

TEXT ANALYSIS



RELATED ACTIVITY
DOWNLOAD THE
“Exercise – Text Analysis.zip”
NOTEBOOK & DATASET

RELATIONAL DATA MODEL

Represent data as a **TABLE** (*relation*)

ROW (*tuple*) represents a single record
Each record is a fixed-length tuple

COLUMN (*attribute*) represents a single *variable*
Each has a name and a data type

SCHEMA – table's set of names and data types

DATABASE - a collection of tables

Month	Treatment	Pressure
March	Control	165
March	Placebo	163
March	300 mg	166
March	450 mg	168
April	Control	162
April	Placebo	159
April	300 mg	161
April	450 mg	163
May	Control	164

Blood Pressure Study (4 treatments, 6 months)

**TEXT IS DIFFERENT
COMMON
UNSTRUCTURED (MOSTLY)
HIGH-DIMENSIONAL (10,000+)
BIG!**

WHY ANALYZE TEXT?

WHY ANALYZE TEXT?

UNDERSTANDING: Examine patterns in word use.
Get the “gist” of a document.

GROUPING: Cluster for overview or classification.

COMPARE: Compare document collections,
or inspect evolution of collection over time.

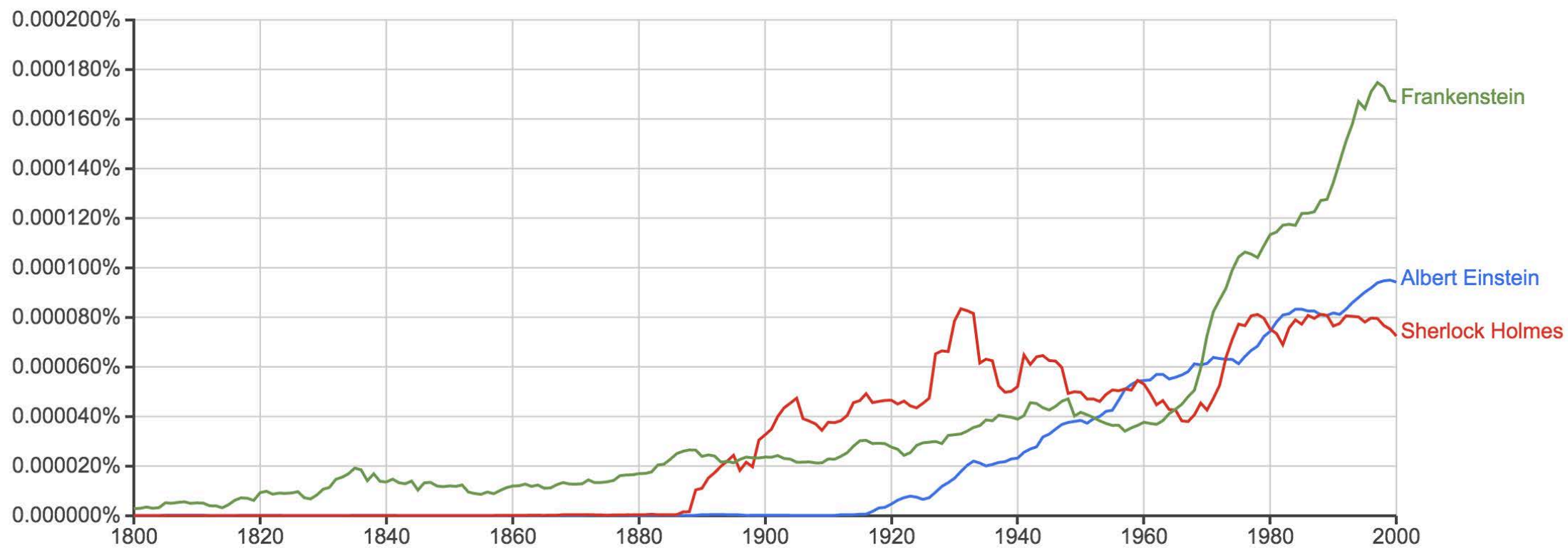
CORRELATE: Compare patterns in text to those in other data – for
example, correlate with social network.

Google Books Ngram Viewer

Graph these comma-separated phrases: ☐ case-insensitive

between and from the corpus with smoothing of

[Search lots of books](#)



(click on line/label for focus)



**Bible Verse
Length, Keim &
Oelke, VAST '07**

Neh 10



Those who sealed it were:

Nehemiah the governor, the son of Hakaliah.

Zedekiah, ² Seraiah, Azariah, Jeremiah,

³ Pashhur, Amariah, Malkijah,

⁴ Hattush, Shebaniah, Malluk,

⁵ Harim, Meremoth, Obadiah,

⁶ Daniel, Ginnethon, Baruch,

⁷ Meshullam, Abijah, Mijamin,

⁸ Maaziah, Bilgai and Shemaiah.

These were the priests.

⁹ The Levites:

Jeshua son of Azaniah, Binnui of the sons of Henadad, Kadmiel,

¹⁰ and their associates: Shebaniah,

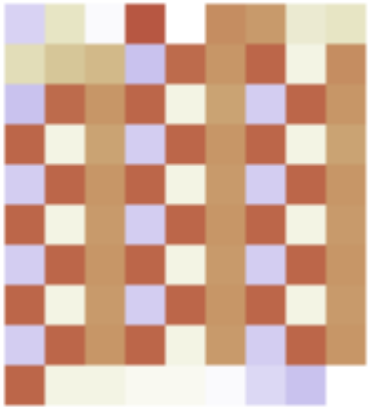
Hodiah, Kelita, Pelaiah, Hanan,

¹¹ Mika, Rehob, Hashabiah,

¹² Zakkur, Sherebiah, Shebaniah,

¹³ Hodiah, Bani and Beninu.

Num 7



The one who brought his offering on the first day was Nahshon son of Amminadab of the tribe of Judah.

¹³ His offering was one silver plate weighing a hundred and thirty shekels^[a] and one silver sprinkling bowl weighing seventy shekels,^[b] both according to the sanctuary shekel, each filled with the finest flour mixed with olive oil as a grain offering; ¹⁴ one gold dish weighing ten shekels,^[c] filled with incense; ¹⁵ one young bull, one ram and one male lamb a year old for a burnt offering; ¹⁶ one male goat for a sin offering^[d]; ¹⁷ and two oxen, five rams, five male goats and five male lambs a year old to be sacrificed as a fellowship offering. This was the offering of Nahshon son of Amminadab.

¹⁸ On the second day Nethanel son of Zuar, the leader of Issachar, brought his offering.

