

CS378: Natural Language Processing

Lecture 11: Word Embeddings



Eunsol Choi

Slides adapted from Greg Durrett, Princeton NLP course



Representation of Words

- ▶ Traditional approaches build one-hot feature vectors.

truck = [0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]

austin = [0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0]

- ▶ The dimension of the vector: the size of the vocabulary!
- ▶ More over, no generalization....



Distributional Semantics

“tejuino”



Context 1: A bottle of ____ is on the table.

Context 2: Everybody likes ____.

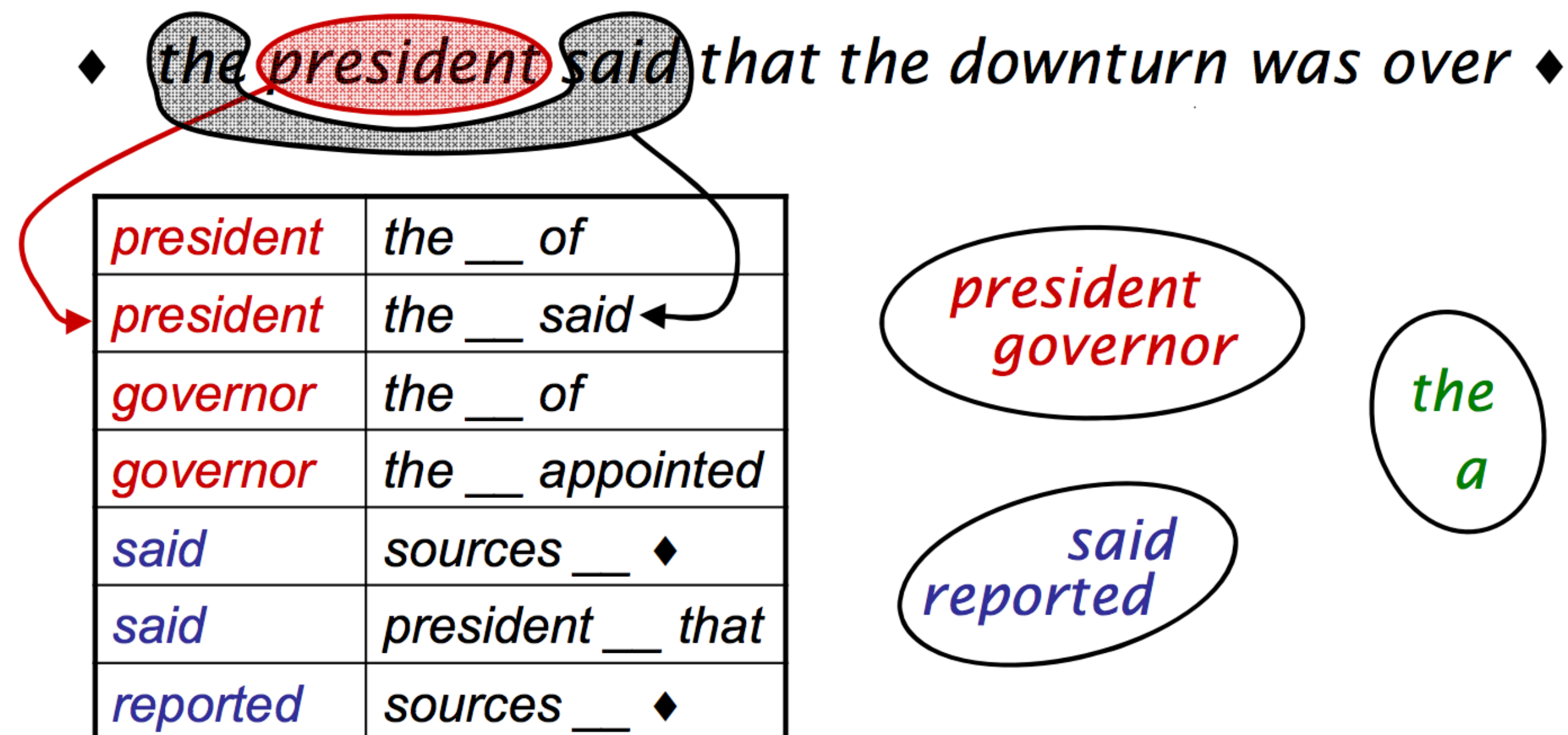
Context 3: Don't have ____ before you drive.

Context 4: We make ____ out of corn.



Distributional Semantics

- ▶ “You shall know a word by the company it keeps” Firth (1957)



[Finch and Chater 92, Shuetze 93, many others]



Goal of Word Embeddings

- ▶ Learn a ***continuous, dense*** vector for each word type.
- ▶ Always the same vector, regardless of in which context the word appears

$$\text{employees} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 10.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$$





- The cats are under the table.

Word ids $w_1, w_2, w_3, w_4, w_5, w_6$

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

```

graph TD
    Root(( )) ---|0| Node0(( ))
    Root ---|1| Node1(( ))
    Node0 ---|0| Node00(( ))
    Node0 ---|1| Node01(( ))
    Node1 ---|0| Node10["is<br/>go"]
    Node1 ---|1| Node11(( ))
    Node00 ---|0| L00[0]
    Node00 ---|1| L01[1]
    Node01 ---|0| L01_0[cat]
    Node01 ---|1| L01_1[fish]
    Node11 ---|0| L11_0[great]
    Node11 ---|1| L11_1[enjoyable]
    L01_0 --- L01_0_1[dog]
    L01_0 --- L01_0_2[...]
    L11_0 --- L11_0_1[good]
  
```

Brown et al. (1992)



