

We'll build a new model of meaning focusing on similarity

Each word = a vector

- Not just "word" or word45.

Similar words are "nearby in space"

A 2D coordinate system illustrating the semantic space for words. The horizontal axis represents a spectrum from negative to positive, with "not good" at the left end and "good" at the right end. The vertical axis represents a spectrum from neutral to strong, with "neutral" at the bottom and "extreme" at the top. Words are plotted based on their semantic features:

- Top row (positive, extreme): not good, dislike, bad, worst
- Middle row (positive, moderate): incredibly bad, worse
- Bottom row (neutral): neutral
- Left column (negative, extreme): not good, dislike, bad, worst
- Middle column (negative, moderate): incredibly bad, worse
- Right column (positive, extreme): good
- Diagonal row (moderate): neutral
- Bottom-left column (negative, extreme): not good, dislike, bad, worst
- Bottom-middle column (moderate): neutral
- Bottom-right column (positive, extreme): good

Words plotted include: to, by, 's, that, now, are, a, i, you, than, with, is.

A 2D coordinate system illustrating the semantic space for words. The horizontal axis represents a spectrum from negative to positive, with "not good" at the left end and "good" at the right end. The vertical axis represents a spectrum from neutral to strong, with "neutral" at the bottom and "extreme" at the top. Words are plotted based on their semantic features:

- Top row (positive, extreme): very good, incredibly good
- Middle row (positive, moderate): amazing, fantastic, wonderful
- Bottom row (neutral): terrific, nice
- Bottom-middle column (moderate): neutral
- Bottom-right column (positive, extreme): good

Words plotted include: very good, incredibly good, amazing, fantastic, wonderful, terrific, nice, good.

We define a word as a vector

Called an "embedding" because it's embedded into a space

The standard way to represent meaning in NLP

Fine-grained model of meaning for similarity

- NLP tasks like sentiment analysis
 - With words, requires **same** word to be in training and test
 - With embeddings: ok if **similar** words occurred!!!

2 kinds of embeddings

tf-idf

- The workhorse of Information Retrieval!
- A common baseline model
- **Sparse** vectors
- Words are represented by a simple function of the counts of nearby words

Word2vec

- **Dense** vectors
- Representation is created by training a classifier to distinguish nearby and far-away words

