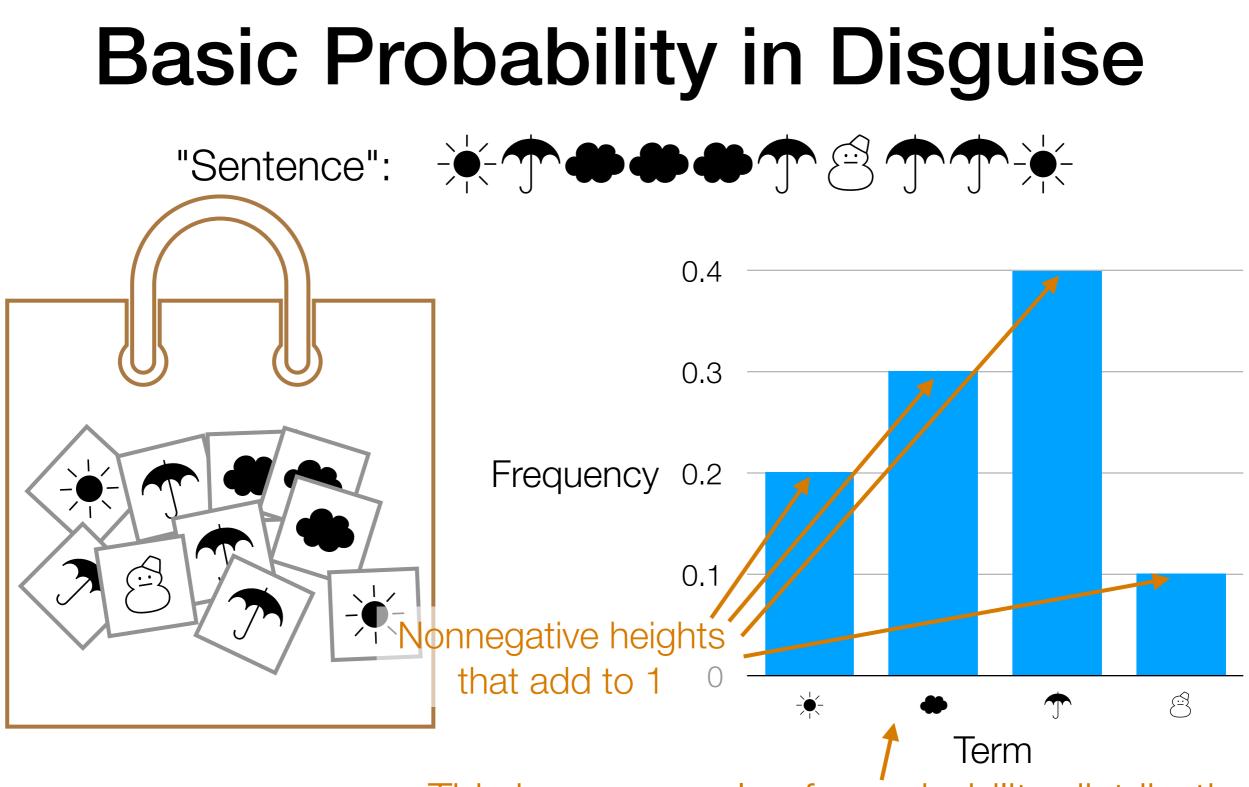Carnegie Mellon University

# HeinzCollege

# 94-775/95-865
# Unstructured Data Analytics
# Lecture 2: Basic Text Analysis
# Wrap-up, Co-occurrence Analysis

George Chen

# Basic Probability in Disguise

"Sentence":
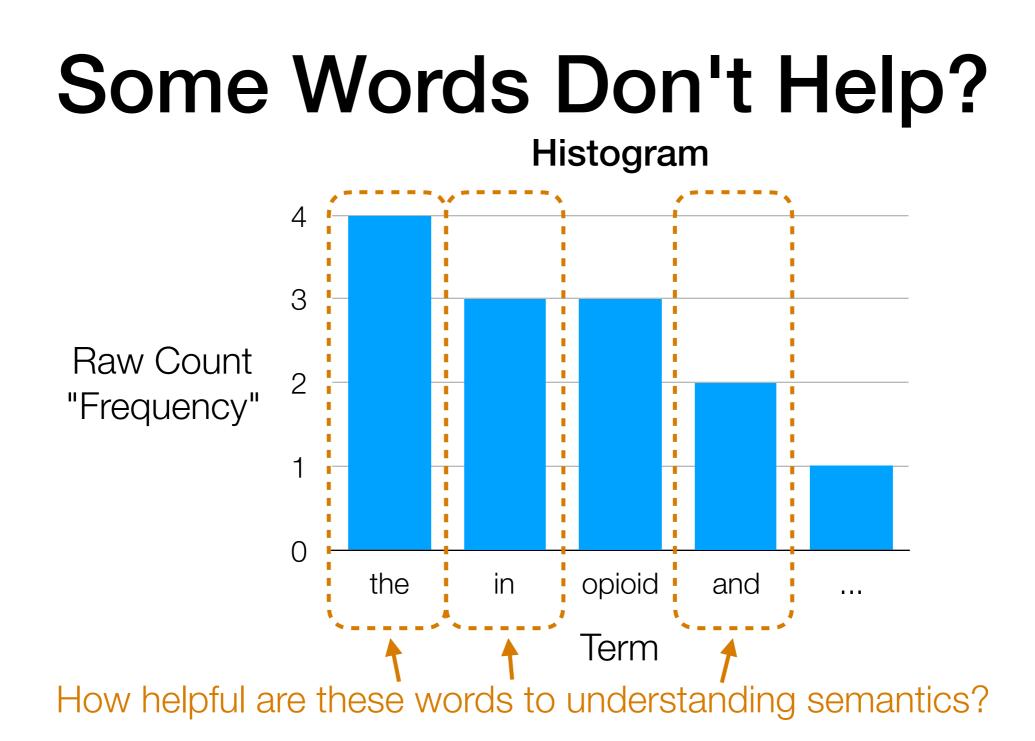
Frequency

Nonnegative heights that add to 1

Term

This is an example of a probability distribution

Probability distributions will appear throughout the course and are a **key component** to the success of many modern AI methods

# Now let's take advantage of properties of text

In other words: natural language humans use has a lot of *structure* that we can exploit

# Some Words Don't Help?

**Histogram**



*Raw Count "Frequency"* (y-axis: 0, 1, 2, 3, 4)

Terms: the (4), in (3), opioid (3), and (2), ... (1)

*Term* (x-axis)

How helpful are these words to understanding semantics?

Bag-of-words models: many frequently occurring words unhelpful

We can remove these words first (remove them from the "bag")
➔ words that are removed are called **stopwords**