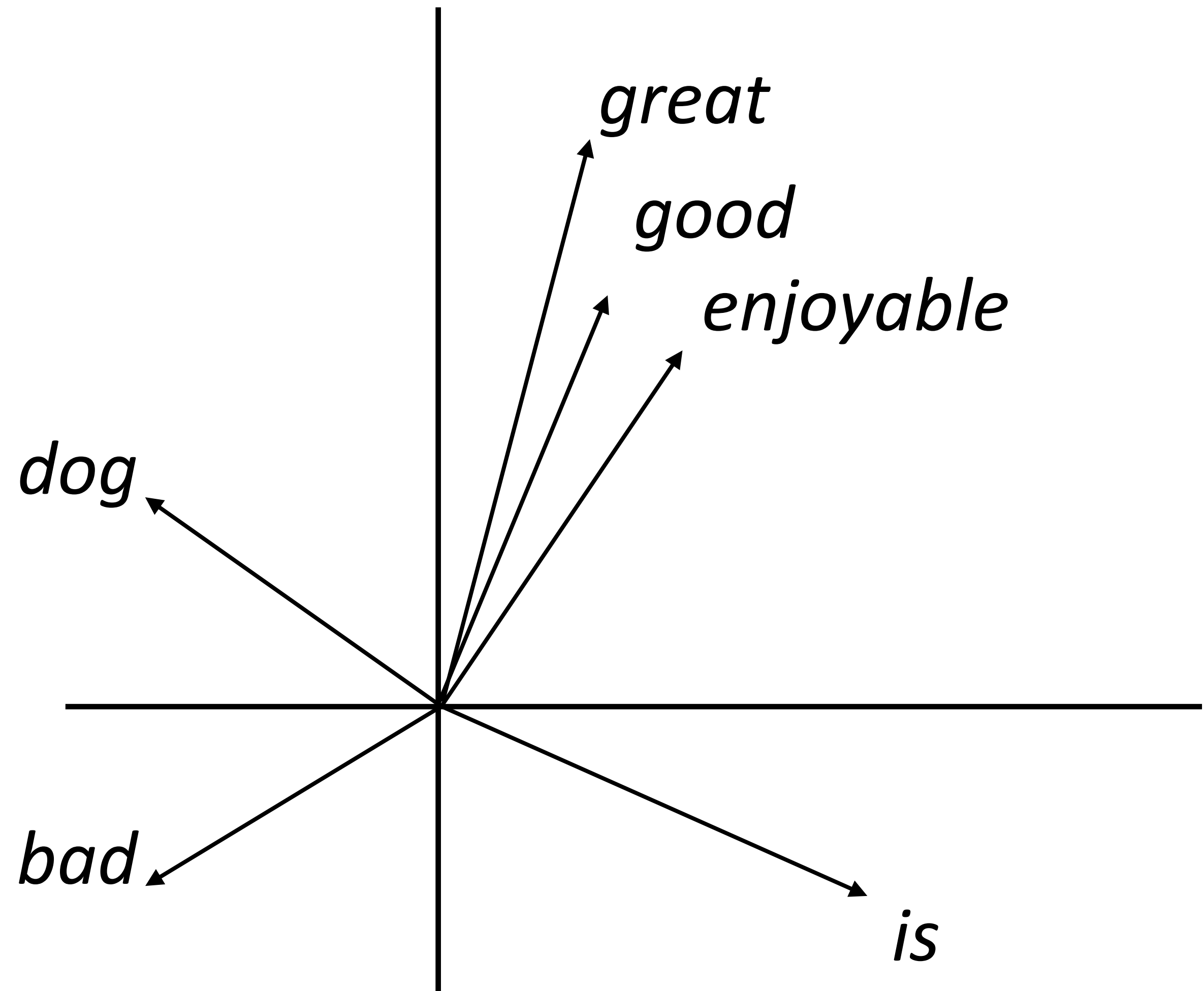
Word Embeddings

the movie was great

\approx

the movie was good

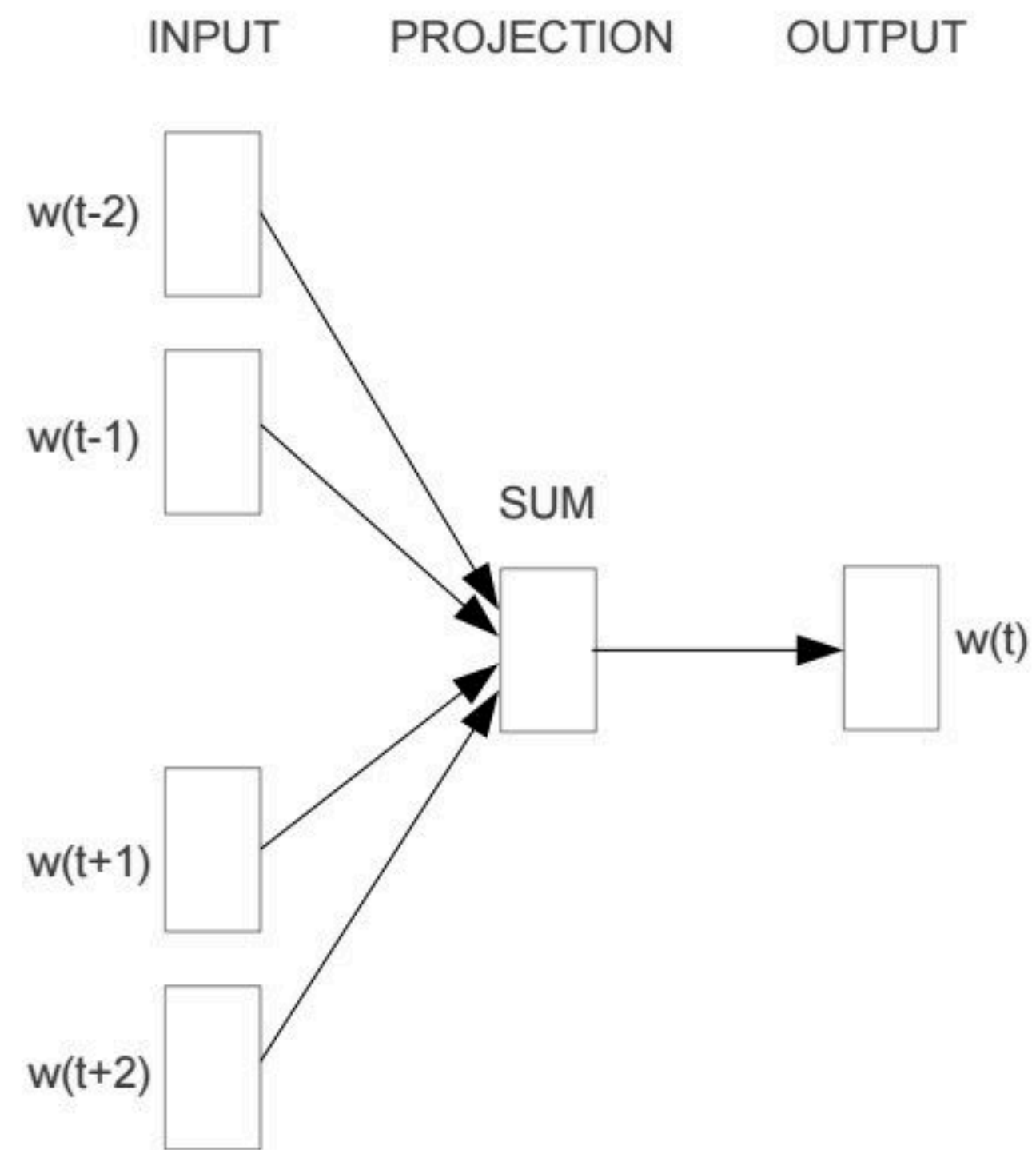
- ▶ Goal: come up with a way to produce these embeddings
- ▶ For each word, “medium” dimensional vector (50-300 dims) representing it



