

Vector Semantics

Dense Vectors

Outline

- Motivation
- Dense Vectors via SVD
- Embeddings by neural language models
 - skip-grams and CBOW
- Brown clustering

Taxonomy vs Context Vector

- Great as resource but missing nuances
- Missing new words (impossible to keep up to date)
- Subjective
- Requires human labor to create and adapt
- Hard to compute accurate word similarity

Sparse vs Dense Vectors

- Discrete representation
 - **long** (length $|V| = 20,000$ to $50,000$)
 - **sparse** (most elements are zero)
- Dense representation
 - **short** (length 200-1000)
 - **dense** (most elements are non-zero)

Problems with the Discrete Representation

- In vector space terms, this is a vector with one 1 and a lot of zeroes
- We call this a “one-hot” representation
 $[0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0]$
- Dimensionality:
 - 20K (speech) – 50K (PTB) – 500K (big vocab) – 13M (Google 1T)

car: $[0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0]$

automobile: $[0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0]$

Sparse vs 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 they 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

Outline

- Motivation
- Dense Vectors via SVD
- Embeddings by neural language models
 - skip-grams and CBOW
- Brown clustering

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

Singular Value Decomposition

Any rectangular $w \times c$ matrix X equals the product of 3 matrices:

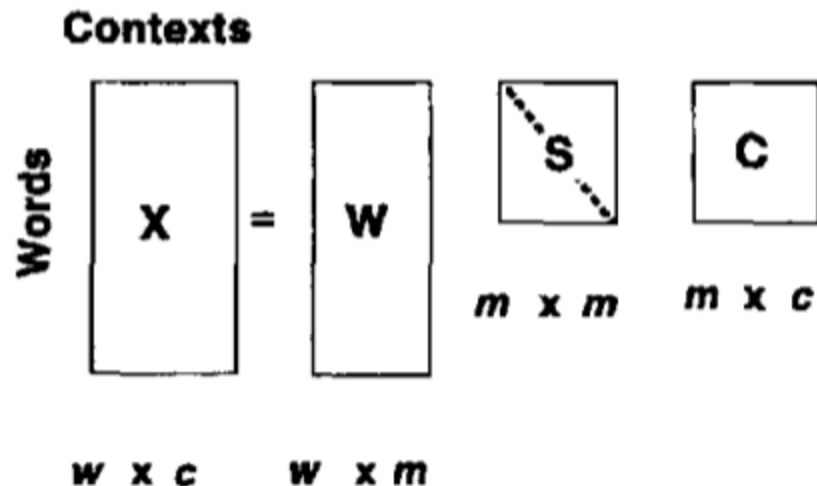
W: rows corresponding to original but m columns represents a dimension in a new latent space, such that

- 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 $m \times m$ matrix of **singular values** expressing the importance of each dimension.

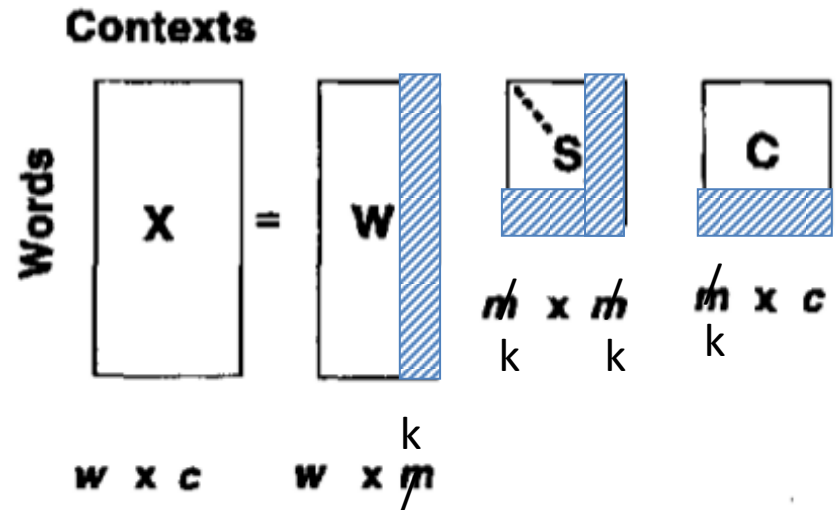
C: columns corresponding to original but m rows corresponding to singular values

Singular Value Decomposition



SVD applied to term-document matrix: Latent Semantic Analysis

- 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?

