

# COMP9444: Neural Networks and Deep Learning Week 7a. Word Vectors

Sonit Singh School of Computer Science and Engineering October 25, 2023

# Outline

- Statistical Language Processing
- *n*-gram models
- co-occurence matrix
- word representation
- Word2Vec
- GloVe
- word relationships



# Word Meaning - Synonyms and Taxonomy?

What is the meaning of meaning?

- dictionary definitions
- synonyms and antonyms
- taxonomy
  - penguin is-a bird is-a mammal is-a vertebrae



## **Statistical Language Processing**

Synonyms for "elegant"

stylish, graceful, tasteful, discerning, refined, sophisticated, dignified, cultivated, distinguished, classic, smart, fashionable, modish, decorous, beautiful, artistic, aesthetic, lovely; charming, polished, suave, urbane, cultured, dashing, debonair; luxurious, sumptuous, opulent, grand, plush, high-class, exquisite

Synonyms, antonyms and taxonomy require human effort, may be incomplete and require discrete choices. Nuances are lost. Words like "king", "queen" can be similar in some attributes but opposite in others.

Could we instead extract some statistical properties automatically, without human involvement?



#### There was a Crooked Man

There was a crooked man, who walked a crooked mile And found a crooked sixpence upon a crooked stile. He bought a crooked cat, who caught a crooked mouse And they all lived together in a little crooked house.



www.kearley.co.uk/images/uploads/JohnPatiencePJ03.gif



# **Counting Frequencies**

word	frequency
а	7
all	1
and	2
bought	1
cat	1
caught	1
crooked	7
found	1
he	1
house	1
in	1
little	1
lived	1
man	1
mile	1
mouse	1
sixpence	1
stile	1
there	1
they	1
together	1
upon	1
walked	1
was	1
who	2

- some words occur frequently in all (or most) documents
- some words occur frequently in a particular document, but not generally
- this information can be useful for document classification



## **Document Classification**

word	doc 1	doc 2	doc X
а			7
all			1
and			2
bought			1
cat			1
caught			1
crooked			7
found			1
he			1
house			1
in			1
little			1
lived			1
man			1
mile			1
mouse			1
sixpence			1
stile			1
there			1
they			1
together			1
upon			1
walked			1
was			1
who			2



#### **Document Classification**

- each column of the matrix becomes a vector representing the corresponding document
- words like "cat", "mouse", "house" tend to occur in children's books or rhymes
- other groups of words may be characteristic of legal documents, political news, sporting results, etc.
- words occurring many times in one document may skew the vector might be better to just have a "1" or "0" indicating whether the word occurs at all



# **Counting Consecutive Word Pairs**

word	a all	and bought			he	house in	little	lived	man	mile	sixpence	stile	there	they	together	uodn	walked was	who
a			6	6			1											
all and				1				1						1				
bought	1			'										'				
cat caught	1																	1
crooked found	1	1				1			1	1	11	1						
he house		1																
in	1																	
little			1															
lived															1			
man mile		1																1
mouse		1																
sixpence																1		
stile					1													
there																	1	
they together	1					1												
upon	1																	
walked	1																	
was	1																	
who			1														1	



## **Predictive 1-Gram Word Model**

word	ß	all	and bought	cat	caught	crooked	found	he	house	.Ľ	little	lived	man	mile	mouse	sixpence	stile	there	they	together	uodn	walked	wds	
а						<u>6</u> 7					7													
all												1												
and							$\frac{1}{2}$												$\frac{1}{2}$					
bought cat	1						_																1	
caught	1																							
crooked				$\frac{1}{7}$					$\frac{1}{7}$				1	1	1	$\frac{1}{7}$	1							
found	1			7					7				7	7	7	7	7							
he			1																					
house																								
in	1																							
little						1																		
lived																				1				
man																							1	
mile			1																					
mouse			1																					
sixpence stile								1													1			
there								1															1	
they		1																						
together										1														
upon	1																							
walked	1																							
was	1																							
who					$\frac{1}{2}$																	$\frac{1}{2}$		