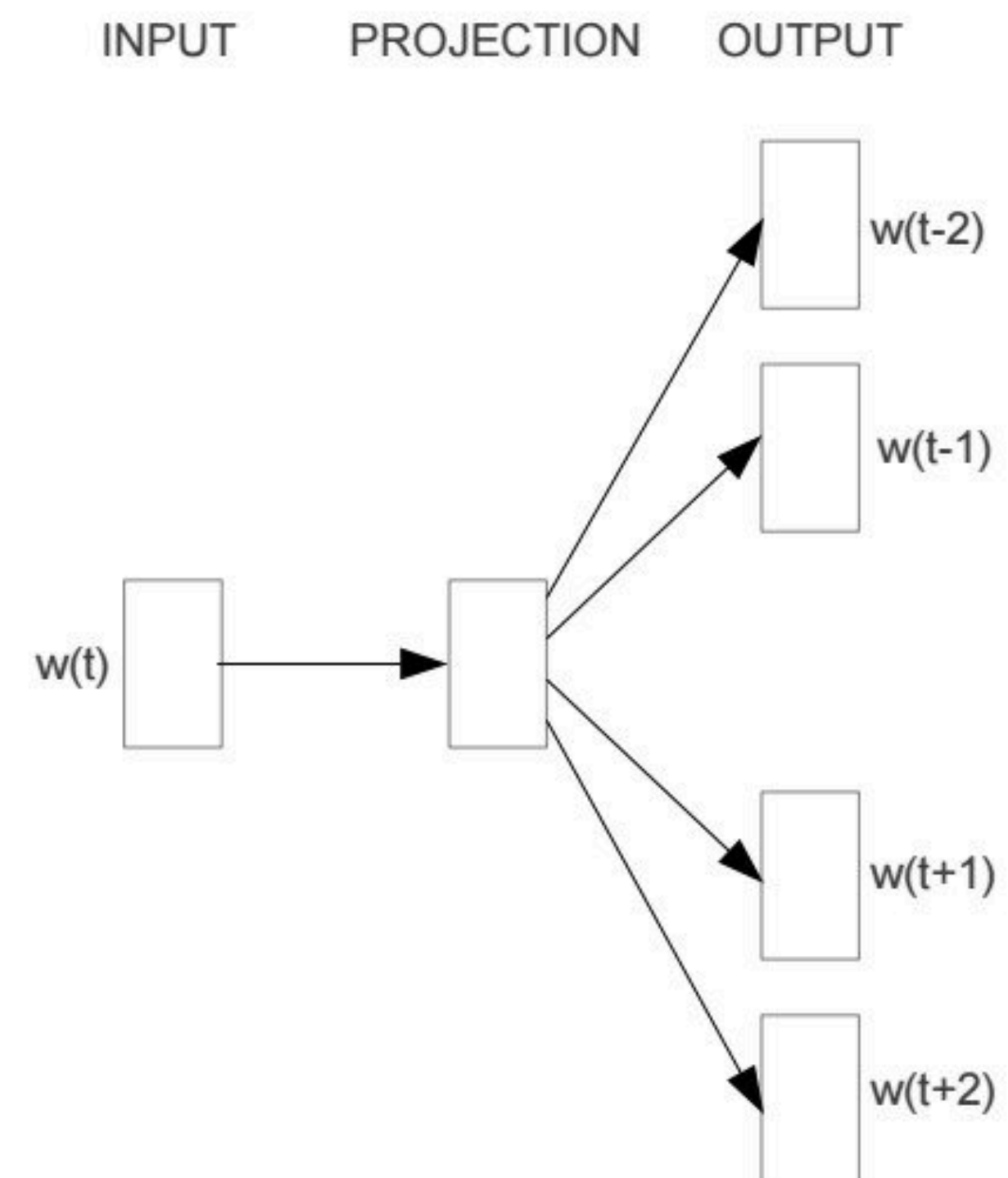
word2vec / GloVe (**Global Vectors**)



word2vec



Continuous Bag of Words (CBOW)



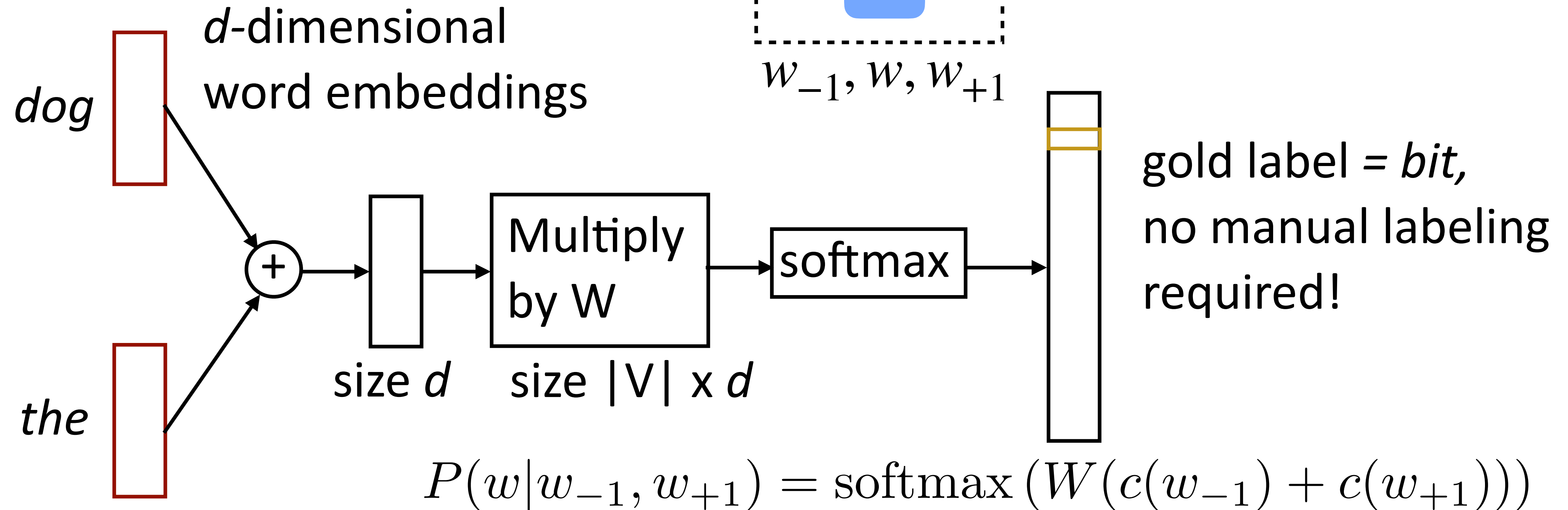
Skip-grams



Continuous Bag-of-Words

- Predict word from context

the dog bit the man
 w_{-1}, w, w_{+1}



- Parameters: $d \times |V|$ (one d -length **context vector per word**),
 $|V| \times d$ output parameters (W)

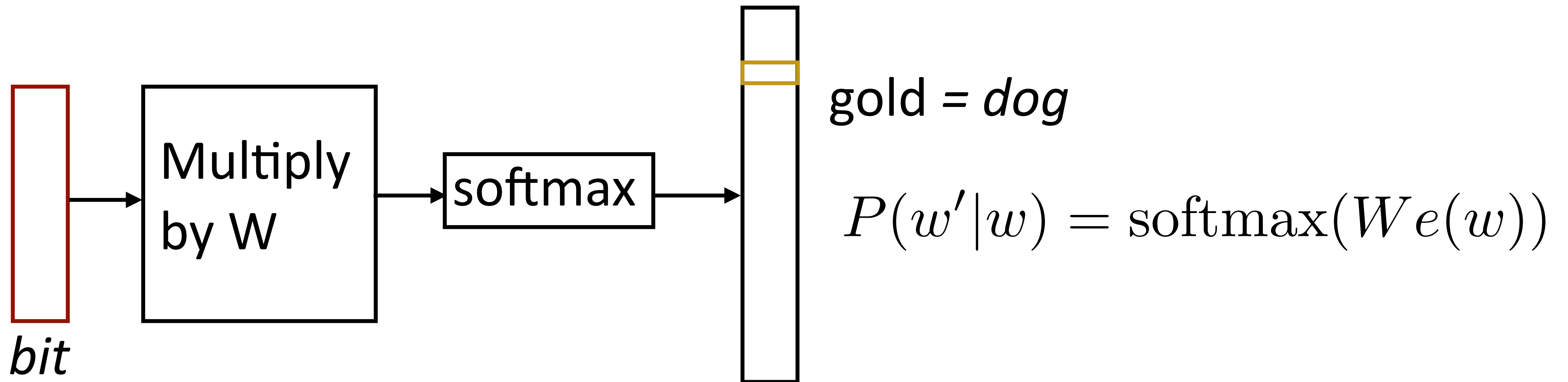
Mikolov et al. (2013)



Skip-Gram

- Predict one word of context from word

the dog bit the man



- Another training example: *bit* -> *the*
- Parameters: $d \times |V|$ word **vector**, $|V| \times d$ output parameters (W) (also usable as vectors!)

