# Data Cleansing

Dr Pierre Le Bras

# About Me

2011 Started to study Software Engineering
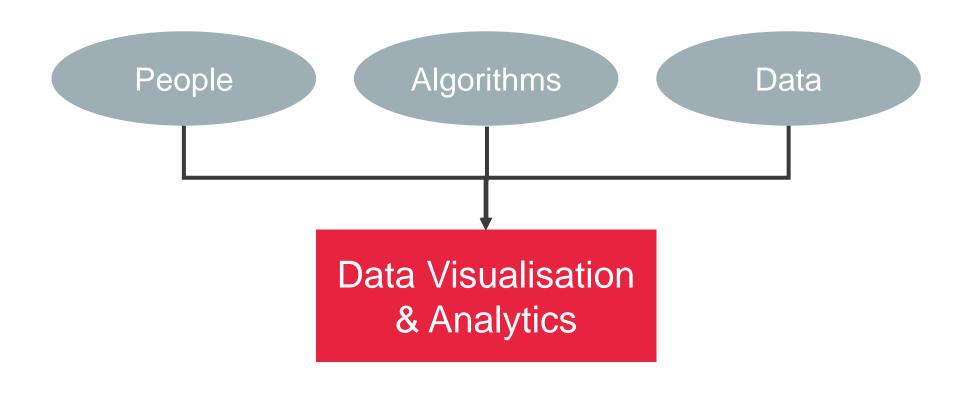
2013 Moved to Edinburgh to study Computer Science

2015 Started my PhD in Data Visualisation and Analytics

2018 Worked as a Research Associate

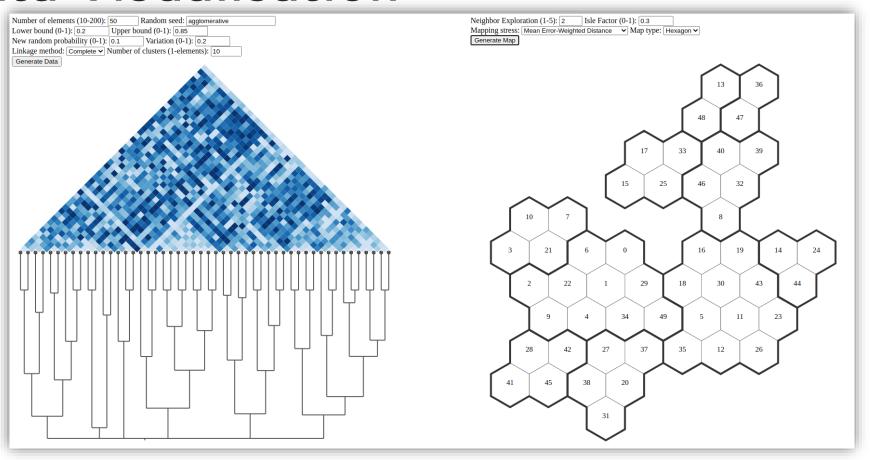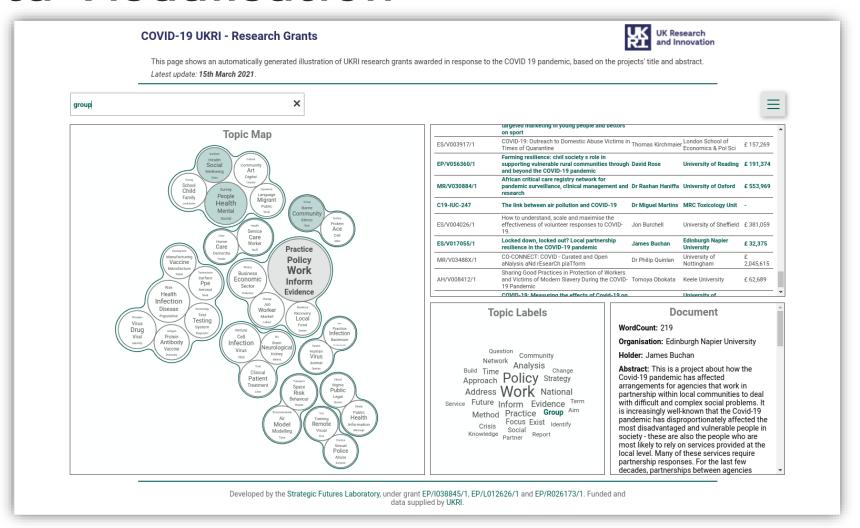2021 Lecturer in Software Engineering and Data Science

