



DIS08 – Data Modeling

05 – Open Data, data file formats, spreadsheets, data cleaning

Philipp Schaer, Technische Hochschule Köln, Cologne, Germany

Version: SS2020

Disclaimer

This lesson is based on the Library Carpentry

- <https://librarycarpentry.org/lc-spreadsheets/>

Chapter 14 of Automating the Boring Stuff

And some slides by Kai Dürkop

- <https://bio.informatik.uni-jena.de/wp/wp-content/uploads/2015/03/web.pdf>

Agenda

Last week

- Regular expressions
- Mining or searching in files

This week

- Open Data and where to get it
- Data formats: CSV, JSON, XML
- Working with CSV data
- Data cleaning in Excel in a nutshell

Next week

- Hands-on Python

Definition of „Openness“

- “Open means **anyone** can **freely access, use, modify, and share** for **any purpose** (subject, at most, to requirements that preserve provenance and openness).”
- “Open data and content can be **freely used, modified, and shared** by **anyone** for **any purpose**”
 - Open Data
 - Open Access
 - Open Science
 - Open Source
 - Open Government ...

Potential Impact of Open Data

	Potential impact: 2011 research	Value captured %	Major barriers
Location-based data	<ul style="list-style-type: none"> \$100 billion+ revenues for service providers Up to \$700 billion value to end users 	<p>50–60</p>	<ul style="list-style-type: none"> Penetration of GPS-enabled smartphones globally
US retail ¹	<ul style="list-style-type: none"> 60%+ increase in net margin 0.5–1.0% annual productivity growth 	<p>30–40</p>	<ul style="list-style-type: none"> Lack of analytical talent Siloed data within companies
Manufacturing ²	<ul style="list-style-type: none"> Up to 50% lower product development cost Up to 25% lower operating cost Up to 30% gross margin increase 	<p>20–30</p>	<ul style="list-style-type: none"> Siloed data in legacy IT systems Leadership skeptical of impact
EU public sector ³	<ul style="list-style-type: none"> ~€250 billion value per year ~0.5% annual productivity growth 	<p>10–20</p>	<ul style="list-style-type: none"> Lack of analytical talent Siloed data within different agencies
US health care	<ul style="list-style-type: none"> \$300 billion value per year ~0.7% annual productivity growth 	<p>10–20</p>	<ul style="list-style-type: none"> Need to demonstrate clinical utility to gain acceptance Interoperability and data sharing

1 Similar observations hold true for the EU retail sector.

2 Manufacturing levers divided by functional application.

3 Similar observations hold true for other high-income country governments.

Open Government

The screenshot shows the official website of the White House under President Barack Obama. At the top, it says "the WHITE HOUSE PRESIDENT BARACK OBAMA". Below that is the official seal of the White House. To the right are links for "Get Email Updates" and "Contact Us". A navigation bar at the bottom includes links for "BLOG", "PHOTOS & VIDEO", "BRIEFING ROOM", "ISSUES", "the ADMINISTRATION", "the WHITE HOUSE", and "our GOVERNMENT".

Barack Obama, 21 January 2009 (his first day in the office): "Memorandum for the Heads of Executive Departments and Agencies"

- **Transparency and Open Government**
 - Government should be transparent.
 - Government should be participatory.
 - Government should be collaborative.
- **Freedom of Information Act (FOIA)**
 - "In the face of doubt, openness prevails."
 - "commitment to accountability and transparency"



Open Data Principles

- **Complete** *All public data is made available.*
- **Primary** *Data is as collected at the source, with the highest possible level of granularity, not in aggregate/modified forms.*
- **Timely** *Data is made available as quickly as necessary to preserve the value of the data.*
- **Accessible** *Data is available to the widest range of users for the widest range of purposes.*
- **Machine processable** *Data is reasonably structured to allow automated processing.*
- **Non-discriminatory** *Data is available to anyone, with no requirement of registration.*
- **Non-proprietary** *Data is available in a format over which no entity has exclusive control.*
- **License-free** *Data is not subject to any copyright, patent, ...*