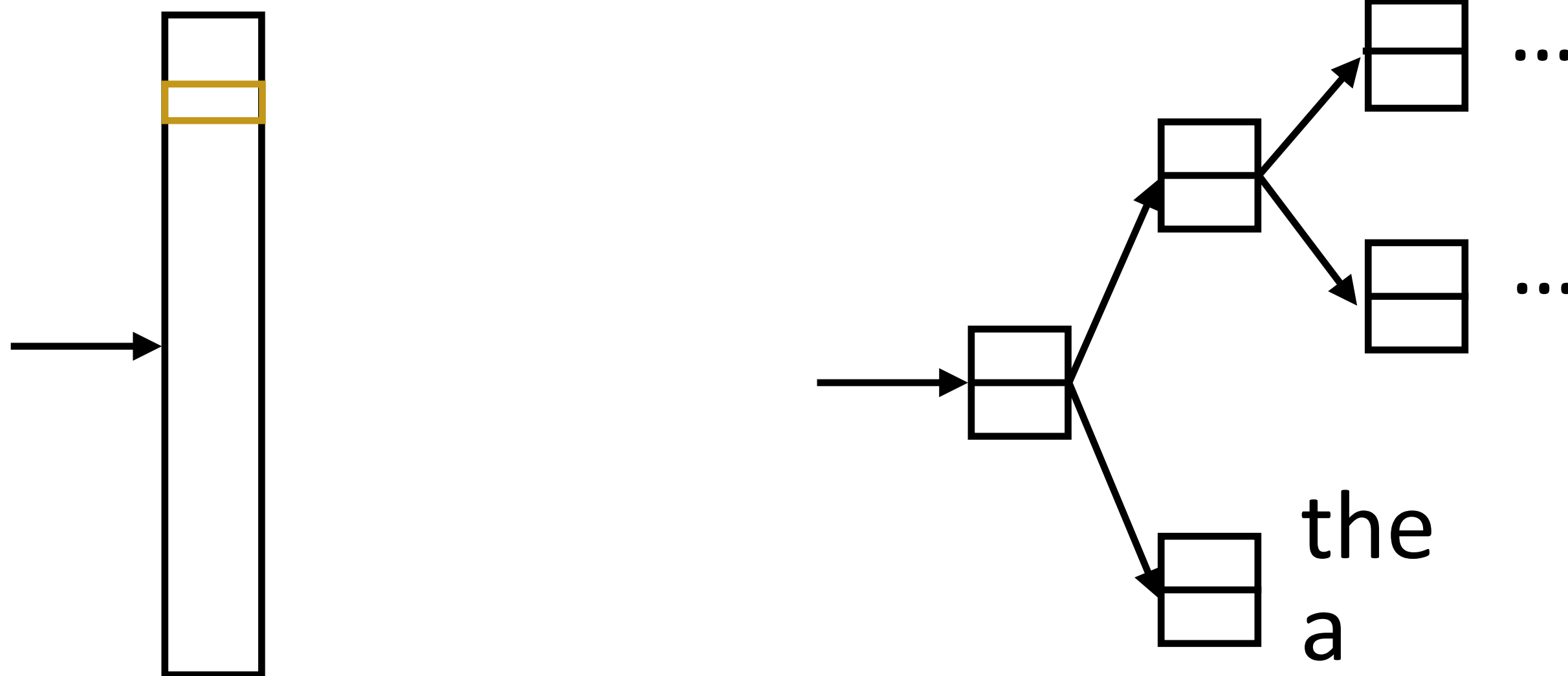
Mikolov et al. (2013)



Hierarchical Softmax

$$P(w|w_{-1}, w_{+1}) = \text{softmax}(W(c(w_{-1}) + c(w_{+1}))) \quad P(w'|w) = \text{softmax}(We(w))$$

- ▶ Matmul + softmax over $|V|$ is very slow to compute for CBOW and SG



- ▶ Binary classifiers to decide which branch to take
- ▶ $\log(|V|)$ binary decisions



Skip-Gram with Negative Sampling

- Take (word, context) pairs and classify them as “real” or not. Create random negative examples by sampling from unigram distribution

(bit, the) => +1

(bit, cat) => -1

(bit, a) => -1

(bit, fish) => -1

$$P(y = 1 | w, c) = \frac{e^{w \cdot c}}{e^{w \cdot c} + 1}$$

words in similar contexts select for similar c vectors

- $d \times |V|$ vectors, $d \times |V|$ context vectors (same # of params as before)

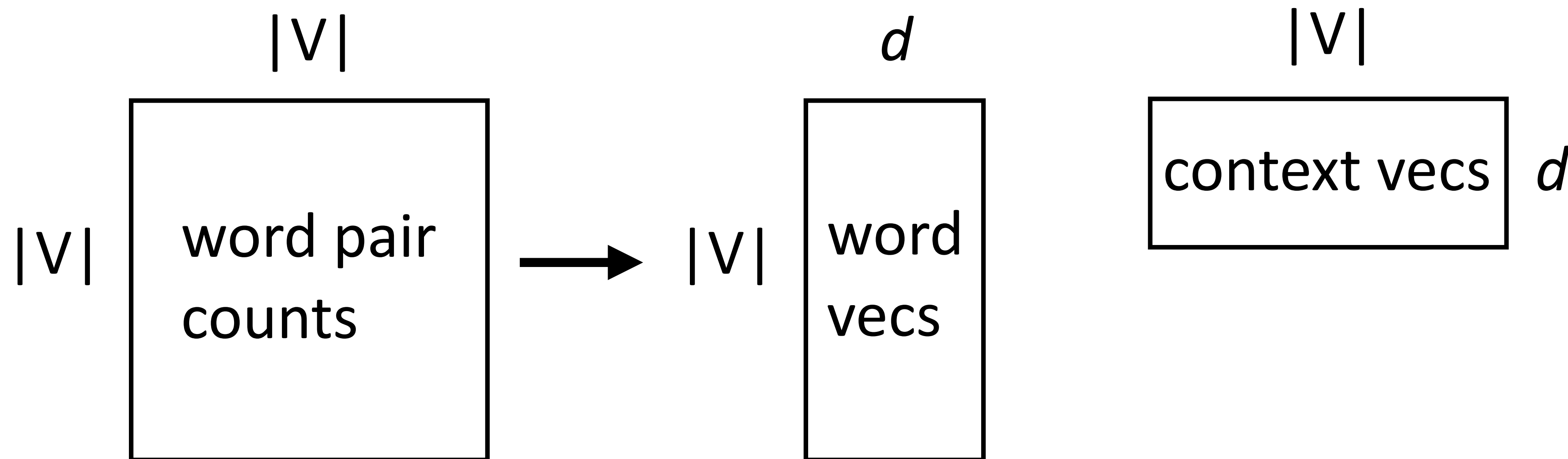
- Objective = $\log P(y = 1 | w, c) + \frac{1}{k} \sum_{i=1}^k \log P(y = 0 | w, c_i)$ ← sampled

Mikolov et al. (2013)



