

# **Introduction to Data Science**

## **Text Analysis**

Kerstin Bunte and **George Azzopardi**

WMCS16002, semester Ia 2018



Mathematician Ted Kaczynski!

# Learning goals

- Understand what text mining is
- Learn how to preprocess text
- Learn how to compare documents
- Compute the bag of words model
- Identify model topics
- Build bi-grams, . . . n-grams

# Search

A large, empty search input field with a microphone icon in the top right corner.

Google Search

I'm Feeling Lucky

Google.nl offered in: [Nederlands](#)

# Extract Knowledge



In 2011, Watson competed against the legendary human competitors and outperformed them! Watson had access to 200 million pages of structured and unstructured content consuming four terabytes of disk storage including the full text of Wikipedia

# Classify



You can thank text classification algorithms for keeping spam out of your inbox

# Detect Emotions



SEARCHED TERM

taxes

POSITIVE TWEETS

151

NEUTRAL TWEETS

1272

NEGATIVE TWEETS

77

TOTAL TWEETS

1500

### 10.07% POSITIVE



speaking of coase and taxes,  
wonder if prof helland is on twitter?  
law and econ = great class.  
learned abt coase and torts and  
injury ([view](#))

at least i am sitting in the sun  
preparing for taxes - that's the  
good news ([view](#))

@johnchow if you have a friend  
into import/export.. you can you  
just need to pay for the import  
taxes and stuff.. but still cheap!  
:d ([view](#))



o\_o apparently, majority of those  
losign jobs are in the private

### 84.80% NEUTRAL



South Fayette to hold the line on  
property taxes: South Fayette  
school administrators last night  
proposed a \$31...  
<http://tinyurl.com/c6vnms> ([view](#))



<http://cliqz.com/us.headlines/c/9221.html> : Congress Looking  
at New Taxes on AIG Bonuses  
([view](#))



South Fayette to hold the line on

### 5.13% NEGATIVE

taxes. pissed that our unofficial  
sxsw spot fell through. jealous of  
those of you going maybe next  
year. feeling the need to paint  
something ([view](#))

my speck package came  
today.... i hate that i have to pay  
30% taxes at the dutyoffice >  
([view](#))

yeah, i know aig execs are  
 greedy, but there's something  
 wrong when govt's create new  
taxes just to hound a handful of  
people. ([view](#))

@codeape did you iust do your

# Trend of publications

Go to Web of Knowledge

# Let us start small...

How do we extract features from text?

## Let us start small...

How do we extract features from text?

1. Tokenization
2. Case folding
3. Removal of short and stop words
4. Stemming
5. Lemmatization

# **Tokenization**

## **Tokenization**

the isolation of word-like units from a text

- base for all further text processing
- accuracy of tokenisation affects higher level processing

# Tokenization

## Tokenization

the isolation of word-like units from a text

- base for all further text processing
- accuracy of tokenisation affects higher level processing

ATTN: “garbage in, garbage out”

# Tokenization

## Tokenization

the isolation of word-like units from a text

- base for all further text processing
- accuracy of tokenisation affects higher level processing  
ATTN: “garbage in, garbage out”
- tokenizing by: spaces, spaces and punctuation, etc.

# **Tokenization is More Difficult than it seems**

# Tokenization is More Difficult than it seems

Ambiguous examples

- state-of-the-art
- Résumé vs. resume

# Tokenization is More Difficult than it seems

Ambiguous examples

- state-of-the-art
- Résumé vs. resume
- computer science, San Francisco

this is called a *noun phrase*

# Tokenization Beyond English

- German/Dutch combine individual words to create longer compound words

searching for Hund should retrieve Windhund while searching for Blau should likely not return Blauzungekrankheit

# Tokenization Beyond English

- German/Dutch combine individual words to create longer compound words

searching for Hund should retrieve Windhund while searching for Blau should likely not return Blauzungekrankheit
- Some Asian languages have no spaces

# Tokenization Beyond English

- German/Dutch combine individual words to create longer compound words

searching for Hund should retrieve Windhund while searching for Blau should likely not return Blauzungekrankheit
- Some Asian languages have no spaces
- Some languages have very complex morphology: e.g.  
Uygarlastiramadiklarimizdanmissinizcasina (tr)

# Tokenization Best Practices

- No single tokenizer that will work in every scenario.

# Tokenization Best Practices

- No single tokenizer that will work in every scenario.
- Common to pair a standard tokenization technique with a lookup table
  - e.g. San Francisco

# Case folding

## Case folding

= reducing all the letters in a word to lowercase

Works

- “**Look**, I took the liberty to **look** at that parrot when I got it home...”

# Case folding

## Case folding

= reducing all the letters in a word to lowercase

Works

- “**Look**, I took the liberty to **look** at that parrot when I got it home...”

Does not work

- General Motors → general motors

# Case folding

## Case folding

= reducing all the letters in a word to lowercase

### Works

- “**Look**, I took the liberty to **look** at that parrot when I got it home...”

### Does not work

- General Motors → general motors
- The Who → the who

Heuristic strategies (e.g. no case fold in mid-sentence)

## **Removal of short and stop words**

- Typically, remove words that have fewer than 3 characters

## Removal of short and stop words

- Typically, remove words that have fewer than 3 characters
- Remove the most common words in a language, known as **stop words**  
E.g. the, and, which, where, ...

# Stemming

## Stemming

= reducing a word to a stem (or root) form *algorithmically*

- Number: stemmer, stemmers
- Tense: stems, stemmed, stemming
- Gender: stemmer, *stemmeress*
- ...

# Algorithmic stemmer: Porter

Rules for chopping the end of some words<sup>1</sup>

Step 1a

SSES	->	SS	caresses	->	caress
IES	->	I	ponies	->	poni
			ties	->	ti
SS	->	SS	caress	->	caress
S	->		cats	->	cat

Step 1b

(m>0) EED	->	EE	feed	->	feed
			agreed	->	agree
(*v*) ED	->		plastered	->	plaster
			bled	->	bled
(*v*) ING	->		motoring	->	motor
			sing	->	sing

Link: The -ing rule Link: The full description of the rules

<sup>1</sup>A good discussion at: <http://snowballstem.org/texts/introduction.html> and a place to try them online: <http://textanalysisonline.com/>

## Porter Stemmer<sup>2</sup> – Example

Call me Ishmael . Some **year** ago **have littl** or no money in my **purs** , and **noth** particular to interest me on shore , I thought I would sail about a **littl** and see the **wateri** part of the world .

vs.

