# Context in Vector Semantics

## James Pustejovsky

Brandeis University

November 1, 2024

# Desiderata

What should a theory of word meaning do for us?

Let's look at some desiderata

From lexical semantics, the linguistic study of word meaning

# Lemmas and senses

**lemma**

mouse (N)

**sense**
1. any of numerous small rodents...
2. a hand-operated device that controls a cursor...

Modified from the online thesaurus WordNet

A sense or "concept" is the meaning component of a word
Lemmas can be polysemous (have multiple senses)

# Relations between senses: Synonymy

Synonyms have the same meaning in some or all contexts.

- ◦ filbert / hazelnut
- ◦ couch / sofa
- ◦ big / large
- ◦ automobile / car
- ◦ vomit / throw up
- ◦ water / $H_2O$

# Relations between senses: Synonymy

Note that there are probably no examples of perfect synonymy.
◦ Even if many aspects of meaning are identical
◦ Still may differ based on politeness, slang, register, genre, etc.

Relation: **Synonymy**?

water/$H_2O$
    "$H_2O$" in a surfing guide?
big/large
    my big sister != my large sister

# The Linguistic Principle of Contrast

Difference in form → difference in meaning

## Abbé Gabriel Girard 1718

Re: "exact" synonyms

"je ne crois pas qu'il y ait de
mot synonime dans aucune
Langue."

[I do not believe that there
is a synonymous word in any
language]

Thanks to Mark Aronoff!

LA JUSTESSE
DE LA
LANGUE FRANÇOISE,
OU
LES DIFFERENTES SIGNIFICATIONS
DES MOTS QUI PASSENT
POUR
SYNONIMES·

Par M. l'Abbé GIRARD C. D. M. D. D. E.

A PARIS,
Chez LAURENT D'HOURY, Imprimeur-
Libraire, au bas de la rue de la Harpe, vis-
à vis la rue S. Severin, au Saint Esprit.

M. DCC. XVIII.
Avec Approbation & Privilege du Roy.

# Relation: **Similarity**

Words with similar meanings. Not synonyms, but sharing some element of meaning

```
car, bicycle
cow, horse
```

# Ask humans how similar 2 words are

| word1 | word2 | similarity |
| --- | --- | --- |
| vanish | disappear | 9.8 |
| behave | obey | 7.3 |
| belief | impression | 5.95 |
| muscle | bone | 3.65 |
| modest | flexible | 0.98 |
| hole | agreement | 0.3 |

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

# Relation: Word relatedness

Also called "word association"

Words can be related in any way, perhaps via a semantic frame or field

- `coffee, tea`: **similar**
- `coffee, cup`: **related**, not similar

# Semantic field

Words that
◦ cover a particular semantic domain
◦ bear structured relations with each other.

**hospitals**
*surgeon*, *scalpel*, *nurse*, *anaesthetic*, *hospital*
**restaurants**
*waiter*, *menu*, *plate*, *food*, *menu*, *chef*
**houses**
*door*, *roof*, *kitchen*, *family*, *bed*

# Relation: Antonymy

Senses that are opposites with respect to only one feature of meaning

Otherwise, they are very similar!

```
dark/light    short/long fast/slow    rise/fall
hot/cold         up/down         in/out
```

More formally: antonyms can
◦ define a binary opposition or be at opposite ends of a scale
◦ long/short, fast/slow
◦ Be *reversives*:
◦ rise/fall, up/down

# Connotation (sentiment)

- Words have **affective** meanings
  - Positive connotations (*happy*)
  - Negative connotations (*sad*)

- Connotations can be subtle:
  - Positive connotation: *copy, replica, reproduction*
  - Negative connotation: *fake, knockoff, forgery*

- Evaluation (sentiment!)
  - Positive evaluation (*great*, *love*)
  - Negative evaluation (*terrible*, *hate*)

