# Contemporary challenges in digital social science methodologies

Eetu Mäkelä

This presentation:
http://j.mp/meth4dss-td

# Li et al., 2014: <u>What a Nasty day: Exploring Mood-Weather Relationship from Twitter</u>
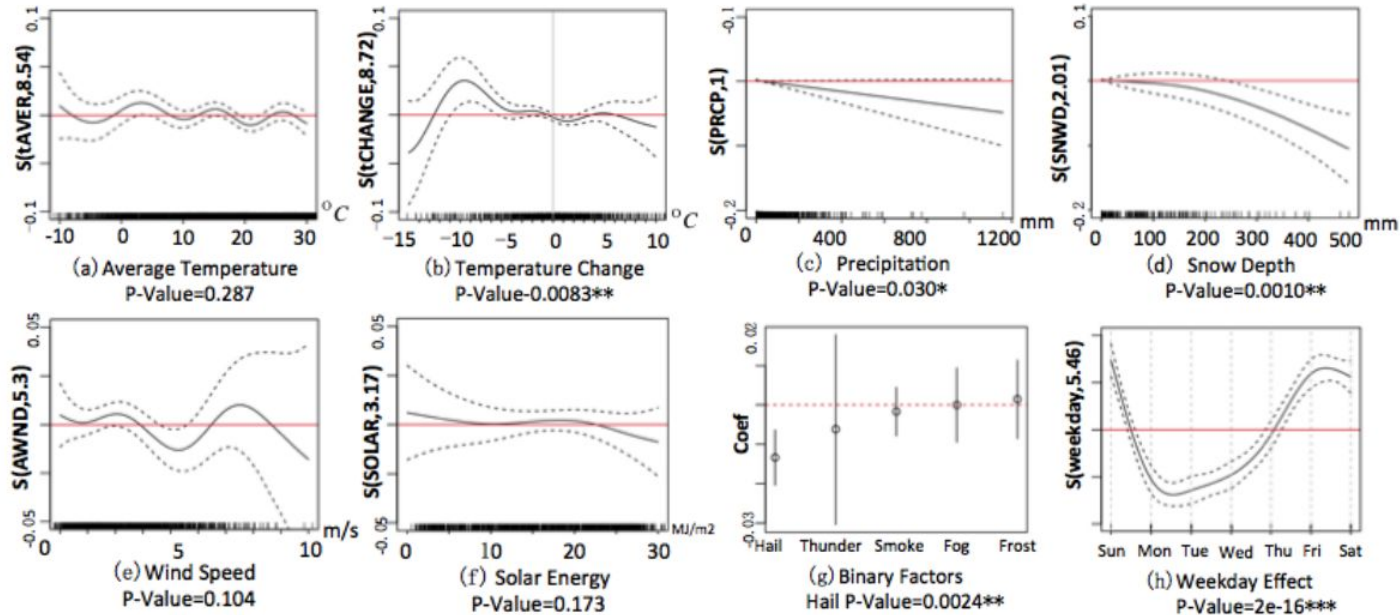


**Figure 3: Positive/Negative mode analysis regarding multiple meteorological factors. Red solid line corresponds to 0 line. Black dotted lines correspond to boundary of confidence interval. Black solid line corresponds to regression curve. y-axis corresponds to smooth regression value from GAM model. Positive value of smooth regression means positive contribution to up-mood state while negative value means the opposite. Label for y-axis corresponds to S(meteorological factor, degree of freedom)**

# Detecting influenza epidemics using search engine query data

Jeremy Ginsberg[1], Matthew H. Mohebbi[1], Rajan S. Patel[1], Lynnette Brammer[2], Mark S. Smolinski[1] & Larry Brilliant[1]

1. Google Inc., 1600 Amphitheatre Parkway, Mountain View, California 94043, USA
2. Centers for Disease Control and Prevention, 1600 Clifton Road, NE, Atlanta, Georgia 30333, USA

# Experimental evidence of massive-scale emotional contagion through social networks

Adam D. I. Kramer[a,1], Jamie E. Guillory[b,2], and Jeffrey T. Hancock[b,c]

[a]Core Data Science Team, Facebook, Inc., Menlo Park, CA 94025; and Departments of [b]Communication and [c]Information Science, Cornell University, Ithaca, NY 14853

## Significance

We show, via a massive ($N = 689,003$) experiment on Facebook, that emotional states can be transferred to others via emotional contagion, leading people to experience the same emotions without their awareness. We provide experimental evidence that emotional contagion occurs without direct interaction between people (exposure to a friend expressing an emotion is sufficient), and in the complete absence of nonverbal cues.

## PSYCHOLOGICAL AND COGNITIVE SCIENCES

PNAS is publishing an Editorial Expression of Concern regarding the following article: "Experimental evidence of massive-scale emotional contagion through social networks," by Adam D. I. Kramer, Jamie E. Guillory, and Jeffrey T. Hancock, which appeared in issue 24, June 17, 2014, of *Proc Natl Acad Sci USA* (111:8788–8790; first published June 2, 2014; 10.1073/pnas.1320040111). This paper represents an important and emerging area of social science research that needs to be approached with sensitivity and with vigilance regarding personal privacy issues.

Questions have been raised about the principles of informed consent and opportunity to opt out in connection with the research in this paper. The authors noted in their paper, "[The work] was consistent with Facebook's Data Use Policy, to which all users agree prior to creating an account on Facebook, constituting informed consent for this research." When the authors prep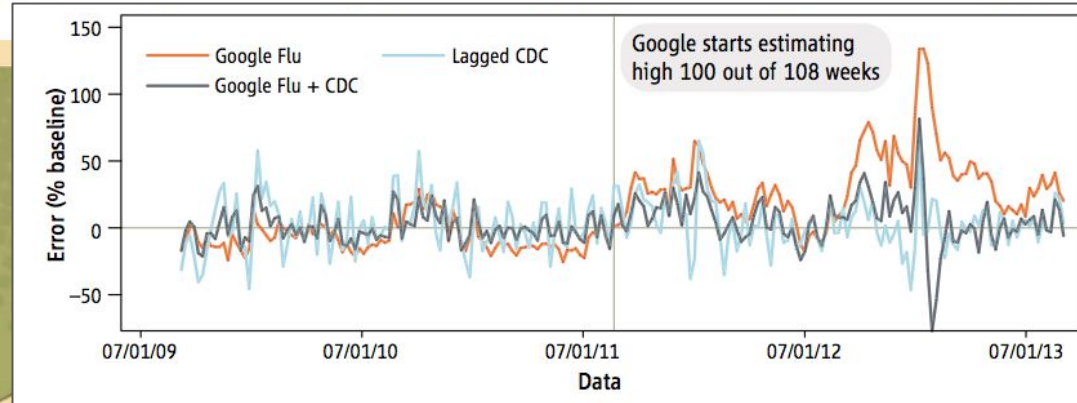ared their paper for publication in PNAS, they stated that: "Because this experiment was conducted by Facebook, Inc. for internal purposes, the Cornell University IRB [Institutional Review Board] determined that the project did not fall under Cornell's Human Research Protection Program." This statement has since been confirmed by Cornell University.

