

Natural Language Processing: Word Semantics

18 February 2022

Naive Bayes: Another Example Clustering

$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(w | c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$$

Priors:

$$P(c) = \frac{3}{4}$$

$$P(j) = \frac{1}{4}$$

Conditional Probabilities:

$$P(\text{Chinese}|c) = (5+1) / (8+6) = 6/14 = 3/7$$

$$P(\text{Tokyo}|c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Japan}|c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Chinese}|j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Tokyo}|j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Japan}|j) = (1+1) / (3+6) = 2/9$$

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

Choosing a class:

$$P(c|d5) \propto 3/4 * (3/7)^3 * 1/14 * 1/14$$

$$\approx 0.0003$$

$$P(j|d5) \propto 1/4 * (2/9)^3 * 2/9 * 2/9$$

$$\approx 0.0001$$

Word Meaning

Previous lectures

- In anything we have seen so far, words were just strings (or indices in a vocabulary list)
- Not very satisfactory..

What do words mean?

First thought: look in a dictionary, e.g.,

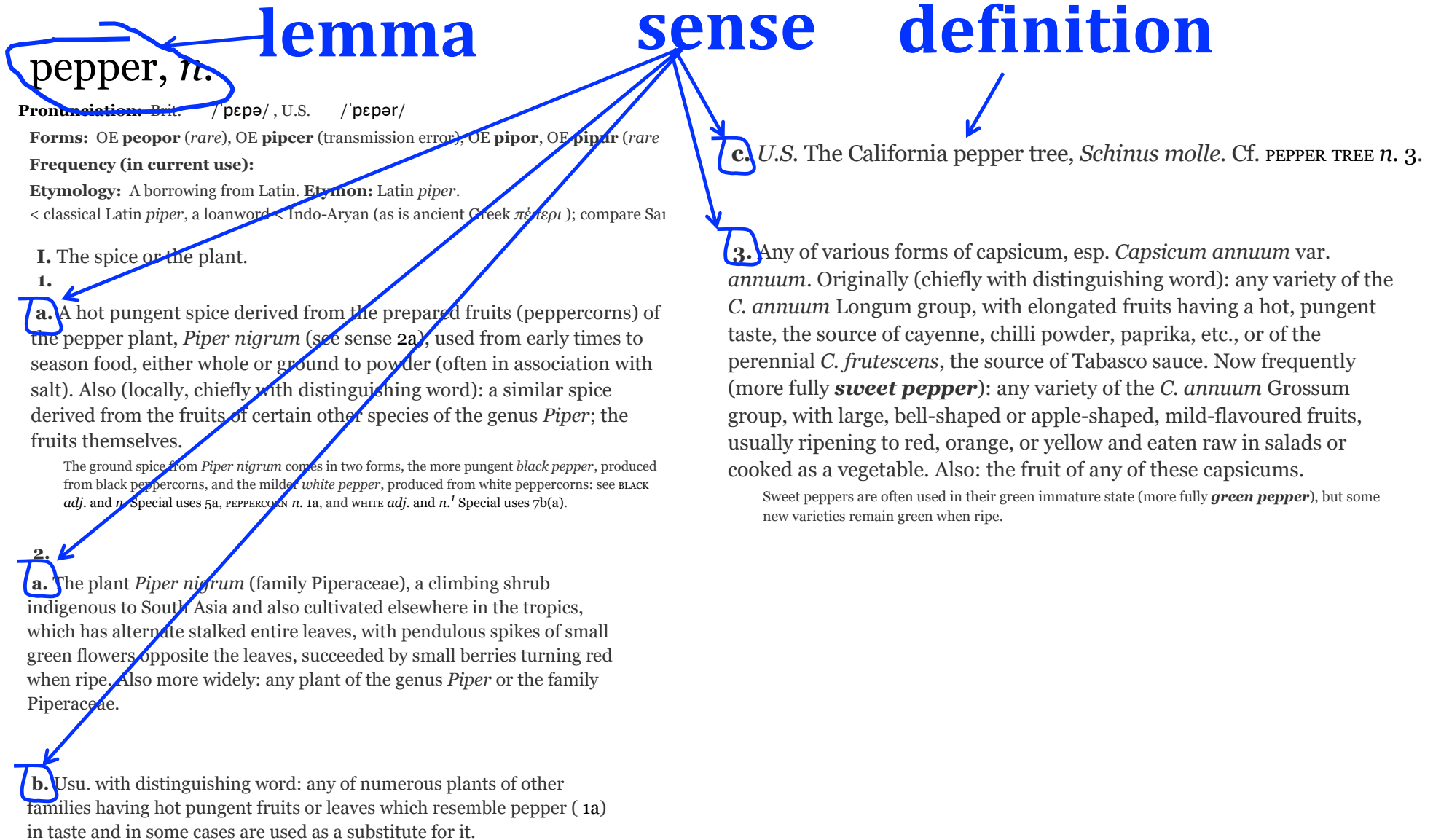
<http://www.oed.com/>

Reminder: lemma and wordform

- A **lemma** or **citation form**
 - Same stem, part of speech, rough semantics
- A **wordform**
 - The “inflected” word as it appears in text

Wordform	Lemma
banks	bank
sung	sing
duermes	dormir

Words, Lemmas, Senses, Definitions



Lemma pepper

Sense 1: spice from pepper plant

Sense 2: the pepper plant itself

Sense 3: another similar plant (Jamaican pepper)

Sense 4: another plant with peppercorns (California pepper)

Sense 5: *capsicum* (i.e. chili, paprika, bell pepper, etc.)

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

Homonymy

Homonyms: words that **share a form but have unrelated, distinct meanings:**

- **bank**₁: financial institution, **bank**₂: sloping land
- **bat**₁: club for hitting a ball, **bat**₂: nocturnal flying mammal

1. Homographs (bank/bank, bat/bat)
2. Homophones:
 1. **Write** and **right**
 2. **Piece** and **peace**

Homonymy causes problems for NLP applications

- Information retrieval
 - “bat care”
- Machine Translation
 - bat: **murciélagos** (animal) or **bate** (for baseball)
- Text-to-Speech
 - bass (stringed instrument) vs. bass (fish)

Polysemy

1. The **bank** was constructed in 1875 out of local red brick.
 2. I withdrew the money from the **bank**
- Are those the same sense?
 - Sense 2: “A financial institution”
 - Sense 1: “The building belonging to a financial institution”
 - A **polysemous** word has **related** meanings
 - Most non-rare words have multiple meanings

Metonymy or Systematic Polysemy:

A systematic relationship between senses

- Many types of polysemy are systematic
 - School, university, hospital, White House, etc.
 - All can mean the institution or the building
- A systematic relationship:
 - Building ↔ Organization
- Other such kinds of systematic polysemy:

Author (Jane Austen wrote Emma)

↔ Works of Author (I love Jane Austen)

Tree (Plums have beautiful blossoms)

↔ Fruit (I ate a preserved plum)

How do we know when a word has more than one sense?

- The “zeugma” test: Two senses of `serve`?
 - Which flights **serve** breakfast?
 - Does Lufthansa **serve** Philadelphia?
 - ?Does Lufthansa serve breakfast and Philadelphia?
- Since this conjunction sounds weird,
 - we say that these are **two different senses of “serve”**

Relations between senses: Synonymy

- Synonyms have the same meaning in some or all contexts
 - couch / sofa
 - big / large
 - automobile / car
 - vomit / throw up
 - water / H₂O
 - filbert / hazelnut

Two lemmas are synonyms

- if they can be substituted for each other in all situations.
- If so they have the same **propositional meaning**

Synonymy

- Note that there are probably no examples of perfect synonymy
 - Even if many aspects of meaning are identical
 - Still may not preserve the acceptability based on notions of politeness, slang, context, genre, etc.

Synonymy?

water/H₂O

big/large

brave/courageous

- "H₂O" in a surfing guide?
- "How big is that plane?" vs. "Would I be flying on a large or small plane?"
- "my big sister" or "my large sister"? (*big has a sense that means being older, or grown up, while large lacks this sense*)

Relation: Similarity

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

car, bicycle

cow, horse

Why word similarity?

- Information retrieval
- Question answering
- Machine translation
- Natural language generation
- Language modeling
- Automatic essay grading
- Plagiarism detection
- Document clustering
- ...

Asking humans how similar two 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)

SimLex-999 vs. WordSim-353

- SimLex-999 provides a way of measuring how well models capture *similarity*, rather than *relatedness* or *association*
- The scores in SimLex-999 therefore differ from other well-known evaluation datasets such as *WordSim-353* [1]

Pair	Simlex-999 rating	WordSim-353 rating
<i>coast - shore</i>	9.00	9.10
<i>clothes - closet</i>	1.96	8.00

Relation: Word relatedness/association

- Also called "word association"
- Words can be related in any way, perhaps via a semantic frame or field
 - car, bicycle: **similar**
 - car, gasoline: **related** (not similar)

Semantic field

- Words that
 - cover a particular semantic domain and
 - 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

Semantic frame

- A semantic frame is a set of words that denote perspectives or participants in a particular type of event
- E.g., a commercial transaction - a kind of event in which one entity trades money to another entity in return for some good or service, after which the good changes hands (or perhaps the service is performed)
 - Lexical encoding:
 - *buy* (the event from the perspective of the buyer), *sell* (from the perspective of the seller), *pay* (focusing on the monetary aspect), or nouns like *buyer*, *seller*
 - Frames have semantic roles (like *buyer*, *seller*, *goods*, *money*), and words in a sentence can take on these roles
 - Knowing the roles of *buy* and *sell* in their semantic frame one can paraphrase:
 - Sam bought the book from Ling → Ling sold the book to Sam

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

Embedding models like word2vec tend to place antonyms near each other in the vector space, thus we may use also thesauri to distinguish antonyms when using such models

Hyponymy and Hypernymy

- One sense is a **hyponym/subordinate** of another if the first sense is more specific, denoting a subclass of the other
 - *car* is a hyponym of *vehicle*
 - *mango* is a hyponym of *fruit*
- Conversely **hypernym/superordinate** (“hyper is super”)
 - *vehicle* is a **hypernym** of *car*
 - *fruit* is a hypernym of *mango*

Superordinate/hyper	vehicle	fruit	furniture
Subordinate/hyponym	car	mango	chair

Hyponymy more formally

- Extensional:
 - The class denoted by the superordinate extensionally includes the class denoted by the hyponym
- Entailment:
 - A sense A is a hyponym of sense B if *being an A* entails *being a B*
- Hyponymy is usually transitive
 - (A hypo B and B hypo C entails A hypo C)
- Another name: the **IS-A hierarchy**
 - A **IS-A** B
 - B **subsumes** A

Hyponyms and Instances

- WordNet has both **classes** and **instances**
- An **instance** is an individual, a proper noun that is a unique entity
 - San Francisco is an **instance** of city
- But city is a class
 - city is a **hyponym** of municipality, location...

The levels are not symmetric..

- One level of category is distinguished from the others
- The "basic level"

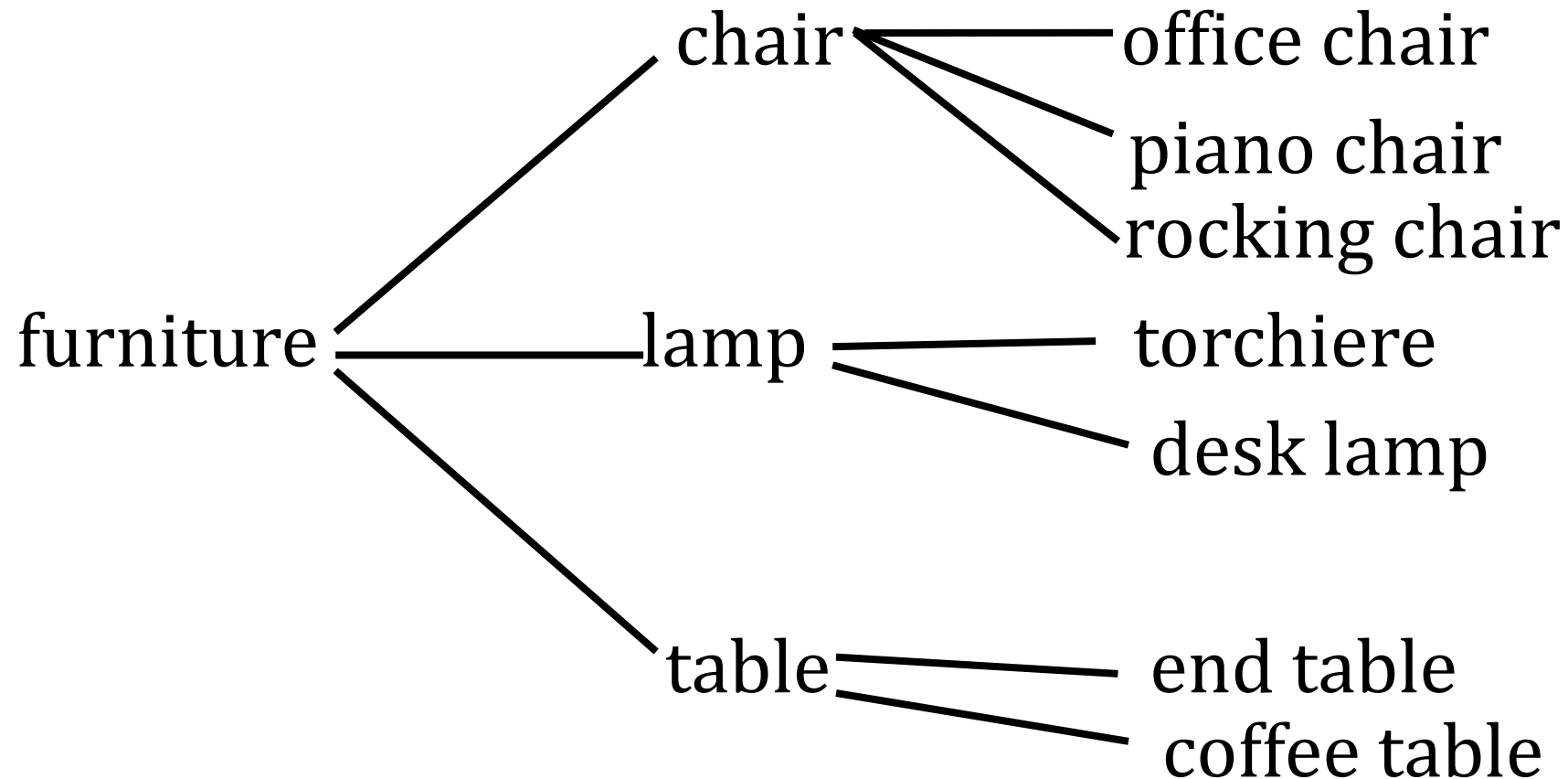
Name these items



Superordinate

Basic

Subordinate



Cluster of Interactional Properties

- Basic level things are “human-sized”
- Consider chairs
 - We know how to interact with a chair (sit)
 - Not so clear for superordinate categories like furniture
 - “Imagine a furniture without thinking of a bed/table/chair/specific basic-level category”

The basic level

- Distinctive actions
- Learned earliest in childhood
- Names are shortest
- Names are most frequent

Connotation (sentiment)

- Words have **affective** meanings
 - positive connotations (*happy*)
 - negative connotations (*sad*)
- fake vs. replica
innocent vs. naïve
- positive evaluation (*great, love*)
 - negative evaluation (*terrible, hate*)

So far

- **Concepts** or word senses
 - Have a complex many-to-many associations with **words** (homonymy, multiple senses)
- Have relations with each other
 - Synonymy
 - Antonymy
 - Similarity
 - Relatedness
 - Superordinate/subordinate, basic level
 - Connotations

WordNet

WordNet 3.0

- A hierarchically organized lexical database
- On-line thesaurus + aspects of a dictionary
 - Some [other languages](#) available or under development
 - (Arabic, Finnish, German, Portuguese...)

