# Representing and Comparing Texts

HSS 510: NLP for HSS

Taegyoon Kim

Mar 13, 2024

# Agenda

Things to be covered

- Unit of analysis

- Tokenization (segmentation)

- Text normalization (= cleaning)

- BoW / vector space models

- Cosine similarity

- TF-IDF weighting

- Guided coding: segmentation, normalization, representation (Python)

# Unit of Analysis

The main element that is being analyzed in a study

- "What" or "who" that is being studied
- Depends on the research question

# Unit of Analysis

Typically, information about one unit is recorded as one row



| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | ID | Q1 | Q2 | Q3 | Q4 | Q5 | Q6 | Q7 | Q8 | Q9 | Q10 |
| 2 | 10010101 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 3 | 10010102 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 4 | 10010103 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 5 | 10010104 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 6 | 10010105 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 7 | 10010106 | Strongly Agree | Strongly Agree | Agree | Strongly Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 8 | 10010107 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 9 | 10010108 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 10 | 10010109 | Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Agree | Strongly Agree | Disagree | Disagree |
| 11 | 10010110 | Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Strongly Agree | Disagree | Disagree |
| 12 | 10010111 | Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Strongly Agree | Disagree | Disagree |
| 13 | 10010112 | Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Strongly Agree | Disagree | Stongly Disagree |
| 14 | 10010113 | Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Strongly Agree | Disagree | Stongly Disagree |
| 15 | 10010114 | Agree | Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Strongly Agree | Disagree | Stongly Disagree |
| 16 | 10010115 | Agree | Strongly Agree | Agree | Disagree | Strongly Disagre | Disagree | Strongly Agree | Strongly Agree | Disagree | Disagree |
| 17 | 10010116 | Agree | Strongly Agree | Agree | Agree | Strongly Disagre | Disagree | Strongly Agree | Strongly Agree | Disagree | Disagree |
| 18 | 10010117 | Agree | Strongly Agree | Agree | Agree | Strongly Disagre | Disagree | Strongly Agree | Strongly Agree | Disagree | Disagree |
| 19 | 10010118 | Agree | Strongly Agree | Agree | Agree | Strongly Disagre | Disagree | Strongly Agree | Strongly Agree | Disagree | Disagree |
| 20 | 10010119 | Strongly Agree | Strongly Agree | Agree | Agree | Strongly Disagre | Disagree | Strongly Agree | Strongly Agree | Disagree | Disagree |
| 21 | 10010120 | Strongly Agree | Strongly Agree | Agree | Agree | Strongly Disagre | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 22 | 10010121 | Strongly Agree | Disagree | Agree | Agree | Strongly Disagre | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 23 | 10010122 | Strongly Agree | Strongly Disagre | Agree | Strongly Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 24 | 10010123 | Strongly Agree | Strongly Disagre | Strongly Agrees | Strongly Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 25 | 10010124 | Disagree | Strongly Disagre | Strongly Agrees | Strongly Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 26 | 10010125 | Disagree | Strongly Agree | Strongly Agrees | Strongly Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 27 | 10010126 | Disagree | Strongly Agree | Agree | Strongly Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 28 | 10010127 | Disagree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Agree |
| 29 | 10010128 | Disagree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 30 | 10010129 | Disagree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 31 | 10010130 | Disagree | Strongly Agree | Agree | Agree | Agree | Disagree | Strongly Agree | Agree | Strongly Disagree | Disagree |
| 32 | 10010131 | Disagree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Agree | Disagree | Disagree |
| 33 | 10010132 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Disagree | Disagree | Disagree |
| 34 | 10010133 | Strongly Agree | Strongly Agree | Agree | Agree | Disagree | Disagree | Strongly Agree | Agree | Disagree | Agree |
| 35 | 10010134 | Strongly Agree | Strongly Agree | Agree | Disagree | Disagree | Disagree | Strongly Agree | Agree | Disagree | Disagree |

# Unit of Analysis

"What are the dominant themes in a corpus of 19th-century British literature?"

