

Word Embeddings

Instructor: Thien Huu Nguyen

Based on slides from: Chris Manning



Words

- The primary elements of natural languages
- Each word carries some unit meaning depending on its context
- The unit meanings of the words are composed/combined to produce new and more complicated meanings/concepts (e.g., sentences, documents)



Word Meanings

- The fundamental of NLP is to be able to allow computers to understand meanings of text

$$\textit{Meaning}(\text{"I have a cat"}) = f(\textit{Meaning}(\text{"I"}), \\ \textit{Meaning}(\text{"have"}), \\ \textit{Meaning}(\text{"a"}), \\ \textit{Meaning}(\text{"cat"}))$$

How do we capture/approximate the composition function f and the *Meaning* function for words?

We will discuss the word meanings in this talk!



What are meanings?

Definition (Webster dictionary)

- The idea that is represented by a word, phrase, etc.
- The idea that a person wants to express by using words, signs, etc.
- The idea that is expressed in a word of writing, art, etc.



How to represent the meanings of a word in computers?

Common solution: Use the sets of synonyms and hypernyms of the word by querying some thesaurus (e.g., WordNet)

e.g. synonym sets containing "good":

```
from nltk.corpus import wordnet as wn
poses = { 'n': 'noun', 'v': 'verb', 's': 'adj (s)', 'a': 'adj', 'r': 'adv' }
for synset in wn.synsets("good"):
    print("{}: {}".format(poses[synset.pos()],
                          ", ".join([l.name() for l in synset.lemmas()])))
```

```
noun: good
noun: good, goodness
noun: good, goodness
noun: commodity, trade_good, good
adj: good
adj (sat): full, good
adj: good
adj (sat): estimable, good, honorable, respectable
adj (sat): beneficial, good
adj (sat): good
adj (sat): good, just, upright
...
adverb: well, good
adverb: thoroughly, soundly, good
```

e.g. hypernyms of "panda":

```
from nltk.corpus import wordnet as wn
panda = wn.synset("panda.n.01")
hyper = lambda s: s.hypernyms()
list(panda.closure(hyper))
```

```
[Synset('procyonid.n.01'),
Synset('carnivore.n.01'),
Synset('placental.n.01'),
Synset('mammal.n.01'),
Synset('vertebrate.n.01'),
Synset('chordate.n.01'),
Synset('animal.n.01'),
Synset('organism.n.01'),
Synset('living_thing.n.01'),
Synset('whole.n.02'),
Synset('object.n.01'),
Synset('physical_entity.n.01'),
Synset('entity.n.01')]
```



Problems with resources like WordNet

- Great as a resource but missing nuance
 - e.g., “proficient” is listed as a synonym for “good”, but this is only true in some contexts.
- Missing new meanings of words
 - e.g., wicked, badass, nifty, wizard, genius, ninja, bombast
 - very challenging to keep up-to-date.
- Subjective
- Require human labor to create and adapt
- Impossible to compute word similarity



Representing words as discrete symbols

- In traditional NLP, words are considered as discrete symbols
- Mathematically, a words are represented by a one-hot vector, where:
 - The dimension of the vector = the number of words in some given vocabulary (e.g., 500,000)
 - Only the bit corresponding to the word is set to 1 (i.e, 0 otherwise)

hotel = [0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0]

motel = [0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]

- This is called the localist representation (to be distinguished with distributed representation in cognitive science later)



Problems with words as discrete symbols

- The size of the vectors is large
- The vectors for any pair of words are orthogonal (i.e., cosine similarity = 0), but for similar words like “hotel” and “motel”, we expect their vectors to exhibit some level of similarity (i.e., the cosine similarity should be non-zero).
 - e.g., in web search, a search for “Seattle hotel” should return documents with “Seattle motel” as well.
- Solution for this?
 - Can we use the idea of synonyms and hyponyms for such one-hot vectors?
 - Not working well in practice (e.g., incompleteness)
 - Learn to explicitly encode similarity in the word vectors themselves, reduce the size of the vectors, go from binary vectors to continuous vectors



Representing words by their contexts

- *Distributional semantics*: a word's meaning is given by the words that frequently appear close-by
 - “You shall know a word by the company it keeps” (J. R. Firth 1957: 11)
 - One of the most successful ideas of modern statistical NLP
- When a word w appears in a text, its **context** is the set of words that appear nearby (within a fixed-size window).
- Use the many contexts of w to build up a representation of w



...government debt problems turning into **banking** crises as happened in 2009...

...saying that Europe needs unified **banking** regulation to replace the hodgepodge...

...India has just given its **banking** system a shot in the arm...



Word vectors

- We will introduce a dense vector for each word, chosen so that it is similar to vectors of words appearing in similar contexts.
- Word vectors are also called word embeddings or word representations. They are a distributed representation, e.g.,

