

Historical Informatics #1-2

@ SDU/Dept. of History

github.com/kln-courses/historical_informatics

Kristoffer L Nielbo

knielbo@sdu.dk

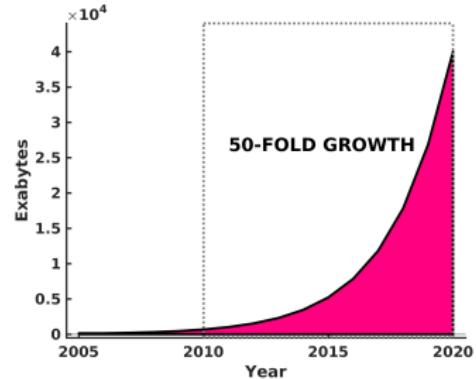
knielbo.github.io

Department of History | University of Southern Denmark

```
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()
```



A need for Human Informatics



the data deluge is transforming knowledge discovery and understanding in every domain of human inquiry

a large part of these data are soft and unstructured ⇒ to get value from these data, humanities (and social sciences) must utilize automation

human informatics - automatic information processing in the humanities

eScience Infrastructure at SDU



14.016 cores, 392 slim nodes, 64 fat nodes, 72 GPU nodes with two Nvidia K40/node

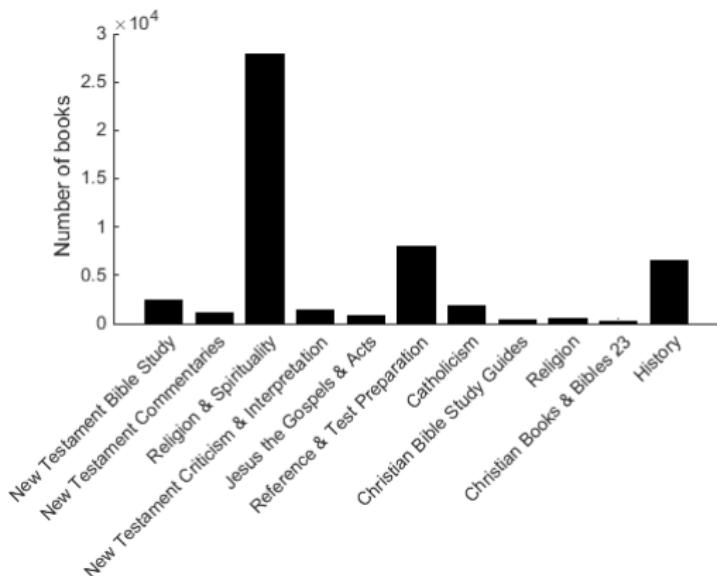
the center is expanding its userbase to include more SSH relevant areas and competencies, e.g., DL, NLP, computer vision

cloud offers a fully GDPR compliant environment for cloud computing with a cumulative codebase



-
- domain knowledge in history, language, literature &c combined with microscopic and (predominantly) qualitative analysis of human cultural manifestations

Gospel of Marc (KJV) ~ 16500 words in 16 chp. on 11 p.

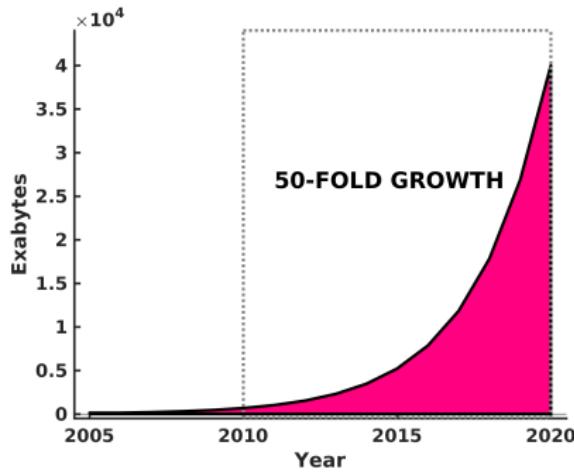


'from the dawn of civilization until 2003, humankind generated five exabytes of data.
Now we produce **five exabytes every two days** ... and the pace is accelerating'

Eric Smith (Google)

