

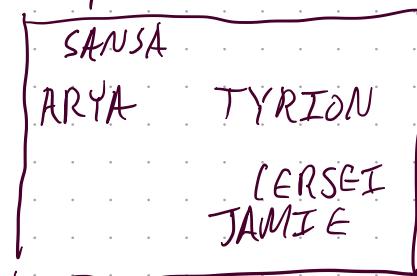
Word Embeddings (aka word2vec)

Word frequency models erase contextual information of words in favor doc-level understanding

Word Embeddings → look at word meaning.

↳ Understanding word meaning directly

↳ Leverage to generate sequences of text



Learning Concepts

Semi-Supervised learning) - includes error term for some part not directly related to the task at hand

Supervised learning - error term directly related to task at hand

Unsupervised learning - no error term at all

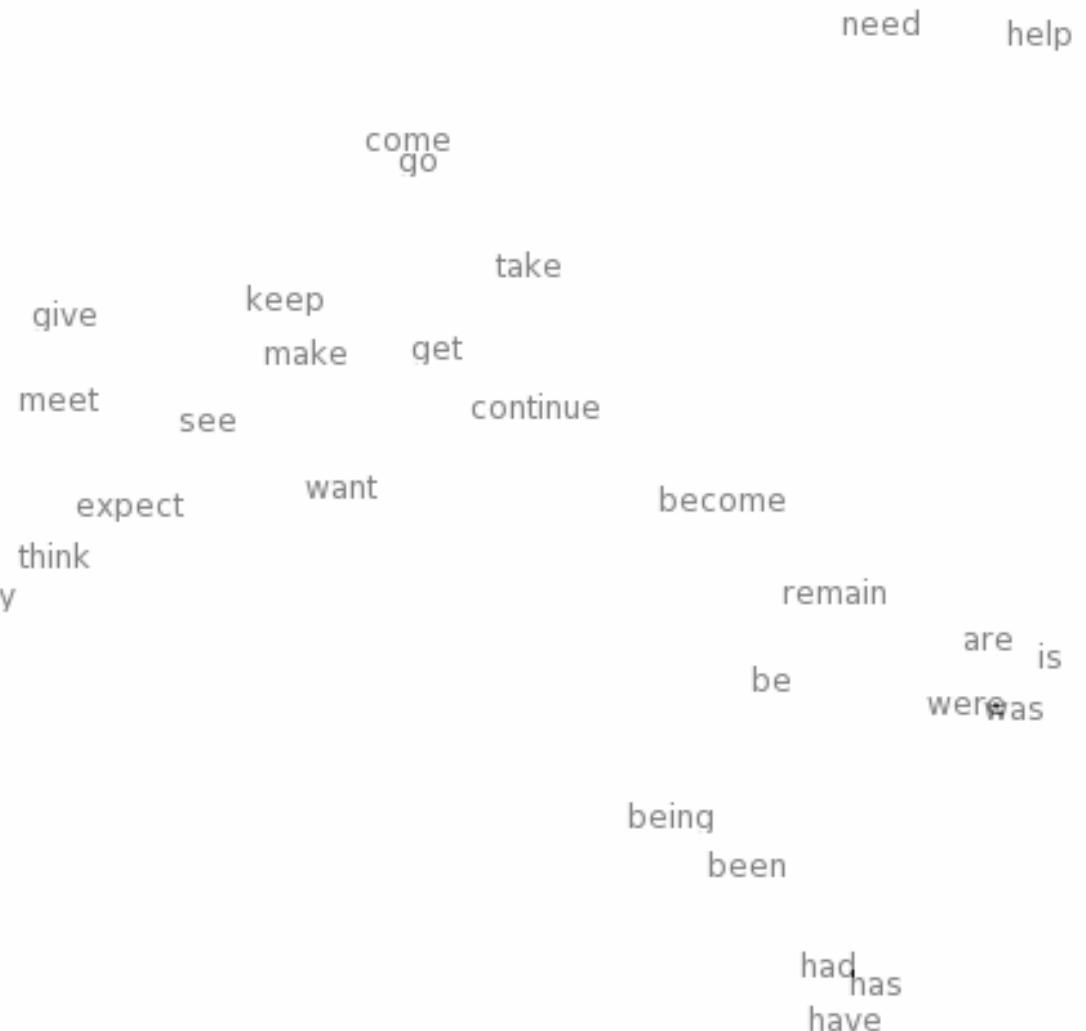
Transfer Learning model trained in part on a different task from the one at hand

VGG

Distributed vs. distributional representations

Distributed: A concept is represented as continuous activation levels in a number of elements

[Contrast: local]



Distributional: Meaning is represented by contexts of use

[Contrast: denotational]

Constructing and Evaluating Word Embeddings

Dr Marek Rei and Dr Ekaterina Kochmar
Computer Laboratory
University of Cambridge

Representing words as vectors

Let's represent words (or any objects) as vectors.

We want to construct them so that similar words have similar vectors.

| Sequence | Count |
|---------------------|-------|
| I live in Cambridge | 19 |
| I live in Paris | 68 |
| I live in Tallinn | 0 |
| I live in yellow | 0 |

- ✗ Tallinn
 - ✗ Cambridge
 - ✗ London
 - ✗ Paris

- ✗ yellow
- ✗ red
- ✗ blue
- ✗ green

Representing words as vectors

Let's represent words (or any objects) as vectors.

We want to construct them, so that similar words have similar vectors.

| Sequence | Count |
|---------------------|-------|
| I live in Cambridge | 19 |
| I live in Paris | 68 |
| I live in Tallinn | 0 |
| I live in yellow | 0 |

- ✗ Tallinn
 - ✗ Cambridge
 - ✗ London
 - ✗ Paris

- ✗ yellow
- ✗ red
- ✗ green
- ✗ blue

1-hot vectors

How can we represent words as vectors?

Option 1: each element represents a different word.

Also known as “1-hot” or “1-of-V” representation.

| | <i>bear</i> | <i>cat</i> | <i>frog</i> |
|-------------|-------------|------------|-------------|
| bear | 1 | 0 | 0 |
| cat | 0 | 1 | 0 |
| frog | 0 | 0 | 1 |

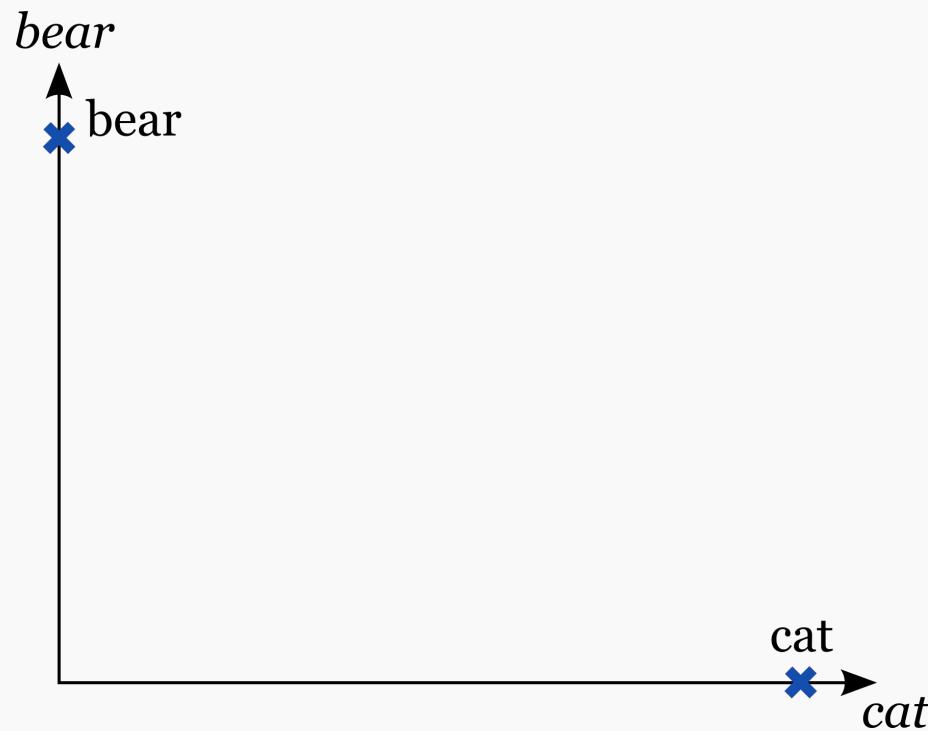
$$\text{bear} = [1.0, 0.0, 0.0]$$

$$\text{cat} = [0.0, 1.0, 0.0]$$

1-hot vectors

When using 1-hot vectors, we can't fit many and they tell us very little.

Need a separate dimension for every word we want to represent.



