GloVe: Global Vectors (Pennington, Socher and Manning 2014)

Erhard Hinrichs

Seminar für Sprachwissenschaft Eberhard-Karls Universität Tübingen

Bibliographical Reference

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation. Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), 1532–1543.

Previous Approaches for Learning Word Vectors

Global matrix factorization methods such as Latent Semantic Analysis (LSA; Deerwester et al. 1990)

- widely used in Information Retrieval (IR) and based term-document matrices
- uses Singular Value Decomposition (SVD) to rerank the dimensions of a matrix from most to least informative
- ► LSA, practicioners assume that only the top 300 or so dimensions (out of tens or even hundreds of thousands) are useful for capturing the meaning of texts.
- downsides of LSA:
 - not suitable for very large corpora
 - does not adequately capture the substructure of the vector space and thus does poorly on analogy tasks.

Previous Approaches for Learning Word Vectors

Local context window methods such as the Skipgram Model

- uses a sliding window of local contexts over a large corpus.
- does not directly capture global information of the corpus.

The GloVe Approach

GloVe is a global log-bilinear regression model. More specifically:

- a weighted least squares model trained on global word-word co-occurrence counts obtained from a large corpus, rather than on
 - sparse term-term matrices
 - a sliding window of local contexts over a large corpus

Distinguishing ratios of target words with discriminative and non-discriminative context words

Probability	k = solid	$k = \mathit{gas}$	k = water	k = fashion
and ratio				
$P(k \mid ice)$				
$P(k \mid steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k \mid ice)/$	8.9	8.5×10^{-2}	1.36	0.96
$P(k \mid steam)$				

with target words: ice, steam with discriminative words: gas, solid

with "noise" words: water, fashion

Loss Function for a Weighted Least Squares Regression Model

$$J = \sum_{i,j=1}^{V} = f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - log X_{ij})^2$$
 (1)

where: f is a weighting function:

$$f(x) = \begin{cases} (x/x_{max})^{\alpha} & \text{if } x < x_{max} \\ 1 & \text{otherwise} \end{cases}$$
 (2)

 $X_{i,j}$ tabulates the number of times word j occurs in the context of word i

Weighting Function f with $\alpha = 3/4$ and

$$x_{max} = 1$$