Category	Unique Strings
Nouns	117,798
Verbs	11,529
Adjectives	22,479
Adverbs	4,481

Senses of “bass” in WordNet

Noun

- [S:](#) [\(n\)](#) **bass** (the lowest part of the musical range)
- [S:](#) [\(n\)](#) **bass**, [bass part](#) (the lowest part in polyphonic music)
- [S:](#) [\(n\)](#) **bass**, [basso](#) (an adult male singer with the lowest voice)
- [S:](#) [\(n\)](#) [sea bass](#), **bass** (the lean flesh of a saltwater fish of the family Serranidae)
- [S:](#) [\(n\)](#) [freshwater bass](#), **bass** (any of various North American freshwater fish with lean flesh (especially of the genus Micropterus))
- [S:](#) [\(n\)](#) **bass**, [bass voice](#), [basso](#) (the lowest adult male singing voice)
- [S:](#) [\(n\)](#) **bass** (the member with the lowest range of a family of musical instruments)
- [S:](#) [\(n\)](#) **bass** (nontechnical name for any of numerous edible marine and freshwater spiny-finned fishes)

Adjective

- [S:](#) [\(adj\)](#) **bass**, [deep](#) (having or denoting a low vocal or instrumental range) *"a deep voice"; "a bass voice is lower than a baritone voice"; "a bass clarinet"*

How is sense defined in WordNet?

- The **synset (synonym set)**, the set of near-synonyms, instantiates a sense or concept, with a **gloss**
- Example: **chump** as a noun with the **gloss**:
“a person who is gullible and easy to take advantage of”
- This sense of “chump” is shared by 9 words:
chump¹, fool², gull¹, mark⁹, patsy¹, fall guy¹, sucker¹, soft touch¹, mug²
- Each of **these** senses have this same gloss
 - (Not **every** sense of a word; e.g., sense 2 of gull is the aquatic bird)

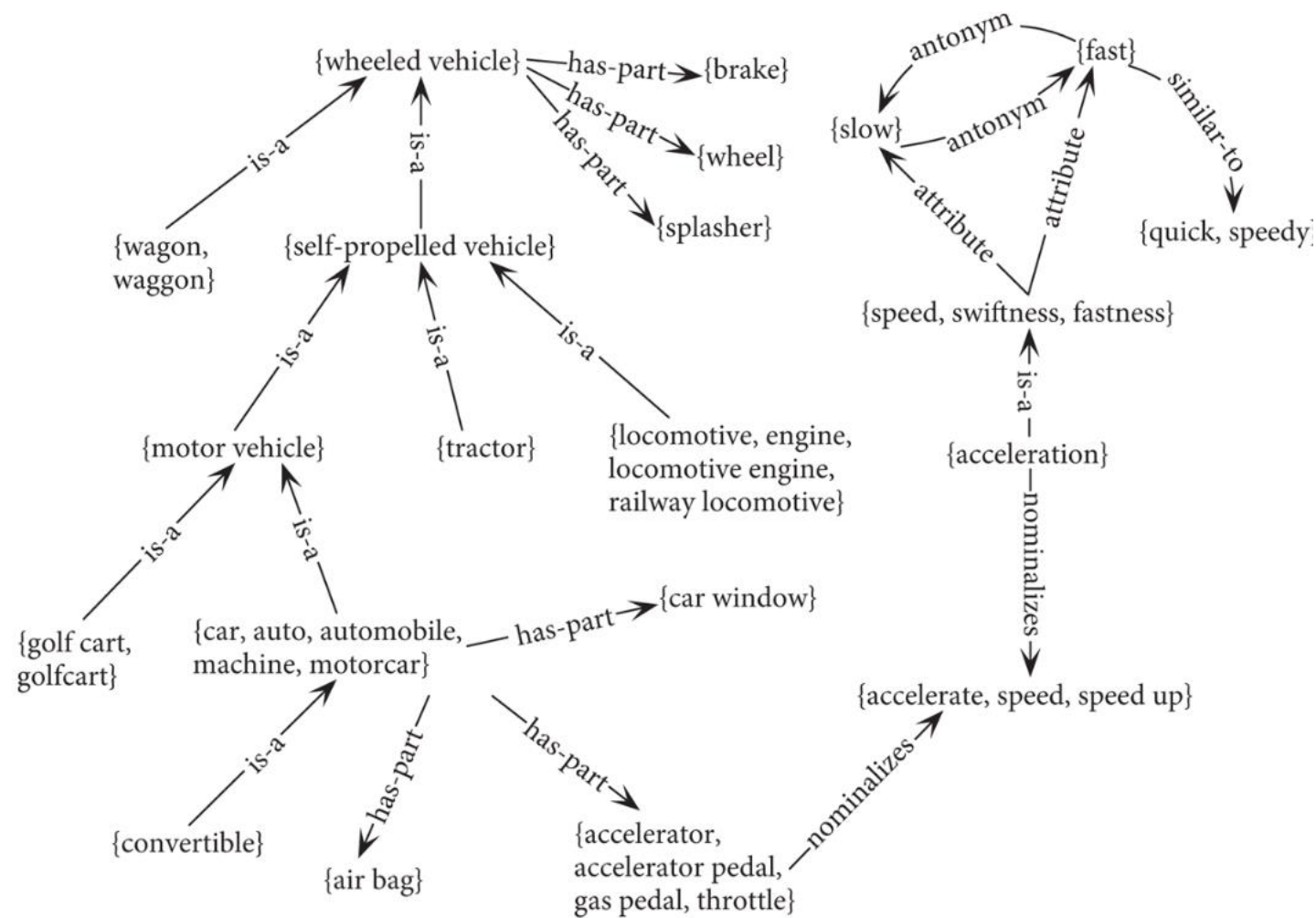
On average noun has 1.23 senses and verb has 2.16 senses

WordNet Hypernym Hierarchy for “bass”

- [S: \(n\) bass](#), [basso](#) (an adult male singer with the lowest voice)
 - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - [S: \(n\) singer](#), [vocalist](#), [vocalizer](#), [vocaliser](#) (a person who sings)
 - [S: \(n\) musician](#), [instrumentalist](#), [player](#) (someone who plays a musical instrument (as a profession))
 - [S: \(n\) performer](#), [performing artist](#) (an entertainer who performs a dramatic or musical work for an audience)
 - [S: \(n\) entertainer](#) (a person who tries to please or amuse)
 - [S: \(n\) person](#), [individual](#), [someone](#), [somebody](#), [mortal](#), [soul](#) (a human being) *"there was too much for one person to do"*
 - [S: \(n\) organism](#), [being](#) (a living thing that has (or can develop) the ability to act or function independently)
 - [S: \(n\) living thing](#), [animate thing](#) (a living (or once living) entity)
 - [S: \(n\) whole](#), [unit](#) (an assemblage of parts that is regarded as a single entity) *"how big is that part compared to the whole?"*; *"the team is a unit"*
 - [S: \(n\) object](#), [physical object](#) (a tangible and visible entity; an entity that can cast a shadow) *"it was full of rackets, balls and other objects"*
 - [S: \(n\) physical entity](#) (an entity that has physical existence)
 - [S: \(n\) entity](#) (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

WordNet Noun Relations

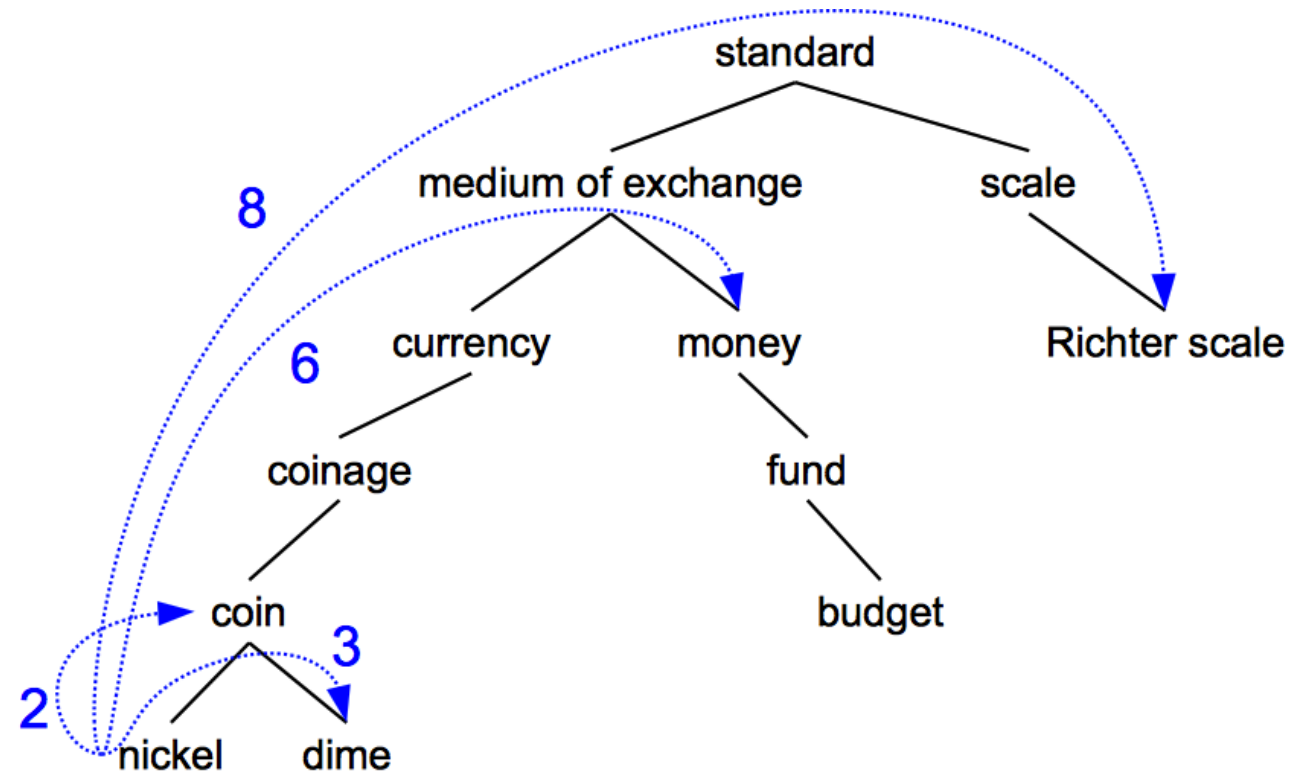
| Relation | Also called | Definition | Example |
|----------------|---------------|---|---|
| Hypernym | Superordinate | From concepts to superordinates | <i>breakfast</i> ¹ → <i>meal</i> ¹ |
| Hyponym | Subordinate | From concepts to subtypes | <i>meal</i> ¹ → <i>lunch</i> ¹ |
| Member Meronym | Has-Member | From groups to their members | <i>faculty</i> ² → <i>professor</i> ¹ |
| Has-Instance | Member-Of | From concepts to instances of the concept | <i>composer</i> ¹ → <i>Bach</i> ¹ |
| Instance | | From instances to their concepts | <i>Austen</i> ¹ → <i>author</i> ¹ |
| Member Holonym | | From members to their groups | <i>copilot</i> ¹ → <i>crew</i> ¹ |
| Part Meronym | | From wholes to parts | <i>table</i> ² → <i>leg</i> ³ |
| Part Holonym | | From parts to wholes | <i>course</i> ⁷ → <i>meal</i> ¹ |
| Antonym | Part-Of | Opposites | <i>leader</i> ¹ → <i>follower</i> ¹ |



(FYI) How to compute similarity of two words?

- Several methods possible (ones listed later tend to be better):
 1. cosine similarity of definitions of two words (e.g., extended Lesk algorithm)
 2. distance of words in the ontology tree of senses (hypernymy tree in WordNet)
 3. cosine similarity measure between two vectors containing tfidf values in term-document matrix
 4. cosine similarity measure between two vectors containing raw co-occurrence counts of words (each word represented as a vector over vocabulary V)
 5. cosine similarity measure between two vectors of PMI associations computed based on word co-occurrence counts
 6. cosine similarity measure between two dense vectors such as ones obtained using deep learning approaches (e.g., word2vec vectors)
 7. ...

Path based similarity



- Two concepts (senses/synsets) are similar if they are near each other in the thesaurus hierarchy
 - =have a short path between them
 - concepts have path 1 to themselves

Path based similarity

$\text{pathlen}(c_1, c_2)$ = 1 + number of edges in the shortest path in the hypernym graph between sense nodes c_1 and c_2