Connections with Matrix Factorization

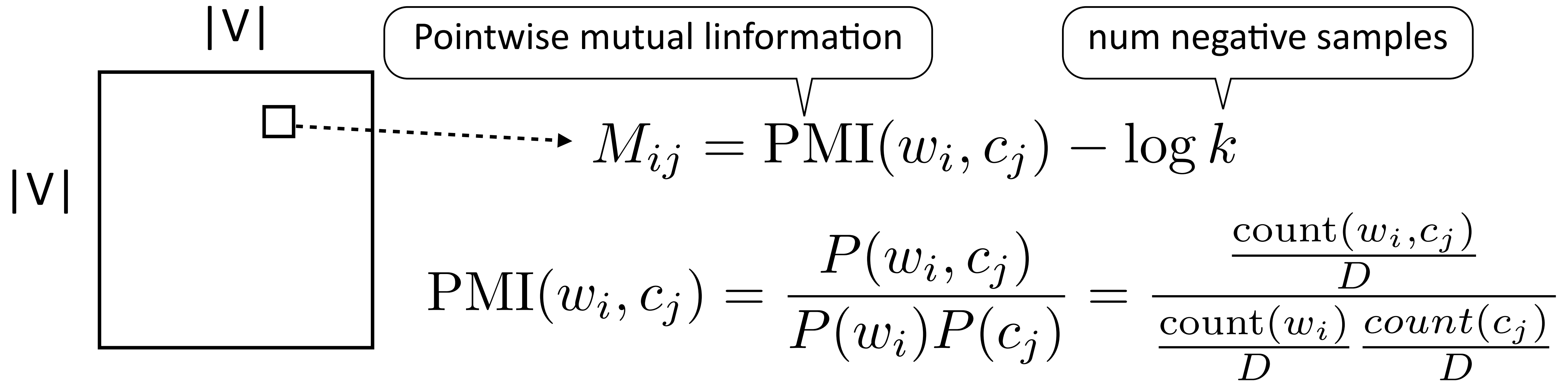
- ▶ Skip-gram model looks at word-word co-occurrences and produces two types of vectors



- ▶ Looks almost like a matrix factorization... can we interpret it this way?



Skip-Gram as Matrix Factorization



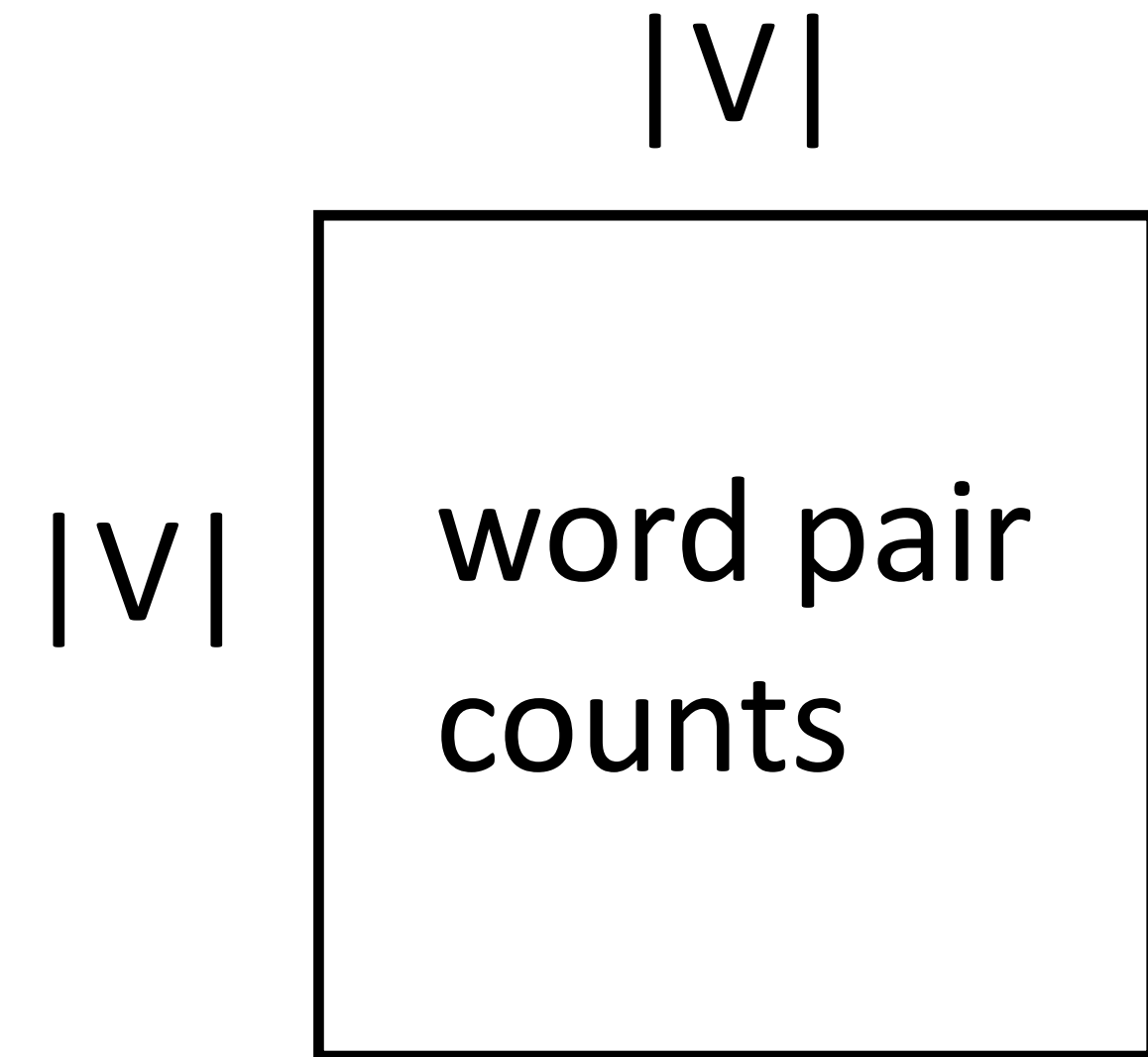
Skip-gram objective *exactly* corresponds to factoring this matrix:

- ▶ If we sample negative examples from the uniform distribution over words
- ▶ ...and it's a *weighted* factorization problem (weighted by word freq)



GloVe (Global Vectors)

- ▶ Also operates on counts matrix, weighted regression on the log co-occurrence matrix



- ▶ Objective = $\sum_{i,j} f(\text{count}(w_i, c_j)) (w_i^\top c_j + a_i + b_j - \log \text{count}(w_i, c_j))^2$
- ▶ Constant in the dataset size (just need counts), quadratic in voc size



fastText: Sub-word Embeddings

- ▶ Same as SGNS, but break words down into n-grams with $n = 3$ to 6

where:

3-grams: <wh, whe, her, ere, re>

4-grams: <whe, wher, here, ere> ,

5-grams: <wher, where, here> ,

6-grams: <where, where>

- ▶ Replace $w \cdot c$ in skip-gram computation with $\left(\sum_{g \in \text{ngrams}} w_g \cdot c \right)$