Call me Ishmael. Some years ago having little or no money in my purse, and nothing particular to interest me on shore, I thought I would sail about a little and see the watery part of the world.

---

<sup>2</sup>The rules of the stemmer:

<http://snowball.tartarus.org/algorithms/porter/stemmer.html>

## Porter Stemmer<sup>2</sup> – Example

Call me Ishmael . Some **year** ago **have littl** or no money in my **purs** , and **noth** particular to interest me on shore , I thought I would sail about a **littl** and see the **wateri** part of the world .

vs.

Call me Ishmael. Some years ago having little or no money in my purse, and nothing particular to interest me on shore, I thought I would sail about a little and see the watery part of the world.

Go to [http://localhost:8888/notebooks/Tokeniser\\_Stemming.ipynb](http://localhost:8888/notebooks/Tokeniser_Stemming.ipynb)

---

<sup>2</sup>The rules of the stemmer:

<http://snowball.tartarus.org/algorithms/porter/stemmer.html>

# Limiting Porter Damage...

```
new_word = PorterStemmer(old_word)
dUS=Enchanted_Dicitonary(en-US)
if new_word != old_word and len(old_word)>3:
    if new_word exist in dUS:
        Replacing the old_word by new_word
```

See: PyEnchant

# Limiting Porter Damage...

```
new_word = PorterStemmer(old_word)
dUS=Enchanted_Dicitonary(en-US)
if new_word != old_word and len(old_word)>3:
    if new_word exist in dUS:
        Replacing the old_word by new_word
```

See: PyEnchant

Surely, Porter was born before Twitter / SMS:

**This is soooooo LOOOVELY, sooo gooood, noooo waaay!!!**

# Another Headache... HTML

"this is <span class='important'>really important</span> text"

A **chunker** is a special tokenizer function that breaks text up into large chunks rather than individual tokens.

```
>>> from enchant.tokenize import get_tokenizer, HTMLChunker
>>>
>>> tknzs = get_tokenizer("en_US")
>>> [w for w in tknzs("this is <span class='important'>really important</span> text")]
[('this', 0), ('is', 5), ('span', 9), ('class', 14), ('important', 21), ('really', 32), ('important', 39),
>>>
>>>
>>> tknzs = get_tokenizer("en_US", chunkers=(HTMLChunker,))
>>> [w for w in tknzs("this is <span class='important'>really important</span> text")]
[('this', 0), ('is', 5), ('really', 32), ('important', 39), ('text', 56)]
```

# Lemmatization

## Lemmatization

= obtaining the single word that allows you to group together a bunch of inflected forms

- Harder than stemming
- Takes context into account
- Uses some kind of dictionary

# Lemmatization

## Lemmatization

= obtaining the single word that allows you to group together a bunch of inflected forms

- Harder than stemming
- Takes context into account
- Uses some kind of dictionary

Pattern = web mining module for Python

<http://www.clips.ua.ac.be/pages/pattern>

## WordNet Lemmatization Example

Call me Ishmael . Some **year** ago **have** little or no money in my purse , and nothing particular to interest me on shore , I **think** I would sail about a little and see the watery part of the world .

vs.

Call me Ishmael. Some years ago having little or no money in my purse, and nothing particular to interest me on shore, I thought I would sail about a little and see the watery part of the world.

# Stemming vs. Lemmatization

- **Common goal:** To reduce inflectional forms and sometimes derivationally related forms of a word to a common base form

# Stemming vs. Lemmatization

- **Common goal:** To reduce inflectional forms and sometimes derivationally related forms of a word to a common base form
- Stemming
  - crude heuristic process that chops off the end of words
  - correct most of the time
  - simple to implement and fast

# Stemming vs. Lemmatization

- **Common goal:** To reduce inflectional forms and sometimes derivationally related forms of a word to a common base form
- Stemming
  - crude heuristic process that chops off the end of words
  - correct most of the time
  - simple to implement and fast
- Lemmatization
  - Use vocabulary but slower than stemming
  - Morphological analysis of words
  - Return the base or dictionary form of a word

# Stemming vs. Lemmatization

- **Common goal:** To reduce inflectional forms and sometimes derivationally related forms of a word to a common base form
- Stemming
  - crude heuristic process that chops off the end of words
  - correct most of the time
  - simple to implement and fast
- Lemmatization
  - Use vocabulary but slower than stemming
  - Morphological analysis of words
  - Return the base or dictionary form of a word
- Examples
  - "better" -> lemma: "good". It is missed from stemming
  - "walking" -> lemma: "walk", stem: "walk"
  - "meeting" -> two meanings. Lemmatization takes context into account

# NLP Word Features

1. Root words
2. Part of speech (POS) tags  
Ex: POSTagging
3. Named entities

# NLP Word Features

1. Root words
2. Part of speech (POS) tags  
Ex: POSTagging
3. Named entities

Traditionally these features were manually annotated in **language corpora**

# Language Corpora

## Corpus

= a public and frozen set of (usually annotated) documents

- Nice for training and benchmarking
- Often NLP libraries come with corpora

Ex: NLTKCorpora

## Word Feature: Information Content

- IC = metric to denote the importance of a term in a corpus

## Word Feature: Information Content

- IC = metric to denote the importance of a term in a corpus
- Combines
  - 1. knowledge of its hierarchical structure from an ontology
  - 2. statistics on its actual usage in text derived from a corpus

# Word Feature: Information Content

- IC = metric to denote the importance of a term in a corpus
- Combines
  - 1. knowledge of its hierarchical structure from an ontology
  - 2. statistics on its actual usage in text derived from a corpus
- e.g.  $IC(\text{necklace}) > IC(\text{jewelry})$  in an English corpus  
jewelry is more general and thus less interesting

# Word Feature: Information Content

- IC = metric to denote the importance of a term in a corpus
- Combines
  - 1. knowledge of its hierarchical structure from an ontology
  - 2. statistics on its actual usage in text derived from a corpus
- e.g.  $IC(\text{necklace}) > IC(\text{jewelry})$  in an English corpus  
jewelry is more general and thus less interesting