Science, March 2014

BIG DATA

# The Parable of Google Flu: Traps in Big Data Analysis

David Lazer,[1,2]* Ryan Kennedy,[1,3,4] Gary King,[3] Alessandro Vespignani[3,5,6]

Large errors in flu prediction were largely avoidable, which offers lessons for the use of big data.

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. *Nature* reported that GFT was predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance reports from laboratories across the United States (*1, 2*). This happened despite the fact that GFT was built to predict CDC reports. Given that GFT is often held up as an exemplary use of big data,
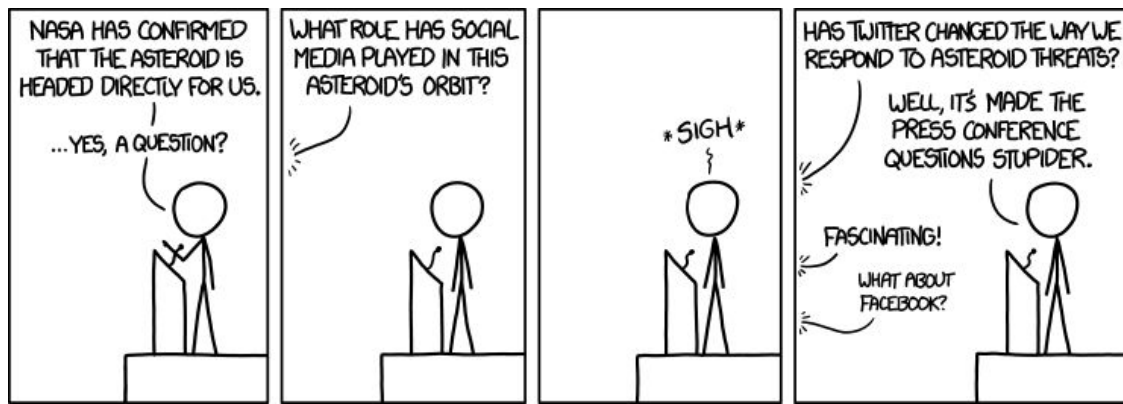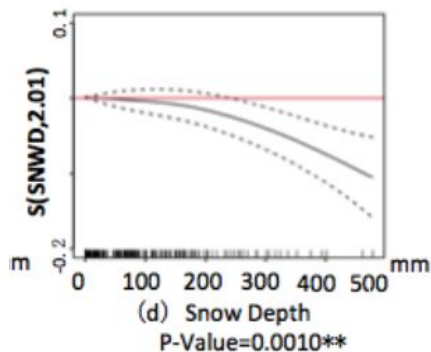


son and has missed high for 100 out of 108 weeks starting with August 2011 (see the graph). These errors are not randomly distributed. For example, last week's errors predict this week's errors (temporal auto-correlation), and the direction and

# Cihon & Yasseri, 2016: A Biased Review of Biases in Twitter Studies on Political Collective Action

This literature offers insight into particular social phenomena on Twitter, but often **fails to** use standardized methods that **permit interpretation beyond individual studies**. Moreover, the literature **fails to ground methodologies** and results in social or political theory, **divorcing empirical research from the theory needed to interpret it**. Rather, investigations focus primarily on methodological innovations for social media analyses, but these too often fail to sufficiently demonstrate the validity of such methodologies.



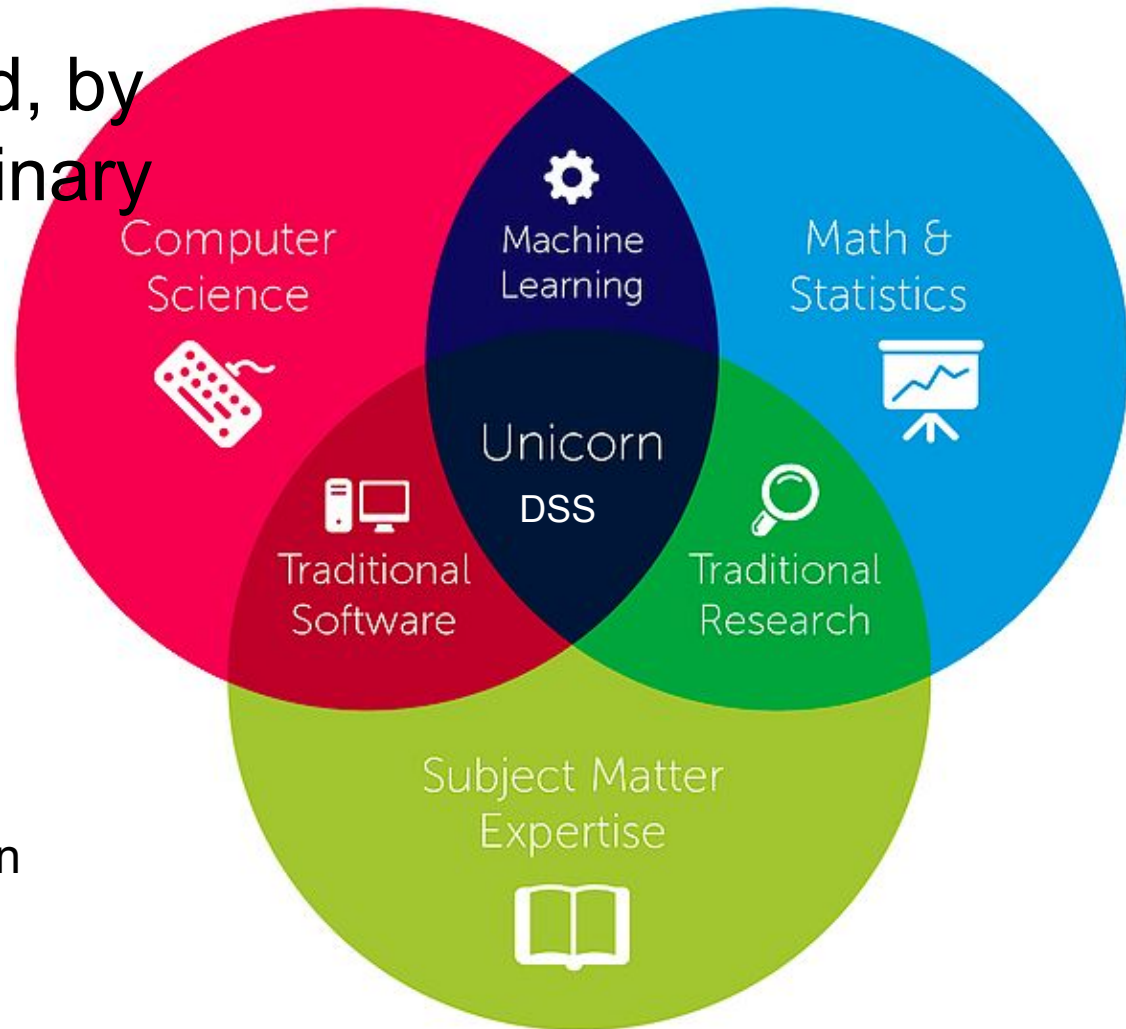(d) Snow Depth
P-Value=0.0010**

# Why does this happen?

# DSS is complex, hard, by necessity interdisciplinary

- Data is big, complex and inaccessible
- CS needed to access, process and explore it
- Knowledge of statistics needed to make reliable conclusions
- Social science subject expertise needed to ground results, provide interpretation and ensure depth