SVD applied to term-term matrix

$$\begin{bmatrix} X \\ |V| \times |V| \end{bmatrix} = \begin{bmatrix} W \\ |V| \times |V| \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_V \\ |V| \times |V| \end{bmatrix} \begin{bmatrix} C \\ |V| \times |V| \end{bmatrix}$$

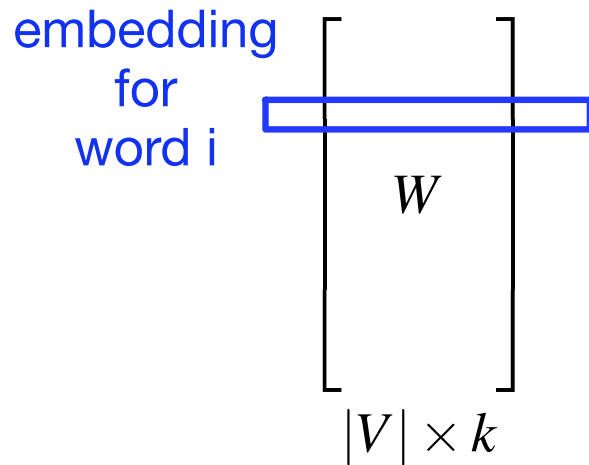
(assuming the matrix has rank $|V|$) 15

Truncated SVD on term-term matrix

$$\begin{bmatrix} X \\ |V| \times |V| \end{bmatrix} = \begin{bmatrix} W \\ |V| \times k \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_k \\ k \times k \end{bmatrix} \begin{bmatrix} C \\ k \times |V| \end{bmatrix}$$

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, A Large Scale Evaluation of Distributional Semantic Models: Parameters, Interactions and Model Selection, 2014).

Embeddings vs 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 co-occurrence.

Problems with SVD

- Computational cost scales quadratically for $n \times m$ matrix:
 $O(mn^2)$ when $n < m$
→ Bad for millions of words or documents
- Hard to incorporate new words or documents

Outline

- Motivation
- Dense Vectors via SVD
- Embeddings by neural language models
 - skip-grams and CBOW
- Brown clustering

Prediction-based models:

An alternative way to get dense vectors

- Idea:
 - Instead of capturing co-occurrence counts directly, predict surrounding words of every word
 - Learn embeddings as part of the process of word prediction.
 - Both are quite similar (see “Glove: Global Vectors for Word Representation” by Pennington et al. (2014) and Levy and Goldberg (2014))
- 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.

Prediction-based models:

An alternative way to get dense vectors

- Advantages:
 - Fast, easy to train (much faster than SVD)
 - can easily incorporate a new sentence/document or add a word to the vocabulary
 - Including sets of pretrained embeddings!
- Available models
 - **Skip-gram** (Mikolov et al. 2013a)
 - **CBOW** (Mikolov et al. 2013b)
- Available online in the `word2vec` package

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:

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

Setup

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

Details of Skip-grams

- Objective function: Maximize the log probability of any context word given the current center word

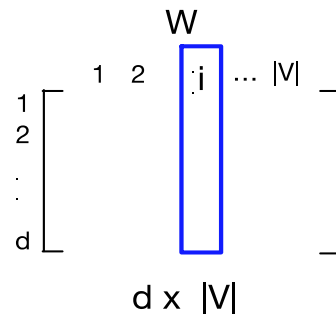
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t)$$