Distributed vectors

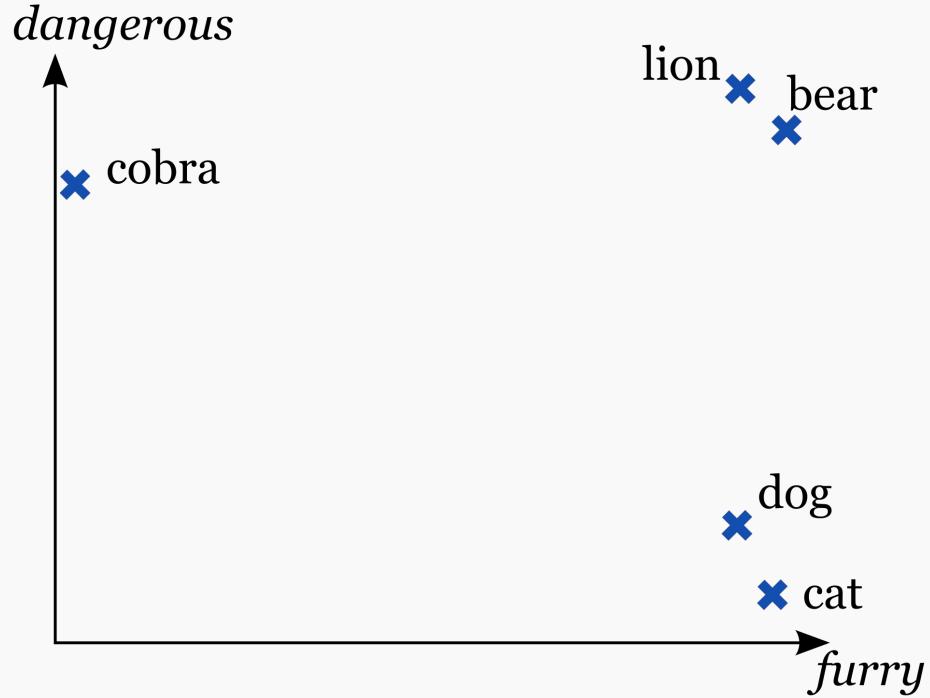
Option 2: each element represents a property, and they are shared between the words.

Also known as “distributed” representation.

| | <i>furry</i> | <i>dangerous</i> | <i>mammal</i> |
|-------------|--------------|------------------|---------------|
| bear | 0.9 | 0.85 | 1 |
| cat | 0.85 | 0.15 | 1 |
| frog | 0 | 0.05 | 0 |

$$\text{bear} = [0.9, 0.85, 1.0] \quad \text{cat} = [0.85, 0.15, 1.0]$$

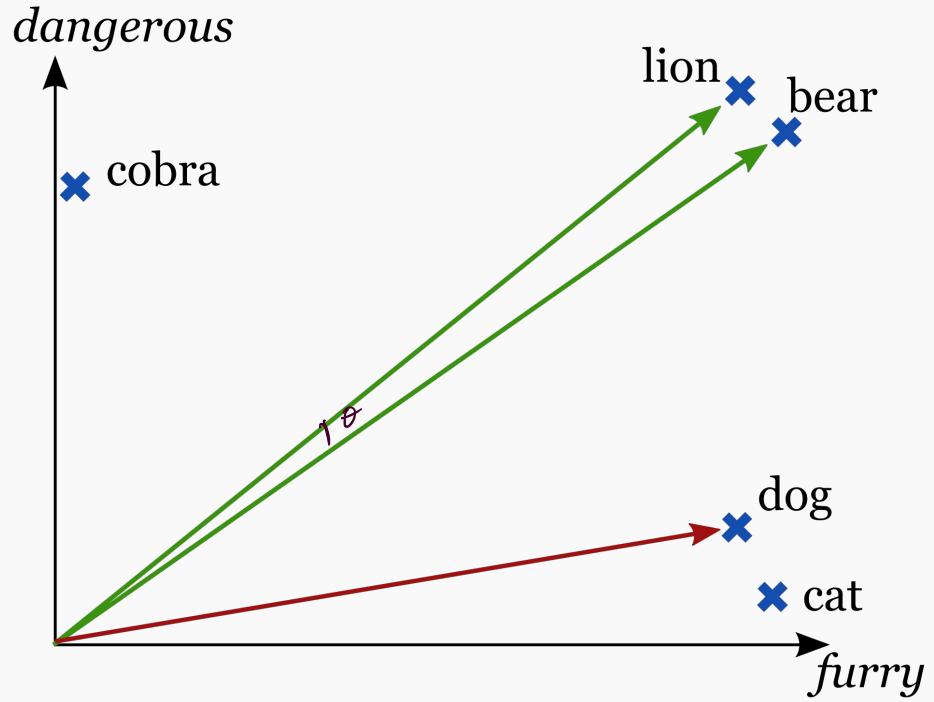
Distributed vectors



| | <i>furry</i> | <i>dangerous</i> |
|--------------|--------------|------------------|
| bear | 0.9 | 0.85 |
| cat | 0.85 | 0.15 |
| cobra | 0.0 | 0.8 |
| lion | 0.85 | 0.9 |
| dog | 0.8 | 0.15 |

Distributed vectors group similar words/objects together

Distributed vectors



$$\cos(a, b) = \frac{\sum_i a_i \cdot b_i}{\sqrt{\sum_i a_i^2} \cdot \sqrt{\sum_i b_i^2}}$$

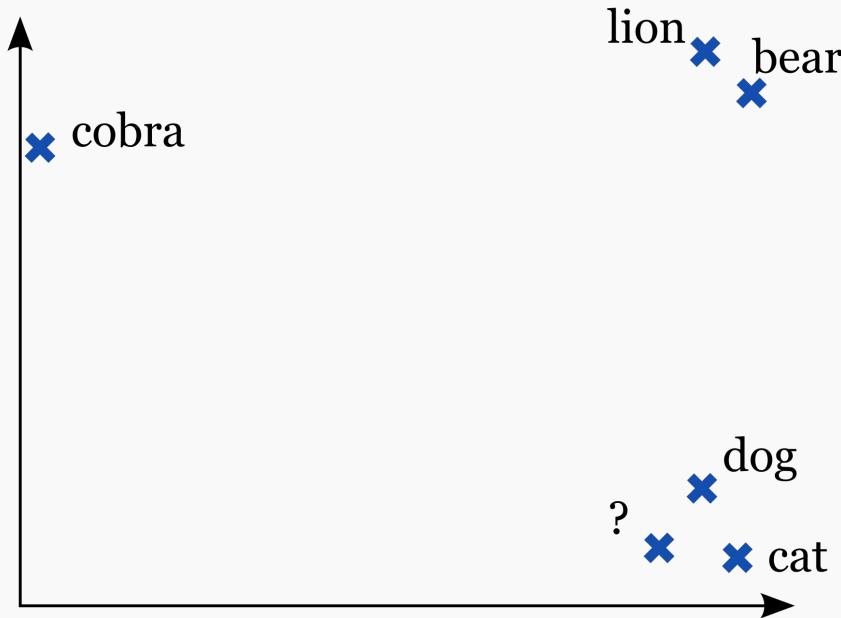
$$\cos(\text{lion}, \text{bear}) = 0.998$$

$$\cos(\text{lion}, \text{dog}) = 0.809$$

$$\cos(\text{cobra}, \text{dog}) = 0.727$$

Can use cosine to calculate similarity between two words

Distributed vectors



We can infer some information, based only on the vector of the word

We don't even need to know the labels on the vector elements

Distributional hypothesis

Words which are similar in meaning occur in similar contexts.

(Harris, 1954)

You shall know a word by the company it keeps

(Firth, 1957)

He is reading a **magazine**

I was reading a **newspaper**

This **magazine** published my story

The **newspaper** published an article

She buys a **magazine** every month

He buys this **newspaper** every day

Count-based vectors

One way of creating a vector for a word:

Let's count how often a word occurs together with specific other words.