*(determined by removing most frequent words or using curated stopword lists)*
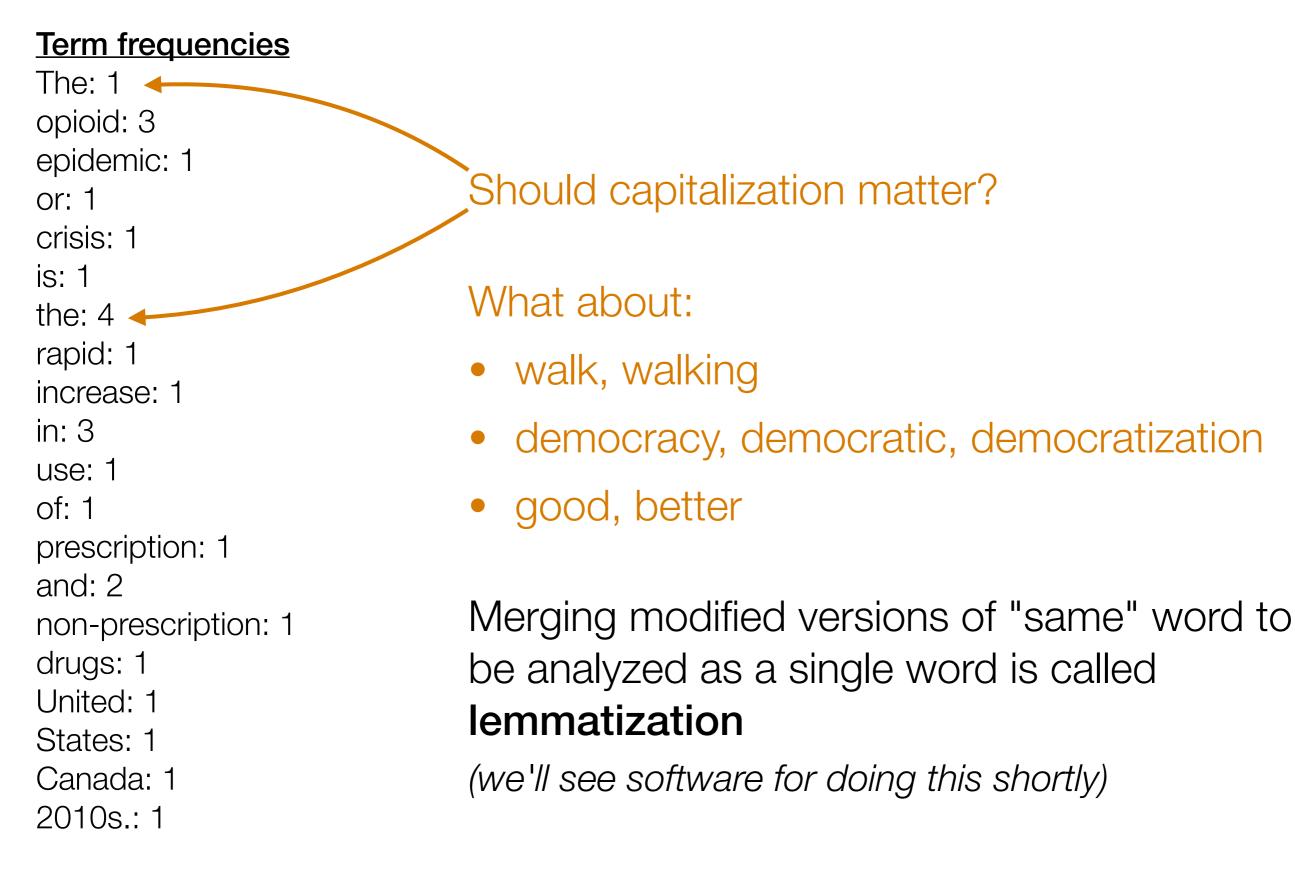
# Example Stopword List (from spaCy)