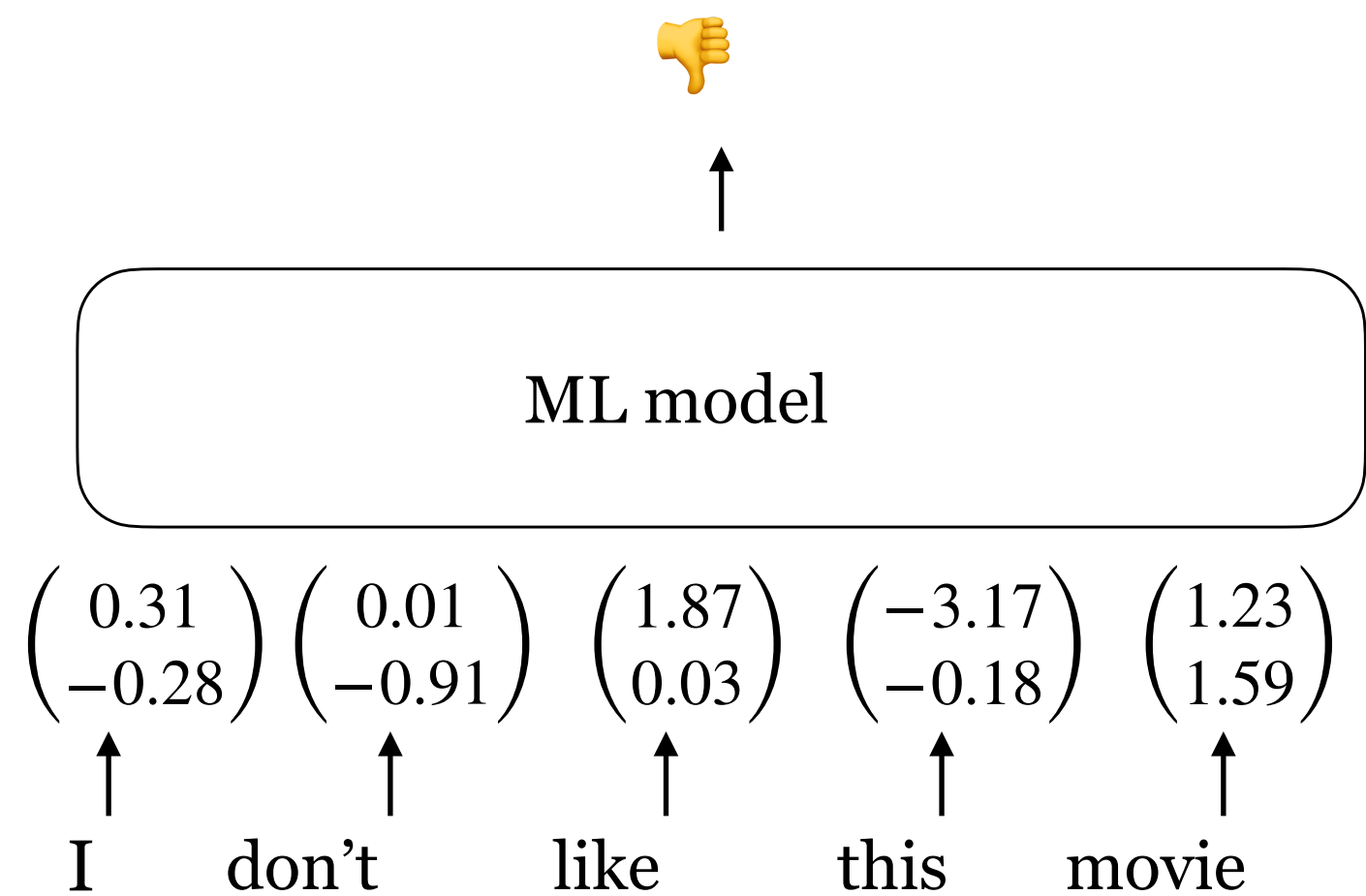
- ▶ Advantages?

Evaluating Word Embeddings



Extrinsinc vs. Intrinsic

- ▶ Plug the vectors into some end-task model, see which does well!
- ▶ Evaluate on some intermediate subtasks (in the coming slides..)





Evaluating Word Embeddings

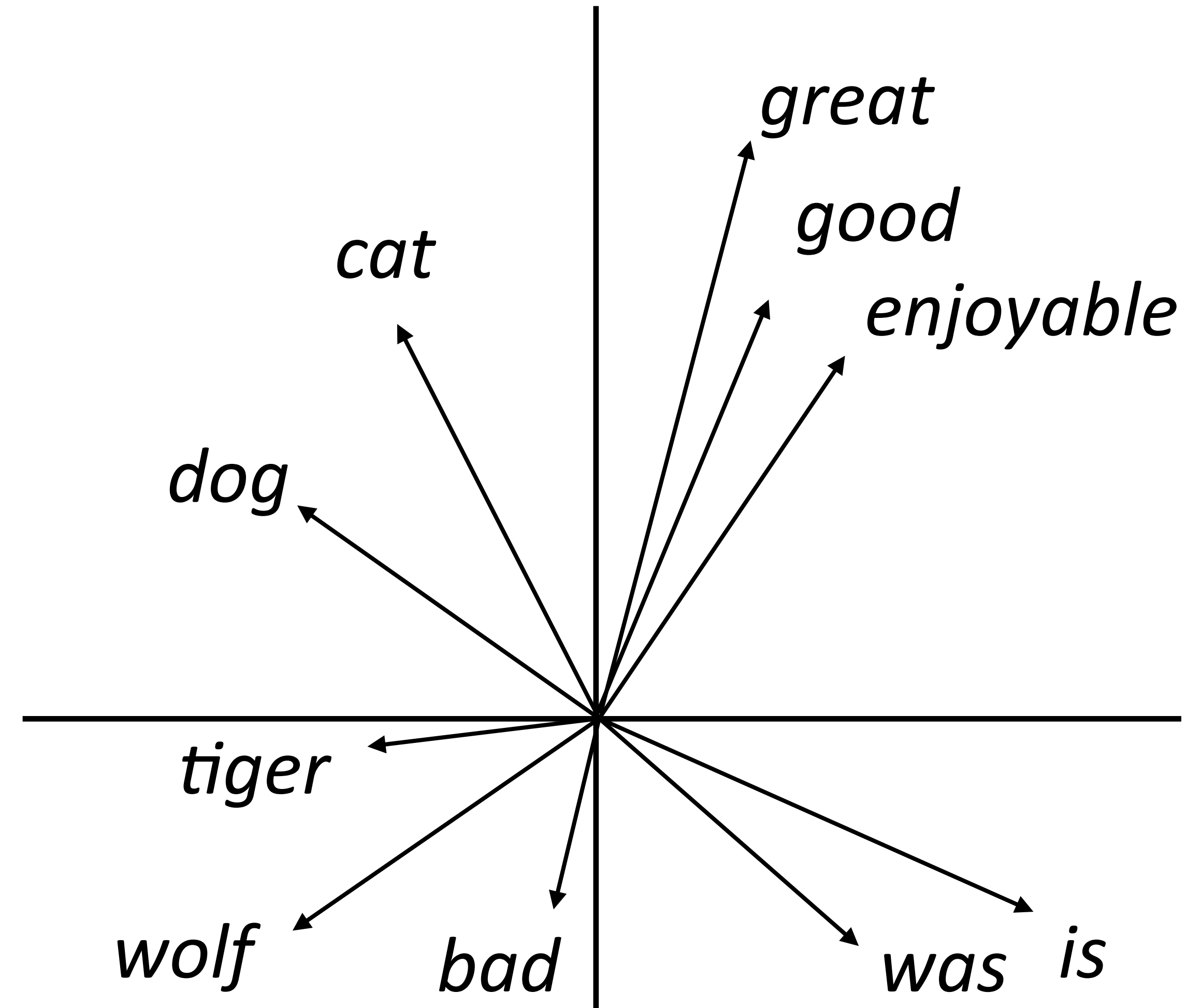
- ▶ What properties of language should word embeddings capture?

- ▶ (1) Similarity: similar words are close to each other

- ▶ (2) Analogy:

good is to best as smart is to ???

Paris is to France as Tokyo is to ???