banking =

$$\begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$

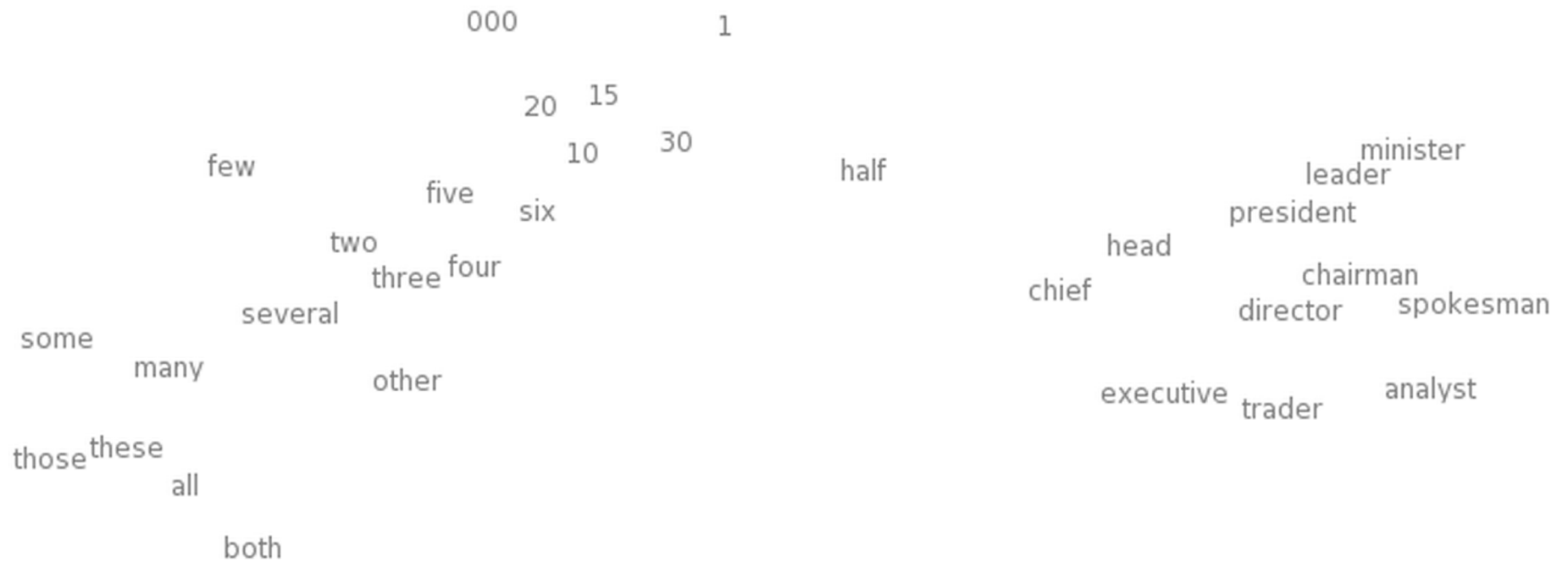

Localist representation vs. distributed representation

- In cognitive science, distributed representation has the following property (Hilton et al., 1986; Plate, 2012):
 - A concept is represented by a pattern of activity over a collection of neurons (i.e., more than one neuron is required to represent a concept.)
 - Each neuron participates in the representation of more than one concept.
- By contrast, in localist representation, each neuron represents a single concept on a stand-alone basis. The critical distinction is that localist units have “meaning and interpretation” whereas units in distributed representation don’t.
 - “These representations are distributed, which typically has the consequence that interpretable information cannot be obtained by examining activity of single hidden units.” – Elman, 1995.

Roy, Asim. “A theory of the brain: localist representation is used widely in the brain.” *Frontiers in psychology* vol. 3 551



Word meaning as a neural word vector



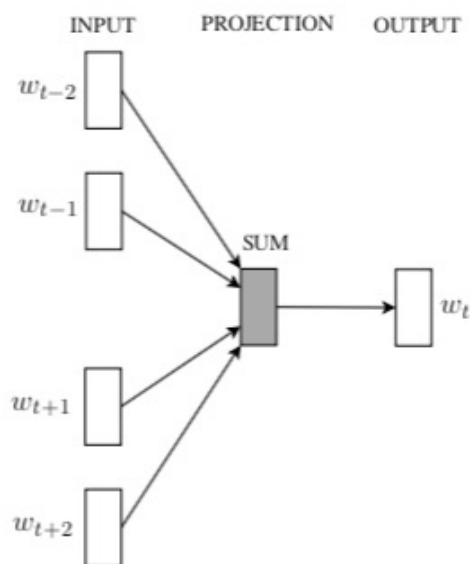
How do we obtain such word vectors?

- Word2vec (Mikolove et al. 2013) is a popular framework to learn word vectors (although many other efforts have been made before it)
- Idea:
 - We start with a large corpus of text
 - Every word in a fixed vocabulary is represented by a **vector**
 - Go through each position t in the text, which has a center word c and context words o (surrounding words)
 - Use the similarity of the word vectors for c and o to compute the probability of c given o ($P(c|o)$) (or vice versa)
 - Keep updating the word vectors to maximize this probability

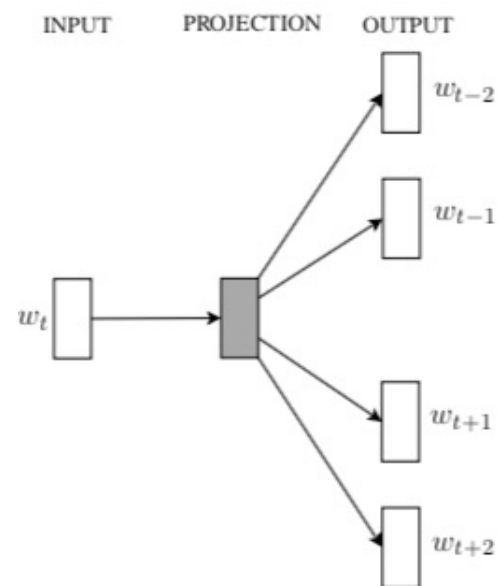


Two variants of word2vec

Context words: windows of size 2 before and after the center word



Continuous Bag of Words (CBOW):
predicting the center words using
the context words ($P(w_t | w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2})$)

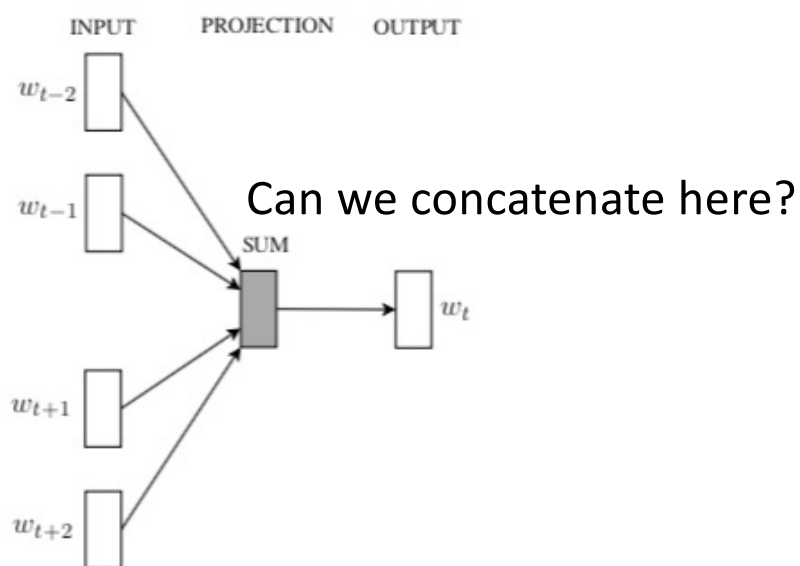


Skip-grams (SG):
predicting the context words using
the center word ($P(w_{t+i} | w_t), i \in \{-2, -1, 1, 2\}$)