Review: words, vectors, and co-occurrence matrices

Term-document matrix

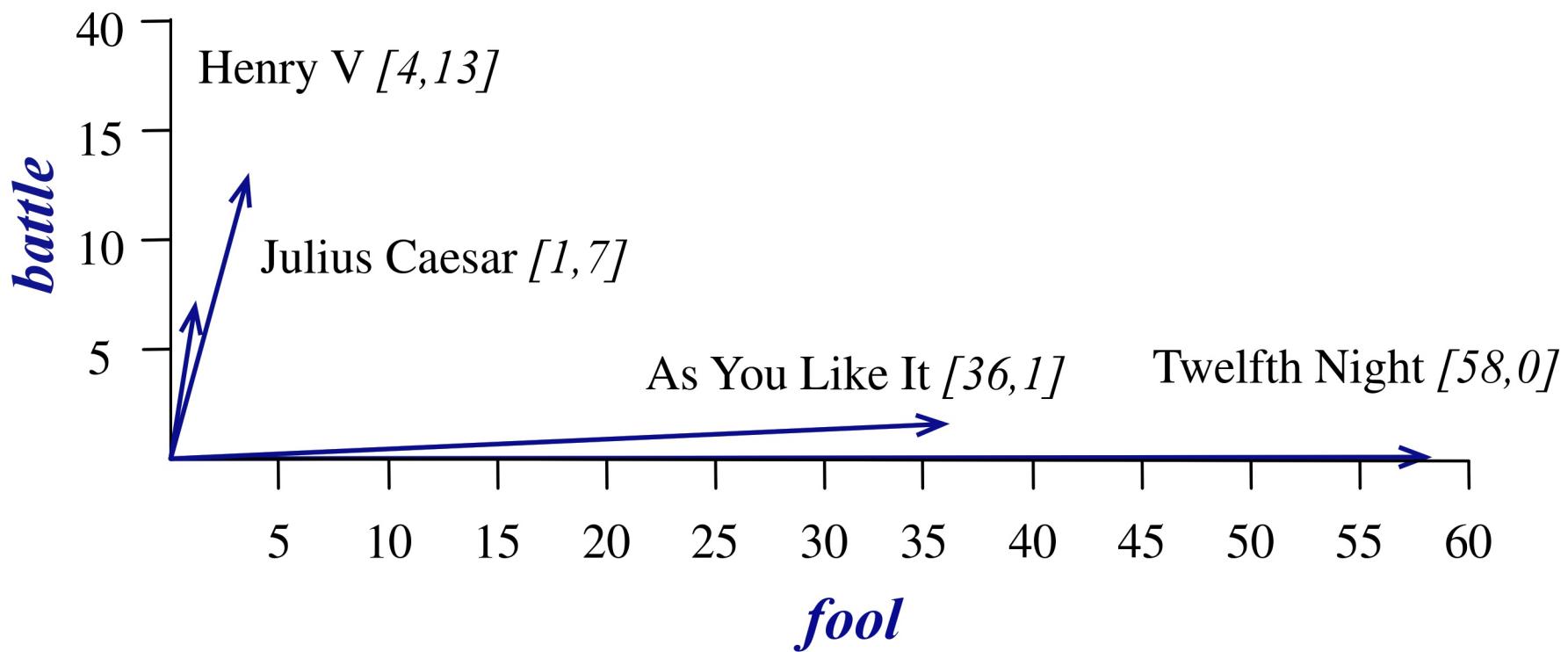
Each document is represented by a vector of words

The diagram illustrates the construction of a term-document matrix. At the top right, a 'corpus' is shown as a vertical stack of documents, each containing words. A window of size ± 1 is highlighted around the word 'battle' in the second document. Below this, a 'content' row shows the words 'battle', 'good', 'fool', and 'wit'. An arrow points from this row to the matrix below, labeled 'content $\rightarrow N$ contents'. The matrix itself is a grid where rows represent words and columns represent documents. The words 'battle', 'good', 'fool', and 'wit' are listed on the left, and the documents 'As You Like It', 'Twelfth Night', 'Julius Caesar', and 'Henry V' are listed at the top. The matrix entries are numerical values representing the count of each word in each document. Red boxes highlight specific entries: 'battle' in 'As You Like It' (1), 'good' in 'Twelfth Night' (80), 'fool' in 'Julius Caesar' (1), and 'wit' in 'Henry V' (3). A legend at the bottom right defines the terms: 'N words' (vertical dashed line), 'Document Vectors' (horizontal arrows), 'N content' (diagonal arrow), 'Word cooccurrence Matrix' (boxed text), and 'word-content mat' (text).

	D1 battle	D2 good	D3 fool	D4 wit
As You Like It	1	0	8	2
Twelfth Night	14	80	62	13
Julius Caesar	36	58	1	89
Henry V	20	15	2	4
				3

Word-word matrix
word-content mat

Visualizing document vectors



Vectors are the basis of information retrieval

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

Vectors are similar for the two comedies
Different than the history

Comedies have more *fools* and *wit* and fewer *battles*.

New idea for word meaning: Words can be vectors too!!!

