



THE UNIVERSITY OF
MELBOURNE

Machine Learning Applications for Health

COMP90089 (2022) - Lecture 3

Dr Brian Chapman
brian.chapman@unimelb.edu.au

Dr Daniel Capurro
dcapurro@unimelb.edu.au





Today's Agenda: Data Sources

- Sources of clinical data
- Some openly available datasets
- MIMIC

Brainstorming: sources of clinical data

what data:

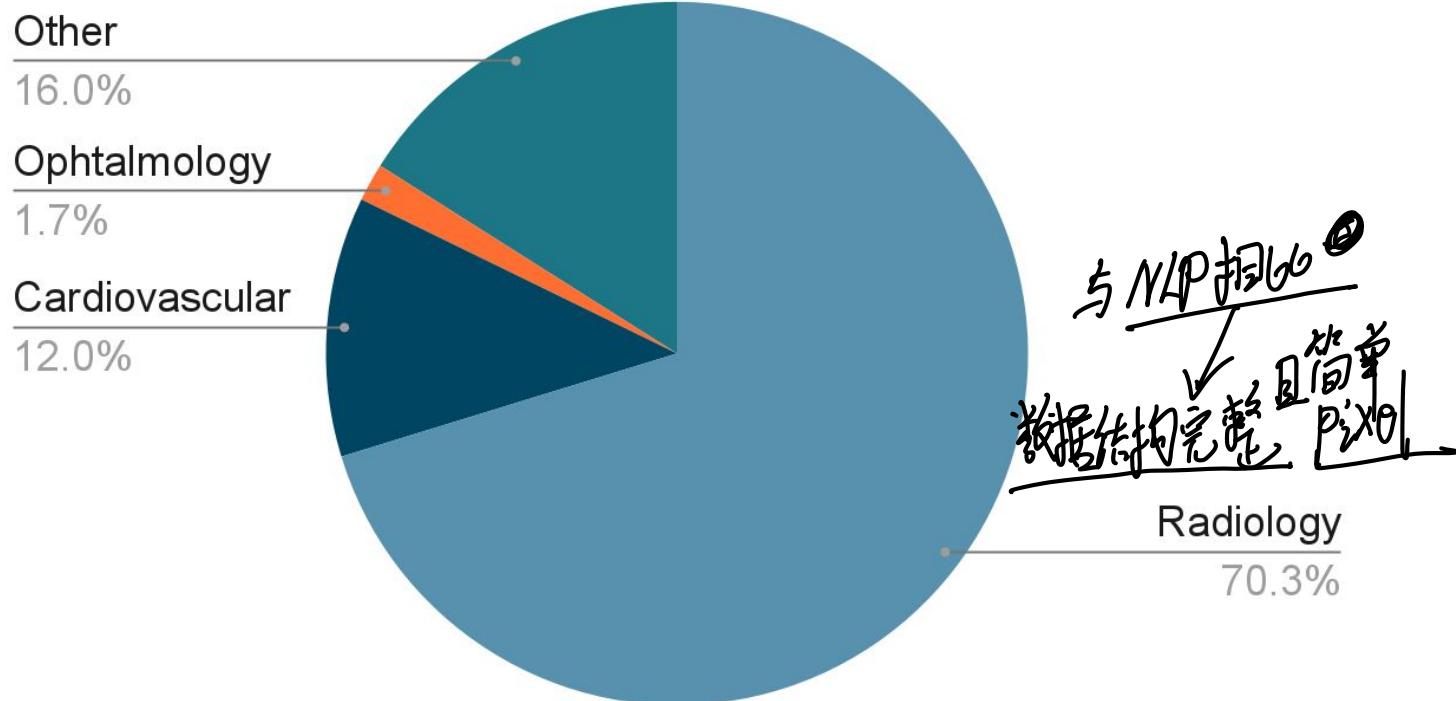
- test
- images
- discharge notes
- insurance claims



Which areas of Medicine have made the greatest progress in developing/adopting AI?

Which areas of Medicine have made the greatest progress in adopting AI?

FDA Approvals

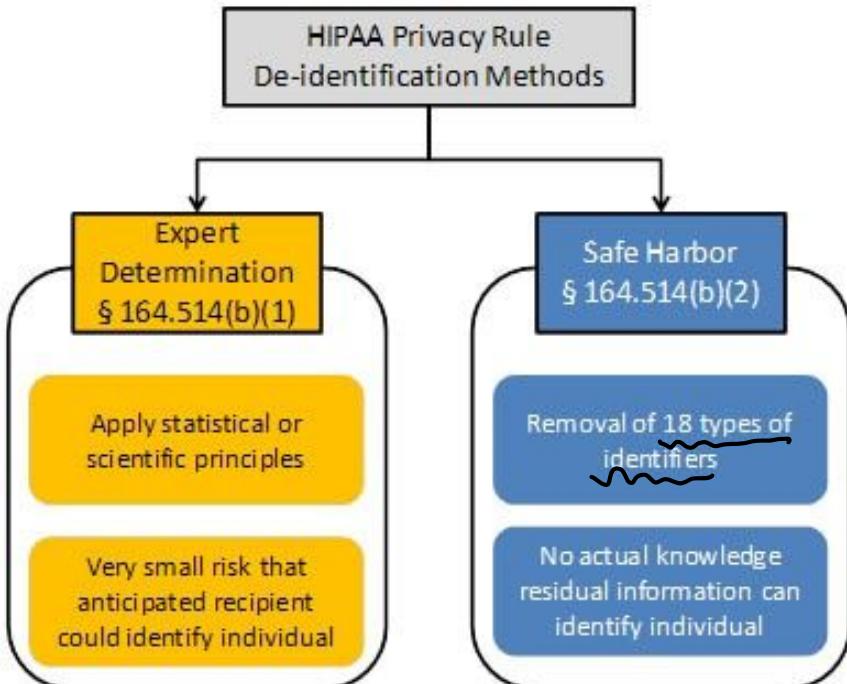




Publicly available data sources

- Not in Australia
- Most publicly datasets are from the USA - HIPAA

In the USA: Health Information Portability and Accountability Act (HIPAA)





In the USA: Health Information Portability and Accountability Act (HIPAA)

- Names
- All geographic subdivisions smaller than a state
- All elements of dates (except year) for dates that are directly related to an individual, including birth date, admission date, discharge date, death date, and all ages over 89 and all elements of dates (including year) indicative of such age, except that such ages and elements may be aggregated into a single category of age 90 or older
- Telephone numbers, Vehicle identifiers and serial numbers, including license plate numbers, Fax numbers, device identifiers and serial numbers, Social security numbers, Account numbers, Health plan beneficiary numbers, certificate/license numbers, Any other unique identifying number, characteristic, or code
- Email addresses, Web Universal Resource Locators (URLs), Internet Protocol (IP) addresses, Medical record numbers,
- Biometric identifiers, including finger and voice prints, Full-face photographs and any comparable images

The De-Identification Decision-Making Framework

Christine M O'Keefe, Stephanie OtoРЕЕС,
Mark Elliot, Elaine Mackey and Kieron O'Hara

18 September 2017



Australian Government

Office of the Australian Information Commissioner

- Full de-identification cannot be guaranteed
- Context dependent



THE SIMPLE PROCESS OF RE-IDENTIFYING PATIENTS IN PUBLIC HEALTH RECORDS

In late 2016, doctors' identities were decrypted in an open dataset of Australian medical billing records. Now patients' records have also been re-identified - and we should be talking about it

By Dr Vanessa Teague, Dr Chris Culnane and Dr Ben Rubinstein, University of Melbourne

ENGINEERING & TECHNOLOGY

Featured

In August 2016, Australia's federal Department of Health published medical billing records of about 2.9 million Australians online. These records came from the Medicare Benefits Scheme (MBS) and the Pharmaceutical Benefits Scheme (PBS)

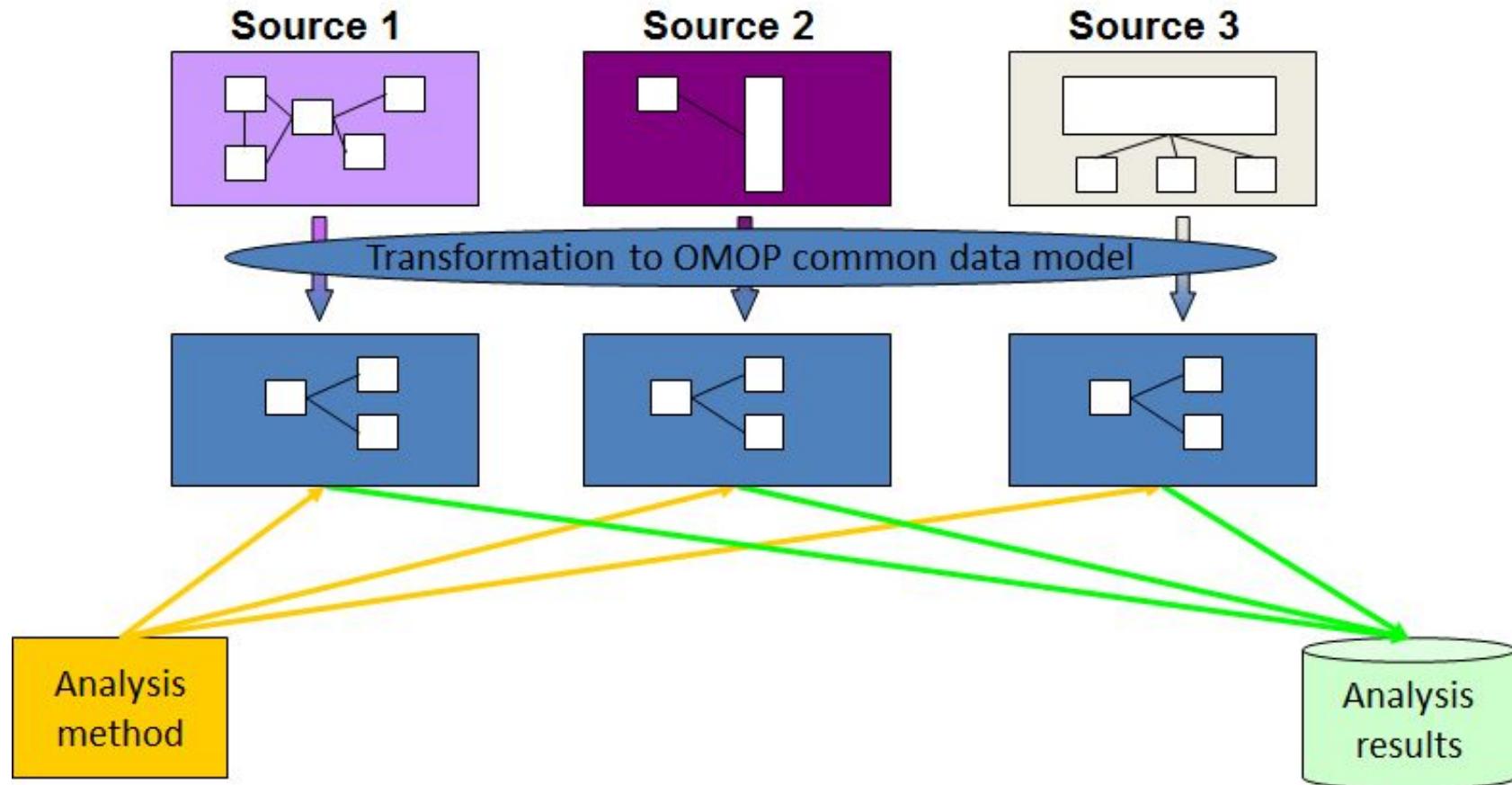
containing 1 billion lines of historical health data from the records of around 10 years