# Connotation

## Words seem to vary along 3 affective dimensions:

◦ **valence**: the pleasantness of the stimulus
◦ **arousal**: the intensity of emotion provoked by the stimulus
◦ **dominance**: the degree of control exerted by the stimulus

|  | Word | Score |  | Word | Score |
|---|---|---|---|---|---|
| **Valence** | love | 1.000 |  | toxic | 0.008 |
|  | happy | 1.000 |  | nightmare | 0.005 |
| **Arousal** | elated | 0.960 |  | mellow | 0.069 |
|  | frenzy | 0.965 |  | napping | 0.046 |
| **Dominance** | powerful | 0.991 |  | weak | 0.045 |
|  | leadership | 0.983 |  | empty | 0.081 |

## So far

**Concepts** or word senses
- Have a complex many-to-many association with **words** (homonymy, multiple senses)

Have relations with each other
- Synonymy
- Antonymy
- Similarity
- Relatedness
- Connotation

# Vector Semantics & Embeddings

Vector Semantics

## Computational models of word meaning

Can we build a theory of how to represent word meaning, that accounts for at least some of the desiderata?

We'll introduce **vector semantics**

The standard model in language processing!

Handles many of our goals!

# Ludwig Wittgenstein

PI #43:
  "The meaning of a word is its use in the language"

# Let's define words by their usages

One way to define "usage":

words are defined by their environments (the words around them)

Zellig Harris (1954):

**If A and B have almost identical environments we say that they are synonyms**.

# What does recent English borrowing *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
  ◦ We could conclude this based on words like "leaves" and "delicious" and "sauteed"

# Ongchoi: *Ipomoea aquatica "Water Spinach"*



空心菜
*kangkong*
rau muống
…

Yamaguchi, Wikimedia Commons, public domain

# The Distributional Hypothesis

- "You shall know a word by the company it keeps." (Firth, 1957)
- "It may be presumed that any two morphemes A and B having different meanings, also differ somewhere in distribution: there are some environments in which one occurs and the other does not." (Harris, 1951)
- "The similarity of the contextual representations of two words contributes to the semantic similarity of those words." (Miller and Charles, 1991) (emphasis mine)

# The Distributional Hypothesis

- "You shall know a word by the company it keeps." (Firth, 1957)
- "It may be presumed that any two morphemes A and B having different meanings, also differ somewhere in distribution: there are some environments in which one occurs and the other does not." (Harris, 1951)
- "The similarity of the contextual representations of two words contributes to the semantic similarity of those words." (Miller and Charles, 1991) (emphasis mine)

- Words can be represented by (abstractions over) their contexts
    - Specifically, linguistic context

Idea 1: Defining meaning by linguistic distribution

Let's define the meaning of a word by its distribution in language use, meaning its neighboring words or grammatical environments.

# Idea 2: Meaning as a point in space (Osgood et al. 1957)

## 3 affective dimensions for a word

◦ **valence**: pleasantness
◦ **arousal**: intensity of emotion
◦ **dominance**: the degree of control exerted

|  | Word | Score |  | Word | Score |
|---|---|---|---|---|---|
| **Valence** | love | 1.000 |  | toxic | 0.008 |
|  | happy | 1.000 |  | nightmare | 0.005 |
| **Arousal** | elated | 0.960 |  | mellow | 0.069 |
|  | frenzy | 0.965 |  | napping | 0.046 |
| **Dominance** | powerful | 0.991 |  | weak | 0.045 |
|  | leadership | 0.983 |  | empty | 0.081 |

NRC VAD Lexicon
(Mohammad 2018)

Hence the connotation of a word is a vector in 3-space

Idea 1: Defining meaning by linguistic distribution

Idea 2: Meaning as a point in multidimensional space

Defining meaning as a point in space based on distribution

Each word = a vector   (not just "good" or "$w_{45}$")

Similar words are "**nearby in semantic space**"

We build this space automatically by seeing which words are **nearby in text**

# We define meaning of a word as a vector

Called an "embedding" because it's embedded into a space (see textbook)

The standard way to represent meaning in NLP

**Every modern NLP algorithm uses embeddings as the representation of word meaning**

Fine-grained model of meaning for similarity

# Intuition: why vectors?

Consider sentiment analysis:

◦ With **words**, a feature is a word identity
  ◦ Feature 5: 'The previous word was "terrible"'
  ◦ requires **exact same word** to be in training and test

◦ With **embeddings**:
  ◦ Feature is a word vector
  ◦ 'The previous word was vector [35,22,17…]
  ◦ Now in the test set we might see a similar vector [34,21,14]
  ◦ We can generalize to **similar but unseen** words!!!