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

battle is "the kind of word that occurs in Julius Caesar and Henry V"

fool is "the kind of word that occurs in comedies, especially Twelfth Night"

More common: word-word matrix (or "term-context matrix")

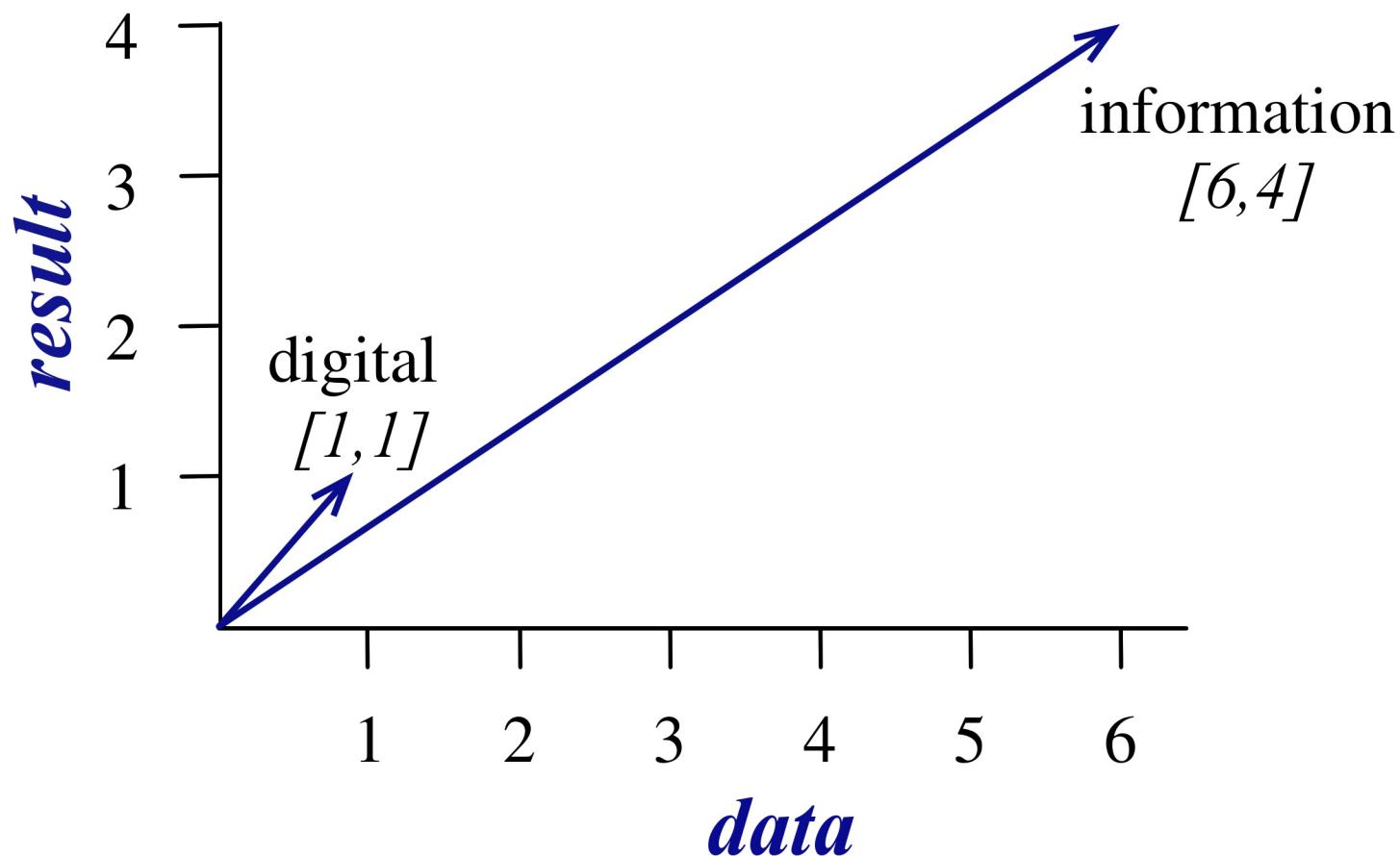
Two **words** are similar in meaning if their context vectors are similar

sugar, a sliced lemon, a tablespoonful of
oyment. Cautiously she sampled her first
vell suited to programming on the digital
for the purpose of gathering data and

apricot
pineapple
computer.
information

jam, a pinch ea
and another fru
In finding the o
necessary for th

	aardvark	computer	data	pinch	result	sugar	...
apricot	0	0	0	1	0	1	
pineapple	0	0	0	1	0	1	
digital	0	2	1	0	1	0	
information	0	1	6	0	4	0	



Cosine for computing word similarity

$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

v_i is the count for word v in context i

w_i is the count for word w in context i .