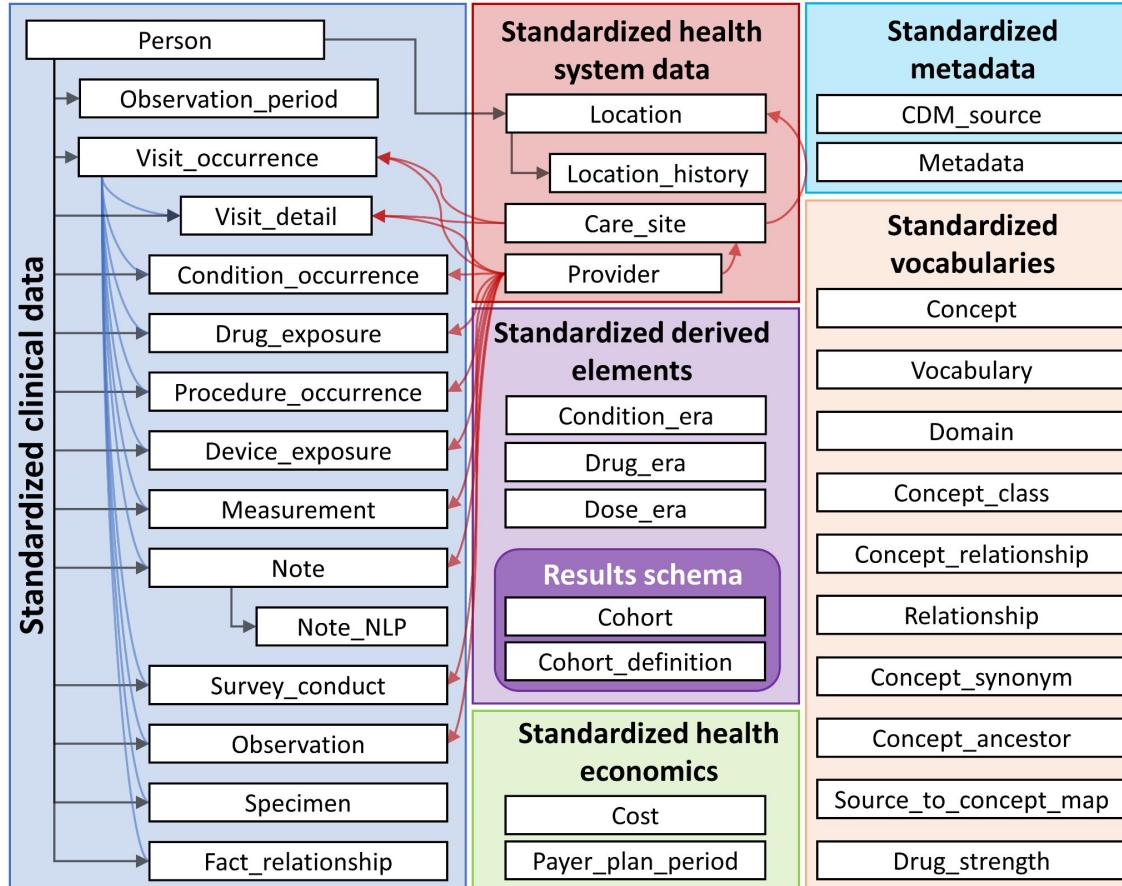
Some datasets available

- MIMIC and its companion datasets in Physionet: MANY!!
<https://physionet.org/about/database/#open>
- US Medicare SynPUF -
<https://www.cms.gov/Research-Statistics-Data-and-Systems/Downloadable-Public-Use-Files/SynPUFs>
- Synthetic data generator - <https://synthetichealth.github.io/synthea/>
- Google public data explorer: <https://www.google.com/publicdata/directory>
- Previous Kaggle challenges:
<https://www.kaggle.com/datasets?search=health>

Common Data Models



OMOP CDM





Project using OMOP

U.S. Department of Health & Human Services > National Institutes of Health > National Center for Advancing Translational Sciences > Enclave Sign In | Create Account | Profile

N3C ABOUT ▾ N3C DATA ENCLAVE ▾ GET INVOLVED ▾ RESOURCES ▾ SUPPORT ▾



National COVID Cohort Collaborative

The **N3C Data Enclave** is a secure platform through which the harmonized clinical data provided by our contributing members is stored. The data itself can only be accessed through a secure cloud portal hosted by NCATS and cannot be downloaded or removed. N3C invites you to begin your journey with the Enclave and join the collaborative efforts of our partners to better understand and address the most pressing COVID-19 clinical questions.

Access the Enclave

Help make science go faster and save lives.

17.3B <i>Total Rows</i>	1,616.3M <i>Clinical Observations</i>	14.8M <i>Persons</i>	5,815,680 <i>COVID+ Cases</i>
----------------------------	--	-------------------------	----------------------------------

Explore the Full Cohort Dashboard ↗

ASK N3C

Data Training Discover Contributors Support & Contact Terms & Conditions Logout



Project using OMOP

U.S. Department of Health & Human Services > National Institutes of Health > National Center for Advancing Translational Sciences > Enclave Sign In | Create Account | Profile

N3C ABOUT ▾ N3C DATA ENCLAVE ▾ GET INVOLVED ▾ RESOURCES ▾ SUPPORT ▾



National COVID Cohort Collaborative

The N3C Data Enclave is a secure platform through which the harmonized clinical data provided by our contributing members is stored. The data itself can only be accessed through a secure cloud portal hosted by NCATS and cannot be downloaded or removed. N3C invites you to begin your journey with the Enclave and join the collaborative efforts of our partners to better understand and address the most pressing COVID-19 clinical questions.

[Access the Enclave](#)

Help make science go faster and save lives.

17.3B <i>Total Rows</i>	1,616.3M <i>Clinical Observations</i>	14.8M <i>Persons</i>	5,815,680 <i>COVID+ Cases</i>
----------------------------	--	-------------------------	----------------------------------

Explore the Full Cohort Dashboard ↗

ASK N3C

Data Training Discover Contributors Support & Contact Terms & Conditions Logout



Project using OMOP

RESEARCHER LOGIN

All of Us Research Hub | NIH National Institutes of Health All of Us Research Program

ABOUT DATA & TOOLS DISCOVER FAQ REGISTER

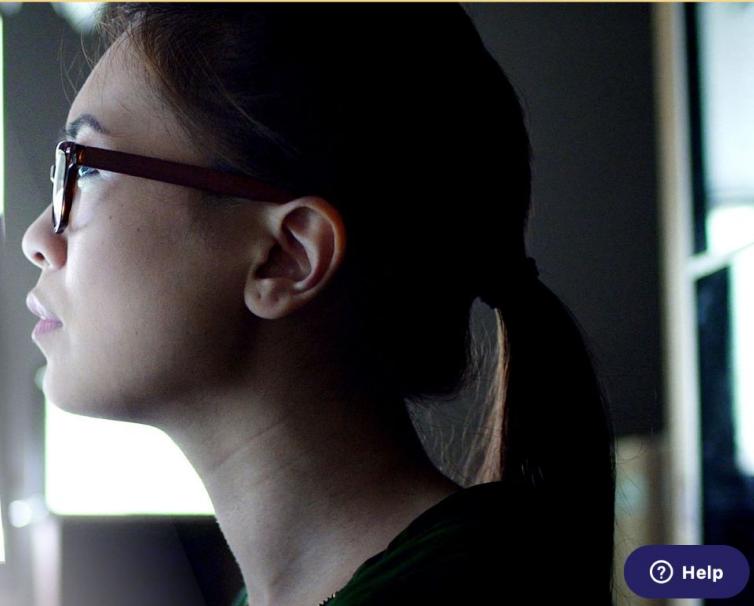
Learn About Our New Controlled Tier >> Explore Genomic Data in the Data Browser >>

Welcome to the *All of Us* Research Hub

The *All of Us* Research Program, led by the National Institutes of Health, is building one of the largest biomedical data resources of its kind. The *All of Us* Research Hub stores health data from a diverse group of participants from across the United States.

Registered researchers can access *All of Us* data and tools to conduct studies to help improve our understanding of human health.

REGISTER FOR ACCESS Help



Project using OMOP



OHDSI

OBSERVATIONAL HEALTH DATA SCIENCES AND INFORMATICS

Welcome to OHDSI Australia!!

The Observational Health Data Sciences and Informatics (or OHDSI, pronounced "Odyssey") program is a multi-stakeholder, interdisciplinary collaborative to bring out the value of health data through large-scale analytics. All our solutions are open-source. OHDSI Australia is a newly formed Australian chapter.

Background

The establishment of OHDSI Australia has been facilitated by close cooperation with the Transformational Data Collaboration (TDC) <https://machaustralia.org/projects/transformational-data-collaboration/>

The TDC is an initiative under the auspices of the Australian Health Research Alliance <https://ahra.org.au/>

Under the 'Data Integration' priority area of the 'AHRA Data Driven Healthcare' activity stream. It has a singular goal:

"To utilise the unique open and collaborative nature of AHRA to help develop and support national data initiatives where an open, inclusive and non-competitive environment is required."

Latest News:

EHDEN Academy free access to all here:
<https://academy.ehdeneu>

Next Event:

27th July 2021 at 1pm AEST

"ETL Framework for the Conversion of Health Databases to OMOP"
by Dr Juan Quiroz

See [events page for details](#)

[Recorded webinars here](#)



Brief Introduction to MIMIC

Some History

- Originated at the Laboratory for Computational Physiology at MIT, originally focused on physiological signal processing.
- Established Physionet, a repository of physiological data available for research.
- Created MIMIC repository: Medical Information Mart for Intensive Care
- MIMIC-II released in 2010
- MIMIC-III released in 2016
- MIMIC-IV released in 2020

What is MIMIC?