- Data: literary works (novels, poems, etc.)
- Unit of analysis: paragraphs, chapters, etc.

# Unit of Analysis

"What are the key scientific topics debated within the scientific community between 2000–2020?"

- Data: scientific publication databases (e.g., Dimension, Web of Science, etc.)
- Unit of analysis: titles, (sentences in) abstracts, paragraphs in full texts, etc.

# Unit of Analysis

The key consideration is our research question

# Unit of Analysis

Barbera et al. (2019)

- Topic models (**L**atent **D**irichlet **A**llocation) on tweets from ordinary users and 500+ legislators in the U.S.

- See if the topics in the former at **t** predicts the latter at **t+1**

# Unit of Analysis

Barbera et al. (2019)

- *"Our definition of "document" is the aggregated total of tweets sent by members of Congress each day"*

- *"Our conceptualization of each day's tweets as the political agenda that each party within each legislative chamber is trying to push for that specific day"*

- *"Conducting an analysis at the tweet level is complex, given its very limited length"*

# Unit of Analysis

## Hammer et al. (2019)

### THREAT: A Large Annotated Corpus for Detection of Violent Threats

1st Hugo L. Hammer
Department of Computer Science
OsloMet – Oslo Metropolitan University
Oslo, Norway
hugo.hammer@oslomet.no

2nd Michael A. Riegler
Simula Metropolitan Center for Digital Engineering
Oslo, Norway

3rd Lilja Øvrelid
4th Erik Velldal
Department of Informatics
University of Oslo
Oslo, Norway

*Abstract*—Understanding, detecting, moderating and in extreme cases deleting hateful comments in online discussions and social media are well-known challenges. In this paper we present a dataset consisting of a total of around 30 000 sentences from around 10 000 YouTube comments. Each sentence is manually annotated as either being a violent threat or not. Violent threats is the most extreme form of hateful communication and is of particular importance from an online radicalization and national security perspective. This is the first publicly available dataset with such an annotation. The dataset can further be useful to develop automatic moderation tools or may even be useful from a social science perspective for analyzing the characteristics of online threats and how hateful discussions evolve.

*Index Terms*—national security, publicly available dataset, social media, threat detection, violent threats

commenting [1], [2], [5], [11], [12], [15], [16], [26]. The methods are mainly based on machine learning and thus require annotated text to learn to separate abominable from harmless online behaviour. Unfortunately, neither of these studies have made the accompanied datasets publicly available. In fact, we are not aware of any publicly available datasets that can be used to develop automatic threat detection.

As a contribution to solve these challenges and to make it possible to perform open and important research on making cyberspace more secure for people we present a large dataset of YouTube comments, where each sentence (manually segmented) is annotated as either being a threat of violence or not.

# Unit of Analysis

## Hammer et al. (2019)

- Supervised learning to detect threatening speech on YouTube comments
- Comments on YouTube videos are split into individual sentences

# Unit of Analysis

## Other considerations

- Rarity of quantity of interest

- Computational power (!?)

- P-value (!?)
    - E.g., legislators' tweets aggregated at the weekly level vs. monthly level

# Models of Text Representation

*"All models are wrong, some are useful" (George Box, 1976)*

# Tokenization

Breaking up a text into discrete words

- Tokenization is a form of segmentation (= word segmentation)
- Token: each individual "word" in the document
  - Possibly including numbers, punctuation, or other symbols
- Tokenization: the process of splitting a document into its constituent words

# Tokenization

"To be or not to be, that is the question"

$\longrightarrow$ "To", "be", "or", "not", "to", "be", "that", "is", "the", "question"

# Tokenization

## Types

- Each token is of a particular "type"

- The set of types is the vocabulary (often denoted as $|V|$)

- "To be or not to be, that is the question"
  $\longrightarrow$ "to" "be" "or" "not" "that" "is" "the", "question" ($|V|$ = 8)

# Tokenization

"Let us learn tokenization."

