

Word Representations & Vector Space Models

Slido: <https://app.sli.do/event/o82YTccjg7nq3LoSr4NzPK>

Yu Meng

University of Virginia

yumeng5@virginia.edu

Sept 10, 2025

Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- **Week 3: Word Senses, Semantics & Classic Word Representations**
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

(Recap) Why Care About Word Semantics?

- Understanding word meanings helps us build better language models!
- Recall the example from N-gram lectures:

[BOS] The cat is on the mat [EOS]

[BOS] I have a cat and a mat [EOS]

[BOS] I like the cat [EOS]

$$p(\text{"cat"}|\text{"the"}) = \frac{2}{3}, \quad p(\text{"mat"}|\text{"the"}) = \frac{1}{3},$$

- Sparsity: many valid bigram counts are zero – count-based measures do not account for word semantics!
- If we know “cat” is semantically similar to “dog”, then $p(\text{"dog"}|\text{"the"}) \approx p(\text{"cat"}|\text{"the"})$

(Recap) What Types of Word Semantics Exist in NLP?

- **Synonyms:** words with similar meanings
 - “happy” & “joyful”
- **Antonyms:** words with opposite meanings
 - “hot” & “cold”
- **Hyponyms & hypernyms:** one word is a more specific instance of another
 - “rose” is a hyponym of “flower”
 - “flower” is a hypernym of “rose”
- **Polysemy:** A single word having multiple related meanings
 - “mouse” can mean small rodents or the device that controls a cursor
- The study of these aspects of word meanings is called **lexical semantics** in linguistics

(Recap) Lemmas

- **Lemma:** the base or canonical form of a word, from which other forms can be derived
 - “run” “runs” “ran” and “running” all share the lemma “run”
 - “better” and “best” share the lemma “good”
- **Lemmatization:** reducing words to their lemma
 - Allows models to recognize that different forms of a word carry the same meaning
 - An important pre-processing step in early NLP models
 - Contemporary LLMs (sort of) perform lemmatization through tokenization (later lectures!)

(Recap) Synonyms

- Word that have the same meaning in some or all contexts
- Two words are synonyms if they can be substituted for each other
- Perfect synonym is very rare!
 - Typically, words are slightly different in notions of politeness, connotation, genre/style...
 - “Child” vs. “kid”: “child” is often more formal/neutral; “kid” is more informal/casual
 - “Slim” vs. “skinny”: “slim” is often more positive in connotation than “skinny”
 - “Big” vs. “Large”: “big sister” is a common phrase but “large sister” is not

(Recap) Antonyms

- Words that have opposite meanings
- Gradable antonyms: exist on the ends of a spectrum or scale
 - “Hot” vs. “cold”
 - “Tall” vs. “short”
- Complementary antonyms: the presence of one directly excludes the other
 - “Alive” vs. “dead”
 - “True” vs. “false”
- Relational antonyms: express a relationship between two dependent entities
 - “Teacher” vs. “student”
 - “Buyer” vs. “seller”

(Recap) Hyponyms & Hypernyms

- Describe hierarchical relationships between words based on specificity and generality
- **Hypernym** is a word that is more general/broader in meaning and can encompass a variety of more specific words
- **Hyponym** is a word that is more specific in meaning and falls under a broader category
- “Vehicle” is a hypernym for “car” “bicycle” “airplane” “boat” etc.
- “Car” “bicycle” “airplane” “boat” are hyponyms of “vehicle”
- **Hypernym/hyponym** relationship is usually transitive
 - A is a hypernym of B; B is a hypernym of C => A is a hypernym of C



(Recap) Polysemy & Senses

- **Polysemy**: a single word has multiple related meanings
 - “**Light**”: “This bag is **light**” / “Turn on the **light**” / “She made a **light** comment”
- **Sense**: a particular meaning or interpretation of a word in a given context
- Word relations (e.g., synonyms, antonyms, hypernyms/hyponyms) are defined between word senses!
- **Word sense disambiguation (WSD)**: determine which sense of a word is being used in a specific context
 - She went to the **bank** to deposit money
 - She lives by the river **bank**
- WSD can be challenging especially when the context is short/insufficient
 - Is the query “mouse info” looking for a pet or a tool?