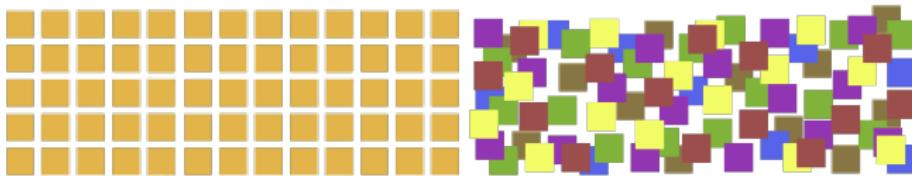
'increasingly, scientific breakthroughs will be powered by advanced computing
capabilities that help researchers manipulate and **explore massive datasets**'

Jim Gray (Fourth Paradigm)



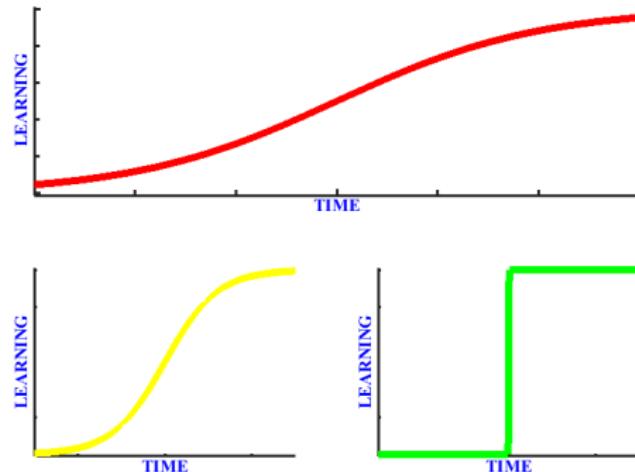
computational sciences are entering the exa-scale era
+
digital technologies are disruptive on a new scale

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
- increased demand for expertise from culture, media and linguistic domains