# Data Visualisation

# Data Visualisation
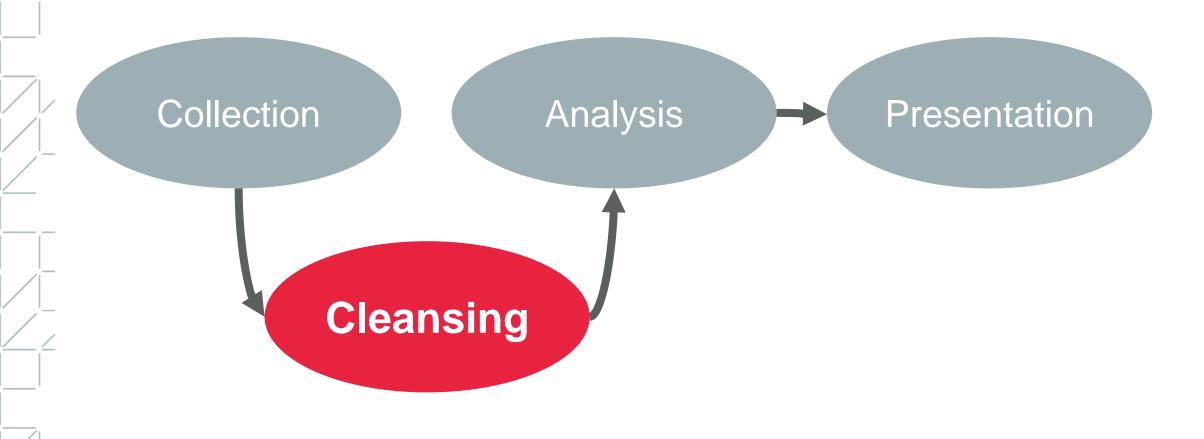
# Data Visualisation

# Processing Data



Collection → Analysis → Presentation

# Processing Data
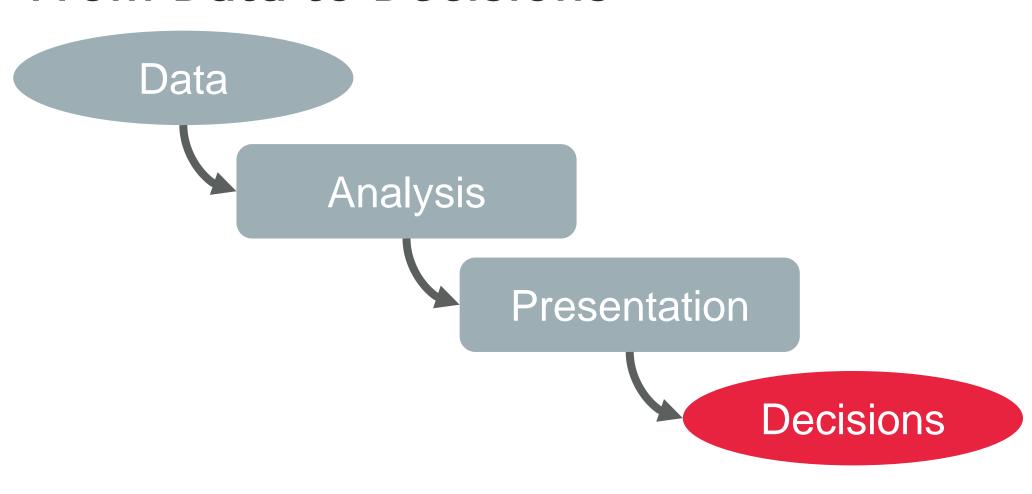
# Why do you want Clean Data?

# Why Investing in Data?

Data is used by organisations or individuals seeking **informed** **decision making**