## **N-Gram Model**

- by normalizing each row (to sum to 1) we can estimate the probability  $prob(w_j|w_i)$  of word  $w_j$  occurring after  $w_i$
- need to aggregrate over a large corpus, so that unusual words like "crooked" will not dominate
- the model captures some common combinations like "there was", "man who", "and found", "he bought", "who caught", "and they", "they all", "lived together", etc.
- this *unigram* model can be generalized to a bi-gram, tri-gram, ..., *n*-gram model by considering the *n* preceding words
- if the vocabulary is large, we need some tricks to avoid exponential use of memory



#### **1-Gram Text Generator**

"Rashly – Good night is very liberal – it is easily said there is – gyved to a sore distraction in wrath and with my king may choose but none of shapes and editing by this , and shows a sea And what this is miching malhecho ; And gins to me a pass , Transports his wit , Hamlet , my arms against the mind impatient , by the conditions that would fain know ; which , the wicked deed to get from a deed to your tutor ."



#### **Co-occurrence Matrix**

- sometimes, we don't necessarily predict the next , but simply a "nearby word" (e.g. a word occurring within an *n*-word window centered on that word)
- we can build a matrix in which each row represents a word, and each column a nearby word
- each row of this matrix could be considered as a vector representation for the corresponding word, but the number of dimensions is equal to the size of the vocabulary, which could be very large ( $\sim 10^5$ )
  - is there a way to reduce the dimensionality while still preserving the relationships between words?



## **Co-occurrence Matrix (2-word window)**

word	g	all	and	bought	cat	caught	crooked	found	he	house	.⊑	little	lived	man	mile	mouse	sixpence	stile	there	they	together	uodn	walked	was	who	
а				1		1	6	1			1	1										1	1	1		
all													1							1						
and								1							1	1				1						
bought	1								1																	
cat							1																		1	
caught crooked	1 6				4					1		1		4	4	4	4	1							1	
found	1		1		'							1		'				1								
he	'			1														1								
house							1											'								
in	1																				1					
little	1						1																			
lived		1																			1					
man							1																		1	
mile			1				1																			
mouse			1				1																			
sixpence							1															1				
stile							1		1																	
there																								1		
they		1	1																							
together											1		1													
upon walked	1																1								4	
walked was	1																		1						1	
was who	'				1	1								1									1			
WIIO														1												



# **Co-occurrence Matrix (10-word window)**

word		в	all	and	bought	cat	caught	crooked	found	he	house	. <u>Ľ</u>	little	lived	man	mile	mouse	sixpence	stile	there	they	together	uodn	walked	was	who
а		10	2	3	2	2	2	13	3	2	1	1	1	1	2	2	1	2	2	1	2	1	2	2	1	4
all		2		1				1				1	1	1			1				1	1				
and		2 3	1				1	3	1			1		1		1	1	1			1	1	1	1		2
bought		2 2 2				1	1	2		1									1				1			1
cat		2			1		1	2		1							1		1							1
caught				1	1	1		2									1				1					1
crooked		13	1	3	2	2	2	10	2	2	1	1	1	2	2	2	1	2	3	1	1	1	2	2	1	4
found		3		1				2								1		1					1	1		
he		2			1	1		2										1	1				1			1
house		1						1				1	1									1				
in		1	1	1				1			1		1	1							1	1				
little		1	1					1			1	1		1								1				
lived		1	1	1				2				1	1				1				1	1				
man		2						2								1				1				1	1	1
mile		2		1				2	1						1			1						1		1
mouse		1	1	1		1	1	1						1							1	1				1
sixpence		2		1				2	1	1						1			1				1			
stile		2			1	1		3		1								1					1			
there		1						1							1										1	1
they		2	1	1			1	1				1		1			1					1				
together		1	1	1				1			1	1	1	1			1				1					
upon		2		1	1			2	1	1								1	1							
walked		2		1				2	1						1	1									1	1
was		1		~				1							1					1				1		1
who	I	4		2	1	1	1	4		1					1	1	1			1				1	1	



#### **Co-occurrence Matrix**

- by aggregating over many documents, pairs (or groups) of words emerge which tend to occur near each other (but not necessarily consecutively)
  - "cat", "caught", "mouse"
    "walked", "mile"

  - "little". "house"
- common words tend to dominate the matrix
  - could we sample common words less often, in order to reveal the relationships of less common words?



## Word Embeddings

"Words that are used and occur in the same contexts tend to purport similar meanings."