He is **reading** a **magazine**

I was **reading** a **newspaper**

This **magazine** published my story

The **newspaper** published an article

She buys a **magazine** every month

He **buys this newspaper** every day

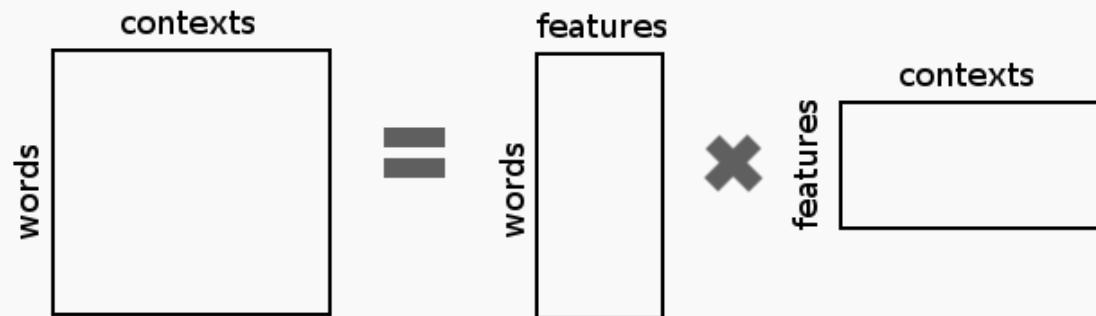
| | reading | a | this | published | my | buys | the | an | every | month | day |
|-----------|---------|---|------|-----------|----|------|-----|----|-------|-------|-----|
| magazine | 1 | 2 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 |
| newspaper | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 |

Count-based vectors

- More frequent words dominate the vectors.
Can use a weighting scheme like PMI or TF-IDF.

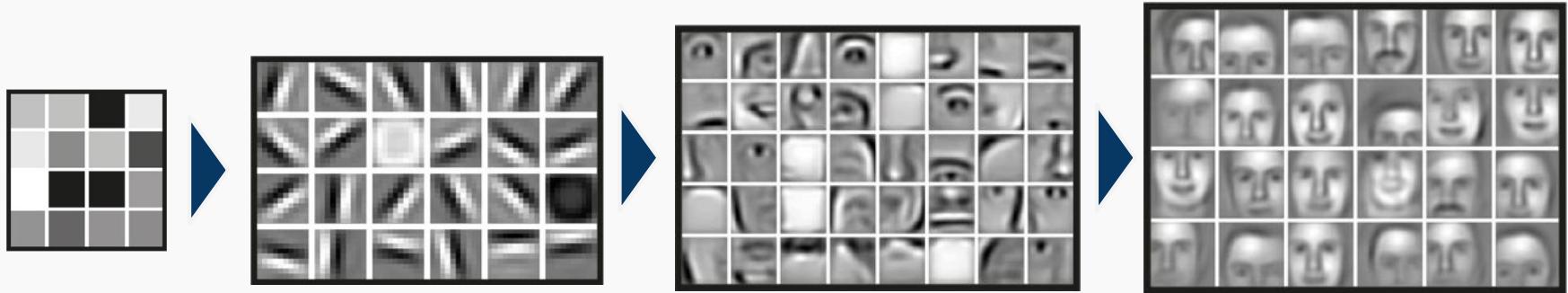
$$pmi(x, z) = \log \frac{p(x, z)}{p(x)p(z)}$$
$$\text{tf-idf}(w, z) = freq_{w,z} \cdot \log \frac{V}{n_z}$$

- Large number of sparse features
Can use matrix decomposition like Singular Value Decomposition (SVD) or Latent Dirichlet Allocation (LDA).



Neural word embeddings

Neural networks will automatically try to discover useful features in the data, given a specific task.



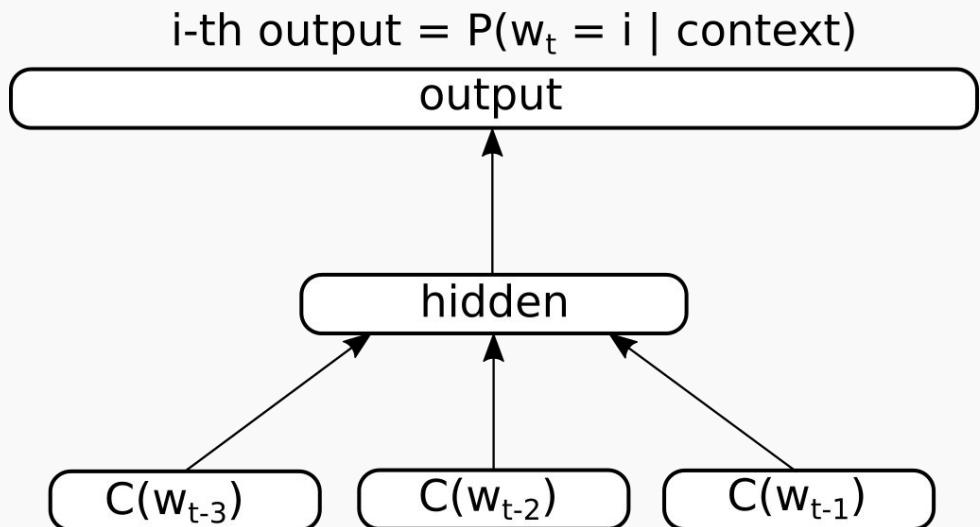
Idea: Let's allocate a number of parameters for each word and allow the neural network to automatically learn what the useful values should be.

Often referred to as “**word embeddings**”, as we are embedding the words into a real-valued low-dimensional space.

Embeddings through language modelling

Predict the next word in a sequence, based on the previous words.

Use this to guide the training for word embeddings.



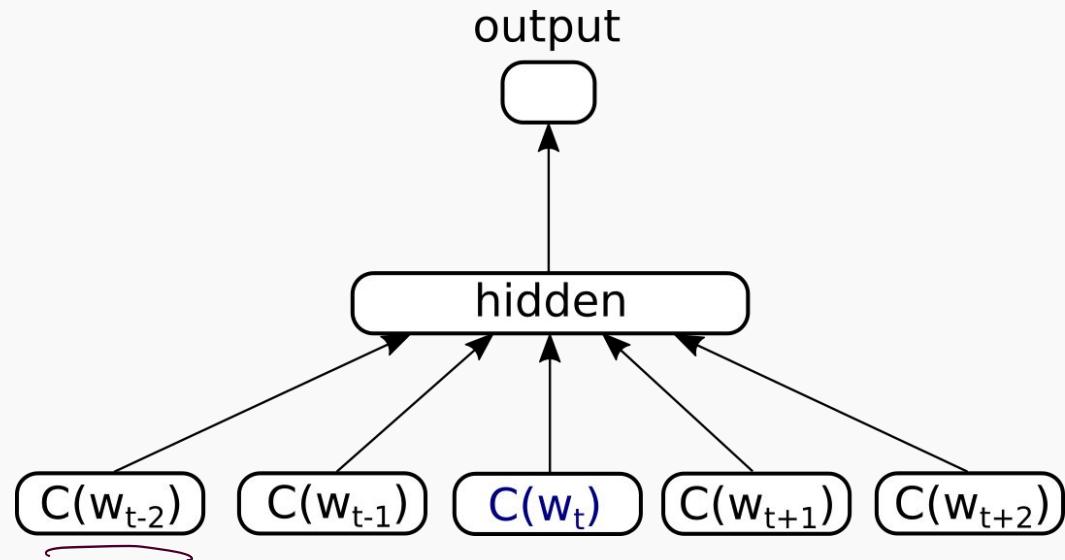
Bengio et. al. 2003. *A Neural Probabilistic Language Model.*

I **read** at my **desk**
I **study** at my **desk**

Embeddings through error detection

Take a grammatically correct sentence and create a corrupted counterpart.

Train the neural network to assign a higher score to the correct version of each sentence.



Collobert et. al. 2011. *Natural Language Processing (Almost) from Scratch.*

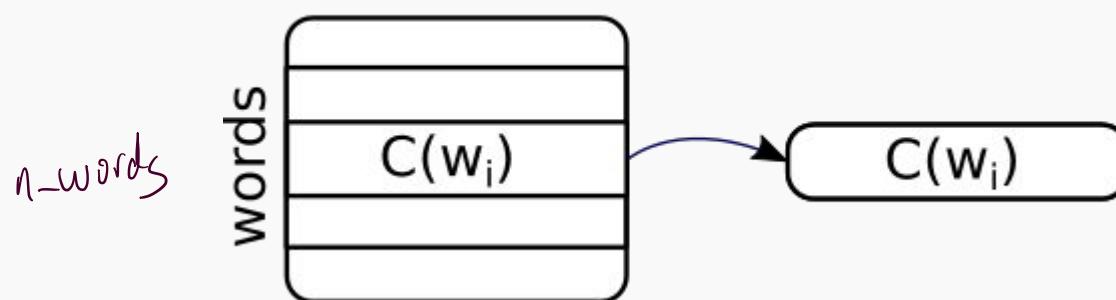
my cat **climbed** a tree

my cat **bridge** a tree

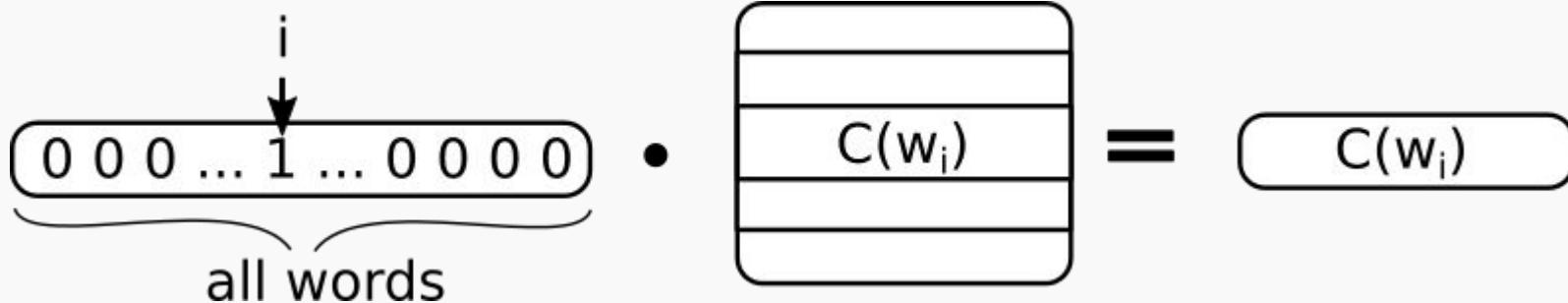
Embedding matrix

Two ways of thinking about the embedding matrix.