'a', 'about', 'above', 'across', 'after', 'afterwards', 'again', 'against', 'all', 'almost', 'alone', 'along', 'already', 'also', 'although', 'always', 'am', 'among', 'amongst', 'amount', 'an', 'and', 'another', 'any', 'anyhow', 'anyone', 'anything', 'anyway', 'anywhere', 'are', 'around', 'as', 'at', 'back', 'be', 'became', 'because', 'become', 'becomes', 'becoming', 'been', 'before', 'beforehand', 'behind', 'being', 'below', 'beside', 'besides', 'between', 'beyond', 'both', 'bottom', 'but', 'by', 'ca', 'call', 'can', 'cannot', 'could', 'did', 'do', 'does', 'doing', 'done', 'down', 'due', 'during', 'each', 'eight', 'either', 'eleven', 'else', 'elsewhere', 'empty', 'enough', 'etc', 'even', 'ever', 'every', 'everyone', 'everything', 'everywhere', 'except', 'few', 'fifteen', 'fifty', 'first', 'five', 'for', 'former', 'formerly', 'forty', 'four', 'from', 'front', 'full', 'further', 'get', 'give', 'go', 'had', 'has', 'have', 'he', 'hence', 'her', 'here', 'hereafter', 'hereby', 'herein', 'hereupon', 'hers', 'herself', 'him', 'himself', 'his', 'how', 'however', 'hundred', 'i', 'if', 'in', 'inc', 'indeed', 'into', 'is', 'it', 'its', 'itself', 'just', 'keep', 'last', 'latter', 'latterly', 'least', 'less', 'made', 'make', 'many', 'may', 'me', 'meanwhile', 'might', 'mine', 'more', 'moreover', 'most', 'mostly', 'move', 'much', 'must', 'my', 'myself', 'name', 'namely', 'neither', 'never', 'nevertheless', 'next', 'nine', 'no', 'nobody', 'none', 'noone', 'nor', 'not', 'nothing', 'now', 'nowhere', 'of', 'off', 'often', 'on', 'once', 'one', 'only', 'onto', 'or', 'other', 'others', 'otherwise', 'our', 'ours', 'ourselves', 'out', 'over', 'own', 'part', 'per', 'perhaps', 'please', 'put', 'quite', 'rather', 're', 'really', 'regarding', 'same', 'say', 'see', 'seem', 'seemed', 'seeming', 'seems', 'serious', 'several', 'she', 'should', 'show', 'side', 'since', 'six', 'sixty', 'so', 'some', 'somehow', 'someone', 'something', 'sometime', 'sometimes', 'somewhere', 'still', 'such', 'take', 'ten', 'than', 'that', 'the', 'their', 'them', 'themselves', 'then', 'thence', 'there', 'thereafter', 'thereby', 'therefore', 'therein', 'thereupon', 'these', 'they', 'third', 'this', 'those', 'though', 'three', 'through', 'throughout', 'thru', 'thus', 'to', 'together', 'too', 'top', 'toward', 'towards', 'twelve', 'twenty', 'two', 'under', 'unless', 'until', 'up', 'upon', 'us', 'used', 'using', 'various', 'very', 'via', 'was', 'we', 'well', 'were', 'what', 'whatever', 'when', 'whence', 'whenever', 'where', 'whereafter', 'whereas', 'whereby', 'wherein', 'whereupon', 'wherever', 'whether', 'which', 'while', 'whither', 'who', 'whoever', 'whole', 'whom', 'whose', 'why', 'will', 'with', 'within', 'without', 'would', 'yet', 'you', 'your', 'yours', 'yourself', 'yourselves'

# Is removing stop words always a good thing?

"To be or not to be"

# Some Words Mean the Same Thing?

**Term frequencies**

The: 1
opioid: 3
epidemic: 1
or: 1
crisis: 1
is: 1
the: 4
rapid: 1
increase: 1
in: 3
use: 1
of: 1
prescription: 1
and: 2
non-prescription: 1
drugs: 1
United: 1
States: 1
Canada: 1
2010s.: 1

Should capitalization matter?

What about:

- walk, walking

- democracy, democratic, democratization

- good, better

Merging modified versions of "same" word to be analyzed as a single word is called **lemmatization**

*(we'll see software for doing this shortly)*

# What about a word that has multiple meanings?

Challenging: try to split up word into multiple words depending on meaning (requires inferring meaning from context)

This problem is called **word sense disambiguation** (WSD)
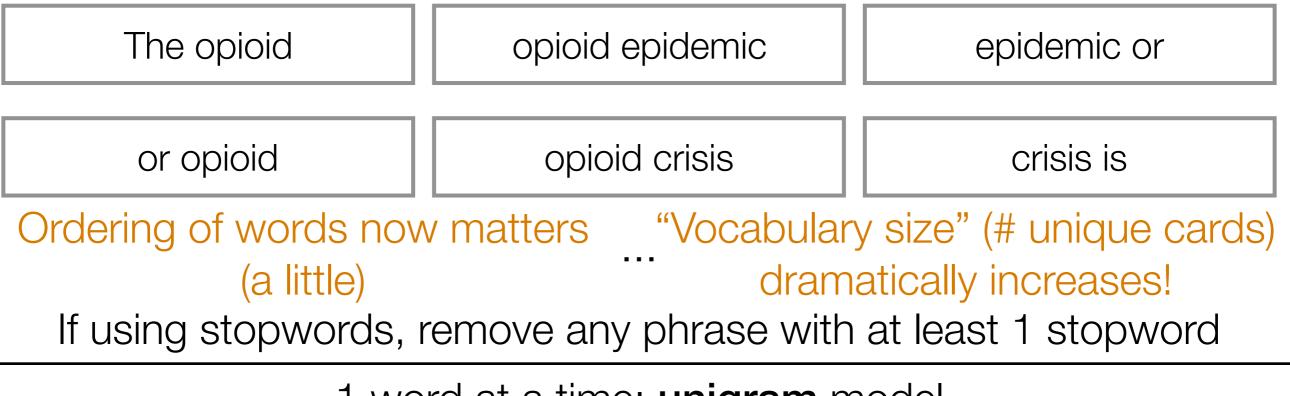
# Treat Some Phrases as a Single Word?