- ranges from 0 to 1 (identity)

$$\text{simpath}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$$

Example: path-based similarity

$$\text{simpath}(c_1, c_2) = 1/\text{pathlen}(c_1, c_2)$$

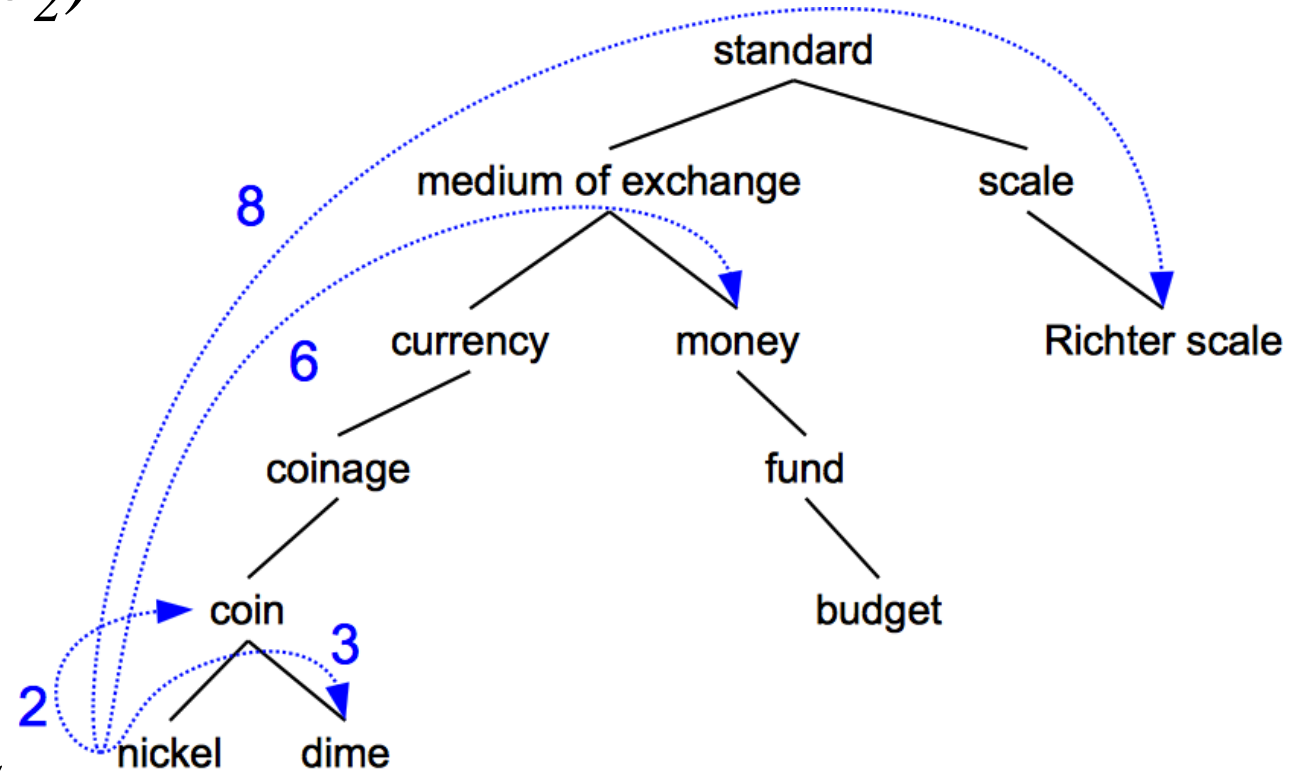
$$\text{simpath}(\textit{nickel}, \textit{coin}) = 1/2 = .5$$

$$\text{simpath}(\textit{fund}, \textit{budget}) = 1/2 = .5$$

$$\text{simpath}(\textit{nickel}, \textit{currency}) = 1/4 = .25$$

$$\text{simpath}(\textit{nickel}, \textit{money}) = 1/6 = .17$$

$$\text{simpath}(\textit{coinage}, \textit{Richter scale}) = 1/6 = .17$$



Problem with basic path based similarity

- Assumes each link represents a uniform distance
 - But *nickel* to *money* seems to us to be closer than *nickel* to *standard*
 - Nodes high in the hierarchy are very abstract
- We instead want a metric that
 - Represents the cost of each edge independently
 - Words connected only through abstract nodes are less similar

The (extended) Lesk Algorithm

- A thesaurus-based measure that looks at **glosses**
- Two concepts are similar if their glosses contain similar words
 - *Drawing paper*: **paper** that is **specially prepared** for use in drafting
 - *Decal*: the art of transferring designs from **specially prepared paper** to a wood or glass or metal surface
- For each n -word phrase that's in both glosses
 - Add a score of n^2
 - **Paper** and **specially prepared** for $1 + 2^2 = 5$
 - Compute the overlap also for other relations
 - glosses of hypernyms and hyponyms

$$\text{sim}_{eLesk}(c_1, c_2) = \sum_{r, q \in RELS} \text{overlap}(\text{gloss}(r(c_1)), \text{gloss}(q(c_2)))$$

WordNet 3.0

- Where it is:
 - <http://wordnetweb.princeton.edu/perl/webwn>
- Libraries
 - R: wordnet package
 - Python: WordNet from NLTK
 - <http://www.nltk.org/Home>
 - Java:
 - JWNL, extJWNL on sourceforge

MeSH: Medical Subject Headings thesaurus from the National Library of Medicine

- **MeSH (Medical Subject Headings)**

- 177,000 entry terms that correspond to 26,142 biomedical “headings”

- **Hemoglobins**

Entry Terms: Eryhem, Ferrous Hemoglobin, Hemoglobin

Definition: The oxygen-carrying proteins of ERYTHROCYTES. They are found in all vertebrates and some invertebrates. The number of globin subunits in the hemoglobin quaternary structure differs between species. Structures range from monomeric to a variety of multimeric arrangements

The MeSH Hierarchy

1. + Anatomy [A]
2. + Organisms [B]
3. + Diseases [C]
4. + Chemicals and Drugs [D]
5. + Analytical, Diagnostic and Therapeutic Techniques and Equipment [E]
6. + Psychiatry and Psychology [F]
7. + Phenomena and Processes [G]
8. + Disciplines and Occupations [H]
9. + Anthropology, Education, Sociology, and Social Sciences [I]
10. + Technology, Industry, Agriculture, and Miscellaneous [J]
11. + Humanities [K]
12. + Information Science [L]
13. + Named Groups [M]
14. + Health Care [N]
15. + Publication Characteristics [O]
16. + Geographicals [Z]

1. + Anatomy [A]
2. + Organisms [B]
3. + Diseases [C]
4. - Chemicals and Drugs [D]
 - [Inorganic Chemicals \[D01\]](#) +
 - [Organic Chemicals \[D02\]](#) +
 - [Heterocyclic Compounds \[D03\]](#) +
 - [Polycyclic Compounds \[D04\]](#) +
 - [Macromolecular Substances \[D05\]](#) +
 - [Hormones, Hormone Substitutes, and Hormone Antagonists \[D06\]](#) +
 - [Enzymes and Coenzymes \[D08\]](#) +
 - [Carbohydrates \[D09\]](#) +
 - [Lipids \[D10\]](#) +
 - [Amino Acids, Peptides, and Proteins \[D12\]](#) +
 - [Nucleic Acids, Nucleotides, and Nucleosides \[D13\]](#) +
 - [Complex Mixtures \[D20\]](#) +
 - [Biological Factors \[D23\]](#) +
 - [Biomedical and Dental Materials \[D25\]](#) +
 - [Pharmaceutical Preparations \[D26\]](#) +
 - [Chemical Actions and Uses \[D27\]](#) +
5. + Analytical, Diagnostic and Therapeutic Techniques and Equipment [E]
6. + Psychiatry and Psychology [F]
7. + Phenomena and Processes [G]

[Amino Acids, Peptides, and Proteins \[D12\]](#)

[Proteins \[D12.776\]](#)

[Blood Proteins \[D12.776.124\]](#)

[Acute-Phase Proteins \[D12.776.124.050\]](#) +

[Anion Exchange Protein 1, Erythrocyte \[D12.776.124.078\]](#)

[Ankyrins \[D12.776.124.080\]](#)

[beta 2-Glycoprotein I \[D12.776.124.117\]](#)

[Blood Coagulation Factors \[D12.776.124.125\]](#) +

[Cholesterol Ester Transfer Proteins \[D12.776.124.197\]](#)

[Fibrin \[D12.776.124.270\]](#) +

[Glycophorin \[D12.776.124.300\]](#)

[Hemocyanin \[D12.776.124.337\]](#)

► [Hemoglobins \[D12.776.124.400\]](#)

[Carboxyhemoglobin \[D12.776.124.400.141\]](#)

[Erythrocyte \[D12.776.124.400.220\]](#)

Uses of the MeSH Ontology

- Provide synonyms (“entry terms”)
 - E.g., glucose and dextrose
- Provide hypernyms (from the hierarchy)
 - E.g., glucose ISA monosaccharide
- Indexing in MEDLINE/PubMED database
 - NLM’s bibliographic database:
 - 20 million journal articles
 - Each article hand-assigned 10-20 MeSH terms

Vector Semantics

How about a radically different approach?

Ludwig Wittgenstein

"The meaning of a word is its use in the language"

Let's define words by their usage

- 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.

“You shall know a word by a company it keeps” (Firth 1957: 11)

What does recent English borrowing *ong choy* mean?

- Suppose you see these sentences:
 - Ong choy is delicious **sautéed with garlic**.
 - Ong choy is superb **over rice**
 - Ong choy **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:
 - Ong choy is a leafy green like spinach, chard, or collard greens

Ong choy: *Ipomoea aquatica* "Water Spinach"

空心菜

kangkong

rau muống

...



Yamaguchi, Wikimedia Commons, public domain

A new model of meaning focusing on distributional similarity

- Each word = a vector
- Similar words are "nearby in space"



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
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 then generalize to **similar but unseen words**

We'll discuss two 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
 - There are extensions such as **contextual embeddings**

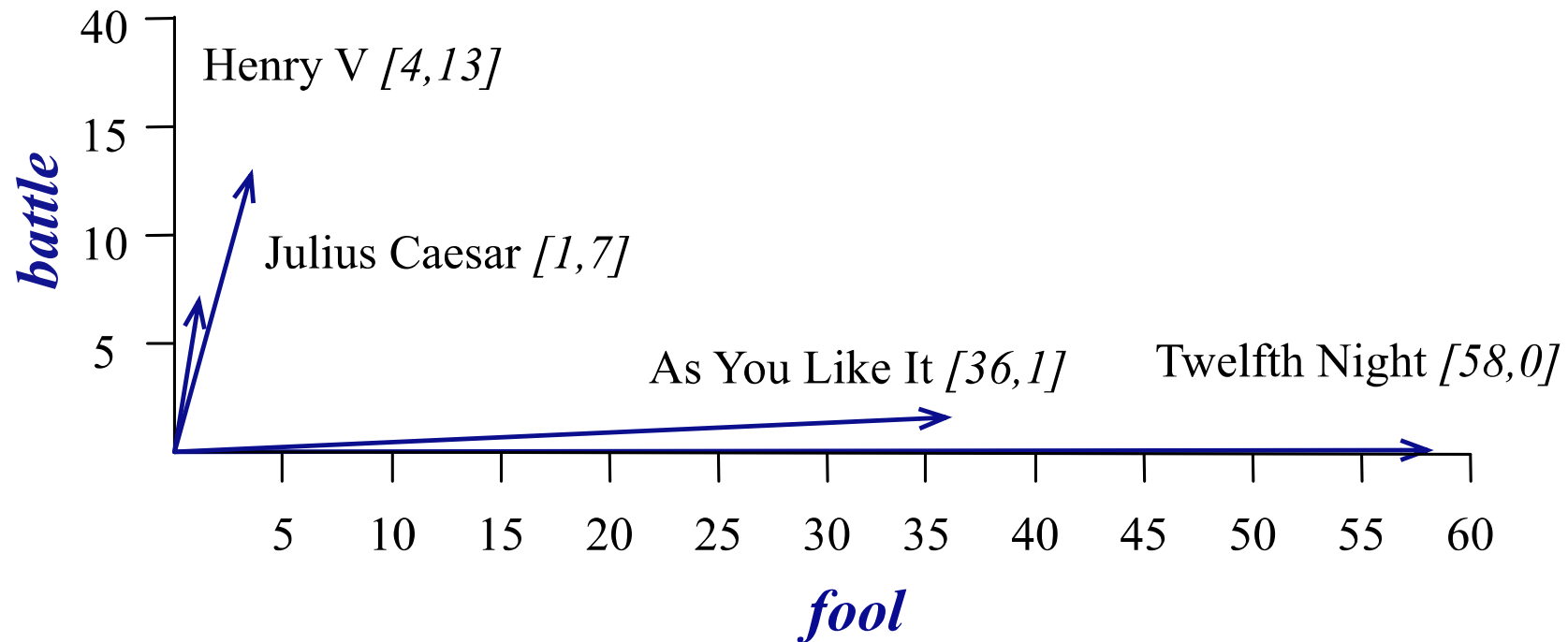
Words and Vectors

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
- Different than the history
- 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"

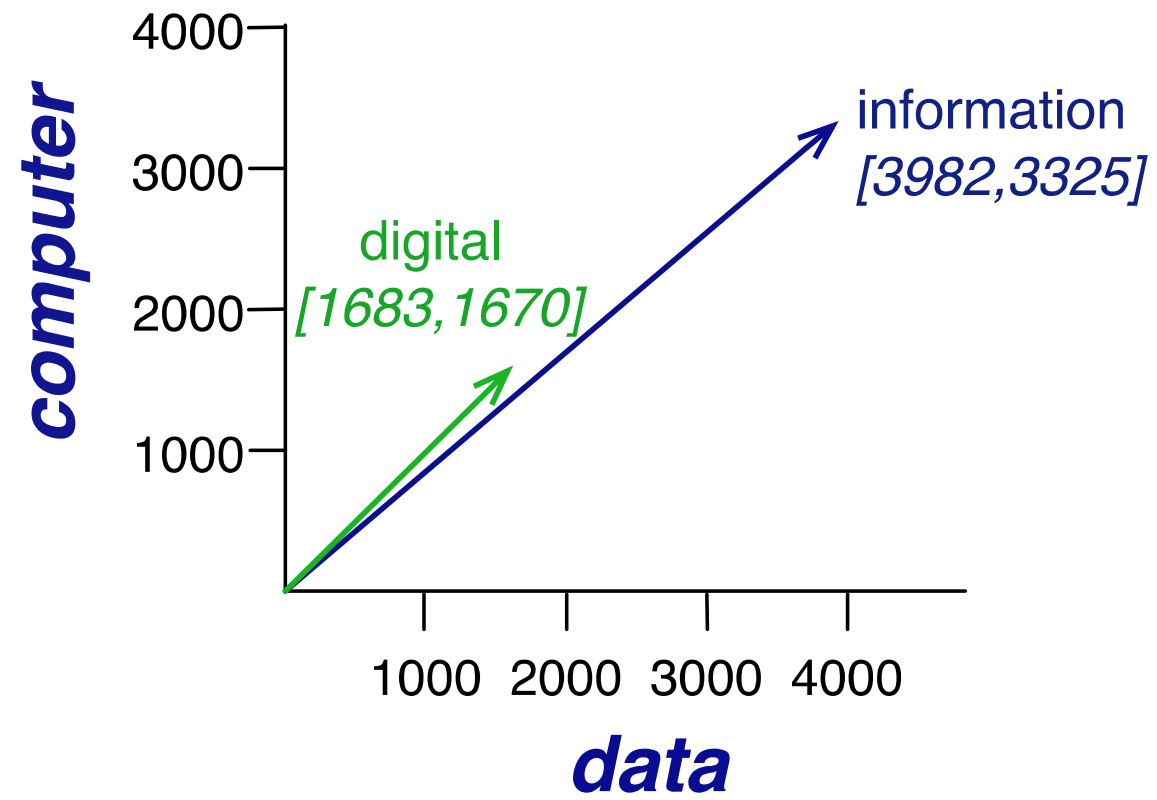
fool is "the kind of word that occurs in comedies, especially Twelfth Night"

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 | ... |



Cosine for computing word similarity

Dot product and cosine

- The dot product between two vectors is a scalar:

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

- The dot product tends to be high when the two vectors have large values in the same dimensions
- Dot product can be a 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 dimensions)
- Vector length:

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

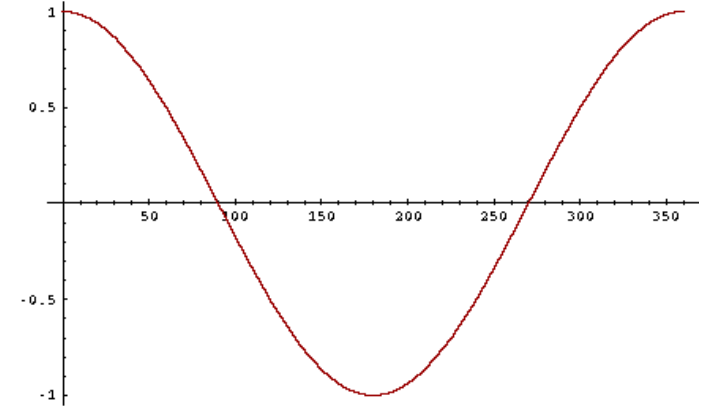
- Frequent words (of, the, you) have long vectors (since they occur many times with other words)
- So dot product overly favors frequent words