(Recap) Word Sense Disambiguation

WSD can be an interesting/challenging test case even for the strong (multimodal) LLMs



Image generated by Nano Banana
under the user prompt: *“generate
an image of a baseball player caring
for his bat in the cave where he lives
with all the other bats”*



(Recap) Word Similarity

- Most words may not have many perfect synonyms, but usually have lots of similar words
 - “cat” is not a synonym of “dog”, but they are similar in meaning

vanish	disappear	9.8
belief	impression	5.95
muscle	bone	3.65
modest	flexible	0.98
hole	agreement	0.3

Word similarity (on a scale from 0 to 10)
manually annotated by humans

- We’ll introduce word embeddings to automatically learn word similarity next week!

(Recap) Word Relatedness & Semantic Field

- **Word relatedness:** the meaning of words can be related in ways other than similarity
 - Functional relationship: “doctor” and “hospital” – doctors work in hospitals
 - Thematic relationship: “bread” and “butter” – often used together in the context of food
 - Conceptual relationship: “teacher” and “chalkboard” – both part of the educational context
- **Semantic field:** a set of words which cover a particular semantic domain and bear structured relations with each other
 - Semantic field of “houses”: door, roof, kitchen, family, bed...
 - Semantic field of “restaurants”: waiter, menu, plate, food, chef...
 - Semantic field of “hospitals”: surgeon, nurse, anesthetic, scalpel...

(Recap) Connotation

- Subjective/cultural/emotional associations that words carry beyond their literal meanings
 - Youthful (positive) vs. childish (negative)
 - Confident (positive) vs. arrogant (negative)
 - Economical (positive) vs. cheap (negative)
- Connotation can be described via three 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

(Recap) Connotation

- Valence: the pleasantness of the stimulus
 - High: “happy” / “satisfied”; low: “unhappy” / “annoyed”
- Arousal: the intensity of emotion provoked by the stimulus
 - High: “excited”; low: “calm”
- Dominance: the degree of control exerted by the stimulus
 - High: “controlling”; low: “influenced”

	Valence	Arousal	Dominance
courageous	8.05	5.5	7.38
music	7.67	5.57	6.5
heartbreak	2.45	5.65	3.58
cub	6.71	3.95	4.24

Earliest work on representing words
with multi-dimensional vectors!

(Recap) WordNet

- Word semantics is complex (multiple senses, various relations)!
- How did people represent word senses and relations in early NLP developments?
- WordNet: A manually curated large lexical database
- Three separate databases: one each for nouns, verbs and adjectives/adverbs
- Each database contains a set of lemmas, each one annotated with a set of senses
- Synset (synonym set): The set of near-synonyms for a sense
- Word relations (hypernym, hyponym, antonym) defined between synsets