**Term frequencies**

The: 1
opioid: 3
epidemic: 1
or: 1
crisis: 1
is: 1
the: 4
rapid: 1
increase: 1
in: 3
use: 1
of: 1
prescription: 1
and: 2
non-prescription: 1
drugs: 1
United: 1
States: 1
Canada: 1
2010s.: 1

First need to detect what are "named entities": called **named entity recognition**
*(we'll see software for doing this shortly)*

Treat as single 2-word phrase "United States"?

# Some Other Basic NLP Tasks

- **Tokenization:** figuring out what are the atomic "words" (including how to treat punctuation)

- **Part-of-speech tagging:** figuring out what are nouns, verbs, adjectives, etc

- **Sentence recognition:** figuring out when sentences actually end rather than there being some acronym with periods in it, etc

# Bigram Model

The opioid epidemic or opioid crisis is the rapid increase in the use of prescription and non-prescription opioid drugs in the United States and Canada in the 2010s.

| The opioid | opioid epidemic | epidemic or |
|---|---|---|
| or opioid | opioid crisis | crisis is |

Ordering of words now matters (a little)  …  "Vocabulary size" (# unique cards) dramatically increases!

If using stopwords, remove any phrase with at least 1 stopword

---

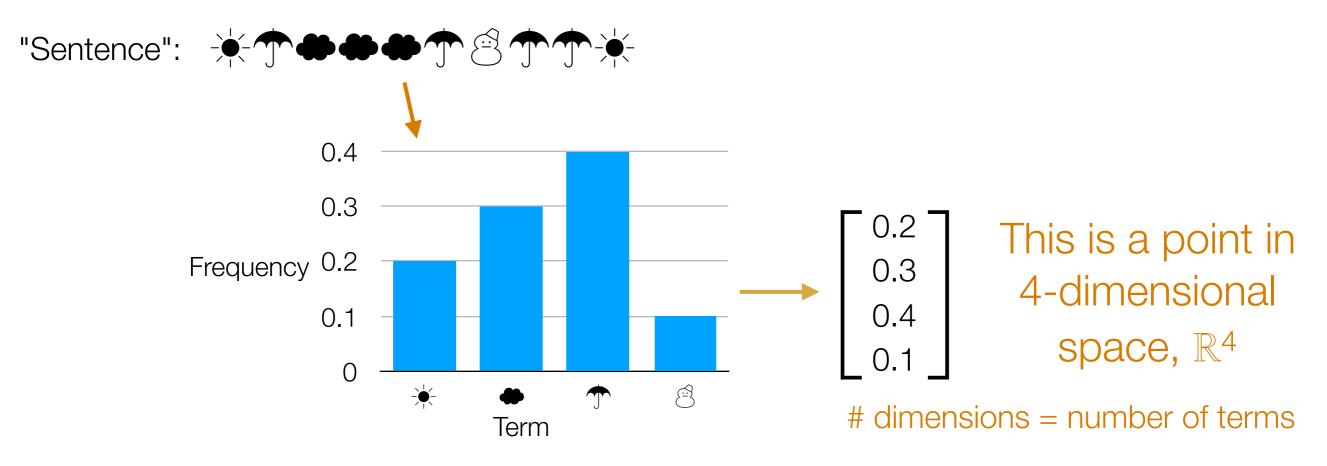1 word at a time: **unigram** model

2 words at a time: **bigram** model

*n* words at a time: ***n*-gram** model

# The spaCy **Python Package**

Demo

# Recap: Basic Text Analysis

- Represent text in terms of "features"
  (such as how often each word/phrase appears)