Open Government Data

Educate Citizens (pedagogy)

Data visualisation -
Display governance process -
Infographics -



Transparency

Open Government

Participation

Collaboration

Consult Citizens

Seek criticisms
suggestions
and ideas



Deliberate with Citizens

Organise public debates



Co-Design Policies with Citizens



Monitor Policies

- Communication strategies
- Dashboard
- Timelines



Break down Silos and Pyramidal Structures

- Inside organisations
- Between organisations

Work Horizontally

- Between organisations
- Through territories
With:
- Service design tools
- Agile methodologies
By:
- Empowering citizens
- Favoring cooperation



Organise Partnerships (inside/between)



Open government diagram
by Armel Le Coz and Cyril Lage
released under Creative Commons
Attribution terms

Open Government in Germany

Der Bundestag hat am **Donnerstag, 18. Mai 2017**, das sogenannte **E-Government-Gesetz** geändert. Dem Gesetzentwurf der Bundesregierung ([□ 18/11614](#)) in der vom Innenausschuss geänderten Fassung ([□ 18/12406](#)) stimmten die Koalitionsfraktionen zu, die Opposition enthielt sich. Mit dem Gesetz werde die Grundlage für die aktive Bereitstellung von elektronischen Daten der Behörden der unmittelbaren Bundesverwaltung geschaffen, schreibt die Regierung. Um dem Anspruch auf eine Vorreiterrolle Deutschlands gerecht zu werden, orientiert sich die Regelung nach eigenen Angaben an international anerkannten Open-Data-Prinzipien, wie sie beispielsweise in der Internationalen Open-Data-Charta (IODC) oder in der Open-Data-Charta der sogenannten G8-Staaten beschrieben werden.

What to do with open public data...

OPEN.NRW
AKTUELLE PROJEKTE

Frühzeitige Beteiligung zu neuer Rheinquerung

Der Landesbetrieb Straßenbau NRW plant derzeit eine Rheinquerung zwischen Köln und Bonn. Der konkrete Verlauf ist noch offen und wird unter Beteiligung der Menschen und Interessengruppen der...

Verwaltung einfacher machen

Wo ist künftig noch eine Unterschrift nötig und wann muss man persönlich erscheinen? Zwischen August und November 2017 konnten Interessierte sich beteiligen und helfen, bürokratische Hürden...

Arbeit im digitalen Wandel – Ihre Ideen für NRW sind gefragt

Die Digitalisierung verändert unsere Arbeitswelt massiv. Doch wie lassen sich digitaler Wandel und die Grundsätze guter und fairer Arbeit in Einklang bringen? Diskutieren und kommentieren Sie im...

Vernetzte Nachbarschaften

Grünflächen, Geschäfte, eine gute Verkehrsanbindung: Wer sein Wohnumfeld gestalten möchte, braucht Unterstützung und Partner. Deshalb hat das Ministerium für Bauen, Wohnen, Stadtentwicklung...

Digitale Lebenswelten in NRW: Wir sind dabei – und Sie?

Naturschutzinformationen NRW: OSIRIS

Ohne vertrauliche und konstruktive Zu-

Where to get Open Data

- <https://www.europeandataportal.eu/>
- <https://www.govdata.de/>
- <https://open.nrw/>
- <https://offenedaten-koeln.de>
- <https://www.kaggle.com/datasets>
- and many more...



kaggle

Data formats

Binary

- not human readable
- memory efficient and fast to parse
- but: platform-dependent
- difficult to convert the format (e.g. open a Word Document in Open Office...)

Text

- (mostly) human readable
- waste more memory and relatively slow to parse
- platform-independent (but: still encoding problems!)
- format can be easily transformed

Common (text) data formats

- **CSV**: simplest format possible: Just strings separated by commas and newlines.
- **XML**: most common text data format. XHTML is the language of the web.
- **JSON**: uses javascript syntax. Much better readable and more sparse than XML.