- Data mart
- One US hospital
- Critical care patient information
- Fully anonymized
- Available for research and quality improvement

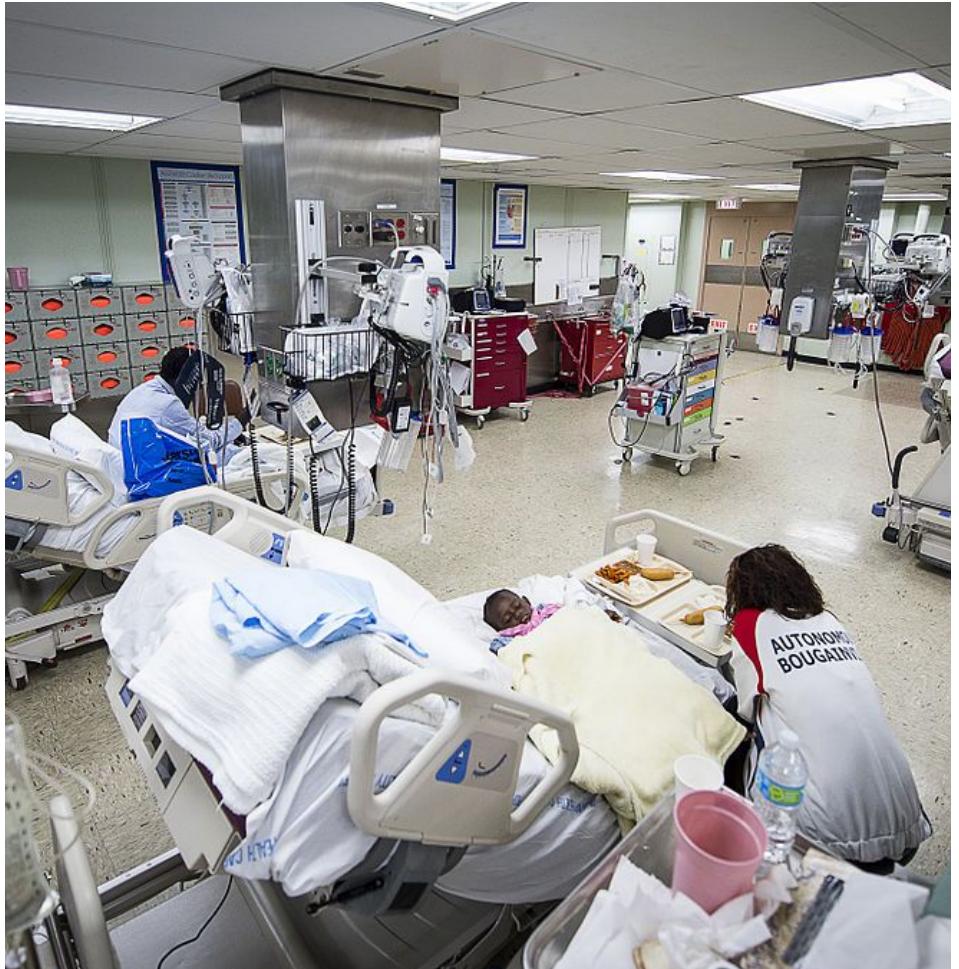
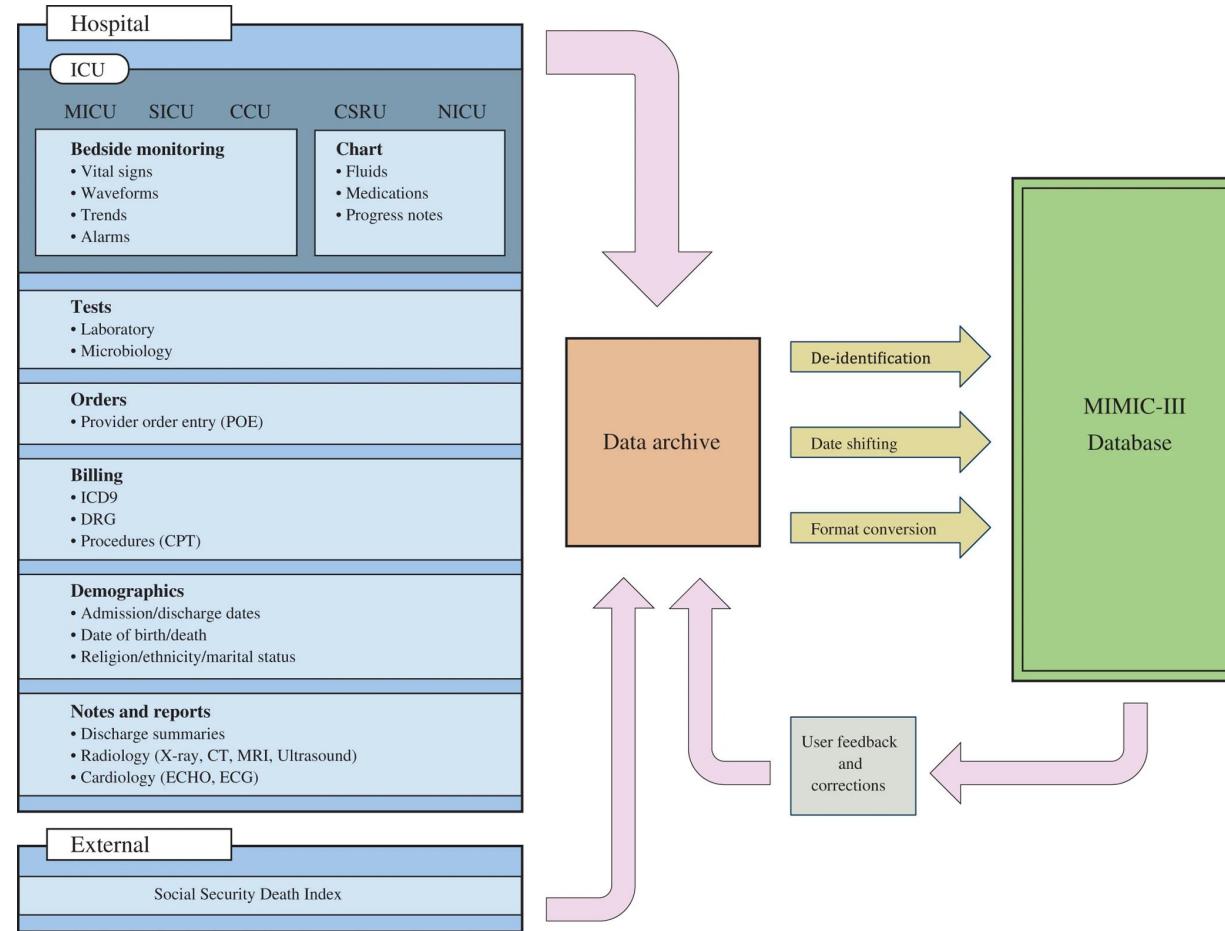


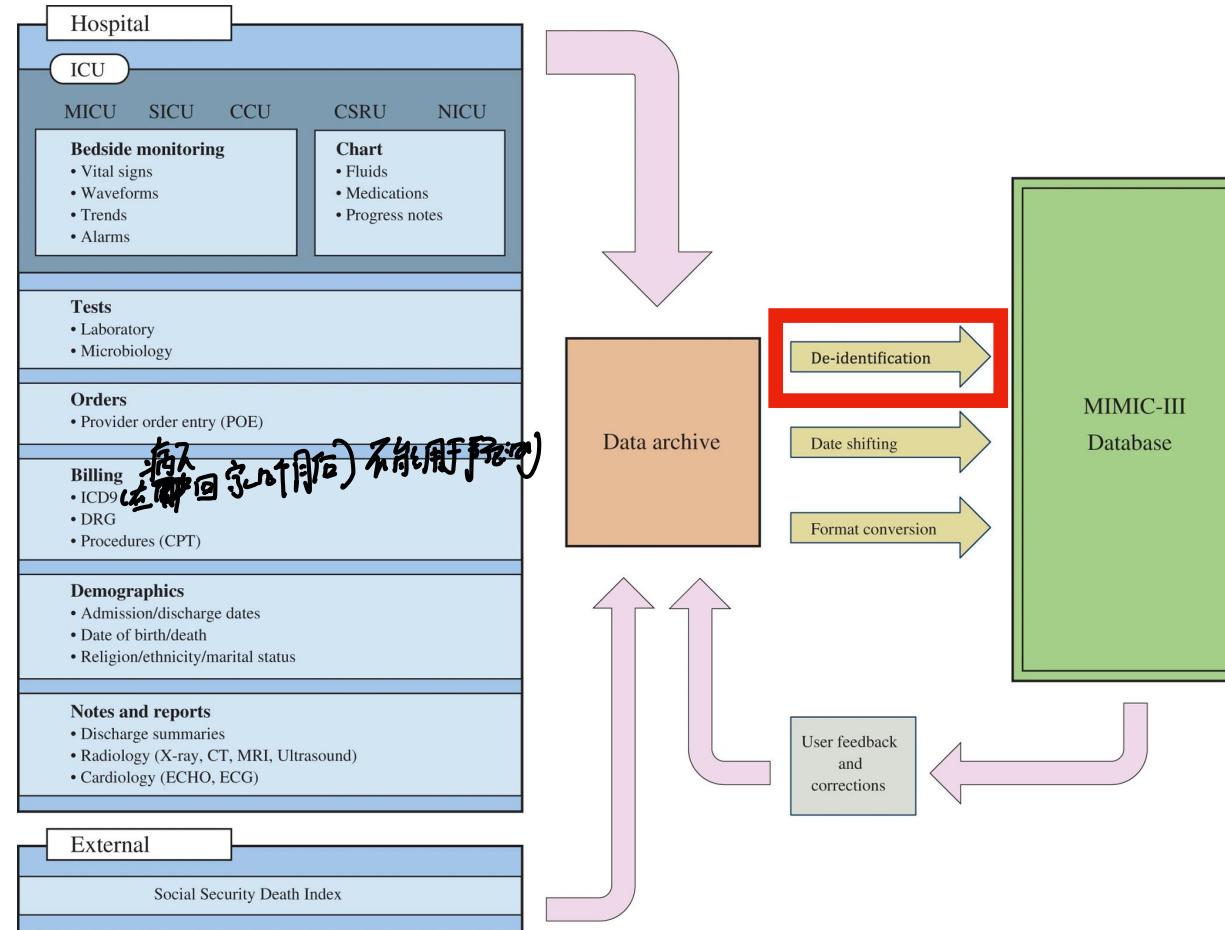
Photo by: Tech. Sgt. Araceli Alarcon,

https://commons.wikimedia.org/wiki/File:Bougainville_locals_returned_after_receiving_critical_care_on_the_Mercy_150701-F-SD522-021.jpg

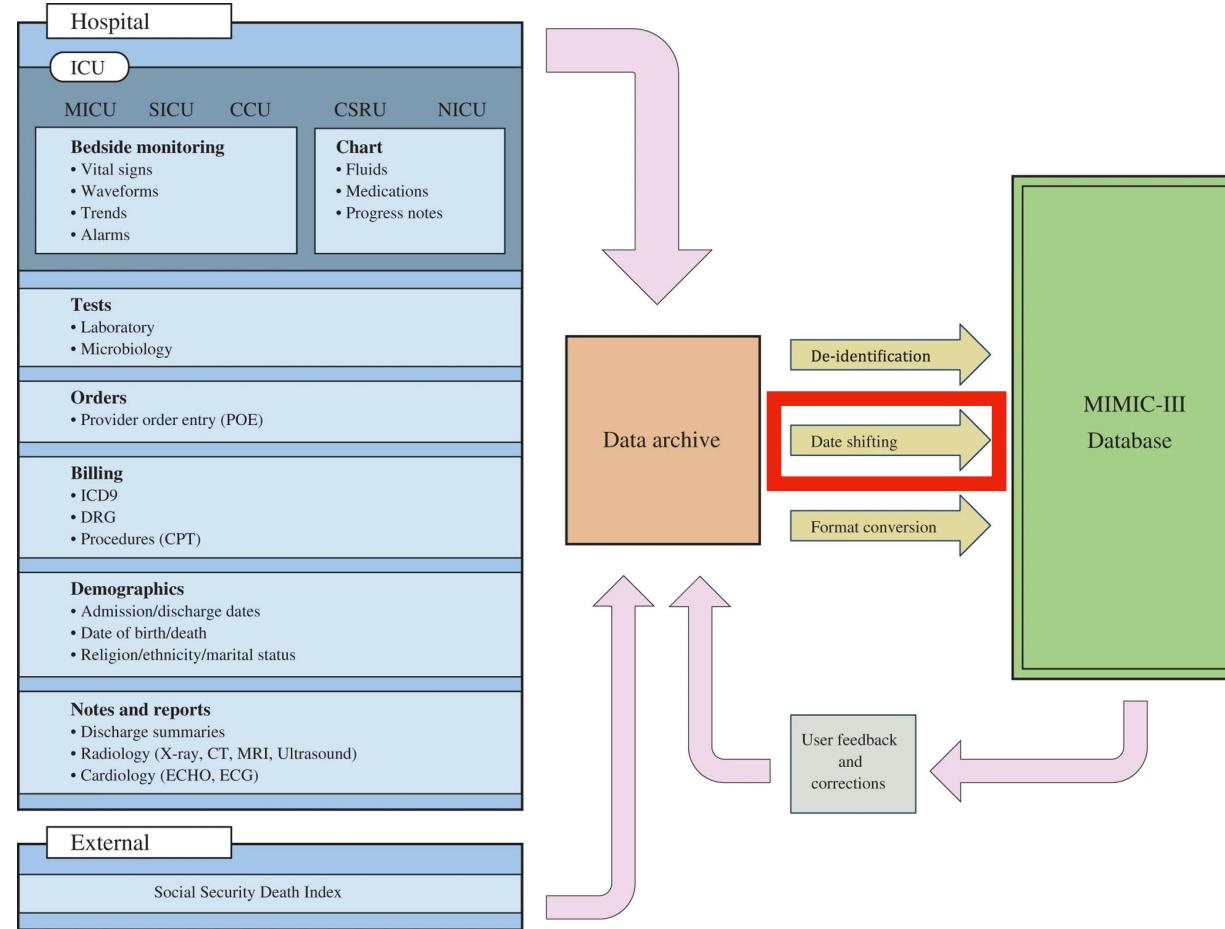
How is the information generated?