¹⁹ The offering he brought was one silver plate weighing a hundred and thirty shekels and one silver sprinkling bowl weighing seventy shekels, both according to the sanctuary shekel, each filled with the finest flour

Text analysis of Trump's tweets confirms he writes only the (angrier) Android half

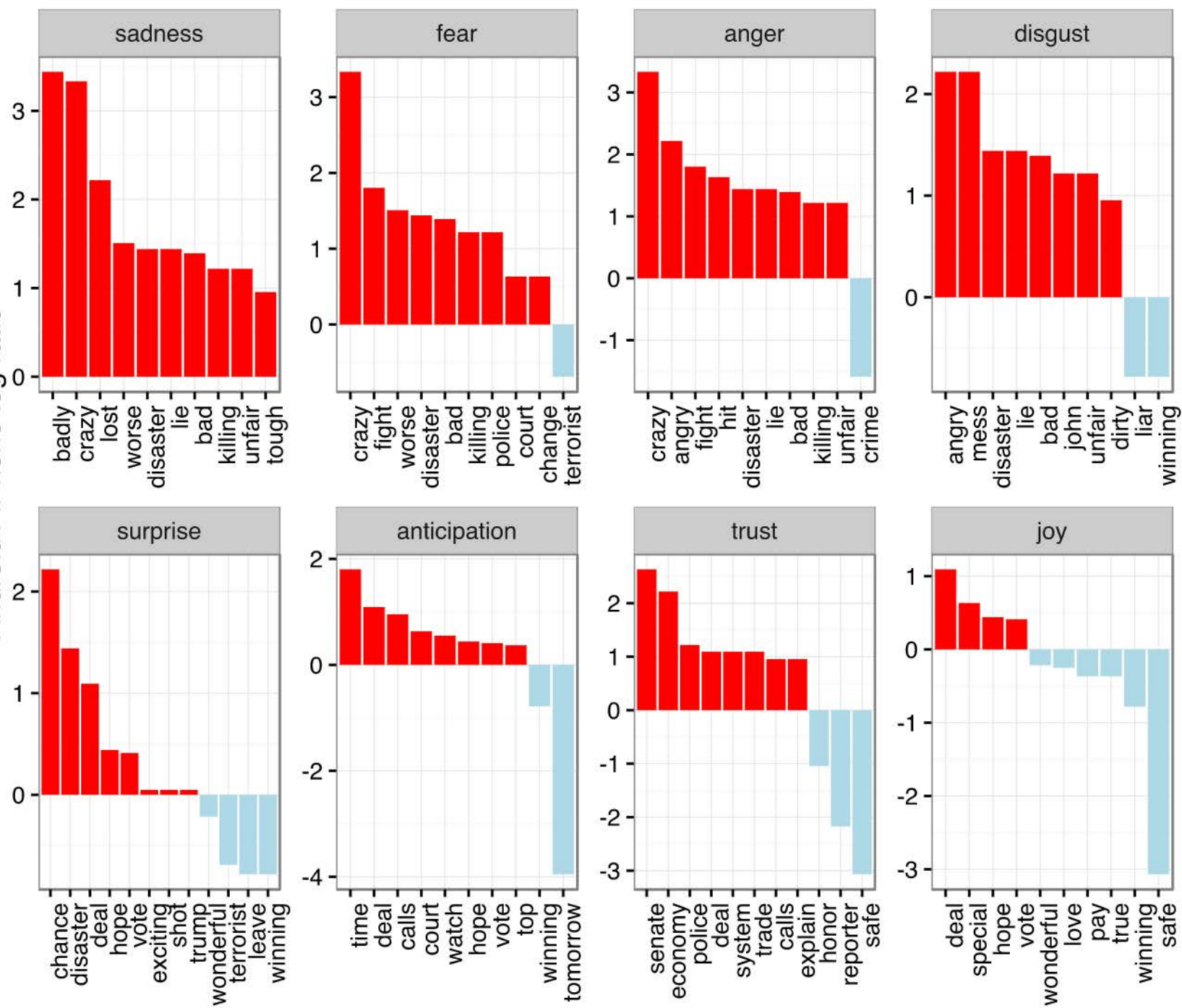
I don't normally post about politics (I'm not particularly savvy about polling, which is where data science has had the largest impact on politics). But this weekend I saw a hypothesis about Donald Trump's twitter account that simply begged to be investigated with data:



David Robinson (2016)

<http://varianceexplained.org/r/trump-tweets/>

Android / iPhone log ratio



WHAT IS TEXT DATA?

DOCUMENTS

Articles, books, and novels

Computer programs

E-mails, web pages, blogs

Messages, posts, tags, comments

COLLECTION OF DOCUMENTS

Messages (e-mail, blogs, tags, comments)

Social networks (personal profiles)

Academic collaborations (publications)

Even whole libraries, websites, social networks

TEXT AS DATA

Words are
the basic
unit of data.

DOCUMENT-LEVEL ATTRIBUTES & METADATA

LENGTH

DATE(S)

AUTHOR(S)

FORMAT

WORD-LEVEL ATTRIBUTES

WORD LENGTH

PART OF SPEECH (noun, verb, adjective, etc.)

FORMAT (*italic*, underline, etc.)

LANGUAGE (English? Latin? Japanese?)

FREQUENCY / DIFFICULTY (is it common?)

SENTIMENT (positive or negative connotation)

SYNONYMS / ANTONYMS / ETYMOLOGY (other meanings? roots?)

ENTITIES (“Calgary”, “Obama”, “Telus”, ...)

... AND MANY MORE

AGGREGATION

REPETITION
PLAGARISM
SHARED ENTITIES
AUTHOR STYLE

COLLECTION

DOCUMENT

SECTION

PAGE

PARAGRAPH

SENTENCE

WORD

TENSE
SENTIMENT
SENTENCE LENGTH
READING LEVEL

WHAT ABOUT THESE WORDS?



automate
automates
automatic
automation



automat

~~a, an, the, to, ...~~

“New York”

“Ban Ki-moon”

“Manchester United”

TEXT PROCESSING PIPELINE

TOKENIZATION: SEGMENT TEXT INTO TERMS

Entities? “San Francisco”, “O’Connor”, “U.S.A.”

Remove stop words? “a”, “an”, “the”, “to”, “be”

N-grams? Can take words in 2-word groups (bi-grams), 3-word (tri-grams), etc.

STEMMING: GROUP TOGETHER DIFFERENT FORMS

Roots: visualization(s), visualize(s), visually visual

lemmatization: goes, went, gone go

For visualization, sometimes need to reverse stemming for labels

Simple solution: map from stem to the most frequent word

RESULT: ORDERED STREAM OF TERMS

TEXT PROCESSING PIPELINE

“The quick brown fox jumps over the lazy dog.”

TOKENIZE (N=1)

[The], [quick], [brown], [fox], [jumps], [over], [the], [lazy], [dog].

TOKENIZE (N=1), REMOVE STOPWORDS, STEM

[quick], [brown], [fox], [jump], [over], [lazy], [dog]

TOKENIZE (N=2)

[the quick], [quick brown], [brown fox], [fox jumps], [jumps over], [over the] □

TOKENIZE (N=5)

[the quick brown fox jumps], [quick brown fox jumps over], [brown fox jumps over] □

PARTS-OF-SPEECH

Tag	Meaning	English Examples
ADJ	adjective	<i>new, good, high, special, big, local</i>
ADP	adposition	<i>on, of, at, with, by, into, under</i>
ADV	adverb	<i>really, already, still, early, now</i>
CONJ	conjunction	<i>and, or, but, if, while, although</i>
DET	determiner, article	<i>the, a, some, most, every, no, which</i>
NOUN	noun	<i>year, home, costs, time, Africa</i>
NUM	numeral	<i>twenty-four, fourth, 1991, 14:24</i>
PRT	particle	<i>at, on, out, over per, that, up, with</i>
PRON	pronoun	<i>he, their, her, its, my, I, us</i>
VERB	verb	<i>is, say, told, given, playing, would</i>
.	punctuation marks	<i>. , ; !</i>
X	other	<i>ersatz, esprit, dunno, gr8, univeristy</i>

NAMED ENTITY RECOGNITION

IDENTIFY AND CLASSIFY NAMED ENTITIES IN TEXT:

JOHN SMITH IS A PERSON

SOVIET UNION IS A COUNTRY

2500 UNIVERSITY DR IS AN ADDRESS

(555) 867-5309 IS A PHONE NUMBER

ENTITY RELATIONS: HOW DO THE ENTITIES RELATE?

