# ANLP

## 04 - Words (Part II)

David Schlangen
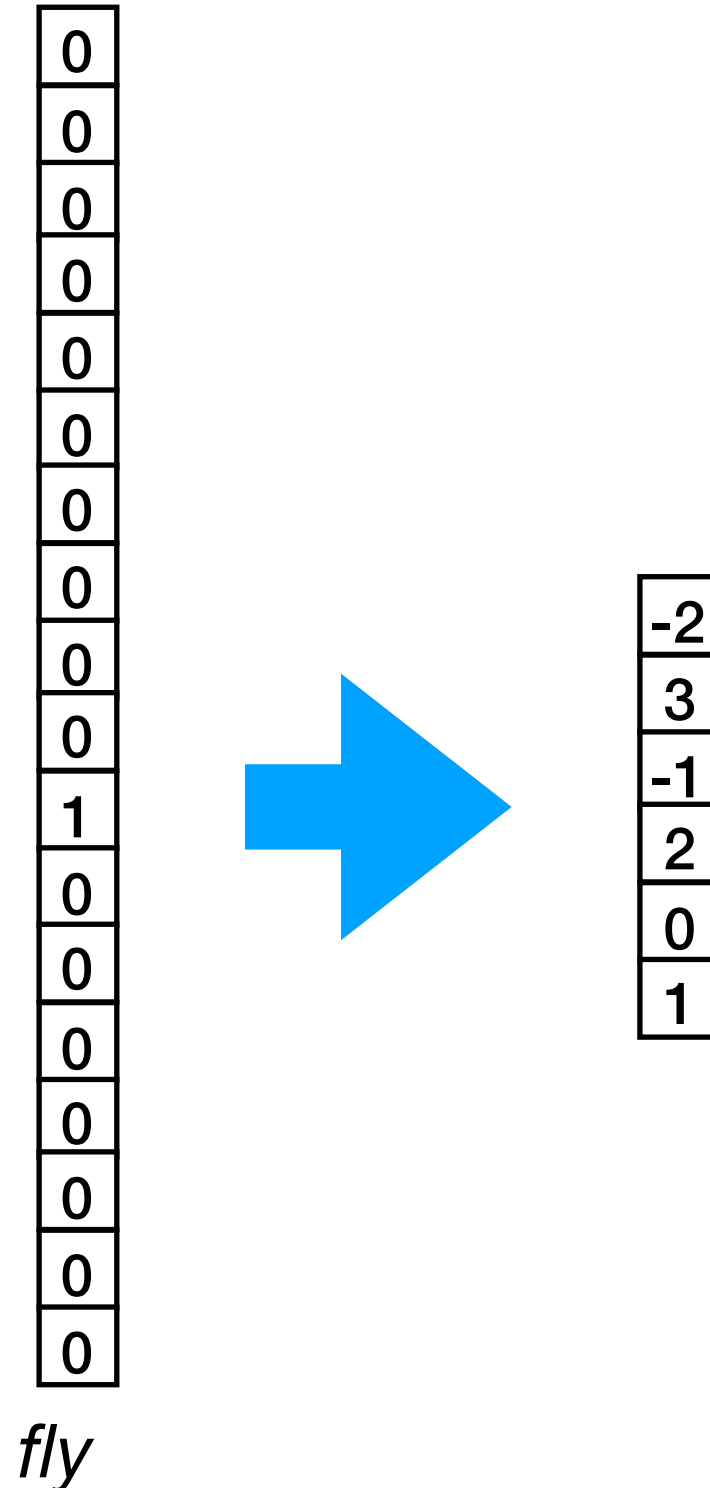University of Potsdam, MSc Cognitive Systems
Winter 2019 / 2020

# taking stock

- to represent words to machine, we now have one-hot vector (of size |V+1|), plus relations (= pairs of such vectors)…

- could now represent pair of words through 2 vectors + one hot vector over relation types

- but wouldn't it be nice to represent word identity and word meaning in the same way?



| 0 | 0 |
|---|---|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 1 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 1 |
| 0 | 0 |
| 0 | 0 |

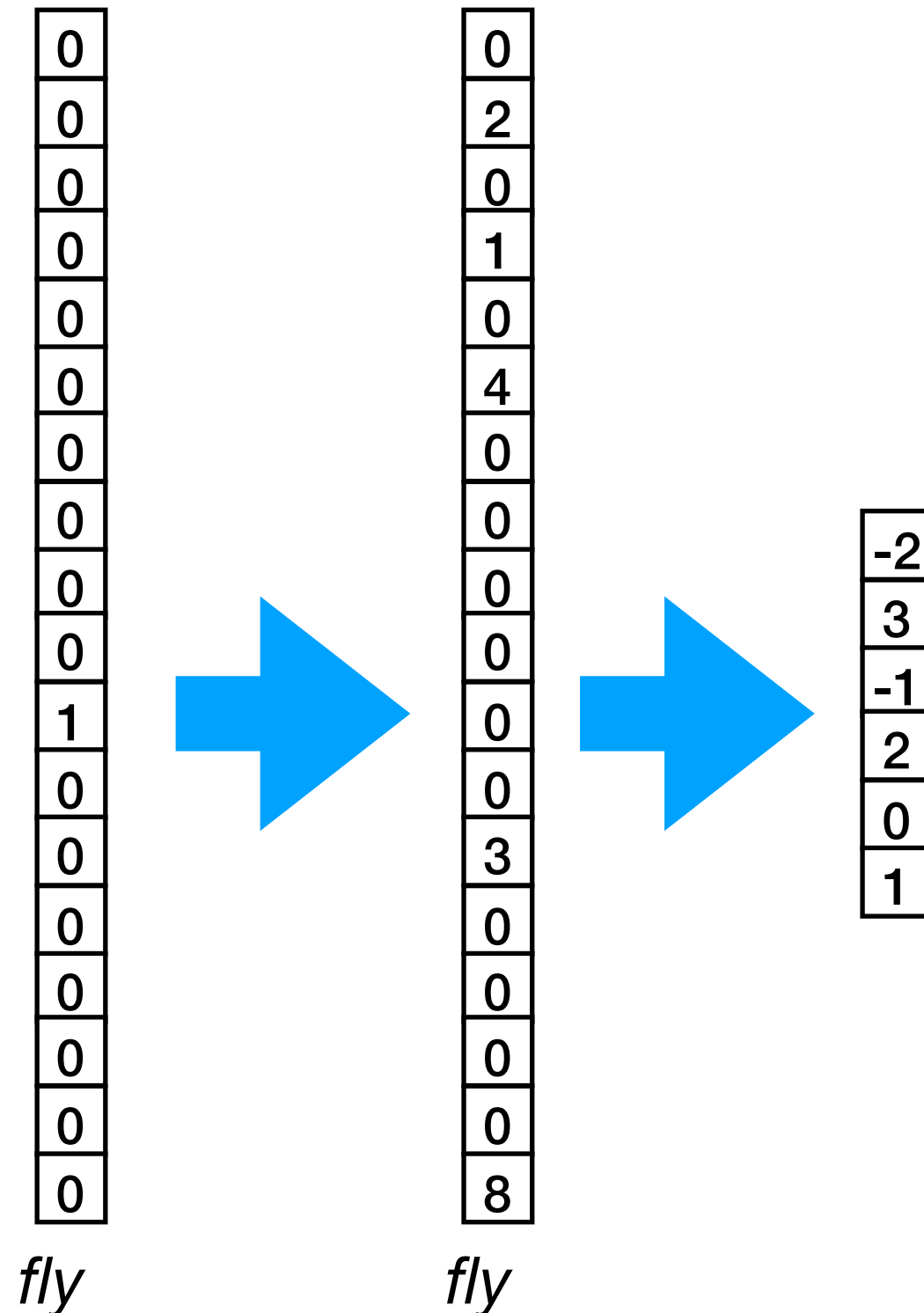is_a (  *fly*  ,  *insect*  )

# outlook

- we'll get this from *distributed representations*

- Hinton (1984): "Each entity is represented by a <span style="color:magenta">pattern of activity</span> distributed over many computing elements, and each computing element is involved in representing many different entities"

- the inherent *similarity* between vectors will (magically?) represent similarity between words!