- Where θ represents all variables we optimize

Skip-grams learn 2 embeddings for each w

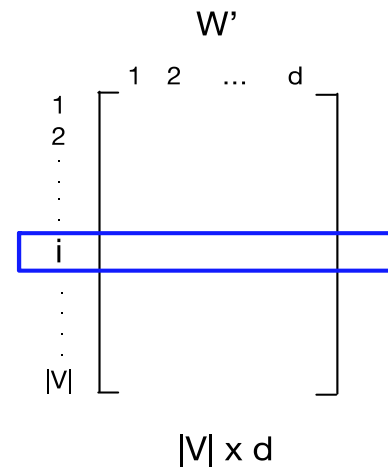
input embedding v , in the input matrix W

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

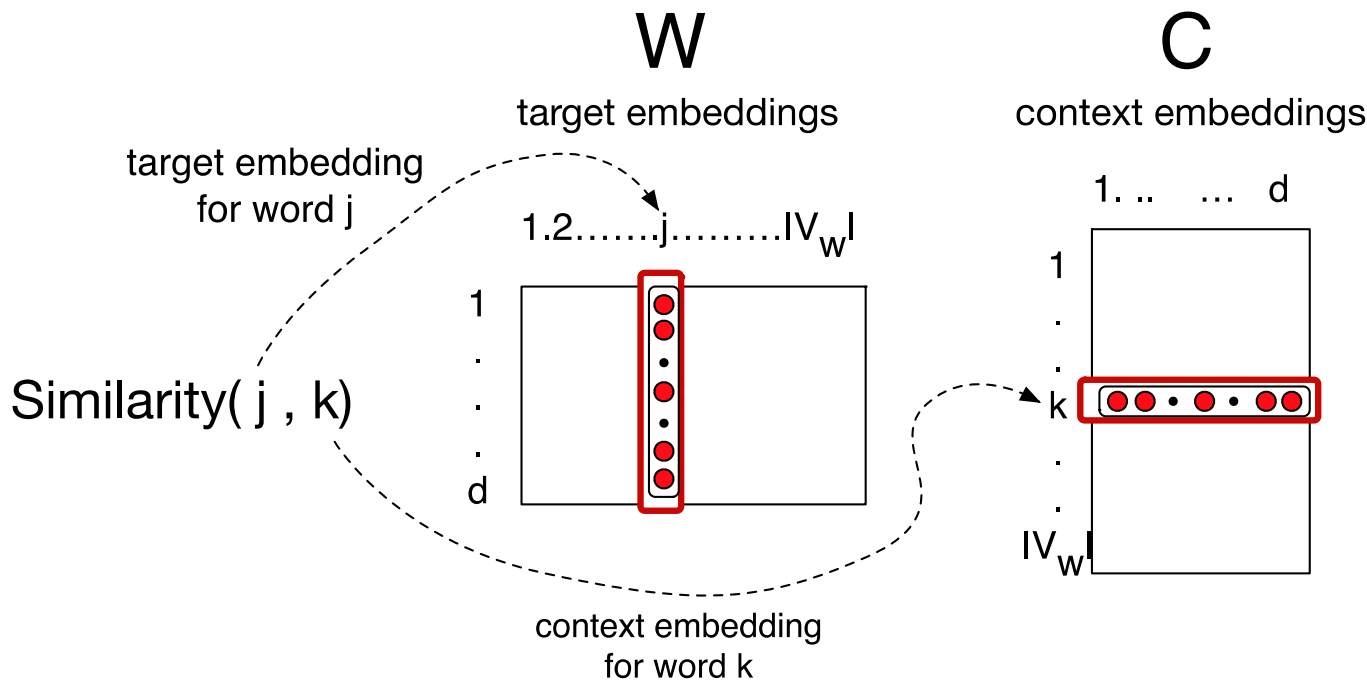


output embedding v' , in output matrix W'

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



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:
 - $\text{Similarity}(s,t) \propto u_s \cdot v_t$
- We'll need to normalize to get a probability

Turning dot products into probabilities

- Predict surrounding words in a window of length m of every word
- Using the softmax function for $p(w_{t+j}|w_t)$ the simplest first formulation is

$$p(w_s|w_t) = \frac{\exp(u_s \cdot v_t)}{\sum_{w \in |V|} \exp(u_w \cdot v_t)}$$