# We'll discuss 2 kinds of embeddings

## tf-idf
◦ Information Retrieval workhorse!
◦ 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 **predict** whether a word is likely to appear nearby
◦ Later we'll discuss extensions called **contextual embeddings**

# Distributed Representations of Words

- ▶ More generally, two approaches to distributed, distributional representations (Baroni et al. 2014):
    - ▶ Count-based
        - ▶ Count occurrences of words in contexts, optionally followed by some mathematical transformation (e.g., tf-idf, PPMI, SVD)
    - ▶ Prediction-based
        - ▶ Given some context vector(s) **c**, predict some word **x** (or vice versa)
        - ▶ a.k.a. language modeling-based

        

        (e.g., word2vec,  ,  )

From now on:
Computing with meaning representations
instead of string representations

荃者所以在鱼，得鱼而忘荃　Nets are for fish;
　　　　　　　　　　　　　　Once you get the fish, you can forget the net.
言者所以在意，得意而忘言　Words are for meaning;
　　　　　　　　　　　　　　Once you get the meaning, you can forget the words
　　　　　　　　　　　　　　　　　　　庄子(Zhuangzi), Chapter 26

# Vector Semantics & Embeddings
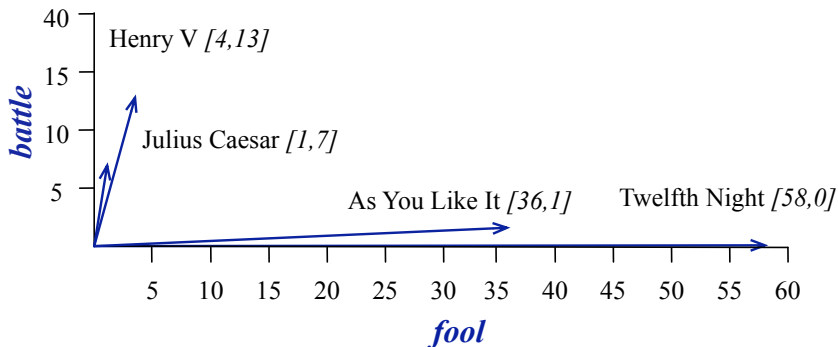
Vector Semantics

# Vector Semantics & Embeddings

Words and Vectors

# 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       |

# 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

But comedies are different than the other two

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

# 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

| | is traditionally followed by | **cherry** | pie, a traditional dessert |
| | often mixed, such as | **strawberry** | rhubarb pie. Apple pie |
| | computer peripherals and personal | **digital** | assistants. These devices usually |
| | a computer. This includes | **information** | available on the internet |

| | aardvark | ... | computer | data | result | pie | sugar | ... |
|---|---|---|---|---|---|---|---|---|
| **cherry** | 0 | ... | 2 | 8 | 9 | 442 | 25 | ... |
| **strawberry** | 0 | ... | 0 | 0 | 1 | 60 | 19 | ... |
| **digital** | 0 | ... | 1670 | 1683 | 85 | 5 | 4 | ... |
| **information** | 0 | ... | 3325 | 3982 | 378 | 5 | 13 | ... |

# Vector Semantics & Embeddings

Words and Vectors

**Vector Semantics & Embeddings**

## Cosine for computing word similarity

Computing word similarity: Dot product and cosine

The dot product between two vectors is a scalar:

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

The dot product tends to be high when the two vectors have large values in the same dimensions

Dot product can thus be a useful similarity metric between vectors

# Problem with raw dot-product

Dot product favors long vectors

Dot product is higher if a vector is longer (has higher values in many dimension)

Vector length:

$$|\mathbf{v}| = \sqrt{\sum_{i=1}^{N} v_i^2}$$

Frequent words (of, the, you) have long vectors (since they occur many times with other words).