$\rightarrow \rightarrow$

$\rightarrow \rightarrow$

$$\vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \theta$$

$$\frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} = \cos \theta$$

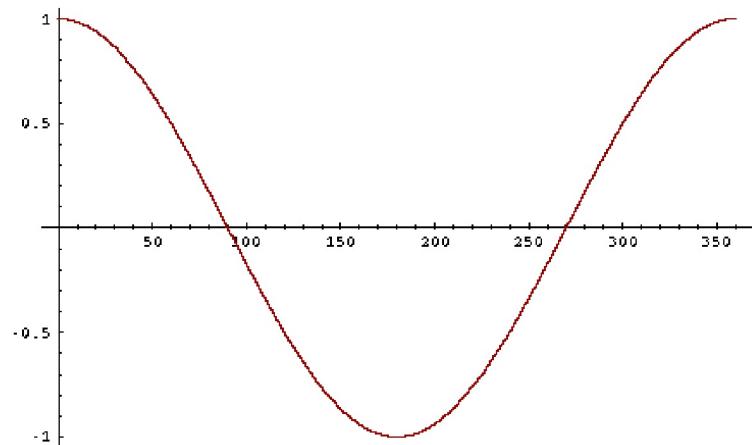
$\text{Cos}(v, w)$ is the cosine similarity of v and w

Cosine as a similarity metric

-1: vectors point in opposite directions

+1: vectors point in same directions

0: vectors are orthogonal



$$\cos(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|} = \frac{\vec{v}}{\|\vec{v}\|} \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

Which pair of words is more similar?

$$\text{cosine(apricot,information)} =$$

	large	data	computer
apricot \leftarrow	1	0	0
digital	0	1	2
information \leftarrow 1	6	1	1