1. **Each row** contains a word embedding, which we need to extract
features *n-features*



2. It is a normal **weight matrix**, working with a 1-hot input vector



Word2vec

A popular tool for creating word embeddings.

Available from: <https://code.google.com/archive/p/word2vec/>

Preprocess the data!

- Tokenise
- Lowercase (usually)

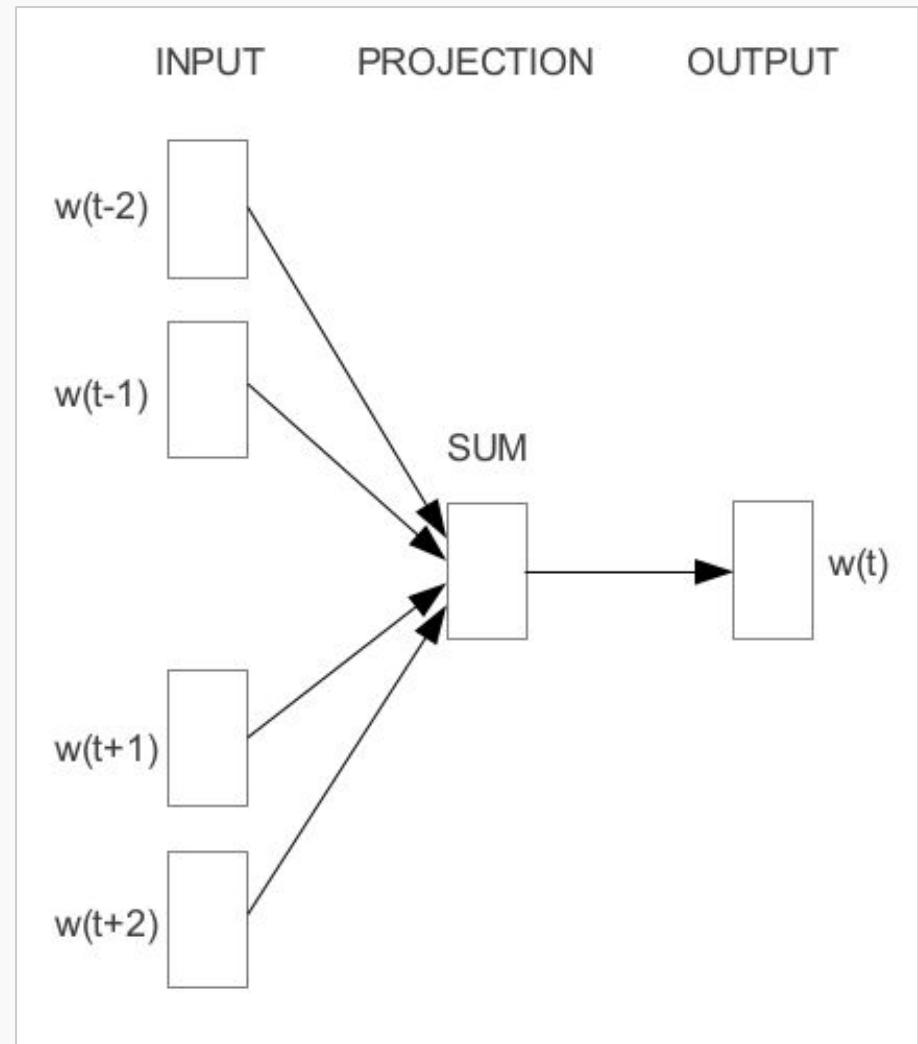
```
./word2vec -train input.txt -output vectors.txt [-cbow] 0 -size 100  
-window 5 -negative 5 -hs 0 -sample 1e-3 -threads 8
```

Continuous Bag-of-Words (CBOW) model

Predict the current word, based on the surrounding words

Mikolov et. al. 2013. *Efficient Estimation of Word Representations in Vector Space.*

my cat climbed a tree
t-2 t t+1 t+2

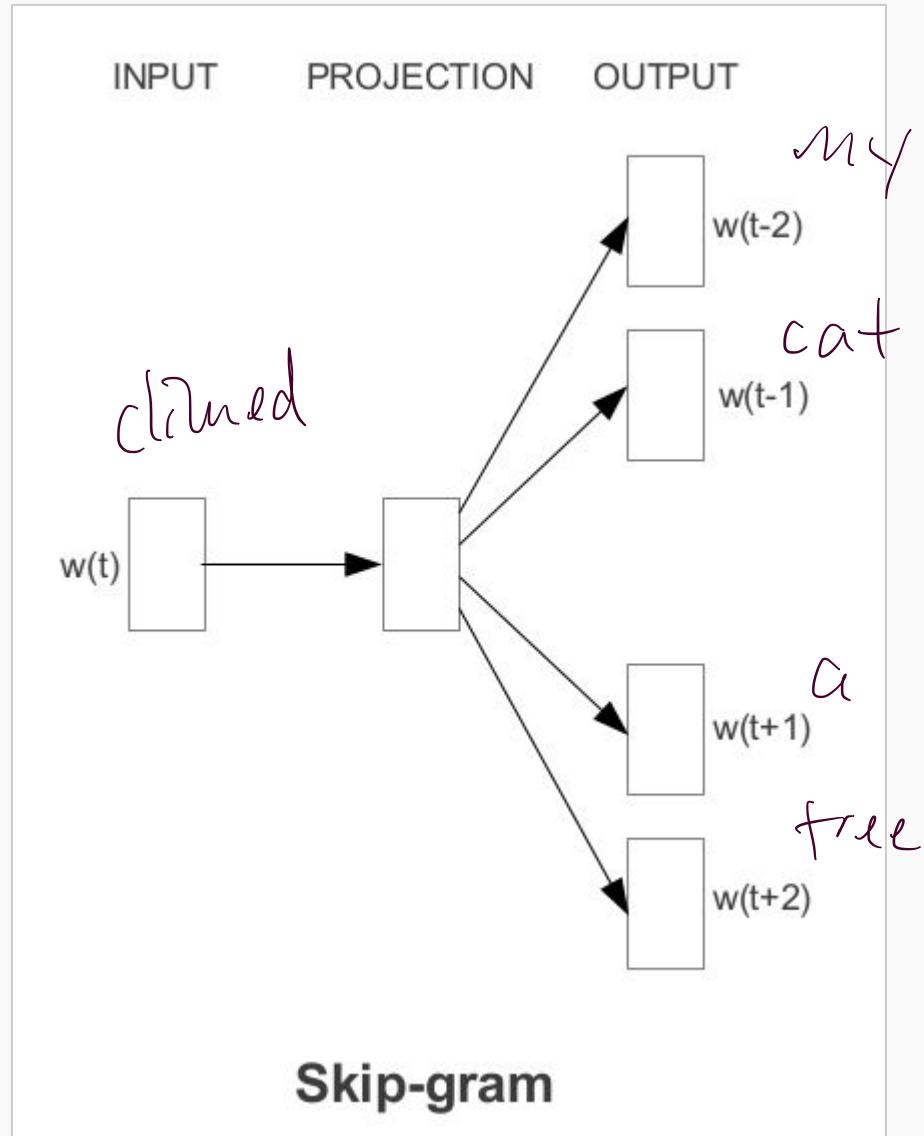


Skip-gram model

Predict the surrounding words, based on the current word.