# DSS is being done without social scientists!

A final challenge for computational social science is that, in spite of many thousands of papers published on topics related to social networks, financial crises, crowdsourcing, influence and adoption, group formation, and so on, **relatively few are published in traditional social science journals or even attempt to engage seriously with social scientific literature**. The result is that much of **computational social science has effectively evolved in isolation** from the rest of social science, largely ignoring much of what social scientists have to say about the same topics, and largely being ignored by them in return.

Duncan J. Watts (Microsoft Research): Computational Social Science: Exciting Progress and Future Directions. The Bridge on **Frontiers of Engineering**, December 20, 2013, Volume 43, Issue 4

Regular Article

# Manifesto of computational social science

R. Conte[1,a], N. Gilbert[2], G. Bonelli[1], C. Cioffi-Revilla[3], G. Deffuant[4], J. Kertesz[5],
V. Loreto[6], S. Moat[7], J.-P. Nadal[8], A. Sanchez[9], A. Nowak[10], A. Flache[11],
M. San Miguel[12], and D. Helbing[13]

[1] ISTC-CNR, Italy
[2] CRESS, University of Surrey, UK
[3] Center for Social Complexity, George Mason University, USA
[4] National Research Institute of Science and Technology for Environment and Agriculture
    (IRSTEA), France
[5] Institute of Physics, Budapest University of Technology and Economics, Hungary
[6] Sapienza University of Rome, Italy
[7] University College London, UK
[8] CNRS, France
[9] GISC, Carlos III University of Madrid, Spain
[10] Center for Complex Systems, University of Warsaw, Poland
[11] ICS, University of Groningen, The Netherlands
[12] IFISC (CSIC-UIB), Campus Universitat Illes Balears, 07071 Palma de Mallorca, Spain
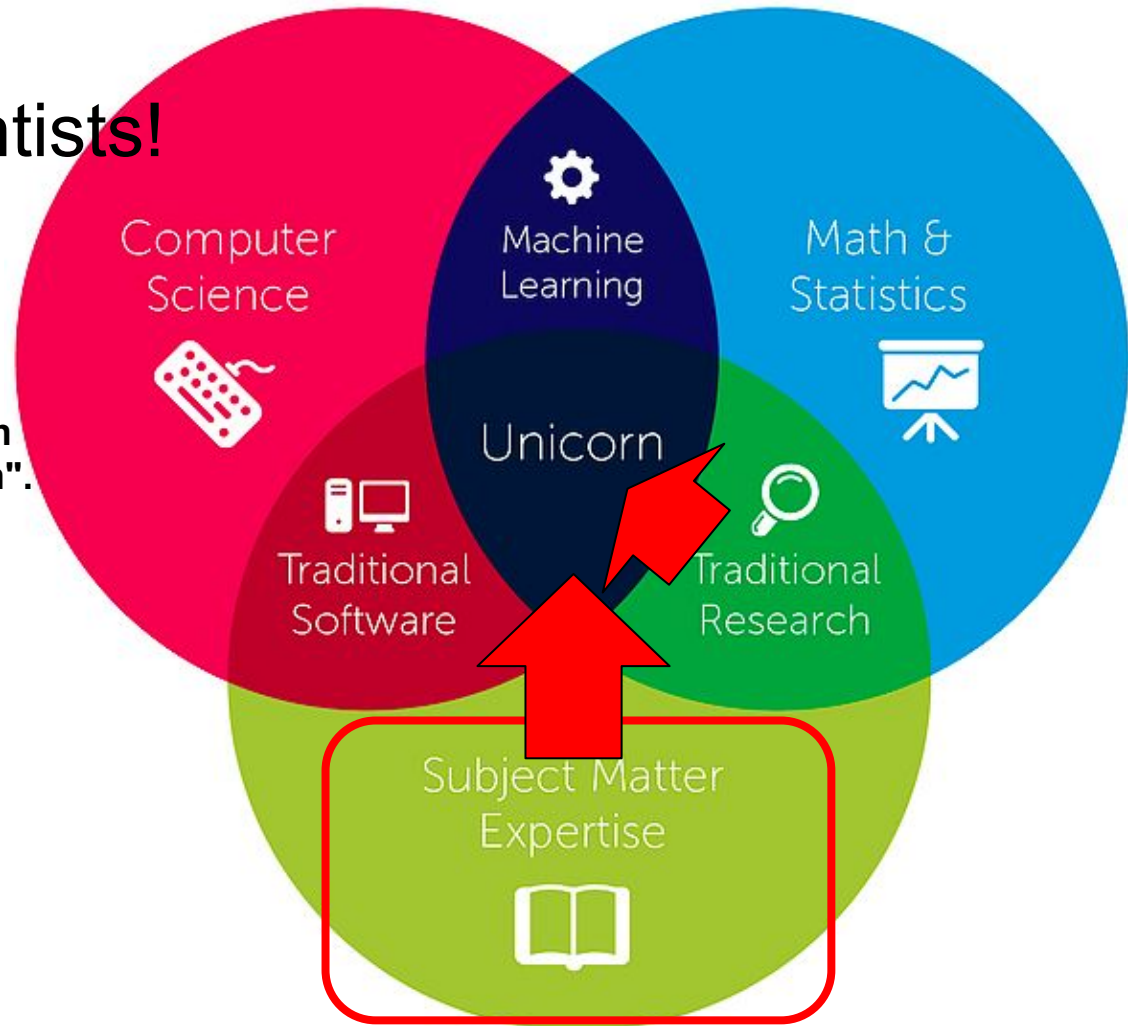[13] ETH Zurich, Switzerland

# AT THE CROSSROADS: LESSONS AND CHALLENGES IN COMPUTATIONAL SOCIAL SCIENCE

EDITED BY : Javier Borge-Holthoefer, Yamir Moreno and Taha Yasseri

# Niche for social scientists!

**"I have the solution, but it works only in the case of spherical cows in a vacuum".**

# And they know they need you!

Olemme **fyysikkotaustaisia** Aalto-yliopiston tutkijoita tekemässä hakemusta MATINE:lle koskien aatteiden ja ideologioiden muodostumista ja kehittymistä agenttipohjaisissa simulaatioissa, ja **etsimme hakemukseen halukkaita yhteistyökumppaneita sosiaalitieteiden puolelta**. Lähestymme tutkimusaihettamme sen oletuksen kautta, että ihmisten pääasiallisena viettinä on maksimoida oma "paremmuutensa" sosiaalisessa ympäristössään. Tämä viitekehys on lähellä Adlerin yksilöpsykologian koulukunnan perusajatuksia, ja siinä ideologioita voidaan kuvata tapoina laittaa asiat ja ihmiset arvojärjestyksiin.