DO THEY CO-OCCUR IN A DOCUMENT? IN A SENTENCE?

BUT THIS CAN BE DIFFICULT

KEEP UP **ON** YOUR **READING** WITH AUDIO **BOOKS**

Vietnam

UK

Louisiana, USA

Audio **books** are highly **popular** with **library** patrons in the **town**

Louisiana, USA

S. Carolina, USA

Pennsylvania, USA

Mass., USA

of **Springfield,** **Greene** County, **MO.** "People are **mobile**

Turkey

Virginia, USA

Maine, USA

Norway

Alabama, USA

and busier, and audio **books** fit into that lifestyle" says **Gary**

Louisiana, USA

Indiana, USA

Sanchez, who oversees the **library's** \$2 **million** budget...

Dominican Republic

Pennsylvania, USA

Kentucky, USA

WORDS AND PHRASES

WORD RELATIONSHIPS

“Concordance” – Words plus the context in which they appear.

The screenshot shows the 'Concordance - Larkin Concordance' window. The interface includes a menu bar (File, Text, Search, Edit, Headwords, Contexts, View, Tools, Help) and a toolbar with various icons. The main area is divided into a left pane showing a list of headwords and their frequencies, and a right pane showing the concordance results for the selected word 'HEART'.

Headword	No.	Context...	Word	...Context	Reference
HEAR	15	That my own	heart	drifts and cries, having no...	Deep Analysis
HEARD	9	By the shout of the	heart	continually at work	And the wave
HEARING	7	Nothing to adapt the skill of the	heart	to, skill	And the wave
HEARS	3	The tread, the beat of it, it is my own	heart	,	Träumerei
HEARSE	1	Because I follow it to my own	heart		Many famous
HEART	25	My	heart	is ticking like the sun:	I am washed u
HEART'S	2	The vague	heart	sharpened to a candid co...	The March Pa
HEART-SHAPED	1	Contract my	heart	by looking out of date.	Lines on a Yo
HEARTH	1	Having no	heart	to put aside the theft	Home is so Sa
HEARTS	7	And the boy puking his	heart	out in the Gents	Essential Bea
HEARTY	1	A harbour for the	heart	against distress.	Bridge for the
HEAT	6	These I would choose my	heart	to lead	After-Dinner F
HEAT-HAZE	1	Time in his little cinema of the	heart		Time and Spa
HEATH	1	This petrified	heart	has taken,	A Stone Churc
HEATS	1	How should they sweep the girl clean...	heart	,	I see a girl dra
HEAVE	1	Hands that the	heart	can govern	Heaviest of flc
HEAVEN	4	For the	heart	to be loveless, and as col...	Dawn
HEAVEN-HOLDING	1	With the unguessed-at	heart	riding	One man walk
HEAVIER-THAN...	1	If hands could free you,	heart	,	If hands could
HEAVIEST	2	That overflows the	heart		Pour away the

At the bottom, a status bar displays summary statistics: Words: 7318, Tokens: 37070, At word: 2990, Deleted lines: 1 [24], Word sort: Asc alpha (string), Context sort: Asc occurrence order.

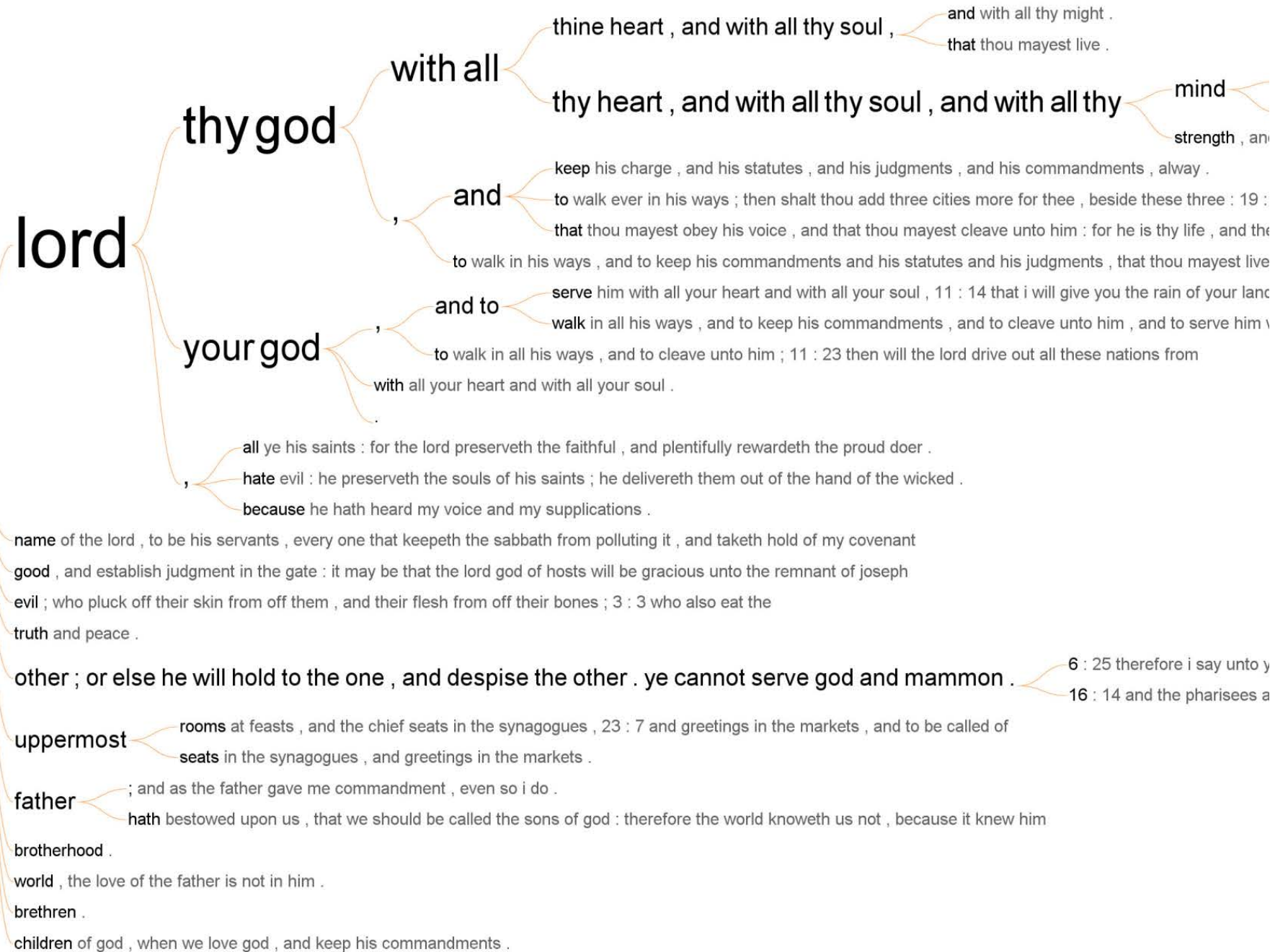
WORD TREES

- cats are better than dogs
- cats eat kibble
- cats are better than hamsters
- cats are awesome
- cats are people too
- cats eat mice
- cats meowing
- cats in the cradle
- cats eat mice
- cats in the cradle lyrics
- cats eat kibble
- cats for adoption
- cats are family
- cats eat mice
- cats are better than kittens
- cats are evil
- cats are weird
- cats eat mice