Mikolov et. al. 2013. *Efficient Estimation of Word Representations in Vector Space.*

my cat climbed a tree
t-2 t-1 t t+1 t+2

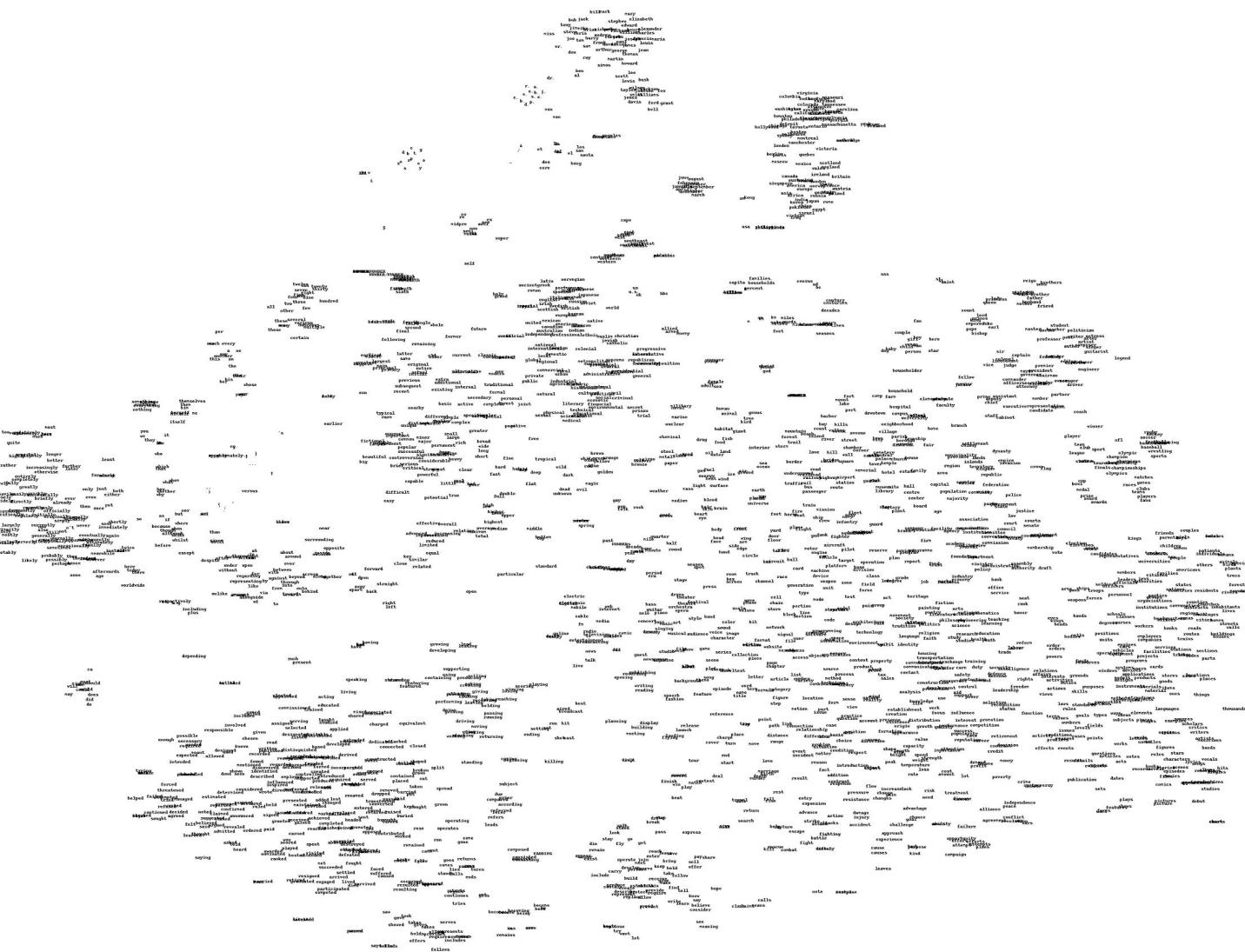


Word similarity

| | | | | | |
|---------------|---------------|--------------|------------------|--------------------|-------------------|
| FRANCE 454 | JESUS 1973 | XBOX 6909 | REDDISH 11724 | SCRATCHED 29869 | MEGABITS 87025 |
| 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 |

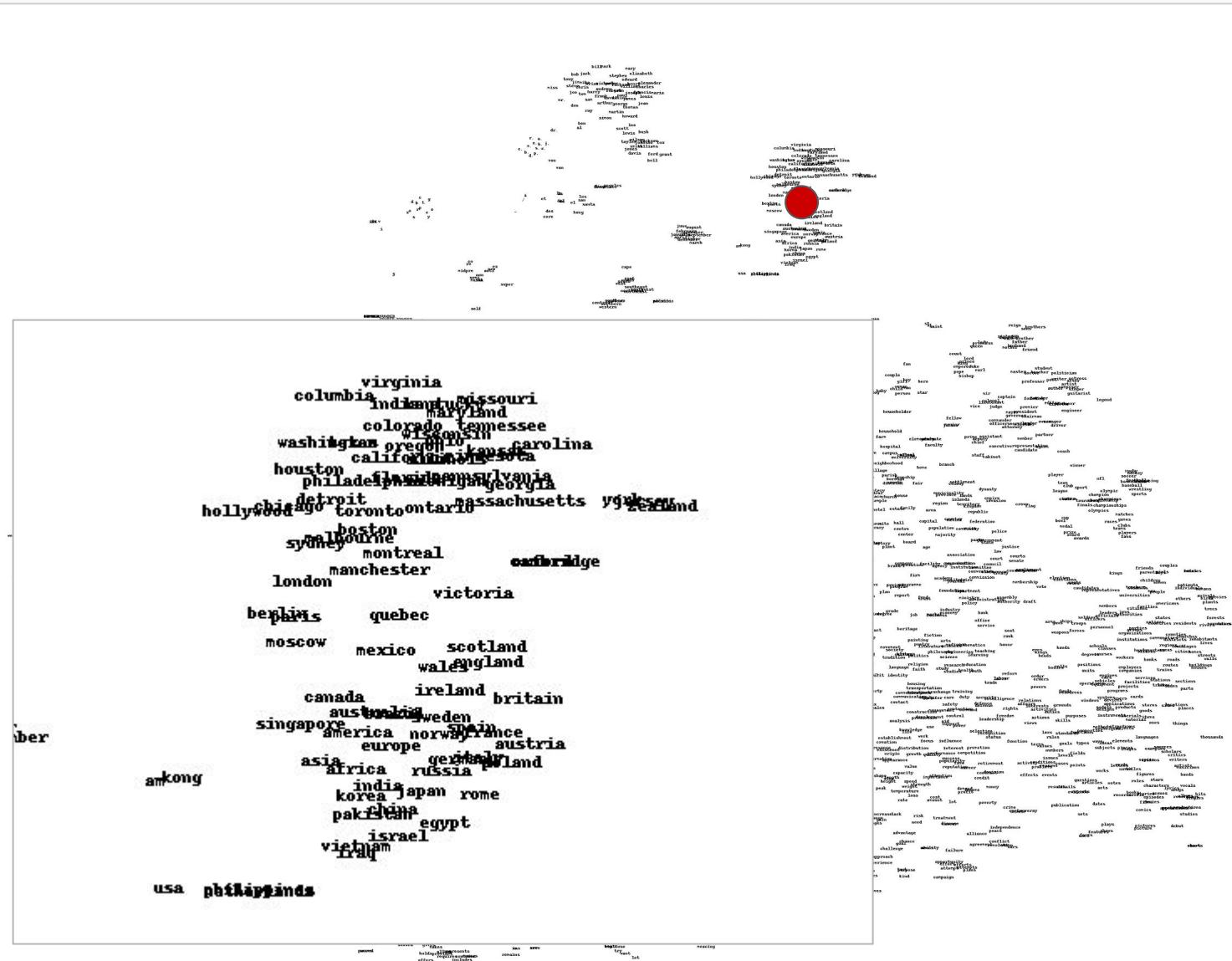
Collobert et. al. 2011. *Natural Language Processing (Almost) from Scratch.*