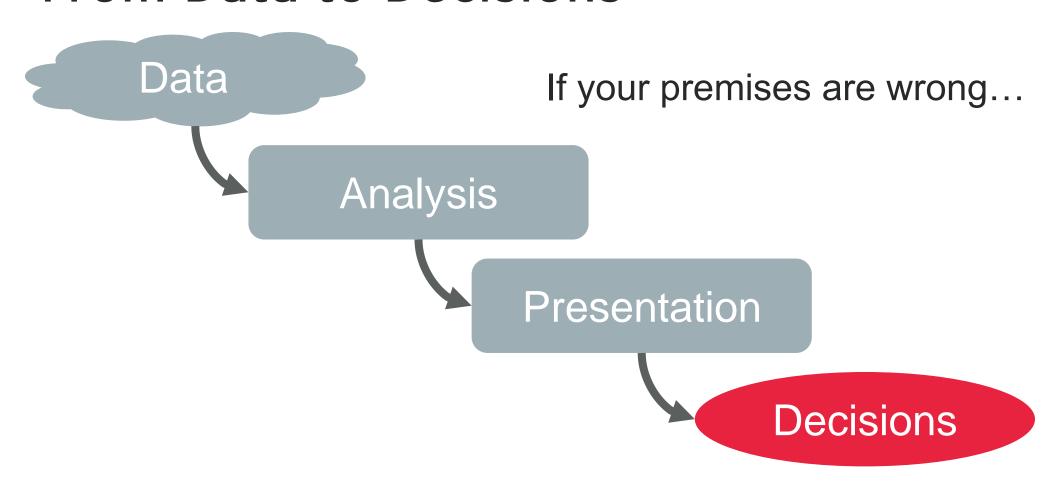# From Data to Decisions

# From Data to Decisions

Data

If your premises are wrong…

Analysis

Presentation

Decisions

# From Data to Decisions

Data

If your premises are wrong…

Analysis

Presentation

… your conclusions will be flawed!

Decisions

# Impacts of Dirty Data

- Bad data lead to inaccurate results and false insights

- ... which lead to the wrong decisions being taken

- ... which can have serious consequences

- For example:
  - Wrong census data and public funds: clickondetroit.com/news/michigan/2021/10/15/census-analysis-shows-black-americans-may-have-been-significantly-undercounted/

# Where does it get dirty?

# Human Error

- At an individual level
    - Incorrect or missed entry
    - Mistake when copying data

- At an organisational level
    - Poor data management policies
    - Poor communication between services
    - Poor documentation

# System Failures

- Hardware and/or software malfunction
  - Data file corruption
  - Loss of data

# Collection Methods

- Collecting and merging data from multiple sources
  - With likely different formats

- Unforeseen exceptions during automatic collections
  - E.g., web crawlers and weird HTML structures

# What is Clean Data?

# Dirty Data can be...

- Duplicated
- Out of date
- Inconsistent
- Incorrect
- Invalid

- Unformatted
- Non-uniform
- Inaccurate
- Missing
- Incomplete

# 5 Rules to Data Quality

- Validity
- Accuracy
- Completeness
- Consistency
- Uniformity

# Validity

- How much the data conforms to a set of rules

- Plays at the individual entries and fields level

|  | Field 1 | Field 2 | Field 3 | Field 4 | … |
|---------|:---:|:---:|:---:|:---:|:---:|
| Entry 1 | ✓ | ✓ | ✓ | ✓ | ✓ |
| Entry 2 | ✓ | ✓ | ✓ | ✓ | ✓ |
| Entry 3 | ✓ | ✓ | ✓ | ✓ | ✓ |
| … | ✓ | ✓ | ✓ | ✓ | ✓ |

# Validity

- Each field has correct types (String, Number, Date, Boolean, …)

- Numbers and dates must fall within the right ranges

- Discrete values must be consistent
  - E.g. a female gender value must be "Female" or "F" across all entries but not both

- Discrete values must be from a finite set
  - E.g. UK regions can only be "England", "Wales", "Scotland" and "Northern Ireland"

# Validity

- Text values must have the right pattern
  - E.g. phone numbers must have a "XXXX-XXX-XXXX" pattern
- Each entry should be unique
  - A field or combination of fields must be defined as unique key
- Mandatory fields are not empty *
- Related fields are consistent *
  - E.g. a birth date and age fields should be coherent (age = current date – birth date)

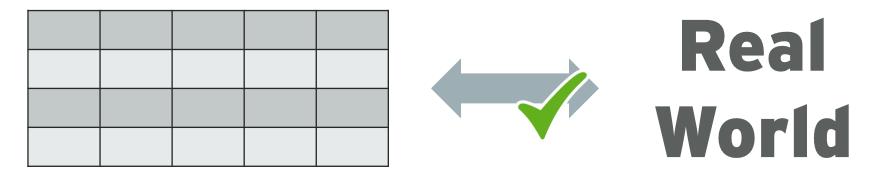\* These requirements fit with other data quality principles

# Accuracy

- How well the data fits with the truth

- Difficult to achieve, requires an external "gold standard" dataset

- Accurate data can become inaccurate (out-of-date)

**Real World**

# Completeness

- How little required data is missing

- Difficult to fix, "*Missing*", "*Unknown*" or "*NA*" value might be needed → **Missing data is still data**

# Consistency

- How much the data values agree or are coherent

- To fix: check which is most up-to-date, which source is most reliable, etc.