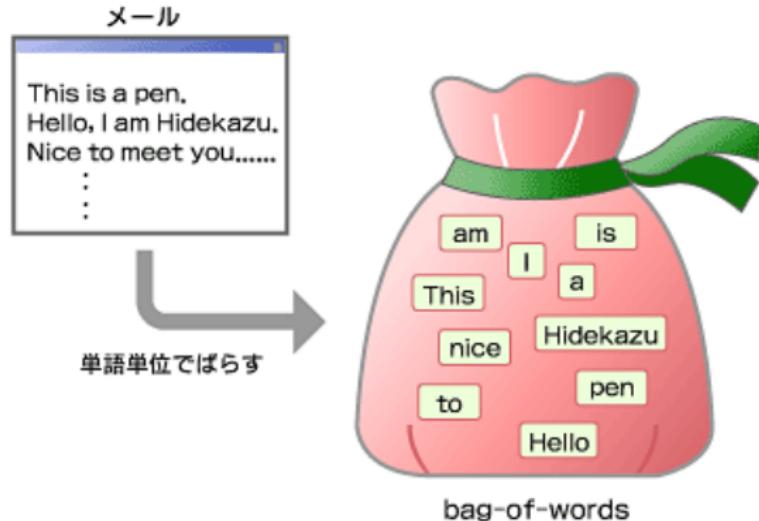
## Problems

- terms not in the corpus automatically receive zero IC  
analyzing U2 with the Brown corpus

Enough about words!

## Documents

# Document Representation: Bag-of-words



- A set of terms ignoring: order, context, etc.
- Each word is considered a term or token

## Clear Limitations but still Powerful

“the dog bites the man” vs. “the man bites the dog”!

{*bites, dog, man, the*}

## Clear Limitations but still Powerful

“the dog bites the man” vs. “the man bites the dog”!

{*bites, dog, man, the*}

“the overwhelming evidence is that the **judicious use of single-term identifiers is preferable to the incorporation of more complex entities extracted from the texts themselves...**” [SB88]

# Term Frequency

- The weight of each term in a document

Proportional to the number of occurrences of the term in that document

# Term Frequency

- The weight of each term in a document

Proportional to the number of occurrences of the term in that document

$$TF_1(term, doc) = \frac{\text{count of term in doc}}{\text{count words in doc}}$$

# Term Frequency

- The weight of each term in a document

Proportional to the number of occurrences of the term in that document

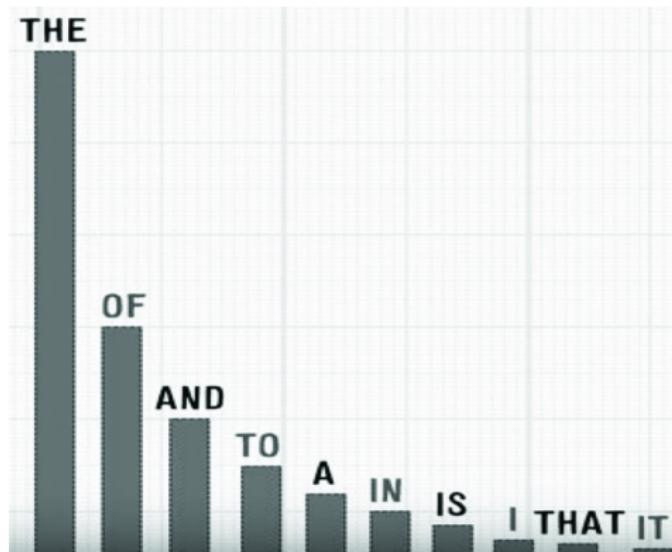
$$TF_1(term, doc) = \frac{\text{count of term in doc}}{\text{count words in doc}}$$

- Roughly follows Zipf's Law

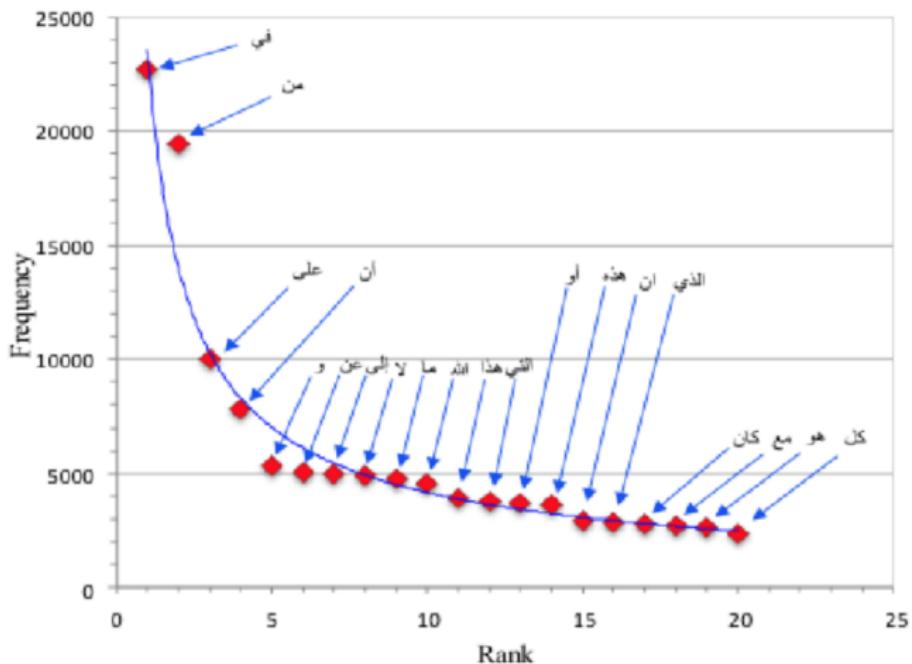
the i-th most common item has freq of approx.  $1/i$  of most frequent – see next slides!

- this is why we often work with  $TF = \log(TF_1)$

# Zipf's Law - Top 10 English Words

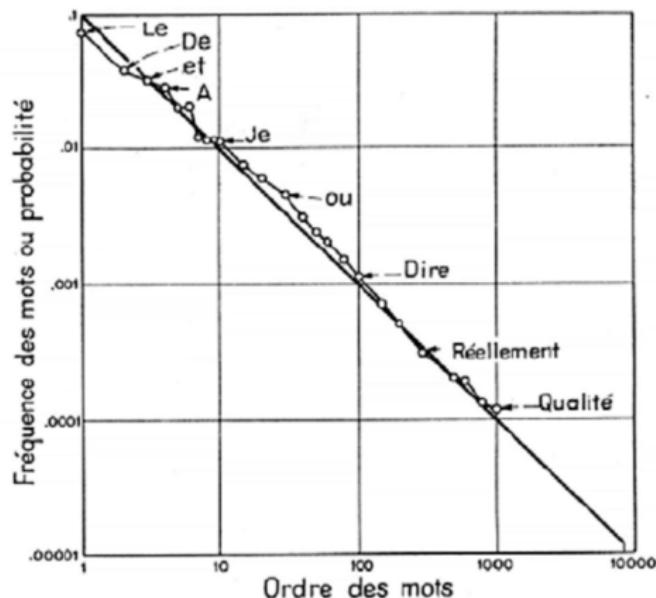


# Zipf's Law - Top 20 Arabic Words



[source]