Word similarity



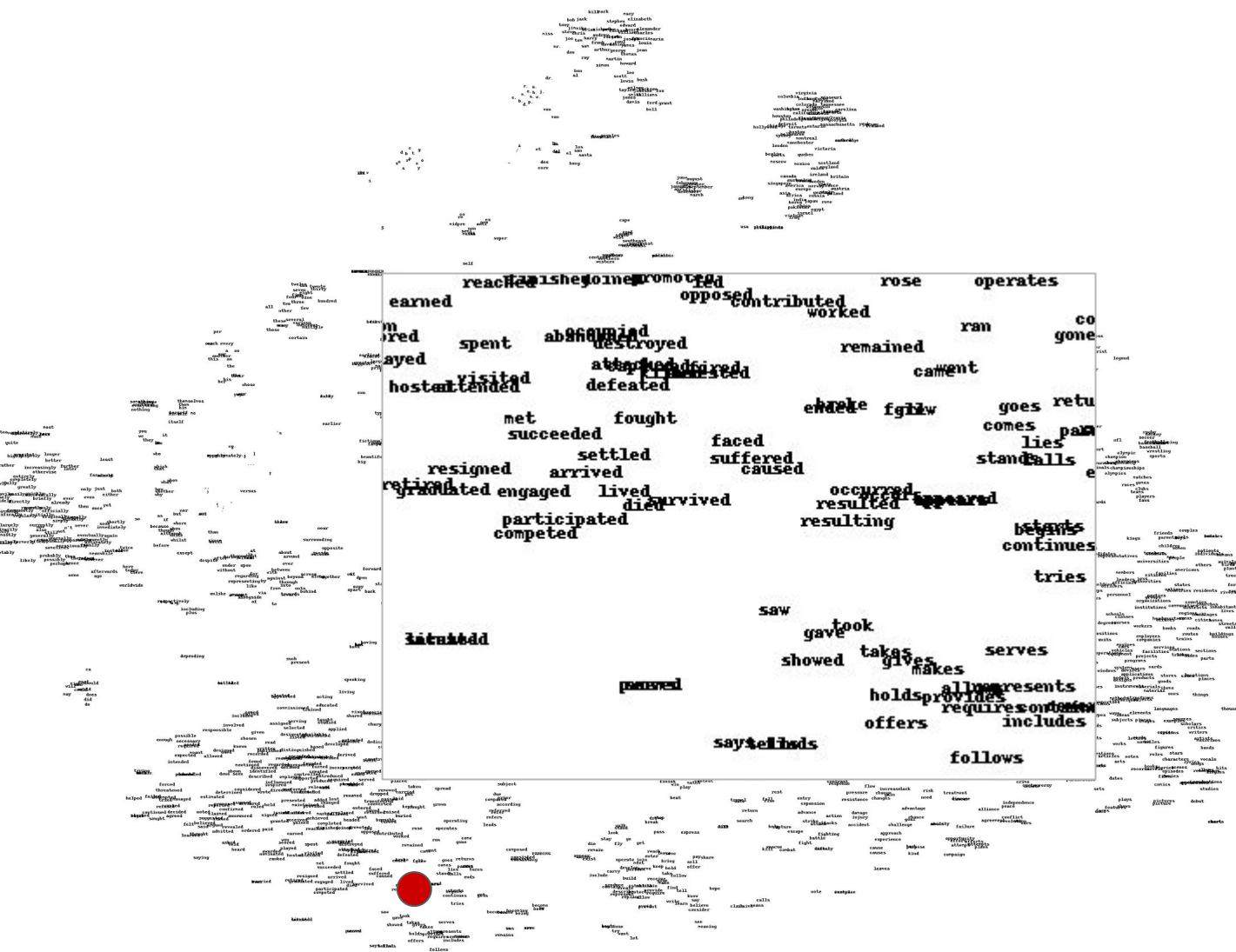
Joseph Turian

Word similarity



Joseph Turian

Word similarity



Joseph Turian

Analogy recovery

The task of **analogy recovery**. Questions in the form:

a is to *b* as *c* is to *d*

The system is given words *a*, *b*, *c*, and it needs to find *d*. For example:

‘apple’ is to ‘apples’ as ‘car’ is to ?

or

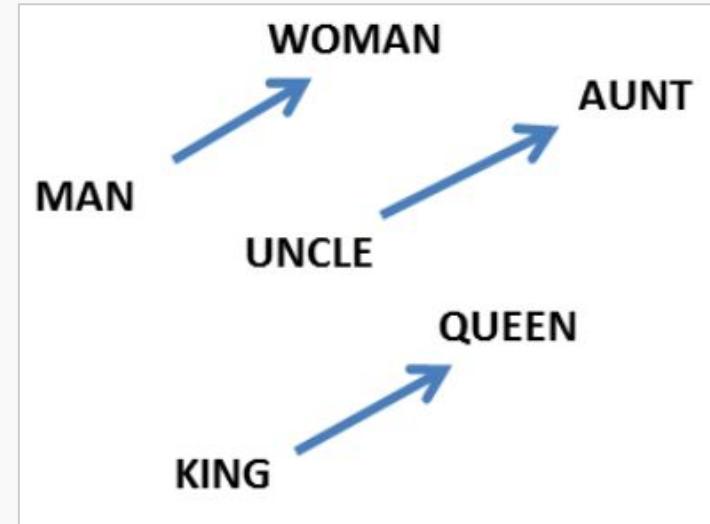
‘man’ is to ‘woman’ as ‘king’ is to ?

Mikolov et. al. 2013. *Efficient Estimation of Word Representations in Vector Space*.

Analogy recovery

Task: a is to b as c is to d

Idea: The direction of the relation should remain the same.



$$a - b \approx c - d$$

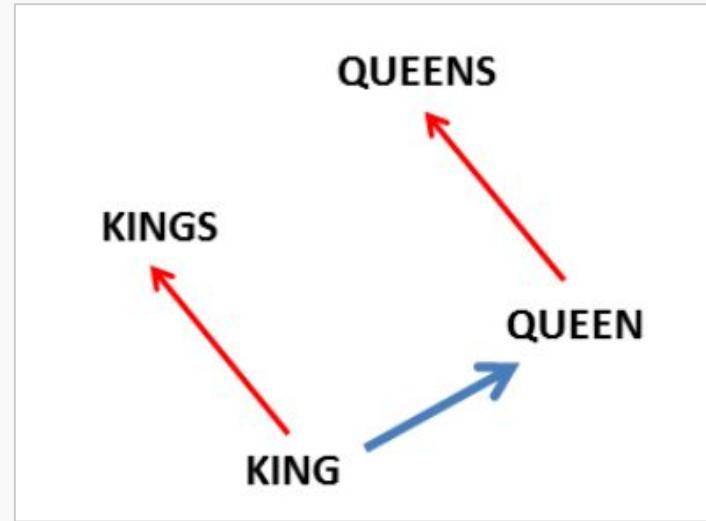
$$man - woman \approx king - queen$$

$$d_w = \operatorname{argmax}_{d'_w \in V} (\cos(a - b, c - d'))$$

Analogy recovery

Task: a is to b as c is to d

Idea: The offset of vectors should reflect their relation.



query $a - b \approx c - d$

$d \approx c - a + b$

$queen \approx king - man + woman$

$d_w = \operatorname{argmax}_{d'_w \in V} (\cos(d', c - a + b))$

relationship

Analogy recovery

| Relationship | Example 1 | Example 2 | Example 3 |
|--------------------------------|------------------------------|------------------------------|--|
| France - Paris big - bigger | Italy: Rome small: larger | Japan: Tokyo cold: colder | Florida: Tallahassee 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 |

Example output using word2vec vectors.

Word embeddings in practice

Word2vec is often used for pretraining.

- It will help your models start from an **informed** position
- Requires only **plain text** - which we have a lot
- Is very **fast** and easy to use
- Already **pretrained** vectors also available (trained on 100B words)



However, for best performance it is important to continue training (fine-tuning).

~~fine-tuning~~

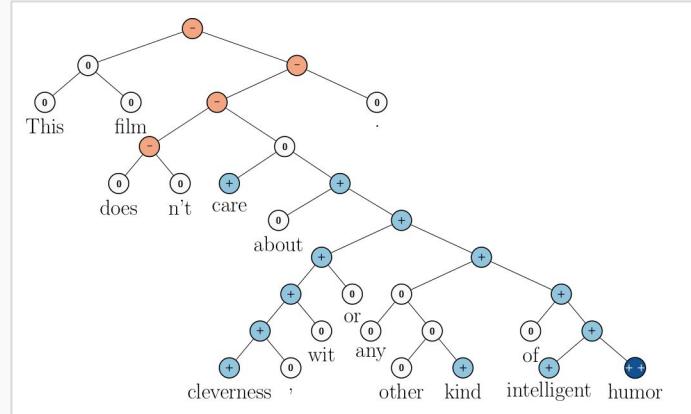
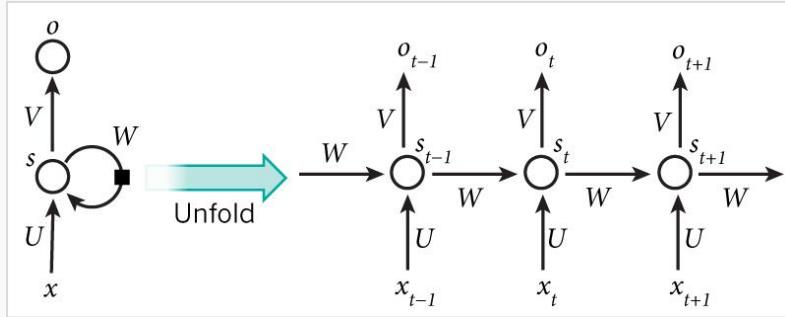
Raw word2vec vectors are good for predicting the surrounding words, but not necessarily for your specific task.