So dot product overly favors frequent words

# Alternative: cosine for computing word similarity

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

Based on the definition of the dot product between two vectors a and b

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

## Cosine as a similarity metric

-1: vectors point in opposite directions

+1: vectors point in same directions

0: vectors are orthogonal



But since raw frequency values are non-negative, the cosine for term-term matrix vectors ranges from 0–1

# Cosine examples

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

|  | pie | data | computer |
|---|---|---|---|
| cherry | 442 | 8 | 2 |
| digital | 5 | 1683 | 1670 |
| information | 5 | 3982 | 3325 |

$\cos(\text{cherry}, \text{information}) =$

$$\frac{442 * 5 + 8 * 3982 + 2 * 3325}{\sqrt{442^2 + 8^2 + 2^2}\sqrt{5^2 + 3982^2 + 3325^2}} = .017$$

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

$$\frac{5 * 5 + 1683 * 3982 + 1670 * 3325}{\sqrt{5^2 + 1683^2 + 1670^2}\sqrt{5^2 + 3982^2 + 3325^2}} = .996$$

## Visualizing cosines
## (well, angles)

**Vector Semantics & Embeddings**

Cosine for computing word similarity

# Vector Semantics & Embeddings

TF-IDF

# But raw frequency is a bad representation

- The co-occurrence matrices we have seen represent each cell by word frequencies.
- 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
- It's a paradox! How can we balance these two conflicting constraints?

# Two common solutions for word weighting

**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 "it" have very low idf

**PMI:** (Pointwise mutual information)

- $PMI(w_1, w_2) = log \dfrac{p(w_1, w_2)}{p(w_1)p(w_2)}$

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

# Term frequency (tf)

$$tf_{t,d} = count(t,d)$$

Instead of using raw count, we squash a bit:

$$tf_{t,d} = \log_{10}(count(t,d)+1)$$

# Document frequency (df)

df$_t$ *is* the number of documents $t$ occurs in.

(note this is not collection frequency: total count across all documents)

"*Romeo*" is very distinctive for one Shakespeare play:

|  | Collection Frequency | Document Frequency |
|---|---|---|
| Romeo | 113 | 1 |
| action | 113 | 31 |

# Inverse document frequency (idf)

$$\mathrm{idf}_t \;=\; \log_{10}\left(\frac{N}{\mathrm{df}_t}\right)$$

N is the total number of documents in the collection

| Word | df | idf |
|------|-----|-------|
| Romeo | 1 | 1.57 |
| salad | 2 | 1.27 |
| Falstaff | 4 | 0.967 |
| forest | 12 | 0.489 |
| battle | 21 | 0.246 |
| wit | 34 | 0.037 |
| fool | 36 | 0.012 |
| good | 37 | 0 |
| sweet | 37 | 0 |

# What is a document?

Could be a play or a Wikipedia article

But for the purposes of tf-idf, documents can be **anything**; we often call each paragraph a document!

# Final tf-idf weighted value for a word

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

Raw counts:

|        | 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       |

tf-idf:

|        | As You Like It | Twelfth Night | Julius Caesar | Henry V |
|--------|----------------|---------------|---------------|---------|
| **battle** | 0.074      | 0             | 0.22          | 0.28    |
| **good**   | 0          | 0             | 0             | 0       |
| **fool**   | 0.019      | 0.021         | 0.0036        | 0.0083  |
| **wit**    | 0.049      | 0.044         | 0.018         | 0.022   |

# Vector Semantics & Embeddings

TF-IDF

# Vector Semantics & Embeddings

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 $c_j$

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

|  | computer | data | result | pie | sugar | count(w) |
|---|---|---|---|---|---|---|
| **cherry** | 2 | 8 | 9 | 442 | 25 | 486 |
| **strawberry** | 0 | 0 | 1 | 60 | 19 | 80 |
| **digital** | 1670 | 1683 | 85 | 5 | 4 | 3447 |
| **information** | 3325 | 3982 | 378 | 5 | 13 | 7703 |
| **count(context)** | 4997 | 5673 | 473 | 512 | 61 | 11716 |

$$pmi_{ij} = \log_2 \frac{p_{ij}}{p_{i*} p_{*j}} \qquad ppmi_{ij} = \begin{cases} pmi_{ij} & \text{if } pmi_{ij} > 0 \\ 0 & \text{otherwise} \end{cases}$$