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

*fly*

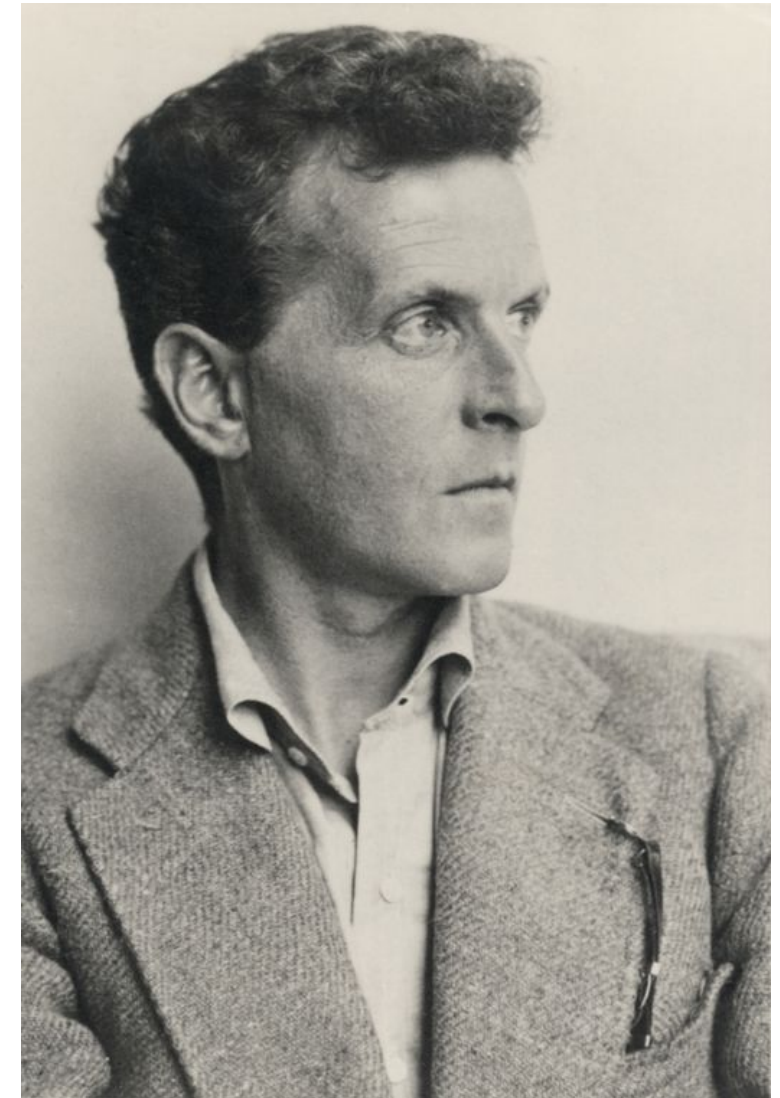| -2 |
|----|
| 3 |
| -1 |
| 2 |
| 0 |
| 1 |

# today

- we'll learn these representations from observing *use* of words

- in two steps: first, sparse representations (but still less sparse than one-hot), then dense representations

- in two different ways

# meaning and use

- "You say: the point isn't the word, but its meaning, and you think of the meaning as a thing of the same kind as the word, though also different from the word. Here the word, there the meaning. The money, and the cow that you can buy with it. (But contrast: money, and its use.)"

*(Wittgenstein 1953, Philosophical Investigations, §120)*

# meaning and use

- how can we get at *use* of a word?

- we can look at the *contexts* in which it was used

- impoverished, but easy to get form of record of such contexts: corpora, and the *other words* that form the context

- Harris (1954): "If A and B have almost identical environments we say that they are synonyms."

- Firth (1957): "You shall know a word by the company it keeps!"

# What does ongchoi mean?

Suppose you see these sentences:

- Ong choi is delicious **sautéed with garlic**.
- Ong choi is superb **over rice**
- Ong choi **leaves** with salty sauces

And you've also seen these:

- …spinach **sautéed with garlic over rice**
- Chard stems and **leaves** are **delicious**
- Collard greens and other **salty** leafy greens

Conclusion:

◦ Ongchoi is a leafy green like spinach, chard, or collard greens

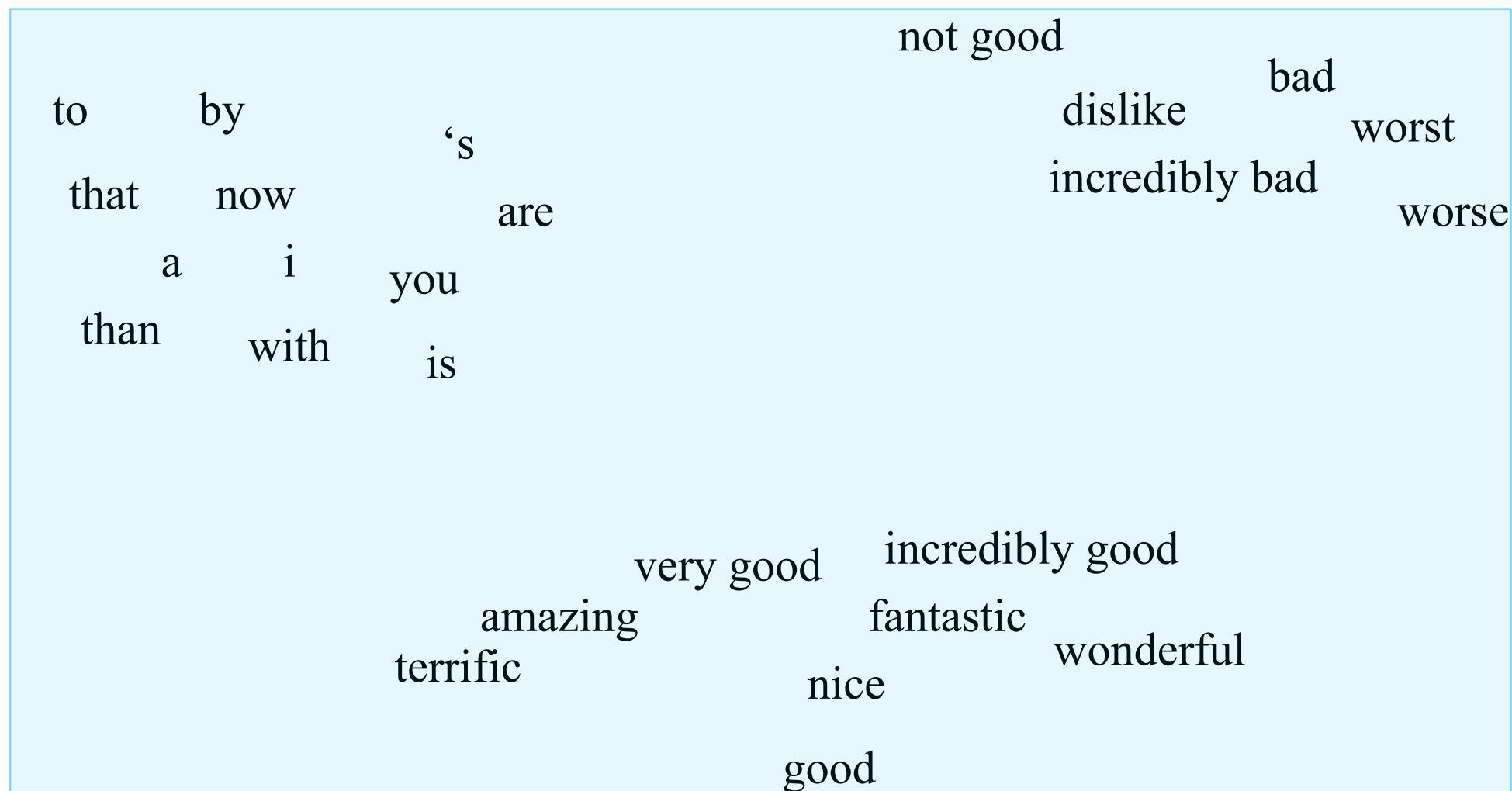# Ong choi: *Ipomoea aquatica* *"Water Spinach"*