Hakemus on jätettävä viimeistään 14.6.2017, joten toivomme yhteistyötarjouksia mahdollisimman pian, ja pahoittelemme tiukasta aikataulusta mahdollisesti aiheutuvaa vaivaa.
Yhteystiedot: Prof. Kimmo Kaski, kimmo.kaski'at'aalto.fi, FT Jan Snellman, jan.snellman'at'aalto.fi

# What to learn?

1. Knowledge of easy to use end-user data processing and exploration tools
   - Easy to use for their intended purpose, but limited
2. Knowledge of the fundamentals concepts of programming
   - Frees you to process your data more efficiently
   - Allows you to more freely apply analyses etc based on ready libraries and tutorials on the Internet
3. High-level understanding of what types of things can be accomplished with advanced CS methods
   - To be able to communicate in **collaborative projects**

For computer scientists, DSS offers:
- complex, meaningful challenges
- both in terms of data as well as use cases

# Longer term = HELDIG

Deep and significant progress in social science, in other words, will require not only new data and methods but also **new institutions** that are designed from the ground up to foster **long-term, large-scale, multidisciplinary, multimethod, problem-oriented social science research**. To succeed, such an institution will **require substantial investment**, on a par with existing institutes for mind, brain, and behavior, genomics, or cancer, as well as the active cooperation of industry and government partners.

Duncan J. Watts (Microsoft Research): Computational Social Science: Exciting Progress and Future Directions. The Bridge on Frontiers of Engineering. December 20, 2013, Volume 43, Issue 4

## What kinds of data do you want us to use for examples?

Cookbooks

New sources of data (eg. Facebook feeds)

The dead and the dying

linguistic data

Tools for transcribing speech

Public documents

Music

Text documents

Skeletons

Musical instruments

Maps

archaeological grave databases etc?

## METH4DH - what are you interested in?

Translating tools

Tools for analyzing the use & meaning of concepts

## What should the methods make possible?

Analyze video games

Quantitative methods through coding and digitalization

Quantitative methods as a supporting toool for historical research

Distant reading for history studies

Methods of investigating musics

How to apply these methods to social sciences

What kind of tools and methods there are?

how to visualize data

sorting through audio

Methods for archaeology

how to manage resources

For a beginner - give an overview of what is possible already (I feel I cannot imagine it yet)

search through a large amount of articles

will definitely do!

Understand pitfalls of dh tools

Ability to analyze larger amounts of data reduces the risk that something goes unnoticed

Digital tools may increase the efficiency of research

Thanks! Will make me feel less lost ;)

important part of the course

this is what we're after, yes..

# METH4DH background questionnaire

## Pertinent background information

If you want to tell us more deeply about your study subject or interests, as they relate to the course

Your answer

## Why are you taking this course?

Your answer

## What would you especially like to learn during this course / where would you like us to focus on?

Your answer

SUBMIT

This presentation:
http://j.mp/meth4dss-td

Unused slides follow →

# Challenge 1 - access to data

- One of the biggest problems cited by researchers doing big data research was **getting access to commercial or proprietary data**, suggesting that more needs to be done to unlock data sets for social science research.



**Figure 10** Data types used by respondents in most recent research involving big data ($n = 3077$)

| Data type | Count |
| --- | --- |
| Administrative data | 1690 |
| Commercial or proprietary data | 697 |
| Other social media | 533 |
| Photographs, video, or audio | 515 |
| Facebook | 460 |
| Twitter | 358 |
| Sensor data | 299 |
| Survey data | 228 |
| Mobile data | 221 |
| Medical/scientific data | 70 |
| Media/press | 44 |
| Bibliographical data | 34 |
| Census | 24 |

Percent of cases

# Challenge 2 - complexity of data

- In the social sciences, the new sources of data … derive overwhelmingly from mixed sources (e.g., social media, unstructured text, digital sensors, financial and administrative transactions) **not designed to produce valid and reliable data** for social scientific analysis (Lazer, Kennedy, King, & Vespignani, 2014), resulting in the **challenge of harmonizing and extracting meaningful features**
- …, **social scientific "big data" are notable less for absolute size per se than for the complexity** that renders conventional methods inadequate (Doorn, 2014).



● Aloituksen jälkeisiä kommentteja    ● Peräkkäisten päivien kommentti-ID:n erotus

# Challenge 3 - complex methods

- Our survey respondents listed finding **collaborators with the right skills** and the **amount of time required to learn** a new field as the biggest barriers to entry.
- A characteristic of researchers doing big data research is that they are more likely to **collaborate with other academics** (79 percent of big data researchers in our survey). Considering that **a large number of social science papers are single authored** (about 40 percent, according to Thomson Reuters (King, 2013), this information is significant.



Computer Science

Math & Statistics

Machine Learning

Unicorn

Traditional Software

Traditional Research

Subject Matter Expertise

**Constrain**

Search

[kc]an* no?t??

| | |
|---|---|
| can not | 52/612 |
| canne nott | 2 |
| can not | 2/27 |
| kan not | 2/7 |
| canne not | 2/9 |
| can note | 0/1 |
| can nott | 0/1 |
| can nowt | 0/1 |
| cane nott | 0/1 |
| kan notte | 0/1 |
| canst not | 0/4 |
| can not | 0/9 |

Education ♦

| | |
|---|---|
| Higher | 3/564 |
| Higher (Foreign) | 0/43 |
| Higher (Inns of Court) | 1/152 |
| Higher (Oxford) | 1/216 |
| Higher (Cambridge) | 2/226 |
| Apprenticed | 0/64 |
| Secondary | 0/30 |
| Private/Self: Classical | 0/48 |
| Elementary | 0/1 |