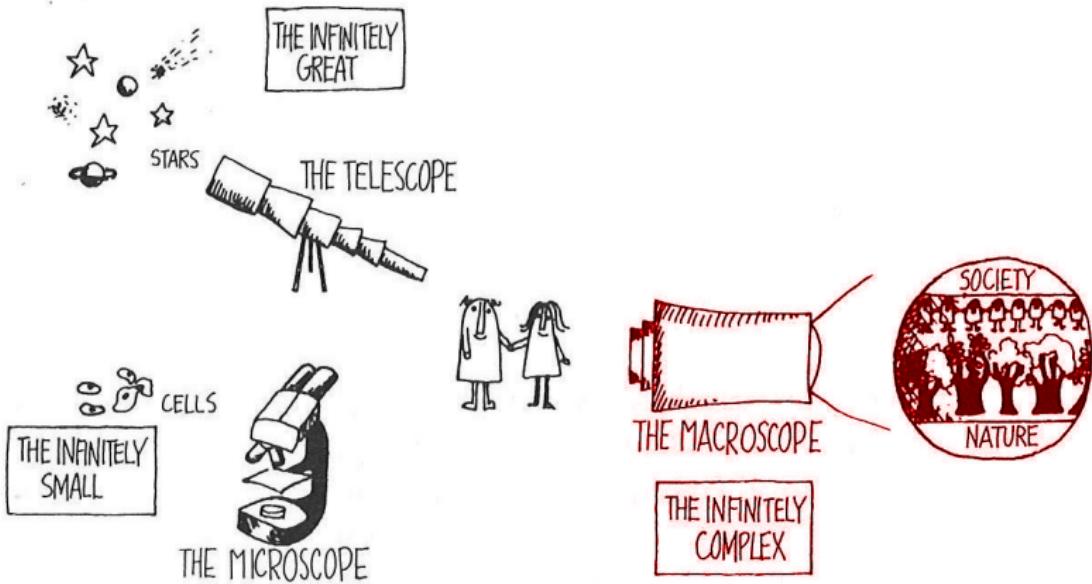


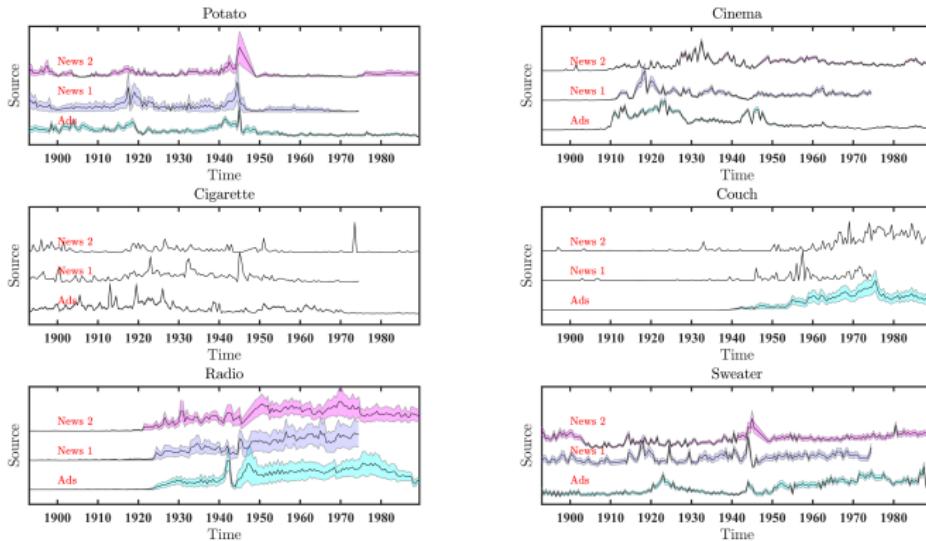
every knowledge-intensive industry have to “break” the learning curve



-
- domain knowledge in history, language, literature &c needs to scale, if we want to maintain our cultural knowledge base

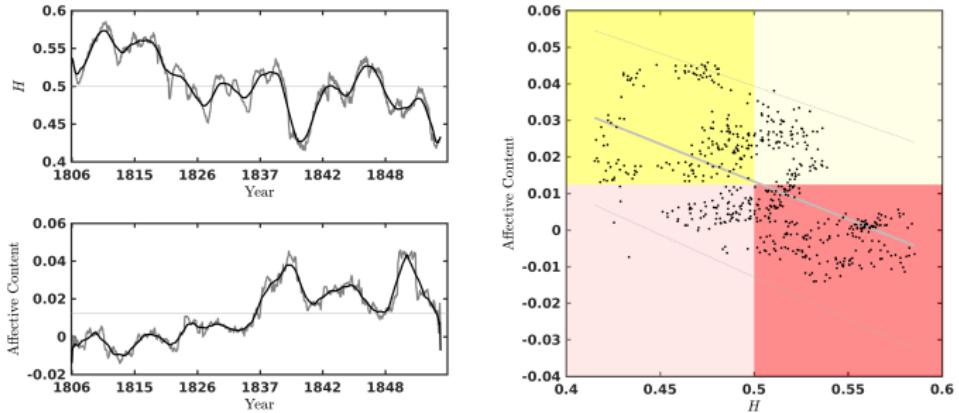
a (historical) macroscope





Digital history and media studies

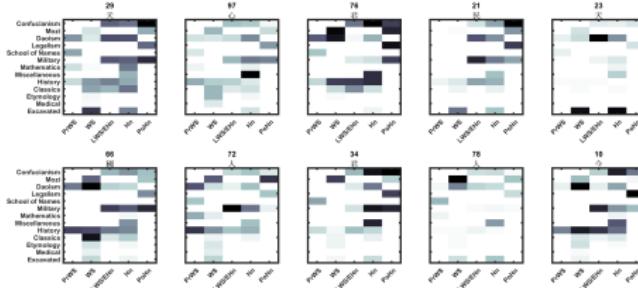
- prerequisite: humanistic domain experts that use content analysis
- source digitization (newspapers) og super computing change resolution and scale
- technologies create new standards for the domains involved
- share technology, but not data!



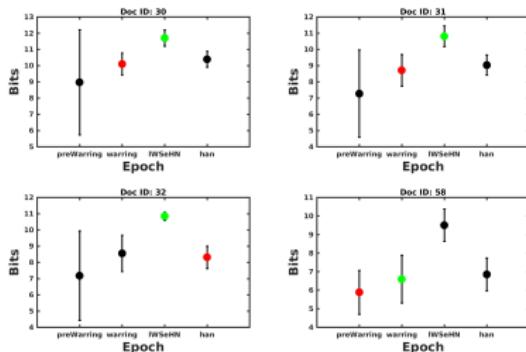
Computational literary history

- prerequisite: humanistic domain experts that study writers and literary periods
- high quality digitization of writers, annotation and NLP changes perspective and scale
- technologies that are creating new standards
- sharing of technology and data

