

text analytics and generic tools

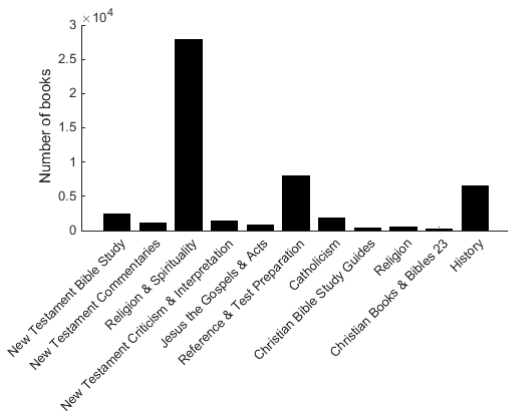
Kristoffer L. Nielbo

kln@cas.au.dk

DTL|Digital Arts Initiative
Interacting Minds Centre|Aarhus University



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

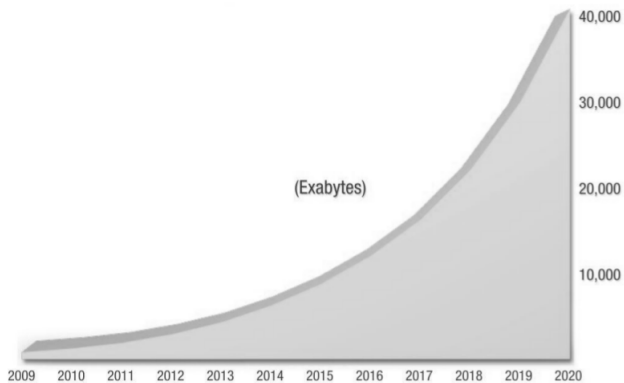


‘from the dawn of civilization until 2003, humankind generated five exabytes of data. Now we produce **five exabytes every two days** ... and the pace is accelerating’

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)