Similarity

Method	WordSim Similarity	WordSim Relatedness	Bruni et al. MEN	Radinsky et al. M. Turk	Luong et al. Rare Words	Hill et al. SimLex
PPMI	.755	.697	.745	.686	.462	.393
SVD	.793	.691	.778	.666	.514	.432
SGNS	.793	.685	.774	.693	.470	.438
GloVe	.725	.604	.729	.632	.403	.398

- ▶ SVD = singular value decomposition on PMI matrix
- ▶ GloVe does not appear to be the best when experiments are carefully controlled, but it depends on hyperparameters + these distinctions don't matter in practice

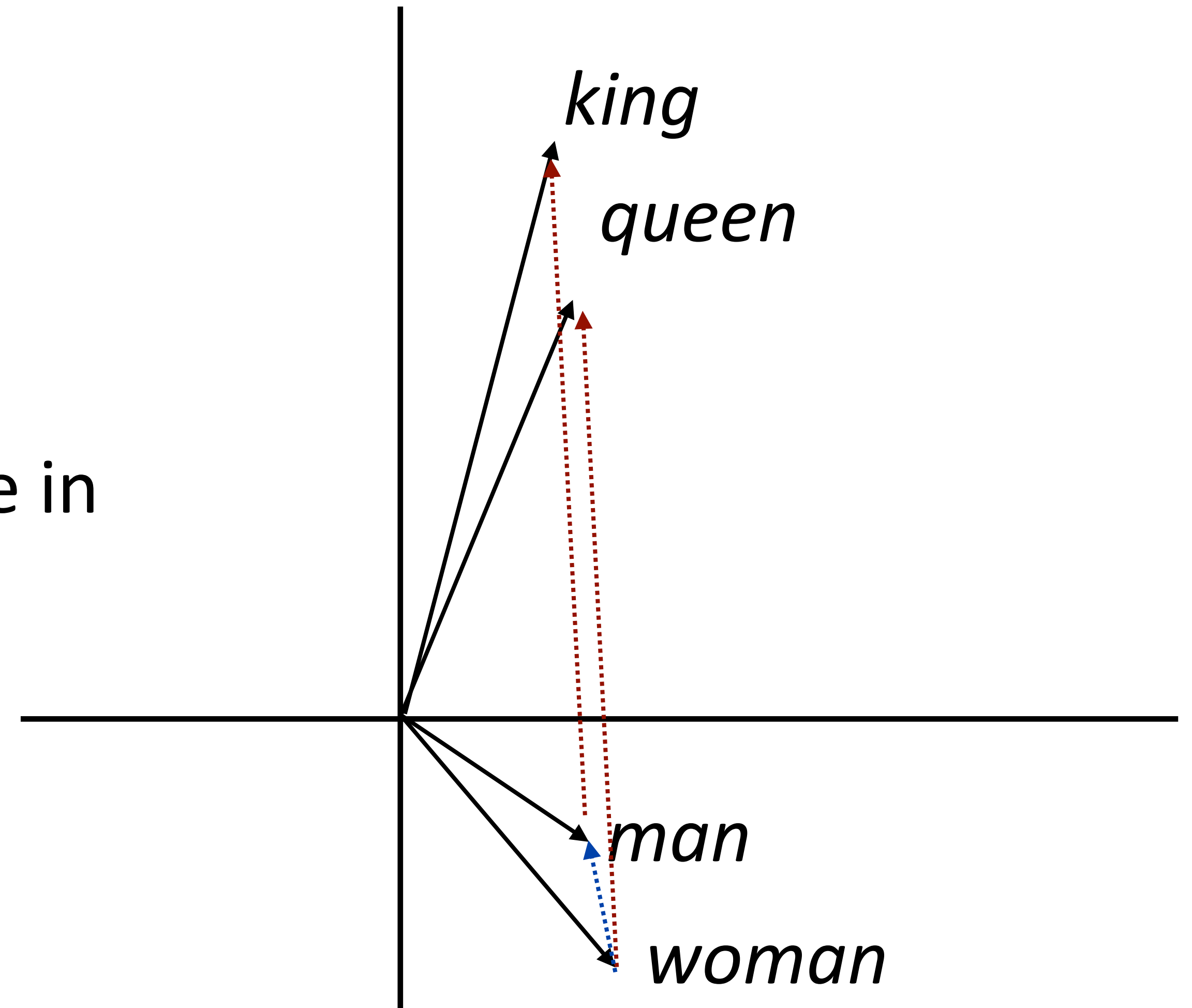


Analogies

$(king - man) + woman = queen$

$king + (woman - man) = queen$

- ▶ Why would this be?
- ▶ woman - man captures the difference in the contexts that these occur in
- ▶ Often used to evaluate word embeddings





What can go wrong with word embeddings?

- ▶ What's wrong with learning a word's "meaning" from its usage?
- ▶ What data are we learning from?
- ▶ What are we going to learn from this data?



What do we mean by bias?

- Identify *she* - *he* axis in word vector space, project words onto this axis

Extreme *she* occupations

- | | | |
|-----------------|-----------------------|------------------------|
| 1. homemaker | 2. nurse | 3. receptionist |
| 4. librarian | 5. socialite | 6. hairdresser |
| 7. nanny | 8. bookkeeper | 9. stylist |
| 10. housekeeper | 11. interior designer | 12. guidance counselor |

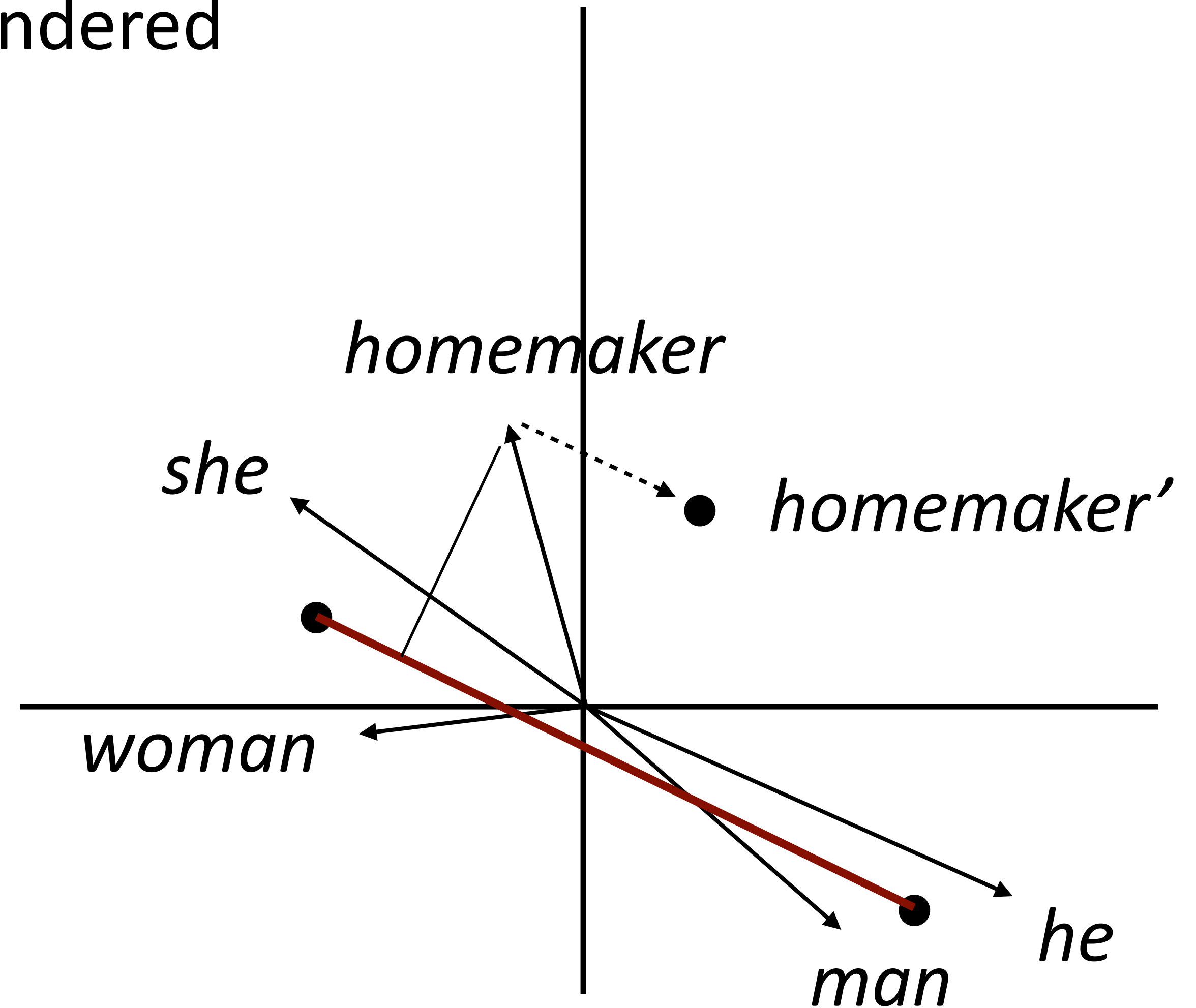
Extreme *he* occupations

- | | | |
|----------------|-------------------|----------------|
| 1. maestro | 2. skipper | 3. protege |
| 4. philosopher | 5. captain | 6. architect |
| 7. financier | 8. warrior | 9. broadcaster |
| 10. magician | 11. fighter pilot | 12. boss |



