

how to *do* computationally assisted research
digital literacy @ comwell

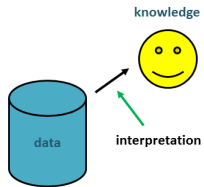
Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io/

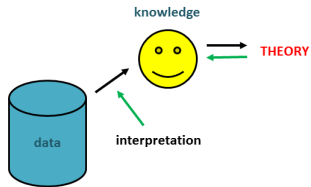
March 22, 2018

```
1 class Person(object):
2     def __init__(self, name):
3         self.name = name
4     def says_hello(self):
5         print 'Hello, my name is', self.name
6
7 class Researcher(Person):
8     def __init__(self, title=None, areas=None, **kwargs):
9         super(Researcher, self).__init__(**kwargs)
10        self.title = title
11        self.areas = areas
12
13 KLN = Researcher(name = 'Kristoffer L Nielbo', \
14                 title = 'Associate professor', \
15                 areas = ['Humanities Computing', 'Culture Analytics', 'eScience'])
16
17 KLN.says_hello()
```



evolution of workflows

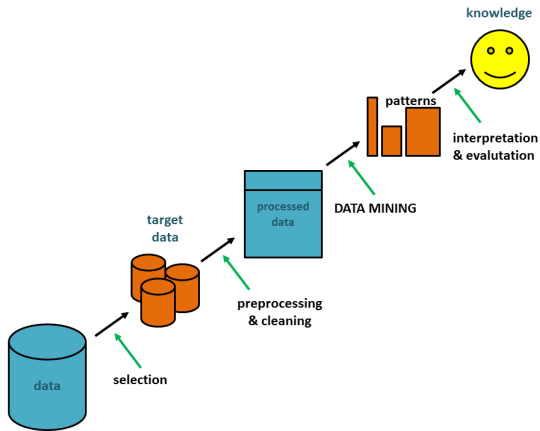


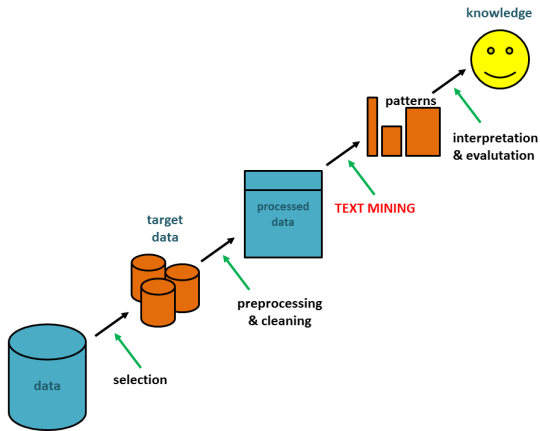


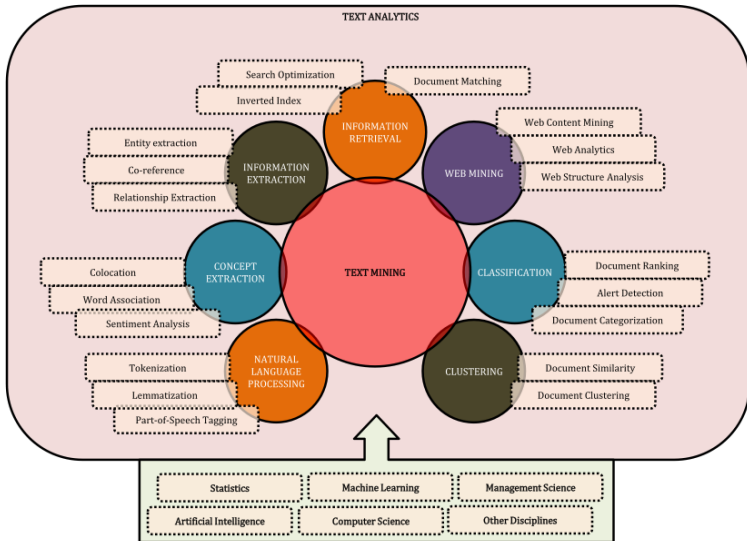
knowledge

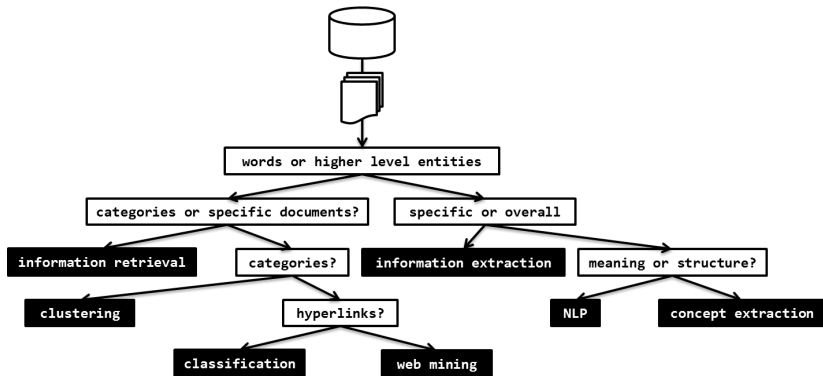


THEORY



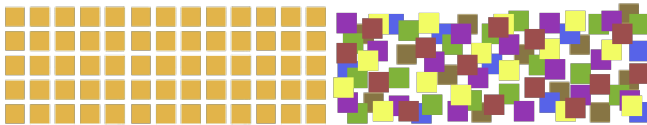






data

data objects that are described over a set of (qualitative or quantitative) features



- fundamental difference between structured data and **unstructured* data**
- word processing files, pdfs, emails, social media posts, digital images, video, and audio
- today > 80% of all data are unstructured
- unstructured data require expertise in culture, media, linguistic ...

data|access and sampling

select (sample*) a set of documents (target data) relevant to your research question from a data collection

>> online databases and research libraries are excellent resources

- proprietary issues
- data protection acts
- ethical concerns
- availability (e.g., historical sources)

>> sample requirements