Z. Harris (1954)

"You shall know a word by the company it keeps."

J.R. Firth (1957)

Aim of Word Embeddings:

Find a vector representation of each word, such that words with nearby representations are likely to occur in similar contexts.

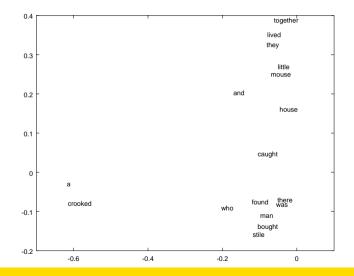


## **History of Word Embeddings**

- Structuralist Linguistics (Firth, 1957)
- Recurrent Networks (Rumelhart, Hinton & Williams, 1986)
- Latent Semantic Analysis (Deerwester et al., 1990)
- Hyperspace Analogue to Language (Lund, Burgess & Atchley, 1995)
- Neural Probabilistic Language Models (Bengio, 2000)
- NLP (almost) from Scratch (Collobert et al., 2008)
- word2vec (Mikolov et al., 2013)
- GloVe (Pennington, Socher & Manning, 2014)



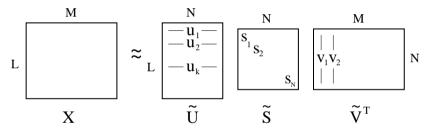
## Word Embeddings





## **Singular Value Decomposition**

Co-occurrence matrix  $X_{(L \times M)}$  can be decomposed as  $X = U S V^T$  where  $U_{(L \times L)}$ ,  $V_{(M \times M)}$  are unitary (all columns have unit length) and  $S_{(L \times M)}$  is diagonal, with diagonal entries  $s_1 \ge s_2 \ge \ldots \ge s_M \ge 0$ 



We can obtain an approximation for X of rank N < M by truncating U to  $\tilde{U}_{(L \times N)}$ , S to  $\tilde{S}_{(N \times N)}$  and V to  $\tilde{V}_{(N \times M)}$ . The *k*th row of  $\tilde{U}$  then provides an *N*-dimensional vector representing the  $k^{\text{th}}$  word in the vocabulary.



# Word2Vec and GloVe

For language processing tasks, typically, *L* is the number of words in the vocabulary (about 60,000) and *M* is either equal to *L* or, in the case of document classification, the number of documents in the collection. SVD is computationally expensive, proportional to  $L \times M^2$  if  $L \ge M$ . Can we generate word vectors in a similar way but with less computation, and incrementally?

- Word2Vec
  - predictive model
  - maximize the probability of a word based on surrounding words
- GloVe
  - count-based model
  - reconstruct a close approximation to the co-occurrence matrix X



### **Eigenvalue vs. Singular Value Decomposition**

Eigenvalue Decomposition:

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \Omega D \Omega^{-1}, \quad \text{where} \quad \Omega = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}, \quad D = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$
$$\begin{bmatrix} 0 - 1 \\ 1 & 0 \end{bmatrix} = \Omega D \Omega^{-1}, \quad \text{where} \quad \Omega = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -i & i \end{bmatrix}, \quad D = \begin{bmatrix} i & 0 \\ 0 & -i \end{bmatrix}$$

Singular Value Decomposition:

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \mathrm{USV}^{\mathrm{T}}, \quad \mathrm{U} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathrm{S} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathrm{V} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$
$$\begin{bmatrix} 0 - 1 \\ 1 & 0 \end{bmatrix} = \mathrm{USV}^{\mathrm{T}}, \quad \mathrm{U} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathrm{S} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathrm{V} = \begin{bmatrix} 1 & 0 \\ 0 - 1 \end{bmatrix}$$



## **Eigenvalue vs. Singular Value Decomposition**

- if X is symmetric and positive semi-definite, eigenvalue and singular value decompositions are the same.
- in general, eigenvalues can be negative or even complex, but singular values are always real and non-negative.
- even if X is a square matrix, singular value decompositon treats the source and target as two entirely different spaces.
- the word co-occurrence matrix is symmetric but not positive semi-definite; for example, if the text consisted entirely of two alternating letters ...ABABABABABABABA... then A would be the context for B, and vice-versa.