Debiasing

- ▶ Identify gender subspace with gendered words
- ▶ Project words onto this subspace
- ▶ Subtract those projections from the original word

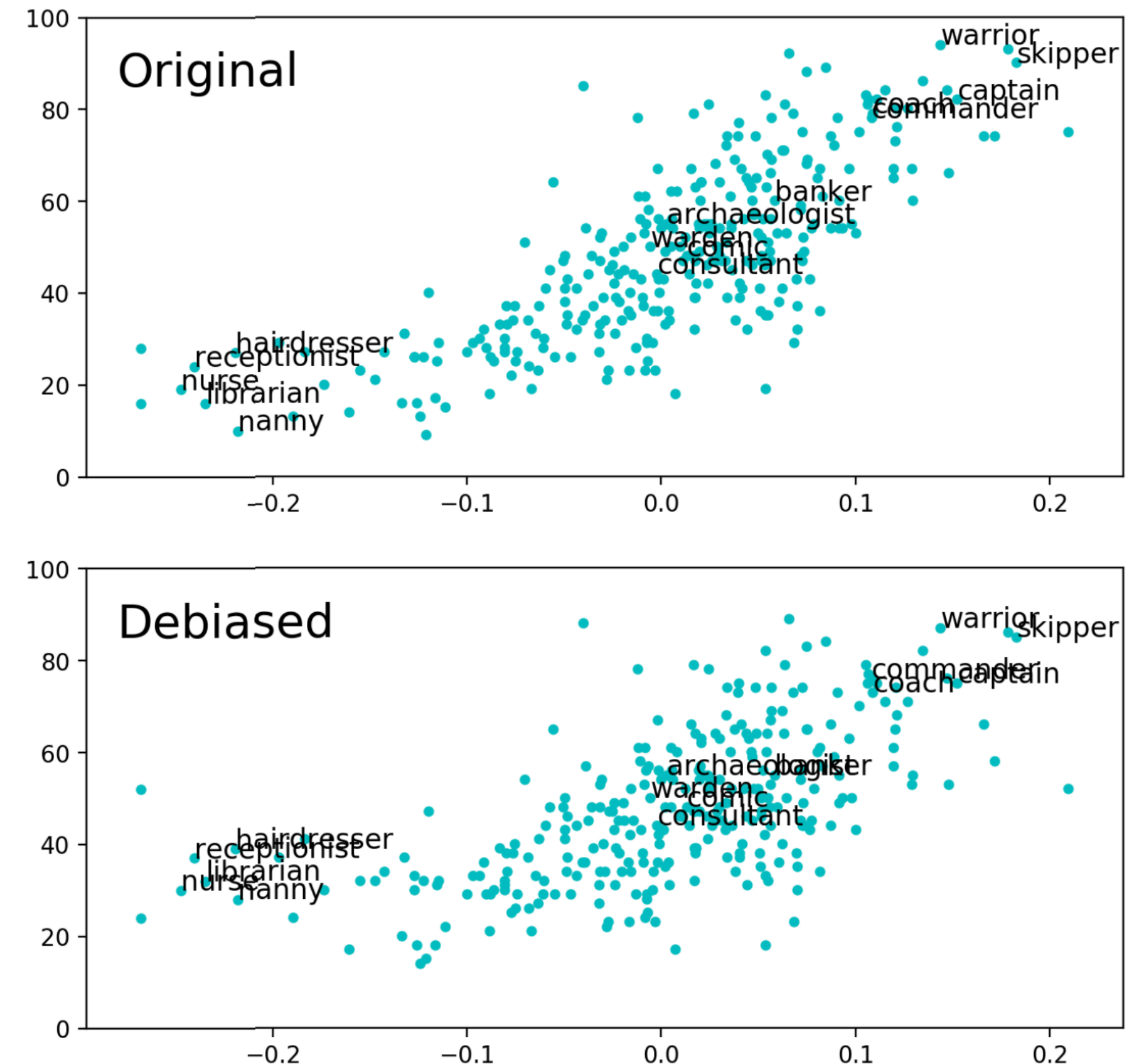


Bolukbasi et al. (2016)



Hardness of Debiasing

- ▶ Not that effective...and the male and female words are still clustered together
- ▶ Bias pervades the word embedding space and isn't just a local property of a few words



(a) The plots for HARD-DEBIASED embedding, before (top) and after (bottom) debiasing.

Gonen and Goldberg (2019)



Trained Word Embeddings

- ▶ word2vec: <https://code.google.com/archive/p/word2vec/>
- ▶ GloVe: <https://nlp.stanford.edu/projects/glove/>
- ▶ FastText: <https://fasttext.cc/>

Download pre-trained word vectors

- Pre-trained word vectors. This data is made available under the [Public Domain Dedication and License](http://www.opendatacommons.org/licenses/pddl/1.0/) v1.0 whose full text can be found at: <http://www.opendatacommons.org/licenses/pddl/1.0/>.
 - [Wikipedia 2014](#) + [Gigaword 5](#) (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download): [glove.6B.zip](#)
 - Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download): [glove.42B.300d.zip](#)
 - Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download): [glove.840B.300d.zip](#)
 - Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors, 1.42 GB download): [glove.twitter.27B.zip](#)
- Ruby [script](#) for preprocessing Twitter data



Using Word Embeddings

- ▶ Approach 1: learn embeddings as parameters from your data
 - ▶ Often works pretty well
- ▶ Approach 2: initialize using pretrained word vectors (e.g., GloVe), keep fixed
 - ▶ Faster because no need to update these parameters
- ▶ Approach 3: initialize using pretrained word vectors (e.g., GloVe), fine-tune
 - ▶ Works best for some tasks



Preview: Context-dependent Embeddings

- ▶ How to handle different word senses? One vector for *balls*

they dance at balls they hit the balls

- ▶ Train a neural language model to predict the next word given previous words in the sentence, use its internal representations as word vectors
- ▶ *Context-sensitive* word embeddings: depend on rest of the sentence
- ▶ *Huge* improvements across nearly all NLP tasks over static word embeddings

Peters et al. (2018)



Compositional Semantics

- ▶ What if we want embedding representations for whole sentences?
- ▶ Skip-*thought* vectors (Kiros et al., 2015), similar to skip-gram generalized to a sentence level
- ▶ Is there a way we can compose vectors to make sentence representations?



Summary

- ▶ Lots of pretrained embeddings work well in practice, they capture some desirable properties
- ▶ Even better: context-sensitive word embeddings (ELMo)
- ▶ Next time: Language Models!