- Word-level: ["Let", "us", "learn", "tokenization."]

- Subword-level: ["Let", "us", "learn", "token", "ization."]

- Character-level: ["L", "e", "t", "u", "s", "l", "e", "a", "r", "n", "t", "o", "k", "e", "n", "i", "z", "a", "t", "i", "o", "n", "."]

# Tokenization

Why word-level?

- Words: most common for many downstream analyses
  - Word embeddings, topic models, etc.
- Subwords
  - Now prevalent in neural NLP
  - Handling of OOV (out-of-vocabulary) words
  - More efficiency (consider "tokenization")
  - E.g., Byte Pair Encoding (BPE), WordPiece
- Character: no meaning (although computationally very efficient)
- Sentences: too many types

# Tokenization

Subword tokenization in recent GPTs

# Tokenization

## How to tokenize?

- In English (and many other languages, including Korea), we can rely heavily on white space
  - Many algorithms build not only on white space but also on various patterns
  - E.g., appstrophies: ["don", "'", "t"] vs. ["do", "n't"]
  - E.g., punctuations: ['vehicle?'] vs. ['vehicle', '?']
- Tools include NLTK, spaCy, Keras, etc.
- In some languages, words cannot be separated deterministically, and they need models (e.g., Chinese, Japanese, etc.)

# Tokenization

*n*-grams

- A sequence of *n* adjacent tokens

- Unigrams, bigrams, trigrams, etc.

- Why would we need multi-grams?
  - E.g., "White House", "look after", "take care of", etc.

# Tokenization

*n*-grams

- Be aware of the computational cost
  - Consider the number of all consecutive sets of two words in the corpus
- Alternatively, we can compile a list of particular bi-grams or tri-grams

# Segmenting Sentences/paragraphs

## Sentence segmentation

- Useful cues: periods, question marks, or exclamation marks
- Prone to errors (the example of `"."`)
  - Abbreviations and initials: "Ph.D.", "J.K. Rowling", etc.
  - Decimal numbers: "3.14"
  - Websites and email addresses: "www.kaist.ac.kr") and email addresses
  - Quotations within a sentence: "He said, 'Stop.' Then he left."
- Rule-based/deterministic or ML-based approaches (part of `nltk` and `spaCy`)

# Segmenting Sentences/paragraphs

## Paragraph segmentation

- Not as commonly addressed
- Few specialized libraries or algorithms in Python
- Useful cues: newline characters (`\n`) or double newline characters (`\n\n`)

# Text Normalization

A set of approaches to reducing complexity in text

- The output from tokenization will contain too many words
- Putting words in a standard form can be useful for information retrieval
  - Findings a pattern in a corpus (e.g., Penny, Pennies, penny, pennies, etc.)

# Text Normalization

We will discuss five approaches

- Lowercasing
- Removing punctuation
- Removing stop words
- Lemmatization/stemming
- Filtering by frequency

# Text Normalization

## Lowercasing

- We often replace all capital letter with lowercase letters
- It is assumed that there is no (semantic) difference
- Is it?

# Text Normalization

## Lowercasing

- Compare "NOW" and "now" in terms of sentiment
- Capital letters also signal the start of of a sentence
- Proper nouns (May vs. may. US vs. us)