(Recap) WordNet Relations

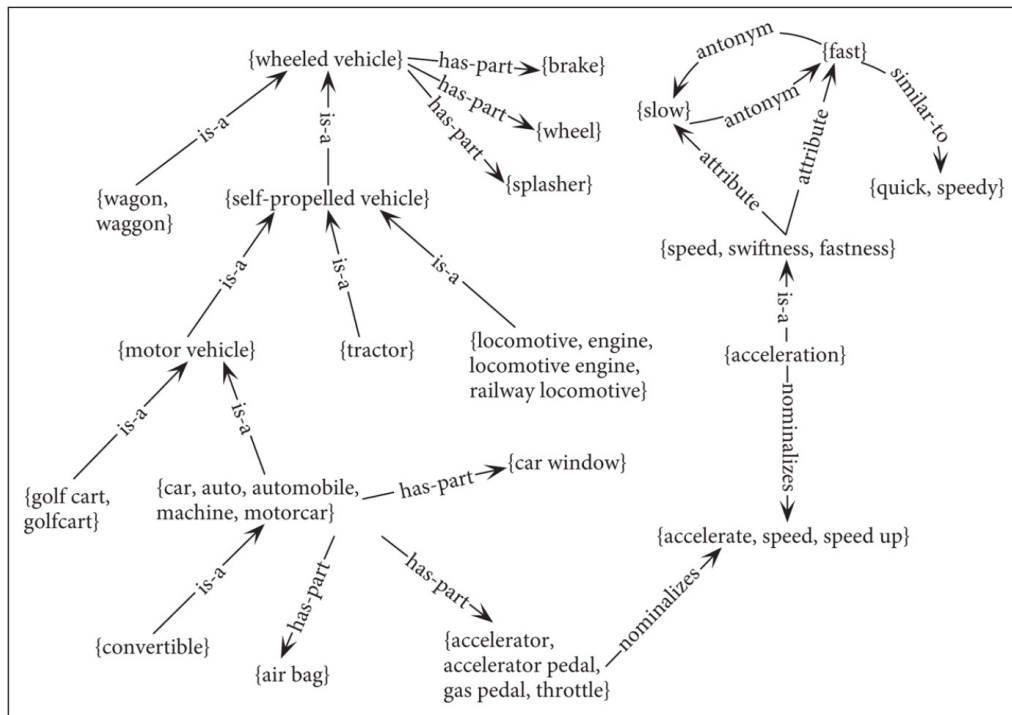
Relation	Also Called	Definition	Example
Hypernym	Superordinate	From concepts to superordinates	<i>breakfast</i> ¹ → <i>meal</i> ¹
Hyponym	Subordinate	From concepts to subtypes	<i>meal</i> ¹ → <i>lunch</i> ¹
Instance Hypernym	Instance	From instances to their concepts	<i>Austen</i> ¹ → <i>author</i> ¹
Instance Hyponym	Has-Instance	From concepts to their instances	<i>composer</i> ¹ → <i>Bach</i> ¹
Part Meronym	Has-Part	From wholes to parts	<i>table</i> ² → <i>leg</i> ³
Part Holonym	Part-Of	From parts to wholes	<i>course</i> ⁷ → <i>meal</i> ¹
Antonym		Semantic opposition between lemmas	<i>leader</i> ¹ ⇔ <i>follower</i> ¹
Derivation		Lemmas w/same morphological root	<i>destruction</i> ¹ ⇔ <i>destroy</i> ¹

Noun relations

Relation	Definition	Example
Hypernym	From events to superordinate events	<i>fly</i> ⁹ → <i>travel</i> ⁵
Troponym	From events to subordinate event	<i>walk</i> ¹ → <i>stroll</i> ¹
Entails	From verbs (events) to the verbs (events) they entail	<i>snore</i> ¹ → <i>sleep</i> ¹
Antonym	Semantic opposition between lemmas	<i>increase</i> ¹ ⇔ <i>decrease</i> ¹

Verb relations

(Recap) WordNet as a Graph





(Recap) WordNet Demo

Category	Unique Strings
Noun	117798
Verb	11529
Adjective	22479
Adverb	4481

Figure source: <https://lm-class.org/lectures/04%20-%20word%20embeddings.pdf>

Word to search for:

Display Options:

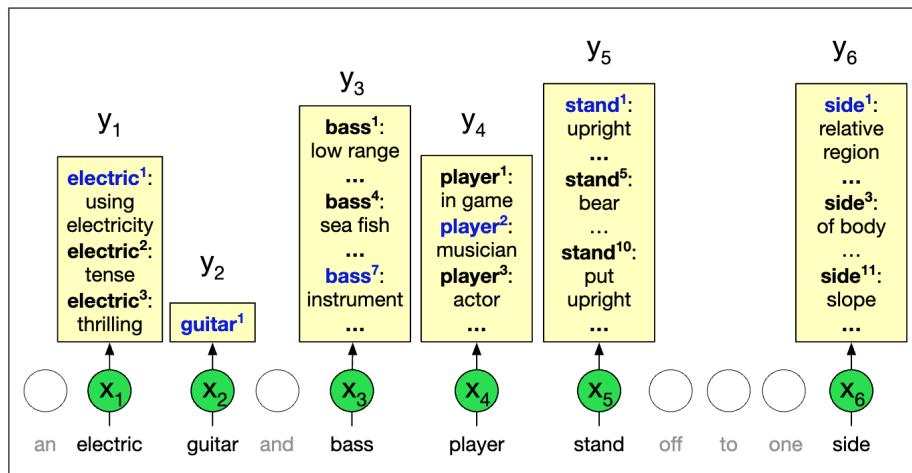
Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations
Display options for sense: (gloss) "an example sentence"

