

# Big Data & Automated Content Analysis

## Week 1 – Monday: »Introduction«

---

Anne Kroon

a.c.kroon@uva.nl

@annekroon

29 March 2021

Afdeling Communicatiewetenschap  
Universiteit van Amsterdam

# Today

Getting to know each other

Setting the stage

Defining “Big Data”

Are we doing Big Data research?

Defining Computational (Social|Communication) Science

The toolbox

The role of software in CSS

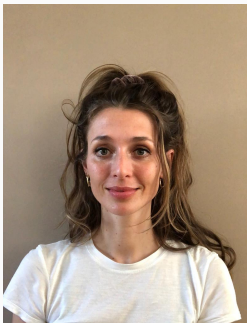
This course

Which techniques?

Which tools?

Python: A language, not a program

# Anne



dr. Anne Kroon

Assistant Professor Corporate Communication

- studied Journalism (HU) and Communication Science at the UVA
- PhD candidate @ ASCoR 2014–2017
- bias in media representations of minorities, algorithmic bias in hiring and selection
- computational research methods

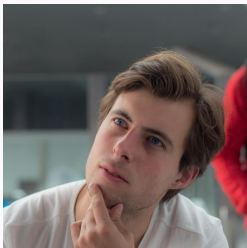
@annekroon    a.c.kroon@uva.nl

REC-C 7<sup>th</sup> floor

# Vlad

Vladislav Petkevich

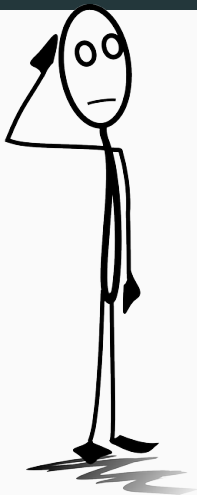
Junior Lecturer Communication Science



- MSc in Communication Research (2020)
- Interested in political communication, especially election campaigns
- Even more interested in applying computational research methods (e.g. NLP, machine vision) to studying social phenomena

[v.petkevich@uva.nl](mailto:v.petkevich@uva.nl)

## You



Your name?

Where are you atm?

Your background?

Your reason to follow this course?

## Setting the stage

---

Defining “Big Data”

Big data is like teenage sex:  
everyone talks about it,  
nobody really knows how to do it,  
everyone thinks everyone else is  
doing it, so everyone claims they  
are doing it...

(Dan Ariely)

Getting to know each other  
ooo

Setting the stage  
oo●oooo  
oo  
oooooo

The toolbox  
oooo

This course  
ooo  
oo  
oo  
oooo  
ooo  
oooo

References

## The “pragmatic” definition

Everything that needs so much computational power and/or storage that you cannot do it on a regular computer.



Getting to know each other  
○○○

Setting the stage  
○○○●○○○  
○○  
○○○○○○

The toolbox  
○○○○

This course  
○○○  
○○  
○○○○  
○○○  
○○○○

References

## The “commercial” definition

### Gartner (n.d.)

“Big data is high-volume, high-velocity and/or high-variety information assets that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, and process automation.”

## The “critical” definition

### Boyd and Crawford (2012)

“

1. Technology: maximizing computation power and algorithmic accuracy to gather, analyze, link, and compare large data sets.
2. Analysis: drawing on large data sets to identify patterns in order to make economic, social, technical, and legal claims.
3. Mythology: the widespread belief that large data sets offer a higher form of intelligence and knowledge that can generate insights that were previously impossible, with the aura of truth, objectivity, and accuracy.

”

## APIs, researchers and tools *make* Big Data

### Vis, 2013

Inevitable influences of:

- APIs
- filtering, search strings, ...
- changing services over time
- organizations that provide the data



*Do you think we are doing Big Data analysis?*

## Setting the stage

---

Are we doing Big Data research?

## Are we doing Big Data research in this course?

### Depends on the definition

- Not if we take a definition that *only* focuses on computing power and the amount of data
- **But:** We are using the same techniques. And they *scale* well.
- Oh, and about that high-performance computing in the cloud: We actually *do* have access to that, so if someone has a really great idea...

## Setting the stage

---

Defining Computational  
(Social|Communication) Science

## A very young field

### Lazer et al. (2009)

“The capacity to collect and analyze massive amounts of data has transformed such fields as biology and physics. But the emergence of a data-driven ‘computational social science’ has been much slower.”



## Epistemologies and paradigm shifts

### Kitchin (2014)

- (Reborn) empiricism: purely inductive, correlation is enough
- Data-driven science: knowledge discovery guided by theory
- Computational social science and digital humanities: employ Big Data research within existing epistemologies
  - DH: descriptive statistics, visualizations
  - CSS: prediction and simulation

## CCS as a subset of CSS

### Hilbert et al. (2019)

“...our definition of computational communication science as an application of computational science to questions of human and social communication. As such, it is a natural subfield of computational social science” (followed by references to CSS definitions)

## Data, analysis, theory

### van Atteveldt and Peng (2018)

“...computational communication science studies generally involve: (1) large and complex data sets; (2) consisting of digital traces and other “naturally occurring” data; (3) requiring algorithmic solutions to analyze; and (4) allowing the study of human communication by applying and testing communication theory.”