Simply treat the embeddings the same as other parameters in your model and keep updating them during training.

Word embeddings in practice

Word embeddings are the building blocks for higher-level models

RNNs



Count based vs direct prediction

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

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

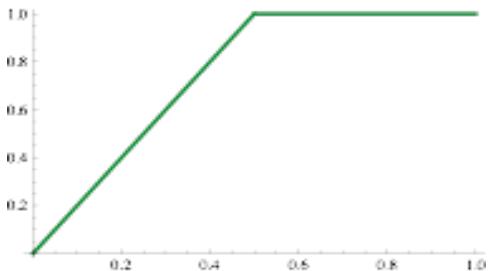
• NNLM, HLBL, RNN, Skip-gram/CBOW, (Bengio et al; Collobert & Weston; Huang et al; Mnih & Hinton; Mikolov et al; Mnih & Kavukcuoglu)

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

Combining the best of both worlds: GloVe

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

$f \sim$



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

Glove results

Nearest words to
frog:

1. frogs
2. toad
3. litoria
4. leptodactylidae
5. rana
6. lizard
7. eleutherodactylus



litoria



rana



leptodactylidae



eleutherodactylus

What to do with the two sets of vectors?

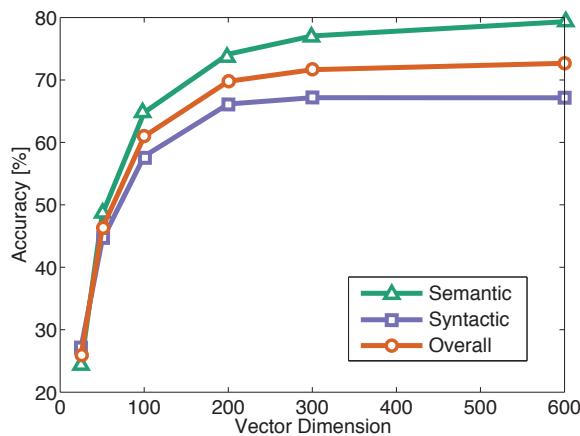
- We end up with U and V from all the vectors u and v (in columns)
- Both capture similar co-occurrence information. It turns out, the best solution is to simply sum them up:

$$X_{\text{final}} = U + V$$

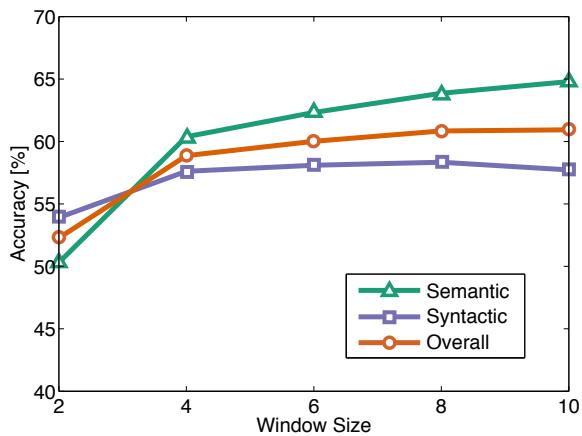
- One of many hyperparameters explored in *GloVe: Global Vectors for Word Representation* (Pennington et al. (2014))

Analogy evaluation and hyperparameters

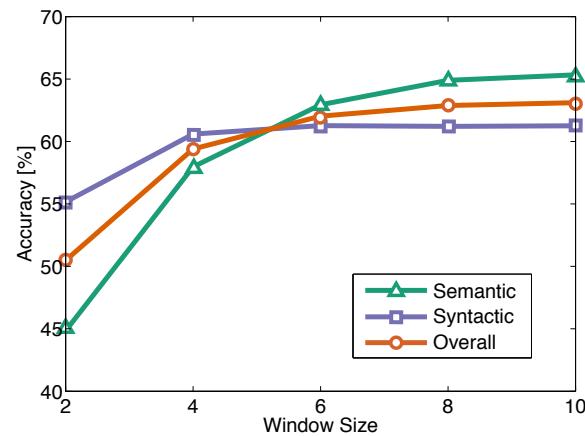
- Asymmetric context (only words to the left) are not as good



(a) Symmetric context



(b) Symmetric context



(c) Asymmetric context

- Best dimensions ~ 300 , slight drop-off afterwards
- But this might be different for downstream tasks!
- Window size of 8 around each center word is good for Glove vectors

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 10.00

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

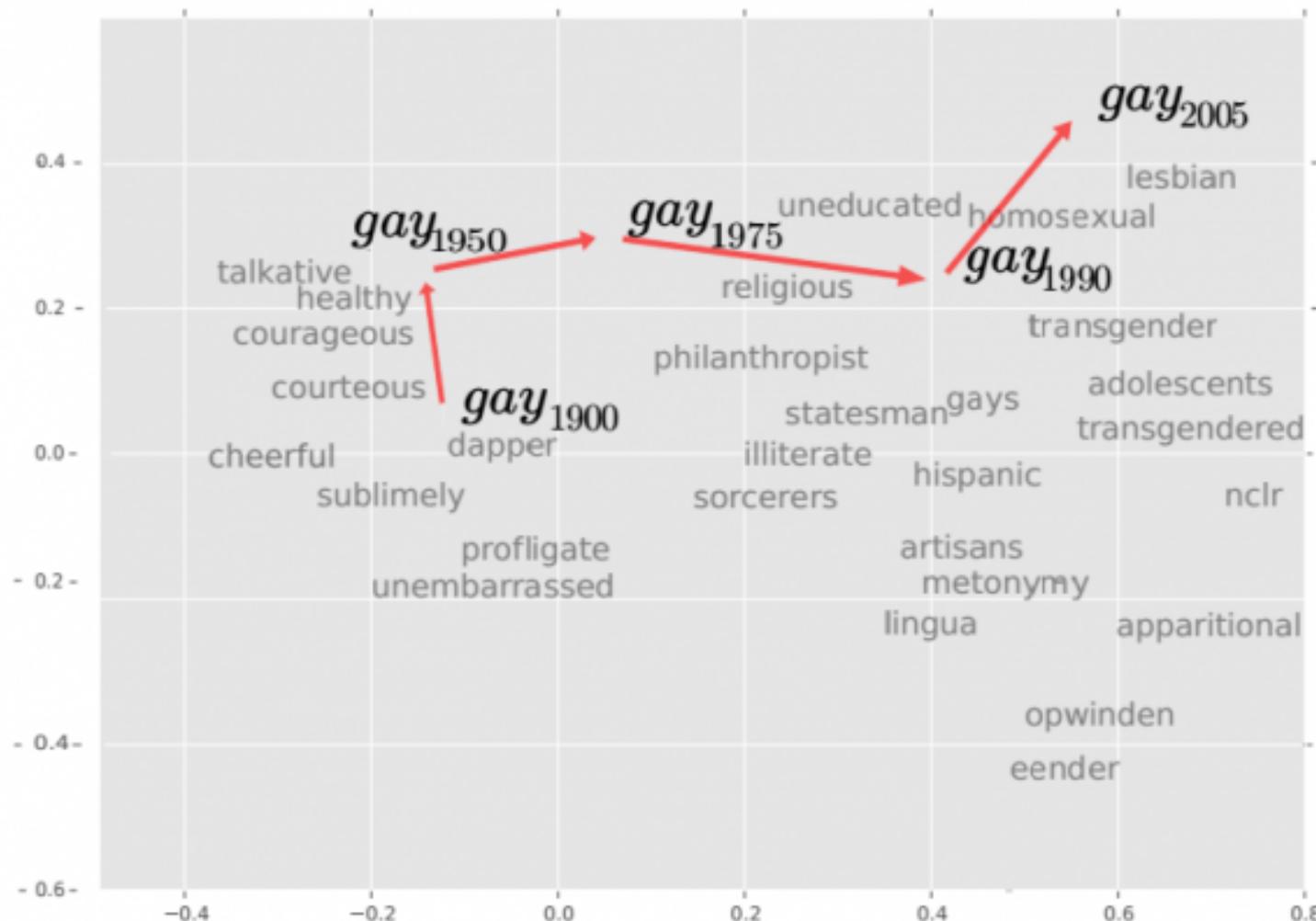
- Word vector distances and their correlation with human judgments

| 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 | 75.9 | 83.6 | 82.9 | 59.6 | 47.8 |
| CBOW* | 100B | 68.4 | 79.6 | 75.4 | 59.4 | 45.5 |