**Filter** ↔

Context

CLEMENT PASTON to JOHN I PASTON on 25 August, 1461

<Q A 1461 FN CPASTON> <X CLEMENT PASTON> <P I,199> {} {\116. TO JOHN PASTON I 1461, 25 AUGUST\] }] {^To hijs rythe reuerent and worchypfwll broder John Paston.^} Rythe reuerent and worchypfwll broder, I recomawnde me to +gowre good broderhood, desieryng to herre of +goure welfare and good prosperite`, the qwyche I pray God encresse to his pleswre and +gowre hertys hesse; certyfyyng +gow +tat I haue spok wyth John Rwsse, and Playter spak wyth him bothe, on Fryday be-fore Seynt ...ce sum were, +tat +tan +ge wold haue hym hom, +te qwyche xwld cause hym not to be hadde in fauore; and also men wold thynke +tat he were put owte of seruice. Also W. Pekoln tellythe me +tat hijs mony is spent, and not ryotesly but wysly and discretly, fore +te costys is gretter in +te Kyngys howse qwen he rydythe +tan +ge wend it hadde be, as Wyllam Pekok can tell +gow. And +tere wee mwst gett hym i c s. at +te lest, as by Wyllam Pekolys seyyng, and +get +tat will be to lityll. And I wot well we kan not get xl d. of Cristofyre Hanswm, so I xall be fayn to lend it him of myn owne siluer. If I knew verily +gour entent were +tat he xwld cum hom I wold send hym non. There I wyll doo as me thynkithe +ge xwld be best plesyd, and +tat, me thynkythe, is to send him +te siluer. +Tere-fore I pray +gow as hastely as +ge may send me a-+gen v mark, and +te remnawnte I trow I xall gett vp-on Cristofire Hanswm and Lwket. I pray +gow send me it as hastely as +ge may, fore I xall leue my-selfe rythe bare; and I pray +gow send me a letter how +ge woll +tat he xall be demenyd. Wrytyn on Twsday after Seynt Barthelmwe, &c. (\Christus vos obseruet.\} by Cle[|ment Paston[]

**View** ↔

Partition:

gender ♦

Sample length (years):

20

Graph type:

individual | compare as area | compare as scatterplot

☐ Hide totals   ☐ Show as accumulation chart
☐ Partitions as charts

Calculate values as:   total | text average | author average

http://purl.org/linked-data/sdmx/2009/code#sex-F

http://ldf.fi/ceec/person_BELIZABETH

39

97

percentage of cannot   Lin ▾

1659

Doctor
Patient

Time

Doctor
Patient

Time

ideas

people

animals

doing

Utterance Pair: 2, 12
Channel: Peter
Similarity: 0.284
Similar Concepts: doing, life, whole, people

Are you happy (SPause)? (GiveTime)

Well, I suppose I am

Yeah (VblAck)

To a point

Yeah (VblAck)

Mmm (affirm)

You got friends in this- in the, *hostel*

*Here?*

Yeah

No

No? (VblAck)

No

1. Conversation lacks conceptual richness, and consequently we see little engagement between CS and PWD manifested as white space under the diagonal

No? (VblAck) But you go walking a lot, don't you (SPause)? (PWDKnowl) (GiveTime)

Well I used to go walking a lot

2. Use of PWDKnowl strategy gives rise to engagement between CS and PWD.

I've seen you walking a lot here (PWDKnowl)

Mmm (affirm)

Around the place

Oh yes you walk wherever you have to

Yeah (VblAck)

# International Conference on Computational Social Science Luminaries

- Santo Fortunato is **Professor of Complex Systems** at the Department of Biomedical Engineering and Computational Science
- Lada Adamic is a **computational social scientist** at Facebook and previously an **associate professor at the School of Information and the Center for the Study of Complex Systems**
- Albert-László Barabási directs the **Center for Complex Network Research**, and holds appointments in the **Departments of Physics and College of Computer and Information Science**

- Nicholas Christakis MD, PhD, MPH, is a **social scientist and physician**
- Alessandro Vespignani is the Sternberg Distinguished **Professor of Physics, Computer Science and Health Sciences**
- Dirk Helbing is **Professor of Sociology**, in particular of Modeling and Simulation, at the Department of Humanities, Social and Political Sciences and member of the Computer Science Department at ETH Zurich. He earned a **PhD in physics**…

# Indaco & Manovich, 2016: Urban Social Media Inequality: Definition, Measurements, and Application

# Indaco & Manovich, 2016: Urban Social Media Inequality: Definition, Measurements, and Application



- Social media inequality of visitors' images in Manhattan (Gini = 0.669) is larger than income inequality of most unequal country in the world (Seychelles where Gini = 0.658).
- On the other hand, social media shared by locals has a Gini coefficient similar to countries that rank between 25 and 30 in the list of countries by income inequality. These are countries like Costa Rica (0.486), Mexico (0.481) and Ecuador (0.466). (The World Bank, 2015).

Since Instagram did not support downloading large volumes of historical data, we had to download data and images continuously during the period we wanted to cover. A single iMac computer running 24/7 continuously was used for downloading this data.

# Solutions to data issues

- Be at Facebook
- Do local stuff
- Make the peculiarity of the data an asset, a part of the research
- Be opportunistic