Alternative: cosine for computing word similarity

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

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 to 1



Cosine examples

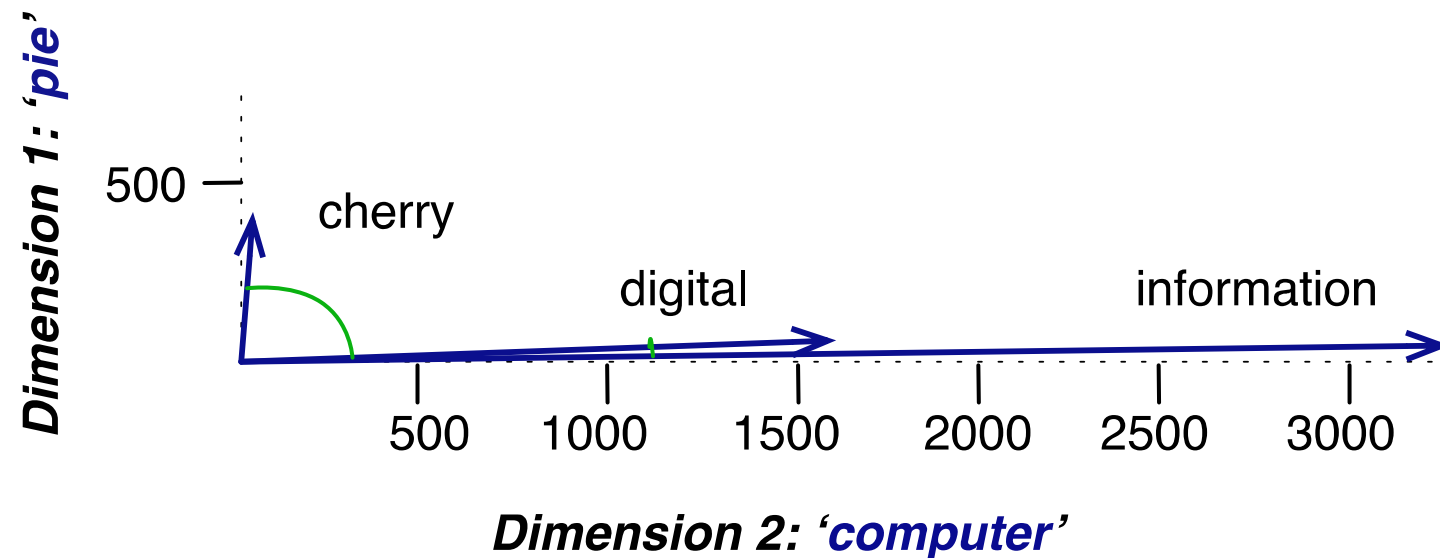
$$\cos(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\vec{v}}{|\vec{v}|} \cdot \frac{\vec{w}}{|\vec{w}|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

| | 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)



TF-IDF

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
- We need a function that resolves this frequency paradox!

Term frequency (tf)

- $\text{tf}_{t,d} = \text{count}(t,d)$

Instead of using the raw count, we may squash a bit:

- $\text{tf}_{t,d} = \log_{10}(\text{count}(t,d)+1)$

Document frequency (df)

- df_t is the number of documents that t occurs in
(note: this is not a [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 |

- Important: documents can be **anything**; we can call each paragraph a document

Inverse document frequency (idf)

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

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 |

Final tf-idf weighted value for a word

- 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:

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

| | 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 |

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)

Word Representations: High vs. Low Dimension Vectors (Dense vs. Sparse)

Sparse Vector Space representation, FCM (term co-occurrence, etc.)

Size of Vocabulary $\approx 10,000 - 400,000$

$$\begin{aligned} \vec{\text{ipod}} &= \begin{bmatrix} 5 & 150 & \dots & 200 & \dots & 5 & \dots & 100 & \dots & 0 & \dots & 0 \end{bmatrix} \\ \vec{\text{car}} &= \begin{bmatrix} 60 & 18 & \dots & 50 & \dots & 200 & \dots & 350 & \dots & 0 & \dots & 0 \end{bmatrix} \end{aligned}$$

(Note: The vectors above are simplified representations of the sparse matrix structure shown in the image, where rows represent words and columns represent terms in the vocabulary.)

Dense Vector Space representation, e.g., LSA, Word2vec, GLoVe

Size of Dimensions 50 - 1000

$$\begin{aligned} \vec{\text{ipod}} &= \begin{bmatrix} 0.1 & -0.2 & \dots & 0.3 & \dots & 0.8 \end{bmatrix} \\ \vec{\text{car}} &= \begin{bmatrix} 0.1 & 0.1 & \dots & -0.8 & \dots & 1.1 \end{bmatrix} \end{aligned}$$

(Note: The vectors above are simplified representations of the dense matrix structure shown in the image, where rows represent words and columns represent dimensions in the vector space.)

- Captures semantic meaning better
- Semantically similar words are located close in vector spaces

Sparse versus dense vectors

- Why dense vectors?
 - Short vectors may be easier to use as **features** in machine learning (fewer weights to tune)
 - Dense vectors may **generalize** better than 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**

Common methods for getting short dense vectors

- “Neural Language Model”-inspired models
 - Word2vec (skipgram, CBOW), Glove
- Singular Value Decomposition (SVD)
 - A special case of this is called LSA – Latent Semantic Analysis
- Alternative to these "static embeddings":
 - Contextual Embeddings (ELMo, BERT)
 - Compute distinct embeddings for a word in its context
 - Separate embeddings for each token of a word

Word2vec: The classifier

Trained Embeddings that can be downloaded

- 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**
 - Big idea: **self-supervision**:
 - A word c that occurs near *apricot* in the corpus acts as the gold "correct answer" for supervised learning
 - No need for human labels
- Bengio et al. (2003); Collobert et al. (2011)

Word2Vec: Skip-Gram Task

Word2vec provides a variety of options. We'll cover:

skip-gram with negative sampling (SGNS)

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

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

Skip-Gram Training Data

- Assume a ± 2 word window, given training sentence:

...lemon, a [tablespoon of apricot jam, a] pinch...

c1 c2 [target] c3 c4

Skip-Gram Training Data

Assume a ± 2 word window, given training sentence:

...lemon, a [tablespoon of apricot jam, a] pinch...

c1 c2 [target] c3 c4

Goal: train a classifier that is given a candidate (**w**ord, **c**ontext) pair

(apricot, tablespoon)

(apricot, aardvark)

...

And assigns each pair a probability: $P(+|w, c)$

Similarity is computed from dot product

Two vectors are similar if they have a high dot product

- Cosine is just a normalized dot product
- So:
 - $\text{Similarity}(w,c) \propto w \cdot c$
- We'll need to normalize to get a probability
 - (cosine isn't a probability either)

Turning dot products into probabilities

- $\text{Sim}(w, c) \approx w \cdot c$
- To turn this into a probability
- We'll use the sigmoid from logistic regression:

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

$$\begin{aligned} P(-|w, c) &= 1 - P(+|w, c) \\ &= \sigma(-c \cdot w) = \frac{1}{1 + \exp(c \cdot w)} \end{aligned}$$

How Skip-Gram Classifier computes $P(+|w, c)$

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

This is for one context word, but we have lots of context words.
We assume independence and just multiply them:

$$P(+|w, c_{1:L}) = \prod_{i=1}^L \sigma(c_i \cdot w)$$

$$\log P(+|w, c_{1:L}) = \sum_{i=1}^L \log \sigma(c_i \cdot w)$$

Skip-gram classifier: summary

A probabilistic classifier that,

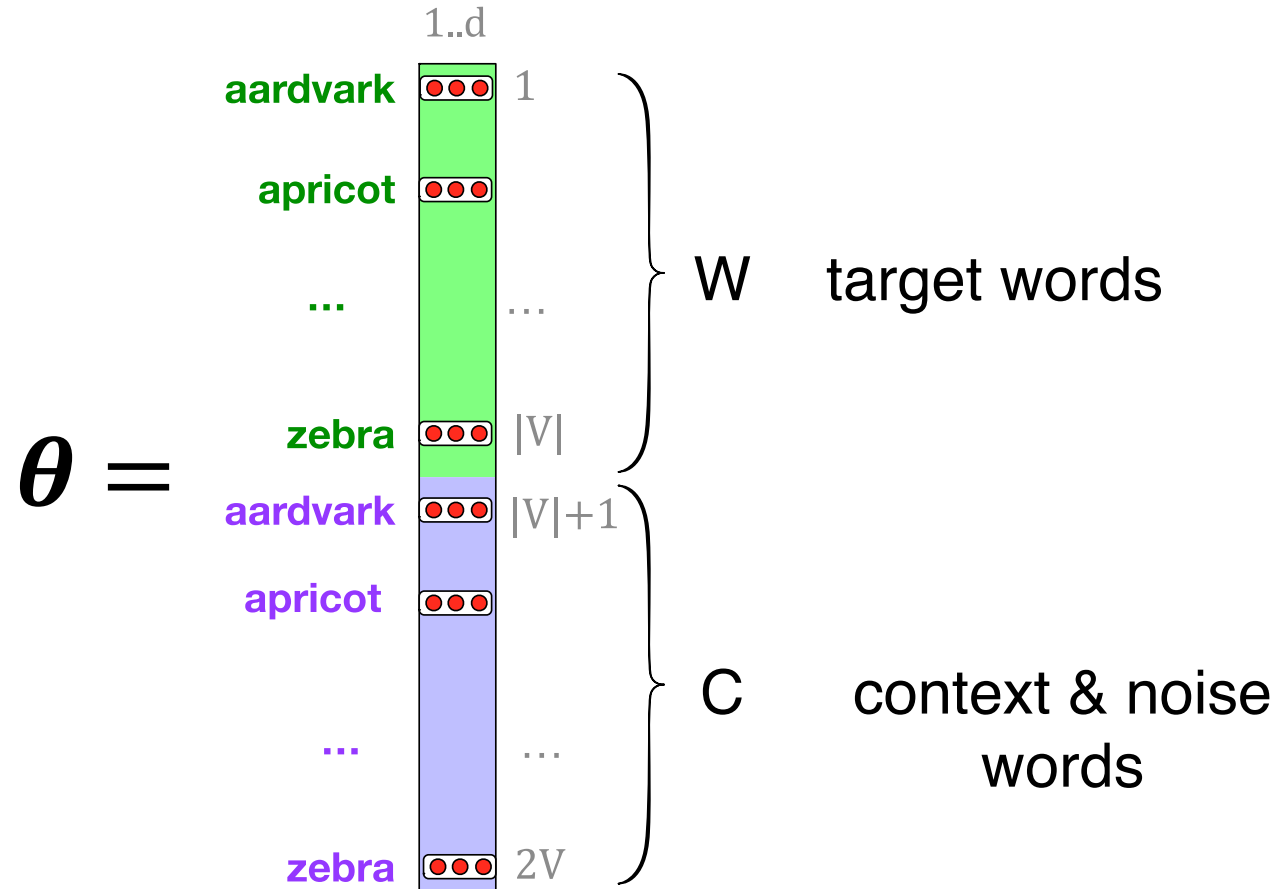
given a test target word w and

its context window of L words $c_{1:L}$,

assigns a probability that w occurs in this window

- To compute this, we just need embeddings for all the words

These embeddings we'll need: a set for w , a set for c

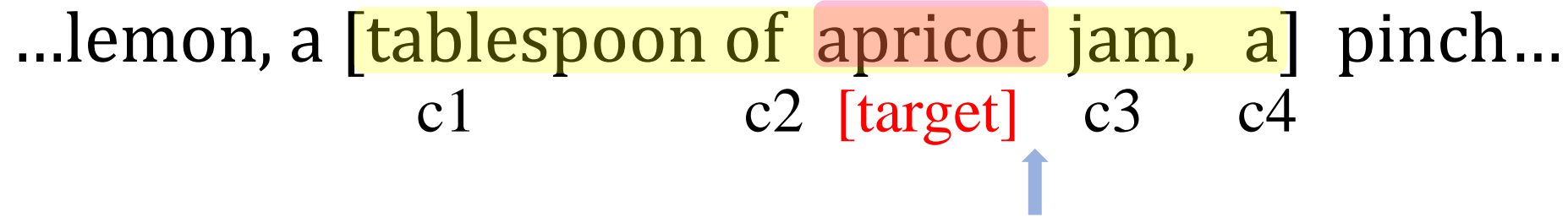


Word2vec: Learning the embeddings

Skip-Gram Training data

...lemon, a [tablespoon of apricot jam, a] pinch...

c1 c2 [target] c3 c4



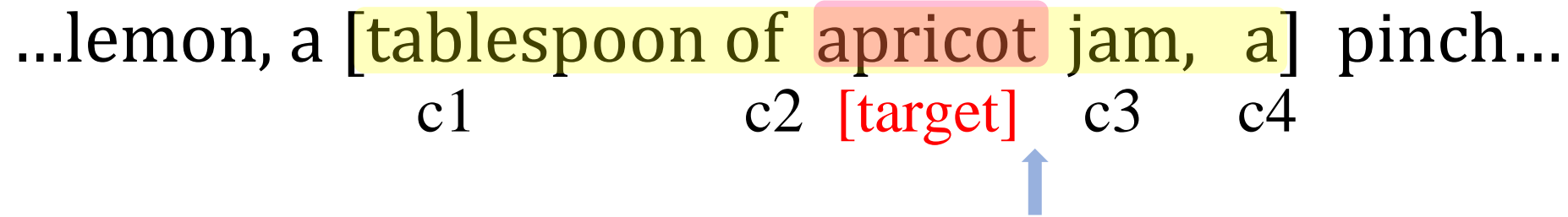
positive examples +

| t | c |
|---------|------------|
| apricot | tablespoon |
| apricot | of |
| apricot | jam |
| apricot | a |

Skip-Gram Training data

...lemon, a [tablespoon of apricot jam, a] pinch...

c1 c2 [target] c3 c4



positive examples +

t

c

apricot tablespoon

apricot of

apricot jam


apricot a

For each positive example
we'll grab k negative
examples, sampling by
frequency

Skip-Gram Training data

...lemon, a [tablespoon of apricot jam, a] pinch...

c1 c2 [target] c3 c4



positive examples +

| t | c |
|---------|------------|
| apricot | tablespoon |
| apricot | of |
| apricot | jam |
| apricot | a |

negative examples -

| t | c | t | c |
|---------|----------|---------|---------|
| apricot | aardvark | apricot | seven |
| apricot | my | apricot | forever |
| apricot | where | apricot | dear |
| apricot | coaxial | apricot | if |

Word2vec: how to learn vectors

- Given the set of positive and negative training instances, and an initial set of embedding vectors:
- The goal of learning is to adjust those word vectors such that we:
 - **Maximize** the similarity of the **target word**, **context word** pairs (w, c_{pos}) drawn from the positive data
 - **Minimize** the similarity of the (w, c_{neg}) pairs drawn from the negative data