1. *What do you think? What is the essence of Big Data/CSS/CCS?*
2. *How will what we do here relate to theories and methods from other courses?*

# The toolbox

---

The role of software in CSS

## Why program your own tool?

### Vis (2013)

“Moreover, the tools we use can limit the range of questions that might be imagined, simply because they do not fit the affordances of the tool. Not many researchers themselves have the ability or access to other researchers who can build the required tools in line with any preferred enquiry. This then introduces serious limitations in terms of the scope of research that can be done.”

## Some considerations regarding the use of software in science

Assuming that science should be *transparent* and *reproducible* by *anyone*, we should

**use tools that are**

- platform-independent
- free (as in beer and as in speech, gratis and libre)
- which implies: open source

This ensures it can our research (a) can be reproduced by anyone, and that there is (b) no black box that no one can look inside. ⇒ ongoing open-science debate! (van Atteveldt, Strycharz, Trilling, & Welbers, 2019)

## Why program your own tool?

### Vis (2013)

“[...] these [commercial] tools are often unsuitable for academic purposes because of their cost, along with the problematic ‘black box’ nature of many of these tools.”

### Mahrt and Scharkow (2013)

“[...] we should resist the temptation to let the opportunities and constraints of an application or platform determine the research question [...]”



# This course

---

## Which techniques?

## What we will learn the next weeks

### 1. How to collect and deal with data

APIs, unstructured data-files, databases, storage in different file formats

### 2. How to analyze data

Automated content analysis, regular expressions, natural language processing, cluster analysis, machine learning

# The ACA toolbox

	Methodological approach		
	<i>Counting and Dictionary</i>	<i>Supervised Machine Learning</i>	<i>Unsupervised Machine Learning</i>
<b>Typical research interests and content features</b>	visibility analysis sentiment analysis subjectivity analysis	frames topics gender bias	frames topics
<b>Common statistical procedures</b>	string comparisons counting	support vector machines naive Bayes	principal component analysis cluster analysis latent dirichlet allocation semantic network analysis

deductive

inductive

Boumans, J.W., & Trilling, D. (2016). Taking stock of the toolkit: An overview of relevant automated content analysis approaches and techniques for digital journalism scholars. *Digital Journalism*, 4, 1. 8–23.

# This course

---

## Which tools?

Getting to know each other  
ooo

Setting the stage  
ooooooo  
oo  
oooooo

The toolbox  
oooo

**This course**  
ooo  
o●  
oooo  
ooo  
oooo

References

## The methods

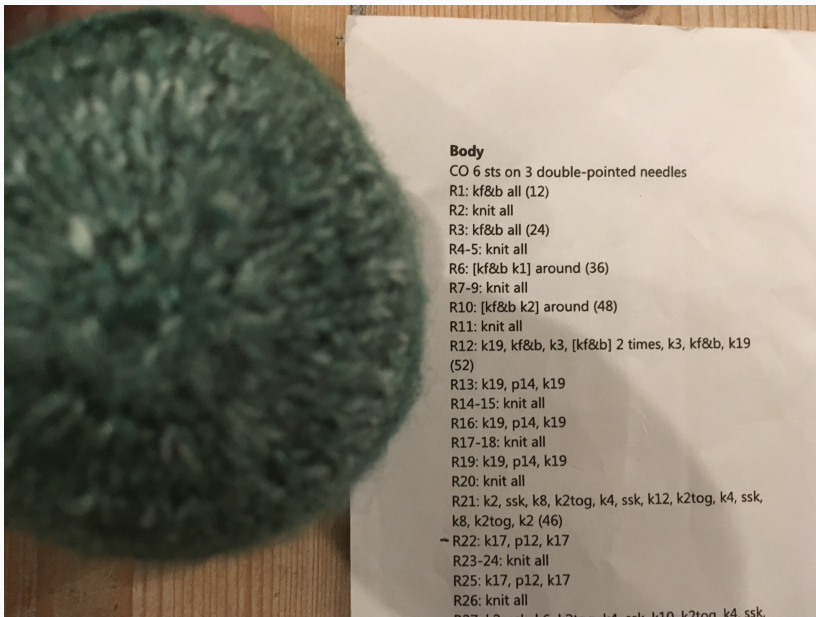
### The tools we use for this

The programming language Python (and the huge amount of Python modules others already wrote).

