

Vector Semantics

Dense Vectors

ntics



Sparse versus dense vectors

- PPMI vectors are
 - **long** (length |V|= 20,000 to 50,000)
 - sparse (most elements are zero)
- Alternative: learn vectors which are
 - short (length 200-1000)
 - dense (most elements are non-zero)



Sparse versus dense vectors

- Why dense vectors?
 - Short vectors may be easier to use as features in machine learning (less weights to tune)
 - Dense vectors may generalize better than storing explicit counts
 - They may do better at capturing synonymy:
 - car and automobile are synonyms; but are represented as distinct dimensions; this fails to capture similarity between a word with *car* as a neighbor and a word with *automobile* as a neighbor





Three methods for getting short dense vectors

- Singular Value Decomposition (SVD)
 - A special case of this is called LSA Latent Semantic Analysis
- "Neural Language Model"-inspired predictive models
 - skip-grams and CBOW
- Brown clustering





Vector Semantics Dense Vectors via SVD



Intuition

- Approximate an N-dimensional dataset using fewer dimensions
- By first rotating the axes into a new space
- In which the highest order dimension captures the most • variance in the original dataset
- And the next dimension captures the next most variance, etc.
- Many such (related) methods:
 - PCA principle components analysis
 - Factor Analysis
 - SVD

Dan Jurafsky



Dimensionality reduction

6





Singular Value Decomposition

Any rectangular w x c matrix X equals the product of 3 matrices:

W: rows corresponding to original but m columns represents a dimension in a new latent space, such that

- M column vectors are orthogonal to each other
- Columns are ordered by the amount of variance in the dataset each new dimension accounts for
- S: diagonal m x m matrix of singular values expressing the importance of each dimension.

C: columns corresponding to original but m rows corresponding to **Šingular** values

Dan Jurafsky



Singular Value Decomposition



w xm WXC

Landuaer and Dumais 1997

Dan Jurafsky



SVD applied to term-document matrix: **Latent Semantic Analysis** Deerwester et al (1988)

- If instead of keeping all m dimensions, we just keep the top k singular values. Let's say 300.
- The result is a least-squares approximation to the original X
- But instead of multiplying, • we'll just make use of W.
- Each row of W:
 - A k-dimensional vector
 - Representing word W











LSA more details

- 300 dimensions are commonly used
- The cells are commonly weighted by a product of two weights
 - Local weight: Log term frequency
 - Global weight: either idf or an entropy measure



Let's return to PPMI word-word matrices

• Can we apply to SVD to them?

Dan Jurafsky



SVD applied to term-term matrix



(I'm simplifying here by assuming the matrix has rank |V|)

|V| imes |V|

Dan Jurafsky



Truncated SVD on term-term matrix





Truncated SVD produces embeddings

- Each row of W matrix is a k-dimensional representation of each word w
- K might range from 50 to 1000
- Generally we keep the top k dimensions, but some experiments suggest that getting rid of the top 1 dimension or even the top 50 dimensions is helpful (Lapesa and Evert 2014).



for



Embeddings versus sparse vectors

- Dense SVD embeddings sometimes work better than sparse PPMI matrices at tasks like word similarity
 - Denoising: low-order dimensions may represent unimportant information
 - Truncation may help the models generalize better to unseen data.
 - Having a smaller number of dimensions may make it easier for classifiers to properly weight the dimensions for the task.
 - Dense models may do better at capturing higher order cooccurrence.



Vector Semantics

Embeddings inspired by neural language models: skip-grams and CBOW

Dan Jurafsky



Prediction-based models: An alternative way to get dense vectors

- **Skip-gram** (Mikolov et al. 2013a) **CBOW** (Mikolov et al. 2013b)
- Learn embeddings as part of the process of word prediction.
- Train a neural network to predict neighboring words
 - Inspired by **neural net language models**.
 - In so doing, learn dense embeddings for the words in the training corpus.
- Advantages:
 - Fast, easy to train (much faster than SVD)
 - Available online in the word2vec package
 - Including sets of pretrained embeddings!



Skip-grams

- Predict each neighboring word
 - in a context window of 2C words
 - from the current word.
- So for C=2, we are given word w_t and predicting these 4 words:

$$[w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2}]$$





Skip-grams learn 2 embeddings for each w

input embedding v, in the input matrix W

Column *i* of the input matrix W is the $1 \times d$ embedding v_i for word *i* in the vocabulary.

output embedding v', in output matrix W'

• Row *i* of the output matrix W' is a $d \times 1$ vector embedding v'_i for word *i* in the vocabulary.



2

V x d





Setup

- Walking through corpus pointing at word w(t), whose index in the vocabulary is j, so we'll call it w_j (1 < j < |V|).
- Let's predict w(t+1), whose index in the vocabulary is k (1 < k < |V|). Hence our task is to compute $P(w_k|w_i)$.