  - Can repeat this for different documents:
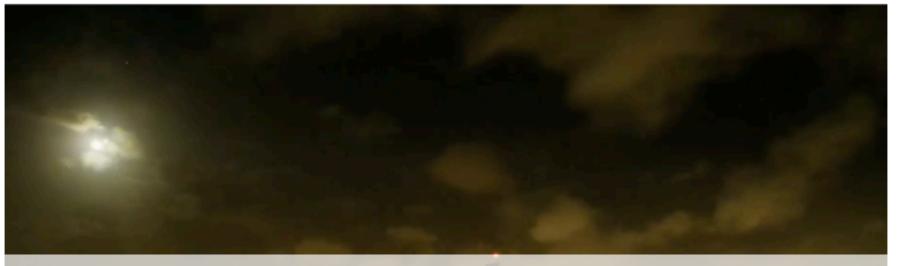    *represent each document as a "feature vector"*

"Sentence":  ☀️🌂☁️☁️☁️🌂☃️🌂🌂☀️

$$\begin{bmatrix} 0.2 \\ 0.3 \\ 0.4 \\ 0.1 \end{bmatrix}$$

This is a point in 4-dimensional space, $\mathbb{R}^4$

# dimensions = number of terms

In general (not just text): first represent data as feature vectors

# Finding Possibly Related Entities

**Elon Musk's Tesla Powerwalls Have Landed in Puerto Rico**

The solar batteries have reportedly been spotted in San Juan's airport.

By John Patrick Pullen  October 16, 2017

Exactly one week after Tesla CEO Elon Musk suggested his company could help with Puerto Rico's electricity crisis in the aftermath of Hurricane Maria, more of the company's Powerwall battery packs have arrived on the island, according to a photo snapped at San Juan airport Friday, Oct. 13.

# How to automatically figure out Elon Musk and Tesla are related?

Source: http://fortune.com/2017/10/16/elon-musks-tesla-powerwalls-have-landed-in-puerto-rico/

# Co-Occurrences

For example: count # news articles that have different named entities co-occur

|  | Apple | Facebook | Tesla |
|---|---|---|---|
| Elon Musk | 10 | 15 | 300 |
| Mark Zuckerberg | 500 | 10000 | 500 |
| Tim Cook | 200 | 30 | 10 |

Big values ➔ *possibly* related named entities

# Different Ways to Count

- Just saw: for all doc's, count # of doc's in which two named entities co-occur

  - This approach ignores # of co-occurrences *within a specific document* (e.g., if 1 doc has "Elon Musk" and "Tesla" appear 10 times, we count this as 1)

  - Could instead add # co-occurrences, not just whether it happened in a doc

- Instead of looking at # doc's, look at co-occurrences within a *sentence*, or a *paragraph*, etc

---

### Bottom Line

- There are many ways to count co-occurrences
- You should think about what makes the most sense/is reasonable for the problem you're looking at

# Co-Occurrences

For example: count # news articles that have different named entities co-occur

|  | Apple | Facebook | Tesla |
|---|---|---|---|
| Elon Musk | 10 | 15 | 300 |
| Mark Zuckerberg | 500 | 10000 | 500 |
| Tim Cook | 200 | 30 | 10 |