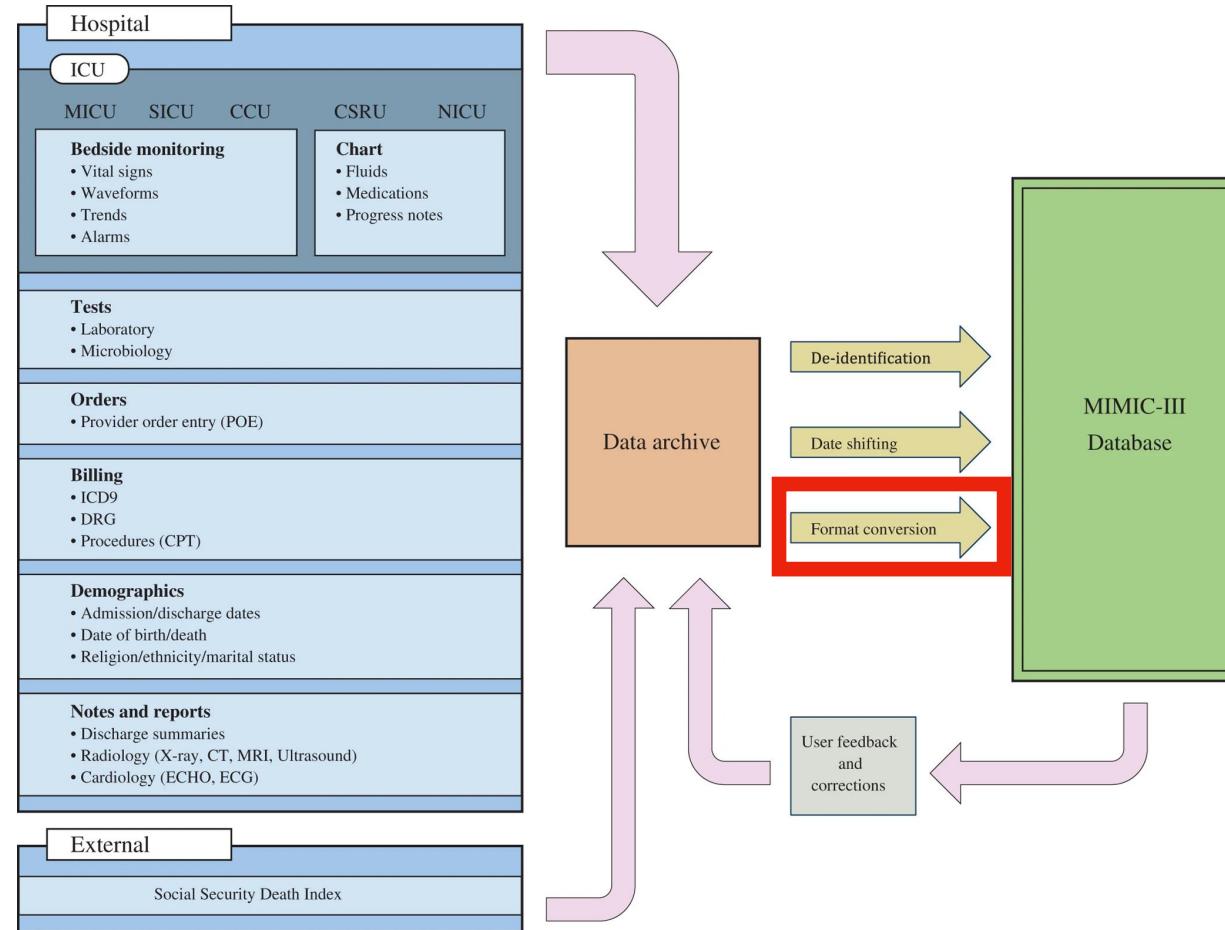
How is the information generated?



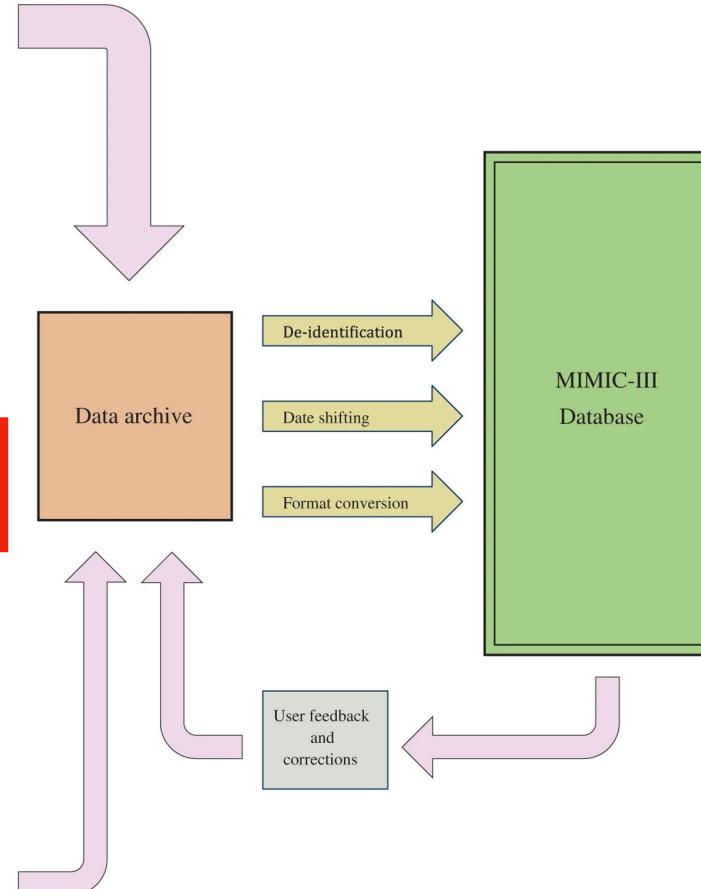
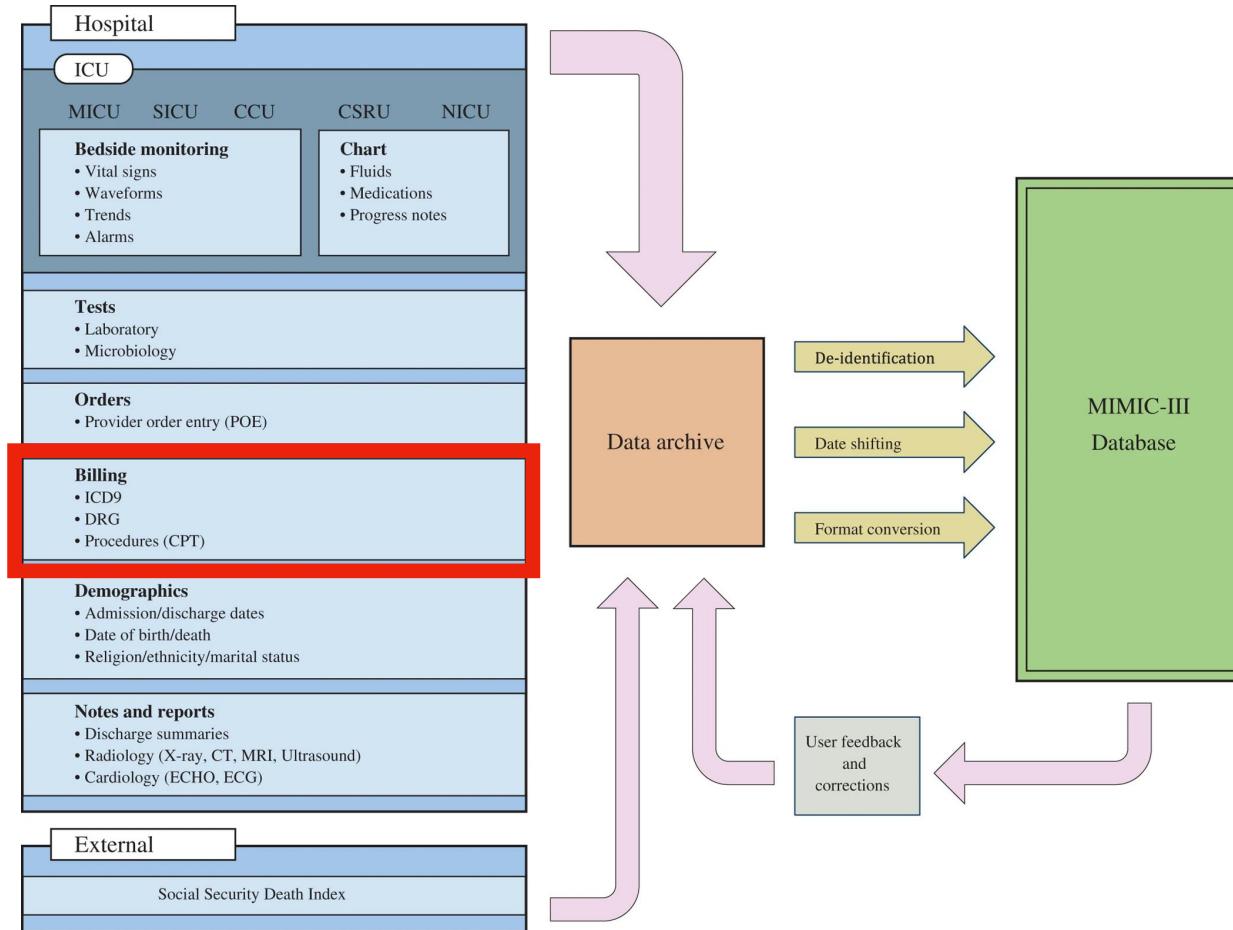
How is the information generated?



How is the information generated?



How is the information generated?