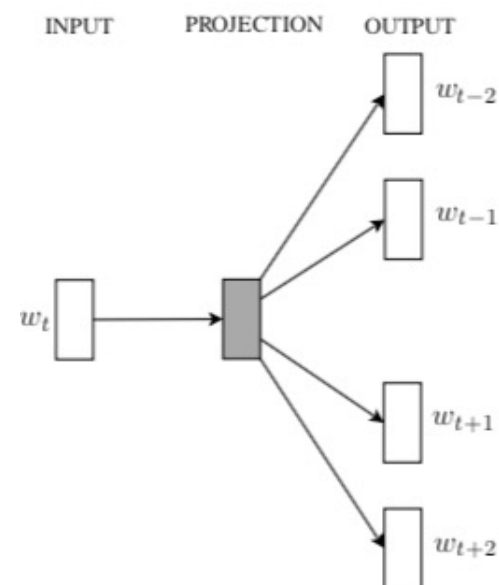


Two variants of word2vec

Context words: windows of size 2 before and after the center word



Continuous Bag of Words (CBOW):
predicting the center words using
the context words ($P(w_t | w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2})$)



Skip-grams (SG):
predicting the context words using
the center word ($P(w_{t+i} | w_t), i \in \{-2, -1, 1, 2\}$)



Wovd2vec: SG objective function

- For each position $i = 1, \dots, N$, predict the context words within a window of fixed size m , given the the center word w_i :

$$\text{Likelihood} = L(\theta) = \prod_{i=1}^N \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{i+j} | w_i; \theta)$$

- The objective/loss function is the (average) negative log likelihood:

$$\text{loss} = J(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{i+j} | w_i; \theta)$$

- θ is the parameter used to define $P(w_{i+j} | w_i; \theta)$. It is the **model parameters**
- Minimizing the loss function amounts to maximizing the predictive accuracy



Wovd2vec: SG objective function

- How do we compute $P(w_{i+j}|w_i; \theta)$?
- We will use two vectors per word w :
 - v_w when w is a center word
 - u_w when w is a context word
 - Using two vectors makes the later optimization easier, average both at the end to obtain final word vectors
 - Although using one vector per word is possible too

Dot product measures similarity of two vectors

$$u^T v = u \cdot v = \sum u_i v_i$$

- Then:

$$P(w_{i+j}|w_i; \theta) = \frac{\exp(u_{w_{i+j}}^T v_{w_i})}{\sum_{w \in V} \exp(u_w^T v_{w_i})}$$



Wovd2vec: SG objective function

- We compute the probability using the softmax function that maps a vector $x = [x_1, \dots, x_n]$ of arbitrary values into a probability distribution:

Exponentiation makes anything positive

Normalization over the dimensions to produce a probability distribution

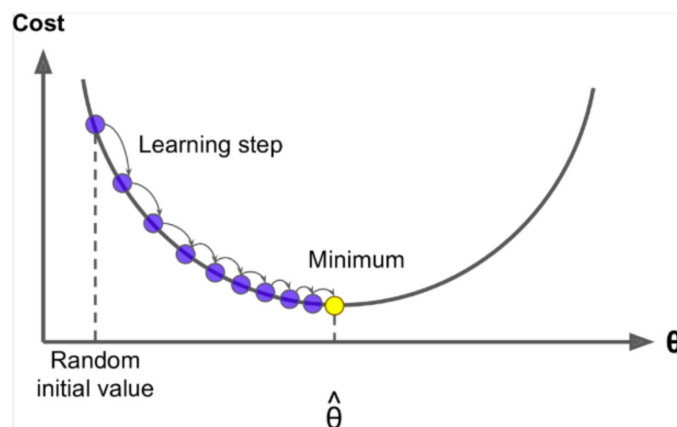
$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)}$$

- “max” because it amplifies the probability for the largest x_i
- “soft” because it still assigns some probability to the smaller x_i
- frequently used in deep learning



Optimization: Gradient Descent

- So, we have a loss function $J(\theta)$ with the word vectors as the parameters. We want to find the parameters (word vectors) that can optimize (minimize) this loss/objective function.
- This is an optimization problem, often solved by (stochastic) gradient descent in deep learning
- Idea: for the current value θ , calculate gradient of $J(\theta)$, then take small step in the direction of negative gradient. Repeat until some convergence condition is met.



The loss functions in practice are often more complicated (i.e., non-convex) than this.



Gradient Descent

- Update rule (in matrix notation):

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$

– where α is the **step size** or **learning rate**

- Update rule (for a single parameter):

$$\theta_j^{new} = \theta_j^{old} - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

- Algorithm:

```
while True:  
    theta_grad = evaluate_gradient(J, corpus, theta)  
    theta = theta - alpha * theta_grad
```



Stochastic Gradient Descent

- Problem: $J(\theta)$ is a function of all windows in the corpus (potentially billions!)
 - So, $\nabla_{\theta}J(\theta)$ is **very expensive** to compute
 - We might need to wait for a very long time before making a single update!
 - We want to be able to update the models more frequently.
- In practice, for deep learning we use **Stochastic Gradient Descent**:
 - Repeatedly sample windows, and update the model after each sampling
 - Just an approximation of batch gradient descent, but have the potential ability to escape local minima (good for non-convex functions)
 - Algorithm:

```
while True:  
    window = sample_window(corpus)  
    theta_grad = evaluate_gradient(J, window, theta)  
    theta = theta - alpha * theta_grad
```



Negative Sampling

- $$P(w_{i+j} | w_i; \theta) = \frac{\exp(u_{w_{i+j}}^T v_{w_i})}{\sum_{w \in V} \exp(u_w^T v_{w_i})}$$
- The normalization factor needs to enumerate over all the words in the vocabulary that can be very large!
- We can instead obtain only a sample of the vocabulary to estimate the normalization factor. This is called **Negative Sampling** as every word other than w_{i+j} is considered as negative in this case.