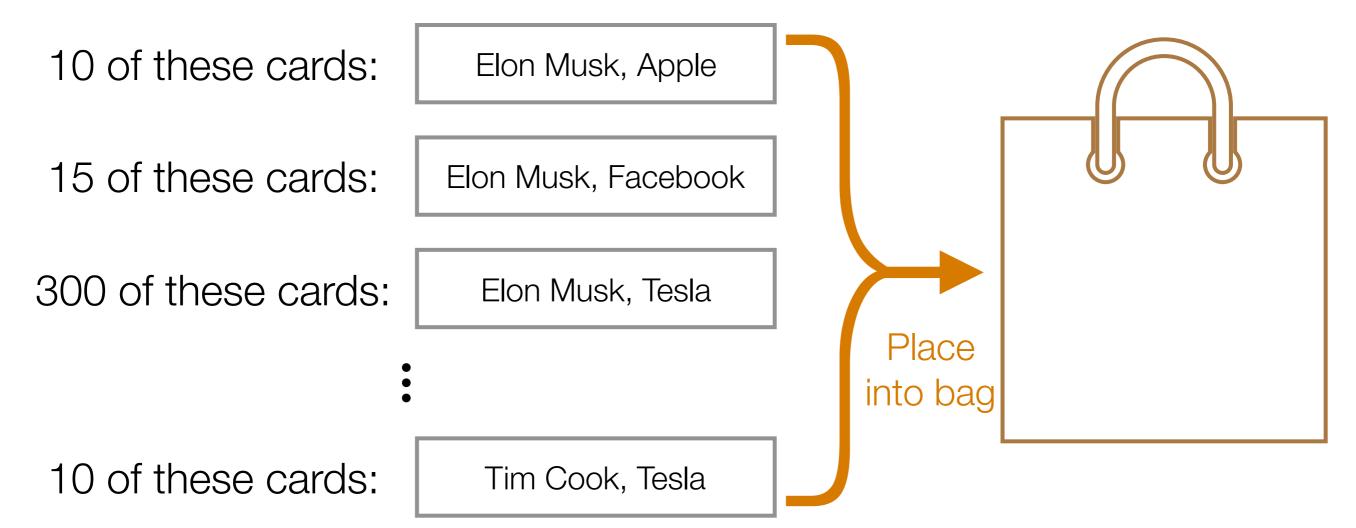Big values ➜ *possibly* related named entities

How to downweight "Mark Zuckerberg" if there are just way more articles that mention him?

Key idea: what would happen if people and companies had nothing to do with each other?

| | Apple | Facebook | Tesla |
|---|---|---|---|
| Elon Musk | 10 | 15 | 300 |
| Mark Zuckerberg | 500 | 10000 | 500 |
| Tim Cook | 200 | 30 | 10 |

Probability of drawing "Elon Musk, Apple"?

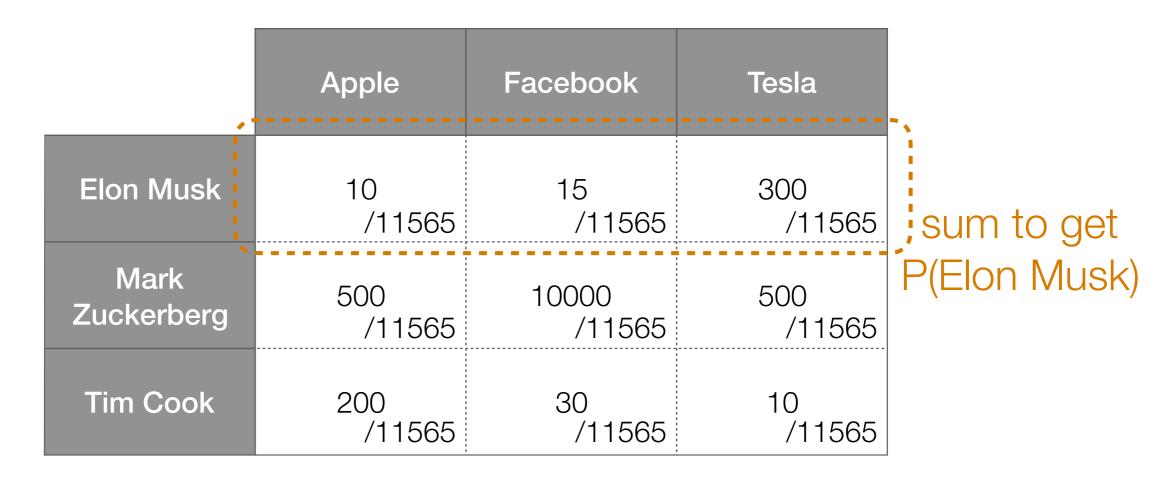Probability of drawing a card that says "Apple" on it?