- “all the data”
- balancing and stratification
- bias reduction



we will focus on documents stored locally in a *plain text* without markup

```
1  """The First Book of Moses, called Genesis
2
3      {1:1} In the beginning God created the heaven and the earth. {1:2}
4  And the earth was without form, and void; and darkness was upon the
5  face of the deep. And the Spirit of God moved upon the face of the
6  waters.
7
8      {1:3} And God said, Let there be light: and there was light. {1:4}
9  And God saw the light, that it was good: and God divided the light"""
```

BUT with a bit of code everything is possible

```
1  import urllib2
2  from HTMLParser import HTMLParser
3
4  class html_parser(HTMLParser):
5      def handle_starttag(self, tag, attrs):
6          print "start tag:", tag
7      def handle_endtag(self, tag):
8          print "end tag :", tag
9      def handle_data(self, data):
10         print "data :", data
11
12  url = "https://knielbo.github.io/"
13  response = urllib2.urlopen(url)
14  webpage = response.read()
15  parser = html_parser()
16  parser.feed(webpage)
```

preprocessing

example

- normalization by reducing inflected words to their stem, base or root form
- the stem need *not* be identical to the morphological root
- sufficient that related words map to the same stem (stem \neq valid root)
- search engines treat words with the same stem as synonyms (conflation)

Porter stemming algorithm - step 1a

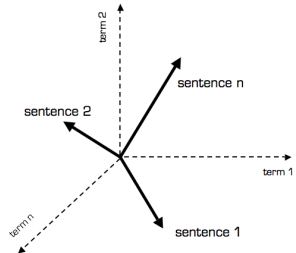
```
1 SSES -> SS      caresses -> caress
2 IES  -> I       ponies   -> poni
3      ties       -> ti
4 SS   -> SS      caress   -> caress
5 S    ->         cats     -> cat
```


example

any collection of m documents can be represented in the vector space model by a document-term matrix of m documents and n terms