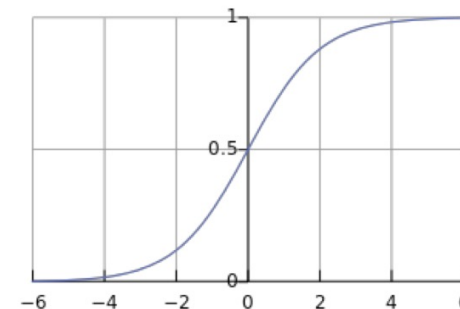


Negative Sampling in the Original Paper

- Paper: “Distributed Representations of Words and Phrases and their Compositionality” (Mikolov et al., 2013).
- Train binary logistic regression for a true pair (a center word and a word in its context window) versus several noise pairs (the center word paired with a random word)
- Overall objection function to maximize:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J_t(\theta)$$

$$J_t(\theta) = \log \sigma(u_o^T v_c) + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-u_j^T v_c)]$$



- The sigmoid function (very popular in machine learning): $\sigma(x) = \frac{1}{1+e^{-x}}$
- In the loss function, we basically maximize the probability of two words co-occurring in the first log

The skip-gram model with negative sampling (implementation)

$$J_{neg-sample}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

- We take K negative samples (using word probabilities)
- Maximize the probability that a real outside word appears and minimize the probability that random words appear around the center word
- $P(w) = U(w)^{3/4} / Z$: the unigram distribution $U(w)$ raised to the $3/4$ power.
- The power makes frequently words be sampled more often



Co-occurrence counts

- Word2Vec capture the co-occurrence of words via the prediction tasks.
- A simpler approach to capture word co-occurrence is via the direct co-occurrence counts between words and X
- Two options for X: words in windows and full documents
 - Window: Counts are done between pairs of words. Similar to Word2Vec, use window around each word -> capturing both syntactic (POS) and semantic information
 - Document: The co-occurrence counts are done between words and documents, encoding the general topics and leading to “Latent Semantic Analysis”

https://en.wikipedia.org/wiki/Latent_semantic_analysis



Example: Window based co-occurrence matrix

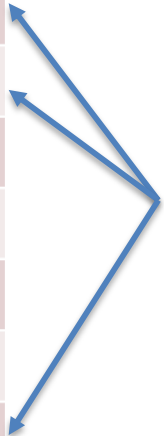
- Window length 1 (although 5-10 are more common)
- Symmetric (don't distinguish left or right context)
- Example corpus:
 - I like deep learning.
 - I like NLP.
 - I enjoy flying.



Example: Window based co-occurrence matrix

- Example corpus:
 - I like deep learning.
 - I like NLP.
 - I enjoy flying.

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0



These are the word vectors!



Problems with simple co-occurrence vectors

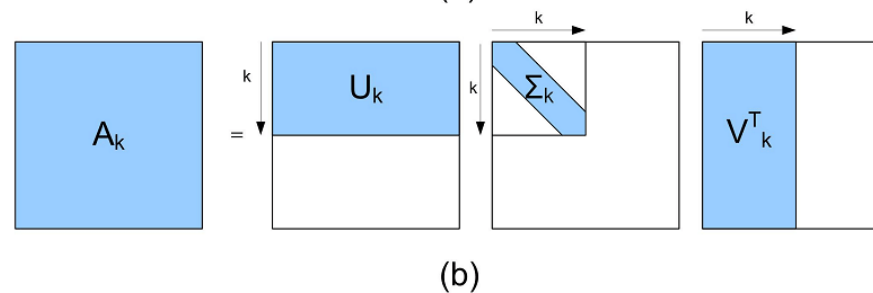
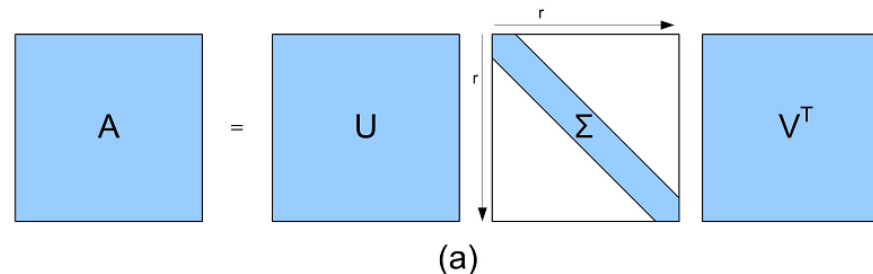
- Increase in size with vocabulary
- Very high dimensional: need a lot of storage
- Subsequent classification models have sparsity issues
- Thus, models are less robust

- Solution: Low dimensional vectors
 - Idea: store most of the important information in a fixed, small number of dimensions: a dense vector
 - Usually 25-1000 dimensions (like Word2Vec)
 - Main question: How to reduce the dimensionality?



Method 1: Dimensionality Reduction

- Singular Value Decomposition (SVD) of the co-occurrence matrix A
- Factorize A into $U\Sigma V^T$ where U and V are orthonormal.



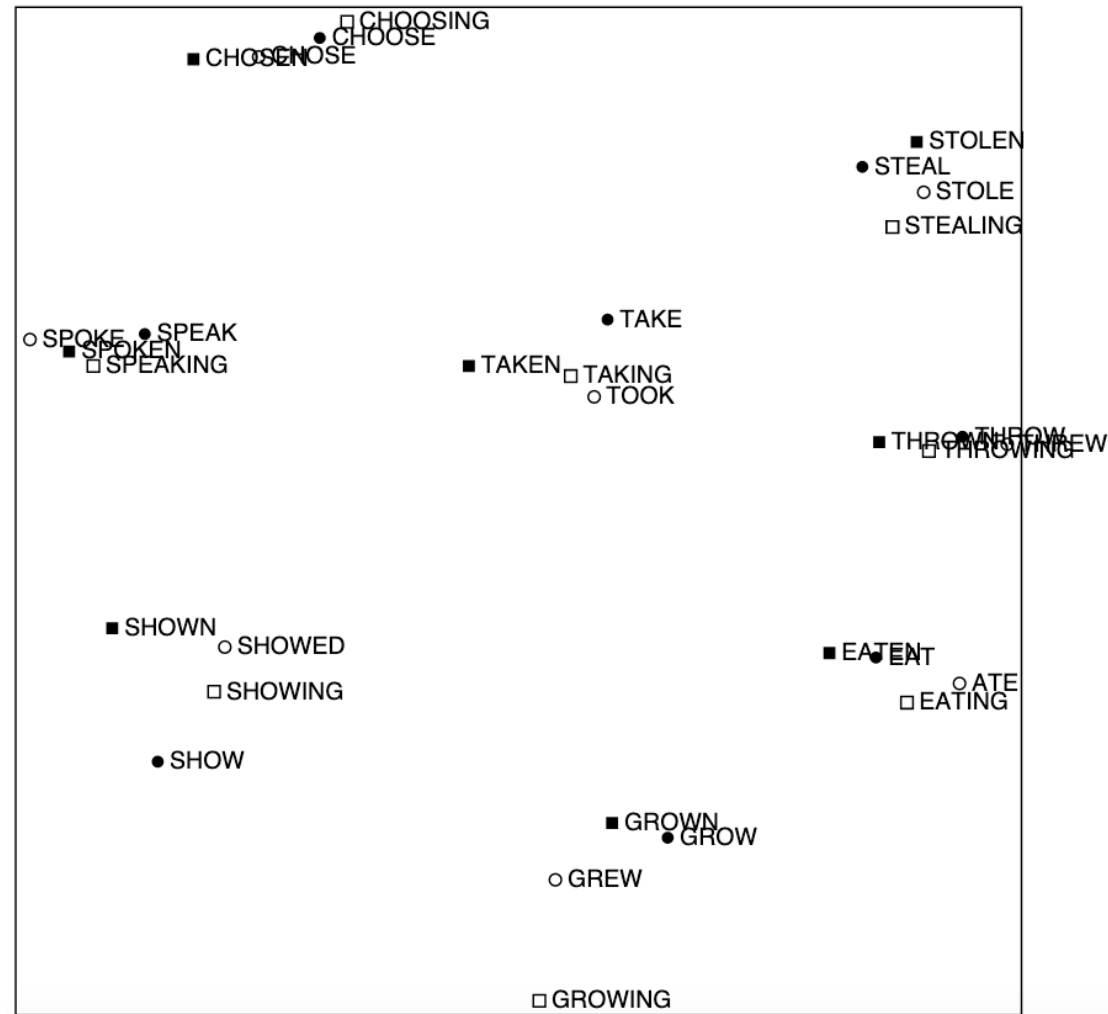
- Retain only k singular values, in order to generalize.
- A_k is the best rank k approximation to A , in terms of least squares.
- Classic linear algebra result. Very expensive to compute for large matrices.