# Research process

1. Have data
2. Magic (?)
3. Something interesting shows up
4. Profit!

# Research process

1. Have data
2. Magic (?)
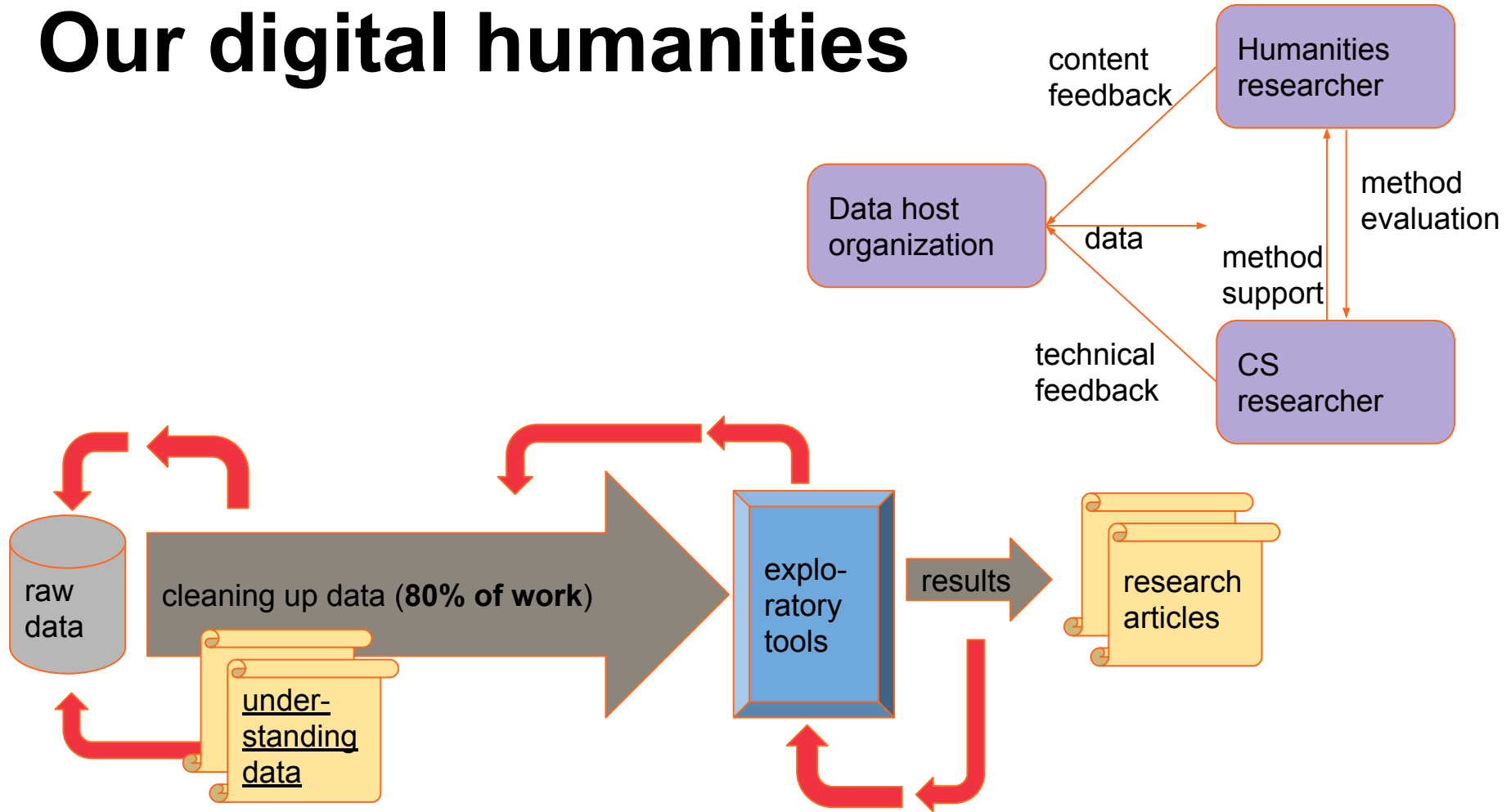3. Something interesting shows up
4. Profit!

*"Any sufficiently advanced technology is indistinguishable from magic."*

-   Arthur C. Clarke

# Research process

1. Have data
2. Magic (?)
   a. Hedge magic (spreadsheets, Excel graphs)
   b. Common ritual magic (statistics: correlation, ANOVA, PCA)
      ■ Relatively simple, commonly understood formulae you could mostly go through with pen and paper if you wanted to
   c. Higher ritual magic (SVM, LSA, LDA, SnE)
      ■ More complex, harder to follow formulae, impossible to work through manually
      ■ Well-grounded black box oracles (e.g. you feed a machine learning algorithm stuff, it processes it based on complex but well-defined rules, out comes results)
   d. Black magic (Deep learning)
      ■ True black box oracles (you feed a neural network both an input and a desired output, it derives mostly unintelligible black box rules that link the two)
   e. Flashy magic (proper visualizations)
3. Something interesting shows up
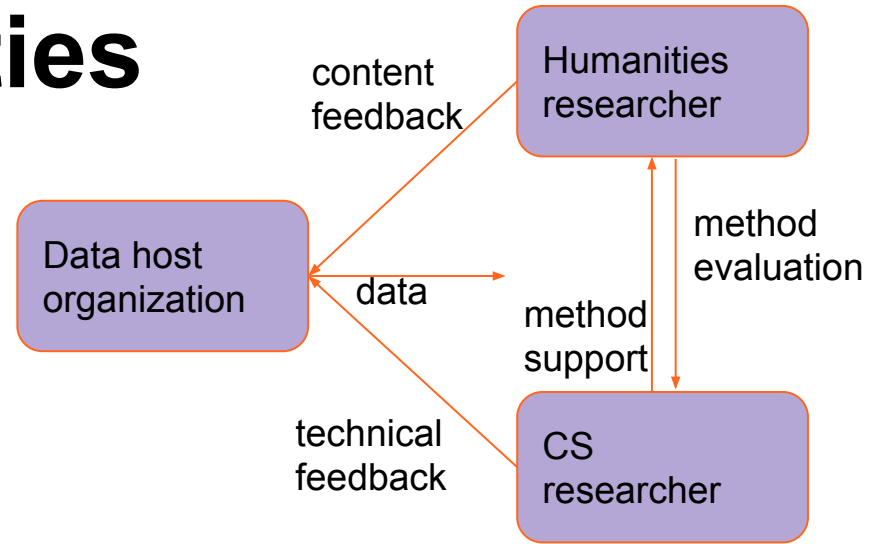
# Our digital humanities

# Our digital humanities

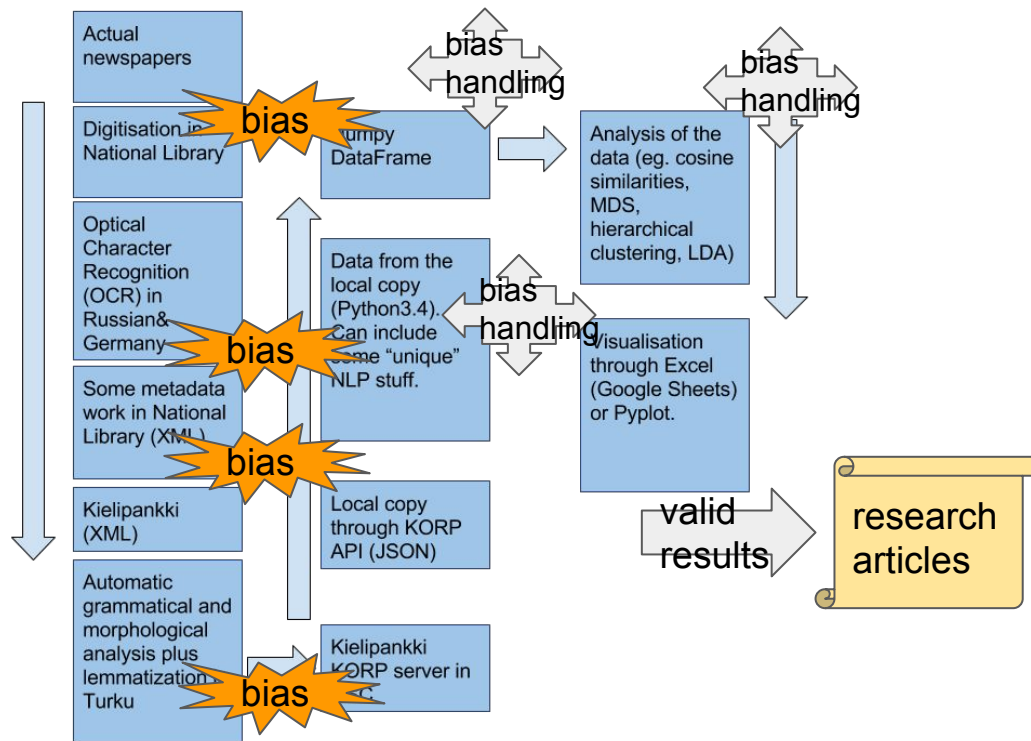At its best, such close collaboration offers **benefits for everyone involved**

- scholars in the humanities are able to tackle **questions too labour-intensive for manual study**
- computer scientists encounter **new and challenging use cases** for the tools and algorithms they develop
- data providers **gain insigh**t into their own data

# Don't get carried away by fancy methods!

1. Your dataset must be applicable to the methods you choose. Complex methods often make presuppositions about the data they apply to - if you don't understand these deeply, you'll end up with invalid results
2. In typical DH research, 90% of your time will go to gathering and understanding the data and transforming it into a form you can use - using complex methods, another 90% of your time may go to altering them to fit your data, and it'll run out
3. Complex methods are often unnecessary for DH work. On the contrary, often simpler methods are actually better.