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

|  | computer | data | result | pie | sugar | count(w) |
|---|---|---|---|---|---|---|
| **cherry** | 2 | 8 | 9 | 442 | 25 | 486 |
| **strawberry** | 0 | 0 | 1 | 60 | 19 | 80 |
| **digital** | 1670 | 1683 | 85 | 5 | 4 | 3447 |
| **information** | 3325 | 3982 | 378 | 5 | 13 | 7703 |
| **count(context)** | 4997 | 5673 | 473 | 512 | 61 | 11716 |

p(w=information,c=data) = 3982/11716 = .3399

p(w=information) = 7703/11716 = .6575

p(c=data) = 5673/11716 = .4842

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

|  | p(w,context) | | | | | p(w) |
|---|---|---|---|---|---|---|
|  | computer | data | result | pie | sugar | p(w) |
| **cherry** | 0.0002 | 0.0007 | 0.0008 | 0.0377 | 0.0021 | 0.0415 |
| **strawberry** | 0.0000 | 0.0000 | 0.0001 | 0.0051 | 0.0016 | 0.0068 |
| **digital** | 0.1425 | 0.1436 | 0.0073 | 0.0004 | 0.0003 | 0.2942 |
| **information** | 0.2838 | 0.3399 | 0.0323 | 0.0004 | 0.0011 | 0.6575 |
| **p(context)** | 0.4265 | 0.4842 | 0.0404 | 0.0437 | 0.0052 |  |

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

| | **p(w,context)** | | | | | **p(w)** |
|---|---|---|---|---|---|---|
| | **computer** | **data** | **result** | **pie** | **sugar** | **p(w)** |
| **cherry** | 0.0002 | 0.0007 | 0.0008 | 0.0377 | 0.0021 | 0.0415 |
| **strawberry** | 0.0000 | 0.0000 | 0.0001 | 0.0051 | 0.0016 | 0.0068 |
| **digital** | 0.1425 | 0.1436 | 0.0073 | 0.0004 | 0.0003 | 0.2942 |
| **information** | 0.2838 | 0.3399 | 0.0323 | 0.0004 | 0.0011 | 0.6575 |
| **p(context)** | 0.4265 | 0.4842 | 0.0404 | 0.0437 | 0.0052 | |

pmi(information,data) = $\log_2$ (.3399 / (.6575*.4842) ) = .0944

Resulting PPMI matrix (negatives replaced by 0)

| | **computer** | **data** | **result** | **pie** | **sugar** |
|---|---|---|---|---|---|
| **cherry** | 0 | 0 | 0 | 4.38 | 3.30 |
| **strawberry** | 0 | 0 | 0 | 4.10 | 5.51 |
| **digital** | 0.18 | 0.01 | 0 | 0 | 0 |
| **information** | 0.02 | 0.09 | 0.28 | 0 | 0 |

# 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 $\alpha = 0.75$:

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

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

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

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

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

# Distributed Representations of Words

- More generally, two approaches to distributed, distributional representations (Baroni et al. 2014):
    - Count-based
        - Count occurrences of words in contexts, optionally followed by some mathematical transformation (e.g., tf-idf, PPMI, SVD)
    - Prediction-based
        - Given some context vector(s) $\mathbf{c}$, predict some word $\mathbf{x}$ (or vice versa)
        - a.k.a. language modeling-based

            

            (e.g., word2vec,  ,  )

# Language Models

- Given some context vector(s) $\mathbf{c}$, predict some word $\mathbf{x}$ (or vice versa)
- Two approaches to language models:
    - Generative models
        - Model the joint probability distribution $P(\mathbf{x}, \mathbf{c})$
        - Examples: n-gram language models
            - Unigram: predict $P(\mathbf{x}_i)$
            - Bigram: predict $P(\mathbf{x}_i | \mathbf{x}_{i-1})$
            - Trigram: predict $P(\mathbf{x}_i | \mathbf{x}_{i-2}, \mathbf{x}_{i-1})$

# Language Models

▶ Given some context vector(s) **c**, predict some word **x** (or vice versa)

▶ Two approaches to language models:
  ▶ Discriminative models
    ▶ Predict the conditional probability $P(\mathbf{x}|\mathbf{c})$ (or $P(\mathbf{c}|\mathbf{x})$) directly
    ▶ Examples: neural network language models
      ▶ Feedforward: word2vec (Mikolov et al., 2013a, 2013b)