# 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"

not good

bad

to        by

dislike        worst

's

incredibly bad

that        now                worse

are

a        i

you

than        with

is

very good        incredibly good

amazing        fantastic

terrific        wonderful

nice

good

# We'll introduce 2 kinds of vector spaces

## Tf-idf / count-based

◦ A common baseline model

◦ Sparse vectors

◦ Words are represented by a simple function of the counts of nearby words

## Word2vec / prediction-based

◦ Dense vectors

◦ Representation is created by training a classifier to distinguish nearby and far-away words
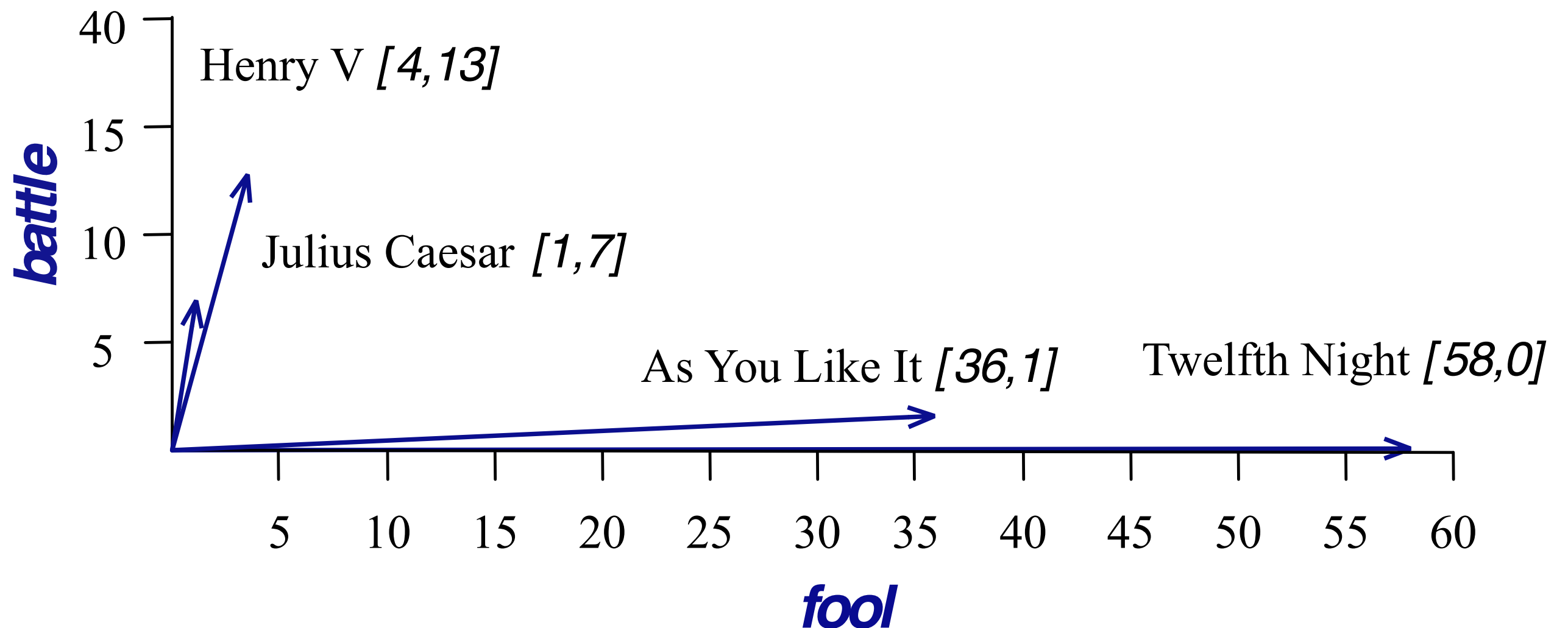
# count-based vector spaces

- idea: count events in context

- what kinds of contexts?

# Term-document matrix

Each document is represented by a vector of words

| | 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 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.

# 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"

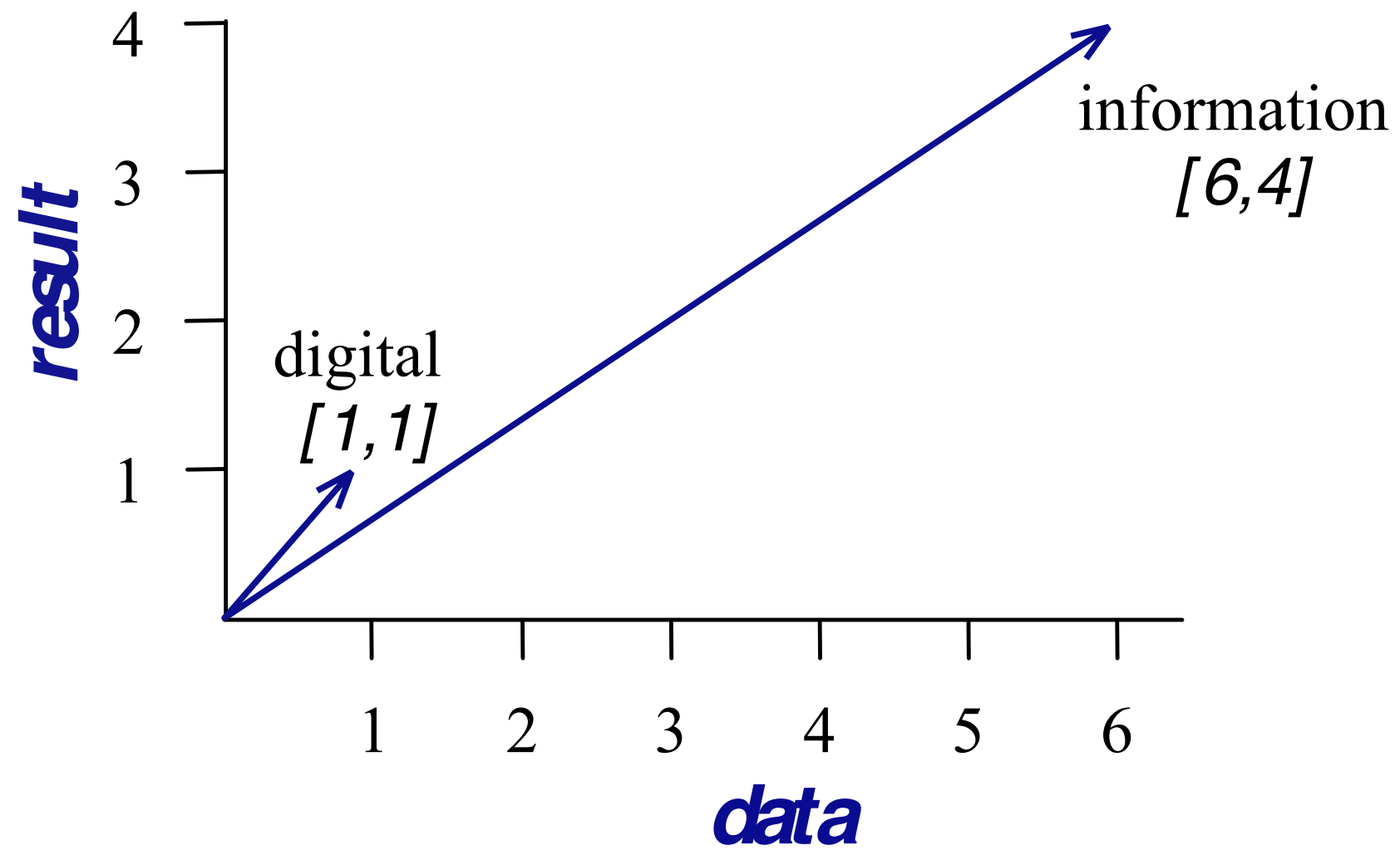# More common: word-word matrix (or "term-context matrix")