# Text Normalization

Removing punctuation

- Period ( **.** ), comma ( **,** ), apostrophe ( **'** ), quotation ( **""** ), question ( **?** ), exclamation ( **!** ), dash ( **−** ), ellipsis ( **...** ), colon ( **:** ), semicolon ( **;** ), etc.
- In many cases, these are (considered) unimportant
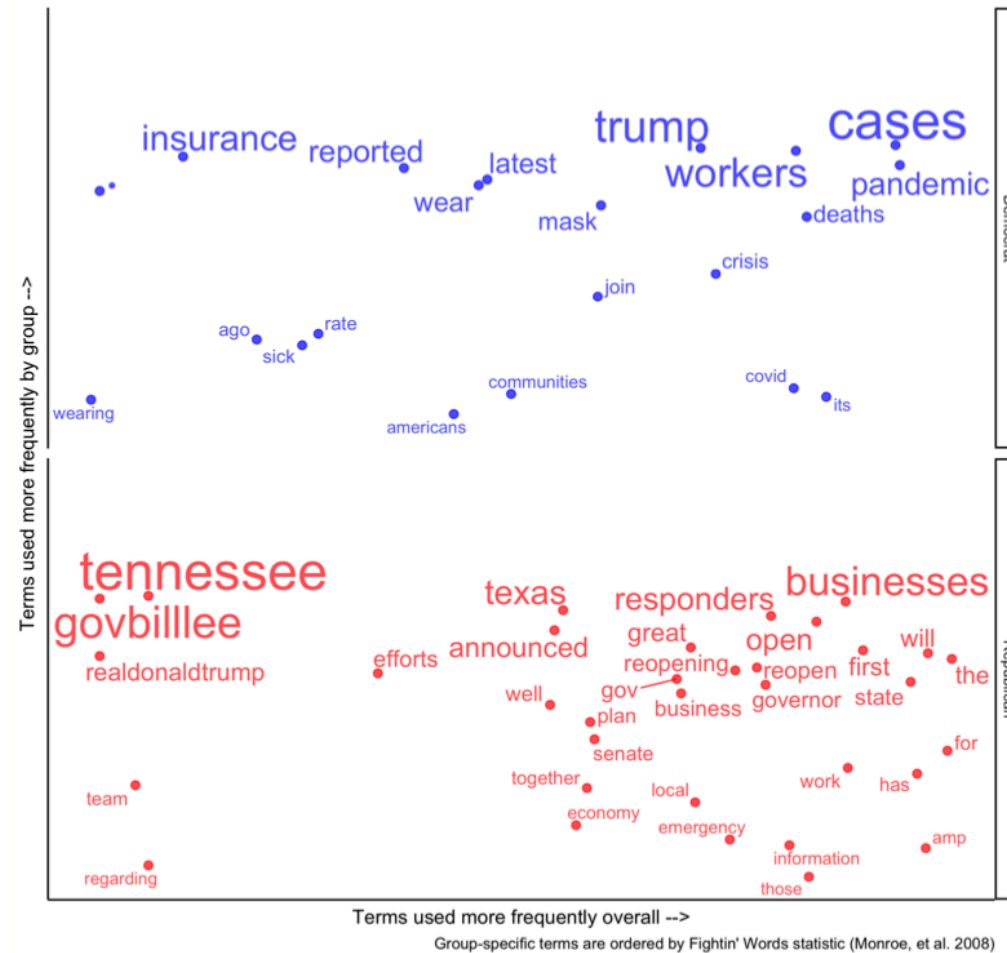- Are they?

# Text Normalization

## Removing punctuation

- Punctuation carries important information
    - Exclamation mark (`!!!`), hashtags (`#metoo`), emojis (`<3`, `:)`, `^^`, `-_-;`), etc.
- Punctuation itself can be of interest (studying writing styles)

# Text Normalization



Comparison of Terms by Groups

Terms used more frequently by group -->

Terms used more frequently overall -->

Group-specific terms are ordered by Fightin' Words statistic (Monroe, et al. 2008)

# Text Normalization

Removing stop words

- Common words used across documents that do not give much information
- E.g., "and", "the", or "that"

# Text Normalization

Removing stop words can spare much computational power

- C.f., Heaps' Law
- However, under what circumstances are these words *not* stop words?