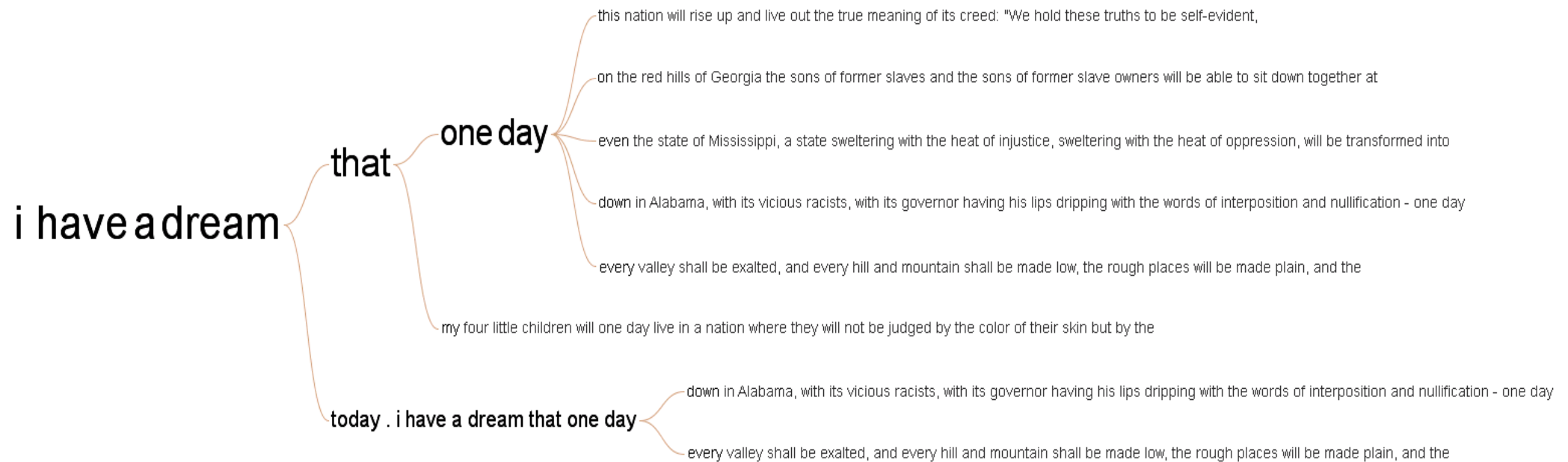


WATTENBERG & VIÉGAS 2008

love the



RECURRENT THEMES IN SPEECH



GLIMPSES OF STRUCTURE

Concordances show local, repeated structure,
but what about other types of patterns?

FOR EXAMPLE

LEXICAL: <A> at

SYNTACTIC: <Noun> <Verb> <Object>

PHRASE NETS

LOOK FOR SPECIFIC LINKING PATTERNS IN THE TEXT:

‘A **AND** B’, ‘A **AT** B’, ‘A **OF** B’, ETC

Could be output of regexp or parser

VISUALIZE EXTRACTED PATTERNS IN A NODE-LINK VIEW

OCCURRENCES = NODE SIZE

PATTERN POSITION = EDGE DIRECTION

Select a phrase

word1 and word2

word1 's word2

word1 of the word2

word1 the word2

word1 a word2

word1 at word2

word1 is word2

word1 [space] word2

or enter your own

* and *

Submit

Filters

Show top: 100

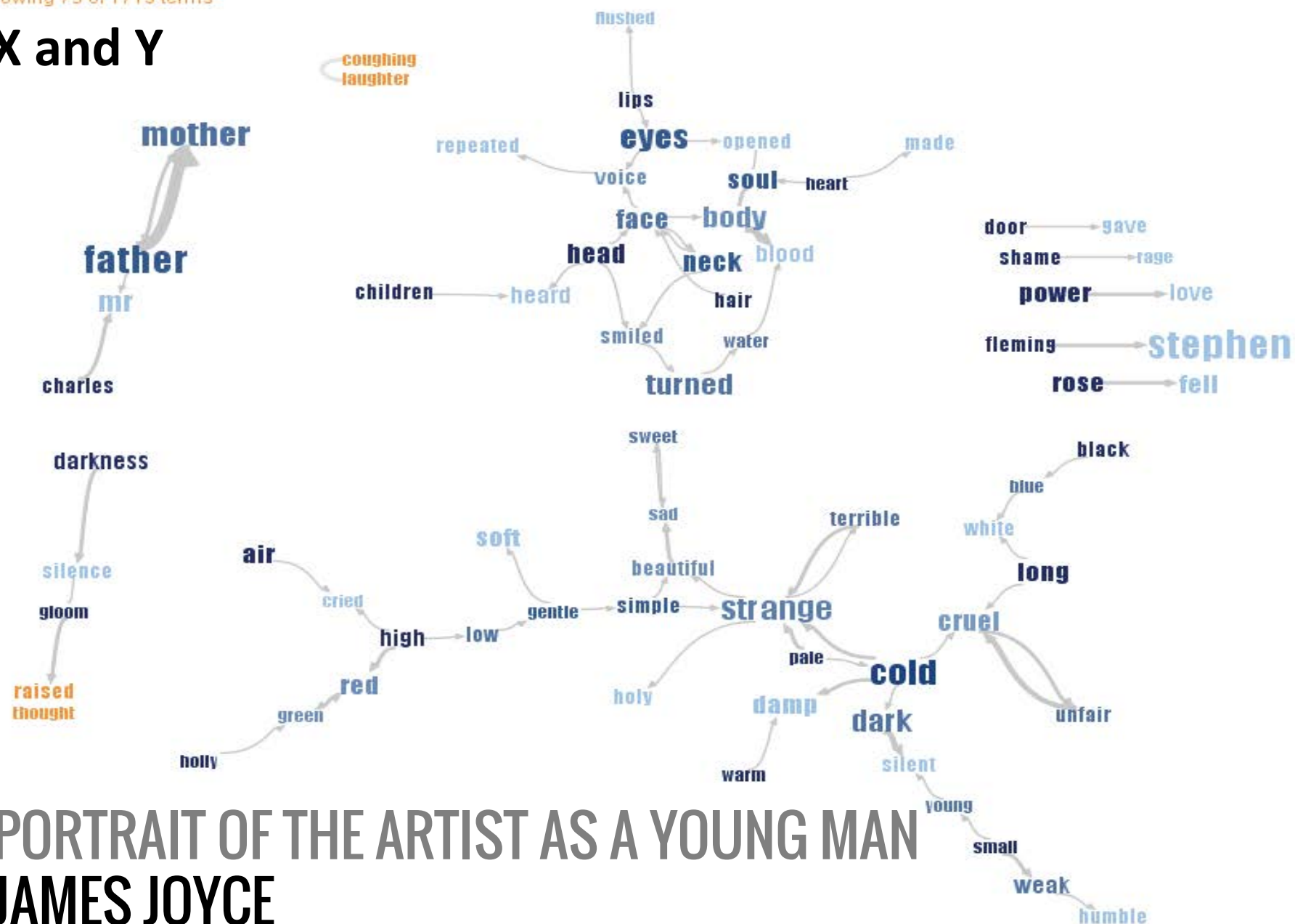
Hide common words ☒

Zoom

In Out Reset

Showing 73 of 1719 terms

X and Y

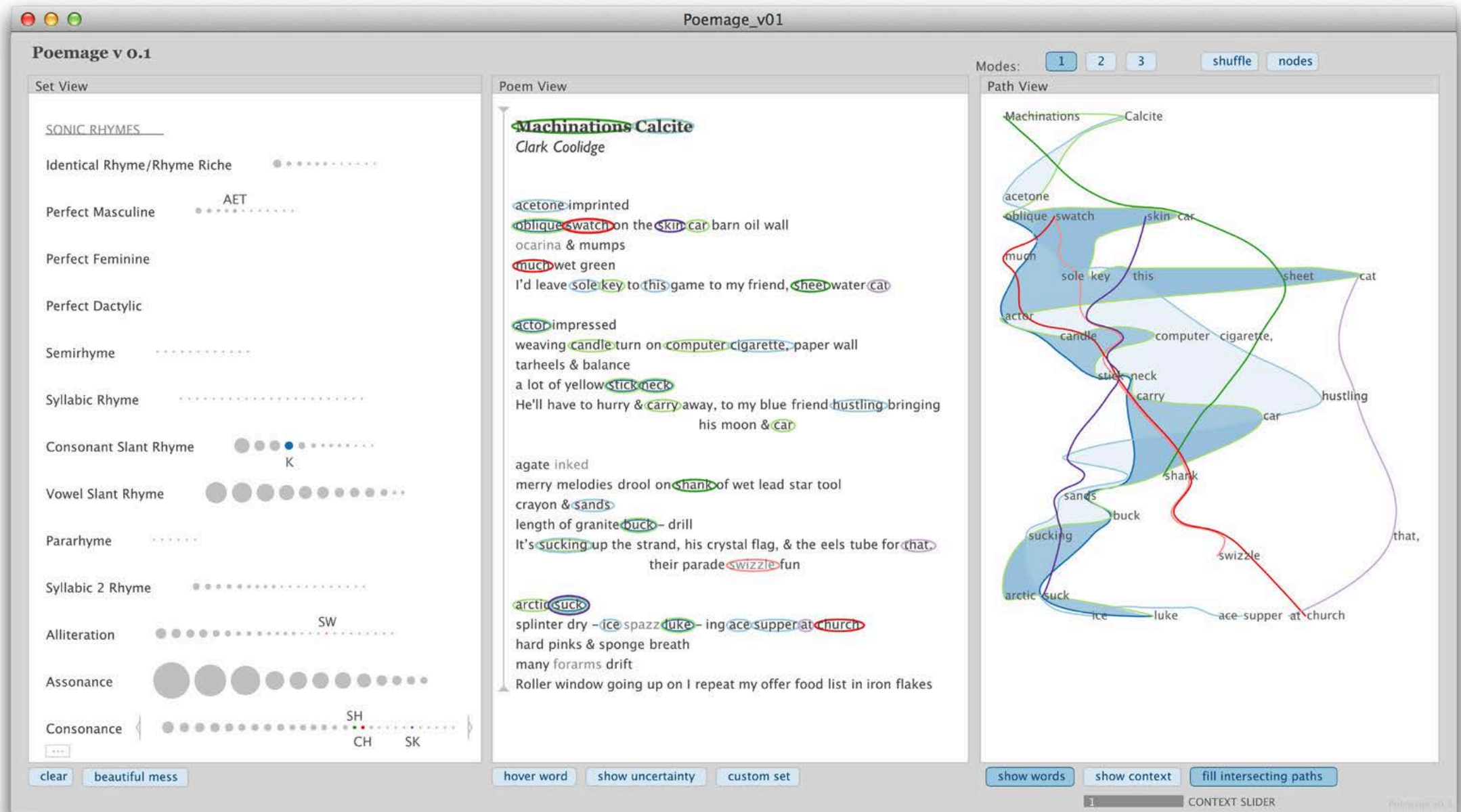


PORTRAIT OF THE ARTIST AS A YOUNG MAN

JAMES JOYCE

RHYME, SPEECH, ETC.

POEMAGE McCurdy Et al. 2016



NLTK (Natural Language ToolKit)

Tokenize and tag some text:

```
>>> import nltk
>>> sentence = """At eight o'clock on Thursday morning
... Arthur didn't feel very good."""
>>> tokens = nltk.word_tokenize(sentence)
>>> tokens
['At', 'eight', "o'clock", 'on', 'Thursday', 'morning',
 'Arthur', 'did', "n't", 'feel', 'very', 'good', '.']
>>> tagged = nltk.pos_tag(tokens)
>>> tagged[0:6]
[('At', 'IN'), ('eight', 'CD'), ("o'clock", 'JJ'), ('on', 'IN'),
 ('Thursday', 'NNP'), ('morning', 'NN')]
```

Identify named entities:

```
>>> entities = nltk.chunk.ne_chunk(tagged)
>>> entities
Tree('S', [
  (('At', 'IN'), ('eight', 'CD'), ("o'clock", 'JJ'),
   ('on', 'IN'), ('Thursday', 'NNP'), ('morning', 'NN')),
  Tree('PERSON', [
    ('Arthur', 'NNP')]),
  ('did', 'VBD'), ("n't", 'RB'), ('feel', 'VB'),
  ('very', 'RB'), ('good', 'JJ'), ('.', '.')]])
```