CSV - Comma seperated values #1

4/5/2015 13:34,Apples,73

4/5/2015 23:41,Cherries,85

4/6/2015 12:46,Pears,14

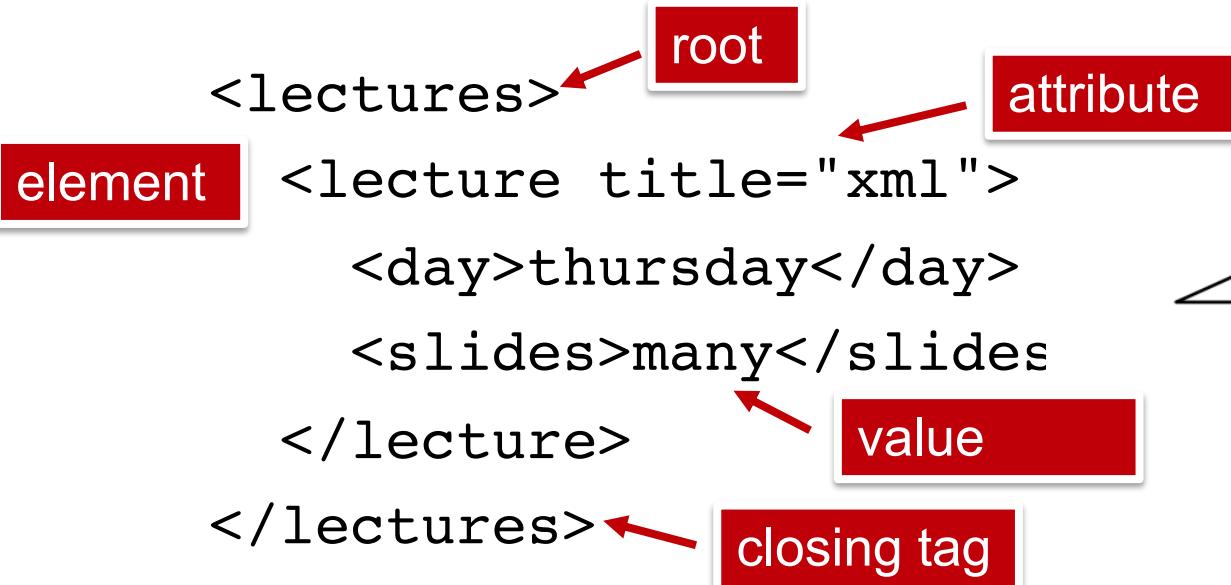
4/8/2015 08:59,Oranges,52

- Column separator:
 - default is comma (","), but also tabulator ("\t")
 - very common row separator: usually newline ("\n")
- Strings can be enclosed in quotation marks to escape special characters
- Quotation marks are escaped by another quotation marks

CSV - Comma seperated values #2

- The advantage of CSV files is **simplicity**.
- CSV files are **widely supported** by many types of programs, can be viewed in text editors, and are a straightforward way to represent data.
- Whenever your data has **no nested structure**: **USE CSV!**
- Can be easily read in every programming language.
- Can be opened in text editors and in **Excel!**
- Comma is default, but tabs are often better as you rarely have to escape strings in tab separated files.

XML - Extensible Markup Language



- Very common data format
- Solves a lot of problems
(namespaces, encoding,
embedded data)
- Not very readable, very verbose

JSON - JavaScript Object Notation

```
[  
  {  
    "name": "Glucose",  
    "formula": "C6H12O6",  
    "similarTo": [ "Hexose", "Fructose" ]  
  }  
]
```

- Just Javascript data. Human readable but not easy to parse.
- Popular use: Client and web server often use JSON to communicate. Javascript can naturally work with json.
- Also almost compatible to Pythons syntax booleans, numbers, strings, arrays, objects (dictionaries).

Question time!

- How many people have used spreadsheets in their work?
- What kind of operations do you do in spreadsheets?
- Which ones do you think spreadsheets are good for?

Why do we care about spreadsheets?