## Zipf's Law - Top 10K French



"a good approximation to what is demonstrably a much more complicated distribution of word frequencies" (S. Piantadosi)

# Demo

TF in Shakespeare  
(Ex: TFShake)

# **Term Frequency Inverse Document Frequency (TFxIDF)**

## **TF.IDF**

= numerical statistic reflecting the importance of a word in a doc in a corpus

Considering a **word important for a document** if it appears

- often within that document (TF), and

# **Term Frequency Inverse Document Frequency (TFxIDF)**

## **TF.IDF**

= numerical statistic reflecting the importance of a word in a doc in a corpus

Considering a **word important for a document** if it appears

- often within that document (TF), and
- rarely in other documents of the corpus (IDF)

# Computing TF.IDF

$$N = |docs|$$

# Computing TF.IDF

$$N = |docs|$$

$DF(term, docs) = \langle \text{count docs in which term occurs} \rangle$

# Computing TF.IDF

$$N = |docs|$$

$DF(term, docs)$  = <count docs in which term occurs>

$$IDF(term, docs) = \log\left(\frac{N + 1}{DF(term, docs) + 1}\right) + 1$$

# Computing TF.IDF

$$N = |docs|$$

$DF(term, docs)$  = <count docs in which term occurs>

$$IDF(term, docs) = \log\left(\frac{N + 1}{DF(term, docs) + 1}\right) + 1$$

$$TF.IDF(term, doc, docs) = TF(term, doc) * IDF(term, docs)$$

# Demo

Simple TF.IDF Example

Ex: TF.IDF.Example

# Vector Space Models<sup>3</sup>

## Vector Space Model

= the representation of a set of documents as vectors in a common space

---

<sup>3</sup>parts of some following slides are based on slides by Malvina Nissim

# Vector Space Model

- Each term  $t$  of the dictionary is considered as a *dimension*
- A document  $d$  can be represented by the weight of each vocabulary term:  $w(term, doc)$ :

$$\mathbf{v}(d) = (w(t_1, d), w(t_2, d), \dots, w(t_n, d))$$

- Note that also raw frequency could be used (that's what we are doing)

# Vector Space Model

- Each term  $t$  of the dictionary is considered as a *dimension*
- A document  $d$  can be represented by the weight of each vocabulary term:  $w(term, doc)$ :

$$\mathbf{v}(d) = (w(t_1, d), w(t_2, d), \dots, w(t_n, d))$$

- Note that also raw frequency could be used (that's what we are doing)

representation of three documents using term raw frequencies:

$d_1, d_2, d_3$

	$d_1$	$d_2$	$d_3$
affection	115	58	20
jealous	10	7	11
gossip	2	0	6

# Vector Space Model

- Each term  $t$  of the dictionary is considered as a *dimension*
- A document  $d$  can be represented by the weight of each vocabulary term:  $w(term, doc)$ :

$$\mathbf{v}(d) = (w(t_1, d), w(t_2, d), \dots, w(t_n, d))$$

- Note that also raw frequency could be used (that's what we are doing)

representation of three documents using term raw frequencies:  
 $d_1, d_2, d_3$

	$d_1$	$d_2$	$d_3$
affection	115	58	20
jealous	10	7	11
gossip	2	0	6

- **Question:** how do we compute the **similarity between documents?**

# Vector normalization and similarity

- Similarity between vectors  
→ inner product  $\mathbf{V}(d_1) \cdot \mathbf{V}(d_2)$

# Vector normalization and similarity

- Similarity between vectors  
→ inner product  $\mathbf{V}(d_1) \cdot \mathbf{V}(d_2)$
- What about the length of a vector?  
Longer documents will be represented with longer vectors, but that does not mean they are more important
- Euclidian normalization (vector length normalization):

$$\mathbf{v}(d) = \frac{\mathbf{V}(d)}{\|\mathbf{V}(d)\|} \quad \text{where } \|\mathbf{V}(d)\| = \sqrt{\sum_{i=1}^n x_i^2}$$

# Vector normalization and similarity

- Similarity between vectors  
→ inner product  $\mathbf{v}(d_1) \cdot \mathbf{v}(d_2)$
- Similarity given by the *cosine* measure between normalized vectors:

$$sim(d_1, d_2) = \mathbf{v}(d_1) \cdot \mathbf{v}(d_2)$$

## Example (Manning et al<sup>4</sup>)

- $\text{sim}(d1, d2) = \mathbf{v}(d1) \cdot \mathbf{v}(d2)$

let's backtrack this:

- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$

- normalising for length:

$$\mathbf{v}(d_i) = \frac{\mathbf{v}(d_i)}{\|\mathbf{v}(d)\|}$$

- Euclidean length:

$$\|\mathbf{v}(d)\| = \sqrt{\sum_{i=1}^n \mathbf{v}_i^2(d)}$$

---

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

- $sim(d1, d2) = \mathbf{v}(d1) \cdot \mathbf{v}(d2)$

- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$

- normalising for length:

$$\mathbf{v}(d_i) = \frac{\mathbf{V}(d_i)}{\|\mathbf{V}(d)\|}$$

vocabulary	d1	d2	d3
1: affection	115	58	20
2: jealous	10	7	11
3: gossip	2	0	6

- Euclidean length:

$$\|\mathbf{V}(d)\| = \sqrt{\sum_{i=1}^n \mathbf{V}_i^2(d)}$$

---

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

- $sim(d1, d2) = \mathbf{v}(d1) \cdot \mathbf{v}(d2)$

- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$

- normalising for length:

$$\mathbf{v}(d_i) = \frac{\mathbf{V}(d_i)}{\|\mathbf{V}(d)\|}$$

- Euclidean length:

vocabulary	$d1$	$d2$	$d3$
1: affection	115	58	20
2: jealous	10	7	11
3: gossip	2	0	6

$$\|\mathbf{V}(d1)\| = \sqrt{115^2 + 10^2 + 2^2}$$

$$\|\mathbf{V}(d)\| = \sqrt{\sum_{i=1}^n \mathbf{V}_i^2(d)}$$

---

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

- $\text{sim}(d1, d2) = \mathbf{v}(d1) \cdot \mathbf{v}(d2)$
- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$
- normalising for length:

$$\mathbf{v}(d_i) = \frac{\mathbf{V}(d_i)}{\|\mathbf{V}(d)\|}$$

- Euclidean length:

$$\|\mathbf{V}(d)\| = \sqrt{\sum_{i=1}^n \mathbf{V}_i^2(d)}$$

vocabulary	d1	d2	d3
1: affection	115	58	20
2: jealous	10	7	11
3: gossip	2	0	6

$$\|\mathbf{V}(d1)\| = \sqrt{115^2 + 10^2 + 2^2}$$

$$\mathbf{v}(d1_1) = \frac{115}{\sqrt{115^2+10^2+2^2}} = 0.996$$

$$\mathbf{v}(d1_2) = \frac{10}{\sqrt{115^2+10^2+2^2}} = 0.087$$

$$\mathbf{v}(d1_3) = \frac{2}{\sqrt{115^2+10^2+2^2}} = 0.017$$

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

- $\text{sim}(d1, d2) = \mathbf{v}(d1) \cdot \mathbf{v}(d2)$
- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$
- normalising for length:

$$\mathbf{v}(d_i) = \frac{\mathbf{V}(d_i)}{\|\mathbf{V}(d)\|}$$

vocabulary	d1	d2	d3
1: affection	115	58	20
2: jealous	10	7	11
3: gossip	2	0	6

$$\|\mathbf{V}(d2)\| = \sqrt{58^2 + 7^2 + 0}$$

- Euclidean length:

$$\|\mathbf{V}(d)\| = \sqrt{\sum_{i=1}^n \mathbf{V}_i^2(d)}$$

$$\mathbf{v}(d2_1) = \frac{58}{\sqrt{58^2+7^2+0}} = 0.993$$

$$\mathbf{v}(d2_2) = \frac{7}{\sqrt{58^2+7^2+0}} = 0.120$$

$$\mathbf{v}(d2_3) = 0$$

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

- $\text{sim}(d1, d2) = \mathbf{v}(d1) \cdot \mathbf{v}(d2)$
- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$
- normalising for length:

$$\mathbf{v}(d_i) = \frac{\mathbf{V}(d_i)}{\|\mathbf{V}(d)\|}$$

vocabulary	d1	d2	d3
1: affection	115	58	20
2: jealous	10	7	11
3: gossip	2	0	6

$$\|\mathbf{V}(d3)\| = \sqrt{20^2 + 11^2 + 6^2}$$

- Euclidean length:

$$\|\mathbf{V}(d)\| = \sqrt{\sum_{i=1}^n \mathbf{V}_i^2(d)}$$

$$\mathbf{v}(d3_1) = \frac{20}{\sqrt{20^2+11^2+6^2}} = 0.847$$

$$\mathbf{v}(d3_2) = \frac{11}{\sqrt{20^2+11^2+6^2}} = 0.466$$

$$\mathbf{v}(d3_3) = \frac{6}{\sqrt{20^2+11^2+6^2}} = 0.254$$

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

- $\text{sim}(d1, d3) = \mathbf{v}(d1) \cdot \mathbf{v}(d3)$

- $\mathbf{v}(d1) \cdot \mathbf{v}(d3) = \sum_{i=1}^n d1_i d3_i$

$$\mathbf{v}(d1_1) = \frac{115}{\sqrt{115^2+10^2+2^2}} = 0.996$$

$$\mathbf{v}(d1_2) = \frac{10}{\sqrt{115^2+10^2+2^2}} = 0.087$$

$$\mathbf{v}(d1_3) = \frac{2}{\sqrt{115^2+10^2+2^2}} = 0.017$$

$$[i=1] \quad 0.996 * 0.847 +$$

$$[i=2] \quad 0.087 * 0.466 +$$

$$[i=3] \quad 0.017 * 0.254 =$$

$$= 0.888$$

$$\mathbf{v}(d3_1) = \frac{20}{\sqrt{20^2+11^2+6^2}} = 0.847$$

$$\mathbf{v}(d3_2) = \frac{11}{\sqrt{20^2+11^2+6^2}} = 0.466$$

$$\mathbf{v}(d3_3) = \frac{6}{\sqrt{20^2+11^2+6^2}} = 0.254$$

---

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

- $\text{sim}(d1, d2) = \mathbf{v}(d1) \cdot \mathbf{v}(d2)$        $\mathbf{v}(d1_1) = \frac{115}{\sqrt{115^2+10^2+2^2}} = 0.996$
- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$        $\mathbf{v}(d1_2) = \frac{10}{\sqrt{115^2+10^2+2^2}} = 0.087$
- $\mathbf{v}(d1) \cdot \mathbf{v}(d2) = \sum_{i=1}^n d1_i d2_i$        $\mathbf{v}(d1_3) = \frac{2}{\sqrt{115^2+10^2+2^2}} = 0.017$

$$\begin{aligned} [\text{i}=1] \quad & 0.996 * 0.993 + \\ [\text{i}=2] \quad & 0.087 * 0.120 + \\ [\text{i}=3] \quad & 0.017 * 0.000 = \\ & = 0.999 \end{aligned}$$

$$\begin{aligned} \mathbf{v}(d2_1) &= \frac{58}{\sqrt{58^2+7^2+0}} = 0.993 \\ \mathbf{v}(d2_2) &= \frac{7}{\sqrt{58^2+7^2+0}} = 0.120 \\ \mathbf{v}(d2_3) &= 0 \end{aligned}$$

---

<sup>4</sup>See [MRS08]

## Example (Manning et al<sup>4</sup>)

summing up:

dictionary	$v(d_1)$	$v(d_2)$	$v(d_3)$
affection	0.996	0.993	0.847
jealous	0.087	0.120	0.466
gossip	0.017	0	0.254

$$sim(d_1, d_2) = 0.999$$

$$sim(d_1, d_3) = 0.888$$

---

<sup>4</sup>See [MRS08]