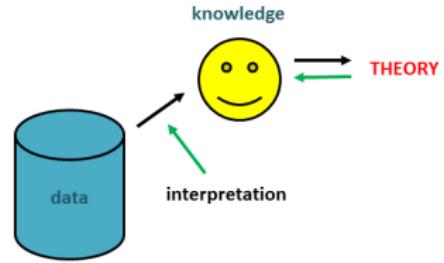
Dating Classical Chinese Documents

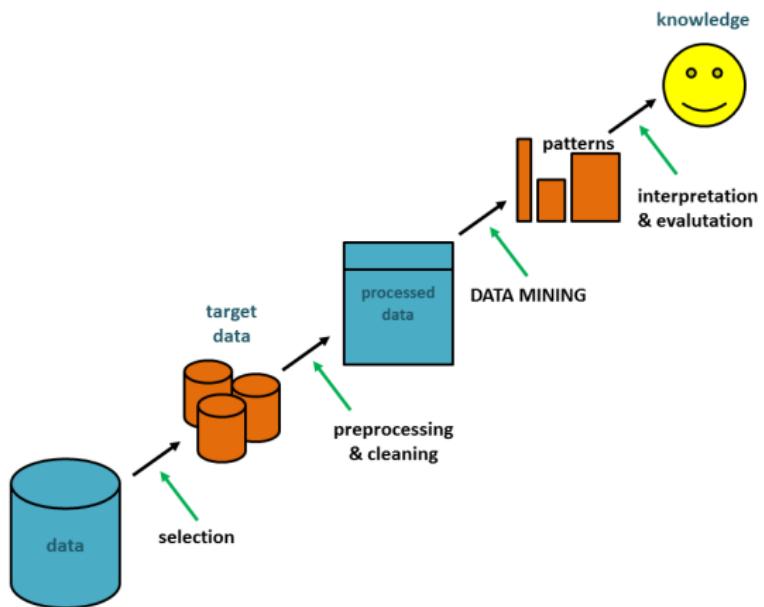


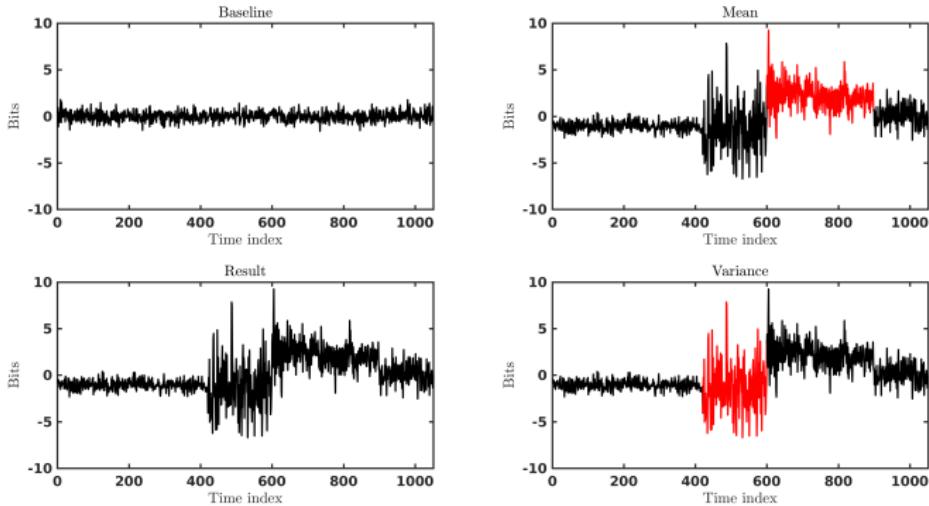
Thematic structure in CTEXT corpus conditioned on period and genre



Average distance from each period, **classical period**, and **alternate period**







Perspective: “Predictive History” ~ Culture Analytics

- given enough data, we can use past knowledge to predict future trends
- linked archives, news databases, social media ...
- knowledge of these technologies become imperative for critical use and assessment
- BUT we need free data access and mobility → OPEN SCIENCE

```
1 if questions:
2     try:
3         answer()
4     except RunTimeError:
5         pass
6     else:
7         print "break"
```