Some tricks for dimensionality reduction

- Scaling the counts in the cells of A can help a lot
 - Problem: function words (*the, he, has*) are too frequent, so syntax has too much impact. Some fixes:
 - $\text{Min}(A, t)$ with $t \approx 100$
 - Ignore them all
- Use Pearson correlations instead of counts, then set negative values to 0
- ...



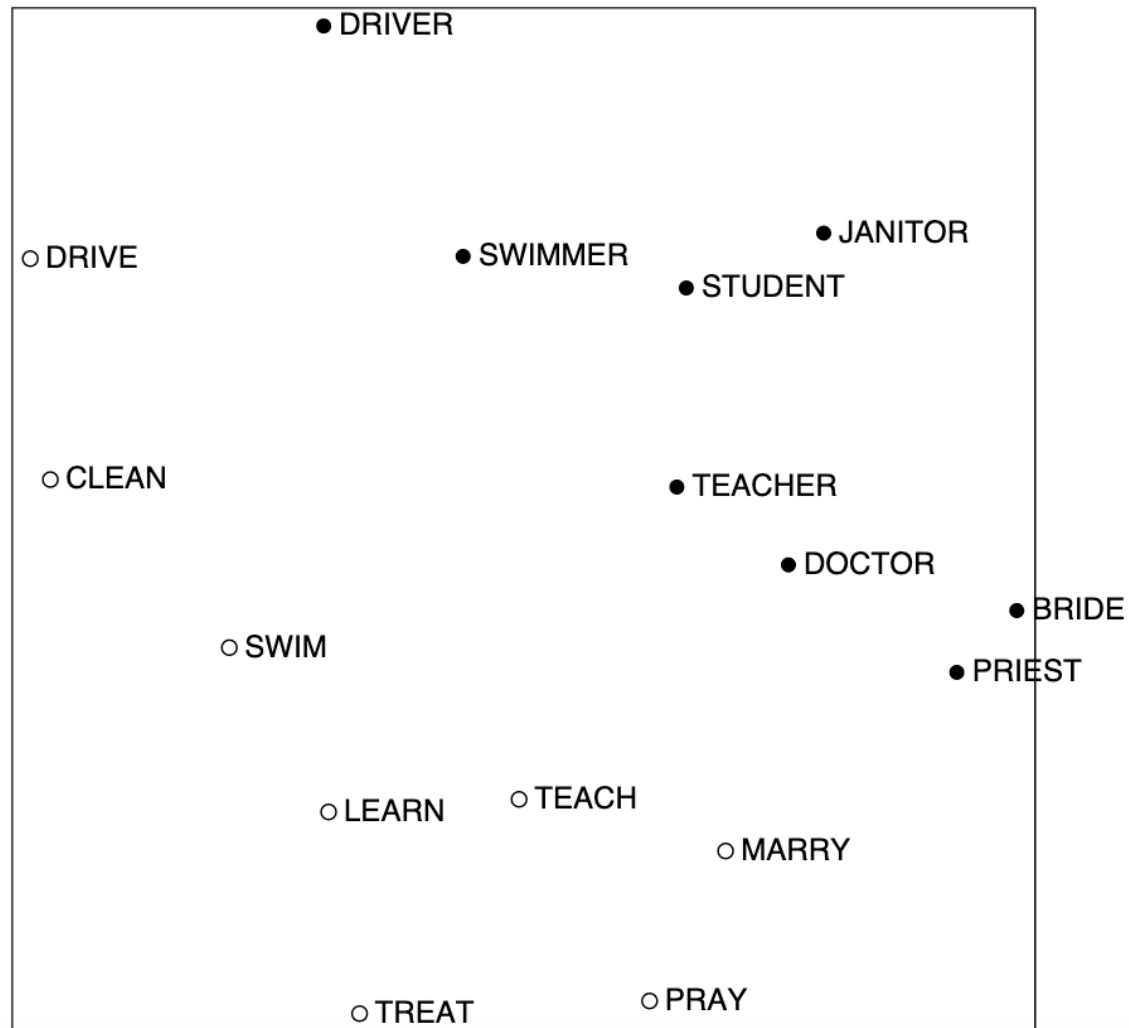
Interesting syntactic patterns emerging in word vectors



COALS model from: [An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence \(Rohde et al., 2005\)](#)



Interesting semantic patterns emerging in word vectors



COALS model from: [An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence \(Rohde et al., 2005\)](#)



Count based vs. direct prediction

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

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

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

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



Method 2: GloVe (Pennington et al., EMNLP 2014)