- Spreadsheets are **good for data entry**, but in reality we **tend to use spreadsheet programs for much more** than data entry. We use them to create data tables for publications, to generate summary statistics, and make figures.
- Generating **tables for reports** in a spreadsheet is not optimal. We advise you to do this sort of operation within your document editing software.
- **HOWEVER**, there are circumstances where you might want to use a spreadsheet program to produce “quick and dirty” calculations or figures, and some of these features can be used in **data cleaning**, prior to importation into a statistical analysis program. We will show you how to use some features of spreadsheet programs to check your **data quality** along the way and produce preliminary summary statistics.

Spreadsheet Software

- Microsoft Excel
- Libre/Open Office
- Apple Numbers
- OnlyOffice (in Sciebo)
- ...

- **Google Spreadsheets**
 - <http://spreadsheets.google.com>

Unbenannte Tabelle Freigeben P

Datei Bearbeiten Ansicht Einfügen Format Daten Tools Add-ons Hilfe Alle Änderungen in Drive gespeichert

100% Arial 10 B I A H J K L M

A B C D E F G H I J K L M

18	2004	GHJ	Handel										
19	2004	JK	Finanzierung										
20	2004	L-O	Öffentliche und p	49385	35716	39572	45402	51580	49073	51663			
21	2004	C	Bergbau und Ge	57350	53982	44850	57172	.	.	57981			
22	2004	D	Verarbeitendes C	51065	40095	43481	48140	51524	53383	60912			
23	2004	E	Energieversorgu	79869	49030	59071	67855	.	.	88152			
24	2004	F	Baugewerbe	42586	39348	42501	45861	60951	57281	54784			
25	2004	G	Handel										
26	2004	H	Gastgewerbe	27341	24982	27100	24371	31671	28463	30896			
27	2004	I	Verkehr und Nac	43831	32257	35829	37278	44376	41601	49919			
28	2004	J	Kredit- und Versi	64488	63845	63934	62967	60502	62487	66410			
29	2004	K	Grundst.-										
30	2004	L	Öffentl.Verwaltung										
31	2004	M	Erziehung und U	58640	-	-	-	-	-	58640			
32	2004	N	Gesundheitswesen										
33	2004	O	Erbringung sonst	51302	36780	44438	54367	/	50234	62353			
34	2004	CA	Kohlenbergbau										
35	2004	CB	Erzbergbau										
36	2004	10	Kohlenbergbau										
37	2004	11	Gew. v. Erdöl u. Erdgas										
38	2004	12	Bergbau auf Urai -	-	-	-	-	-	-	-			
39	2004	13	Erzbergbau	-	-	-	-	-	-	-			
40	2004	14	Gewinnung von Steinen und Erden										
41	2004	15	Ernährungsgewe	39285	275								
42	2004	16	Tabakverarbeitur	57232	326								
43	2004	17	Textilgewerbe	41956	326								
44	2004	18	Bekleidungsgew	44016	394								
45	2004	19	Ledergewerbe	39190	34383	36699	32107	46975	-	-			
46	2004	20	Holzgewerbe (oh	38810	32445	36678	40448	40154	.g	.			
47	2004	21	Papiergewerbe	48919	44128	40797	45795	49228	57797	52610			

To be honest...It's not the best
solution, but always available!

Intermission! Using regex to change data

The screenshot shows the 'Suchen und ersetzen' (Search and Replace) dialog in Google Sheets. The search term is `^(\\d)` and the replacement term is `$1$1`. A red callout box highlights these terms with arrows pointing to them. The search scope is set to 'Alle Tabellenblätter'. Under 'Suchen', the following options are checked:

- Groß-/Kleinschreibung beachten
- Gesamten Zelleninhalt vergleichen
- Suche mithilfe regulärer Ausdrücke [Hilfe](#)
- Auch in Formeln suchen

At the bottom, it shows the result: `"^(\\d)" an 4 Stellen durch "$1$1 " ersetzt`. The buttons at the bottom are 'Finden' (Find), 'Ersetzen' (Replace), 'Alle ersetzen' (Replace All), and 'Fertig' (Done).