# Text Normalization

## Lemmatization

- Lemma: the base form
  - E.g., "run"
- Wordform: various forms derived the lemma
  - E.g., "runs", "ran", "running"
- Lemmatizatoin is the process of mapping words to their lemma

# Text Normalization

## Lemmatization

- Not always straightforward
  - Irregular variations E.g., "see-saw-seen"
  - Same token but different lemmas
    - E.g., he is "writing" an email vs. a nice piece of "writing"
- Necessitates a dictionary and POS (**p**art **o**f **s**peech) tagging

# Text Normalization

Stemming is a popular approximation to lemmatization

- Simply discards the end of a word
    - E.g., family: famili
- Errors
    - E.g., "leav" for both "leaves" (as in "He leaves the room") and "leaves" (as in parts of a plant)
- Various algorithms: *Porter*, *Lancaster*, etc.

# Text Normalization

## Filtering by frequency

- Too (in)frequent words across documents
  - E.g., stop words
- The rationale
  - Discriminatory power
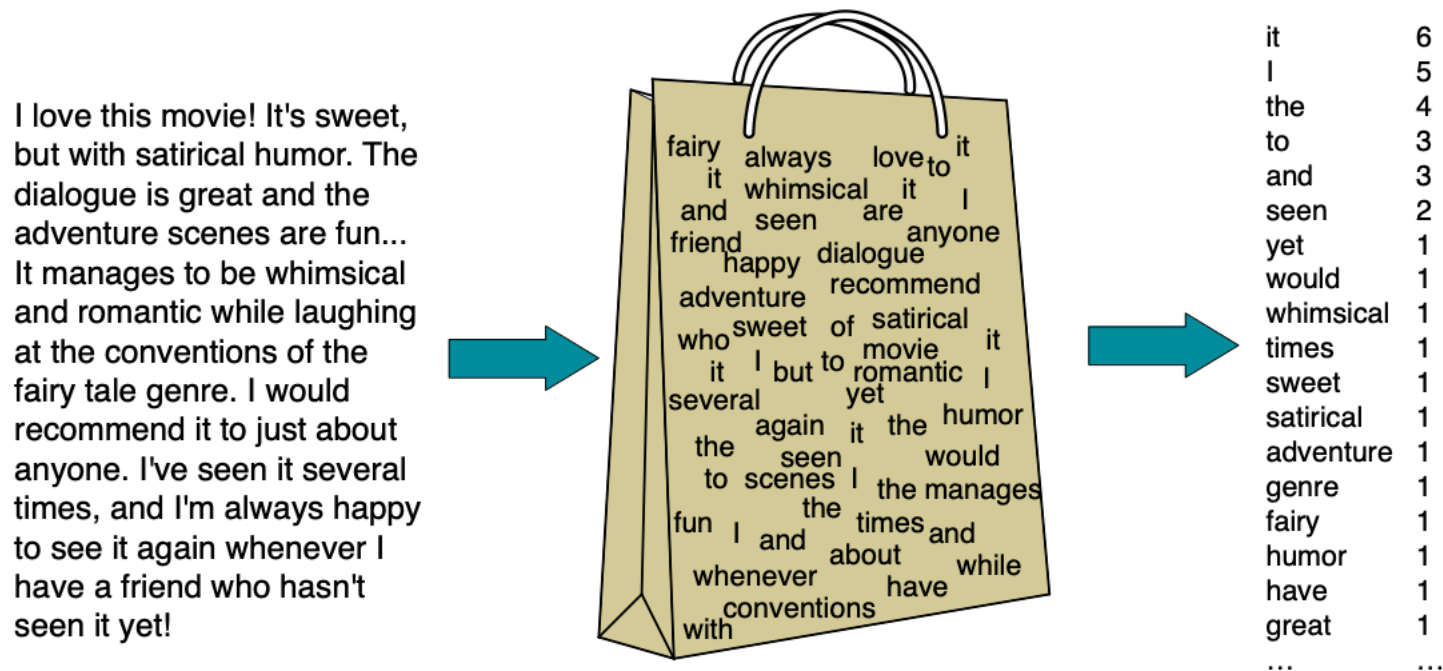  - Computational savings