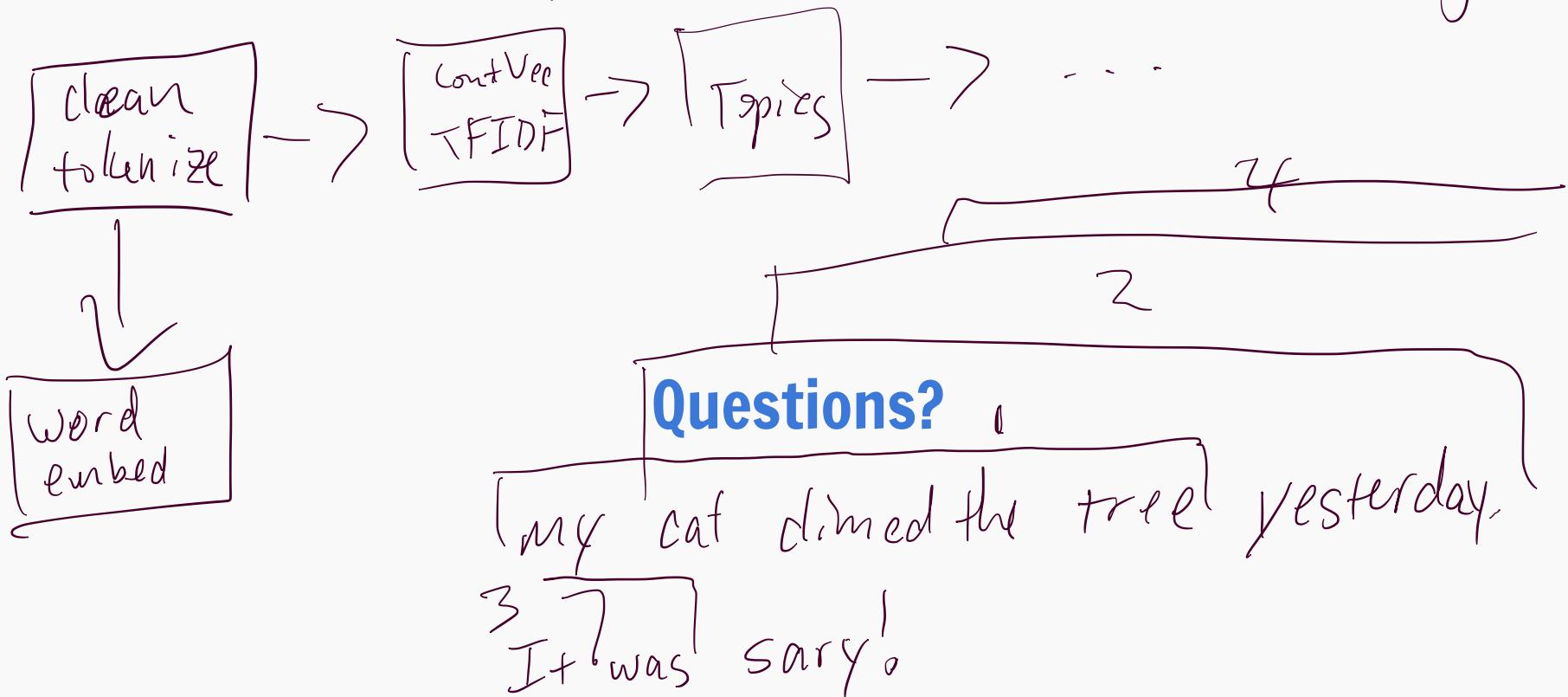
- Some ideas from Glove paper have been shown to improve skip-gram (SG) model also (e.g. sum both vectors)

pick out word + tag with year

Kulkarni, Al-Rfou, Perozzi, & Skiena (2015)



Pipeline ⚡ Separate from topic modeling



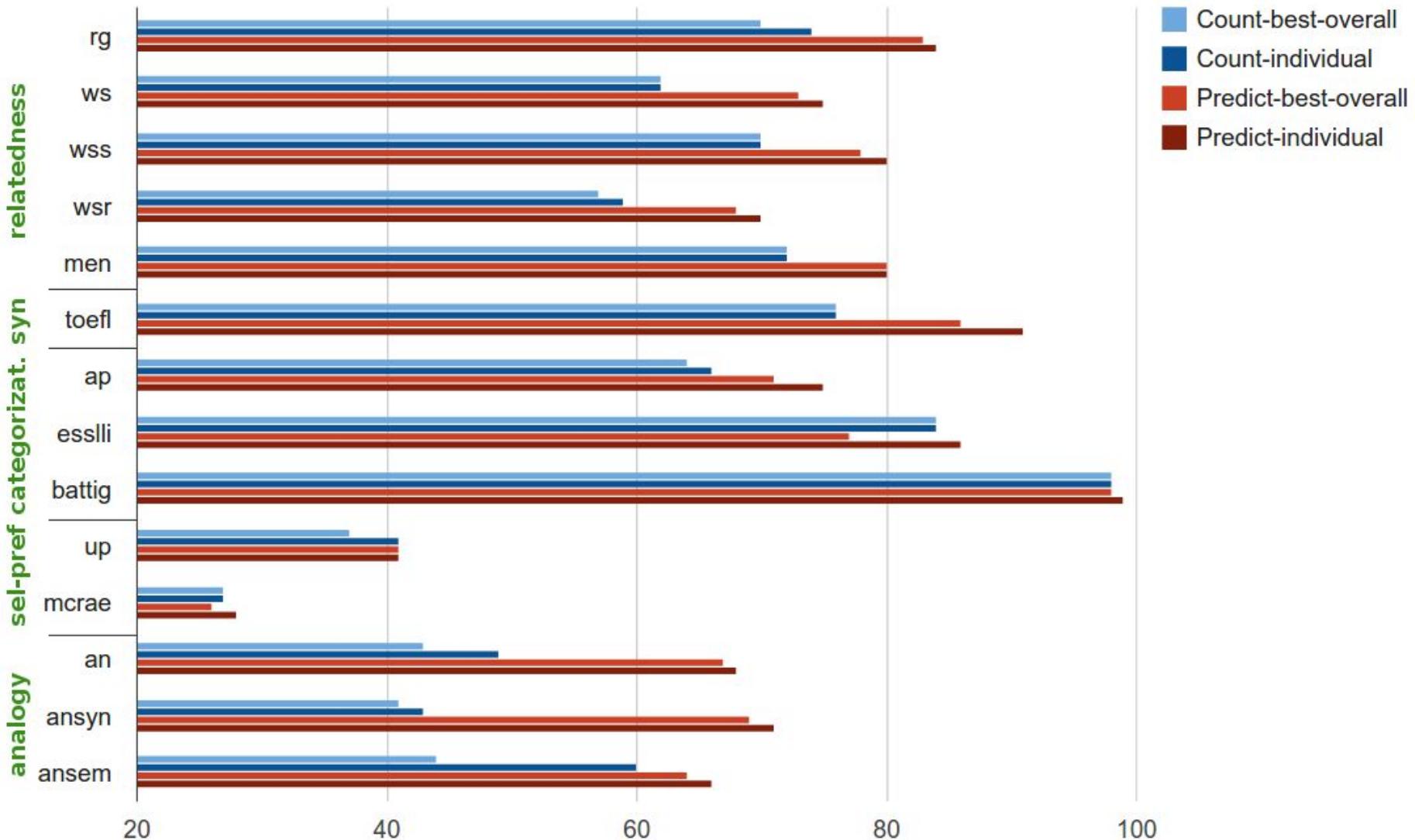
Count-based vs neural

Comparison of **count-based** and **neural** word vectors on 5 types of tasks and 14 different datasets:

1. Semantic relatedness
2. Synonym detection
3. Concept categorization
4. Selectional preferences
5. Analogy recovery

Baroni et. al. 2014. *Don't count, predict! A systematic comparison of context-counting vs. context-predicting semantic vectors.*

Count-based vs neural



Some of these conclusions are challenged by:

Levy et. al. 2015. *Improving Distributional Similarity with Lessons Learned from Word Embeddings*.

Count-based vs neural

The best parameter choices for counting models:

- window size 2 (bigger is not always better)
- weighted with PMI, not LMI
- no dimensionality reduction (not using SVD or NNMF)

The best parameters for the neural network model:

- window size 5
- negative sampling (not hierarchical softmax)
- subsampling of frequent words
- dimensionality 400