**LET'S TRY SOME
TEXT ANALYSIS**

DOCUMENTS AND COLLECTIONS

COMMON DOCUMENT-LEVEL METRICS

Measuring text complexity/author identity

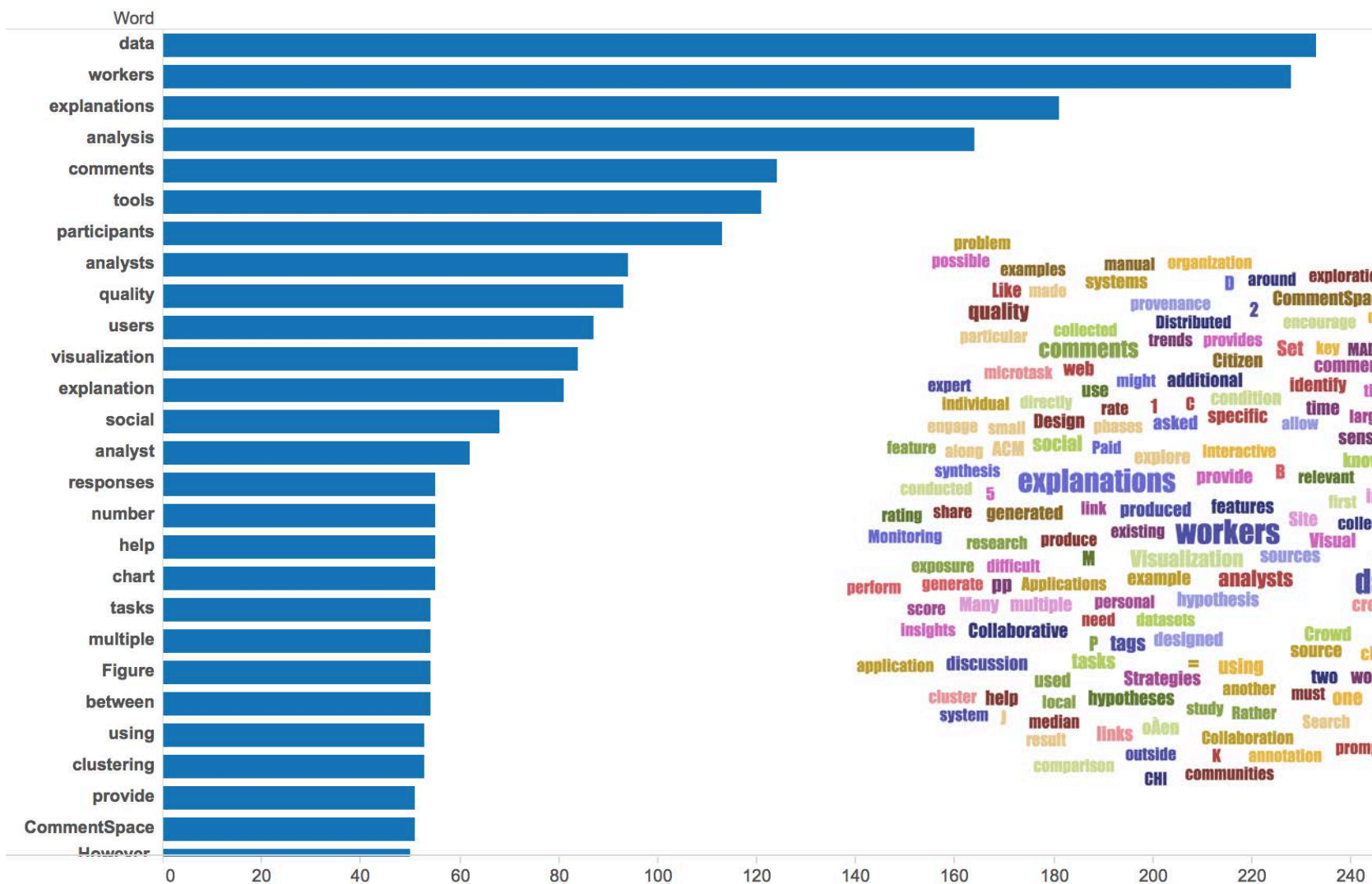
(often used in literary analysis)

- **Word length**
- **Syllables** per word
- **Average sentence length**
- **Percentage by parts of speech** (nouns, verbs, etc.)
- **Frequencies** of specific words
- **Language diversity** (number of words used)
- **Hapax Legomena** (words that appear only once)

HOW OFTEN DOES A GIVEN WORD
APPEAR IN A DOCUMENT?

WHAT TERMS ARE THE **MOST**
REPRESENTATIVE OF OR **MOST**
UNIQUE TO THIS DOCUMENT?

WORD COUNTS



WORD CLOUDS

... are usually a bad idea

- Size and word length are conflated
- Hard to visually search/compare
- Color skews salience
- No context

“The mullets of the internet.”

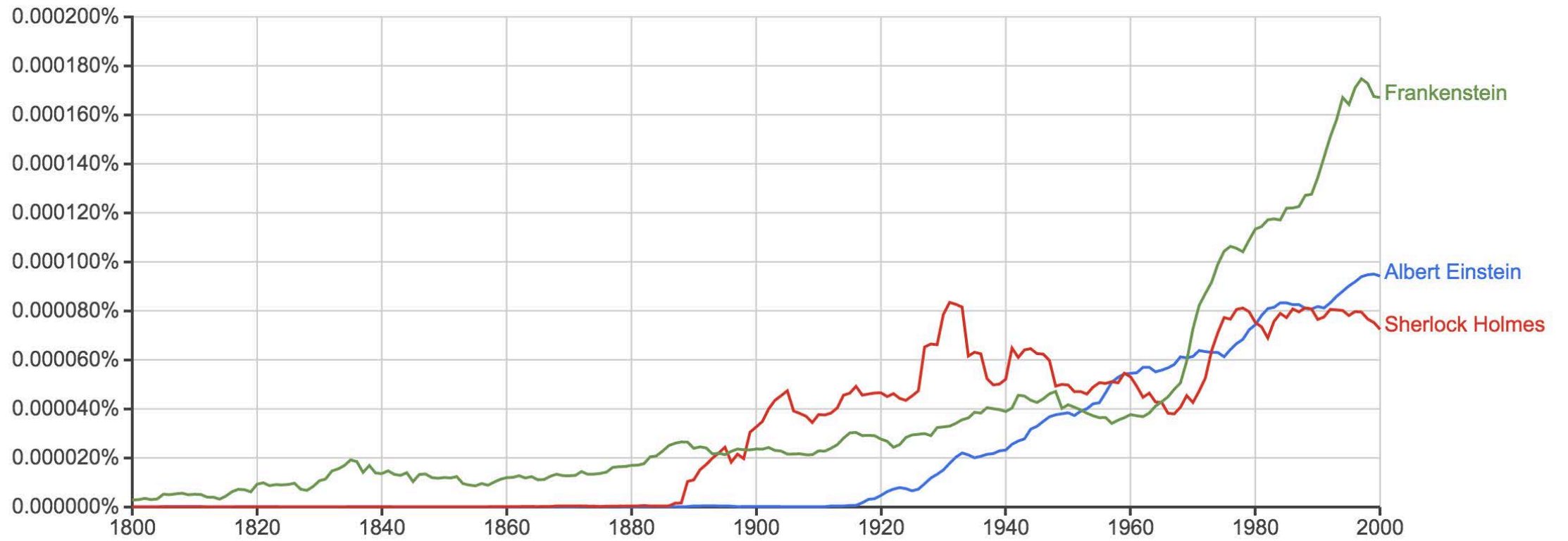


Google Books Ngram Viewer

Graph these comma-separated phrases: ☐ case-insensitive

between and from the corpus with smoothing of

[Search lots of books](#)



(click on line/label for focus)

IDENTIFYING IMPORTANT WORDS

Raw frequency isn't always appropriate

Stop words can pollute counts.

We often want to consider how counts compare to our expectations.

Weighting schemes

TF-IDF

weight by number of times word appears in collection

Probabilities

weight by probability of seeing the word

Etc.