Encoding meaning in vector differences

- Crucial insight:** Ratios of co-occurrence probabilities can encode meaning components (i.e., relationships of words)

Probe words

	$x = \text{solid}$	$x = \text{gas}$	$x = \text{water}$	$x = \text{random}$
$P(x \text{ice})$	large	small	large	small
$P(x \text{steam})$	small	large	large	small
$\frac{P(x \text{ice})}{P(x \text{steam})}$	large	small	~ 1	~ 1



Method 2: GloVe (Pennington et al., EMNLP 2014)

Encoding meaning in vector differences

- Crucial insight:** Ratios of co-occurrence probabilities can encode meaning components (i.e., relationships of words)

Probe words

	$x = \text{solid}$	$x = \text{gas}$	$x = \text{water}$	$x = \text{fashion}$
$P(x \text{ice})$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(x \text{steam})$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$\frac{P(x \text{ice})}{P(x \text{steam})}$	8.9	8.5×10^{-2}	1.36	0.96



Encoding meaning in vector differences (GloVe)

- **Question:** How can we capture ratios of co-occurrence probabilities as linear meaning components in a word vector space?
- **Answer:**

– Log-bilinear model $w_i \cdot w_j = \log P(i|j)$

with vector differences $w_x \cdot (w_a - w_b) = \log \frac{P(x|a)}{P(x|b)}$



Encoding meaning in vector differences (GloVe)

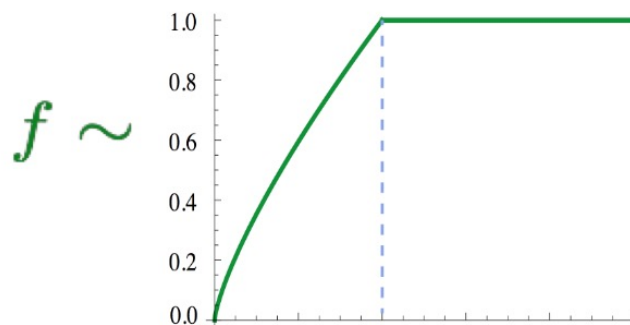
- The loss function:

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

The weighting function

Number of times word j occur
in the context of word i

- Fast training
- Scalable to huge corpora
- Good performance even with small corpus and small vectors
- Word2Vec and GloVe were very popular in NLP before the advent of large language models. Which one is better depends on their specific applications.



GloVe results

- Nearest words to frog:
 - frogs
 - toad
 - litoria
 - leptodactylidae
 - rana
 - lizard
 - eleutherodactylus



litoria



leptodactylidae



rana



eleutherodactylus

How to evaluate word vectors?

- Related to general evaluation in NLP: Intrinsic vs extrinsic
- **Intrinsic:**
 - Evaluation on a specific/intermediate subtask
 - Fast to compute
 - Helps to understand that system
 - Unclear if really helpful unless correlation to real tasks is established
- **Extrinsic:**
 - Evaluation on a real task (things that we will study in this class)
 - Can take a long time to evaluate the accuracy
 - If a problem occurs, unclear if it is due to the word vectors the system for the real task or their interactions
 - If replacing exactly one system for the real task with another improves accuracy -> Winning!



Intrinsic word vector evaluation

- Word Vector Analogies

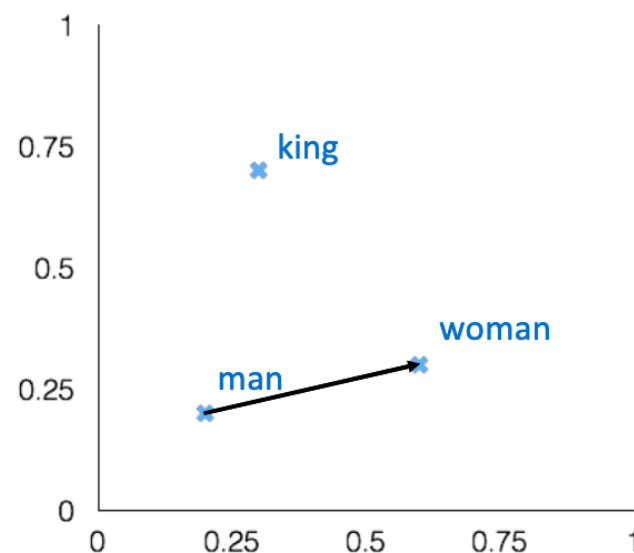
$a:b :: c:?$



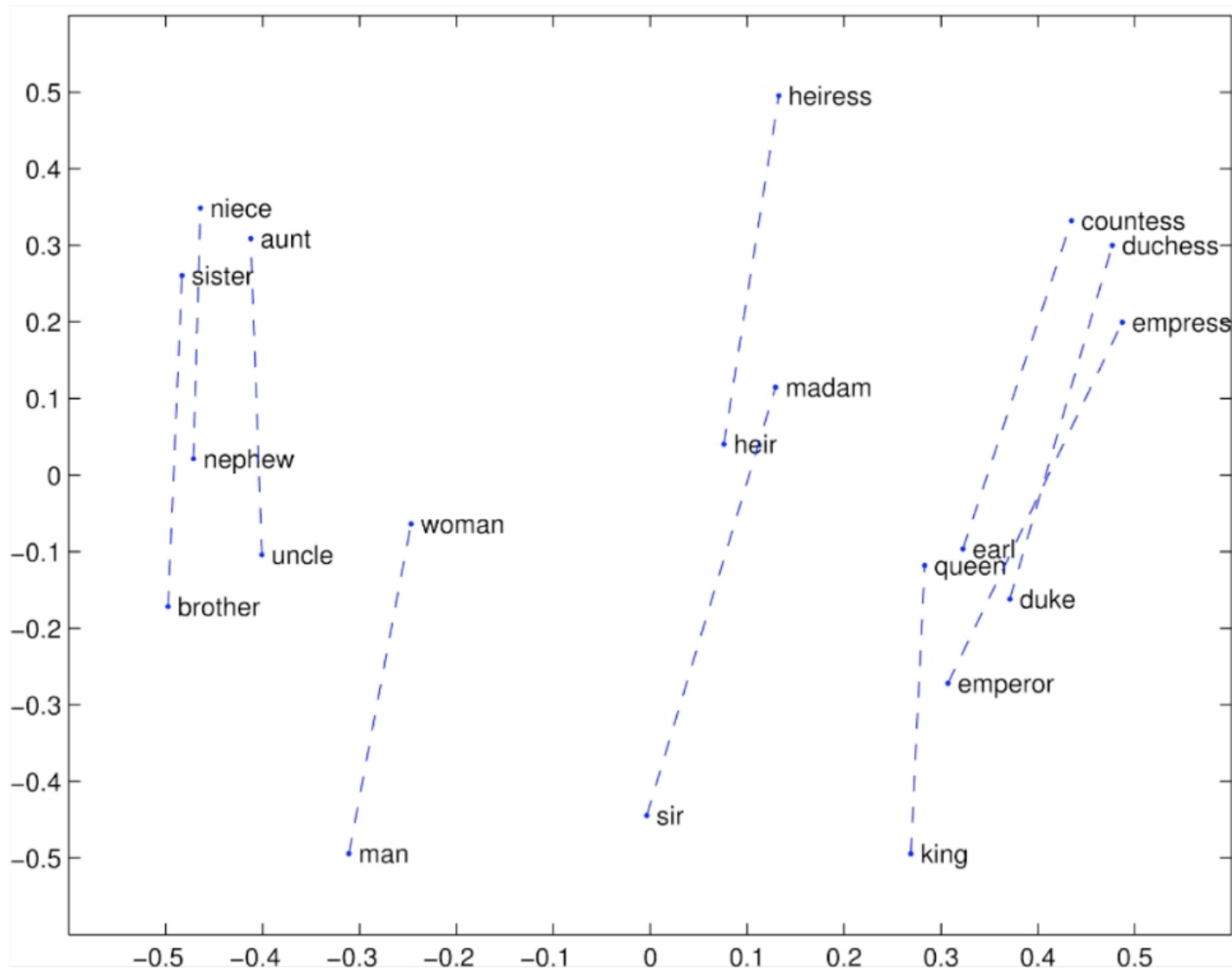
$$d = \arg \max_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|}$$