Data



Information



Presentation



Knowledge



Data



Information



Presentation



Knowledge



Data



Information

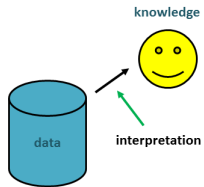


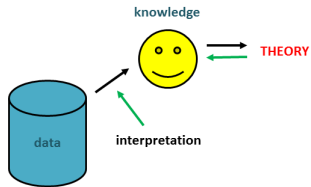
Presentation



Knowledge





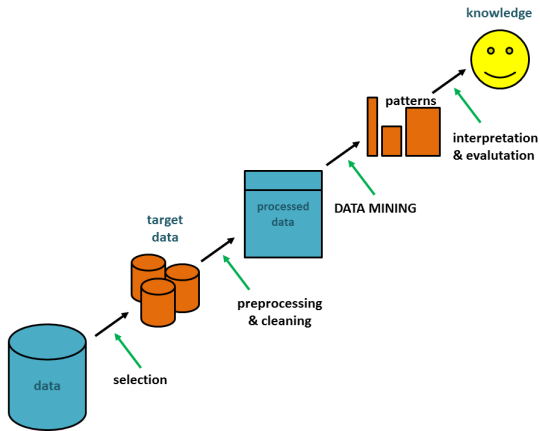


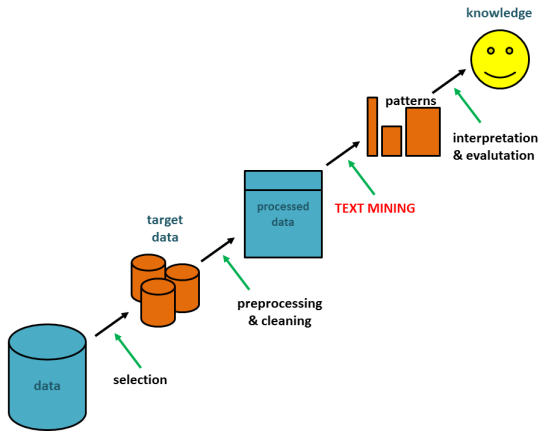
knowledge



THEORY







text analytics ~ text mining ~ automated text analysis

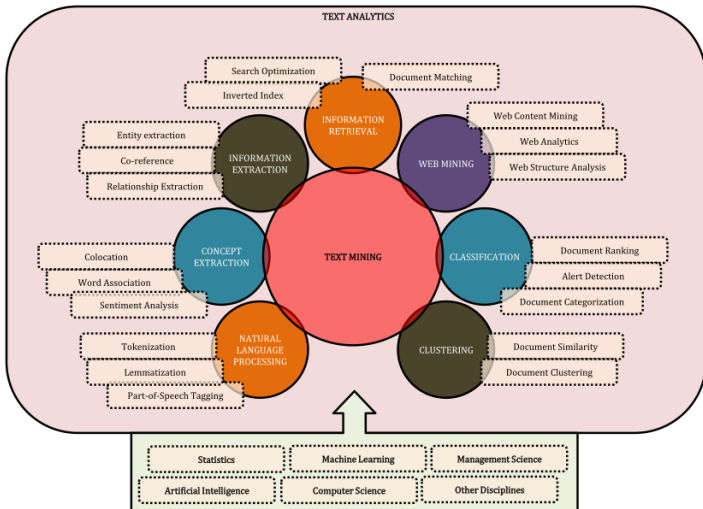
set of data mining¹ techniques for extracting high quality information from **large scale text-heavy** (unstructured) data sets

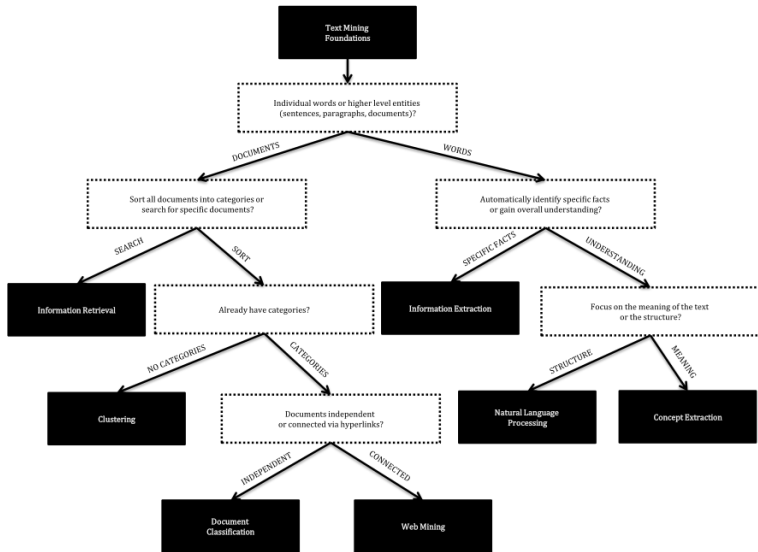
(~ Miner et al 2012)

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

(~ O'Connor, Bamman & Smith 2011)

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





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 → 'ground truth'

unsupervised learning identify latent classes in the data → lack 'ground truth'

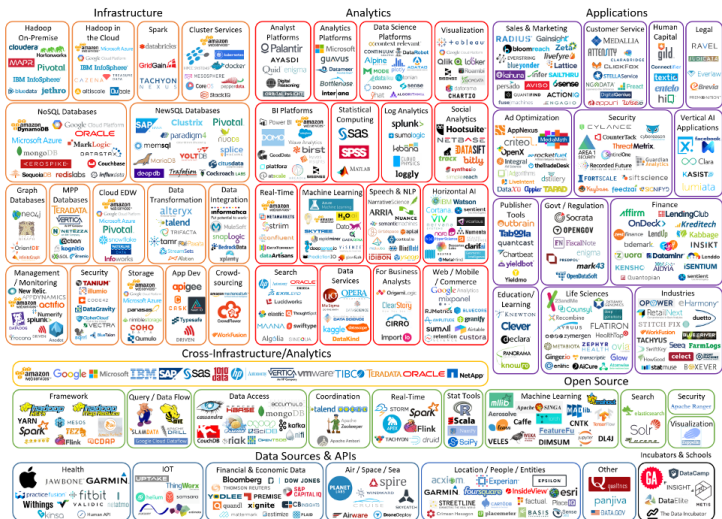
adequate problem solution requires that we test a range of approaches (algorithms, (hyper-)parameter estimation) - the validation of an approach is an **experiment**

experiment input: code, data sets, hyperparameter values

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

a complex and error-prone process

⇒ systematically comment your work and process and use **version control and source code management**









repository

Platform	Cost	Exclude	License
<input type="text" value="- Any -"/>	<input type="text" value="- Any -"/>	<input type="text" value="- Any -"/>	<input type="text" value="- Any -"/>

Research objects

Sort by Order

What kind of data should the tool work with?

BASE

BASE (Bielefeld Academic Search Engine) is a search engine for academic open access web resources that searches materials stored in OAI-PMH (Open Archives Initiative Protocol for Metadata Harvesting) enabled repositories.

Website: <https://www.base-search.net/>

Last updated: 19 Apr 2016

CONTENTdm

LANGUAGES

- English
- Español

“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