Properties of Word Embeddings

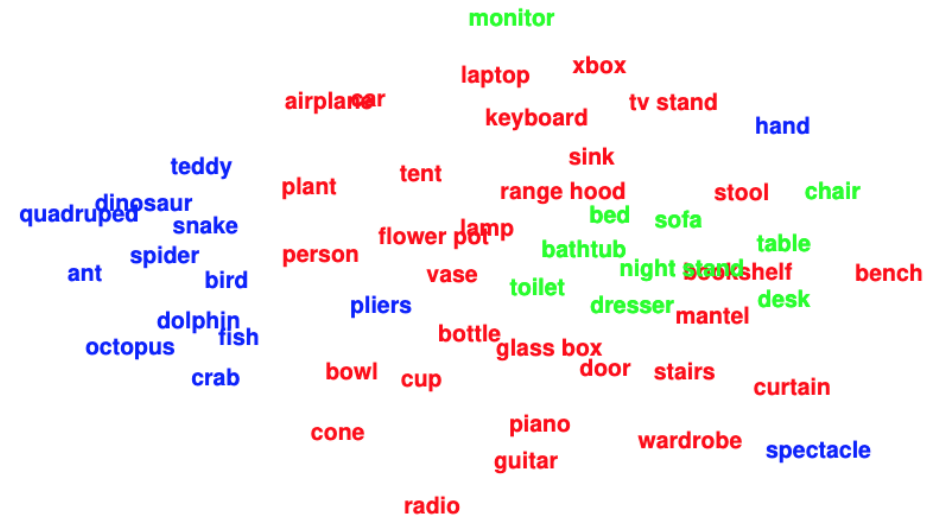
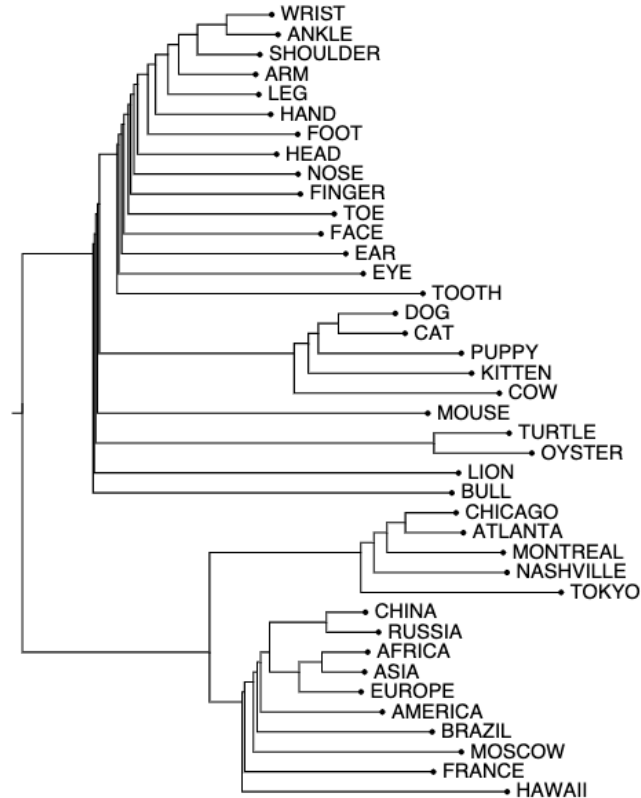
Effect of window size in dense embedding training

- **Shorter context windows** tend to lead to more syntactic representations as the information comes from immediately nearby words
 - The most similar words to a target word w tend to be **semantically similar words** with the **same parts of speech**
- For **longer context windows**, the most similar words to a target word w tend to be words that are **topically related but not similar**

Examples

- **Small windows** ($C = +/- 2$) : nearest words are similar nouns, words in the **same taxonomy**, e.g.:
 - *Hogwarts* nearest neighbors are other fictional schools *Sunnydale*, *Evernight*, *Blandings*
- **Large windows** ($C = +/- 5$): nearest words are related words in the **same semantic field**, e.g.:
 - *Hogwarts* nearest neighbors are terms from the Harry Potter world: *Dumbledore*, *Half-blood*, *Malfoy*

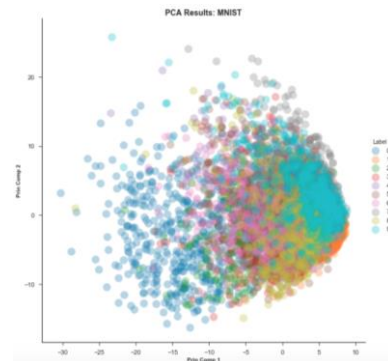
Visualizing Word Embeddings



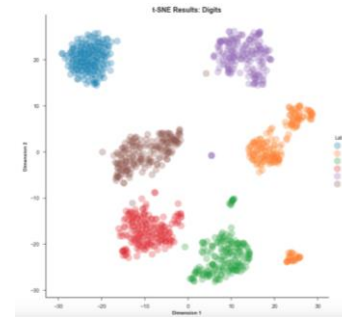
t-SNE

- **t-distributed stochastic neighbor embedding (t-SNE):** statistical method for visualizing high-dimensional data by giving each datapoint a location in a two or three-dimensional map
- A non-linear dimensionality reduction technique well-suited for embedding high-dimensional data
- Goal is to take a set of points in a high-dimensional space and find a faithful representation of those points in a lower-dimensional space, typically the 2D plane
- Implemented in Scikit-learn

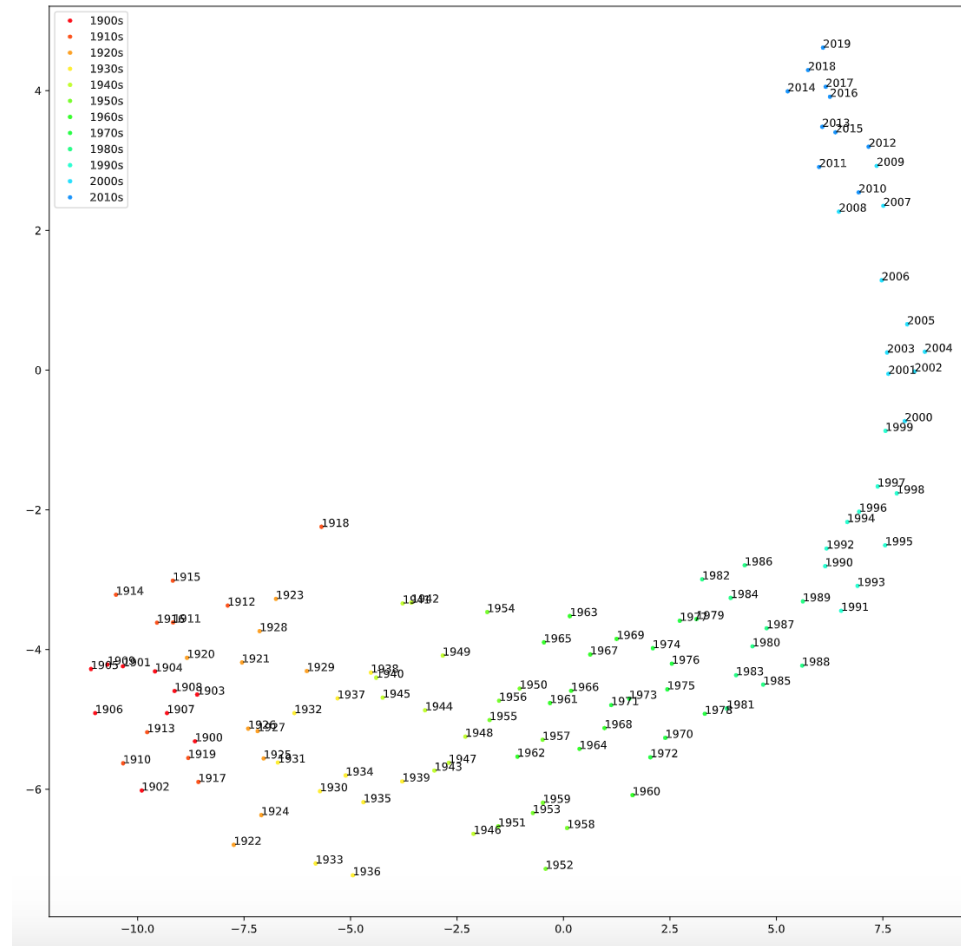
MNIST dataset by PCA



MNIST dataset by T-SNE



Example of visualizing year embeddings using t-SNE



Target: year references extracted from New York Times articles that were published from 1981 to 2013

Online Demos

[Turku NLP Group]

Models

Select one of the available models

English GoogleNews Negative300 ▼

Nearest words

Given a word, this demo shows a list of other words that are similar to it, i.e. nearby in the vector space.

cat Show nearest Case sensitive: ☒ Top N: 5 ▼

cats
dog
kitten
feline
beagle

Similarity of two words

Given two words, this demo gives the similarity value between 1 and -1.

Type in a word Type in a word Show similarity

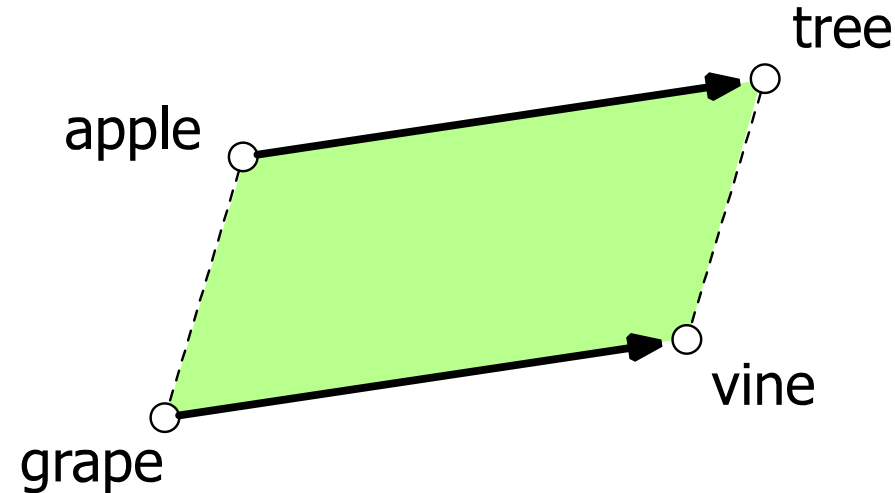
Word analogy

This demo computes word analogy: the first word is to the second word like the third word is to which word? Try for example *ilma - lintu - vesi* (air - bird - water) which would expect to return *kala* (fish) because fish is to water like birds is to air. Other cases could be for example *sammakko - hyppää - kala*. This is however only a toy to show what is possible - most of the time the analogy does not work particularly well (at least for the Finnish data).

Type in a word Type in a word Type in a word Show Top N: 2 ▼

Analogue relations & capturing relational meanings

- The classic parallelogram model of analogue reasoning (Rumelhart and Abrahamson 1973)
- To solve: *"apple is to tree as grape is to ____"*
- Add *apple - tree* to *grape* to get *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)