TF-IDF (Term Frequency – Inverse Document Frequency)

Term Frequency

How common is this term in this doc.

for a term t in a document d in a corpus of N documents

$$tf_{td} = \text{count}(t) \text{ in } d$$

$$\log(1 + tf_{td}) \quad \leftarrow \text{log frequency to smooth}$$

$$tf_{td} / \sum_t tf_{td} \quad \leftarrow \text{normalize based on total number of terms in doc}$$

TF-IDF

How unique is this term to this doc.

$$tf.idf_{td} = \log(1 + tf_{td}) \times \log(N/df_t) \quad \leftarrow \text{weight by the inverse document frequency}$$

(fraction of the total docs that contain the term)

OTHER FREQUENCY STATISTICS

[Chuang et al 2012]

Table III. Frequency Statistics

Statistic	Definition
$\log(\text{tf})$	$\log(t_{\text{Doc}})$
tf.idf	$(t_{\text{Doc}}/t_{\text{Ref}}) \cdot \log(N/D)$
G^2	$2 \left(t_{\text{Doc}} \log \frac{t_{\text{Doc}} \cdot T_{\text{Ref}}}{T_{\text{Doc}} \cdot T_{\text{Doc}}} + t_{\overline{\text{Doc}}} \log \frac{t_{\overline{\text{Doc}}} \cdot T_{\text{Ref}}}{T_{\overline{\text{Doc}}} \cdot T_{\text{Doc}}} \right)$
BM25	$3 \cdot t_{\text{Doc}} / (t_{\text{Doc}} + 2(0.25 + 0.75 \cdot T_{\text{Doc}}/r)) \cdot \log(N/D)$
WordScore	$(t_{\text{Doc}} - t_{\text{Ref}}) / (T_{\overline{\text{Doc}}} - T_{\overline{\text{Ref}}})$
log-odds ratio (weighted)	$\left(\log \frac{t'_{\text{Doc}}}{t'_{\overline{\text{Doc}}}} - \log \frac{T'_{\text{Doc}}}{T'_{\overline{\text{Doc}}}} \right) / \sqrt{\frac{1}{t'_{\text{Doc}}} + \frac{1}{t'_{\overline{\text{Doc}}}}}$

FINDING SIMILAR DOCUMENTS... **AND IDENTIFYING GROUPS**

MANY APPROACHES EXIST (WE'LL LOOK AT 2)

(1) COMPUTE SIMILARITY BETWEEN DOCUMENTS

Based on the words they share.

For example, use words' TF-IDF scores to compute similarity.

(2) TOPIC MODELING

Assume documents mixtures of multiple topics.

Model topics based on co-occurring terms

SIMILARITY - TREAT WORDS AS FEATURES



Rex O'Saurus @UCRex · Oct 20

HUGE good luck, positive vibes, major butt kicking feels to our Lady Dinos Rugby team - competing for the CanWest title this w/end!

#GoDinos



Rex O'Saurus @UCRex · Oct 29

MAJOR GOOD LUCK to ⚽⚽⚽ in playoffs today! Get out and support them! They're wayyyy better than this I promise 👉 **#GoDinos**

#WeAreAllDinos



Rex O'Saurus @UCRex · Sep 25

Your @UCDinos record this weekend was 6-3! Looking to go 8-0 next weekend! **#DinosPride** **#WeAreAllDinos** **#GoDinos** **#ucalgary50** 🏈⚽



CAN USE BINARY 0/1,
COUNTS, NORMALIZED,
TF-IDF SCORES, ETC.

IF YOU WANT CLUSTERS,
JUST APPLY YOUR
FAVORITE CLUSTERING
TECHNIQUE!

Tweet #	good	luck	Dinos	Rugby	weekend	...
1	1	1	1	1	0	...
2	1	1	0	0	0	...
3	0	0	0	0	2	

SOME CHALLENGES

Space of all words is really large.

Similarity matrices can be very sparse.

Synonyms, homonyms, and multiple definitions (mean, lead, etc.).

Language is nuanced! Discrete clusters may not make sense.

TOPIC MODELING

Identify a set of “**topics**” that describe documents.

(Weighted combinations of terms.)

$CAT = (0.25 * \text{“cat”} + 0.15 * \text{“meow”} + 0.004 * \text{“toy”} + \dots)$

$TOYS = (0.10 * \text{“toy”} + 0.04 * \text{“child”} + 0.002 * \text{“game”} + \dots)$

Lost Cat! -----

30% CAT

...

New Fall Toys --

40% TOYS

...

Frolicat-----

18% CAT

22% TOYS

...



LOTS OF TECHNIQUES

LSA – Latent Semantic Analysis

NMF – Non-Negative Matrix Factorization

LDA – Latent Dirichlet Allocation

...

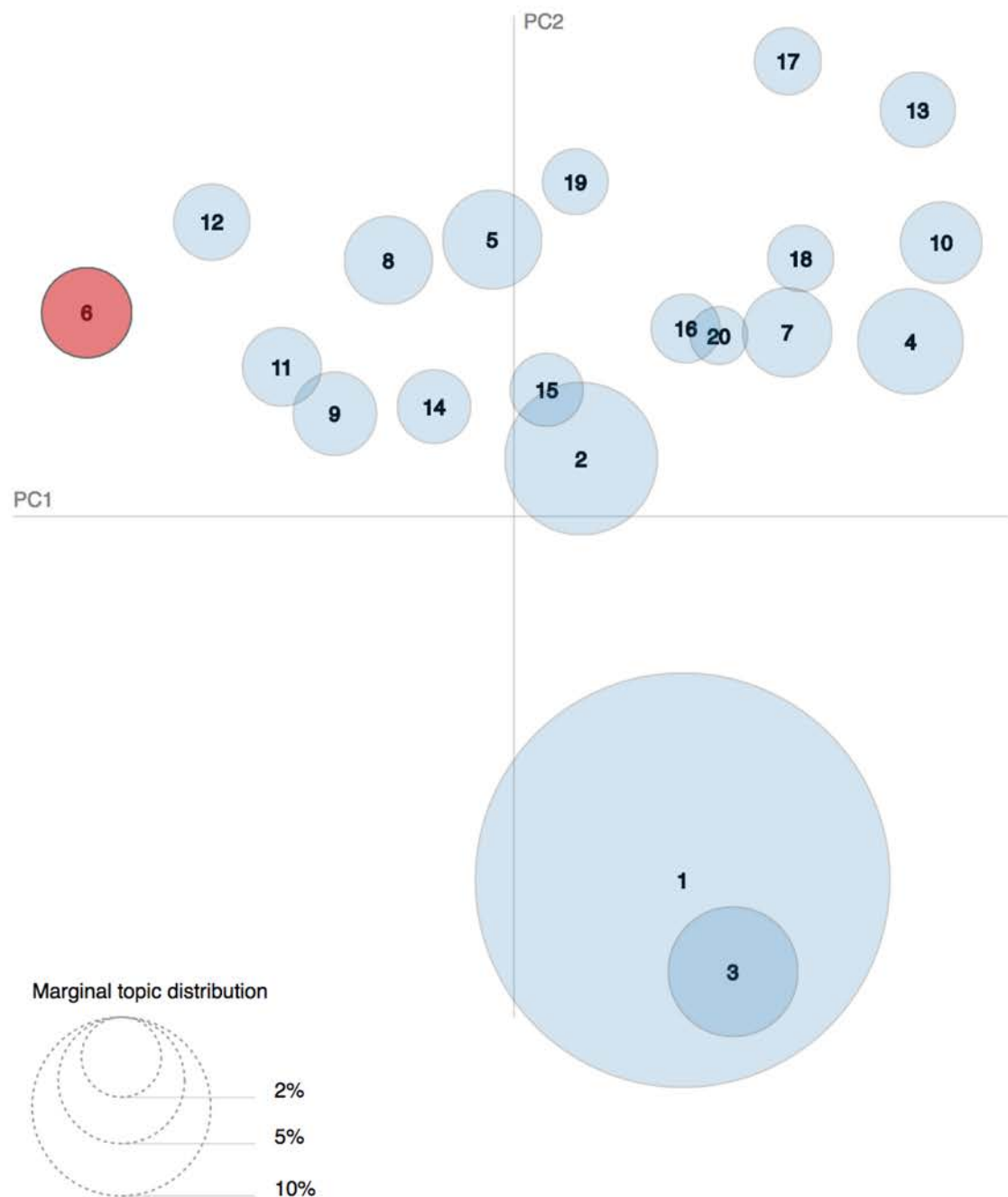
Less rigid than other classification methods. Documents can match **multiple topics simultaneously** and to **varying degrees**.