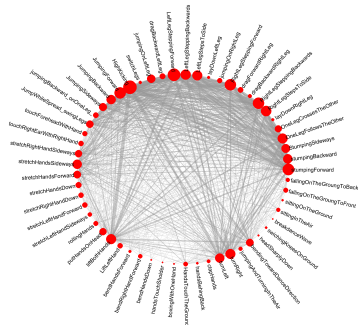
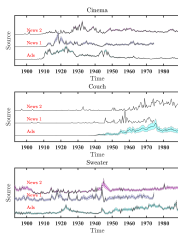
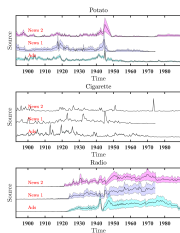
a vector space model is a basic modeling mechanism for a word- or document-space (whether we look at rows or columns)

- a document vector with only one word is collinear to the vocabulary word axis
- a document vector that does not contain a specific word is orthogonal/perpendicular to the word axis
- two documents are identical if they contain the same words in a different order (BOW assumption)



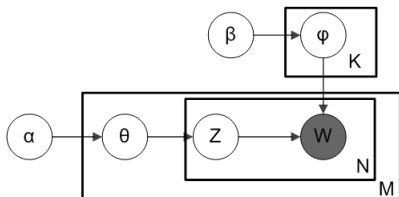
| Document space | t_1 | t_2 | t_3 | ... | t_n | ← Term vector space |
|----------------|----------|----------|----------|-----|----------|---------------------|
| D_1 | a_{11} | a_{12} | a_{13} | ... | a_{1n} | |
| D_2 | a_{21} | a_{22} | a_{23} | ... | a_{2n} | |
| D_3 | a_{31} | a_{32} | a_{33} | ... | a_{3n} | |
| ... | | | | | | |
| D_m | a_{m1} | a_{m2} | a_{m3} | ... | a_{mn} | |
| Q | b_1 | b_2 | b_3 | ... | b_n | |

analysis



- describe basic properties of the data, e.g., simple distributions and relations
- result in themselves or input to more advanced analysis
- the value depends critically on domain knowledge

beauty lies in simplicity



- α Dirichlet prior for per-doc topic dist
- proportions parameter
- β Dirichlet prior for per-topic word dist
- topic parameter
- θ_i word dist for topic
- per-document topic proportions
- ϕ_k word dist for topic k
- topics
- Z_{ij} topic for j^{th} word in doc i
- per-word topic assignment
- W_{ij} the observed word

Procedure 1 Generative Model

- 1: **choose** $\theta_i \sim \text{Dir}(\alpha)$, i is a document
 - 2: **choose** $\phi_k \sim \text{Dir}(\beta)$, k is a topic
 - 3: **for** each word position **do**
 - 4: **choose** a topic $z_{ij} \sim \text{Multinomial}(\theta_i)$
 - 5: **choose** a word $w_{ij} \sim \text{Multinomial}(\phi_{z_{ij}})$
 - 6: **end for**
-

The joint distribution defines a posterior probability: $P(\theta, z, \phi)$

use posterior to:

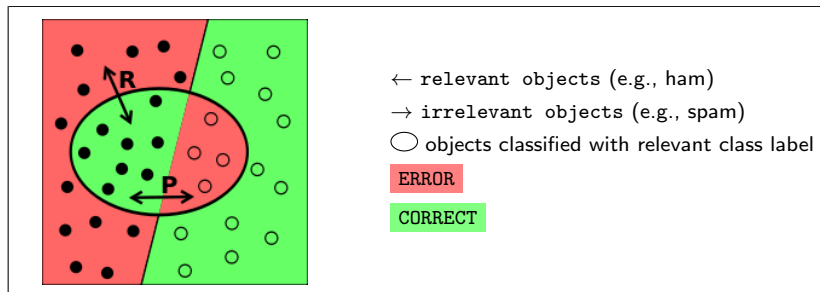
Train on a corpus: Bayesian inference on θ and ϕ

Train on a new documents d: fix $P(w | z)$ to infer $P(z | d)$

– Multiple inference algorithms available (expectation-maximization/VEM and Gibbs sampling/GIBBS)

interpretation and evaluation

evaluation|is our model valid?