(a bit about the) Database Schema

Details available in: <https://mimic.mit.edu/docs/iv/modules/>

Patient

SUBJECT_ID	GENDER	DOB	
DOD	DOD_HOSP	DOD_SSN	EXPIRE_FLAG

Admissions

SUBJECT_ID	HADM_ID	
ADMITTIME	DISCHTIME	DEATHTIME
ADMISSION_TYPE	ADMISSION_LOCATION	
INSURANCE	LANGUAGE	RELIGION
MARITAL_STATUS	ETHNICITY	
EDREGTIME	EDOUTTIME	DIAGNOSIS
HOSPITAL_EXPIRE_FLAG		

ICU Stays

SUBJECT_ID	HADM_ID	ICUSTAY_ID	
DBSOURCE			
FIRST_CAREUNIT	LAST_CAREUNIT		
FIRST_WARDID	LAST_WARDID		
INTIME	OUTTIME	LOS	

Important to know



Some relevant tables

Diagnoses_ICD

Procedures_ICD

SUBJECT_ID , HADM_ID SEQ_NUM
ICD9_CODE

Labevents

SUBJECT_ID , HADM_ID ITEMID
CHARTTIME VALUE , VALUENUM
VALUEUOM FLAG

Complex \rightarrow sensitivity of bacteria

Mircobiology Events

SUBJECT_ID , HADM_ID
CHARTDATE , CHARTTIME
SPEC_ITEMID , SPEC_TYPE_DESC
ORG_ITEMID , ORG_NAME ISOLATE_NUM
AB_ITEMID , AB_NAME
DILUTION_TEXT , DILUTION_COMPARISON ,
DILUTION_VALUE
INTERPRETATION

Prescriptions *describe doctor Presay to patients*

SUBJECT_ID , HADM_ID , ICUSTAY_ID
STARTDATE , ENDDATE DRUG_TYPE
DRUG , DRUG NAME POE ,
DRUG_NAME_GENERIC
FORMULARY_DRUG_CD , GSN , NDC
PROD_STRENGTH
DOSE_VAL_RX , DOSE_UNIT_RX
FORM_VAL_DISP , FORM_UNIT_DISP ROUTE

Codes and Descriptions

diagnoses_ICD

row_id [PK] integer	subject_id integer	hadm_id integer	seq_num integer	icd9_code character varying (10)
12350	1062	105525	1	85200
12351	1062	105525	2	43491
12352	1062	105525	3	0389
12353	1062	105525	4	51881
12354	1062	105525	5	2639
12355	1062	105525	6	48283
12356	1062	105525	7	99592

one patient many diagnoses

d_diagnoses_icd

row_id [PK] integer	icd9_code character varying (10)	short_title character varying (50)	long_title character varying (255)
1	174	01166	TB pneumonia-oth test
2	175	01170	TB pneumothorax-unspec
3	176	01171	TB pneumothorax-no exam
4	177	01172	TB pneumothorax-exam unkn
5	178	01173	TB pneumothorax-micro dx
6	179	01174	TB pneumothorax-outl dx

row_id integer	subject_id integer	hadm_id integer	seq_num integer	icd9_code character varying (10)	row_id integer	icd9_code character varying (10)	short_title character varying (50)	long_title character varying (255)
1	12350	1062	105525	1	85200	8394	85200	Traum subarachnoid hem
2	12351	1062	105525	2	43491	5053	43491	Crbl art ocl NOS w infrc
3	12352	1062	105525	3	0389	660	0389	Septicemia NOS
4	12353	1062	105525	4	51881	5279	51881	Acute respiratory failure
5	12354	1062	105525	5	2639	1642	2639	Protein-cal malnutr NOS
6	12355	1062	105525	6	48283	5514	48283	Pneumo oth arm-peo bact

The same happens with procedures, labs, chartevents, inputevents, and others

Some issues with time

- Dates are anonymized, that's why you see Star Trek years
- Dates are consistent within patients, not between patients
- Some elements are generated at the END of a hospitalization (Diagnoses, some procedure codes, DRGs)
- MIMIC now includes patients in the Emergency Department that were not in the ICU (no hadm_id)

↓
in hospital, not in ICU



THE UNIVERSITY OF
MELBOURNE

Machine Learning Applications for Health

COMP90089 (2022) - Lecture 3

Dr Brian Chapman

brian.chapman@unimelb.edu.au

Dr Daniel Capurro

dcapurro@unimelb.edu.au



Introduction to Data Wrangling

BE Chapman, PhD

2023-08-02

Housekeeping

- MIMIC access
 - Have you finished your training?
- Groups
 - Have you submitted your signed group contract?
 - Please let us know of any issues/difficulties connecting with peers
- Student representatives
 - Clemence Mottez (cmmottez@student.unimelb.edu.au)
 - Tanvesh Takawale (ttakawale@student.unimelb.edu.au)

Generating research-ready clinical data

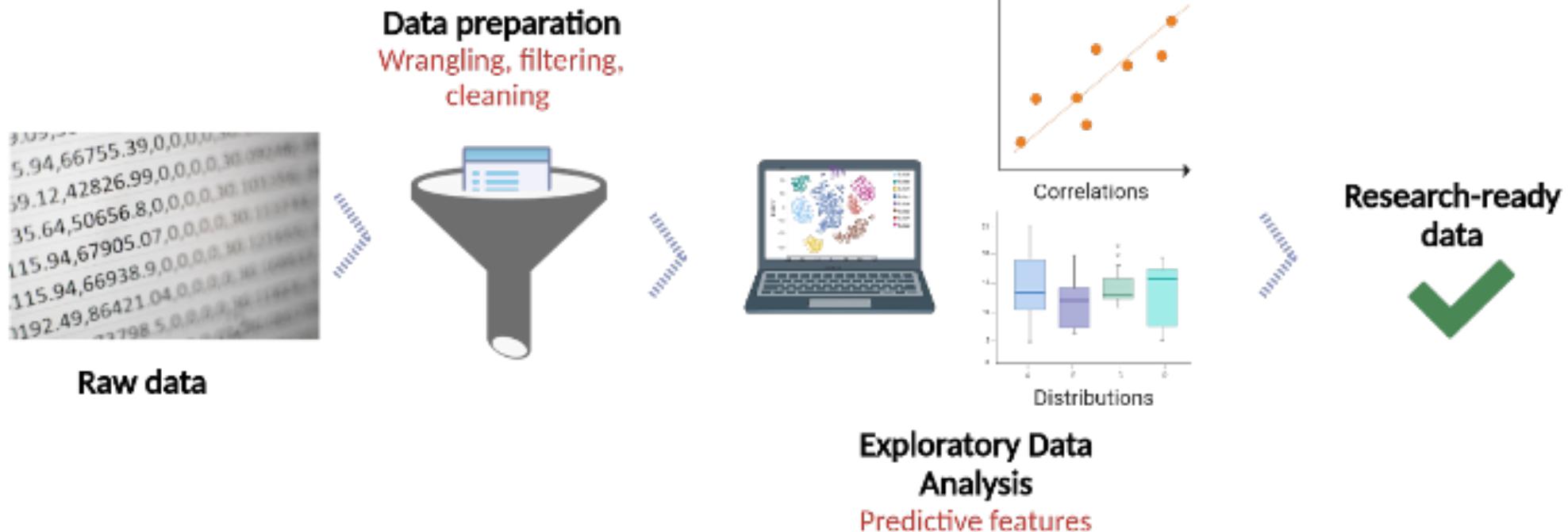


Figure 1: wrangling pipeline

Data Wrangling: a Simple Example

- What are all the “data” issues you might imagine with this tab delimited file?

Rank	State	Highest elevation	Lowest elevation	Average elevation
1	Colorado	"14,440 feet"	"3,315 feet"	"6,800 feet"
2	Wyoming	"13,804 feet"	"3,099 feet"	"6,700 feet"
3	Utah	"13,528 feet"	"2,000 feet"	"6,100 feet"
4	New Mexico	"13,161 feet"	"2,842 feet"	"5,700 feet"
5	Nevada	"13,140 feet"	"479 feet"	"5,500 feet"
6	Idaho	"12,662 feet"	"710 feet"	"5,000 feet"
7	Arizona	"12,633 feet"	"70 feet"	"4,100 feet"
8	Montana	"12,799 feet"	"1,800 feet"	"3,400 feet"
9	Oregon	"11,239 feet"	"Sea level"	"3,300 feet"
10	Hawaii	"13,796 feet"	"Sea level"	"3,030 feet"
11	California	"14,494 feet"	"282 feet"	"2,900 feet"
12	Nebraska	"5,424 feet"	"840 feet"	"2,600 feet"
13	South Dakota	"7,242 feet"	"966 feet"	"2,200 feet"

Figure 2: elevation

Data Wrangling: Another Example

- About how long is 3,175 mm?

“What is the point?”

According to the ISO, the decimal sign is written either as a comma or a point (a period), but in English the decimal sign is usually, although not always, written as a point. . . . The 22nd Conference Generale des Poids et Measures in 2003 repeated that “the symbol for the decimal marker shall be either the point on the line or the comma on the line.” Further, it reaffirmed that when numbers are divided in groups of three in order to facilitate reading, “neither dots nor commas are ever inserted in the spaces between groups.” (Grimvall 2011)

Data wrangling demo



https://github.com/chapmanbe/data_wrangling_demo

Generating research-ready clinical data

Preparing your data for analysis



Figure 3: data prep

Data quality: common issues

When poll is active respond at PollEv.com/brianchapman270

Send **brianchapman270** to **22333**



Figure 4: q1

Data quality: Lack of standardisation

- “In talking with our clinical experts, we learned that normal blood glucose levels are between 3.9 and 5.5 and a value above 7 would indicate a diagnosis of diabetes.”

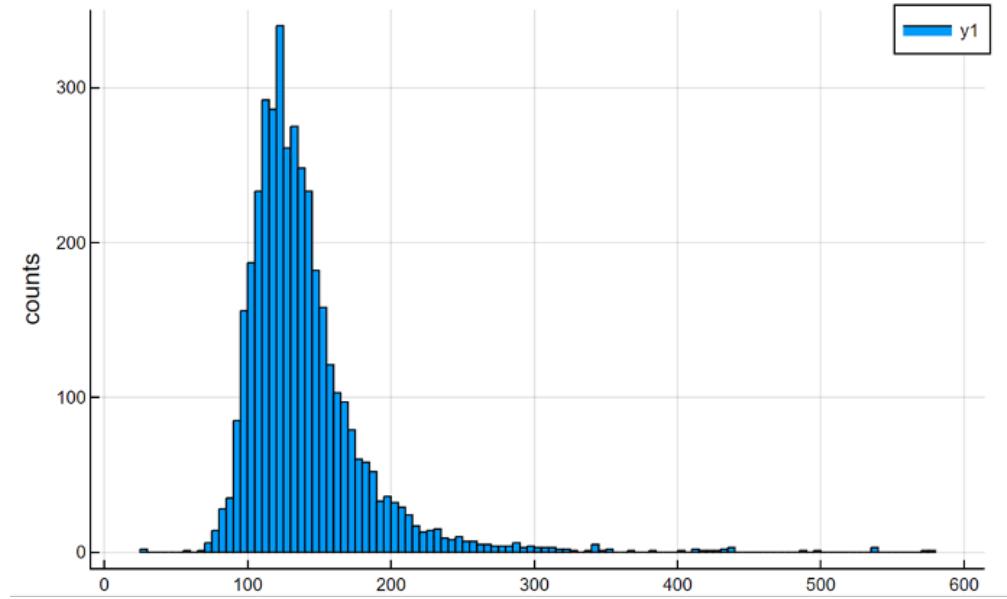


Figure 5: elevation

Data quality: Lack of standardisation

- Between 3.9 and 5.5 what?