Suchen und ersetzen

Suchen

Ersetzen durch

Suchen

Alle Tabellenblätter

Groß-/Kleinschreibung beachten

Gesamten Zelleninhalt vergleichen

Suche mithilfe regulärer Ausdrücke [Hilfe](#)

Auch in Formeln suchen

`"^(\\d)" an 4 Stellen durch "$1$1 " ersetzt`

Finden Ersetzen Alle ersetzen Fertig

Structuring data in spreadsheets

The cardinal rules of using spreadsheet programs for data:

- Put all your **variables in columns** - the thing you're measuring, like 'weight' or 'temperature'.
- Put each **observation in its own row**.
- **Don't combine multiple pieces of information in one cell.** Sometimes it just seems like one thing, but think if that's the only way you'll want to be able to use or sort that data.
- **Leave the raw data raw** - don't mess with it!
- Export the cleaned data to a **text based format** like CSV. This ensures that anyone can use the data, and is the format required by most data repositories.

An example - Messy data?

- For instance, we have data from a **survey of small mammals** in a desert ecosystem. **Different people** have gone to the field and **entered data in to a spreadsheet**. They keep track of things like species, plot, weight, sex and date collected.

Date collected	Plot	Species-Sex	Weight
1/9/78	1	DM-M	40
1/9/78	1	DM-F	36
1/9/78	1	DS-F	135
1/20/78	1	DM-F	39
1/20/78	2	DM-M	43
1/20/78	2	DS-F	144
3/13/78	2	DM-F	51
3/13/78	2	DM-F	44
3/13/78	2	DS-F	146

An example - Messy data!

- The problem is that **species** and **sex** are in the same field.
- If you want to look at all of one species or look at different weight distributions by sex, it would be hard to do so.
- Put sex and species in **different columns!**

Date collected	Plot	Species	Sex	Weight
1/9/78	1	DM	M	40
1/9/78	1	DM	F	36
1/9/78	1	DS	F	135
1/20/78	1	DM	F	39
1/20/78	2	DM	M	43
1/20/78	2	DS	F	144
3/13/78	2	DM	F	51
3/13/78	2	DM	F	44
3/13/78	2	DS	F	146

Common Spreadsheet Errors #1

Multiple tables

- A common strategy is creating multiple data tables within one spreadsheet. **This confuses the computer!**
 - When you create multiple tables within one spreadsheet, you're drawing false associations between things for the computer, which sees each row as an observation. You're also potentially using the same field name in multiple places!

Common Spreadsheet Errors #2

Using formatting to convey information

Plot: 2			
Date collected	Species	Sex	Weight
1/8/14	NA		
1/8/14	DM	M	44
1/8/14	DM	M	38
1/8/14	OL		
1/8/14	PE	M	22
1/8/14	DM	M	38
1/8/14	DM	M	48
1/8/14	DM	M	43
1/8/14	DM	F	35
1/8/14	DM	M	43
1/8/14	DM	F	37
1/8/14	PF	F	7
1/8/14	DM	M	45
1/8/14	OT		
1/8/14	DS	M	157
1/8/14	OX		
2/18/14	NA	M	218
2/18/14	PF	F	7
2/18/14	DM	M	52
measurement device not calibrated			



Plot: 2				
Date collected	Species	Sex	Weight	Calibrated
1/8/14	NA			
1/8/14	DM	M	44	Y
1/8/14	DM	M	38	Y
1/8/14	OL			
1/8/14	PE	M	22	Y
1/8/14	DM	M	38	Y
1/8/14	DM	M	48	Y
1/8/14	DM	M	43	Y
1/8/14	DM	F	35	Y
1/8/14	DM	M	43	Y
1/8/14	DM	F	37	Y
1/8/14	PF	F	7	Y
1/8/14	DM	M	45	Y
1/8/14	OT			
1/8/14	DS	M	157	N
1/8/14	OX			
2/18/14	NA	M	218	N
2/18/14	PF	F	7	Y
2/18/14	DM	M	52	Y

Common Spreadsheet Errors #3

Using formatting to make the data sheet look pretty

- **Example:** merging cells.
- **Solution:** If you're not careful, formatting a worksheet to be more aesthetically pleasing can compromise your computer's ability to see associations in the data. Merged cells are an absolute formatting NO-NO if you want to make your data readable by statistics software. Consider restructuring your data in such a way that you will not need to merge cells to organize your data.