Precision: fraction of retrieved instances that are relevant

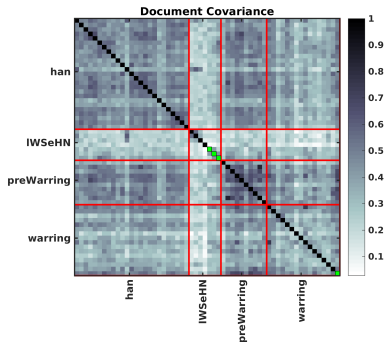
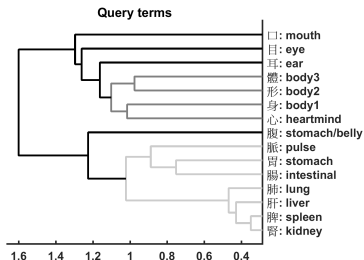
$$P = \frac{TP}{TP + FP}$$

Recall: fraction of relevant instances that are retrieved

$$R = \frac{TP}{TP + FN}$$

P and R are inversely related. Identify balance through a Precision-Recall curve.

interpretation|what does our model mean?

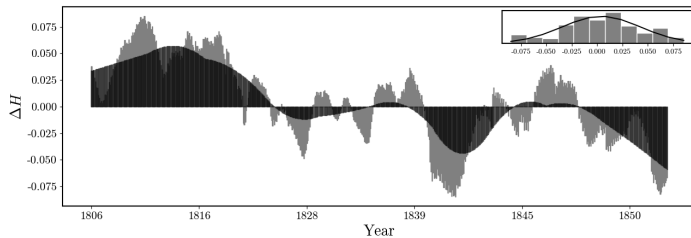
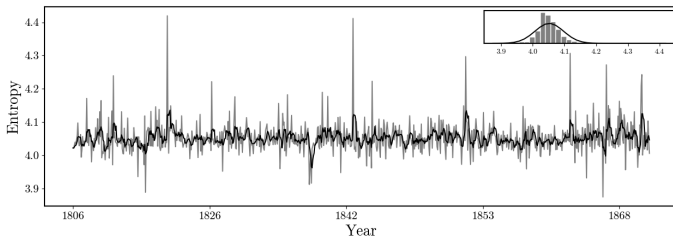


- philosophers and sinologists have been debating the existence of mind-body dualism in classical Chinese philosophy
- with domain experts, unsupervised learning was used to identify a multi-level dualistic semantic space
- one model (LDA) was further utilized to predict class of origin for controversial texts slices

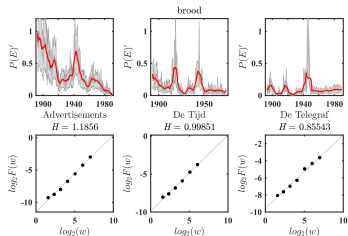
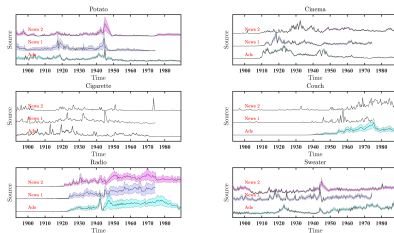
knowledge



YOUR
GAME
HERE



History|Predictive Causality & Slow Decay



- historians and media researchers theorize about the causal dependencies between public discourse and advertisement
- time series analysis of keyword frequencies (from seedlists) indicated that for some categories ‘ads shape society’, while other categories merely ‘reflect’
- advertisements show a faster decay (on-off intermittent behavior) than public discourse (long-range dependencies)