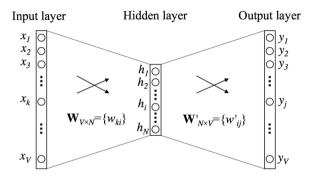


# word2vec (Mikolov et al., 2013)

- Idea: predict rather than count
- Instead of counting how often each word *w* occurs near "university" train a classifier on a **binary prediction task**:
  - Is w likely to show up near "university"?
- We don't actually care about this task
  - $\circ~$  But we'll take the learned classifier weights as the word embeddings
- Use running text as implicitly supervised training data
- No need for hand-labeled supervision



# Word2Vec 1-word Context Model



The  $k^{\text{th}}$  row  $\mathbf{v}_k$  of  $\mathbf{W}$  is a representation of word k. The  $j^{\text{th}}$  column  $\mathbf{v}'_j$  of  $\mathbf{W}'$  is an (alternative) representation of word j. If the (1-hot) input is k, the linear sum at each output will be  $u_j = \mathbf{v}_j^{T} \mathbf{v}_k$ 



#### **Cost Function**

Softmax can be used to turn these linear sums  $u_j$  into a probability distribution estimating the probability of word *j* occurring in the context of word *k* 

$$\operatorname{prob}(j|k) = \frac{\exp(u_j)}{\sum_{j'=1}^{V} \exp(u_{j'})} = \frac{\exp(\mathbf{v}_j^{T} \mathbf{v}_k)}{\sum_{j'=1}^{V} \exp(\mathbf{v}_{j'}^{T} \mathbf{v}_k)}$$

We can treat the text as a sequence of numbers  $w_1, w_2, ..., w_T$  where  $w_i = j$  means that the *i*<sup>th</sup> word in the text is the *j*<sup>th</sup> word in the vocabulary. We then seek to maximize the log probability

$$\frac{1}{T} \sum_{t=1}^{T} \sum_{-c \leq r \leq c, r \neq 0} \log \operatorname{prob}(w_{t+r}|w_t)$$

where c is the size of training context (which may depend on  $w_t$ )



#### Word2Vec issues

- Word2Vec is a linear model in the sense that there is no activation function at the hidden nodes
- this 1-word prediction model can be extended to multi-word prediction in two different ways:
  - Continuous Bag of Words
  - Skip-Gram
- need a computationally efficient alternative to Softmax (Why?)
  - Hierarchical Softmax
  - Negative Sampling
- · need to sample frequent words less often



#### Word2Vec Weight Updates

If we assume the full softmax, and the correct output is  $j^*$ , then the cost function is

$$E = -u_{j^*} + \log \sum_{j'=1}^V \exp(u_{j'})$$

the output differentials are

$$e_j = \frac{\partial E}{\partial u_j} = -\delta_{jj^*} + \frac{\partial}{\partial u_j} \log \sum_{j'=1}^{V} \exp(u_{j'})$$

where

$$\delta_{jj^*} = \begin{cases} 1, & \text{ if } j = j^*, \\ 0, & \text{ otherwise.} \end{cases}$$



## Word2Vec Weight Updates

hidden-to-output differentials

$$\frac{\partial E}{\partial w'_{ij}} = \frac{\partial E}{\partial u_j} \frac{\partial u_j}{\partial w'_{ij}} = \mathbf{e}_j \, \mathbf{h}_i$$

hidden unit differentials

$$\frac{\partial E}{\partial h_i} = \sum_{j=1}^{V} \frac{\partial E}{\partial u_j} \frac{\partial u_j}{\partial h_i} = \sum_{j=1}^{V} e_j w_{ij}'$$

input-to-hidden differentials

$$\frac{\partial E}{\partial w_{ki}} = \frac{\partial E}{\partial h_i} \frac{\partial h_i}{\partial w_{ki}} = \sum_{j=1}^{V} e_j w_{ij}' x_k$$



#### **CBOW vs Skip-Gram**

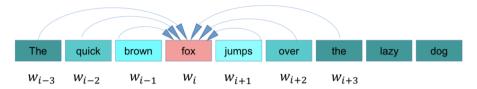


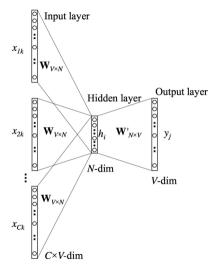
Figure: Continous Bag of Words (CBOW)