Noun

- **S: (n) light, visible light, visible radiation** ((physics) electromagnetic radiation that can produce a visual sensation) *"the light was filtered through a soft glass window"*
 - [direct hyponym](#) / [full hyponym](#)
 - [domain category](#)
 - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - [part holonym](#)
 - [derivationally related form](#)
- **S: (n) light, light source** (any device serving as a source of illumination) *"he stopped the car and turned off the lights"*
- **S: (n) light** (a particular perspective or aspect of a situation) *"although he saw it in a different light, he still did not understand"*
- **S: (n) luminosity, brightness, brightness level, luminance, luminousness, light** (the quality of being luminous; emitting or reflecting light) *"its luminosity is measured relative to that of our sun"*
- **S: (n) light** (an illuminated area) *"he stepped into the light"*
 - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - [derivationally related form](#)
- **S: (n) light, illumination** (a condition of spiritual awareness; divine illumination) *"follow God's light"*
- **S: (n) light, lightness** (the visual effect of illumination on objects or scenes as created in pictures) *"he could paint the lightest light and the darkest dark"*
- **S: (n) light** (a person regarded very fondly) *"the light of my life"*
- **S: (n) light, lighting** (having abundant light or illumination) *"they played as long as it was light"; "as long as the lighting was good"*
- **S: (n) light** (mental understanding as an enlightening experience) *"he finally saw the light"; "can you shed light on this problem?"*
- **S: (n) sparkle, twinkle, spark, light** (merriment expressed by a brightness or gleam or animation of countenance) *"he had a sparkle in his eye"; "there's a perpetual twinkle in his eyes"*
- **S: (n) light** (public awareness) *"it brought the scandal to light"*
- **S: (n) Inner Light, Light, Light Within, Christ Within** (a divine presence)

(Recap) WordNet for Word Sense Disambiguation

- All words WSD task: map all input words (nouns/verbs/adjectives/adverbs) to WordNet senses
- Strong baseline: map to the first sense in WordNet (most frequent)
- Modern approaches: sequence modeling architectures (later lectures!)



(Recap) WordNet Limitations

- Require significant efforts to construct and maintain/update
 - Hard to keep up with rapidly evolving language usage
- Limited coverage of domain-specific terms & low-resource language
 - No coverage of specialized, domain-specific terms (e.g., medical, legal, or technical)
- Only support individual words and their meanings
 - Do not account for idiomatic expressions, phrasal verbs, or collocations

A more automatic, scalable, and contextualized word semantic learning approach is needed!

Agenda

- Introduction to Word Senses & Semantics
- Classic Word Representations
- Vector Space Model Basics

Motivation: Representing Texts with Vectors

- Word similarity computation is important for understanding semantics

Word similarity (on a scale from 0 to 10)
manually annotated by humans

vanish	disappear	9.8
belief	impression	5.95
muscle	bone	3.65
modest	flexible	0.98
hole	agreement	0.3

Word semantics can be multi-faceted

	Valence	Arousal	Dominance
courageous	8.05	5.5	7.38
music	7.67	5.57	6.5
heartbreak	2.45	5.65	3.58
cub	6.71	3.95	4.24

- How to represent words numerically? Using multi-dimensional vectors!

Vector Semantics

- Represent a word as a point in a multi-dimensional semantic space
- A desirable vector semantic space: words with similar meanings are nearby in space



2D visualization of a desirable high-dimensional vector semantic space

Vector Space Basics

- Vector notation: an N-dimensional vector $\mathbf{v} = [v_1, v_2, \dots, v_N] \in \mathbb{R}^N$
- Vector dot product/inner product:

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

- Vector length/norm:

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

Other (less commonly-used) vector norms:
Manhattan norm, p -norm, infinity norm...