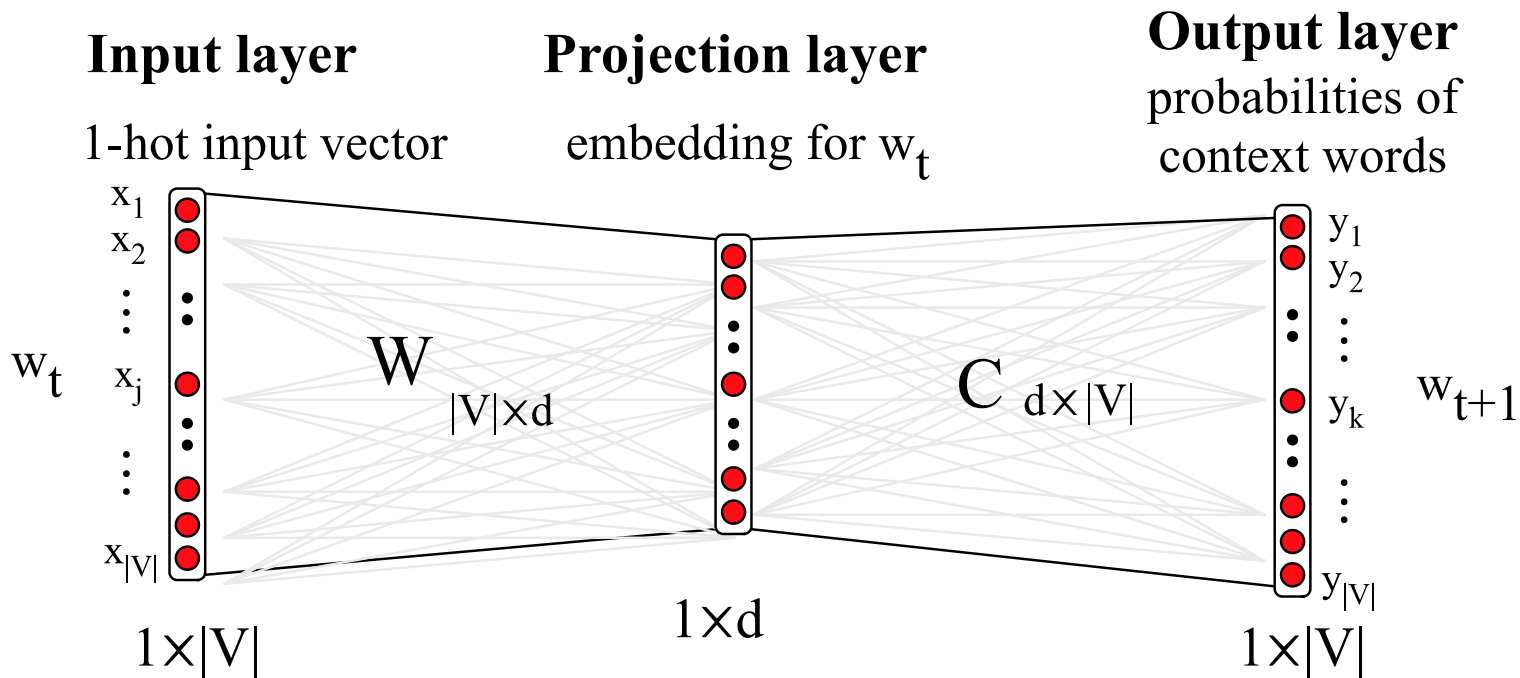
Embeddings from W and W'