Intuition: similarity as dot-product between a target vector and context vector







Similarity is computed from dot product

- Remember: two vectors are similar if they have a high dot product
 - Cosine is just a normalized dot product
- So:
 - Similarity(j,k) $\propto c_k \cdot v_j$
- We'll need to normalize to get a probability

roduct ave a high

for word k

• We use softmax to turn into probabilities

$$p(w_k|w_j) = \frac{exp(c_k \cdot v_j)}{\sum_{i \in |V|} exp(c_i \cdot v_j)}$$



Embeddings from W and W'

- Since we have two embeddings, v_i and c_i for each word w_i
- We can either:
 - Just use v_i
 - Sum them
 - Concatenate them to make a double-length embedding



Learning

- Start with some initial embeddings (e.g., random)
- iteratively make the embeddings for a word
 - more like the embeddings of its neighbors
 - less like the embeddings of other words.

Dan Jurafsky



Visualizing W and C as a network for doing error backprop



У₁ y_2 w_{t+1} $\mathbf{y}_{\mathbf{k}}$

 $y_{|V|}$

Dan Jurafsky



One-hot vectors

- A vector of length |V|
- 1 for the target word and 0 for other words
- So if "popsicle" is vocabulary word 5
- The one-hot vector is
- [0,0,0,0,1,0,0,0,0,.....0]







Skip-gram

Input layer

$$h = v_j$$

Projection layer





o = Ch $o_k = c_k h$ $o_k = c_k V_j$ Output layer probabilities of



The denominator: have to compute over every word in vocab

$$p(w_k|w_j) = \frac{exp(c_k \cdot v_j)}{\sum_{i \in |V|} exp(c_i \cdot v_j)}$$

Instead: just sample a few of those negative words

Dan Jurafsky



Goal in learning

Make the word like the context words

a [tablespoon of apricot preserves or] jam lemon, **c**1 c2 w c3 c4 • We want this to be high:

 $\sigma(c1 \cdot w) + \sigma(c2 \cdot w) + \sigma(c3 \cdot w) + \sigma(c4 \cdot w)$

- And not like k randomly selected "noise words" [cement metaphysical dear coaxial apricot attendant whence forever puddle] n1 n2 n3 **n**5 n6 n4
- We want this to be low:

 $\sigma(n1 \cdot w) + \sigma(n2 \cdot w) + \ldots + \sigma(n8 \cdot w)$ 31



n7 n8





Skipgram with negative sampling: Loss function





Relation between skipgrams and PMI!

- If we multiply WW'^{T}
- We get a |V|x|V| matrix M, each entry m_{ii} corresponding to some association between input word *i* and output word *j*
- Levy and Goldberg (2014b) show that skip-gram reaches its optimum just when this matrix is a shifted version of PMI: $WW'^{T} = M^{PMI} - \log k$
- So skip-gram is implicitly factoring a shifted version of the PMI matrix into the two embedding matrices.



Properties of embeddings

Nearest words to some embeddings (Mikolov et al. 20131)

target:	Redmond	Havel	ninjutsu	gra
	Redmond Wash.	Vaclav Havel	ninja	spr
	Redmond Washington	president Vaclav Havel	martial arts	gra
	Microsoft	Velvet Revolution	swordsmanship	tag

ffiti ay paint fitti gers

capitulate capitulation capitulated capitulating



Embeddings capture relational meaning!

vector('king') - vector('man') + vector('woman') \approx vector('queen') vector('*Paris'*) - vector('*France'*) + vector('*Italy'*) \approx vector('Rome')



Cross-lingual Embeddings

 Skip-gram allows us learning embeddings for words in a single language



Cross-lingual Embeddings

- Skip-gram allows us learning embeddings for words in a single language
- But what if we want to work with multiple languages?



General Schema for Cross-lingual Embeddings



General Schema for Cross-lingual Embeddings



Sources of Cross-Lingual Supervision



BiSparse - Sparse Bilingual Embeddings

- A method to learn embeddings, that are
 - Bilingual
 - Sparse
 - Non-negative
- Starting from
 - Monolingual embeddings in two languages
 - A "seed" dictionary

• Method based on matrix factorization



Method based on matrix factorization



Method based on matrix factorization



Method based on matrix factorization



Building the S Matrix



•

Interpreting Embeddings

French Dimensions	English Dimensions	
logiciel, fichiers, web, microsoft	files, web, microsoft, www	
université, collège, lycée, conseil de administration	university, college, graduate, faculty	
virus informatique, virus, infection, cellules	virus, viruses, infection, cells	
doigts, genoux, jambes, muscles	bruises, fingers, toes, knees	
budapest, stockholm, copenhague, buenos	lahore, dhaka, harare, karachi	

Summary

- Vector Semantics with Dense Vectors
 - Singular Value Decomposition
 - Skip-gram embeddings
- Cross-lingual embeddings
 - BiSparse model