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

sugar, a sliced lemon, a tablespoonful of **apricot** jam, a pinch each of,
their enjoyment. Cautiously she sampled her first **pineapple** and another fruit whose taste she likened
well suited to programming on the digital **computer**. In finding the optimal R-stage policy from
for the purpose of gathering data and **information** necessary for the study authorized in the

|  | 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 | |

# Reminders from linear algebra

$$\text{dot-product}(\vec{v}, \vec{w}) = \vec{v} \cdot \vec{w} = \sum_{i=1}^{N} v_i w_i = v_1 w_1 + v_2 w_2 + \ldots + v_N w_N$$

vector length $\quad |\vec{v}| = \sqrt{\sum_{i=1}^{N} v_i^2}$

# Cosine for computing similarity

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

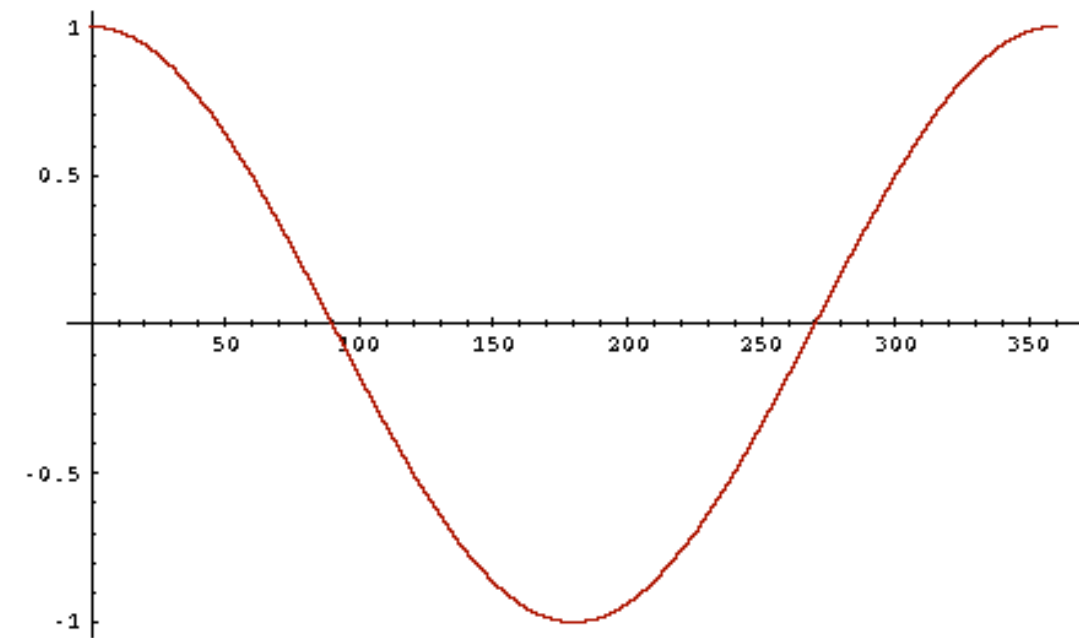$v_i$ is the count for word $v$ in context $i$

$w_i$ is the count for word $w$ in context $i$.

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

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

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

# Cosine as a similarity metric

-1: vectors point in opposite directions

+1:  vectors point in same directions

0: vectors are orthogonal

Frequency is non-negative, so  cosine range 0-1

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

|  | large | data | computer |
|---|---|---|---|
| apricot | 1 | 0 | 0 |
| digital | 0 | 1 | 2 |
| information | 1 | 6 | 1 |

Which pair of words is more similar?

cosine(apricot,information) =

cosine(digital,information) =

cosine(apricot,digital) =

$$\sqrt{}$$

# 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!

# tf-idf: combine two factors

**tf: term frequency**. frequency count (usually log-transformed):

$$\text{tf}_{t,d} = \begin{cases} 1 + \log_{10} \text{count}(t,d) & \text{if } \text{count}(t,d) > 0 \\ 0 & \text{otherwise} \end{cases}$$

**Idf: inverse document frequency: tf-**

Total # of docs in collection

$$\text{idf}_i = \log \left( \frac{N}{\text{df}_i} \right)$$

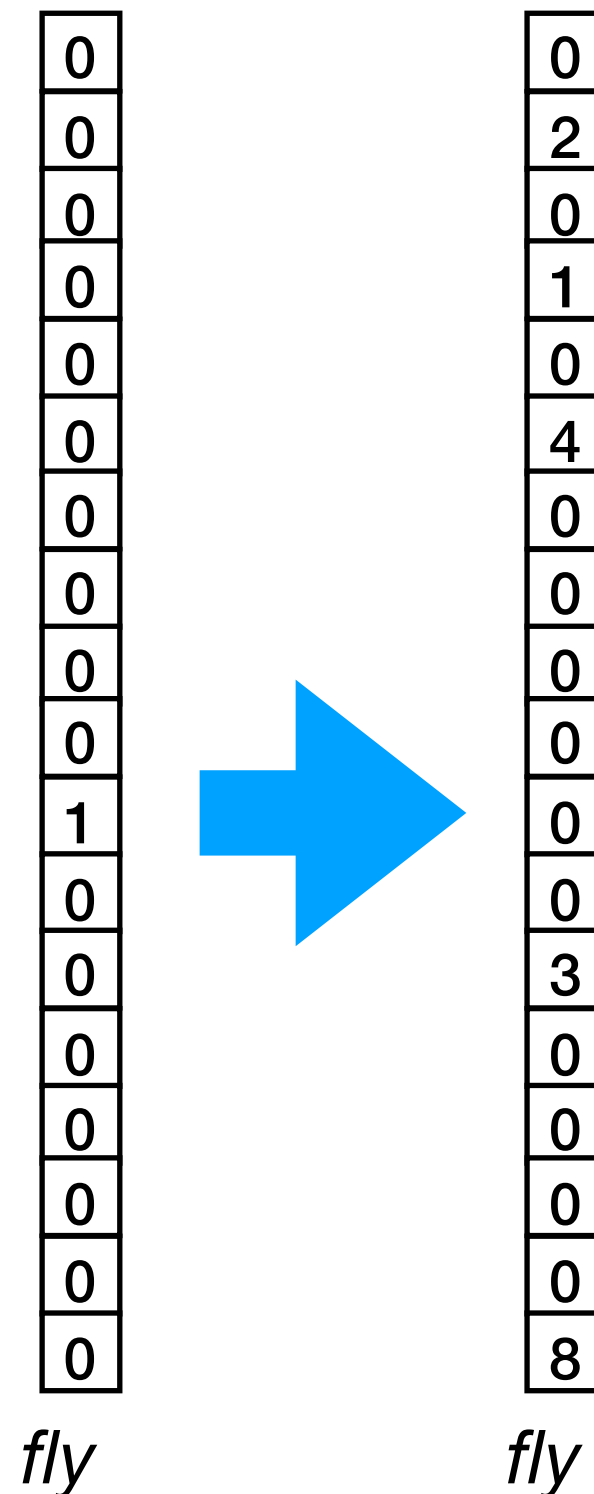Words like "the" or "good" have very low idf

# of docs that have word i

tf-idf value for word t in document d:

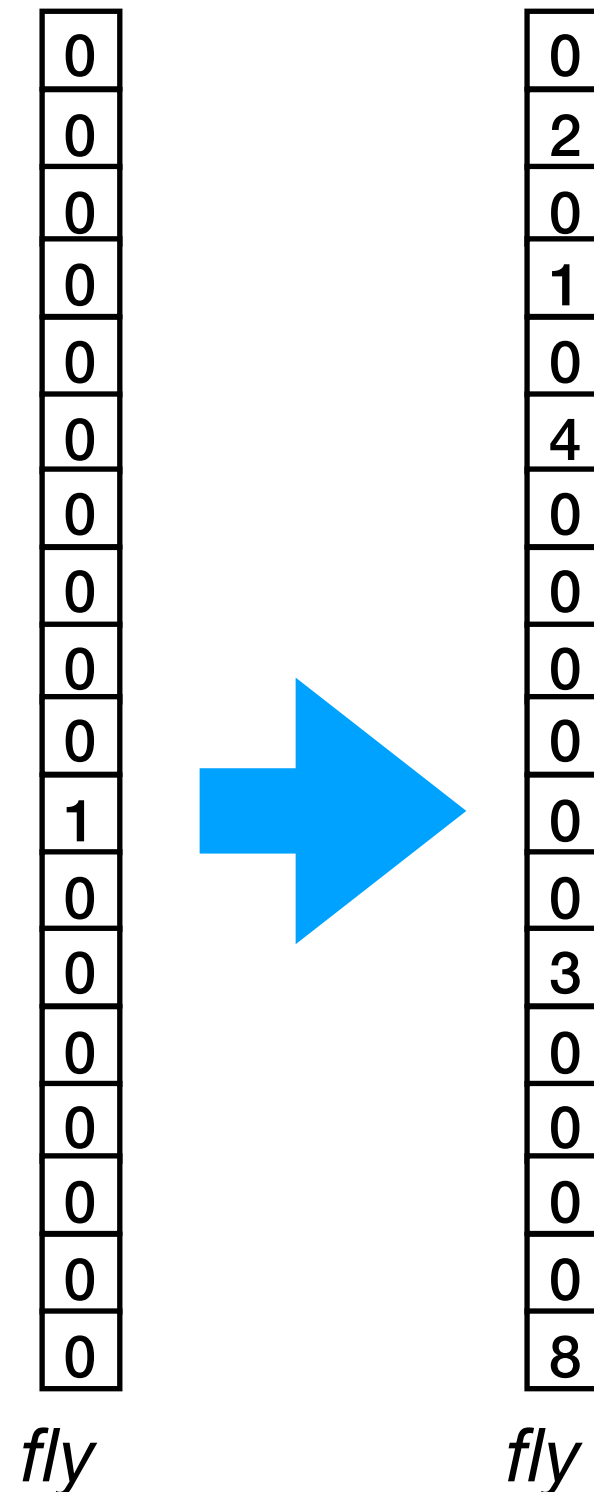$$w_{t,d} = \text{tf}_{t,d} \times \text{idf}_t$$

# summary count-based

- we now have less sparse representation for word, where dimensions are tf-idf transformed counts of context words.

- (alternative method instead of tf-idf: pointwise mutual information. read slides / background on your own.)

- we now have a relation defined on representation space (cosine), which captures (to a degree) relation between objects (semantic similarity)!

- addition still makes sense. can represent document as *centroid* vector (average)

| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

*fly*

| 0 |
| 2 |
| 0 |
| 1 |
| 0 |
| 4 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 3 |
| 0 |
| 0 |
| 0 |
| 0 |
| 8 |

*fly*

# summary count-based

- majority of components of word vector will still be…?

  - 0   (most other words don't co-occur with given word)

- can we make this yet less sparse?

- yes. Use techniques of *dimensionality reduction.* Popular choice: *singular value decomposition* over word/context matrix.

- or, use prediction methods (next)

| fly | | fly |
|---|---|---|
| 0 | | 0 |
| 0 | | 2 |
| 0 | | 0 |
| 0 | | 1 |
| 0 | | 0 |
| 0 | | 4 |
| 0 | | 0 |
| 0 | | 0 |
| 0 | | 0 |
| 0 | | 0 |
| 1 | | 0 |
| 0 | | 0 |
| 0 | | 3 |
| 0 | | 0 |
| 0 | | 0 |
| 0 | | 0 |
| 0 | | 8 |

# An alternative to tf-idf

Ask whether a context word is **particularly informative** about the target word.

◦ Positive Pointwise Mutual Information (PPMI)

# Pointwise Mutual Information

**Pointwise mutual information**:

Do events x and y co-occur more than if they were independent?

$$\text{PMI}(X, Y) = \log_2 \frac{P(x,y)}{P(x)P(y)}$$

**PMI between two words**: (Church & Hanks 1989)

Do words x and y co-occur more than if they were independent?

$$\text{PMI}(word_1, \ word_2) = \log_2 \frac{P(word_1, \ word_2)}{P(word_1)P(word_2)}$$

# Positive Pointwise Mutual Information

- PMI ranges from $-\infty$ to $+\infty$
- But the negative values are problematic
  - Things are co-occurring **less than** we expect by chance
  - Unreliable without enormous corpora
    - Imagine w1 and w2 whose probability is each $10^{-6}$
    - Hard to be sure p(w1,w2) is significantly different than $10^{-12}$
  - Plus it's not clear people are good at "unrelatedness"
- So we just replace negative PMI values by 0
- Positive PMI (PPMI) between word1 and word2:

$$\text{PPMI}(word_1,\ word_2) = \max\left(\log_2\frac{P(word_1,\ word_2)}{P(word_1)P(word_2)}, 0\right)$$

# Computing PPMI on a term-context matrix

Matrix $F$ with $W$ rows (words) and $C$ columns (contexts)

$f_{ij}$ is # of times $w_i$ occurs in context $c_j$

| | 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 |

$$p_{ij} = \frac{f_{ij}}{\sum\limits_{i=1}^{W}\sum\limits_{j=1}^{C} f_{ij}}$$

$$p_{i*} = \frac{\sum\limits_{j=1}^{C} f_{ij}}{\sum\limits_{i=1}^{W}\sum\limits_{j=1}^{C} f_{ij}}$$

$$p_{*j} = \frac{\sum\limits_{i=1}^{W} f_{ij}}{\sum\limits_{i=1}^{W}\sum\limits_{j=1}^{C} f_{ij}}$$

$$pmi_{ij} = \log_2 \frac{p_{ij}}{p_{i*}\, p_{*j}}$$

$$ppmi_{ij} = \begin{cases} pmi_{ij} & \text{if } pmi_{ij} > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$p_{ij} = \frac{f_{ij}}{\displaystyle\sum_{i=1}^{W}\sum_{j=1}^{C} f_{ij}}$$

**Count(w,context)**

| | computer | data | pinch | result | sugar |
|---|---|---|---|---|---|
| apricot | 0 | 0 | 1 | 0 | 1 |
| pineapple | 0 | 0 | 1 | 0 | 1 |
| digital | 2 | 1 | 0 | 1 | 0 |
| information | 1 | 6 | 0 | 4 | 0 |

p(w=information,c=data) =  6/19  = .32

p(w=information) =  11/19  = .58

p(c=data) =  7/19  = .37

$$p(w_i) = \frac{\displaystyle\sum_{j=1}^{C} f_{ij}}{N} \qquad p(c_j) = \frac{\displaystyle\sum_{i=1}^{W} f_{ij}}{N}$$

**p(w,context)**        **p(w)**

| | computer | data | pinch | result | sugar | |
|---|---|---|---|---|---|---|
| apricot | 0.00 | 0.00 | 0.05 | 0.00 | 0.05 | 0.11 |
| pineapple | 0.00 | 0.00 | 0.05 | 0.00 | 0.05 | 0.11 |
| digital | 0.11 | 0.05 | 0.00 | 0.05 | 0.00 | 0.21 |
| information | 0.05 | 0.32 | 0.00 | 0.21 | 0.00 | 0.58 |
| **p(context)** | 0.16 | 0.37 | 0.11 | 0.26 | 0.11 | |

$$pmi_{ij} = \log_2 \frac{p_{ij}}{p_{i*}\, p_{*j}}$$

**p(w,context)**     **p(w)**

| | computer | data | pinch | result | sugar | |
|---|---|---|---|---|---|---|
| apricot | 0.00 | 0.00 | 0.05 | 0.00 | 0.05 | 0.11 |
| pineapple | 0.00 | 0.00 | 0.05 | 0.00 | 0.05 | 0.11 |
| digital | 0.11 | 0.05 | 0.00 | 0.05 | 0.00 | 0.21 |
| information | 0.05 | 0.32 | 0.00 | 0.21 | 0.00 | 0.58 |
| **p(context)** | 0.16 | 0.37 | 0.11 | 0.26 | 0.11 | |

pmi(information,data) = $\log_2$ (    .32 / (.37*.58) ) = .58

*(.57 using full precision)*

## PPMI(w,context)

| | computer | data | pinch | result | sugar |
|---|---|---|---|---|---|
| apricot | - | - | 2.25 | - | 2.25 |
| pineapple | - | - | 2.25 | - | 2.25 |
| digital | 1.66 | 0.00 | - | 0.00 | - |
| information | 0.00 | 0.57 | - | 0.47 | - |