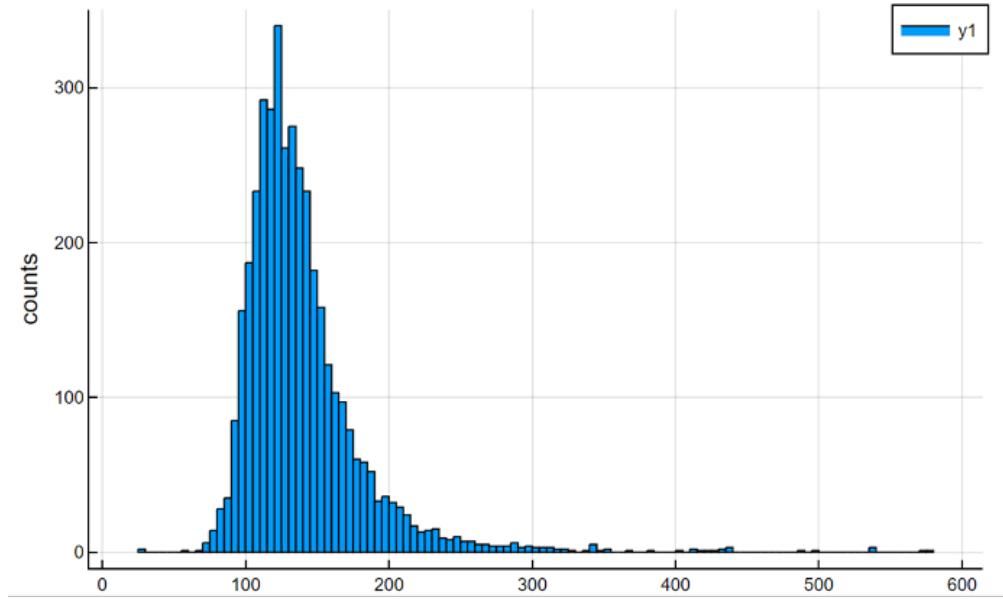


Figure 6: elevation

Data quality: Lack of standardisation

- Between 3.9 and 5.5 what?
- mmol/L

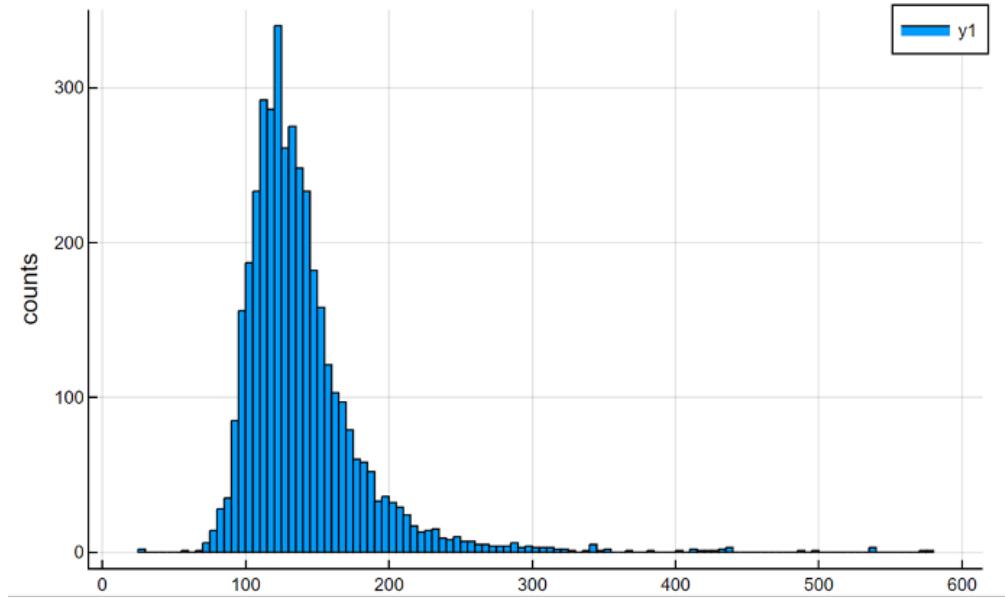


Figure 6: elevation

Data quality: Lack of standardisation

- Between 3.9 and 5.5 what?
- mmol/L
- MIMIC data units: mg/dL

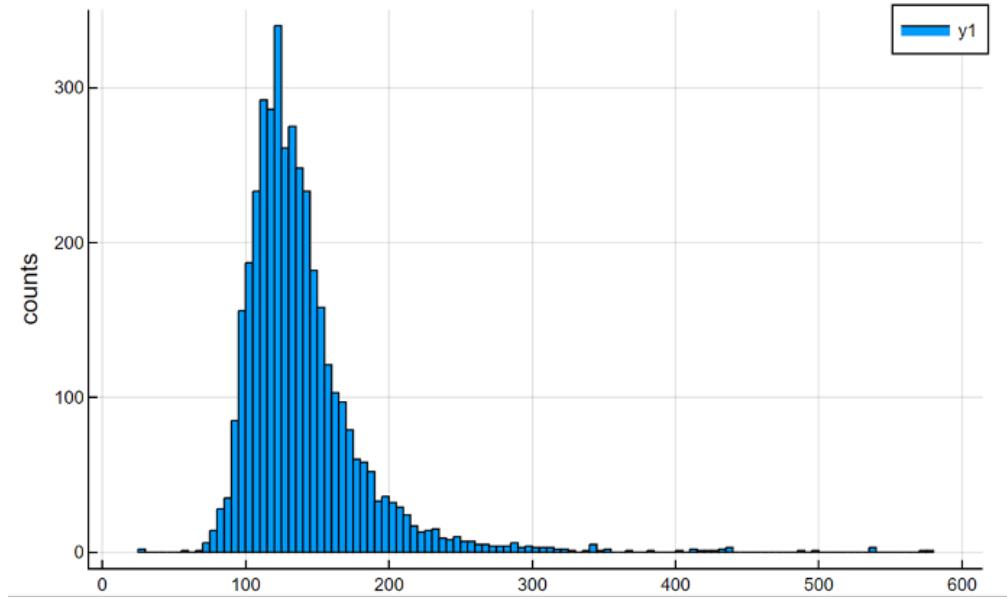


Figure 6: elevation

Data quality: Lack of standardisation

- Between 3.9 and 5.5 what?
- mmol/L
- MIMIC data units: mg/dL
- Convert!

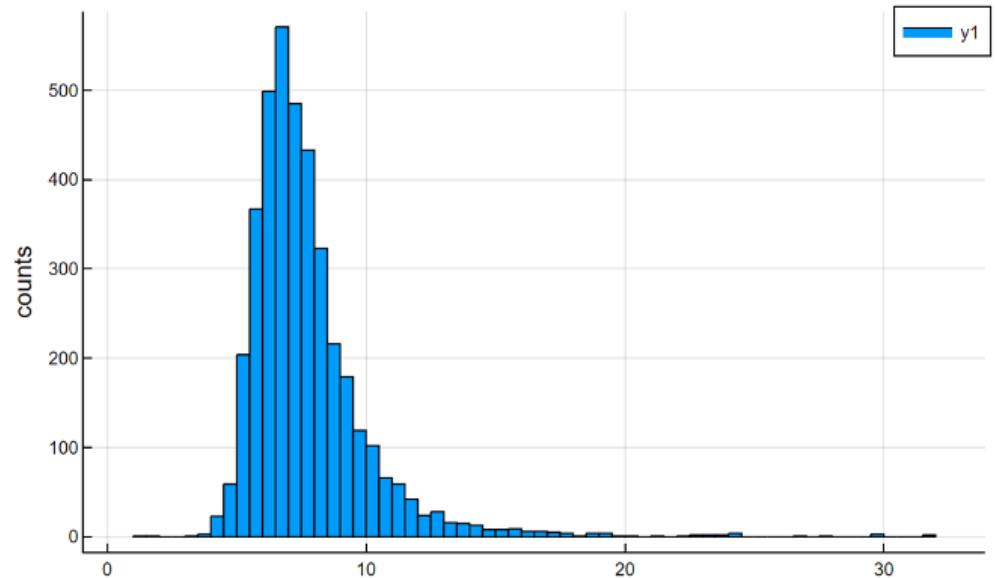


Figure 7: elevation

Does your machine learning algorithm care... .

. . . whether it is mmol/L or mg/dL?

Data Wrangling: an Example¹

2.11 Systolic Blood Pressure Histogram

```
select bucket, count(*) from (
    select width_bucket(value1num, 0, 300, 300) as bucket
        from mimic2v26.chartevents ce,
            mimic2v26.d_patients dp
    where itemid in (6, 51, 455, 6701)
        and ce.subject_id = dp.subject_id
        and months_between(ce.charttime, dp.dob)/12 > 15
) group by bucket order by bucket;
```

!

¹*MIMIC II SQL Cookbook*, Daniel J. Scott and Ikaro Silva

Data quality: outliers and errors

- What is an outlier?
- An observation that diverges or is distant from an overall pattern of samples/observations
- Data points lying outside the overall distribution of a data set

