>1</sup>, context distribution smoothing of .75, learning rate of .025, no subword information, and negative distribution power of .75. We train on a lowercased, alphanumerical 2015 Wikipedia dump with 3.8B tokens, discarding tokens with frequency  $< 100$ , for a vocabulary size of 303,517 words.

For comparison, we include results for a randomly initialized, non-trained embedding to establish task baselines.

**Semantics:** To evaluate word-level semantics, we use the SimLex (Hill et al., 2015) and Rare Word (RW) (Luong et al., 2013) word similarity datasets, and the Google Semantic (GSem) analogies (Mikolov et al., 2013). We evaluate sentence-level semantics using averaged bag of vectors (BoV) representations on the Semantic Textual Similarity (STSb) task (Cer et al., 2017) and Word Content<sup>2</sup> (WC) probing task (identify from a list of words which is contained in the sentence representation) from SentEval (Conneau et al., 2018).

**Syntax:** Similarly, we use the Google Syntactic analogies<sup>3</sup> (GSyn) (Mikolov et al., 2013) to evaluate word-level syntactic information, and Depth (Dep) and Top Constituent (TopC) (of the input sentence’s constituent parse tree) probing tasks from SentEval (Conneau et al., 2018) for sentence-level syntax. Classifiers for all SentEval probing tasks are multilayer perceptrons with a single hidden layer of 100 units and dropout of .1. Our final syntactic task is part-of-speech (POS) tagging using the same BiLSTM-CRF<sup>4</sup> setup as Huang et al. (2015) but using only word embeddings (no hand-engineered features) as input, trained on the WSJ section of the Penn Treebank (Marcus et al., 1993).

## 5 Results

All results are shown in table 1.

<sup>1</sup>Positional contexts account for the position of a context word relative to the target word— e.g.,  $M_{wc_{-1}}$  is the number of occurrences of  $c$  immediately to the left of  $w$ .

<sup>2</sup>By construction most probe words are content words, thus recovery relies on semantic information.

<sup>3</sup>Google Syntactic analogies are in fact morphological but many categories test for POS relations and are therefore syntactic in nature.

<sup>4</sup>Using <https://github.com/zalando-research/flair>

Model	Semantic					Syntactic			
	SimLex	RW	GSem	STSB	WC	GSyn	POS	Dep	TopC
+ <i>PPMI</i>	<b>.377</b>	.352	56.1	<u>.622</u>	<b>74.1</b>	50.3	92.2	30.5	<u>34.6</u>
- <i>CPMI</i> <sub>-2</sub>	.164	.231	3.6	.402	22.7	7.1	89.6	<b>32.7</b>	<b>34.7</b>
- <i>NNEGPMI</i>	.142	.232	3.3	.366	16.6	6.3	88.8	<u>32.4</u>	34.1
<i>PPMI</i>	<u>.363</u>	<b>.459</b>	<u>80.3</u>	.618	69.6	62.2	<u>92.5</u>	29.0	30.5
<i>CPMI</i> <sub>-2</sub>	.355	.432	<u>80.3</u>	.621	69.9	<b>65.1</b>	<b>92.6</b>	28.5	31.1
<i>NPMI</i>	.322	.437	63.5	.578	58.0	58.0	92.1	29.2	31.4
<i>NNEGPMI</i>	.360	<u>.439</u>	<b>80.7</b>	<b>.629</b>	70.0	<u>64.2</u>	<b>92.6</b>	27.2	30.3
Random	-.018	-.026	0.0	.453	0.3	0.0	55.2	17.9	5.0

Table 1: SimLex and RW word similarity: Spearman rank correlation. STSB: Pearson correlation. GSem/GSyn word analogy, POS tagging and WC, Dep, TopC probing tasks: % accuracy. Best result for each column in bold, second best underlined.

**Negative PMI:** We observe that using only  $-PMI$  (rows  $-CPMI_{-2}$  and  $-NNEGPMI$ ) performs similarly to all other models in POS tagging and both syntactic probing tasks, but very poorly on all semantic tasks, strongly supporting our main claim that  $-PMI$  mostly encodes syntactic information.

Our hypothesis for this is that the grammar that generates language implicitly creates negative cooccurrence and so  $-PMI$  encodes this syntactic information. Interestingly, this idea creates a bridge between distributional semantics and the argument by Regier and Gahl (2004) that indirect negative evidence might play an important role in human language acquisition of grammar.

**Positive PMI:** The  $+PPMI$  model performs as well or better as the full spectrum models on nearly all tasks, clearly indicating that  $+PMI$  encodes both semantic and syntactic information.

**Why incorporate -PMI?**  $+PPMI$  only falters on the RW and analogy tasks, and we hypothesize this is where  $-PMI$  is useful: in the absence of positive information, negative information can be used to improve rare word representations and word analogies. Analogies are solved using nearest neighbor lookups in the vector space, and so accounting for negative cooccurrence effectively repels words with which no positive cooccurrence was observed. In future work, we will explore incorporating  $-PMI$  only for rare words (where it is most needed).