# Uniformity

- How units of measures are similar across fields

# Other Data Quality Aspects

- Comparability
- Relevance
- Credibility
- Currency (up-to-date)
- Confidentiality

http://www.dama-nl.org/data_quality/

# How to Clean Data?
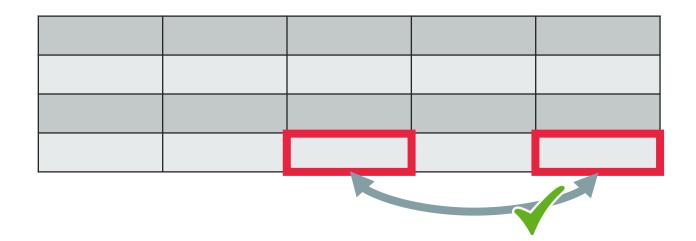
# Step 1 – Inspection

- Exploring the data and detecting unexpected, incorrect and/or inconsistent entries or fields.

- Data Profiling:
  - Are values corresponding to their field's format and pattern specification?
  - Are there missing fields?
  - Are the mean, median, range, sum, etc. coherent?
  - Is the distribution of values coherent?
  - Are there outliers?
  - Is there duplicate entries?

- Visualisation

# Step 2 - Cleansing

- Clear irrelevant data
- Handle duplicates
- Handle missing values
- Convert types
- Fix the content
- Fix the format
- Scale and/or normalise numbers

# Step 2 - Cleansing

**Drop irrelevant data**: remove fields or entries which are not needed for our analysis

- *Column-wise*: fields which provide information outside of our task domain, e.g., phone numbers in an analysis of student grades
- *Row-wise*: entries with attributes outside of our task domain, e.g., postgraduate students in an analysis of undergraduate student grades

- **Only delete data if you are sure about it!**
  - You might need it later
  - It might make some interesting correlation analysis

# Step 2 - Cleansing

**Remove duplicates**

- *Column-wise*: a field repeating another field
- *Row-wise*: duplicate entries
  - Identify a field or combination of fields as unique keys
  - Flag and remove duplicates based on the keys

# Step 2 - Cleansing

**Check and correct missing values**

- Drop: if a field or entry is mostly empty, it might be worth removing
- Complete: some values can be filled, using another dataset, doing manual search, etc.
- Infer: some values may be inferred from other observation (e.g., with linear regression)
- Flag: **missing data can be data!**
  - Numerical data can be set to absurd values: *-1* or *0*
  - Categorical data can be set to a new category: *"missing"* or *"NA"*
  - *"missing"*, *"NA"* or *"unknown"* have different meaning

# Step 2 - Cleansing

**Convert values to their correct types**

- Dates should be date objects (Unix timestamp)
- Numerical data should be numbers
- Categorical data can be strings or numbers
- Text data should be strings
- Boolean values (true/false, yes/no) should be Booleans

# Step 2 - Cleansing

**Fix the content:** Correct syntax, mistakes and conversion errors

- Strings: trim whitespaces, fix typos
    - Easy on categorical data
    - Difficult with free text data (tools can help)

- Numbers: fix errors (from your inspection)
    - Can use the same approaches as with missing data (drop / complete / infer / flag)

# Step 2 - Cleansing

**Fix the format:** Standardize the values across your dataset

- Strings:
  - Patterns (e.g., XXXX-XXX-XXXX for phone numbers)
  - Format (e.g., lower case)

- Numbers:
  - Units (e.g., distances, weights, currencies)
  - Format (e.g., zero-padding, decimal digits)

- Date:
  - Format (e.g. YYYY-MM-DD or DD/MM/YY)

- Categories:
  - Values (e.g. *"F"* or *"female"*, *"Sat."* or *"saturday"*)

# Step 2 - Cleansing

**Scale and/or normalise number ranges**

- Scaling
    - Making your numbers comparable, on the same scale
- Normalisation
    - Making your number distribution normally distributed (for later statistical analysis)

# Step 3 - Verification

- Re-inspect the data to make sure of its correct cleaning

- Correct mistakes (manually) if needed

# Step 4 - Report

- Reporting your cleansing is important

- It helps measure data quality, highlight issues and identify common problems

- It can inform future data practices
  - Collection (e.g., survey question often unanswered, web crawler cropping text, etc.)
  - Storage (e.g., mismatch between datasets)

# Data Cleansing Tools

# Programming

# Software

# Software


OpenRefine

- Free
- Open Source
- Cross Platform
- Design for Cleansing
- Powerful

# Next…

- Questions?

- Break

- Live demonstration