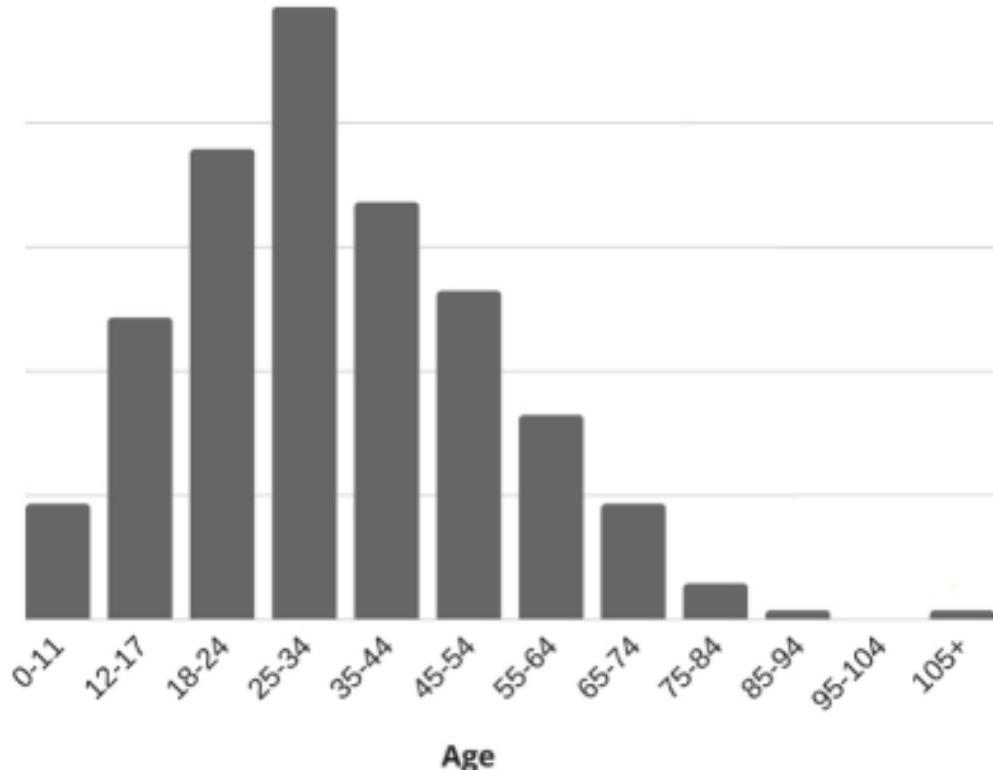
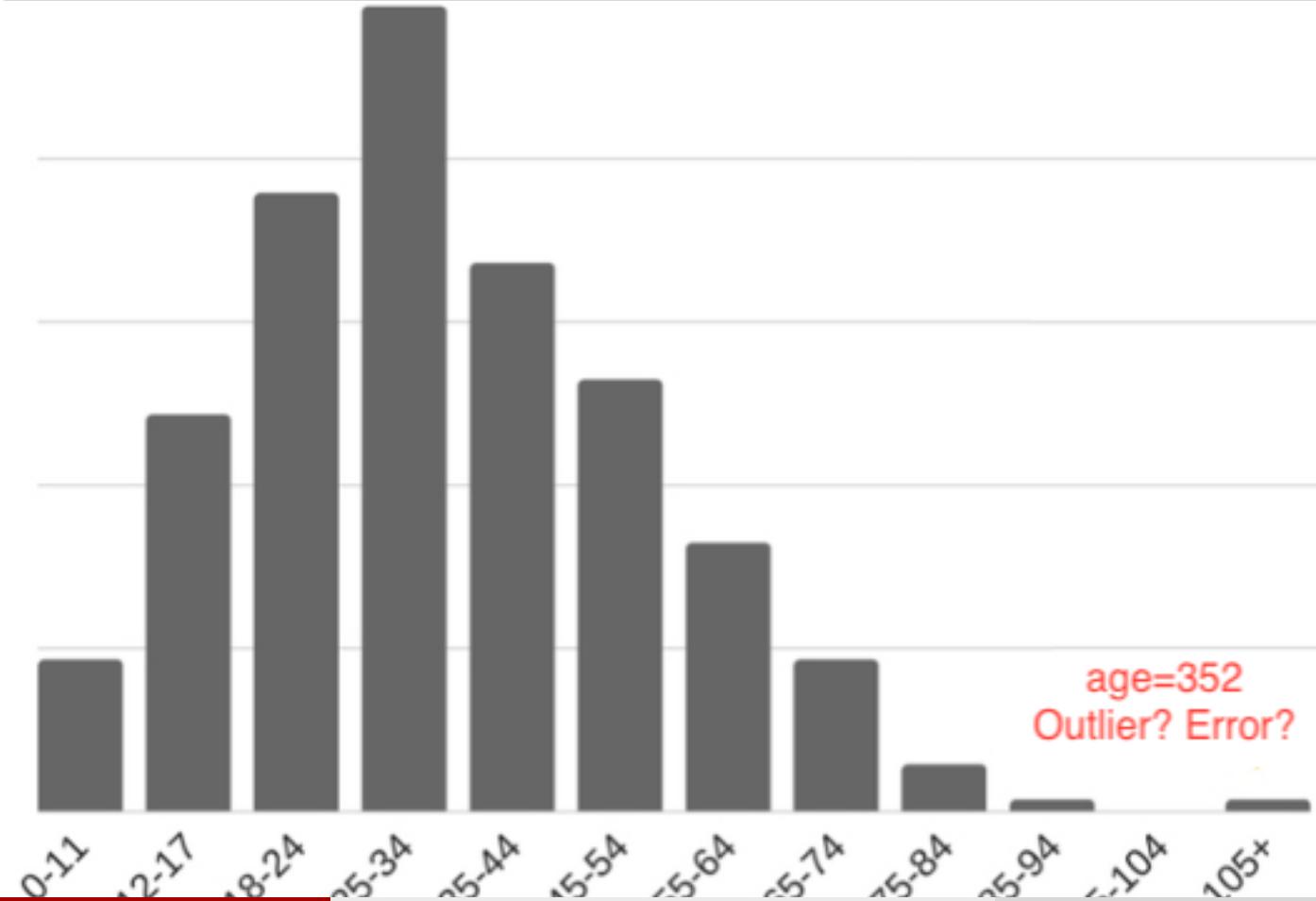
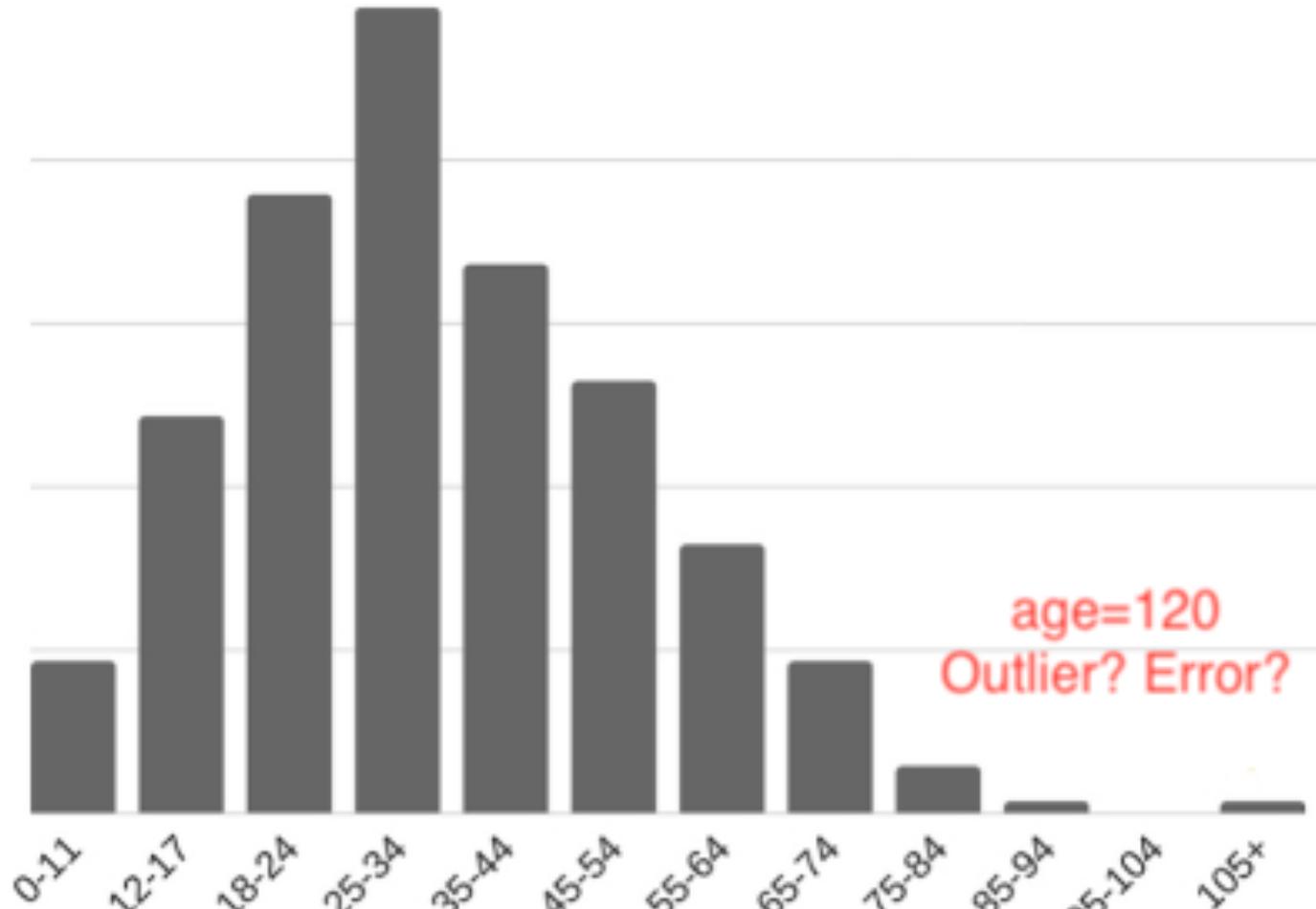


Figure 8: example distribution

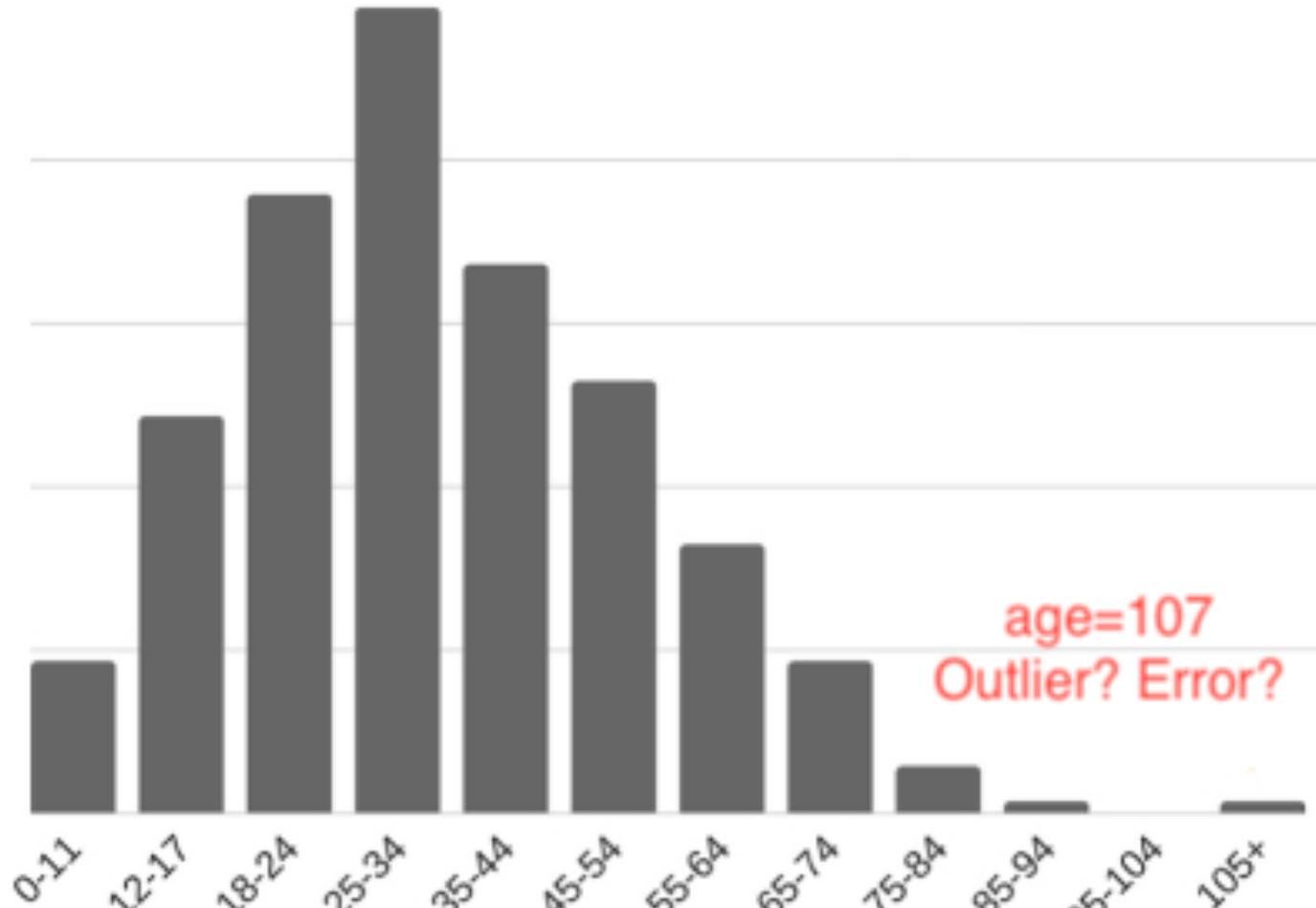
Data quality: outliers and errors



Data quality: outliers and errors



Data quality: outliers and errors



Outliers and implications for statistics & ML

- Girls under 16 basketball teams 2019



Outliers and implications for statistics & ML

- Girls under 16 basketball teams 2019
- Is the photo a mistake?