$$\text{cosine(digital,information)} =$$

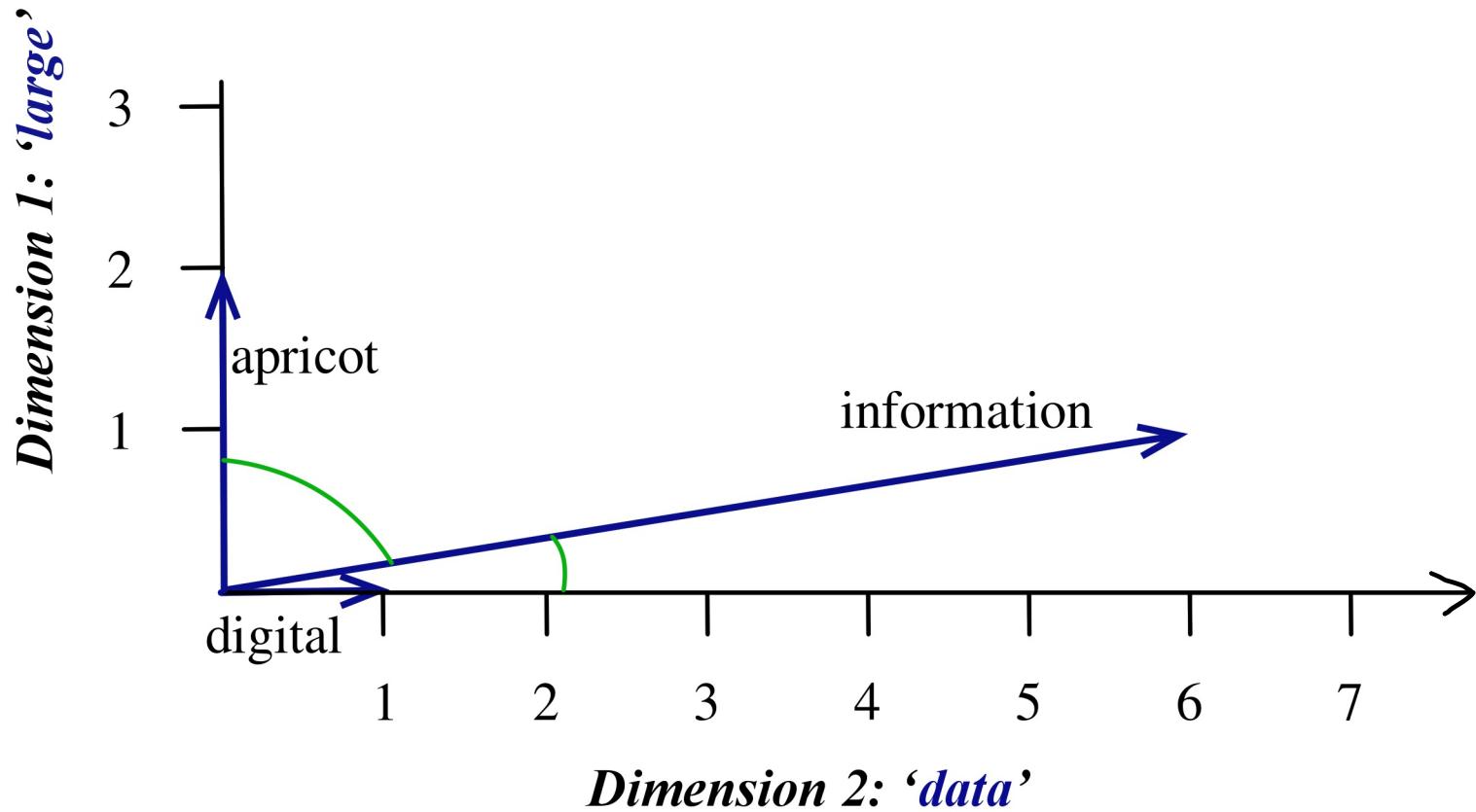
$$\text{cosine(apricot,digital)} =$$

$$\frac{1+0+0}{\sqrt{1+0+0} \sqrt{1+36+1}} = \frac{1}{\sqrt{38}} = .16$$

$$\frac{0+6+2}{\sqrt{0+1+4} \sqrt{1+36+1}} = \frac{8}{\sqrt{38} \sqrt{5}} = .58$$

$$\frac{0+0+0}{\sqrt{1+0+0} \sqrt{0+1+4}} = 0$$

Visualizing cosines (well, angles)



But raw frequency is a bad representation

- Frequency is clearly useful; if *sugar* appears a lot near *apricot*, that's useful information.
- But overly frequent words like *the*, *it*, or *they* are not very informative about the context
- Need a function that resolves this frequency paradox!

Many functions can help in reweighting the counts

tf-idf:

tf-idf value for word t in document d:

$$w_{t,d} = \text{tf}_{t,d} \times \text{idf}_t$$

Words like "the" or "good" have very low idf

Pointwise mutual information

- $\text{PMI}(w_1, w_2) = \log \frac{p(w_1, w_2)}{p(w_1)p(w_2)} = \frac{\text{Probability of Occuring together}}{\text{Probability of occurring independently}}$

See if words like "good" appear more often with "great" than we would expect by chance



Vector Semantics

Dense Vectors

Sparse versus dense vectors

- tf-idf vectors are
 - **long** (length $|V| = 20,000$ to $50,000$)
 - **sparse** (most elements are zero)
- Alternative: learn vectors which are
 - **short** (length 50-1000)
 - **dense** (most elements are non-zero)

Sparse versus 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 are distinct dimensions
 - a word with *car* as a neighbor and a word with *automobile* as a neighbor should be similar, but aren't
- **In practice, they work better**

Three methods for getting short dense vectors

- “Neural Language Model”-inspired models
 - Word2vec (skip-grams, CBOW), Glove
- Singular Value Decomposition (SVD)
 - A special case of this is called LSA – Latent Semantic Analysis
- Brown clustering



Vector Semantics

Embeddings inspired by neural language models: word2vec

Embeddings you can download!