king – man + woman is close to queen

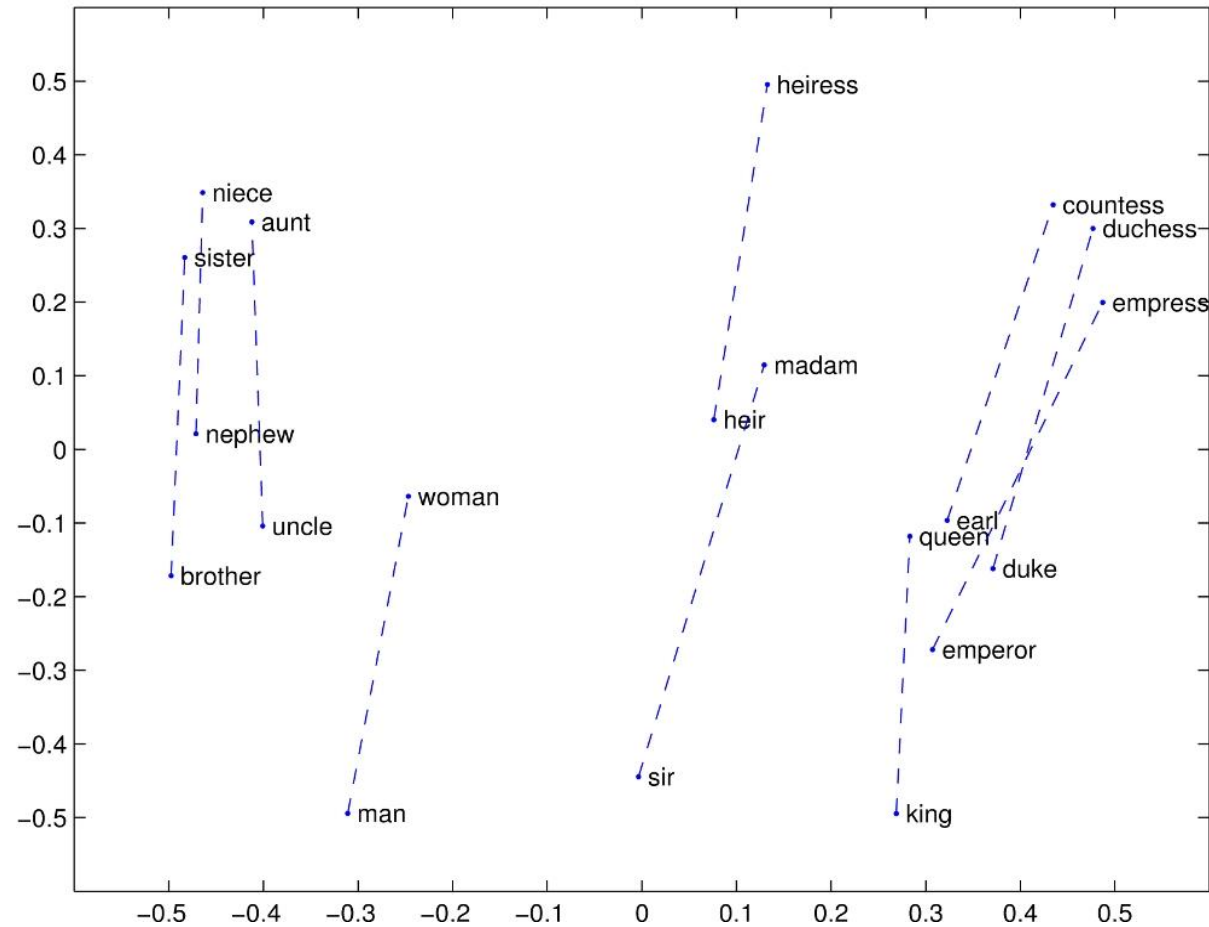
Paris – France + Italy is close to Rome

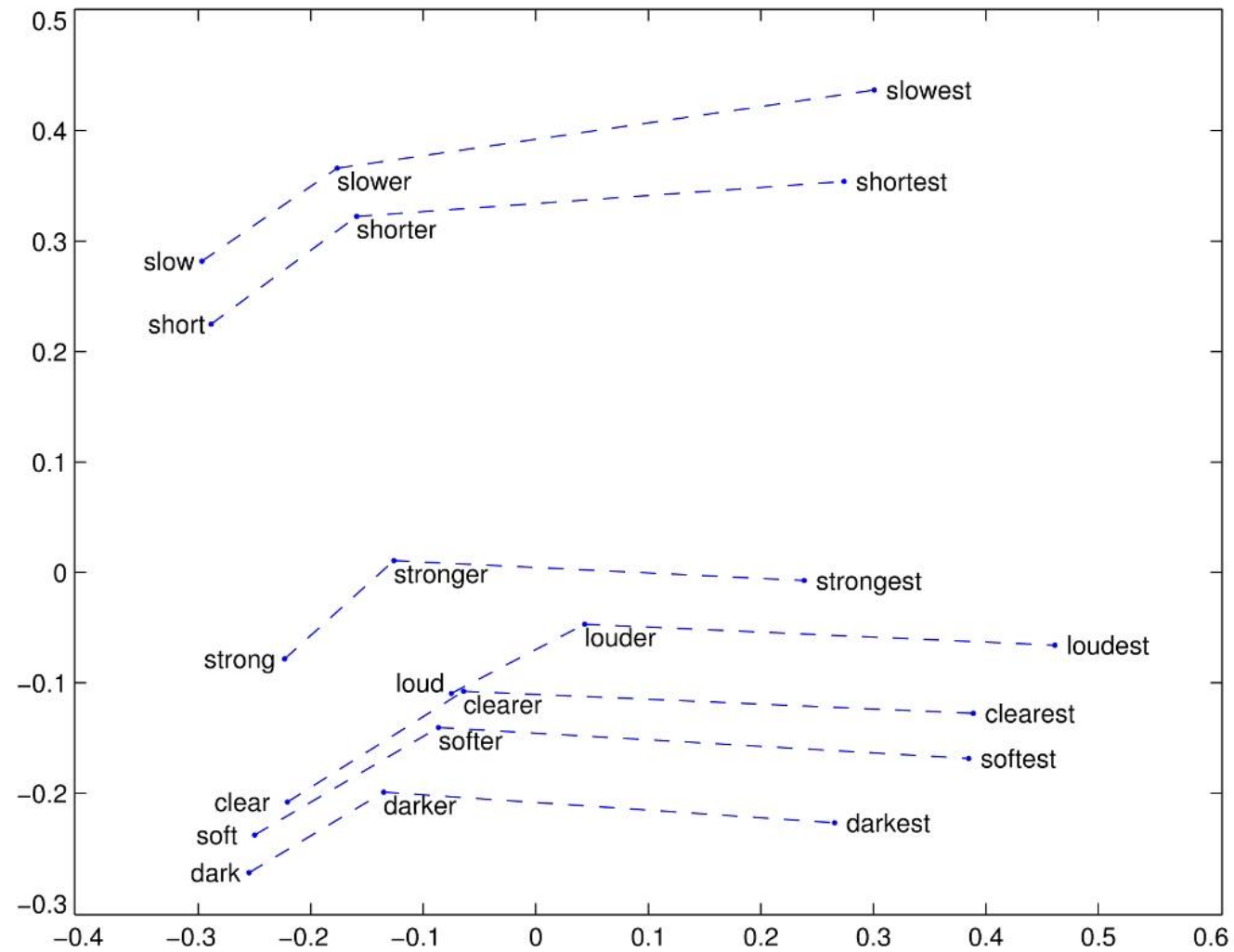
- For a problem $a:a^*::b:b^*$, the parallelogram method is:

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

↑
Cosine or Euclidean
distance

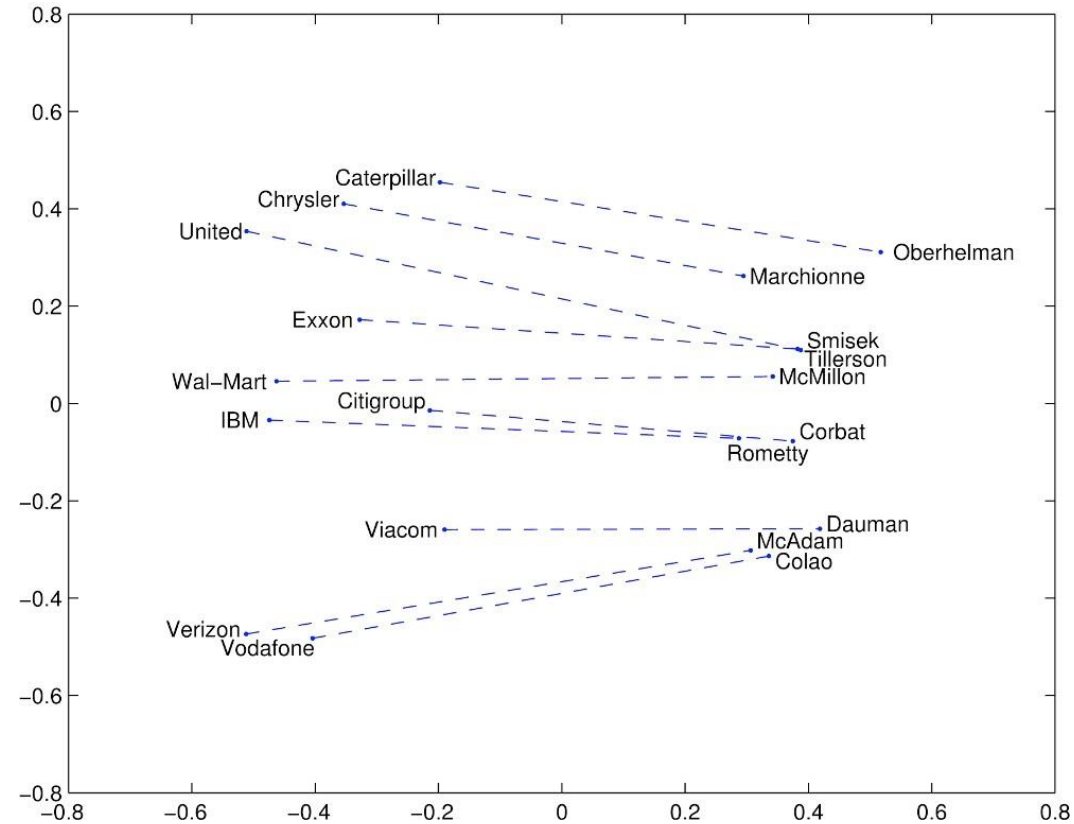
Structure Example in GloVe Embedding space





Comparative and superlative morphology

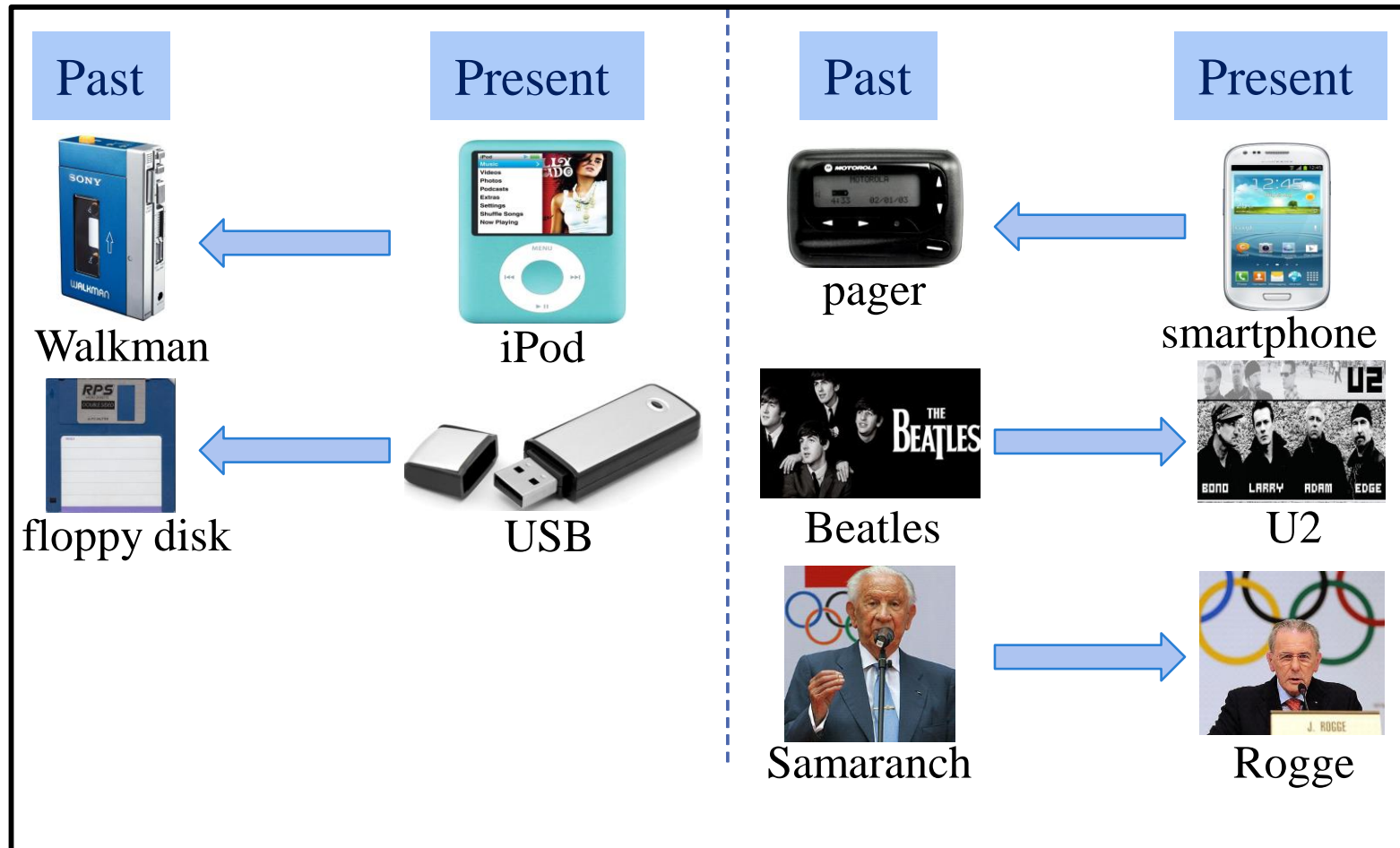
Glove Visualizations: Company - CEO



Caveats with the parallelogram method

- It only seems to work for frequent words, small distances and certain relations (e.g., relating countries to capitals, or parts of speech), but not others (Linzen 2016, Ethayarajh et al. 2019a)
- Often morphological variants are returned, e.g., *cherry:red :: potato:x* returns *potatoes* instead of brown so they must be explicitly excluded
- Understanding analogy is an open area of research (Peterson et al. 2020)...

Finding Temporal Analogy



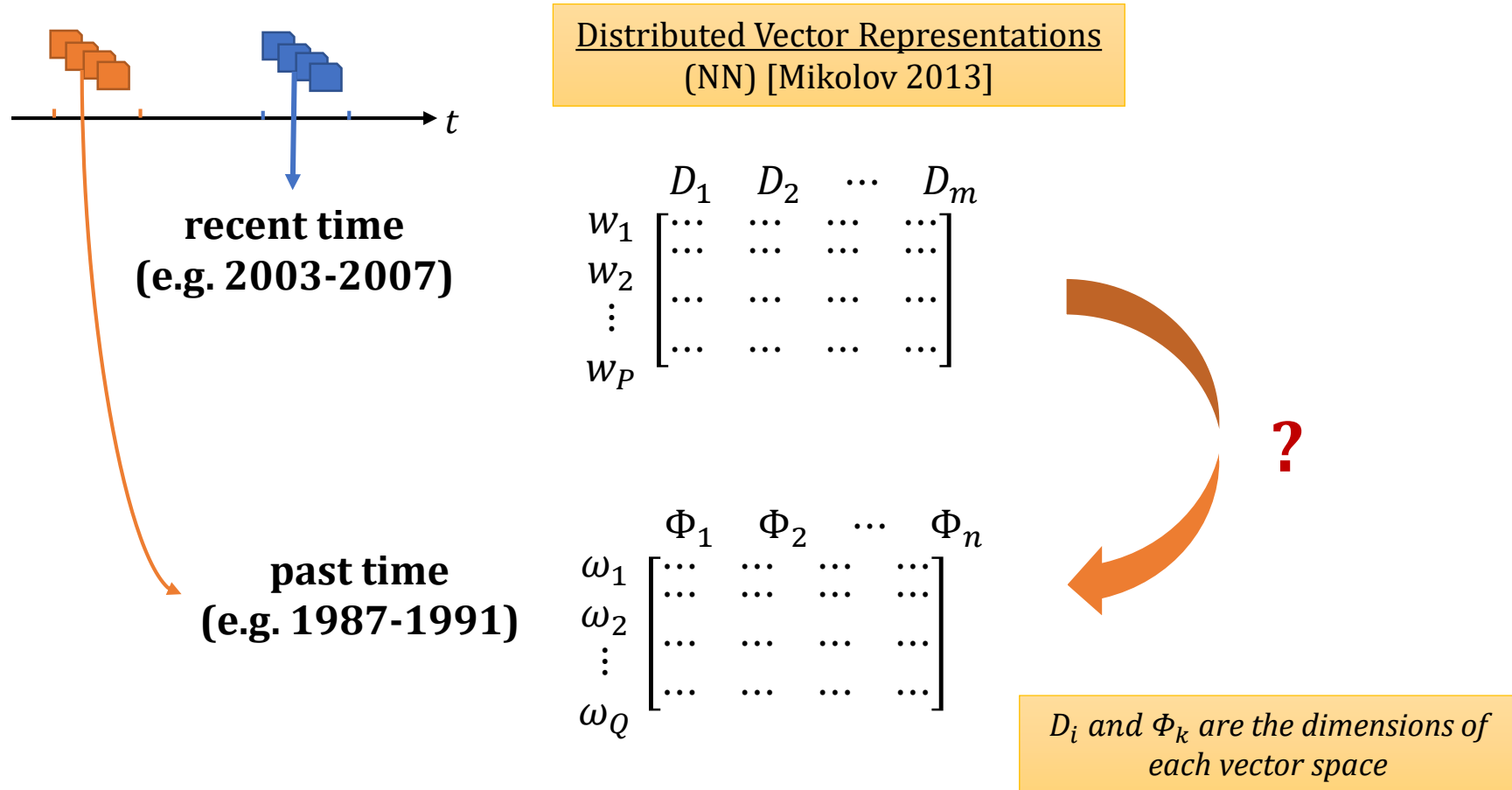
Finding Temporal Analogy: Panta Rei

- **Everything changes:** thus contexts surrounding *temporal analogs* are different