Figure: Skip-Gram model



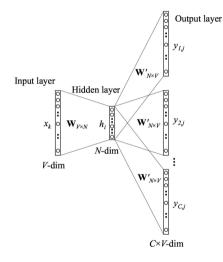
# **Continuous Bag of Words (CBOW)**



- If several context words are each used independently to predict the center word, the hidden activation becomes a sum (or average) over all the context words
- Note the difference between this and NetTalk – in word2vec (CBOW) all context words share the same input-to-hidden weights



# Word2Vec Skip-Gram Model

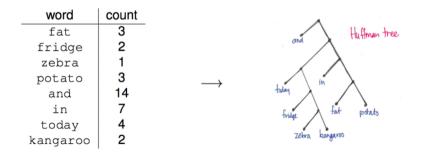


- try to predict the context words, given the center word
- this skip-gram model is similar to CBOW, except that in this case a single input word is used to predict multiple context words
- all context words share the same hidden-to-output weights



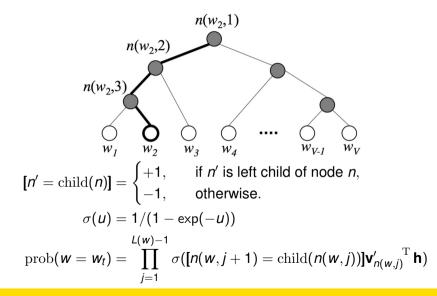
# **Hierarchical Softmax**

- target words are organized in a Huffman-coded Binary Tree
- · each output of the network corresponds to one branch point in the tree
- only those nodes that are visited along the path to the target word are evaluated (which is log<sub>2</sub>(V) nodes on average)





#### **Hierarchical Softmax**





# **Negative Sampling**

- Consider for each word *w* binary classifier: if given word *C* is good context for *w*, or not
- The idea of negative sampling is that we train the network to increase its estimation of the target word  $j^*$  and reduce its estimate not of all the words in the vocabulary but just a subset of them  $W_{neg}$ , drawn from an appropriate distribution.

$$oldsymbol{\mathcal{E}} = -\log \sigma({oldsymbol{v}'_{j^*}}^{\mathrm{T}}oldsymbol{\mathsf{h}}) \ - \sum_{j\in\mathcal{W}_{\mathrm{neg}}}\log \sigma(-{oldsymbol{v}'_{j}}^{\mathrm{T}}oldsymbol{\mathsf{h}})$$

• This is a simplified version of Noise Constrastive Estimation (NCE). It is not guaranteed to produce a well-defined probability distribution, but in practice it does produce high-quality word embeddings.



# **Negative Sampling**

- The number of samples is 5-20 for small datasets, 2-5 for large datasets.
- Empirically, a good choice of the distribution from which to draw the negative samples is  $P(w) = U(w)^{3/4}/Z$  where U(w) is the unigram distribution determined by the previous word, and Z is a normalizing constant.



## Sub-sampling of Frequent Words

In order to diminish the influence of more frequent words, each word in the corpus is discarded with probability

$$\mathsf{P}(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

where  $f(w_i)$  is the frequency of word  $w_i$  and  $t \sim 10^{-5}$  is an empirically determined threshold.



#### **Global Vectors (GloVe)**

#### **Co-occurrence probabilities**

Given two words i and j that occur in text, their co-occurrence probability is defined as the probability of seeing i in the context of j

$$P(j/i) = \frac{\text{count}(j \text{ in context of } i)}{\sum_{k}(\text{count}(k \text{ in context if } i))}$$

**Claim**: If we want to distinguish between two words, it is not enough to look at their co-occurrences, we need to look at the ratio of their co-occurrences with other words.