# (How) Should We Normalize?

Difficult to know its consequences *a priori*

- Before analysis: carefully think about the pros and cos in each of the steps
- After analysis: conduct robustness check

# The Bag of Words (BoW) model

The most common text representation model

- A text is represented as a set of words that appear in it

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

fairy always love to it
it whimsical it I
and seen are
friend happy dialogue anyone
adventure recommend
who sweet of satirical it
it I but to romantic I
several yet
again it the humor
the seen would
to scenes I the manages
the times
fun I and and
whenever about while
have
conventions
with

| | |
|---|---|
| it | 6 |
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| … | … |

# The Bag of Words (BoW) model

Document-Feature Matrix (or Document-Term Matrix)

- Columns record features/terms (all types or $|V|$)
- Rows record documents
- Cells can be binary vectors or count vectors

# The Bag of Words (BoW) model

An example corpus

- Doc 1: "The clever fox cleverly jumps over the lazy dog, showcasing its cleverness."

- Doc 2: "Magic and mysteries mingle in the wizard's daily musings, revealing mysteries unknown."

- Doc 3: "Sunny days bring sunshine and sunsets, making sunny parks the best for sunny strolls."

# The Bag of Words (BoW) model

An example DFM

| Document | clever | jumps | lazy | dog | magic | mysteries | |
|----------|--------|-------|------|-----|-------|-----------|---|
| Doc 1 | 3 | 1 | 1 | 1 | 0 | 0 | . |
| Doc 2 | 0 | 0 | 0 | 0 | 1 | 2 | . |
| Doc 3 | 0 | 0 | 0 | 0 | 0 | 0 | . |

# The Bag of Words (BoW) model

An example DFM

| Document | clever | jumps | lazy | dog | magic | mysteries |
|----------|--------|-------|------|-----|-------|-----------|
| Doc 1 | 3 | 1 | 1 | 1 | 0 | 0 |
| Doc 2 | 0 | 0 | 0 | 0 | 1 | 2 |
| Doc 3 | 0 | 0 | 0 | 0 | 0 | 0 |

# The Bag of Words (BoW) model

An example corpus

- Doc 1: "The `clever` fox `cleverly` jumps over the lazy dog, showcasing its `cleverness`."

- Doc 2: "Magic and mysteries mingle in the wizard's daily musings, revealing mysteries unknown."

- Doc 3: "Sunny days bring sunshine and sunsets, making sunny parks the best for sunny strolls."

# The Vector Space Model

What is the vector space model?

- Each row (representing a text) in a DFM is a vector (an array of numbers) in a high-dimensional space
- The size of the dimension (the number of columns) is $|V|$
- Originates from IR (**I**nformation **R**etrieval)
  - See Turney and Pantel (2010) for details

# Comparing Texts

With some form of DFM, we are ready to compare different documents

- "Similar" can mean different things
  - Sentiments, stances, themes, etc.
- There is no "correct" notion of similarity
- Yet there are metrics that are more of less effective across contexts

## Cosine Similarity

We have two vectors (representing two documents), $\vec{A}$ and $\vec{B}$:

$$\vec{A} = [a_1, a_2, \ldots, a_n]$$

$$\vec{B} = [b_1, b_2, \ldots, b_n]$$

The inner product:

$$\vec{A} \cdot \vec{B} = (a_1 \times b_1) + (a_2 \times b_2) + \ldots + (a_n \times b_n)$$

## Cosine Similarity

Cosine similarity between vectors $\vec{A}$ and $\vec{B}$ is given by:

$$\text{Cosine Similarity}(\vec{A}, \vec{B}) = \frac{\vec{A} \cdot \vec{B}}{\|\vec{A}\| \|\vec{B}\|}$$