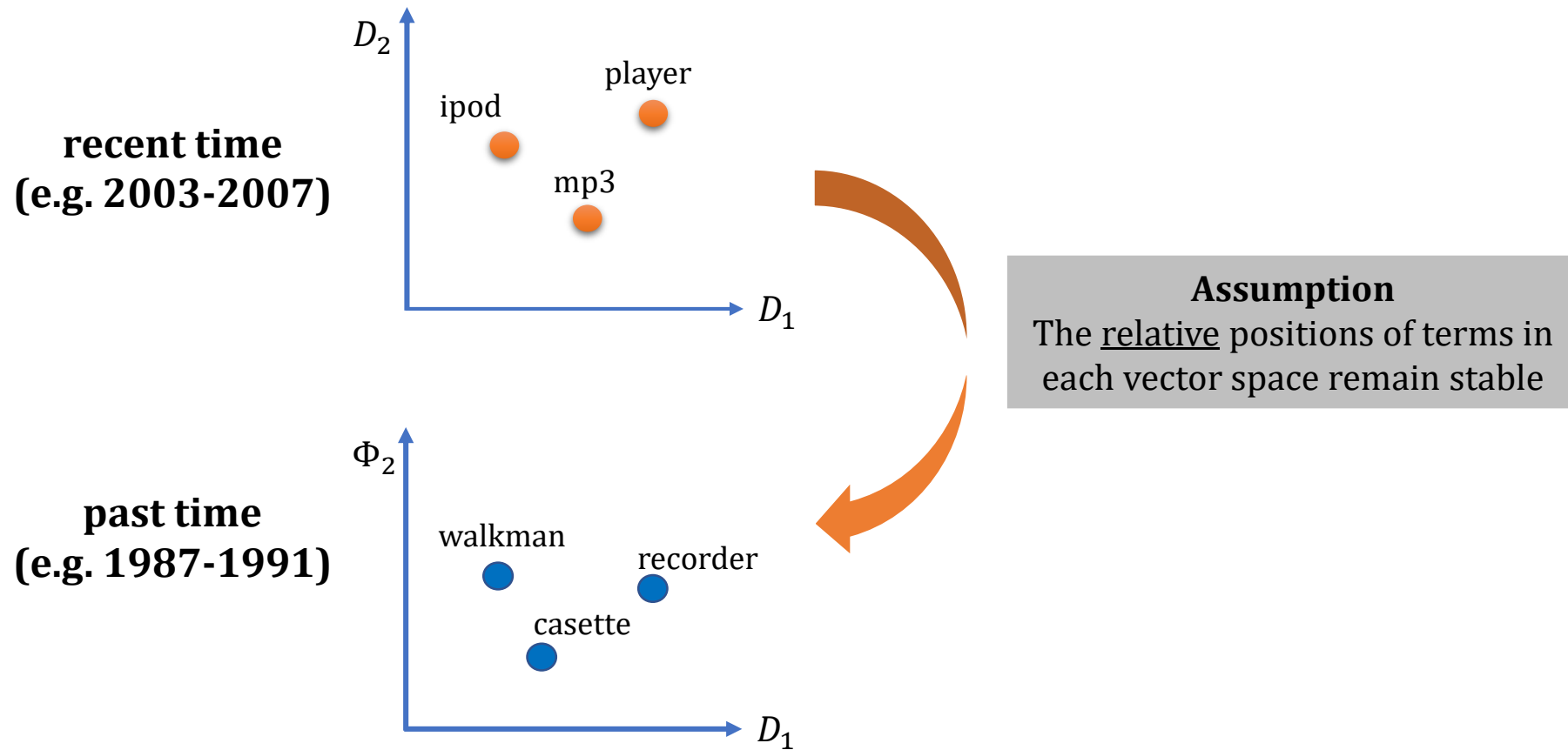
| Walkman (1980s) | iPod (2010s) |
|-----------------|--------------|
| cassette | apple |
| audio | mp3 |
| video | roqit |
| tape | player |
| music | music |
| sony | geeks |
| digital | jukebox |
| stereo | portable |
| earphone | macintosh |
| recorder | dlink |

* Contexts in the New York Times corpus

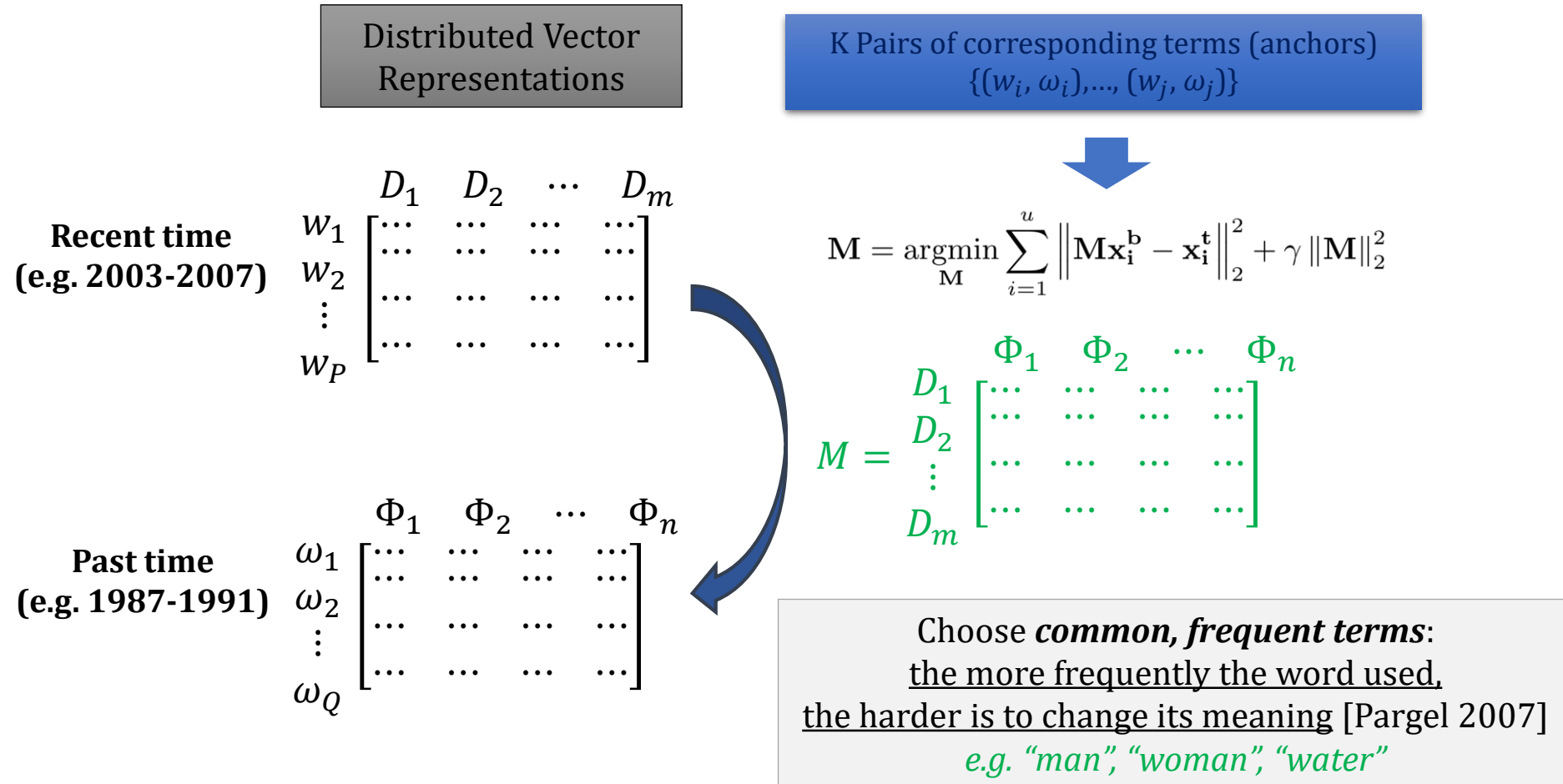
Finding Temporal Analogy: Making NN-based Term Embedding



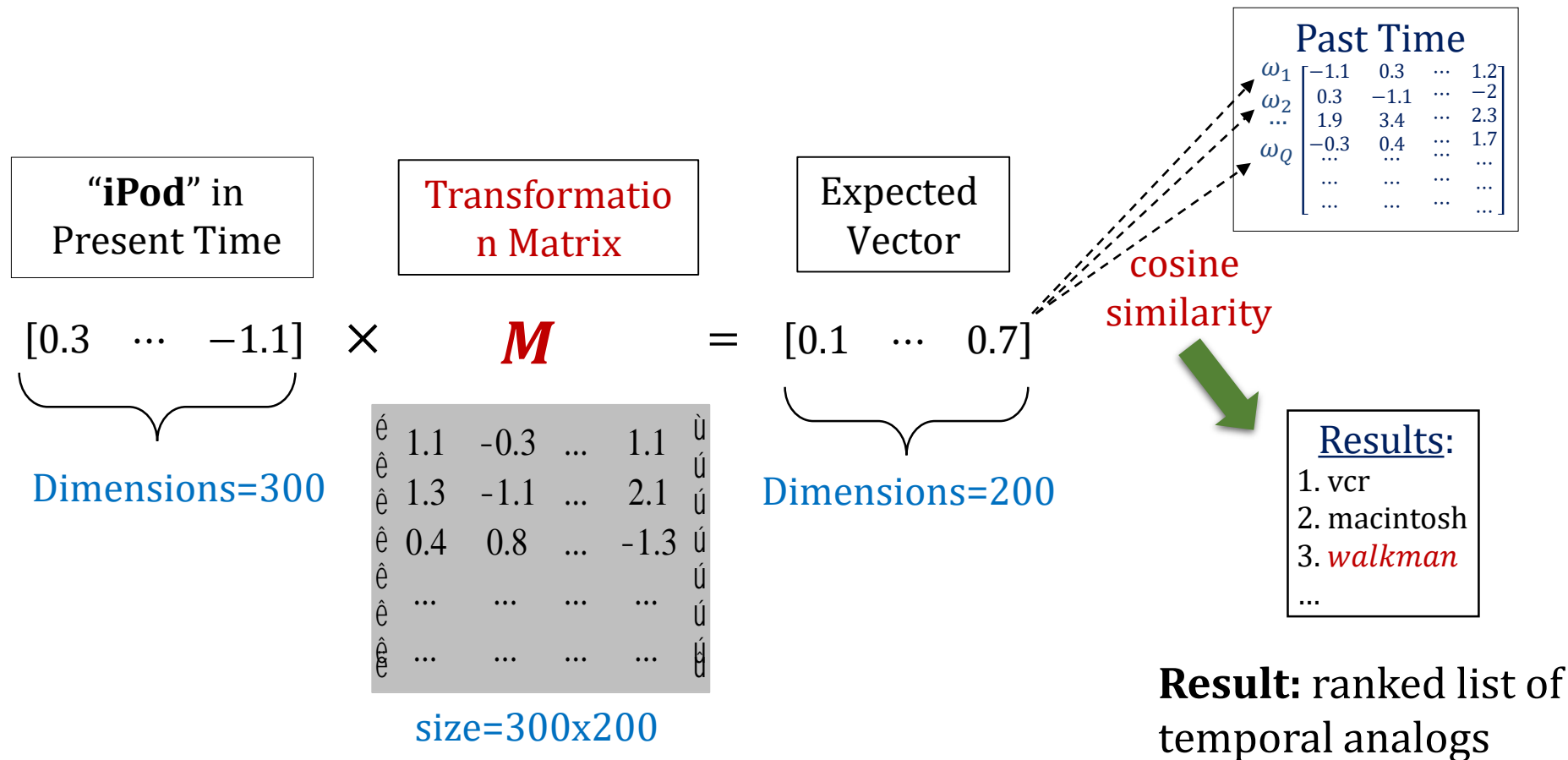
Finding Temporal Analogy: Assumption behind Proposed Approach



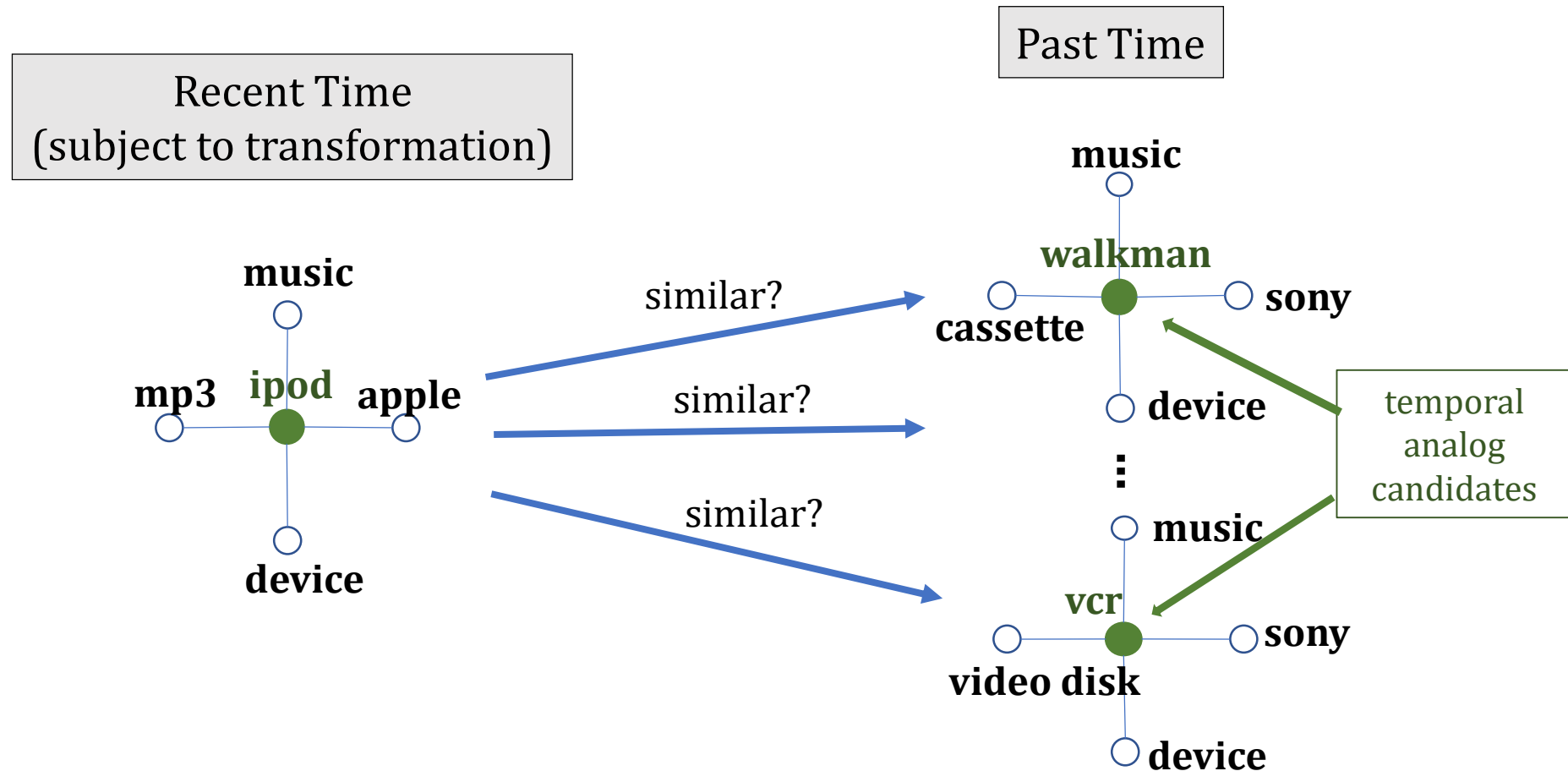
Finding Temporal Analogy: Constructing Transformation Matrix



Finding Temporal Analogy: Global Term Transformation Approach



Finding Temporal Analogy: Transformation Using Local Graph by Using Reference Points