- Formalizing this intuition gives us an optimization problem



## **GloVe Objective**

#### Notation:

- *i*: word, *j*: a context word
- w<sub>i</sub>: The word embedding for i
- $c_j$ : The context embedding for j
- $b_i^w$ ,  $b_i^c$ : Two bias terms: word and context specific
- $X_{ij}$ : The number of times word *i* occurs in the context of *j*

#### Intuition:

- 1. Construct a word-content matrix whose  $(i, j)^{th}$  entry is  $log(X_{ij})$
- 2. Find vectors  $w_i$ ,  $c_j$  and the biases  $b_i$ ,  $c_j$  such that the dot product of the vectors added to the biases approximates the matrix entries



### **GloVe Objective**

#### Notation:

- *i*: word, *j*: a context word
- w<sub>i</sub>: The word embedding for i
- $c_j$ : The context embedding for j
- $b_i^w$ ,  $b_i^c$ : Two bias terms: word and context specific
- $X_{ij}$ : The number of times word *i* occurs in the context of *j*

Objective:

$$J = \sum_{i,j=1}^{|V|} (w_i^T c_j + b_i + b_j - \log X_{ij})^2$$

**Problem**: Pairs that frequently co-occur tend to dominate the objective **Solution**: Correct for this by adding an extra term that prevents this



## **GloVe Objective**

#### Notation:

- *i*: word, *j*: a context word
- w<sub>i</sub>: The word embedding for i
- $c_j$ : The context embedding for j
- $b_i^w$ ,  $b_i^c$ : Two bias terms: word and context specific
- $X_{ij}$ : The number of times word *i* occurs in the context of *j*

Objective:

$$J = \sum_{i,j=1}^{|V|} f(X_{ij}) (w_i^T c_j + b_i + b_j - \log X_{ij})^2$$

*f*: A weighting function that assigns lower relative importance to frequent co-occurrences



#### GloVe

Table 2: Results on the word analogy task, given as percent accuracy. Underlined scores are best within groups of similarly-sized models; bold scores are best overall. HPCA vectors are publicly available<sup>2</sup>; (i)vLBL results are from (Mnih et al., 2013); skip-gram (SG) and CBOW results are from (Mikolov et al., 2013a,b); we trained SG<sup>†</sup> and CBOW<sup>†</sup> using the word2vec tool<sup>3</sup>. See text for details and a description of the SVD models.