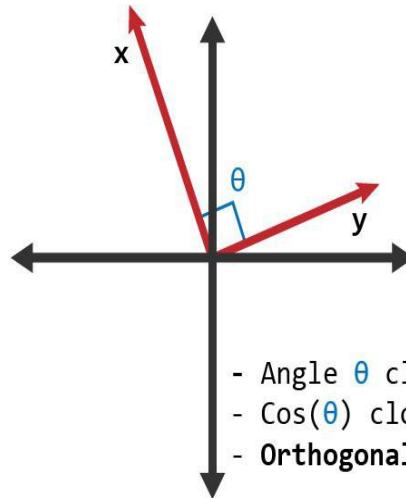$\vec{A} \cdot \vec{B}$ is the inner product, and $\|\vec{A}\|$ and $\|\vec{B}\|$ are defined as

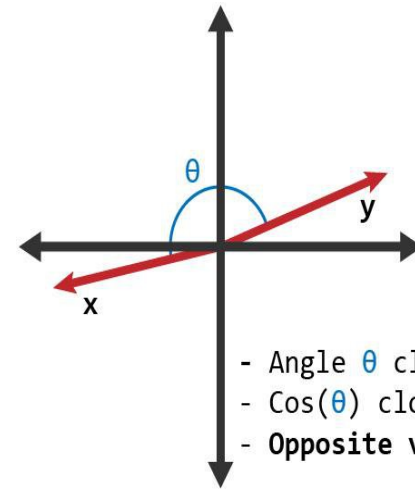$$\|\vec{A}\| = \sqrt{a_1^2 + a_2^2 + \ldots + a_n^2} \quad \|\vec{B}\| = \sqrt{b_1^2 + b_2^2 + \ldots + b_n^2}$$

# Cosine Similarity



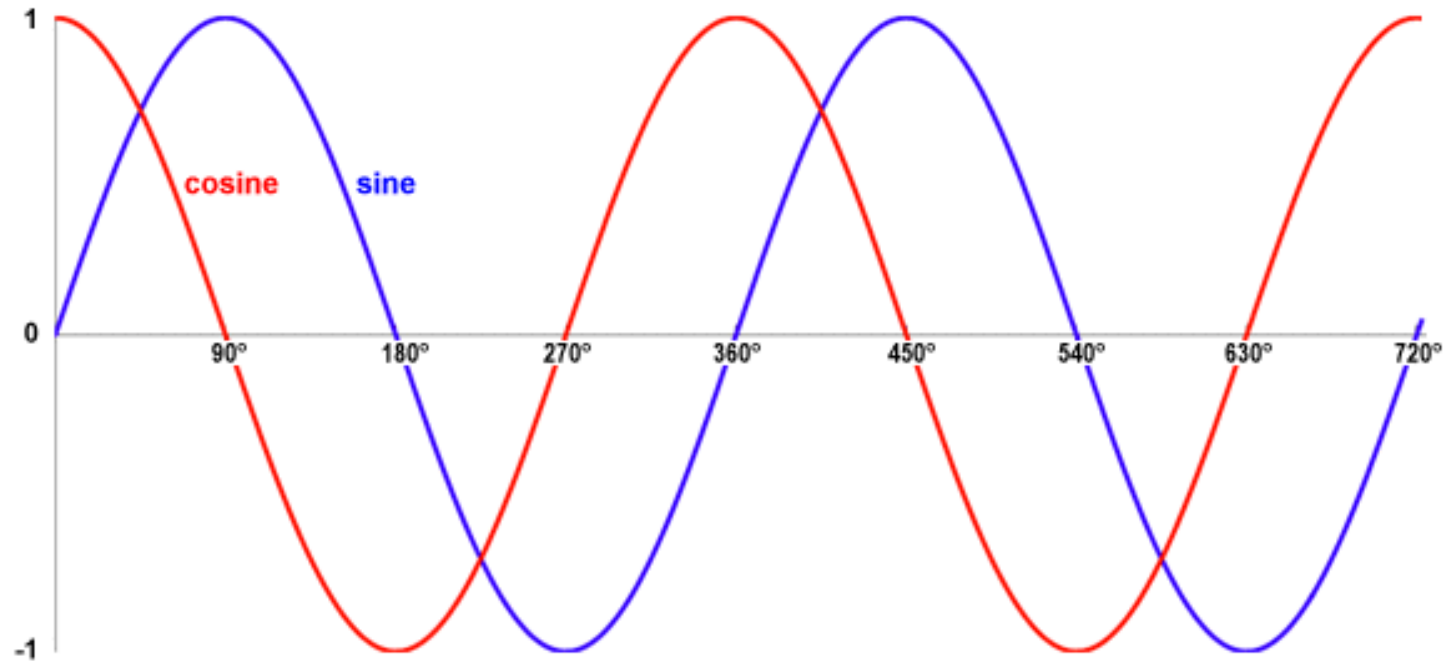- Angle θ close to 0
- Cos(θ) close to 1
- **Similar vectors**