- Cosine similarity between vectors:

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



Vector Space Basics: Example

- Consider two 4-dimensional vectors $\mathbf{v} = [1, 0, 1, 0] \in \mathbb{R}^4$ $\mathbf{w} = [0, 1, 1, 0] \in \mathbb{R}^4$
- Vector dot product/inner product:

$$\mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^N v_i w_i = 1$$

- Vector length/norm:

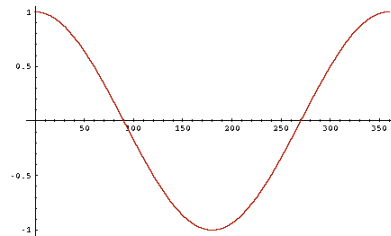
$$|\mathbf{v}| = \sqrt{\sum_{i=1}^N v_i^2} = \sqrt{2} \quad |\mathbf{w}| = \sqrt{\sum_{i=1}^N w_i^2} = \sqrt{2}$$

- Cosine similarity between vectors:

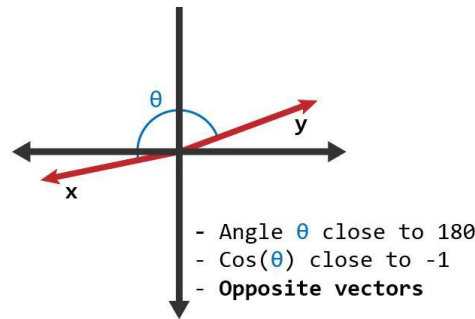
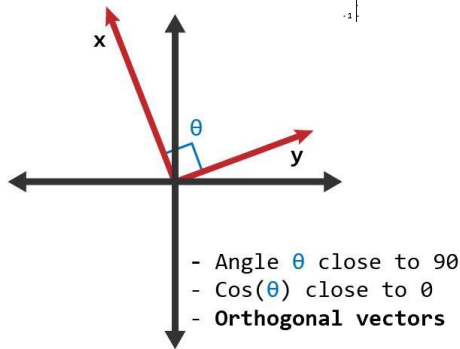
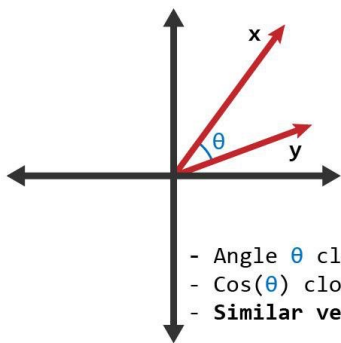
$$\cos(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v} \cdot \mathbf{w}}{|\mathbf{v}| |\mathbf{w}|} = \frac{1}{2}$$

Vector Similarity

- Cosine similarity is the most commonly used metric for similarity measurement
 - Symmetric: $\cos(\mathbf{v}, \mathbf{w}) = \cos(\mathbf{w}, \mathbf{v})$
 - Not influenced by vector length
 - Has a normalized range: $[-1, 1]$
 - Intuitive geometric interpretation



Cosine function values under different angles



How to Represent Words as Vectors?

- Given a vocabulary $\mathcal{V} = \{\text{good, feel, I, sad, cats, have}\}$
- Most straightforward way to represent words as vectors: use their indices
- One-hot vector: only one high value (1) and the remaining values are low (0)
- Each word is identified by a unique dimension