      ▶ Recurrent:  (Peters et al., 2018)

      ▶ Transformer:  (Devlin et al., 2019)

# word2vec

- ▶ Based on a feedforward neural network language model



CBOW

Skip-gram

# Neural Networks



- Output layer
- Hidden layer(s)
- Input layer

# Neural Networks



- Output layer
- Hidden layer(s)
- Input layer
- $\mathbf{x}$ is the input
- $\mathbf{h}$ is the hidden layer output
  - Can be seen as intermediate representation of the input
- $\hat{\mathbf{y}}$ is the predicted output
  - $\hat{} =$ predicted

# Neural Networks



- ▶ Output layer
- ▶ Hidden layer(s)
- ▶ Input layer
- ▶ $\mathbf{h} = g(\mathbf{x} \cdot \mathbf{W})$
- ▶ $\hat{\mathbf{y}} = f(\mathbf{h} \cdot \mathbf{C})$
    - ▶ $\mathbf{W}$ and $\mathbf{C}$ are weight (or parameter) matrices
        - ▶ May or may not include a bias term
    - ▶ $g$ and $f$ are activation functions

# word2vec

- Based on a feedforward neural network language model



CBOW

Skip-gram

- Continuous bag of words (CBOW): use context to predict current word
- Skip-gram: use current word to predict context

# CBOW



- Input layer: one-hot word vectors
    - $\begin{bmatrix} 0 & \cdots & 0 & 1 & 0 & \cdots & 0 \end{bmatrix}$
    - Context words within some window

# CBOW



- ▶ Hidden (projection) layer: identity activation function, no bias
  - ▶ Weight matrix shared for all context words
  - ▶ Input → hidden = table lookup (in weight matrix)
  - ▶ Context word vectors are averaged

# CBOW



- ▶ Output layer: softmax activation function
  - ▶ Numbers → probabilities

# Skip-gram



- Input layer: one-hot word vectors
  - $\begin{bmatrix} 0 & \cdots & 0 & 1 & 0 & \cdots & 0 \end{bmatrix}$

# Skip-gram



- ▶ Hidden (projection) layer: identity activation function, no bias
    - ▶ Input → hidden = table lookup (in weight matrix)

# Skip-gram



- ▶ Output layer: softmax activation function
    - ▶ Predict context words within some window
    - ▶ Separate classification for each context word
    - ▶ Closer context words sampled more than distant context words

# word2vec

- Skip-gram model: for each word, word2vec learns two word embeddings
  - Target word vector $\mathbf{w}$ (row of $\mathbf{W}$, = output of hidden layer)
  - Context word vector $\mathbf{c}$ (column of $\mathbf{C}$)
- Common final word embeddings
  - Add $\mathbf{w} + \mathbf{c}$
  - Just $\mathbf{w}$ (throw away $\mathbf{c}$)

Vector Semantics & Embeddings

Properties of Embeddings

## The kinds of neighbors depend on window size

**Small windows** (C= +/- 2) : nearest words are syntactically similar words in same taxonomy
   ◦*Hogwarts* nearest neighbors are other fictional schools
   ◦*Sunnydale, Evernight, Blandings*

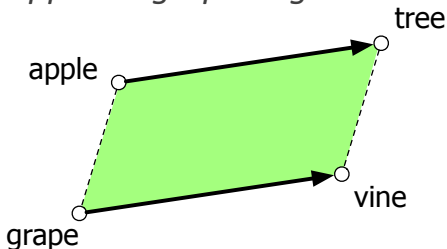**Large windows** (C= +/- 5) :  nearest words are related words in same semantic field
   ◦*Hogwarts* nearest neighbors are Harry Potter world:
   ◦*Dumbledore, half-blood,  Malfoy*

## Analogical relations

The classic parallelogram model of analogical reasoning
(Rumelhart and Abrahamson 1973)

To solve: *"apple is to tree as grape is to _____"*

*Add $\overrightarrow{tree} - \overrightarrow{apple}$ to $\overrightarrow{grape}$ to get $\overrightarrow{vine}$*