Common Spreadsheet Errors #4

Field name problems

- Choose **descriptive field names**, but be careful not to include: spaces, numbers, or special characters of any kind.
- **Spaces** can be misinterpreted by parsers that use whitespace as delimiters and some programs don't like field names that are text strings that start with numbers.
- **Underscores** (_) are a good alternative to spaces and consider writing names in camel-case to improve readability.
- Remember that **abbreviations** that make sense at the moment may not be so obvious in 6 months but don't overdo it with names that are excessively long. Including the units in the field names avoids confusion and enables others to readily interpret your fields.

Common Spreadsheet Errors #4

Good Name	Good Alternative	Avoid
Max_temp_C	MaxTemp	Maximum Temp (°C)
Precipitation_mm	Precipitation	precmm
Mean_year_growth	MeanYearGrowth	Mean growth/year
sex_mf	sex	M/F
weight_kg	weight	w.
cell_type	CellType	Cell Type
Observation_01	first_observation	1st Obs

Common Spreadsheet Errors #5

Special characters in data

- **Example:** You treat Excel as a word processor, even copying data directly from Word or other applications.
- **Solution:** This is a **common strategy**. For example, when writing longer text in a cell, people often include line breaks, em-dashes, et al in their spreadsheet.
- Worse yet, when copying data in from applications such as Word, formatting and fancy non-standard characters are included.
- **General best practice is to avoid adding characters such as newlines, tabs, and vertical tabs.** In other words, treat a text cell as if it were a simple web form that can only contain text and spaces.

Common Spreadsheet Errors #6

Not filling in zeroes

- It might be that when you're measuring something, it's usually a zero, say the number of times an elephant is observed in the object or the survey. Why bother writing in the number zero in that column, when it's mostly zeros?
- There's a **difference between a zero and a blank cell** in a spreadsheet. To the computer, a zero is actually data. You measured or counted it. A blank cell means that it wasn't measured and the computer will interpret it as a null value.

Table 1. Commonly used null values, limitations, compatibility with common software and a recommendation regarding whether or not it is a good option. Null values are indicated as compatible with specific software if they work consistently and correctly with that software. For example, the null value "NULL" works correctly for certain applications in R, but does not work in others, so it is not presented in the table as R compatible.

Null values	Problems	Compatibility	Recommendation
0	Indistinguishable from a true zero		Never use
Blank	Hard to distinguish values that are missing from those overlooked on entry. Hard to distinguish blanks from spaces, which behave differently.	R, Python, SQL	Best option
-999, 999	Not recognized as null by many programs without user input. Can be inadvertently entered into calculations.		Avoid
NA, na	Can also be an abbreviation (e.g., North America), can cause problems with data type (turn a numerical column into a text column). NA is more commonly recognized than na.	R	Good option
N/A	An alternate form of NA, but often not compatible with software		Avoid
NULl	Can cause problems with data type	SQL	Good option
None	Uncommon. Can cause problems with data type	Python	Avoid
No data	Uncommon. Can cause problems with data type, contains a space		Avoid
Missing	Uncommon. Can cause problems with data type		Avoid
-,+,. .	Uncommon. Can cause problems with data type		Avoid

Let's try it out!

- Download a free online open data set
 - <https://open.nrw/dataset/ldbnrw-service-62411-14i>
 - Choose the CSV format
- Import it into Google Spreadsheets
- Audit the data format and find common mistakes
- Correct the data format
- Export as CSV