TOPIC MODELING WITH LDA

(A SIMPLIFIED EXPLANATION)

- Decide how many topics you want to model.
- For each document, randomly assign each word to a topic.
(Gives initial word/topic and doc/topic mappings ... but really bad ones.)
- For each word in each document compute:
 - $P(\text{topic} \mid \text{document})$ = fraction of words in a doc assigned to that topic
 - $P(\text{word} \mid \text{topic})$ = fraction of assignments to topic due to that word
 - Assign word a new topic based on $P(\text{topic} \mid \text{document}) * P(\text{word} \mid \text{topic})$
(Basically assumes other word-topic mappings are right and we're just fixing this one.)
- Repeat until topics converge.

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 6 (2.5% of tokens)



20 topics based on
the text of 2000
movie reviews.

1. $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$ for topics t ; see Chuang et. al (2012)
2. $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

SUMMARY

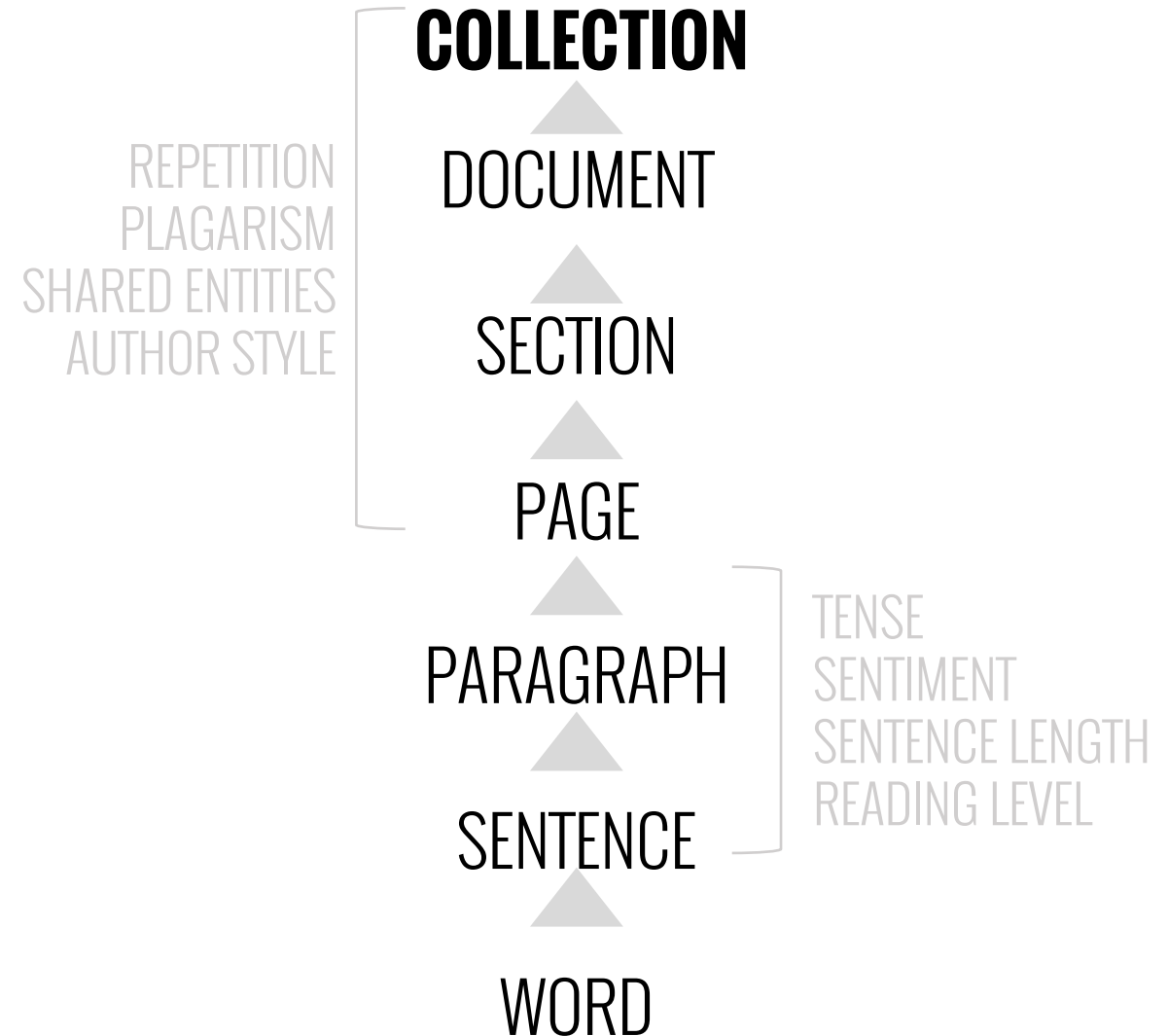
Text can be a really **rich** data source.

(Mostly) **unstructured** nature means analysis is sometimes more complex.

Can work at **many different levels**

Word→Sentence→Doc→Collection

Lots of **analysis approaches** and lots of **visualization techniques**.



LOTS OF TEXT VISUALIZATION TECHNIQUES

Text Visualization Browser

A Visual Survey of Text Visualization Techniques (IEEE PacificVis 2015 short paper)

Provided by ISOVIS group

About Summary Add entry

Techniques displayed:

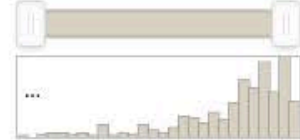
400

Search:

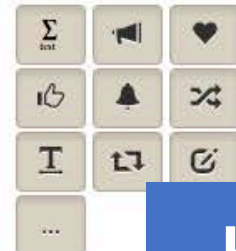
Time filter:

1976

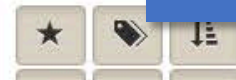
2017



Analytic Tasks



Visualization



<http://textvis.lnu.se/>