$\begin{array}{c c c c c c c c c c c c c c c c c c c $						
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	Model	Dim.	Size	Sem.	Syn.	Tot.
Glove         100         1.6B         67.5         54.3         60.3           SG         300         1.B         61         61         61           CBOW         300         1.6B         16.1         52.6         36.1           VLBL         300         1.5B         54.2         64.8         60.0           ivLBL         300         1.5B         65.2         63.0         64.0           GloVe         300         1.6B         80.8         61.5         70.3           SVD         300         6B         63.         8.1.         7.3           SVD-S         300         6B         63.6         67.4         65.7           SVD-S         300         6B         56.6         63.0         60.1           CBOW <sup>+</sup> 300         6B         63.6         67.4         65.7           SVD-S         300         6B         73.0         60.0         69.1           GBOW <sup>+</sup> 300         6B         73.0         60.0         69.1           Glove         300         6B         73.4         67.0         71.7           CBOW <sup>+</sup> 300         6B         73.3 <t< td=""><td>ivLBL</td><td>100</td><td>1.5B</td><td>55.9</td><td>50.1</td><td>53.2</td></t<>	ivLBL	100	1.5B	55.9	50.1	53.2
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	HPCA	100	1.6B	4.2	16.4	10.8
CBOW         300         1.6B         16.1         52.6         36.1           VLBL         300         1.5B         54.2         64.8         60.0           viLBL         300         1.5B         54.2         64.8         60.0           GloVe         300         1.5B         65.2         63.0         64.0           GloVe         300         1.6B         80.8         61.5         70.3           SVD         300         6B         6.3         81.7         7.3           SVD-L         300         6B         63.6         63.0         60.1           CBOW <sup>1</sup> 300         6B         63.6         67.4         66.6         50.1           GloVe         300         6B         63.6         67.4         65.0         50.1         61.0         50.1           GloVe         300         6B         63.6         67.4         65.1         61.6         63.1         61.1         51.4         50.4         53.1         53.0         53.7         56.9         53.7         56.9         53.7         56.9         53.7         55.1         65.0         53.0         53.2         49.2         49.2         49.2	GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
vLBL         300         1.5B         54.2         64.8         60.0           ivLBL         300         1.5B         65.2         63.0         64.0           GloVe         300         1.6B         80.8         61.5         70.3           SVD         300         6B         63.         8.1         7.3           SVD-S         300         6B         63.         8.1         7.3           SVD-S         300         6B         63.6         67.4         65.7           SVD-S         300         6B         63.6         67.4         65.7           SVD-S         300         6B         73.0         66.0         69.1           GBW <sup>+</sup> 300         6B         73.0         66.0         69.1           GloVe         300         6B         77.4         67.0         71.7           CBOW <sup>+</sup> 300         6B         75.3         68.9         63.7           GloVe         300         6B         75.3         68.9         63.7           SG         1000         6B         66.1         65.1         65.6           SVD-L         300         42B         38.4         5	SG	300	1B	61	61	61
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	CBOW	300	1.6B	16.1	52.6	36.1
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	ivLBL	300	1.5B	65.2	63.0	64.0
SVD-S         300         6B         36.7         46.6         42.1           SVD-L         300         6B         56.6         63.0         60.1           CBOW <sup>†</sup> 300         6B         63.6         67.4         65.7           SG <sup>†</sup> 300         6B         73.0         66.0         69.1           GloVe         300         6B         77.4         67.0         71.7           CBOW         1000         6B         57.3         68.9         63.7           SG         1000         6B         56.1         65.1         65.5           SVD-L         300         42B         38.4         58.2         49.2	GloVe	300	1.6B	80.8	61.5	<u>70.3</u>
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	SVD	300	6B	6.3	8.1	7.3
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	SVD-S	300	6B	36.7	46.6	42.1
SG <sup>†</sup> 300         6B         73.0         66.0         69.1           GloVe         300         6B         77.4         67.0         71.7           CBOW         1000         6B         57.3         68.9         65.1         55.5           SG         1000         6B         66.1         65.1         65.6         5VD-L         300         42B         38.4         58.2         49.2	SVD-L	300	6B	56.6	63.0	60.1
Glove         300         6B         77.4         67.0         71.7           CBOW         1000         6B         57.3         68.9         63.7           SG         1000         6B         66.1         65.1         65.6           SVD-L         300         42B         38.4         58.2         49.2	$CBOW^{\dagger}$	300	6B	63.6	<u>67.4</u>	65.7
CBOW         1000         6B         57.3         68.9         63.7           SG         1000         6B         66.1         65.1         65.6           SVD-L         300         42B         38.4         58.2         49.2	$SG^{\dagger}$	300	6B	73.0	66.0	69.1
SG         1000         6B         66.1         65.1         65.6           SVD-L         300         42B         38.4         58.2         49.2	GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>
SVD-L 300 42B 38.4 58.2 49.2	CBOW	1000	6B	57.3	68.9	63.7
	SG	1000	6B	66.1	65.1	65.6
GloVe 300 42B 81.9 69.3 75.0	SVD-L	300	42B	38.4	58.2	49.2
	GloVe	300	42B	<u>81.9</u>	<u>69.3</u>	<u>75.0</u>

Table 3: Spearman rank correlation on word similarity tasks. All vectors are 300-dimensional. The CBOW\* vectors are from the word2vec website and differ in that they contain phrase vectors.

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	<u>72.7</u>	75.1	56.5	37.0
CBOW <sup>†</sup>	6B	57.2	65.6	68.2	57.0	32.5
SG <sup>†</sup>	6B	62.8	65.2	69.7	58.1	37.2
GloVe	6B	<u>65.8</u>	<u>72.7</u>	<u>77.8</u>	53.9	<u>38.1</u>
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	75.9	<u>83.6</u>	<u>82.9</u>	<u>59.6</u>	<b>47.8</b>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

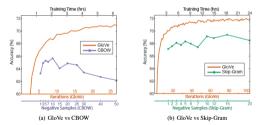


Figure 4: Overall accuracy on the word analogy task as a function of training time, which is governed by the number of iterations for GloVe and by the number of negative samples for CBOW (a) and skip-gram (b). In all cases, we train 300-dimensional vectors on the same 6B token corpus (Wikipedia 2014 + Gigaword 5) with the same 400,000 word vocabulary, and use a symmetric context window of size 10.



