# text analytics

DTL SU @ AU

kristoffer l nielbo
`kln@cas.au.dk`
`github.com/kln-courses/tmgu17`
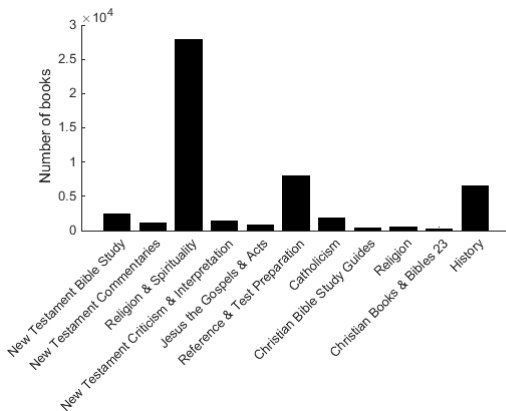`tmgu17.slack.com`

DAI|IMC|AARHUS UNIVERSITY

AARHUS UNIVERSITET

INTERACTING MINDS CENTRE

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

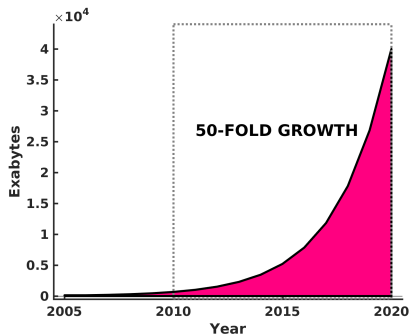# Gospel of Marc (KJV) $\sim$ 16500 words in 16 chp. on 11 p.

'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'

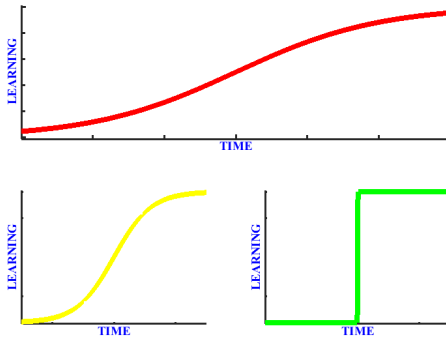Eric Smith (Google)

'increasingly, scientific breakthroughs will be powered by advanced computing capabilities that help researchers manipulate and explore massive datasets'
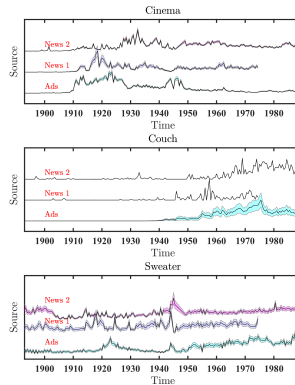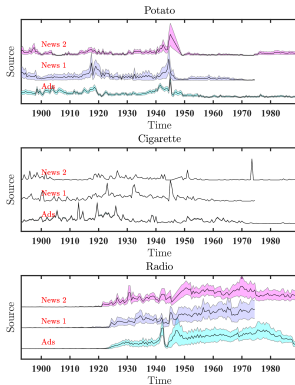
Jim Gray (Fourth Paradigm)

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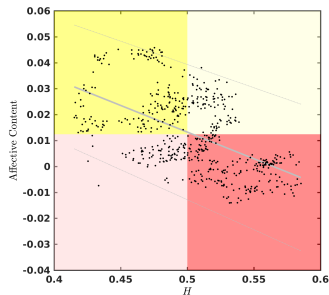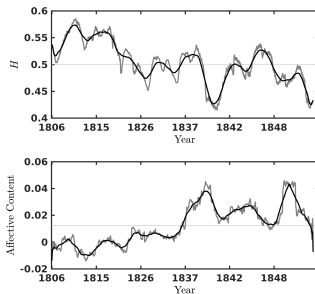+
digital technologies are disruptive on a new scale

every knowledge-intensive industry have to "break" the learning curve

**INTERVENTION|from the console**

———

**GUI → CLI**

- novice-friendly visual approach to computer interaction w. a fast learning curve <span style="color:red">ERROR</span>
- expert-friendly text-based approach to computer interaction w. ++freedom <span style="color:green">VALID</span>
- <span style="color:orange">CONFLICT</span> break the learning curve through training intensive, non-intuitive, and specialized tools

**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 tehcnology and data

**Data**  **Information**  **Presentation**  **Knowledge**

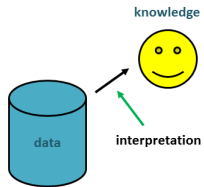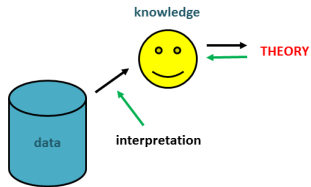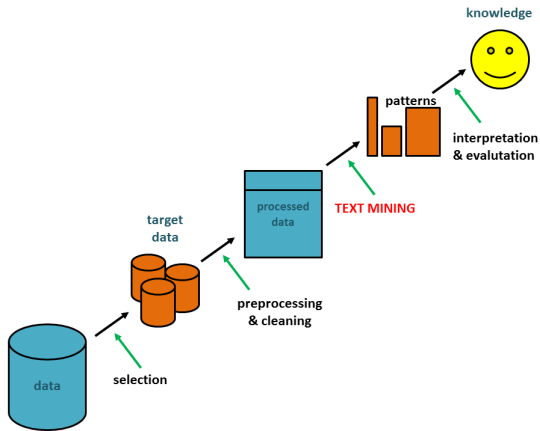Data       Information       Presentation       Knowledge