# Weighting PMI

PMI is biased toward infrequent events
◦ Very rare words have very high PMI values

Two solutions:
◦ Give rare words slightly higher probabilities
◦ Use add-one smoothing (which has a similar effect)

# Weighting PMI: Giving rare context words slightly higher probability

Raise the context probabilities to $\alpha = 0.75$:

$$\text{PPMI}_\alpha(w,c) = \max\left(\log_2 \frac{P(w,c)}{P(w)P_\alpha(c)}, 0\right)$$

$$P_\alpha(c) = \frac{count(c)^\alpha}{\sum_c count(c)^\alpha}$$

This helps because $P_\alpha(c) > P(c)$ for rare *c*

Consider two events, P(a) = .99 and P(b)=.01

$$P_\alpha(a) = \frac{.99^{.75}}{.99^{.75} + .01^{.75}} = .97 \quad P_\alpha(b) = \frac{.01^{.75}}{.01^{.75} + .01^{.75}} = .03$$

**Add-2 Smoothed Count(w,context)**

|  | computer | data | pinch | result | sugar |
|---|---|---|---|---|---|
| apricot | 2 | 2 | 3 | 2 | 3 |
| pineapple | 2 | 2 | 3 | 2 | 3 |
| digital | 4 | 3 | 2 | 3 | 2 |
| information | 3 | 8 | 2 | 6 | 2 |

**p(w,context) [add-2]**                    **p(w)**

|  | computer | data | pinch | result | sugar |  |
|---|---|---|---|---|---|---|
| apricot | 0.03 | 0.03 | 0.05 | 0.03 | 0.05 | 0.20 |
| pineapple | 0.03 | 0.03 | 0.05 | 0.03 | 0.05 | 0.20 |
| digital | 0.07 | 0.05 | 0.03 | 0.05 | 0.03 | 0.24 |
| information | 0.05 | 0.14 | 0.03 | 0.10 | 0.03 | 0.36 |
| **p(context)** | 0.19 | 0.25 | 0.17 | 0.22 | 0.17 |  |

# PPMI versus add-2 smoothed PPMI

## PPMI(w,context)

| | computer | data | pinch | result | sugar |
|---|---|---|---|---|---|
| apricot | - | - | 2.25 | - | 2.25 |
| pineapple | - | - | 2.25 | - | 2.25 |
| digital | 1.66 | 0.00 | - | 0.00 | - |
| information | 0.00 | 0.57 | - | 0.47 | - |

## PPMI(w,context) [add-2]

| | computer | data | pinch | result | sugar |
|---|---|---|---|---|---|
| apricot | 0.00 | 0.00 | 0.56 | 0.00 | 0.56 |
| pineapple | 0.00 | 0.00 | 0.56 | 0.00 | 0.56 |
| digital | 0.62 | 0.00 | 0.00 | 0.00 | 0.00 |
| information | 0.00 | 0.58 | 0.00 | 0.37 | 0.00 |

# alternative: prediction-based representation learning

- counting is easy.. but SVD is a costly operation

- starting from 2013/2014 (but with earlier precursors), a different technique became popular

- idea: train a neural network to be good at *predicting* contexts, thus forcing it to be good to represent similar words similarly.

# 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

Brilliant insight: Use running text as implicitly supervised training data!

- A word *s* near *apricot*
  - Acts as gold 'correct answer' to the question
  - "Is word *w* likely to show up near *apricot*?"
- No need for hand-labeled supervision
- The idea comes from **neural language modeling**
  - Bengio et al. (2003)
  - Collobert et al. (2011)

# Word2Vec: Skip-Gram Task

Word2vec provides a variety of options. Let's do
- "skip-gram with negative sampling" (SGNS)

# Skip-gram algorithm

1. Treat the target word and a neighboring context word as positive examples.

2. Randomly sample other words in the lexicon to get negative samples

3. Use logistic regression to train a classifier to distinguish those two cases

4. Use the weights as the embeddings

# Skip-Gram Training Data

Training sentence:

... lemon, a **tablespoon** of **apricot** jam   a   pinch ...

          c1          c2  target  c3    c4

Asssume context words are those in +/- 2 word window

# Skip-Gram Goal

Given a tuple (t,c)  = target, context
  ◦ (*apricot, jam*)
  ◦ (*apricot, aardvark*)
Return probability that c is a real context word:

P(+|t,c)
$P(-|t,c) = 1 - P(+|t,c)$

# How to compute p(+|t,c)?

Intuition:
- Words are likely to appear near similar words
- Model similarity with dot-product!
- Similarity(t,c) $\propto$ t · c

*Problem:*
- *Dot product is not a probability!*
  - *(Neither is cosine)*

# Turning dot product into a probability

The sigmoid lies between 0 and 1:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

# Turning dot product into a probability

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

$$P(-|t,c) = 1 - P(+|t,c)$$

$$= \frac{e^{-t \cdot c}}{1+e^{-t \cdot c}}$$

# For all the context words:

Assume all context words are independent

$$P(+|t, c_{1:k}) = \prod_{i=1}^{k} \frac{1}{1 + e^{-t \cdot c_i}}$$

$$\log P(+|t, c_{1:k}) = \sum_{i=1}^{k} \log \frac{1}{1 + e^{-t \cdot c_i}}$$

# Skip-Gram Training Data

Training sentence:

... lemon, a **tablespoon of apricot** jam   a   pinch ...

           c1         c2    t      c3    c4

Training data: input/output pairs centering on *apricot*

Asssume a +/- 2 word window

# Skip-Gram Training

Training sentence:

... lemon, a **tablespoon of apricot** jam   a   pinch ...

        c1          c2    t      c3   c4

**positive examples +**

| t | c |
| --- | --- |
| apricot | tablespoon |
| apricot | of |
| apricot | preserves |
| apricot | or |

- For each positive example, we'll create *k* negative examples.
- Using *noise* words
- Any random word that isn't *t*

# Skip-Gram Training

Training sentence:

… lemon, a **tablespoon of apricot** jam   a   pinch …

c1          c2    t      c3   c4

| positive examples + | |
| --- | --- |
| t | c |
| apricot | tablespoon |
| apricot | of |
| apricot | preserves |
| apricot | or |

k=2

| negative examples - | | | |
| --- | --- | --- | --- |
| t | c | t | c |
| apricot | aardvark | apricot | twelve |
| apricot | puddle | apricot | hello |
| apricot | where | apricot | dear |
| apricot | coaxial | apricot | forever |

# Choosing noise words

Could pick w according to their unigram frequency P(w)

More common to chosen then according to $p_\alpha(w)$

$$P_\alpha(w) = \frac{count(w)^\alpha}{\sum_w count(w)^\alpha}$$

α= ¾ works well because it gives rare noise words slightly higher probability

To show this, imagine two events p(a)=.99 and p(b) = .01:

$$P_\alpha(a) = \frac{.99^{.75}}{.99^{.75} + .01^{.75}} = .97$$

$$P_\alpha(b) = \frac{.01^{.75}}{.99^{.75} + .01^{.75}} = .03$$

# 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.

# Learning the classifier

Iterative process.

We'll start with 0 or random weights

Then adjust the word weights to
◦ make the positive pairs more likely
◦ and the negative pairs less likely

over the entire training set:

# 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, and the − label for the pairs sample from the negative data.

# Focusing on one target word t:

$$L(\theta) = \log P(+|t,c) + \sum_{i=1} \log P(-|t,n_i)$$

$$= \log \sigma(c \cdot t) + \sum_{i=1}^{k} \log \sigma(-n_i \cdot t)$$

$$= \log \frac{1}{1+e^{-c \cdot t}} + \sum_{i=1}^{k} \log \frac{1}{1+e^{n_i \cdot t}}$$

W

**apricot**

1.2.......j.........v

1
.
.
.
d

C

1... ... d

1
.
k
.
n
.
v

increase
similarity( apricot , jam)
$w_j \cdot c_k$

"…apricot jam…"

**jam** *neighbor word*

**aardvark** *random noise word*

decrease
similarity( apricot , aardvark)
$w_j \cdot c_n$

# Train using gradient descent

*[ we will look at this in more detail later ]*

Actually learns two separate embedding matrices W and C

Can use W and throw away C, or merge them somehow

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

Start with V random 300-dimensional vectors as initial embeddings

Use logistic regression, the second most basic classifier used in machine learning after naïve bayes
- 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.