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

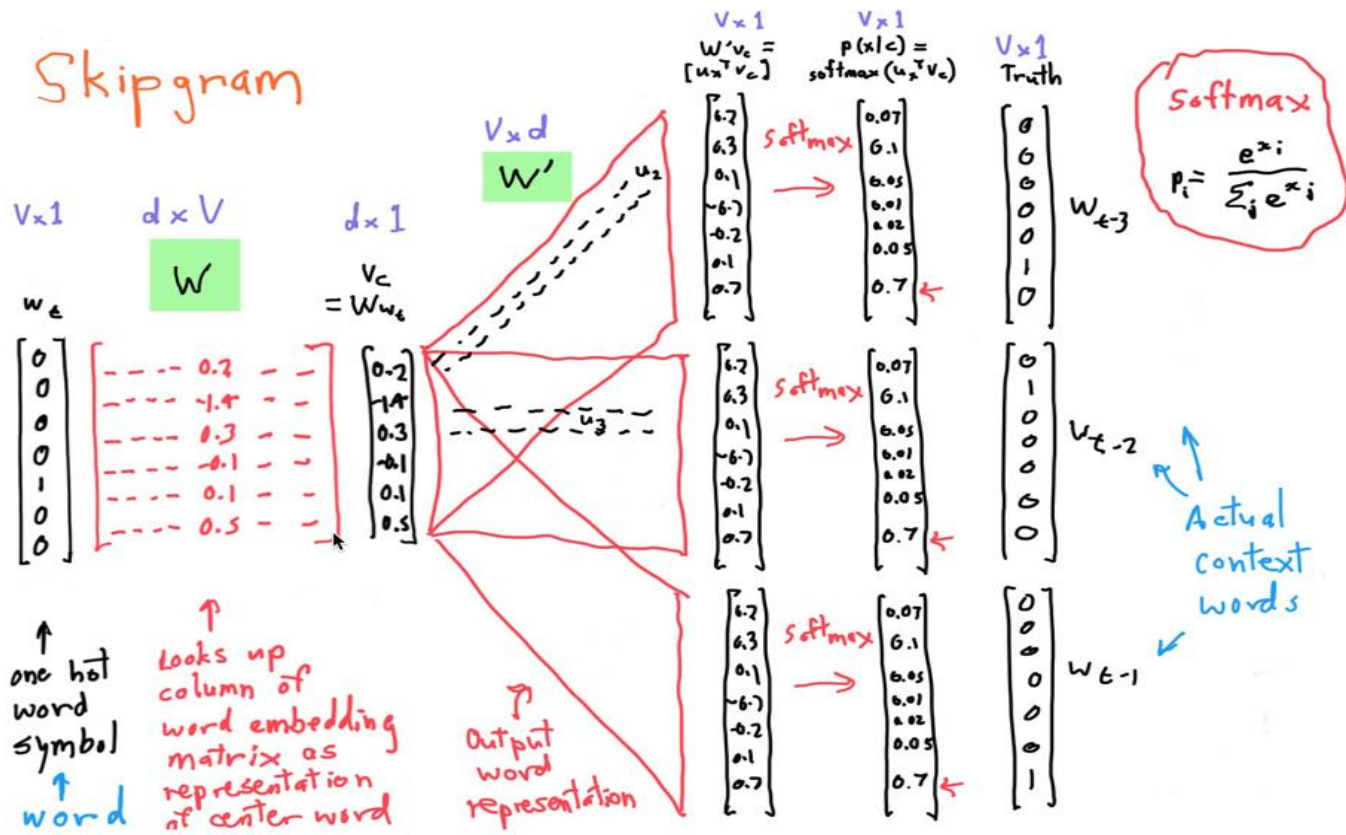
Learning

- 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.

Visualizing W and C as a network for doing error back propagation



Visualizing W and C as a network for doing error back propagation



Problem with the softmax

- The denominator: have to compute over every word in vocab

$$p(w_s|w_t) = \frac{\exp(u_s^T \cdot v_t)}{\sum_{w \in |V|} \exp(u_w^T \cdot v_t)}$$

- Instead: just sample a few of those negative words

Goal in learning

- Make the word like the context words

lemon, a [tablespoon of apricot preserves or] jam
c1 c2 w c3 c4

$$\sigma(x) = \frac{1}{1+e^x}$$

- 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 n4 n5 n6 n7 n8

- We want this to be low:

$$\sigma(n1 \cdot w) + \sigma(n2 \cdot w) + \dots + \sigma(n8 \cdot w)$$

Skipgram with negative sampling: Loss function

$$\log \sigma(c \cdot w) + \sum_{i=1}^{\kappa} \mathbb{E}_{w_i \sim p(w)} [\log \sigma(-w_i \cdot w)]$$