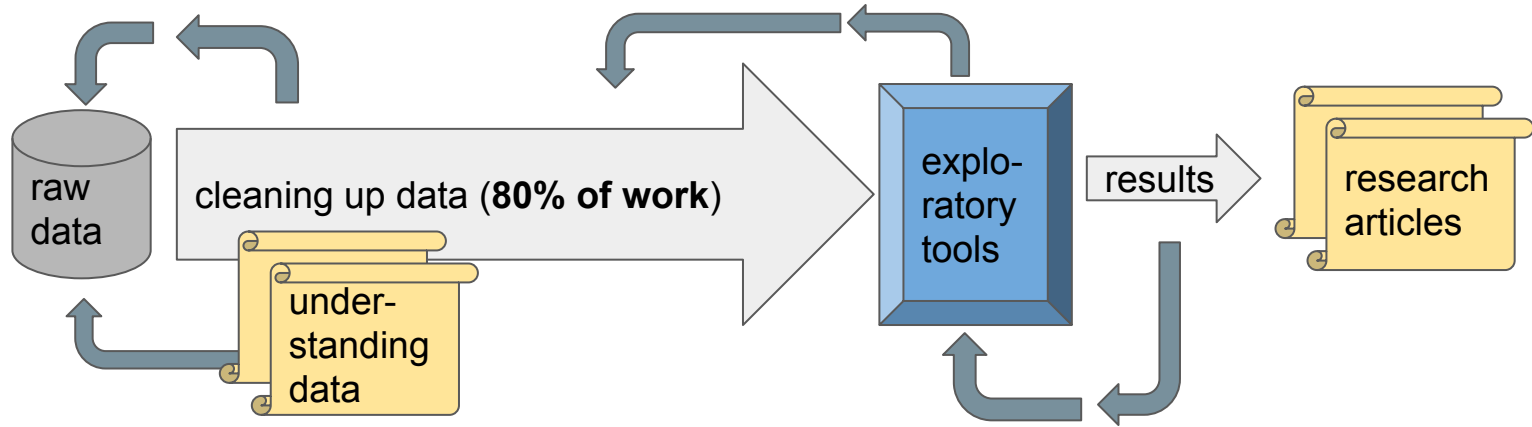# KLK Newspaper Pipeline: from archives to a hypothetical researcher

# Our digital humanities

- Scholars in the humanities and computer sciences collaborating, **applying novel computer science** to solve **humanities research questions**
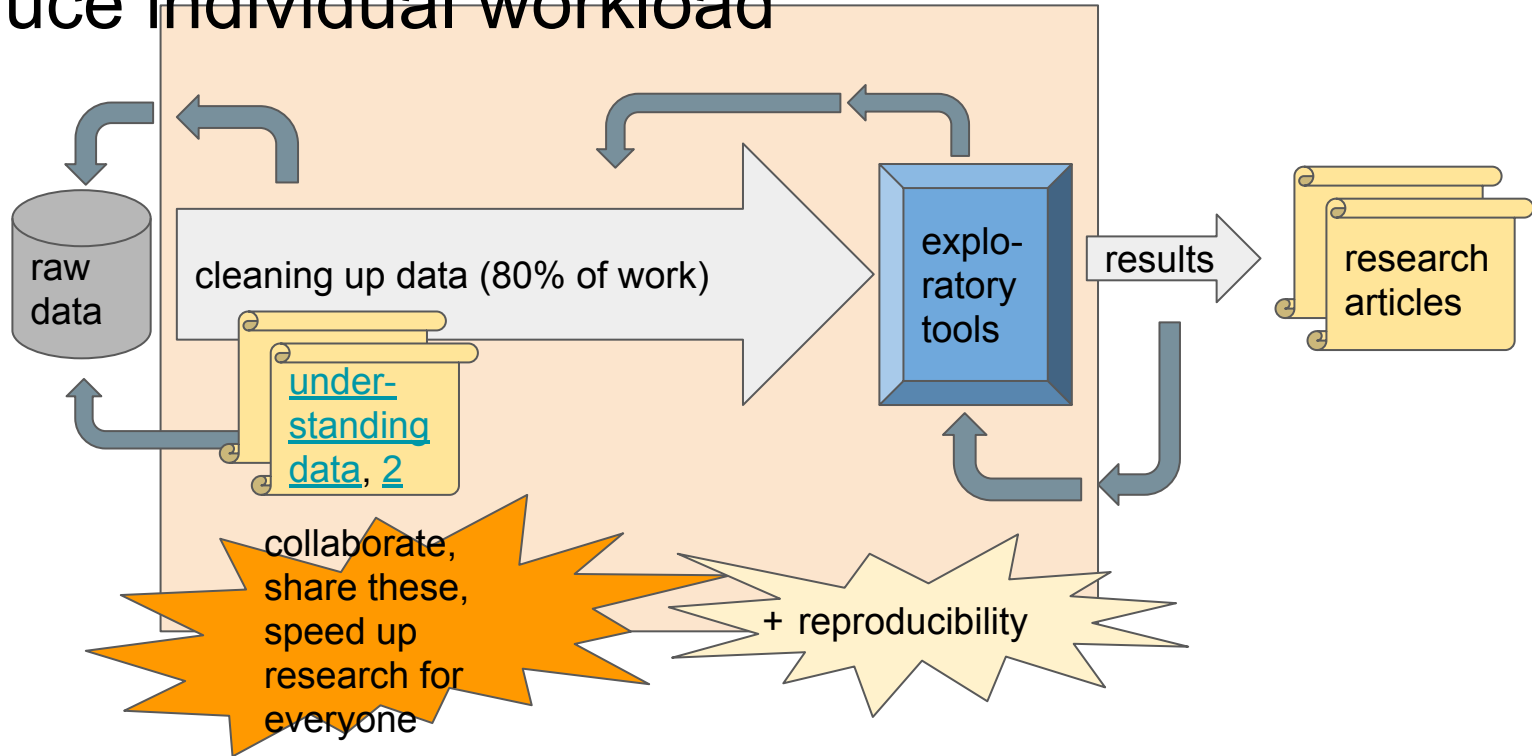
# Digital humanities research process