$$\mathbf{v}_{\text{good}} = [1, 0, 0, 0, 0, 0]$$

$$\mathbf{v}_{\text{feel}} = [0, 1, 0, 0, 0, 0]$$

$$\mathbf{v}_{\text{I}} = [0, 0, 1, 0, 0, 0]$$

$$\mathbf{v}_{\text{sad}} = [0, 0, 0, 1, 0, 0]$$

$$\mathbf{v}_{\text{cats}} = [0, 0, 0, 0, 1, 0]$$

$$\mathbf{v}_{\text{have}} = [0, 0, 0, 0, 0, 1]$$

Represent Sequences by Word Occurrences

- Consider the mini-corpus with three documents

$d_1 = \text{"I feel good"}$

$d_2 = \text{"I feel sad"}$

$d_3 = \text{"I have cats"}$

$$\mathbf{v}_{\text{good}} = [1, 0, 0, 0, 0, 0]$$

$$\mathbf{v}_{\text{feel}} = [0, 1, 0, 0, 0, 0]$$

$$\mathbf{v}_{\text{I}} = [0, 0, 1, 0, 0, 0]$$

$$\mathbf{v}_{\text{sad}} = [0, 0, 0, 1, 0, 0]$$

$$\mathbf{v}_{\text{cats}} = [0, 0, 0, 0, 1, 0]$$

$$\mathbf{v}_{\text{have}} = [0, 0, 0, 0, 0, 1]$$

- Straightforward way of representing documents: look at which words are present

$$\mathbf{v}_{d_1} = [1, 1, 1, 0, 0, 0]$$

$$\mathbf{v}_{d_2} = [0, 1, 1, 1, 0, 0]$$

$$\mathbf{v}_{d_3} = [0, 0, 1, 0, 1, 1]$$

Document vector similarity



$$\cos(\mathbf{v}_{d_1}, \mathbf{v}_{d_2}) = \frac{2}{3}$$

$$\cos(\mathbf{v}_{d_1}, \mathbf{v}_{d_3}) = \frac{1}{3}$$

$$\cos(\mathbf{v}_{d_2}, \mathbf{v}_{d_3}) = \frac{1}{3}$$



Term-Document Matrix

- With larger text collections, word frequencies in documents entail rich information
- Consider the four plays by Shakespeare and obtain the word frequency statistics
- Look at 4 manually-picked words: “battle” “good” “fool” “wit”

	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

There are many more words!

- Document vector representation with word frequencies:

$$\mathbf{v}_{d_1} = [1, 114, 36, 20] \quad \mathbf{v}_{d_2} = [0, 80, 58, 15] \quad \mathbf{v}_{d_3} = [7, 62, 1, 2] \quad \mathbf{v}_{d_4} = [13, 89, 4, 3]$$



Document Similarity

- Document vector representation with word frequencies:

$$\mathbf{v}_{d_1} = [1, 114, 36, 20] \quad \mathbf{v}_{d_2} = [0, 80, 58, 15] \quad \mathbf{v}_{d_3} = [7, 62, 1, 2] \quad \mathbf{v}_{d_4} = [13, 89, 4, 3]$$

	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

- “fool” and “wit” occur much more frequently in d_1 and d_2 than d_3 and d_4
- d_1 and d_2 are comedies $\cos(\mathbf{v}_{d_1}, \mathbf{v}_{d_2}) = 0.95$ $\cos(\mathbf{v}_{d_2}, \mathbf{v}_{d_3}) = 0.81$
- Word frequencies in documents do reflect the semantic similarity between documents!

Words Represented with Documents

- “Battle”: “the kind of word that occurs in Julius Caesar and Henry V (history plays)”
- “Fool”: “the kind of word that occurs in comedies”

	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

- Represent words using their co-occurrence counts with documents:

$$\mathbf{v}_{\text{battle}} = [1, 0, 7, 13]$$

$$\mathbf{v}_{\text{good}} = [114, 80, 62, 89]$$

$$\mathbf{v}_{\text{fool}} = [36, 58, 1, 4]$$

$$\mathbf{v}_{\text{wit}} = [20, 15, 2, 3]$$

Words Represented with Documents

	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

$$\mathbf{v}_{\text{battle}} = [1, 0, 7, 13]$$

$$\mathbf{v}_{\text{good}} = [114, 80, 62, 89]$$

$$\mathbf{v}_{\text{fool}} = [36, 58, 1, 4]$$

$$\mathbf{v}_{\text{wit}} = [20, 15, 2, 3]$$



$$\cos(\mathbf{v}_{\text{fool}}, \mathbf{v}_{\text{wit}}) = 0.93$$

$$\cos(\mathbf{v}_{\text{fool}}, \mathbf{v}_{\text{battle}}) = 0.09$$

Previously:

$$\mathbf{v}_{\text{battle}} = [1, 0, 0, 0]$$

$$\mathbf{v}_{\text{good}} = [0, 1, 0, 0]$$

$$\mathbf{v}_{\text{fool}} = [0, 0, 1, 0]$$

$$\mathbf{v}_{\text{wit}} = [0, 0, 0, 1]$$



$$\cos(\mathbf{v}_{\text{fool}}, \mathbf{v}_{\text{wit}}) = 0$$

$$\cos(\mathbf{v}_{\text{fool}}, \mathbf{v}_{\text{battle}}) = 0$$

Document co-occurrence statistics provide coarse-grained contexts

Fine-Grained Contexts: Word-Word Matrix

Instead of using documents as contexts for words, we can also use words as contexts

4 words to the left	center word	4 words to the right
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



Fine-Grained Contexts: Word-Word Matrix

Count how many times words occur in a ± 4 word window around the center word

context word

center word

	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	...

Counts derived from the Wikipedia corpus



Word Similarity Based on Word Co-occurrence

- Word-word matrix with ± 4 word window

	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	...