man:woman :: king:?

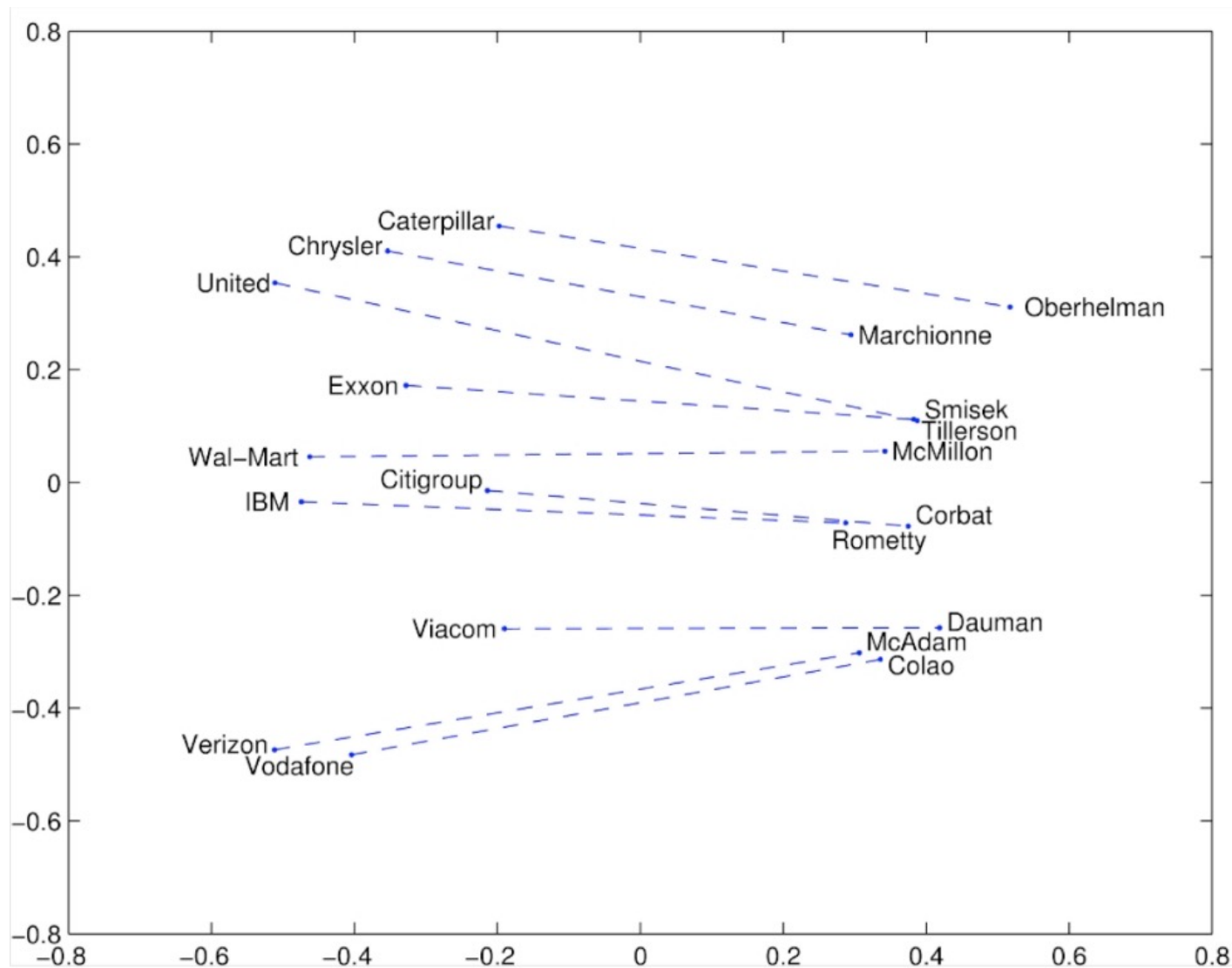
- Evaluate word vectors by how well their cosine distance after addition captures intuitive semantic and syntactic analogy questions
- Discarding the input words from the search! $\arg \max_{x \in V \setminus \{king, man, woman\}} \cos(x, king - man + woman)$
- Problem: What if the information is there but not linear?