# New documents

- each new document  $n$  is represented using vectors in the same way

	$d_1$	$d_2$	$d_3$	$n$
affection	115	58	20	0
jealous	10	7	11	1
gossip	2	0	6	1

- $\text{sim}(n, d) = \mathbf{v}(n) \cdot \mathbf{v}(d)$

# New documents

- each new document  $n$  is represented using vectors in the same way

	$d_1$	$d_2$	$d_3$	$n$
affection	115	58	20	0
jealous	10	7	11	1
gossip	2	0	6	1

- $\text{sim}(n, d) = \mathbf{v}(n) \cdot \mathbf{v}(d)$

# New documents

- each new document  $n$  is represented using vectors in the same way

	$d_1$	$d_2$	$d_3$	$n$
affection	115	58	20	0
jealous	10	7	11	1
gossip	2	0	6	1

- $\text{sim}(n, d) = \mathbf{v}(n) \cdot \mathbf{v}(d)$
- with  $n = \langle \text{jealous}, \text{gossip} \rangle$

we obtain:

$$\mathbf{v}(n) \cdot \mathbf{v}(d_1) = 0.074$$

$$\mathbf{v}(n) \cdot \mathbf{v}(d_2) = 0.085$$

$$\mathbf{v}(n) \cdot \mathbf{v}(d_3) = 0.509$$

# Vector Space Models

Useful for

- Search
- Clustering
- **Discovering the main topics in a corpus**

# Topic Models

## Statistical models that

- Determine the themes over the documents / corpus
- Useful for: organizing, searching, summarizing texts

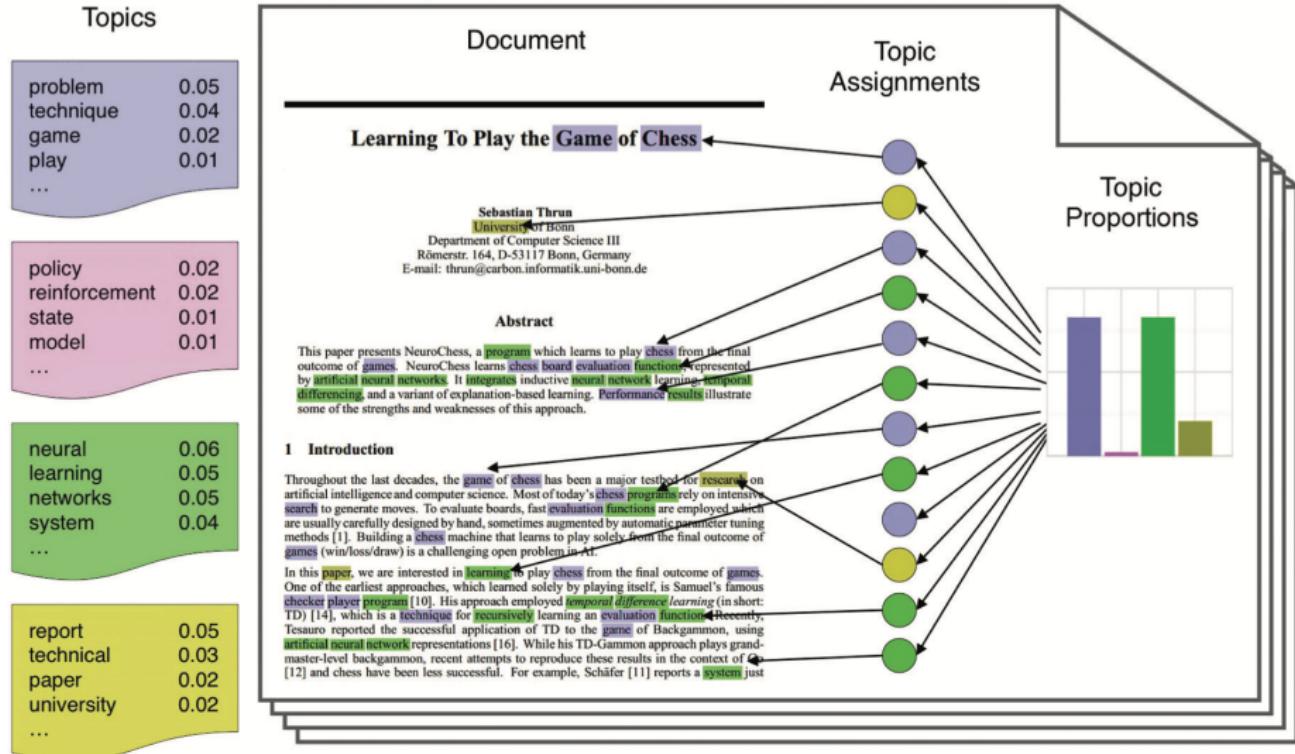
# Topic Models

## Statistical models that

- Determine the themes over the documents / corpus
- Useful for: organizing, searching, summarizing texts

Based on the **distributional hypothesis**: words close in meaning will occur in similar pieces of text.

# Topic Models Example<sup>5</sup>



<sup>5</sup>See [ES15]

# Topic

## Topic

- cluster of words that frequently occur together and share the same theme.
- (formally) a distribution over a fixed vocabulary of words  
Different topics = different distributions over the same vocabulary
- can emerge solely based on text analysis

# Latent Dirichlet Allocation

- Document is represented as a **mix of latent topics**

# Latent Dirichlet Allocation

- Document is represented as a **mix of latent topics**
- Number of latent topics ( $k$ ) is predefined
  - unlike k-means a document can belong to multiple topics
  - unlike GMM features do not need to be normally distributed