**Data**       **Information**       **Presentation**       **Knowledge**

knowledge

data

interpretation

knowledge

THEORY
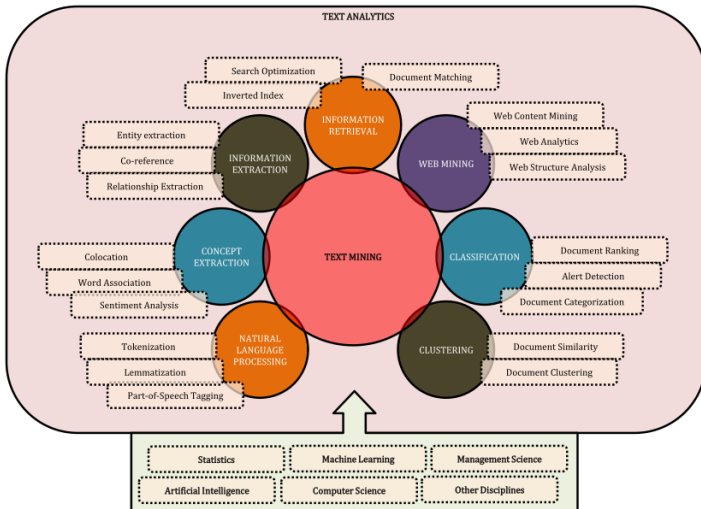
data

interpretation

knowledge

THEORY

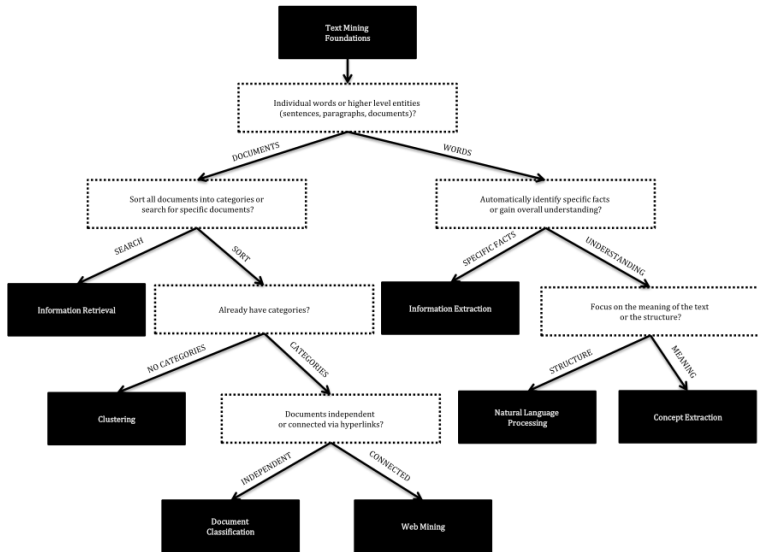**text analytics** $\sim$ text mining $\sim$ automated text analysis

set of data mining[1] techniques for extracting high quality information from large scale text-heavy (unstructured) data sets
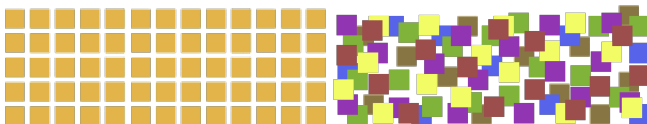
($\sim$ Miner et al 2012)

a tool for discovery and measurement in textual data of prevalent attitudes, concepts, or events

($\sim$ O'Connor, Bamman & Smith 2011)

[1] Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). From data mining to knowledge discovery in databases. AI Magazine, 17(3), 37.

Text Mining Foundations

Individual words or higher level entities (sentences, paragraphs, documents)?

DOCUMENTS — Sort all documents into categories or search for specific documents?

WORDS — Automatically identify specific facts or gain overall understanding?

SEARCH — Information Retrieval

SORT — Already have categories?

SPECIFIC FACTS — Information Extraction

UNDERSTANDING — Focus on the meaning of the text or the structure?

NO CATEGORIES — Clustering

CATEGORIES — Documents independent or connected via hyperlinks?

STRUCTURE — Natural Language Processing

MEANING — Concept Extraction

INDEPENDENT — Document Classification

CONNECTED — Web Mining

AARHUS UNIVERSITET

IMC INTERACTING MINDS CENTRE

**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

the goal of **statistical learning** is to build a machine that can learn from data and automatically make the right decisions

supervised learning infer mapping between data & class-information $\rightarrow$ 'ground truth'

unsupervised learning identify latent classes in the data $\rightarrow$ lack 'ground truth'

adequate problem solution requires that we test a range of approaches (algorithms, (hyper-)parameter estimation) - the validation of an approach is an **experiment**

experiment input: code, data sets, hyperparameter values

experiment output: model definition (weights), metric values (experiment comparison), execution logs

<div align="center">**a complex and error-prone process**</div>

$\Rightarrow$ systematically comment your work and process and use <span style="color:red">version control and source code management</span>

Voyant Tools 2.0 (Corpus View)

"There is no true interpretation of anything; interpretation is a vehicle in the service of human comprehension. The value of interpretation is in enabling others to fruitfully think about an idea"

Andreas Buja