Word2vec (Mikolov et al)

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



Glove (Pennington, Socher, Manning)

<http://nlp.stanford.edu/projects/glove/>



Word2vec

- Popular embedding method
- Very fast to train
- Code available on the web
- Idea: **predict** rather than **count**

Word2vec

- Instead of **counting** how often each word w occurs near "*apricot*"
- Train a classifier on a binary **prediction** task:
 - Is w likely to show up near "*apricot*"?
- We don't actually care about this task
 - But we'll take the learned classifier weights as the word embeddings

Word2Vec: Skip-Gram Task

Word2vec provides a variety of options.

Let's do

"skip-gram with negative sampling" (SGNS)

Approach: predict if candidate word c is a "neighbor"

- Treat the target word t and a true neighboring context word c as **positive examples**.
- Randomly sample other words in the lexicon to get negative examples
- Use logistic regression to train a classifier to distinguish those two cases
- Use the weights as the embeddings

Skip-Gram Training Data

Assume a +/- 2 word window, given training sentence:

- *I would like to go someplace nearby for lunch*

[target]

The training data:

- input/output pairs (centering on *go*)

Skip-Gram Training data

*I would like to go someplace nearby for
lunch*



Positive

{target, context}

{go, like}

{go, to}

{go, someplace}

{go, nearby}

Skip-Gram Training data

I would like to go someplace nearby for lunch



*Positive
{target, context}*

*{go, like}
{go, to}
{go, someplace}
{go, nearby}*

That's fine for positive data. But for training a binary classifier we need negative examples.

Let's sample other words from the lexicon (that don't occur with the target word in this context).

Skip-Gram Training data

I would like to go someplace nearby for lunch



Positive

{+ target, context}

{go, like}

{go, to}

{go, someplace}

{go, nearby}

Negative

{- target, context}

{go, aardvark}

{go, incubate}

{go, twelve}

{go, therefore}

Setup

Let's represent words as vectors of some length (say 300), randomly initialized.

So we start with $300 * V$ random parameters

Over the entire training set, we'd like to adjust those word vectors such that we

- Maximize the similarity of the **target word**, **context word** pairs (t, c) drawn from the positive data
- Minimize the similarity of the (t, c) pairs drawn from the negative data.

The classifier's goal

- Compute the probability that c is a real context word and not a fake noise word
- $P(\text{real} | t, c)$
- $P(\text{fake} | t, c)$

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:
 - $\text{Similarity}(t, c) \propto t \cdot c$
 - We'll need to normalize to get a⁷¹ probability
 - (cosine isn't a probability either)

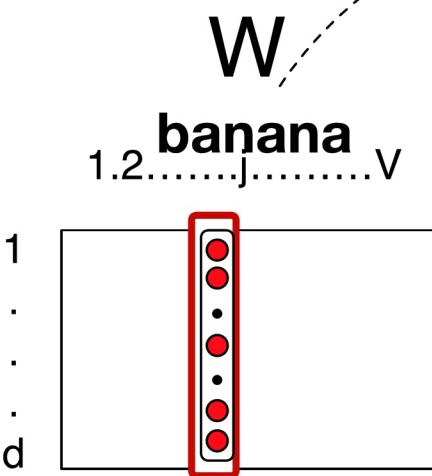
Turning dot products into probabilities

- $\text{Sim}(t,c) = t \cdot c$
- To turn this into a probability.
- We'll use the sigmoid from logistic regression:

$$P(+)|t, c) = \frac{1}{1 + e^{-\text{sim}(t,c)}}$$

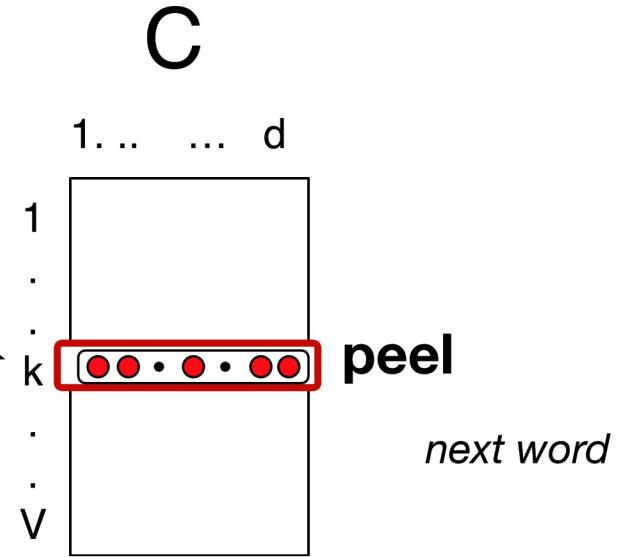
Learning the classifier

- How to learn?
 - Stochastic gradient descent!
- We'll adjust the word weights to
 - make the positive pairs more likely
 - and the negative pairs less likely,
 - over the entire training set.

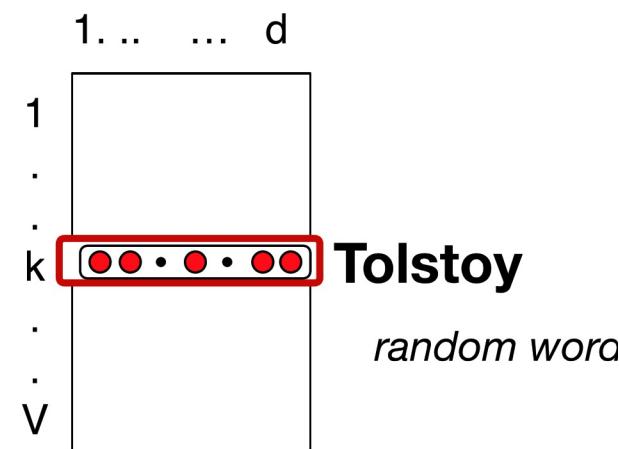


increase
similarity(banana , peel)

“...banana peel...”



decrease
similarity(banana , Tolstoy)



Objective Criteria

We want to maximize...

$$\sum_{(t,c) \in +} \log P(+|t, c) + \sum_{(t,c) \in -} \log P(-|t, c)$$

Maximize

- the + label for the pairs from the positive training data
- the – label for the pairs sample from the negative data.

Train using stochastic gradient descent

Summary: How to learn word2vec (skip-gram) embeddings

- Start with V random 300-dimensional vectors as initial embeddings
- Use logistic regression:
 - Take a corpus and take pairs of words that co-occur as positive examples
 - Take pairs of words that don't co-occur as negative examples
 - Train the classifier to distinguish these by slowly adjusting all the embeddings to improve the classifier performance
 - Throw away the classifier code and keep the embeddings.

Or in other words

- 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.

Conclusion

- **Concepts** or word senses
 - Have a complex many-to-many association with **words** (homonymy, multiple senses)
 - Have relations with each other
 - Synonymy, Antonymy, Superordinate
 - But are hard to define formally (necessary & sufficient conditions)
- **Embeddings** = vector models of meaning
 - More fine-grained than just a string or index
 - Especially good at modeling similarity/analogy
 - Just download them and use cosines!!
 - Can use sparse count models (tf-idf/PPMI) or dense predict models (word2vec, GLoVE)
 - Useful in practice but also encode cultural stereotypes
 - Can **debias** embeddings, and use them to **study bias**