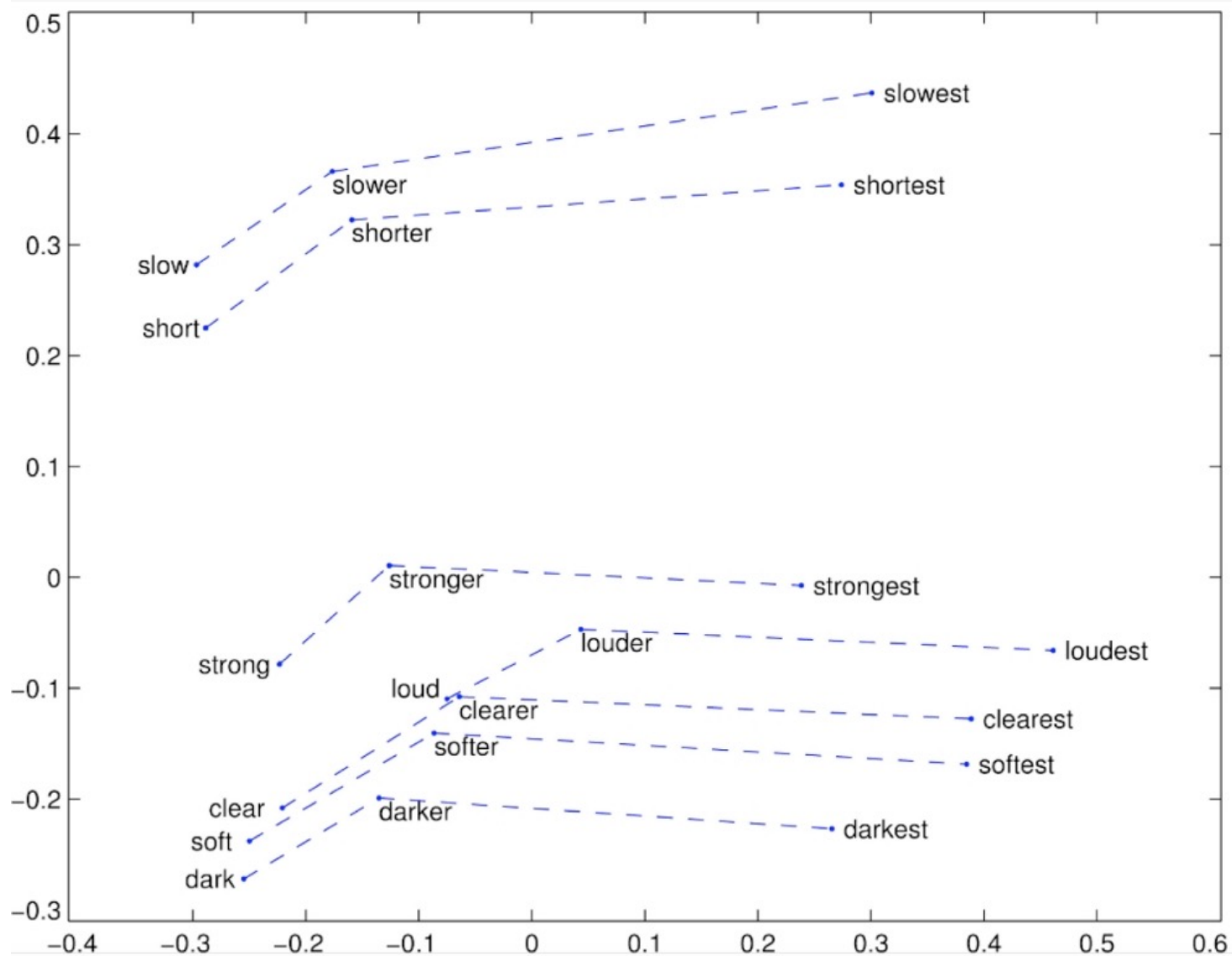
GloVe visualization



GloVe visualization: Company - CEO



GloVe visualization: Superlative



Intrinsic word vector evaluation

- Word Vector Analogies: Syntactic and Semantic examples from:
<https://code.google.com/archive/p/word2vec/source>

: gram4-superlative

bad worst big biggest

bad worst bright brightest

bad worst cold coldest

bad worst cool coolest

bad worst dark darkest

bad worst easy easiest

bad worst fast fastest

bad worst good best

bad worst great greatest

: city-in-state

Chicago Illinois Houston Texas

Chicago Illinois Philadelphia Pennsylvania

Chicago Illinois Phoenix Arizona

Chicago Illinois Dallas Texas

Chicago Illinois Jacksonville Florida

Chicago Illinois Indianapolis Indiana

Chicago Illinois Austin Texas

Chicago Illinois Detroit Michigan

Chicago Illinois Memphis Tennessee

Chicago Illinois Boston Massachusetts



Analogy evaluation and hyperparameters

- Accuracy

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<u>80.8</u>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW [†]	300	6B	63.6	<u>67.4</u>	65.7
SG [†]	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>
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
GloVe	300	42B	<u>81.9</u>	<u>69.3</u>	<u>75.0</u>



Another intrinsic word vector evaluation

- Word vector distances and their correlation with human judgments
- Example dataset: WordSim353
<http://www.cs.technion.ac.il/~gabr/resources/data/wordsim353/>

Word 1	Word 2	Human (mean)
tiger	cat	7.35
tiger	tiger	10
book	paper	7.46
computer	internet	7.58
plane	car	5.77
professor	doctor	6.62
stock	phone	1.62
stock	CD	1.31
stock	jaguar	0.92



Correlation evaluation

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 [†]	6B	57.2	65.6	68.2	57.0	32.5
SG [†]	6B	62.8	65.2	69.7	<u>58.1</u>	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	<u>75.9</u>	<u>83.6</u>	<u>82.9</u>	<u>59.6</u>	<u>47.8</u>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5



Extrinsic word vector evaluation

- Extrinsic evaluation of word vectors: All subsequent tasks in this class
- One example where good word vectors should help directly is Named Entity Recognition (i.e., finding names of persons, organization, or locations in text)

Model	Dev	Test	ACE	MUC7
Discrete	91.0	85.4	77.4	73.4
SVD	90.8	85.7	77.3	73.7
SVD-S	91.0	85.5	77.6	74.3
SVD-L	90.5	84.8	73.6	71.5
HPCA	92.6	88.7	81.7	80.7
HSMN	90.5	85.7	78.7	74.7
CW	92.2	87.4	81.7	80.2
CBOW	93.1	88.2	82.2	81.1
GloVe	93.2	88.3	82.9	82.2

- Word vectors/representations have been a major breakthrough in NLP, enabling a novel approach for NLP based on deep learning, and leading to a new era for NLP with models of better performance, robustness and portability.
- We will study a new generation of word vectors later.