**Full spectrum models:** The  $PPMI$ ,  $CPMI_{-2}$ , and  $NNEGPMI$  models perform similarly, whereas the  $NPMI$  model is significantly worst on nearly all semantic tasks. We thus conclude that accounting for scale in the positive spectrum

is more important than in the negative spectrum. We hypothesize this is because scale helps to uniquely identify words, which is critical for semantics (results on  $WC$  task correlate strongly with performance on semantic tasks), but in syntax, words with the same function should be indistinguishable. Since  $+PMI$  encodes both semantics and syntax, scale must be preserved, whereas  $-PMI$  encodes mostly syntax, and so scale information can be discarded.

**Collapsing the negative spectrum:** The  $PPMI$  model, which collapses the negative spectrum to zero, performs almost identically to the  $CPMI_{-2}$  and  $NNEGPMI$  models that account for the range of negative values. This is justified by 1) Our discussion which shows that  $+PMI$  is far more informative than  $-PMI$  and 2) Looking at fig. 1, we see that collapsed values — interval  $(-5, 0]$  — account for only 11% of samples compared to 41.7% for non-collapsed negative values.

## 6 Conclusions and Future Work

In this paper, we evaluated existing and novel ways of incorporating  $-PMI$  into word embedding models based on explicit weighted matrix factorization<sup>5</sup>, and, more importantly, studied the role that  $-PMI$  and  $+PMI$  each play in distributional semantics, finding that “*a word is not only characterized by the company that it keeps, but also by the company it rejects*”. In future work, we wish to further study the link between our work and language acquisition, and explore the fact the  $-PMI$  is almost purely syntactic to (possibly) subtract syntax from the full spectrum models, study-

<sup>5</sup>Code available at <https://github.com/alexandres/lexvec>

ing the frontier (if there is one) between semantics and syntax.

## Acknowledgments

This research was partly supported by CAPES and CNPq (projects 312114/2015-0, 423843/2016-8, and 140402/2018-7).

## References

- Gerlof Bouma. 2009. Normalized (pointwise) mutual information in collocation extraction. *Proceedings of GSCL*, pages 31–40.
- John A Bullinaria and Joseph P Levy. 2007. Extracting semantic representations from word co-occurrence statistics: A computational study. *Behavior research methods*, 39(3):510–526.
- Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14.
- Kenneth Ward Church and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational linguistics*, 16(1):22–29.
- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single vector: Probing sentence embeddings for linguistic properties. *arXiv preprint arXiv:1805.01070*.
- James R. Curran and Marc Moens. 2002. Improvements in automatic thesaurus extraction. In *Proceedings of the ACL-02 Workshop on Unsupervised Lexical Acquisition - Volume 9*, ULA '02, pages 59–66, Stroudsburg, PA, USA. Association for Computational Linguistics.
- John R Firth. 1957. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*.
- Zellig S Harris. 1954. Distributional structure. *Word*.
- Felix Hill, Roi Reichart, and Anna Korhonen. 2015. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. *CoRR*, abs/1508.01991.
- Dan Jurafsky. 2000. *Speech & language processing*. Pearson Education India.
- Douwe Kiela and Stephen Clark. 2014. A systematic study of semantic vector space model parameters. In *Proceedings of the 2nd Workshop on Continuous Vector Space Models and their Compositionality (CVSC)*, pages 21–30, Gothenburg, Sweden. Association for Computational Linguistics.
- Omer Levy and Yoav Goldberg. 2014. Neural word embedding as implicit matrix factorization. In *Advances in Neural Information Processing Systems*, pages 2177–2185.
- Minh-Thang Luong, Richard Socher, and Christopher D Manning. 2013. Better word representations with recursive neural networks for morphology. *CoNLL-2013*, 104.
- Christopher D Manning, Christopher D Manning, and Hinrich Schütze. 1999. *Foundations of statistical natural language processing*. MIT press.
- Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of english: The penn treebank. *Comput. Linguist.*, 19(2):313–330.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Dmitrijs Milajevs, Mehrnoosh Sadrzadeh, and Matthew Purver. 2016. Robust co-occurrence quantification for lexical distributional semantics. In *Proceedings of the ACL 2016 Student Research Workshop*, pages 58–64, Berlin, Germany. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12.
- Tamara Polajnar and Stephen Clark. 2014. Improving distributional semantic vectors through context selection and normalisation. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 230–238, Gothenburg, Sweden. Association for Computational Linguistics.
- Terry Regier and Susanne Gahl. 2004. Learning the unlearnable: The role of missing evidence. *Cognition*, 93(2):147–155.
- Alexandre Salle, Aline Villavicencio, and Marco Idiart. 2016. Matrix factorization using window sampling and negative sampling for improved word representations. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 2: Short Papers*. The Association for Computer Linguistics.
- Noam Shazeer, Ryan Doherty, Colin Evans, and Chris Waterson. 2016. Swivel: Improving embeddings by noticing what’s missing. *arXiv preprint arXiv:1602.02215*.

Peter D Turney and Michael L Littman. 2003. Measuring praise and criticism: Inference of semantic orientation from association. *ACM Transactions on Information Systems (TOIS)*, 21(4):315–346.

Xin Xin, Yuan Fajie, and He Xiangnan. 2018. Batch is not heavy: Learning word embeddings from all samples. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*.