- “digital” and “information” both co-occur with “computer” and “data” frequently
- “cherry” and “strawberry” both co-occur with “pie” and “sugar” frequently
- Word co-occurrence statistics reflect word semantic similarity!
- Issues? Sparsity!

Is Raw Frequency A Good Representation?

- On the one hand, high frequency can imply semantic similarity
- On the other hand, there are words with universally high frequencies

	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

- Can we reweight the raw frequencies so that distinctively high frequency terms are highlighted?



Term Frequency (TF)

- A word appearing 100 times in a document doesn't make it 100 times more likely to be relevant to the meaning of the document
- Instead of using the raw counts, we squash the counts with log scale

$$\text{TF}(w, d) = \begin{cases} 1 + \log_{10} \text{count}(w, d) & \text{count}(w, d) > 0 \\ 0 & \text{otherwise} \end{cases}$$

Document Frequency (DF)

- Motivation: Give a higher weight to words that occur only in a few documents
 - Terms that are limited to a few documents are more discriminative
 - Terms that occur frequently across the entire collection aren't as helpful
- Document frequency (DF): count how many documents a word occurs in

$$\text{DF}(w) = \sum_{i=1}^N \mathbb{1}(w \in d_i) \longrightarrow \begin{array}{l} \text{Evaluates to 1 if } w \text{ occurs in } d_i \\ \text{otherwise evaluates to 0} \end{array}$$

- DF is NOT defined to be the total count of a word across all documents (collection frequency)!

	Collection Frequency	Document Frequency
Romeo	113	1
action	113	31

Inverse Document Frequency (IDF)

- We want to emphasize discriminative words (with low DF)
- Inverse document frequency (IDF): total number of documents (N) divided by DF, in log scale

$$\text{IDF}(w) = \log_{10} \left(\frac{N}{\text{DF}(w)} \right)$$

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

DF & IDF statistics in the
Shakespeare corpus

TF-IDF Weighting

The TF-IDF weighted value characterizes the “salience” of a term in a document

$$\text{TF-IDF}(w, d) = \text{TF}(w, d) \times \text{IDF}(w)$$

TF-IDF weighted

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	0.246	0	0.454	0.520
good	0	0	0	0
fool	0.030	0.033	0.0012	0.0019
wit	0.085	0.081	0.048	0.054

$$\cos(\mathbf{v}_{d_2}, \mathbf{v}_{d_3}) = 0.10 \quad \cos(\mathbf{v}_{d_3}, \mathbf{v}_{d_4}) = 0.99$$

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

$$\cos(\mathbf{v}_{d_2}, \mathbf{v}_{d_3}) = 0.81 \quad \cos(\mathbf{v}_{d_3}, \mathbf{v}_{d_4}) = 0.99$$

How to Define Documents?

- The concrete definition of documents is usually open to different design choices
 - Wikipedia article/page
 - Shakespeare play
 - Book chapter/section
 - Paragraph/sentence
 - ...
- Larger documents provide broader context; smaller ones provide focused insights
- Depends on the analysis need: interested in global trends across documents (e.g., news articles) vs. more local patterns (e.g., specific sections of a legal document)?



Probability-Based Weighting

- TF-IDF weighting scheme is based on heuristics
- Can we weigh the raw counts with probabilistic approaches?
- Intuition: the association between two words can be reflected by **how much they co-occur more than by chance**

		context word				summed counts	
center word		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
summed counts	count(context)	4997	5673	473	512	61	11716

Word Association Based on Probability

- In probability theory, when two random variables A & B are independent, we have
Joint probability $p(A, B) = p(A)p(B)$
- When two words co-occur by chance, we expect their probabilities to satisfy the independence assumption: $p(w_1, w_2) = p(w_1)p(w_2)$
- When $p(w_1, w_2) > p(w_1)p(w_2)$, two words co-occur more often than would be expected by chance
- How to develop a probabilistic metric to characterize this association?