### Sentence Completion Task

Q1. Seeing the pictures of our old home made me feel ...... and

- nostalgic. A. fastidious
  - B. indignant
  - C. wistful
  - D. conciliatory
- - A. object
  - B. abdicate
  - C. abstain
  - D. compromise

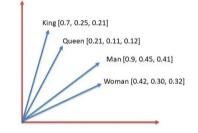
(use model to choose which word is most likely to occur in this context)



## **Linguistic Regularities**

King + Woman - Man  $\simeq$  Queen

More generally, A is to B as C is to ??



$$\boldsymbol{d} = \operatorname{argmax}_{\boldsymbol{x}} \frac{(\boldsymbol{v}_{\mathcal{C}} + \boldsymbol{v}_{\boldsymbol{b}} - \boldsymbol{v}_{\boldsymbol{a}})^{\mathrm{T}} \boldsymbol{v}_{\boldsymbol{x}}}{||\boldsymbol{v}_{\mathcal{C}} + \boldsymbol{v}_{\boldsymbol{b}} - \boldsymbol{v}_{\boldsymbol{a}}||}$$

King – Man + Woman = [0.7, 0.25, 0.21] - [0.9, 0.45, 0.41] + [0.42, 0.30, 0.32] = [0.22, 0.1, 0.12] = Queen



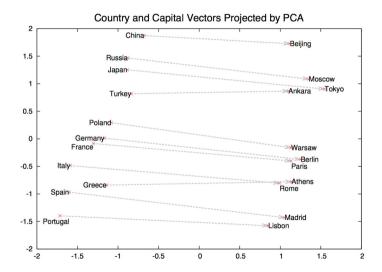
### Word Analogy Task

#### Q1. evening is to morning as dinner is to ..........

- A. breakfast
- B. soup
- C. coffee
- D. time
- Q2. bow is to arrow as ..... is to bullet
  - A. defend
  - B. lead
  - C. shoot
  - D. gun



### **Capital Cities**





# Word Analogies

Type of relationship	Word Pair 1		Word Pair 2	
Common capital city	Athens	Greece	Oslo	Norway
All capital cities	Astana	Kazakhstan	Harare	Zimbabwe
Currency	Angola	kwanza	Iran	rial
City-in-state	Chicago	Illinois	Stockton	California
Man-Woman	brother	sister	grandson	granddaughter
Adjective to adverb	apparent	apparently	rapid	rapidly
Opposite	possibly	impossibly	ethical	unethical
Comparative	great	greater	tough	tougher
Superlative	easy	easiest	lucky	luckiest
Present Participle	think	thinking	read	reading
Nationality adjective	Switzerland	Swiss	Cambodia	Cambodian
Past tense	walking	walked	swimming	swam
Plural nouns	mouse	mice	dollar	dollars
Plural verbs	work	works	speak	speaks



## Word Relationship

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza



## Summary

- *Word vectors*, also sometimes called *word embeddings* or *word representations* are distributed representations of words.
- Two kinds of embeddings
  - Sparse vectors: Words are represented by a simple function of the counts of nearby words.
  - Dense vectors: Representation is created by training a classifier to distinguish nearby and far-away words
- The *contexts* in which a word appears tells us a lot about what it means. Distributional similarities use the set of contexts in which words appear to measure their similarity
- Word2Vec and GloVe are two important dense representations of words.
- Various choice:
  - $\circ~$  dimensionality of embeddings (50, 100, 200, 300, 500)
  - $\circ\;$  scale, quality and type of text to get word embeddings
  - size of the context window



#### References

- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. In Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2 (NIPS'13). Curran Associates Inc., Red Hook, NY, USA, 3111–3119.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global Vectors for Word Representation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Xin Rong. 2014. word2vec Parameter Learning Explained, arXiv link.
- GloVe https://nlp.stanford.edu/projects/glove/
- On Word Embeddings
- Word Embedding Demo: Tutorial



#### References

- Bandyopadhyay, S., Xu, J., Pawar, N. and Touretzky, D. 2022. Interactive Visualizations of Word Embeddings for K-12 Students. Proceedings of the AAAI Conference on Artificial Intelligence. 36, 11 (Jun. 2022), 12713-12720. DOI:https://doi.org/10.1609/aaai.v36i11.21548.
- TensorFlow Word Embeddings
- Word Embeddings in PyTorch
- Word2Vec Demo
- Word Embedding Demo