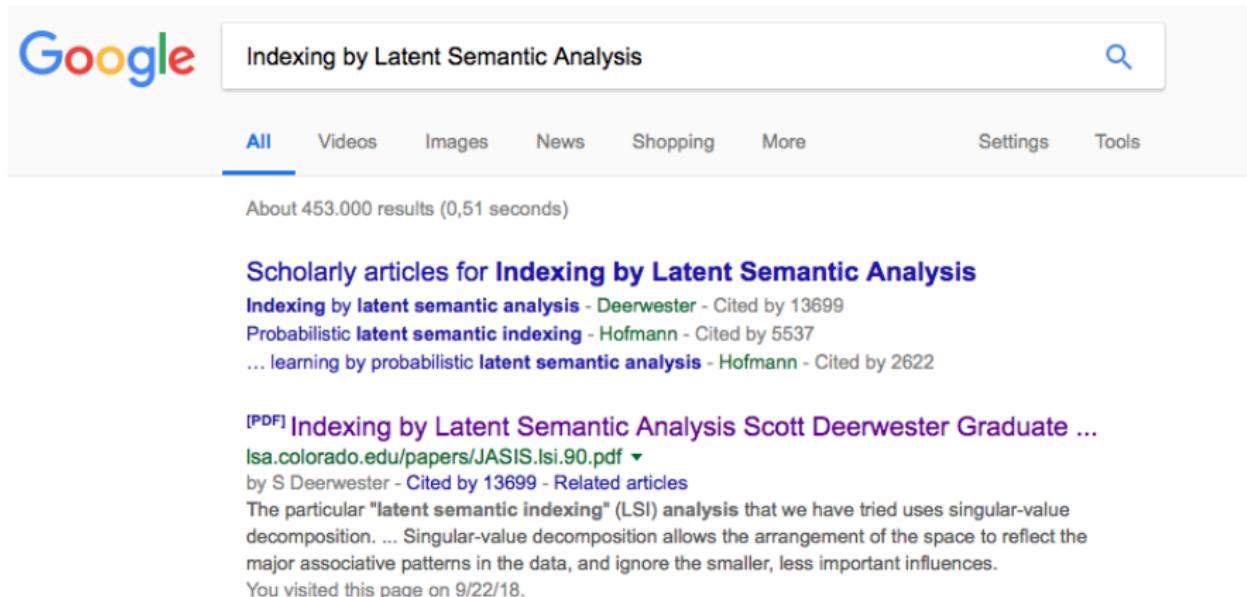
# Latent Dirichlet Allocation

- Document is represented as a **mix of latent topics**
- Number of latent topics ( $k$ ) is predefined
  - unlike k-means a document can belong to multiple topics
  - unlike GMM features do not need to be normally distributed

One of the simplest topic models

implemented in many libraries

# Latent Semantic Analysis



Google search results for "Indexing by Latent Semantic Analysis". The search bar shows the query. Below it, the "All" tab is selected, followed by other categories: Videos, Images, News, Shopping, More, Settings, and Tools. The search results page displays approximately 453,000 results found in 0.51 seconds. A section titled "Scholarly articles for Indexing by Latent Semantic Analysis" lists several top results, including a PDF link to a Scott Deerwester graduate thesis. Below this, a snippet of text from the thesis discusses the LSI analysis and singular-value decomposition.

Indexing by Latent Semantic Analysis

About 453.000 results (0,51 seconds)

### Scholarly articles for Indexing by Latent Semantic Analysis

[Indexing by latent semantic analysis](#) - Deerwester - Cited by 13699

[Probabilistic latent semantic indexing](#) - Hofmann - Cited by 5537

[... learning by probabilistic latent semantic analysis](#) - Hofmann - Cited by 2622

[\[PDF\] Indexing by Latent Semantic Analysis Scott Deerwester Graduate ...](#)  
isa.colorado.edu/papers/JASIS.Lsi.90.pdf ▾  
by S Deerwester - Cited by 13699 - Related articles

The particular "latent semantic indexing" (LSI) analysis that we have tried uses singular-value decomposition. ... Singular-value decomposition allows the arrangement of the space to reflect the major associative patterns in the data, and ignore the smaller, less important influences.

You visited this page on 9/22/18.

# Example

## Titles:

- c1: Human machine *interface* for Lab ABC *computer* applications
- c2: A survey of *user* opinion of *computer system response time*
- c3: The *EPS user interface management system*
- c4: *System* and *human system engineering* testing of *EPS*
- c5: Relation of *user-perceived response time* to error measurement

- m1: The generation of random, binary, unordered *trees*
- m2: The intersection *graph* of paths in *trees*
- m3: *Graph minors IV: Widths of trees* and well-quasi-ordering
- m4: *Graph minors: A survey*

Terms	Documents								
	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

# Example

## Titles:

- c1: Human machine *interface* for Lab ABC *computer* applications
- c2: A survey of *user* opinion of *computer system response time*
- c3: The *EPS user interface management system*
- c4: *System* and *human system engineering* testing of *EPS*
- c5: Relation of *user-perceived response time* to error measurement

- m1: The generation of random, binary, unordered *trees*
- m2: The intersection *graph* of paths in *trees*
- m3: *Graph minors IV: Widths of trees* and well-quasi-ordering
- m4: *Graph minors: A survey*

Terms	Documents								
	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

- Query: "Human Computer Interaction"

# Example

## Titles:

- c1: Human machine *interface* for Lab ABC *computer* applications
- c2: A survey of *user* opinion of *computer system response time*
- c3: The *EPS user interface management system*
- c4: *System* and *human system engineering* testing of *EPS*
- c5: Relation of *user-perceived response time* to error measurement

- m1: The generation of random, binary, unordered *trees*
- m2: The intersection *graph* of paths in *trees*
- m3: *Graph minors IV: Widths of trees* and well-quasi-ordering
- m4: *Graph minors: A survey*

Terms	Documents								
	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

- Query: "Human Computer Interaction"
- Simple term matching techniques would return documents c1, c2 and c4
  - they share one or more terms with the query.

# Example

## Titles:

- c1: Human machine *interface* for Lab ABC *computer* applications
- c2: A survey of *user* opinion of *computer system response time*
- c3: The *EPS user interface management system*
- c4: *System* and *human system engineering* testing of *EPS*
- c5: Relation of *user-perceived response time* to error measurement

- m1: The generation of random, binary, unordered *trees*
- m2: The intersection *graph* of paths in *trees*
- m3: *Graph minors IV: Widths of trees* and well-quasi-ordering
- m4: *Graph minors: A survey*

Terms	Documents								
	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

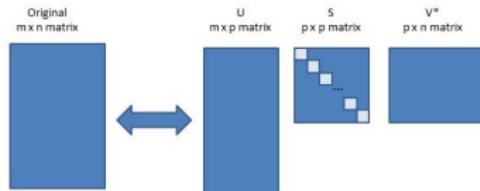
- Query: "Human Computer Interaction"
- Simple term matching techniques would return documents c1, c2 and c4
  - they share one or more terms with the query.
- Documents c3 and c5, which are also relevant, are missed as they have no terms in common with the query.