10 of these cards: Elon Musk, Apple

15 of these cards: Elon Musk, Facebook

300 of these cards: Elon Musk, Tesla

⋮

10 of these cards: Tim Cook, Tesla

Place into bag

# Co-occurrence table

|  | Apple | Facebook | Tesla |
|---|---|---|---|
| Elon Musk | 10 | 15 | 300 |
| Mark Zuckerberg | 500 | 10000 | 500 |
| Tim Cook | 200 | 30 | 10 |

Total: 11565

# Joint probability table

| | Apple | Facebook | Tesla |
|---|---|---|---|
| **Elon Musk** | 10 /11565 | 15 /11565 | 300 /11565 |
| **Mark Zuckerberg** | 500 /11565 | 10000 /11565 | 500 /11565 |
| **Tim Cook** | 200 /11565 | 30 /11565 | 10 /11565 |

sum to get P(Elon Musk)

Total: 11565

# Joint probability table

|  | Apple | Facebook | Tesla |  |
|---|---|---|---|---|
| Elon Musk | 0.00086 | 0.00130 | 0.02594 | **0.02810** |
| Mark Zuckerberg | 0.04323 | 0.86468 | 0.04323 | **0.95115** |
| Tim Cook | 0.01729 | 0.00259 | 0.00086 | **0.02075** |
|  | **0.06139** | **0.86857** | **0.07004** |  |

Recall: if events A and B are independent, P(A, B) = P(A)P(B)

# Joint probability table **if people and companies were independent**

| | Apple | Facebook | Tesla | |
|---|---|---|---|---|
| **Elon Musk** | 0.00173 | 0.02441 | 0.00197 | **0.02810** |
| **Mark Zuckerberg** | 0.05839 | 0.82614 | 0.06662 | **0.95115** |
| **Tim Cook** | 0.00127 | 0.01802 | 0.00145 | **0.02075** |
| | **0.06139** | **0.86857** | **0.07004** | |

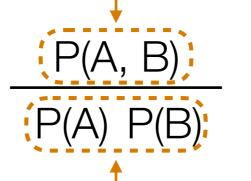Recall: if events A and B are independent, P(A, B) = P(A)P(B)

## What we actually observe

|  | Apple | Facebook | Tesla |
|---|---|---|---|
| Elon Musk | 0.00086 | 0.00130 | 0.02594 |
| Mark Zuckerberg | 0.04323 | 0.86468 | 0.04323 |
| Tim Cook | 0.01729 | 0.00259 | 0.00086 |

## What should be the case if people are companies are independent

|  | Apple | Facebook | Tesla |
|---|---|---|---|
| Elon Musk | 0.00173 | 0.02441 | 0.00197 |
| Mark Zuckerberg | 0.05839 | 0.82614 | 0.06662 |
| Tim Cook | 0.00127 | 0.01802 | 0.00145 |

# Pointwise Mutual Information (PMI)

Probability of A and B co-occurring

$$\frac{P(A,\ B)}{P(A)\ P(B)}$$

if equal to 1
➔ A, B are indep.

Probability of A and B co-occurring *if they were independent*

**PMI(A, B) is defined as the log of the above ratio**

PMI measures (the log of) a ratio that says how
far A and B are from being independent

# Looking at All Pairs of Outcomes

- PMI measures how P(A, B) differs from P(A)P(B) using a **log ratio**

- **Log ratio** isn't the only way to compare!

- Another way to compare:

$$\frac{[\, P(A, B) - P(A)\, P(B)\, ]^2}{P(A)\, P(B)}$$

Phi-square = $\displaystyle\sum_{A,\, B} \frac{[\, P(A, B) - P(A)\, P(B)\, ]^2}{P(A)\, P(B)}$

Chi-square = N × Phi-square

N = sum of all co-occurrence counts

Phi-square is between 0 and 1

0 ➜ pairs are all indep.

Measures how close *all* pairs of outcomes are close to being indep.

# PMI/Phi-Square/Chi-Square Calculation

Demo