- Angle θ close to 90
- Cos(θ) close to 0
- **Orthogonal vectors**

- Angle θ close to 180
- Cos(θ) close to -1
- **Opposite vectors**

# Cosine Similarity

# Cosine Similarity

## Is Cosine-Similarity of Embeddings Really About Similarity?

Harald Steck
hsteck@netflix.com
Netflix Inc.
Los Gatos, CA, USA

Chaitanya Ekanadham
cekanadham@netflix.com
Netflix Inc.
Los Angeles, CA, USA

Nathan Kallus
nkallus@netflix.com
Netflix Inc. & Cornell University
New York, NY, USA

March 11, 2024

# TF-IDF Weighting

TF (Term Frequency) - IDF (Inverse Document Frequency)

- Count vectors consider the frequencies of words
- However, some words are too frequent across different documents
  - E.g., *the*, *a*, *an*, etc.
- We want to weight how unique a word to a document

# TF-IDF Weighting

TF-IDF is a numerical statistic that reflects *how important a word is to a document* in a corpus.

# TF-IDF Weighting

The **TF-IDF** value is obtained by multiplying **TF** (**T**erm **F**requency) and **IDF** (**I**nverse **D**ocument **F**requency) for a term in a document, highlighting the importance of rare terms

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

# TF-IDF Weighting

## Term Frequency

- Reflects how frequently a term occurs in a document, normalized by the document length

$$\text{TF}(t, d) = \frac{\text{Number of times term } t \text{ appears in document } d}{\text{Total number of terms in document } d}$$

# TF-IDF Weighting

## Inverse Document Frequency

- Scales down terms that occur very frequently across the corpus and are less informative

$$\text{IDF}(t, D) = \log\left(\frac{\text{Total number of documents } D}{\text{Number of documents with term } t \text{ in it} + 1}\right)$$

# TF-IDF Weighting

Many versions of TF-IDF: link

# TF-IDF Weighting

# Count Vectors Vs. TF-IDF Vectors

## Count Vectors

| Term | can | you | fly | sleep |
|---|---|---|---|---|
| 'can you fly' | 1 | 1 | 1 | 0 |
| 'can you sleep' | 1 | 1 | 0 | 1 |

## TF-IDF Vectors

| Term | can | you | fly | sleep |
|---|---|---|---|---|
| 'can you fly' | 0.5 | 0.5 | 0.7 | 0 |
| 'can you sleep' | 0.5 | 0.5 | 0 | 0.7 |

# Summary

The process of transforming raw texts into numbers involve a number of important decisions

- Segmentation
- Normalization
- Representation

→ It is worth thinking ahead of and reviewing the potential consequences

# Guided Coding

Normalization, representation, and comparison in Python ([Link](Link))