# Evaluating embeddings

Compare to human scores on word similarity-type tasks:

- WordSim-353 (Finkelstein et al., 2002)

- SimLex-999 (Hill et al., 2015)

- Stanford Contextual Word Similarity (SCWS) dataset (Huang et al., 2012)

- TOEFL dataset: *Levied is closest in meaning to: imposed, believed, requested, correlated*

# Properties of embeddings

Similarity depends on window size C

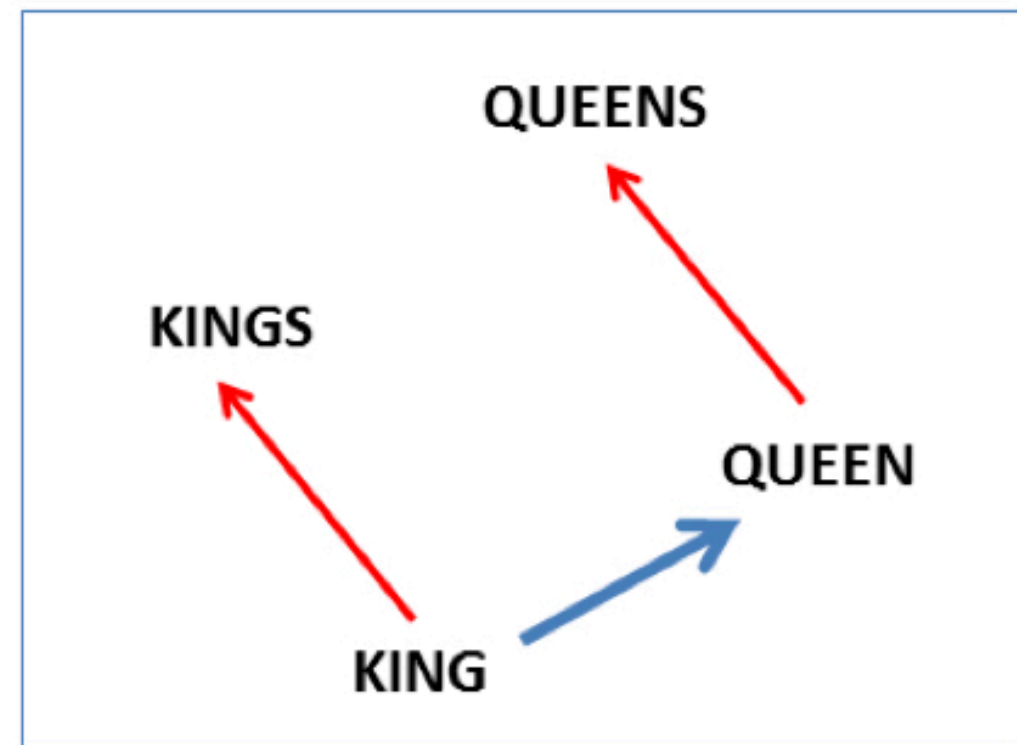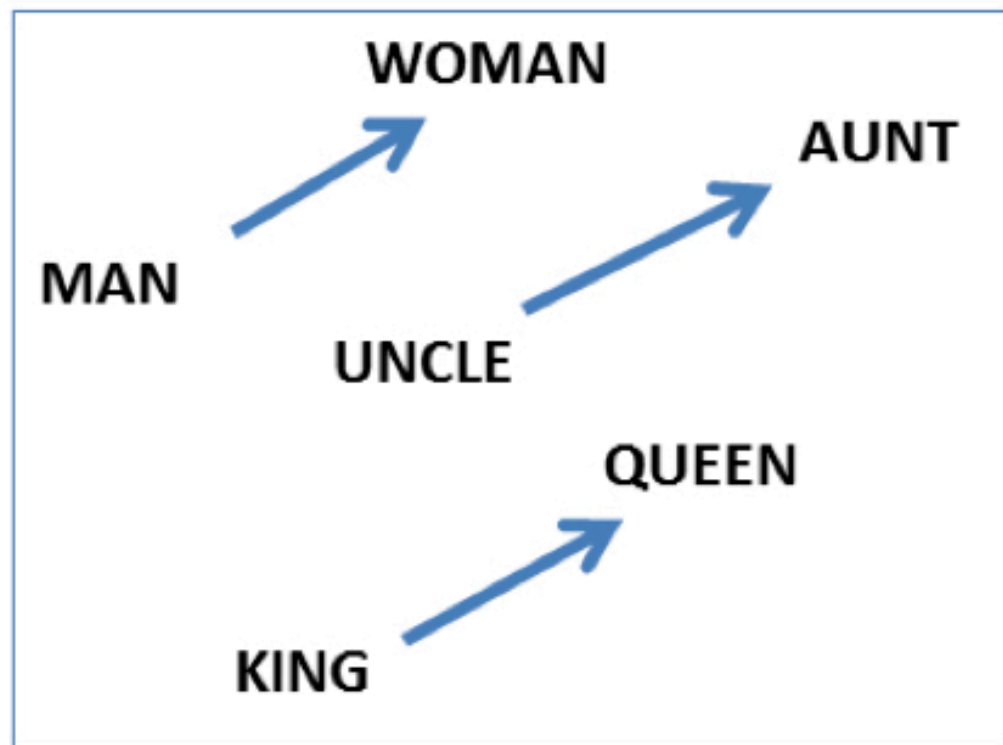C = ±2 The nearest words to *Hogwarts:*
◦ *Sunnydale*
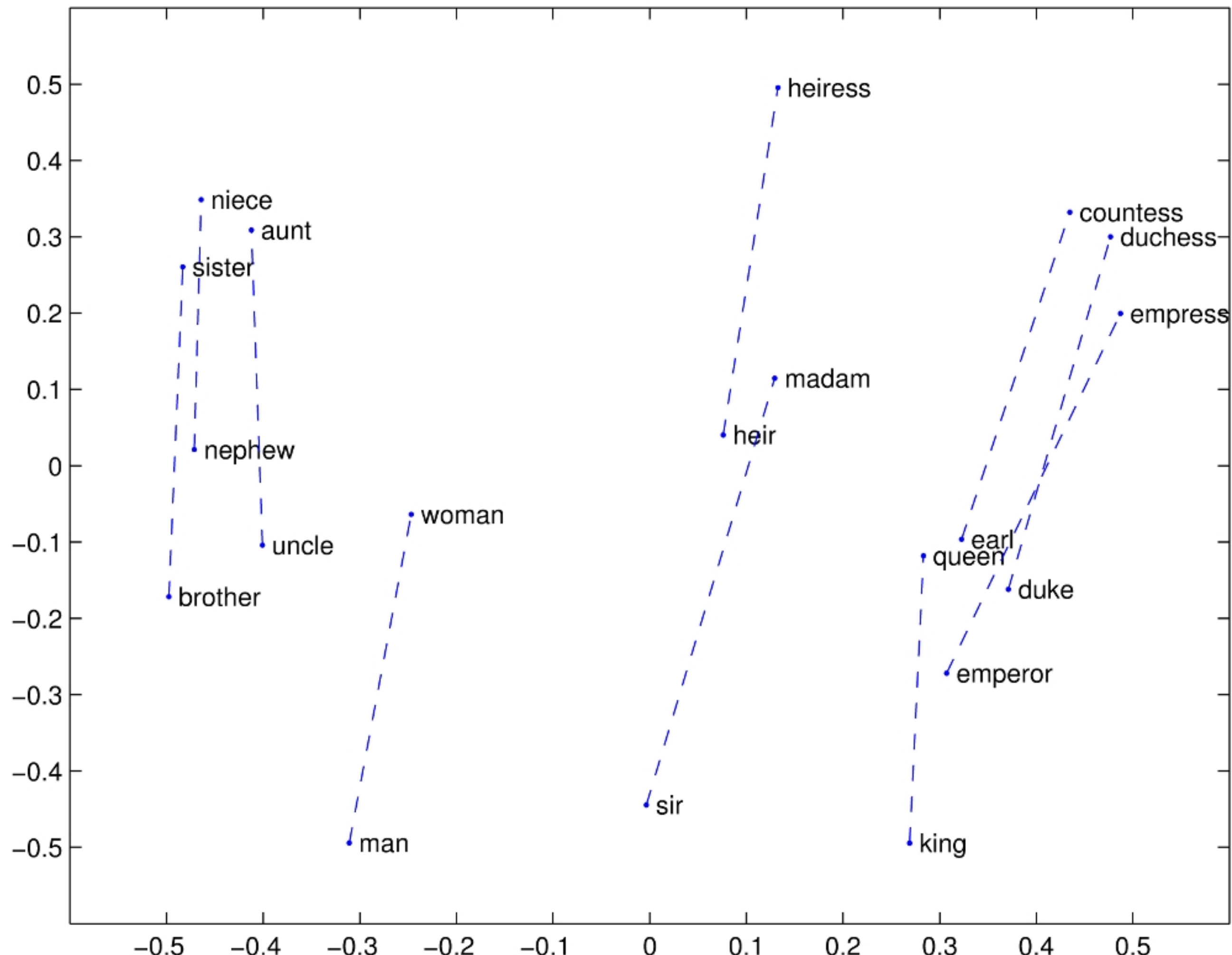◦ *Evernight*