INTERVENTION|from the console

GUI → CLI

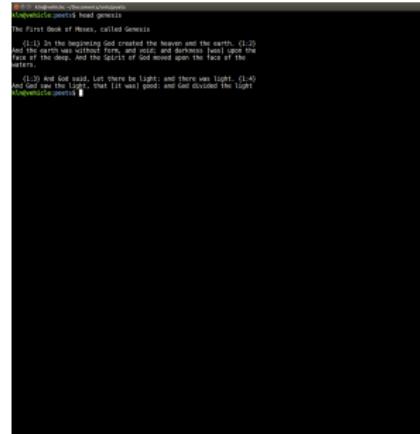
- novice-friendly visual approach to computer interaction w. a fast learning curve **ERROR**
- expert-friendly text-based approach to computer interaction w. ++freedom **VALID**
- **CONFLICT** break the learning curve through training intensive, non-intuitive, and specialized tools

Pannang Curry 1-2-3!

- ① Open your right refrigerator door and remove ingredients from the following locations: Door shelf 2. Spot 1; Crisper drawer 1, Spot 3; Crisper drawer 1, Spot 5.
- ② Open your third kitchen drawer and remove the utensils labeled "1", "3", "4", "9", and "12".
- ③ Use your arms to apply utensil 1 to ingredients 1-3. Place ingredients inside utensil 3.

*Note: This recipe uses ShaKL the Shared Kitchen Layout. To use ShaKL, you'll need to have installed ShaKL shelving, cabinetry, and utensils throughout your kitchen and pantry and have basic understanding of ShaKL managers. To learn more, read *Up and Running with ShaKL* (O'Billy Press, 2015). Want to improve ShaKL? Consider contributing to our team.

- a shell is a program whose primary purpose is to read commands and run other programs
- the shell's main advantages are its high action-to-keystroke ratio, its support for automating repetitive tasks, and its capacity to access networked machines
- the shell's main disadvantages are its primarily textual nature and how cryptic its commands and operation can be



A screenshot of a terminal window titled "Terminal - OpenShift Terminal". The command entered is "curl -s https://www.bible.com/bible/101/genesis.1.1-3". The output shows the beginning of the Book of Genesis in English:

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

`PS1='$ '` sets prompt string in console to

```
1 $
```

prompt indicates that the shell is waiting for input

```
1 $ whoami  
2 kln  
3 $
```

user ID or who the shell thinks you are

`whoami`

- ① finds a program called `whoami`
- ② runs that program
- ③ displays that program's output
- ④ displays a new prompt to tell us that it's ready for more commands

unknown command

```
1 $ somecommand  
2 somecommand: command not found  
3 $
```

- the shell runs other programs, so it does not work if the program does not exist

print working directory - current default directory

```
1 $ pwd  
2 home/kln  
3  
4 $ a=$(pwd)  
5 $ echo "current wd is: $a"  
6 current wd is: /home/knielbo
```

the path to the **home** directory varies between operating systems:

- [linux] /home/yourname
- [mac] /Users/yourname
- [windows] C:\Users\yourname

tokenization - unigrams

```
1 $ tr -sc "A-Za-z" "\n" < 2017-Trump.txt
```

sort in alphabetic order

```
1 $ tr -sc "A-Za-z" "\n" < genesis | sort
```

uniq - lexicon of document

```
1 $ tr -sc "A-Za-z" "\n" < 2017-Trump.txt | sort | uniq
2 $ tr -sc "A-Za-z" "\n" < 2017-Trump.txt | sort | uniq -c
3 $ tr -sc "A-Za-z" "\n" < 2017-Trump.txt | sort | uniq -c > lexicon.txt
4 ...
```

tired of cryptic commands and operations from the command line?

luckily we have:



```
1 >>> import gensim, nltk, polyglot, spacy
2 >>> from adl.util import thefunctionthatrulesthemall
3 >>> thefunctionthatrulesthemall("yourfile.dat")
```

and:



```
1 > libs b<- c("mallet", "tidyverse", "tm", "syuzhet")
2 > lapply(libs, require, character.only = TRUE)
3 > thefunctionthatrulesthemall("yourfile.dat")
```

```
1 if questions:
2     try:
3         answer()
4     except RunTimeError:
5         pass
6     else:
7         print "thank you"
```