80% of your time for data cleanup, another
80% for algorithms, …

# Leverage collaboration, open science workflows to reduce individual workload

# Workflow/Tools

1. Data access
2. Possible preprocessing: <u>R</u>, <u>Python</u>, <u>tm</u> (for texts), <u>OpenRefine</u>, …
3. Zero or more of:
   - Statistics: <u>R</u>, <u>stats</u>, <u>pandas</u>, …
   - Topic modeling: <u>Mallet</u>, <u>topicmodels</u>, <u>LDAvis</u>, <u>gensim</u>, … (for texts)
   - Dimensionality reduction/clustering: <u>stats</u>, <u>lsa</u>, <u>BayesLCA</u>, <u>pvclust</u>, <u>Weka</u>, … (also for texts)
   - Social network analysis: <u>igraph</u>, <u>sna</u>, <u>statnet</u>, <u>sonia</u>, <u>Gephi</u>, …
   - Simulation: <u>NetLogo</u>, …
   - Neural networks: <u>som</u>, <u>TensorFlow™</u>, … (also for texts)
   - Association rule learning: <u>arules</u>, <u>Weka</u>, …
   - Anomaly detection: <u>AnomalyDetection</u>, …
4. Structured visualization: <u>Tableau</u>, <u>Palladio</u>, <u>RAW</u>, <u>nodegoat</u>, <u>matplotlib</u>, <u>ggplot2</u>, <u>iPlots</u>, <u>plot.ly</u>, <u>Leaflet</u>, <u>Gephi</u>, <u>CartoDB</u>, … or text visualization: <u>Voyant Tools</u>, <u>Textexture</u>, <u>Wordsift</u>, …
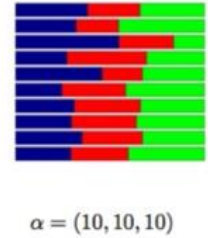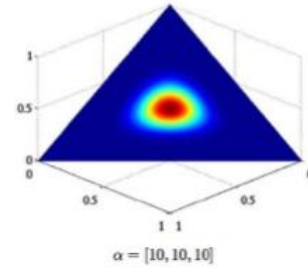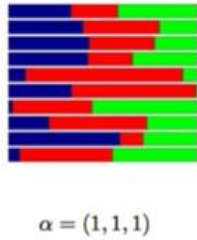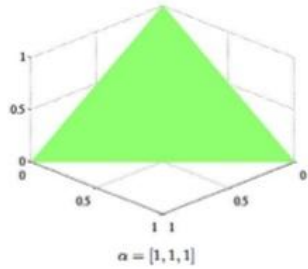
# Voyant tools

# Types of data

- Structured (databases) vs unstructured (text, image, video, audio)
- Clean vs messy
- **Biased? <- incomplete, messy, badly sampled**

# Topic Modeling: LDA - Assumptions

- A document collection contains N topics

- A single document can consist of multiple topics (e.g. 30% war and 70% cooking)

- The N topics are in essence probability distributions over words (e.g. there is a 1,5% chance that a random word from a 'war' topic is 'attack', while only a 0,00001% chance in a 'cooking' topic)

- There are two distributions that give the prior probabilities of:

  a. the probability of topic mixes in documents (e.g. how likely is it that a single document talks about all the topics vs. only a few) , and

  b. the probability mix of words in a topic (e.g. do individual topics mainly contain many words or just a few)

# Topic Modeling: LDA - Role of (symmetric) priors



$\alpha = [1,1,1]$

$\alpha = (1,1,1)$

$\alpha = [10,10,10]$

$\alpha = (10,10,10)$

$\alpha = [.1,.1,.1]$

$\alpha = (0.1, 0.1, 0.1)$
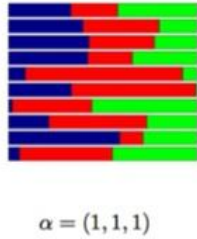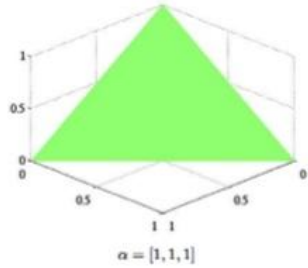
# Topic Modeling: LDA - How it works

- Take all words and documents and randomly assign them to topics (based on the prior distributions)
- Calculate the combined probability of this combination producing the documents we have
- Update the topic assignments as well as the prior distributions so the probability increases
- Repeat many many times until we're happy

# LDA in Practice

```
corpus <-
VCorpus(DirSource("/srv/data/varieng/ceec-subcorpora/scot-17
00-1719/"))
corpus <- tm_map(corpus,content_transformer(tolower))
corpus <- tm_map(corpus,removeNumbers)
corpus <- tm_map(corpus,removePunctuation)
corpus <- tm_map(corpus,removeWords,stopwords("SMART"))
corpus <- tm_map(corpus,stripWhitespace)
numtopics <- 20
lda <- LDA(DocumentTermMatrix(corpus), numtopics)
```

# Topic Modeling: LDA - Role of priors

# Topic Modeling: LDA - Effect of priors

- Traditional LDA supposed uniform priors

- Turns out non-uniform priors make sense for how topics appear in documents, but not for how words appear in topics

  → as-LDA, which also turns out to need less pre-filtering of e.g. stopwords, numbers, because these can be sequestered into a common topic without constraining how other topics appear

**Figure 6** Primary discipline of respondents who have been involved in big data research (*n* = 9195)

| Discipline | Percentage ever involved in big data | Value |
|---|---|---|
| Social Statistics and Research Methods | | 121 |
| Economics | | 155 |
| Demography, Population Studies, and Human Geography | | 49 |
| Health Sciences | | 409 |
| Social Policy and Public Policy | | 80 |
| Marketing | | 80 |
| Management and Business Studies | | 285 |
| Communication and Media Studies | | 188 |
| Political Science and International Studies | | 167 |
| Linguistics | | 64 |
| Sociology | | 208 |
| Other | | 269 |
| History | | 28 |
| Social Work | | 65 |
| Education | | 413 |
| Nursing | | 50 |
| Criminology and Criminal Justice | | 61 |
| Anthropology | | 50 |
| Law and Legal Studies | | 21 |
| Psychology | | 286 |
| Counseling and Psychotherapy | | 27 |

**Figure 15  Challenges facing big data researchers (*n* = 2273)**

| Challenge | Big problem for me | Something of a problem for me | Not a problem for me |
|---|---|---|---|
| Getting funding for my research | 1290 | 1181 | 585 |
| Getting access to commercial or proprietary data for my research | 970 | 1224 | 827 |
| Finding collaborators with the right skills and knowledge | 677 | 1343 | 1039 |
| Learning new software for myself | 672 | 1449 | 944 |
| Learning new analytic methods for myself | 615 | 1485 | 960 |
| Choosing a suitable journal in which to publish my research | 608 | 1339 | 1098 |
| Establishing a successful career in an interdisciplinary field | 554 | 1295 | 1183 |
| Developing effective research designs | 404 | 1402 | 1243 |
| Getting ethical approval for my research | 261 | 824 | 1954 |

**Figure 19** Problems encountered by amount of research using big data (*n* = 2266)

SMELLY MAPS

Q

BLACKFRIARS BRIDGE
PRIMARY

57.9%
EMISSIONS

ANGER
ANTICIPATION
TRUST
DISGUST
SURPRISE
FEAR
SADNESS
JOY

EMISSIONS   NATURE   FOOD   ANIMALS   WASTE