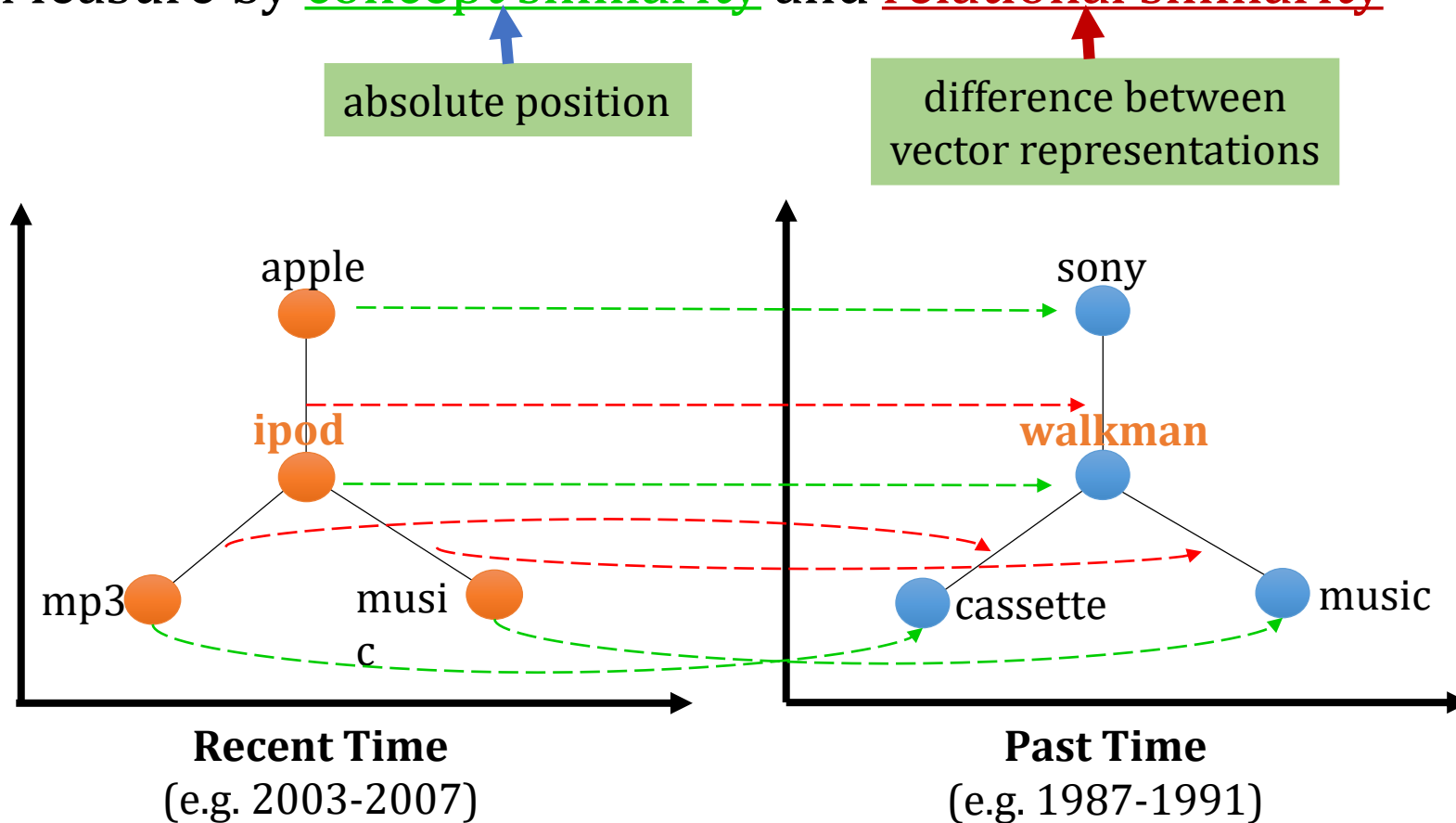
Finding Temporal Analogy: Desired Characteristics of Reference Points

- Reference Points - context terms which help to build effective across-time connection
- Desired criteria:
 - a) have **high relation** with the query
 - b) be sufficiently **general**
 - c) **independent** from each other

Finding Temporal Analogy: Local Graph Similarity Measurement

Approach:

Measure by concept similarity and relational similarity



Example Results

| | queries | correct answers | baselines | | methods | |
|-----|-------------|-----------------|-------------------|-----------------------|--------------|----------------------|
| | [2002,2007] | [1987,1991] | BOW
(baseline) | LSI+Com
(baseline) | Global_Trans | Local_Trans
(Lex) |
| 1 | Putin | Yeltsin | 1000+ | 51 | 24 | 2 |
| 2 | Chirac | Mitterrand | 1000+ | 6 | 7 | 2 |
| 3 | iPod | Walkman | 1000+ | 6 | 3 | 1 |
| 4 | Facebook | Usenet | 1000+ | 1000+ | 1 | 1 |
| 5 | Linux | Unix | 1000+ | 5 | 20 | 1 |
| 6 | spam | junk mail | 1000+ | 1000+ | 5 | 1 |
| 7 | spreadsheet | database | 1000+ | 395 | 3 | 1 |
| 9 | email | messages | 1000+ | 1 | 2 | 7 |
| 10 | email | letters | 1000+ | 1000+ | 1 | 1 |
| 11 | email | mail | 1000+ | 119 | 7 | 6 |
| 12 | email | fax | 1000+ | 1000+ | 3 | 4 |
| 14 | superman | batman | 1000+ | 46 | 5 | 2 |
| 15 | Pixar | Tristar | 1000+ | 110 | 1 | 1 |
| 16 | Pixar | Disney | 1000+ | 1 | 3 | 2 |
| 17 | Euro | Mark | 1000+ | 1000+ | 2 | 1 |
| 19 | Euro | Franc | 1000+ | 1000+ | 7 | 3 |
| 20 | Myanmar | Burma | 1000+ | 3 | 64 | 46 |
| 21 | Koizumi | Kaifu | 1000+ | 66 | 2 | 1 |
| 22 | NATO | NATO | 1000+ | 1 | 304 | 141 |
| 24 | fridge | freezer | 1000+ | 7 | 1 | 1 |
| 25 | fridge | refrigerator | 1000+ | 4 | 2 | 2 |
| 27 | Serbia | Yugoslavia | 1000+ | 12 | 1 | 1 |
| 28 | Kosovo | Yugoslavia | 1000+ | 27 | 14 | 10 |
| 30 | mp3 | compact disk | 1000+ | 44 | 58 | 19 |
| ... | ... | ... | ... | ... | ... | ... |

*Lexico-Syntactic Pattern used to detect reference points

Rank of correct answers

Finding Temporal Analogy: Aspect-based Retrieval System for Temporal Analogs

TempoAnalogus

Query in [2002,2007]:
euro

Past time period:
Select a time period

Method:
Select a method

Aspect term:
currency

Search

Reset

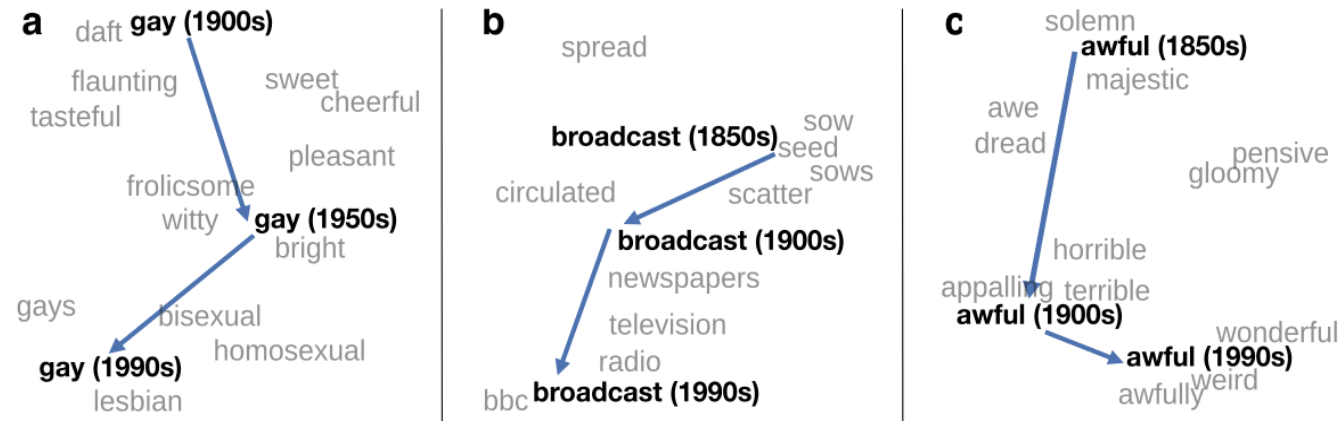
Temporal counterpart of **euro** biased on **currency** in [1987, 1991] is:

- francs : 0.609** ☒
but about nine billion **francs**, or \$250 million, of the aid depends on sabena's obtaining six billion **francs**, or about \$166 million, from a partner.
- belgian_francs : 0.574** ☒
lead: carlo de benedetti doubled his public offer tonight for societe generale de belgique's shares, from 4,000 **belgian francs** a share, or about \$113, to 8,000 **francs** in an attack on the **french-belgian** coalition that claims to have 52 percent of the vast holding company's capital.
- lire : 0.56** ☒
lead: *3*** company reports ** *3* de tomaso industries year to dec 31 1988 1987 sales 207,363,000 201,123,000 net loss 29,443,000 12,822,000 results are translated from italian **lire** at the exchange rate prevailing at dec.
- zloties : 0.544** ☐
the new official rate, which applies only to foreign tourists and foreign trade dealings, is 710 **zloties** to the dollar, compared with 680 on friday.
- lira : 0.538** ☒
lead: european officials were expected to consider devaluing the french franc and italian **lira** against the west german mark this weekend as the german currency's huge rise against the dollar intensified strains within the european monetary system.
- pecent : 0.538** ☐
5 percent stake in mixte to 30 percent, and mixte will cut its 12 **pecent** stake in the bank to 9.
- billion_pesetas : 0.537** ☐
22 **billion**, for the week ended wednesday, the investment company institute said thursday.
- dow_industrials : 0.536** ☐
the **dow** theory provided a bullish confirmation on tuesday, and another one yesterday, as the **dow** jones transportation average moved to record levels, while the **dow** jones industrial average climbed to its highest level since the 1987 crash.
- pound_sterling : 0.534** ☐
but ronald holzer, chief dealer for the harris trust and savings bank in chicago, said the dollar's rise against **sterling** was muted by the british currency's strength against the german mark and a flurry of other trading that helped the japanese yen and hurt the swiss franc.
- volume_shrank : 0.533** ☐
gains in agriculture sector the nation's trade surplus in agriculture jumped sharply despite the drought, the deficit in trade with japan dropped 15 percent and the nation's bill for imported oil declined as **volume shrank** and prices eased.

Feedback

Embeddings as a window onto historical semantics

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



Trained embeddings on different decades of historical text to see meanings shift

Embeddings reflect cultural and other kinds of biases

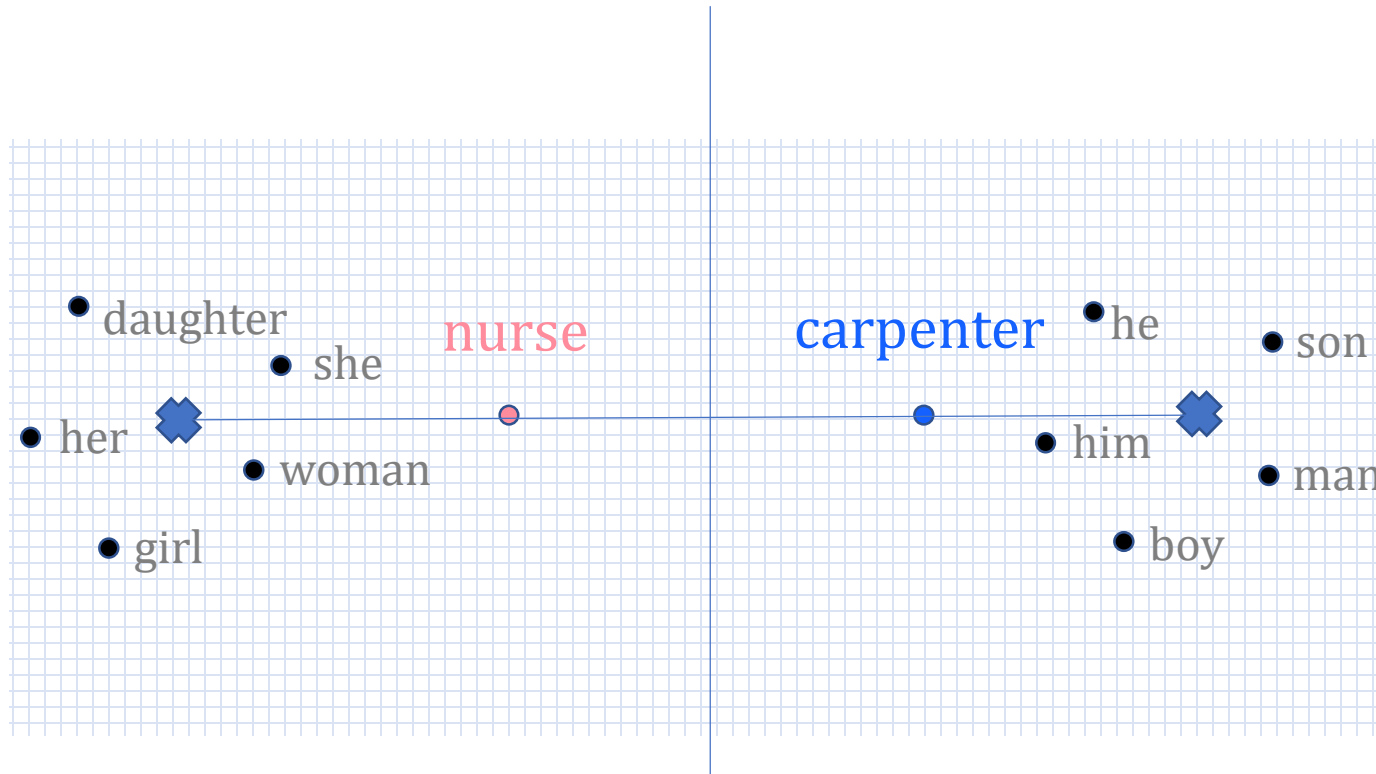
- Ask “Paris : France :: Tokyo : x”
 - x = Japan
- Ask “father : doctor :: mother : x”
 - x = nurse
- Ask “man : computer programmer :: woman : x”
 - x = homemaker

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

| Hispanic | Asian | White |
|--------------|------------|---------------|
| housekeeper | professor | smith |
| mason | official | blacksmith |
| artist | secretary | surveyor |
| janitor | conductor | sheriff |
| dancer | physicist | weaver |
| mechanic | scientist | administrator |
| photographer | chemist | mason |
| baker | tailor | statistician |
| cashier | accountant | clergy |
| driver | engineer | photographer |

The top ten occupations most closely associated with each ethnic group in the Google News embedding (Garg, 2018)

Computing the gender bias of a word



How much closer a word is to "woman" synonyms than "man" synonyms?