# Analogical relations via parallelogram

The parallelogram method can solve analogies with both sparse
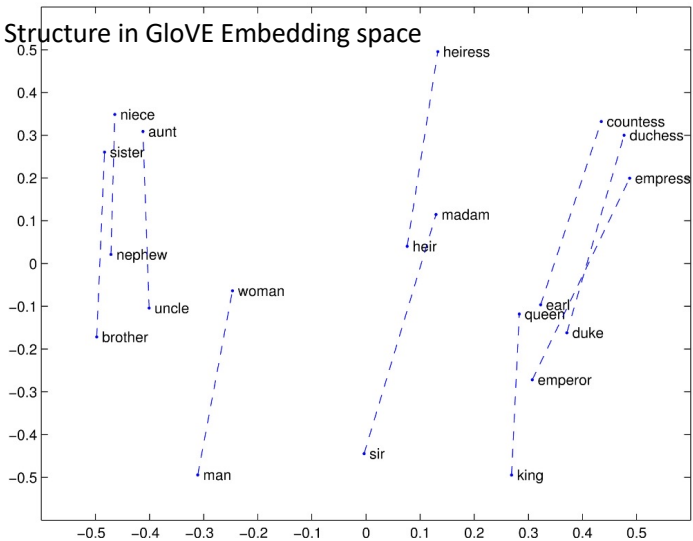and dense embeddings (Turney and Littman 2005, Mikolov et al.
2013b)

$$\overrightarrow{king} - \overrightarrow{man} + \overrightarrow{woman} \text{ is close to } \overrightarrow{queen}$$

$$\overrightarrow{Paris} - \overrightarrow{France} + \overrightarrow{Italy} \text{ is close to } \overrightarrow{Rome}$$

For a problem $a : a^* :: b : b^*$, the parallelogram method is:

$$\hat{b}^* = \underset{x}{\operatorname{argmin}} \operatorname{distance}(x, a^* - a + b)$$

Structure in GloVE Embedding space

# Caveats with the parallelogram method

It only seems to work for frequent words, small distances and certain relations (relating countries to capitals, or parts of speech), but not others. (Linzen 2016, Gladkova et al. 2016, Ethayarajh et al. 2019a)
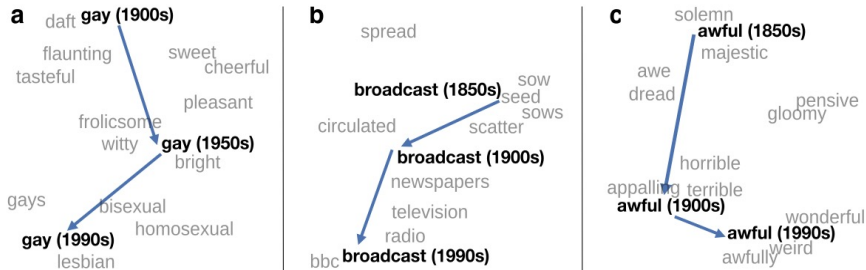
Understanding analogy is an open area of research (Peterson et al. 2020)

# Embeddings as a window onto historical semantics

**Train embeddings on different decades of historical text to see meanings shift**

~30 million books, 1850-1990, Google Books data



William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. Proceedings of ACL.

# 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 *NeurIPS*, 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

Algorithms that use embeddings as part of e.g., hiring searches for programmers, might lead to bias in hiring

# Historical embedding as a tool to study cultural biases

Garg, N., Schiebinger, L., Jurafsky, D., and Zou, J. (2018). Word embeddings quantify 100 years of gender and ethnic stereotypes. Proceedings of the National Academy of Sciences 115(16), E3635–E3644.

- Compute a **gender or ethnic bias** for each adjective: e.g., how much closer the adjective is to "woman" synonyms than "man" synonyms, or names of particular ethnicities
  - Embeddings for **competence** adjective (*smart, wise, brilliant, resourceful, thoughtful, logical)* are biased toward men, a bias slowly decreasing 1960-1990
  - Embeddings for **dehumanizing** adjectives (barbaric, monstrous, bizarre) were biased toward Asians in the 1930s, bias decreasing over the 20th century.
- These match the results of old surveys done in the 1930s