# This course

---

Python: A language, not a program



### Body

CO 6 sts on 3 double-pointed needles

R1: kf&b all (12)

R2: knit all

R3: kf&b all (24)

R4-5: knit all

R6: [kf&b k1] around (36)

R7-9: knit all

R10: [kf&b k2] around (48)

R11: knit all

R12: k19, kf&b, k3, [kf&b] 2 times, k3, kf&b, k19  
(52)

R13: k19, p14, k19

R14-15: knit all

R16: k19, p14, k19

R17-18: knit all

R19: k19, p14, k19

R20: knit all

R21: k2, ssk, k8, k2tog, k4, ssk, k12, k2tog, k4, ssk,  
k8, k2tog, k2 (46)

~ R22: k17, p12, k17

R23-24: knit all

R25: k17, p12, k17

R26: knit all

R27-28: k2, ssk, k8, k2tog, k4, ssk, k10, k2tog, k4, ssk,

An algorithm in a language that's a bit harder than Python

# Python

## What?

- A language, not a specific program
- Huge advantage: flexibility, portability
- One of *the* languages for data analysis. (The other one is R.)  
But Python is more flexible—the original version of Dropbox was written in Python. Some people say: R for numbers, Python for text and messy stuff.

## Which version?

We use Python 3.

<http://www.google.com> or <http://www.stackexchange.com> still may show you some Python2-code, but that can easily be adapted. Most notable difference: In Python 2, you write `print "Hi"`, this has changed to `print ("Hi")`.



## What can we do?

### Cool research, sure, but what can we do?

- Dependency from third parties:  
scraping < API < server-side implementation
- Restrictions (e.g., Twitter: sprinkler, garden hose, fire hose)  $\Rightarrow$   
Vis, 2013: Data are made!
- We can't just trust the numbers. Some tasks require human coders – or a qualitative approach, at least as a pre-study.

# This course

---

## The schedule

Getting to know each other  
ooo

Setting the stage  
ooooooo  
oo  
oooooo

The toolbox  
oooo

**This course**  
ooo  
oo  
oo  
oooo  
o●o  
oooo

References

## The schedule

# The schedule

## Each week

In general: A lecture and a lab session.

Each week one method.

## Examinations

A mid-term take-home exam and an individual research project on which you work during the whole course.

## Self-study

Play around! You really have to *do* it to learn programming. See it as your weekly assignment ;-)

# This course

---

## Next meetings

Getting to know each other  
○○○

Setting the stage  
○○○○○○○  
○○  
○○○○○○

The toolbox  
○○○○

**This course**  
○○○  
○○  
○○  
○○○○  
○○○  
○○  
○●○○

References

Next meetings

Getting to know each other  
○○○

Setting the stage  
○○○○○○○  
○○  
○○○○○○

The toolbox  
○○○○

This course  
○○○  
○○  
○○  
○○○○  
○○○  
○○○  
○○●○

References

## Week 1: Introduction

**Thursday, 1–4**

Lab session.

**Fun with data**

Getting to know each other  
ooo

Setting the stage  
ooooooo  
oo  
oooooo

The toolbox  
oooo

This course  
ooo  
oo  
oo  
oooo  
ooo  
ooo●

References

Make sure you can run Python code and install packages. Otherwise, you won't be able to follow along on Thursday. (See instructions you got. Use Vlad's office hours if you cannot figure it out.)





Boyd, D., & Crawford, K. (2012). Critical questions for Big Data. *Information, Communication & Society*, 15(5), 662–679. doi:10.1080/1369118X.2012.678878



Gartner. (n.d.). Big data. Retrieved from <https://www.gartner.com/en/information-technology/glossary/big-data>



Hilbert, M., Barnett, G., Blumenstock, J., Contractor, N., Diesner, J., Frey, S., ... Zhu, J. J. H. (2019). Computational Communication Science : A Methodological Catalyzer for a Maturing Discipline. *International Journal of Communication*, 13, 3912–3934.



Kitchin, R. (2014). Big Data, new epistemologies and paradigm shifts. *Big Data & Society*, 1(1), 1–12. doi:10.1177/2053951714528481



Lazer, D., Pentland, A., Adamic, L., Aral, S., Barabási, A.-L., Brewer, D., ... van Alstyne, M. (2009). Computational social science. *Science*, 323, 721–723. doi:10.1126/science.1167742



Mahrt, M., & Scharrow, M. (2013). The value of Big Data in digital media research. *Journal of Broadcasting & Electronic Media*, 57(1), 20–33. doi:10.1080/08838151.2012.761700



van Atteveldt, W., Strycharz, J., Trilling, D., & Welbers, K. (2019). Toward Open Computational Communication Science : A Practical Road Map for Reusable Data and Code University of Amsterdam , the Netherlands. *International Journal of Communication*, 13, 3935–3954.



van Atteveldt, W., & Peng, T. Q. (2018). When Communication Meets Computation: Opportunities, Challenges, and Pitfalls in Computational Communication Science. *Communication Methods and Measures*, 12(2-3), 81–92. doi:10.1080/19312458.2018.1458084



Vis, F. (2013). A critical reflection on Big Data: Considering APIs, researchers and tools as data makers. *First Monday*, 18(10), 1–16. doi:10.5210/fm.v18i10.4878