Properties of embeddings

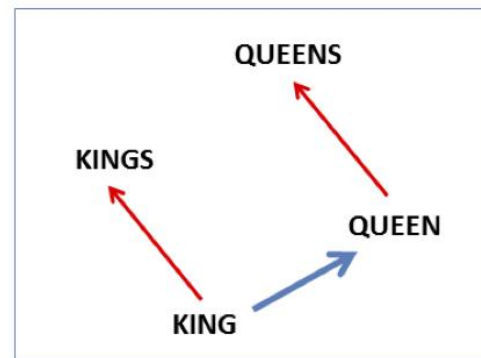
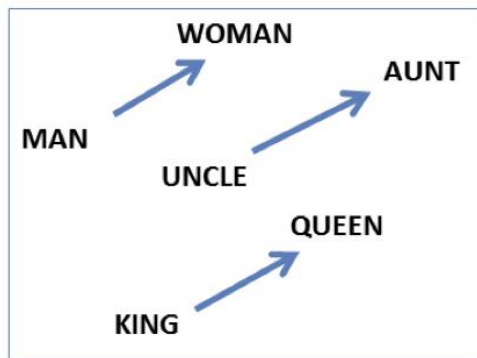
- Nearest words to some embeddings

FRANCE	JESUS	XBOX	REDDISH	SCRATCHED	MEGABITS
AUSTRIA	GOD	AMIGA	GREENISH	NAILED	OCTETS
BELGIUM	SATI	PLAYSTATION	BLUISH	SMASHED	MB/S
GERMANY	CHRIST	MSX	PINKISH	PUNCHED	BIT/S
ITALY	SATAN	IPOD	PURPLISH	POPPED	BAUD
GREECE	KALI	SEGA	BROWNISH	CRIMPED	CARATS
SWEDEN	INDRA	PSNUMBER	GREYISH	SCRAPED	KBIT/S
NORWAY	VISHNU	HD	GRAYISH	SCREWED	MEGAHERTZ
EUROPE	ANANDA	DREAMCAST	WHITISH	SECTIONED	MEGAPIXELS
HUNGARY	PARVATI	GEFORCE	SILVERY	SLASHED	GBIT/S
SWITZERLAND	GRACE	CAPCOM	YELLOWISH	RIPPED	AMPERES

Embeddings capture relational meaning!

$\text{vector}('king') - \text{vector}('man') + \text{vector}('woman') \approx \text{vector}('queen')$

$\text{vector}('Paris') - \text{vector}('France') + \text{vector}('Italy') \approx \text{vector}('Rome')$



Count based vs direct prediction

LSA, HAL (Lund & Burgess)
COALS (Rohde et al)
Hellinger-PCA (Lebret & Collobert)

NNLM, HLBL, RNN, Skip-gram/CBOW
(Bengio et al; Collobert & Weston; Huang et al; Mnih & Hinton; Mikolov et al; Mnih & Kavukcuoglu)

- Fast training
- Efficient usage of statistics
- Primarily used to capture word similarity
- Disproportionate importance given to large counts

- Scales with corpus size
- Inefficient usage of statistics
- Generate improved performance on other tasks
- Can capture complex patterns beyond word similarity

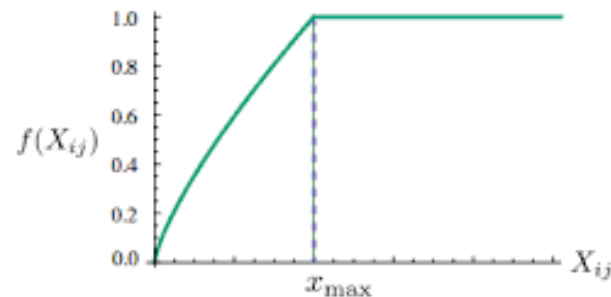
Relation between skipgrams and PMI!