# Latent Semantic Analysis

- Create a term-document matrix ( $M$ )  
word counts = rows / documents = columns

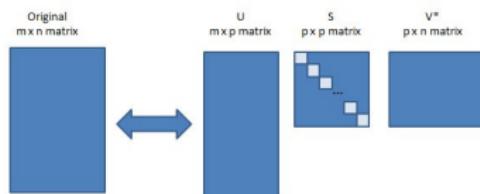
# Latent Semantic Analysis

- Create a term-document matrix ( $M$ )  
word counts = rows / documents = columns
- Perform Singular Value Decomposition on  $M$
- Reduce the dimension of the  $\Sigma$  in the SVD



# Latent Semantic Analysis

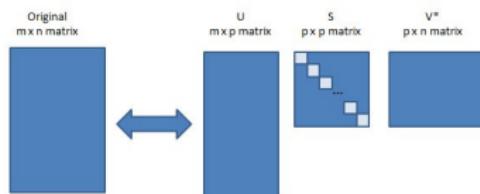
- Create a term-document matrix ( $M$ )  
word counts = rows / documents = columns
- Perform Singular Value Decomposition on  $M$
- Reduce the dimension of the  $\Sigma$  in the SVD



- Create a new matrix  $M'$  based on the reduced  $\Sigma$   
words can be compared using cosine distance

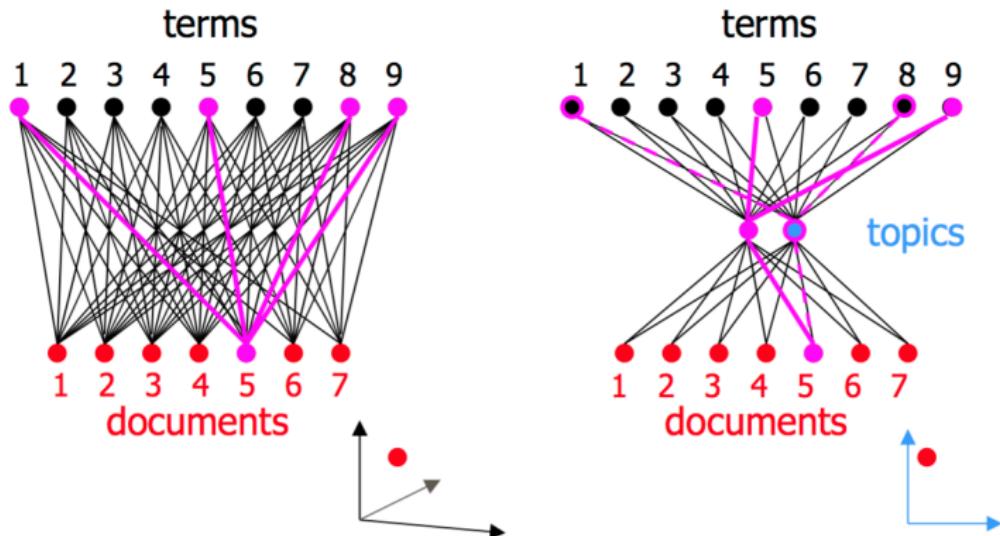
# Latent Semantic Analysis

- Create a term-document matrix ( $M$ )  
word counts = rows / documents = columns
- Perform Singular Value Decomposition on  $M$
- Reduce the dimension of the  $\Sigma$  in the SVD

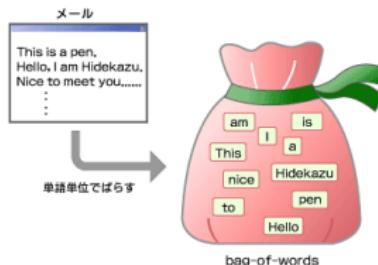


- Create a new matrix  $M'$  based on the reduced  $\Sigma$   
words can be compared using cosine distance
- Now we can
  - Search as we did in vector space models
  - Cluster the document vectors to discover topics

# Singular Value Decomposition

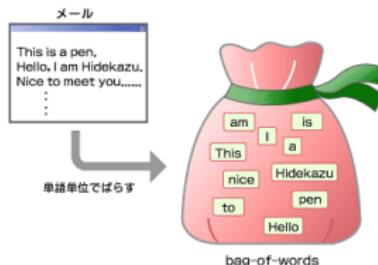


# Limitations of the Bag-of-Words



- Cannot capture phrases and **multi-word expressions**  
e.g. San Juan, a priori, etc.

# Limitations of the Bag-of-Words



- Cannot capture phrases and **multi-word expressions**  
e.g. San Juan, a priori, etc.
- Doesn't account for potential misspellings or word derivations.

# N-grams

- Given a corpus of text, the **n-grams** are the sequences of **n** consecutive words that are in the corpus.
- Can be with words or letters (characters)

# Sliding windows

- the cat that sat on the sofa also sat on the mat.
- the cat that sat on the sofa also sat on the mat.
- the cat that sat on the sofa also sat on the mat.

there are 3-grams, 4-grams, etc...

the sliding window can be larger than  $n!$

# Sliding windows

- the cat that sat on the sofa also sat on the mat.
- the cat that sat on the sofa also sat on the mat.
- the cat that sat on the sofa also sat on the mat.

there are 3-grams, 4-grams, etc...

the sliding window can be larger than  $n!$

- are all combinations relevant?
- are the frequent ones better?

## Bigrams beyond chance: *collocations*

**Collocations:** recurrent combinations of words that co-occur more often than chance, sometimes with non-compositional meaning

- strong tea (NB: ???powerful tea)
- best practice
- ...

## Bigrams beyond chance: *collocations*

What does “more often than chance” mean?

## Bigrams beyond chance: *collocations*

What does “more often than chance” mean?

We want to compare the likelihood of 2 words next to each other being a **chance event** vs. being a **surprise**<sup>67</sup>

- Do the two words appear next to each other more than we might expect, based on what we know about their individual frequencies?

---

<sup>6</sup>[tdunning.blogspot.nl/2008/03/surprise-and-coincidence.html](http://tdunning.blogspot.nl/2008/03/surprise-and-coincidence.html)

<sup>7</sup>See [Dun93]

# Determining strength of collocation

We use **association measures** to determine whether two words occur together by chance or in a **significant way**

- pointwise mutual information (PMI)
- chi-square ( $\chi^2$ )
- log-likelihood ratio (LL)

Ex: Collocations

# References I

- [Dun93] Ted Dunning, *Accurate methods for the statistics of surprise and coincidence*, Comput. Linguist. **19** (1993), no. 1, 61–74.
- [ES15] EMC Education Services, *Data science and big data analytics*, John Wiley and Sons, Inc., Indianapolis, Indiana, 2015.
- [LKGL05] Mircea Lungu, Adrian Kuhn, Tudor Gîrba, and Michele Lanza, *Interactive exploration of semantic clusters*, 3rd International Workshop on Visualizing Software for Understanding and Analysis (VISSOFT 2005), 2005, pp. 95–100.
- [MRS08] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze, *Introduction to information retrieval*, Cambridge University Press, New York, NY, USA, 2008.
- [SB88] Gerard Salton and Christopher Buckley, *Term-weighting approaches in automatic text retrieval*, Inf. Process. Manage. **24** (1988), no. 5, 513–523.