Pointwise Mutual Information (PMI)

- PMI compares the probability of two words co-occurring with the probabilities of the words occurring independently

$$\text{PMI} = \log_2 \frac{p(w_1, w_2)}{p(w_1)p(w_2)} = \log_2 \frac{\#(w_1, w_2) \cdot N}{\#(w_1)\#(w_2)} \quad \text{N: Total word counts}$$

- PMI = 0: Two words co-occur as expected by chance => no particular association
- PMI > 0: Two words co-occur more often than by chance => the higher the PMI, the stronger the association between the words
- PMI < 0: Two words co-occur less often than expected by chance => negative associations; not much actionable insight
- Positive PMI (PPMI): replaces all negative PMI values with zero

$$\text{PPMI} = \max \left(\log_2 \frac{p(w_1, w_2)}{p(w_1)p(w_2)}, 0 \right)$$

PPMI Example

Raw counts

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

PPMI-weighted
matrix

	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

Issue: biased toward infrequent events (rare words tend to have very high PMI values)

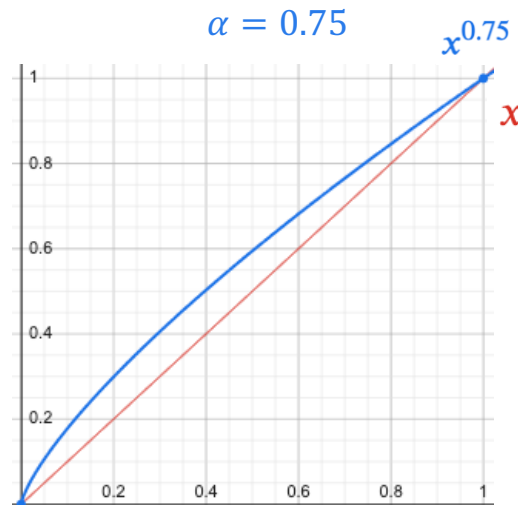
PPMI with Power Smoothing

Power smoothing: Manually boost low probabilities by raising to a power α

$$\text{PPMI} = \max \left(\log_2 \frac{p(w_1, w_2)}{p(w_1)p(w_2)}, 0 \right)$$

Original: $p(w) = \frac{\#(w)}{\sum_{w' \in \mathcal{V}} \#(w')}$

Power smoothed:
($\alpha < 1$) $p_\alpha(w) = \frac{\#(w)^\alpha}{\sum_{w' \in \mathcal{V}} \#(w')^\alpha}$



PPMI with Add- k Smoothing

- Another way of increasing the counts of rare occurrences is to apply add- k smoothing

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

Add a constant k to all counts

- The larger the k (k can be larger than 1), the more we boost the probability of rare occurrences

TF-IDF vs. PMI Weighting

- TF-IDF
 - Measures the importance of a word in a document relative to other documents (corpus)
 - Context granularity: document level
 - Based on heuristics
 - High TF-IDF = frequent in a document but infrequent across the corpus
- PMI:
 - Measures the strength of association between two words
 - Context granularity: word pair level (usually based on local context windows)
 - Based on probability assumptions
 - High PMI = words co-occur more often than expected by chance, a strong association

Summary: Word Semantics & Senses

- Understanding word semantics & senses help us build better language models!
- Word semantics is complex
 - Polysemy: a single word having multiple meanings
 - Multi-faceted: word meanings entail various aspects (e.g., valence, arousal, dominance)
- Many types of word relations: synonyms, antonyms, hyponyms & hypernyms...
- Word relations are usually not binarized (e.g., perfect synonyms are rare); word similarity is usually a more flexible measure

Summary: Classic Word Representations

- Large-scale lexical databases (WordNet) were constructed in early NLP developments
- WordNet consists of manually curated synsets linked by relation edges
- WordNet can be used as a database for word sense disambiguation
- WordNet has significant limitations:
 - Require significant efforts to construct and maintain/update
 - Limited coverage of domain-specific terms & low-resource language
 - Only support individual words and their meanings

Summary: Vector Space Models

- Vector semantic space: use vector representations to reflect word semantics
- Cosine similarity is the most-commonly used metric for vector similarity
- Word-document & word-word co-occurrence statistics provide valuable semantic information – count-based vector representations work decently well
- Raw counts are not good representations (e.g., biased to universally frequent terms)
- TF-IDF highlights the important words in a document relative to other documents
- PMI measures the strength of association between two words based on probabilistic (independence) assumptions



Thank You!

Yu Meng

University of Virginia

yumeng5@virginia.edu