C = ±5 The nearest words to *Hogwarts:*
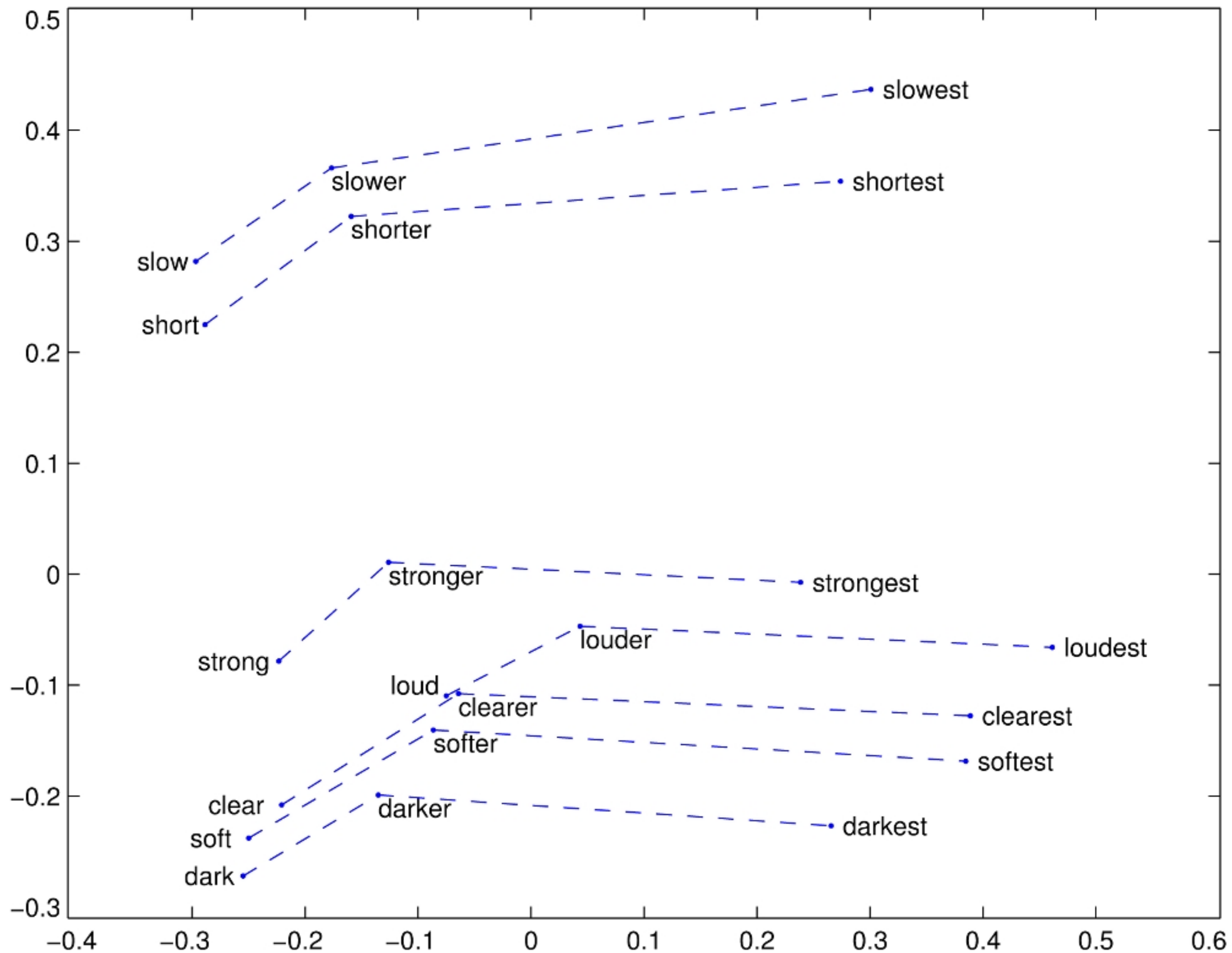◦ *Dumbledore*
◦ *Malfoy*
◦ *halfblood*

# Analogy: Embeddings capture relational meaning!

vector(*'king'*) - vector(*'man'*) + vector(*'woman'*) ≈ vector('queen')

vector(*'Paris'*) - vector(*'France'*) + vector(*'Italy'*) ≈ vector('Rome')

# Embeddings reflect cultural bias

Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings." In *Advances in Neural Information Processing Systems*, pp. 4349-4357. 2016.

Ask "Paris : France :: Tokyo : x"
◦ x = Japan

Ask "father : doctor :: mother : x"
◦ x = nurse

Ask "man : computer programmer :: woman : x"
◦ x = homemaker

# Directions

Debiasing algorithms for embeddings

◦ Bolukbasi, Tolga, Chang, Kai-Wei, Zou, James Y., Saligrama, Venkatesh, and Kalai, Adam T. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in Neural Infor- mation Processing Systems*, pp. 4349–4357.

Use embeddings as a historical tool to study bias

# summary: words

- words map form to meaning

- word meanings can stand in relations

- (aspects of) word meaning can be derived from observing word use in (linguistic) context

- *representation learning:* set up task so that it requires something, let model figure out itself what representation helps best

- we now have representation for both word *identity* and (aspects of) word *meaning* at the same time. (cp. Wittgenstein quote)

- will be very useful for machine learning, helping methods to generalise beyond tokens seen in training data

- missing: linking words to world. (Attempts to integrate this: Lazaridou *et al.* 2014; arguments for representing this separately: Schlangen *et al.* 2016)

# Questions, Queries, Comments?

# slide credits

slides that look like this                    come from

Classical ("Aristotelian") Theory of Concepts

The meaning of a word:

a concept defined by **necessary** and **sufficient** conditions

A **necessary** condition for being an X is a condition C that X must satisfy in order for it to be an X.
  ◦ If not C, then not X
  ◦ "Having four sides" is necessary to be a square.

A **sufficient** condition for being an X is condition such that if something satisfies condition C, then it must be an X.
  ◦ If and only if C, then X
  ◦ The following necessary conditions, jointly, are sufficient to be a square.
    ◦ x has (exactly) four sides
    ◦ each of x's sides is straight
    ◦ x is a closed figure
    ◦ x lies in a plane
    ◦ each of x's sides is equal in length to each of the others
    ◦ each of x's interior angles is equal to the others (right angles)
    ◦ the sides of x are joined at their ends

Example from Norman Swartz, SFU

Dan Jurafsky's slide deck for J&M

and their use is gratefully acknowledged. I try to make any modifications obvious, but if there are errors on a slide, assume that I added them.