- If we multiply WW'^T
- We get a $|V| \times |V|$ matrix M , each entry m_{ij} 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:
$$WW'^T = M^{\text{PMI}} - \log k$$
- So skip-gram is implicitly factoring a shifted version of the PMI matrix into the two embedding matrices.

GloVe: Combining the best of both worlds

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W f(P_{ij}) (u_i^T v_j - \log P_{ij})^2$$

- X_{ij} : the number of times word j occurs in the context of word i
- $X_i = \sum_k X_{ik}$: the number of times any word appears in the context of word i
- $P_{ij} = P(j|i) = X_{ij}/X_i$: probability that word j appear in the context of word i



GloVe: Combining the best of both worlds

- Fast training
- Scalable to huge corpora
- Good performance even with small corpus, and small vectors

Outline

- Motivation
- Dense Vectors via SVD
- Embeddings by neural language models
 - skip-grams and CBOW
- Brown clustering

Brown clustering

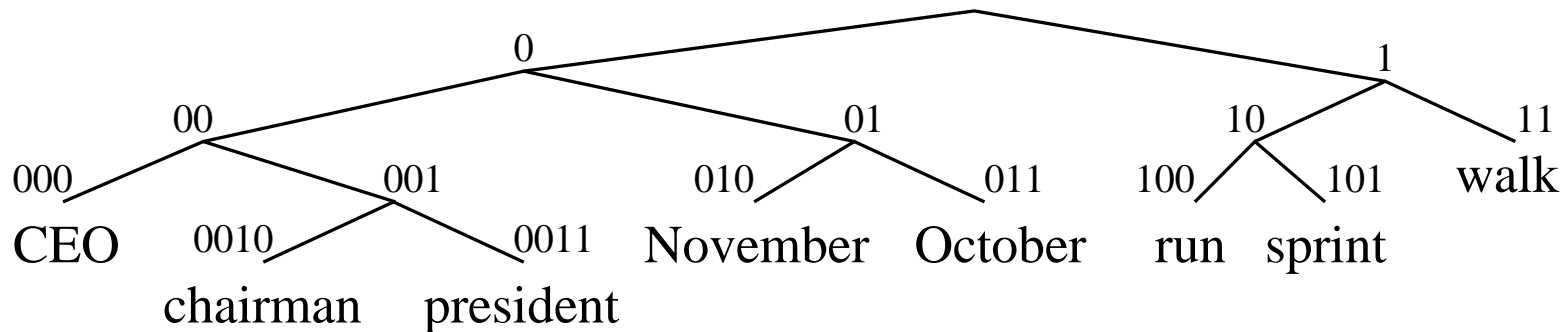
- An agglomerative clustering algorithm that clusters words based on which words precede or follow them
- These word clusters can be turned into a kind of vector

Brown clustering algorithm

- Each word is initially assigned to its own cluster.
- We now consider merging each pair of clusters.
Highest quality merge is chosen.
 - Quality = merges two words that have similar probabilities of preceding and following words
 - More technically quality = smallest decrease in the likelihood of the corpus according to a class-based language model
- Clustering proceeds until all words are in one big cluster.

Brown Clusters as vectors

- By tracing the order in which clusters are merged, the model builds a binary tree from bottom to top.
- Each word represented by binary string: path from root to leaf
- Each intermediate node is a cluster
- Chairman is 0010, “months” = 01, and “verbs” = 1



Brown cluster examples

Friday Monday Thursday Wednesday Tuesday Saturday Sunday weekends Sundays Saturdays
June March July April January December October November September August
pressure temperature permeability density porosity stress velocity viscosity gravity tension
anyone someone anybody somebody
had hadn't hath would've could've should've must've might've
asking telling wondering instructing informing kidding reminding bothering thanking deposing
mother wife father son husband brother daughter sister boss uncle
great big vast sudden mere sheer gigantic lifelong scant colossal
down backwards ashore sideways southward northward overboard aloft downwards adrift

Class-based language model

- Suppose each word was in some class c_i

$$P(w_i|w_{i-1}) = P(c_i|c_{i-1})P(w_i|c_i)$$