Outliers and implications for statistics & ML

- Girls under 16 basketball teams 2019
- Is the photo a mistake?
 - Red/Blue: USA



Outliers and implications for statistics & ML

- Girls under 16 basketball teams 2019
- Is the photo a mistake?
 - Red/Blue: USA
 - Blue/White: El Salvador



Outliers and implications for statistics & ML

- Girls under 16 basketball teams 2019
- Is the photo a mistake?
 - Red/Blue: USA
 - Blue/White: El Salvador
 - USA won 114-19



Outliers and implications for statistics & ML

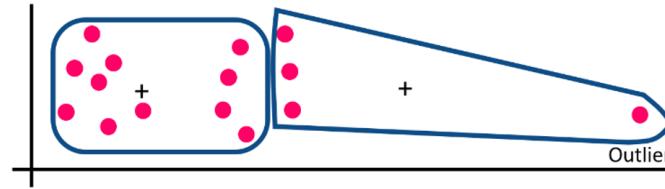


Figure 12: Outlier example 1a

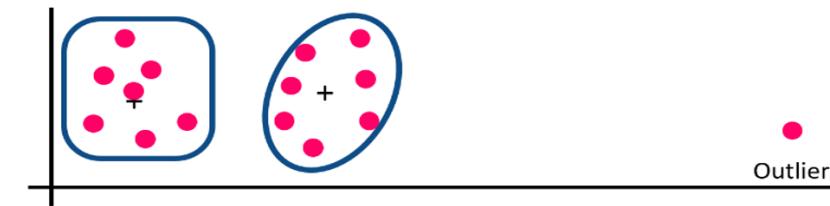


Figure 13: Outlier example 1b

Where do outliers come from?

- Data entry errors (human errors)
- Measurement errors (instrumental errors)
- Experimental errors (data extraction or experiment planning/executing errors)
- Intentional (dummy outliers made to test detection methods)
- Data processing errors (due to data manipulation or unintended mutations in data)
- Sampling errors (extracting or mixing data from wrong or various sources)
- Natural (not an error at all, but rather **novelties** in the data)

How to detect outliers? (non-exhaustive)

- Z-Score or Extreme Value Analysis (parametric)
- Probabilistic and Statistical Modeling (parametric)
- Linear Regression Models (prinicpal component analysis (PCA), least means squares (LMS))
- Proximity Based Models (non-parametric)
- Information Theory Models
- High Dimensional Outlier Detection Methods
- Visualization?
- **Domain knowledge?**

Why do I need to worry about outliers?

- They might reflect errors in our data that we need to correct.
- If they are novelty, they might influence downstream analysis that we carry out.
- Many different summary statistics (e.g., mean, standard deviation) are actually sensitive to outliers, meaning they will be heavily influenced by them.

What are missing values?

Missing data (or missing values) is defined as the data value that is not stored for a variable in the observation of interest. The problem of missing data is relatively common in almost all research and can have a significant effect on the conclusions that can be drawn from the data. Accordingly, some studies have focused on handling the missing data, problems caused by missing data, and the methods to avoid or minimize such in medical research. (Kang 2013)

Where do missing values come from?

Real-world, routinely collected clinical data often has lots of missing values and handling them adequately is a critical step for data cleaning. There can be multitude reasons why they occur:

- Human/data entry errors
- Optional fields/responses in surveys (e.g., patient weight might not be always recorded in all consultations)
- Incorrectly acquired data (e.g., from errors in sensor readings)
- Software bugs in data processing pipelines
- And many others!

Why do we need to handle them adequately?

- Absence of data reduces statistical power
 - Probability that the test will reject the null hypothesis when it is false
- Can cause bias in the estimation of parameters
- Can reduce the representativeness of the samples
- It may complicate the analysis of the study

Handling missing values

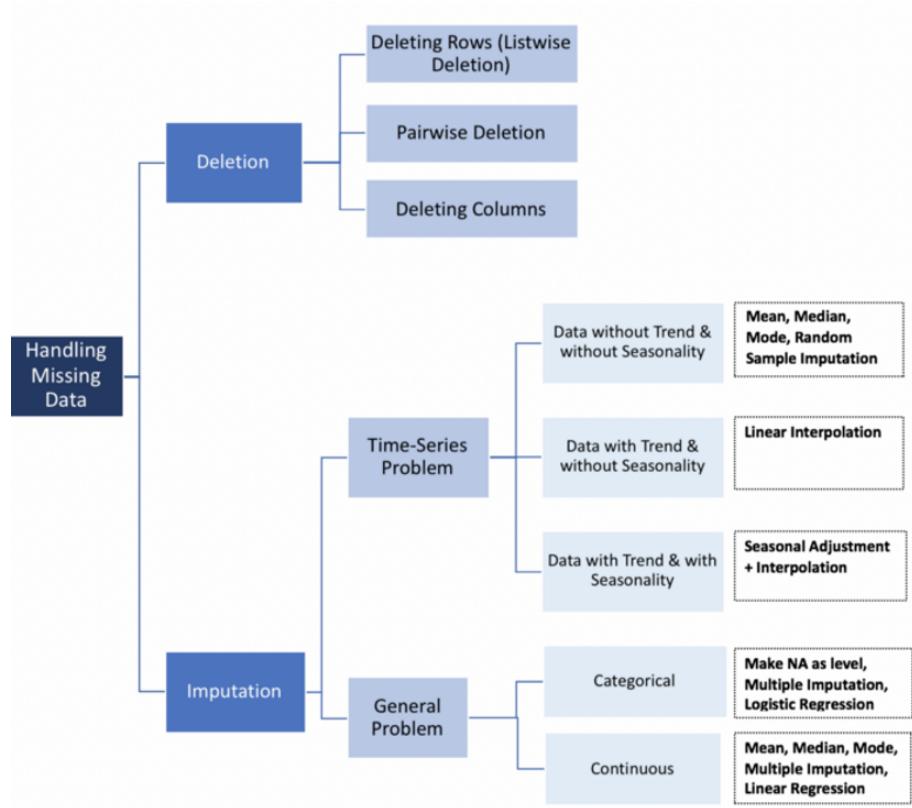


Figure 14: Missing values decision tree

Missing data vs. “Dark data”



Figure 15: David J. Hand

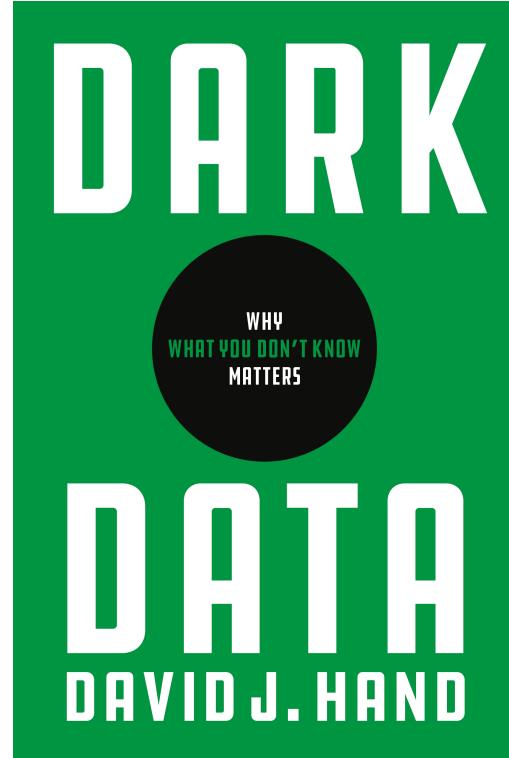


Figure 16: Data Data

Dark Data: A Taxonomy of Dark Data

- DD-Type 1: Data We Know are Missing
- DD-Type 2: Data We Don't Know are Missing
- DD-Type 3: Choosing Just Some Cases
- DD-Type 4: Self-Selection
- DD-Type 5: Missing What Matters
- DD-Type 6: Data Which Might Have Been
- DD-Type 7: Changes with Time
- DD-Type 8: Definitions of Data
- DD-Type 9: Summaries of Data
- DD-Type 10: Measurement Error and Uncertainty
- DD-Type 11: Feedback and Gaming
- DD-Type 12: Information Asymmetry
- DD-Type 13: Intentionally Darkened Data
- DD-Type 14: Fabricated and Synthetic Data
- DD-Type 15: Extrapolating beyond Your Data

(Hand 2020)

The map is not the territory!

Correlation vs. causation

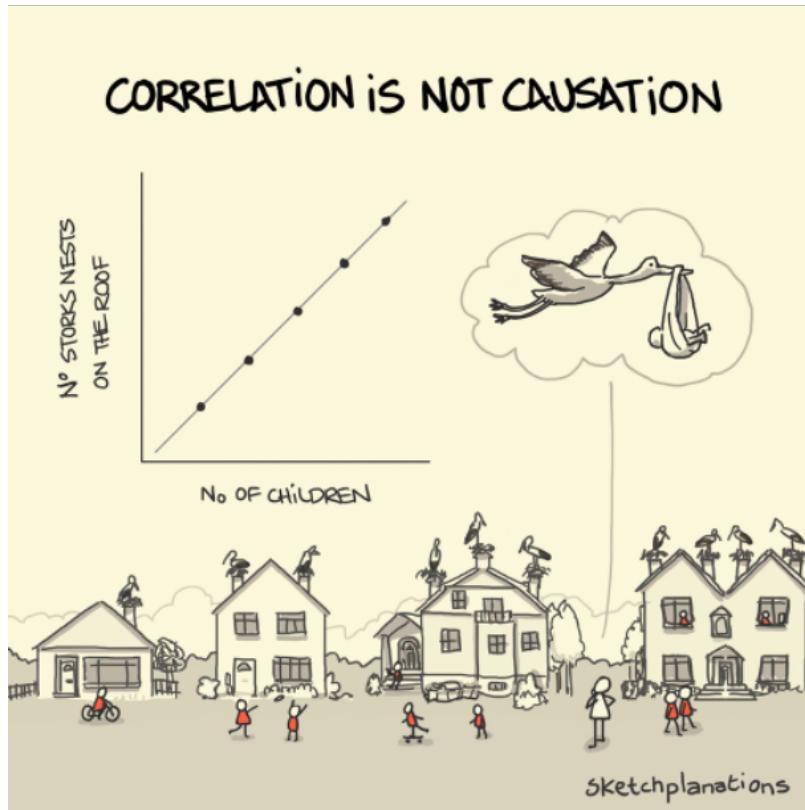


Figure 17: correlation vs. causation

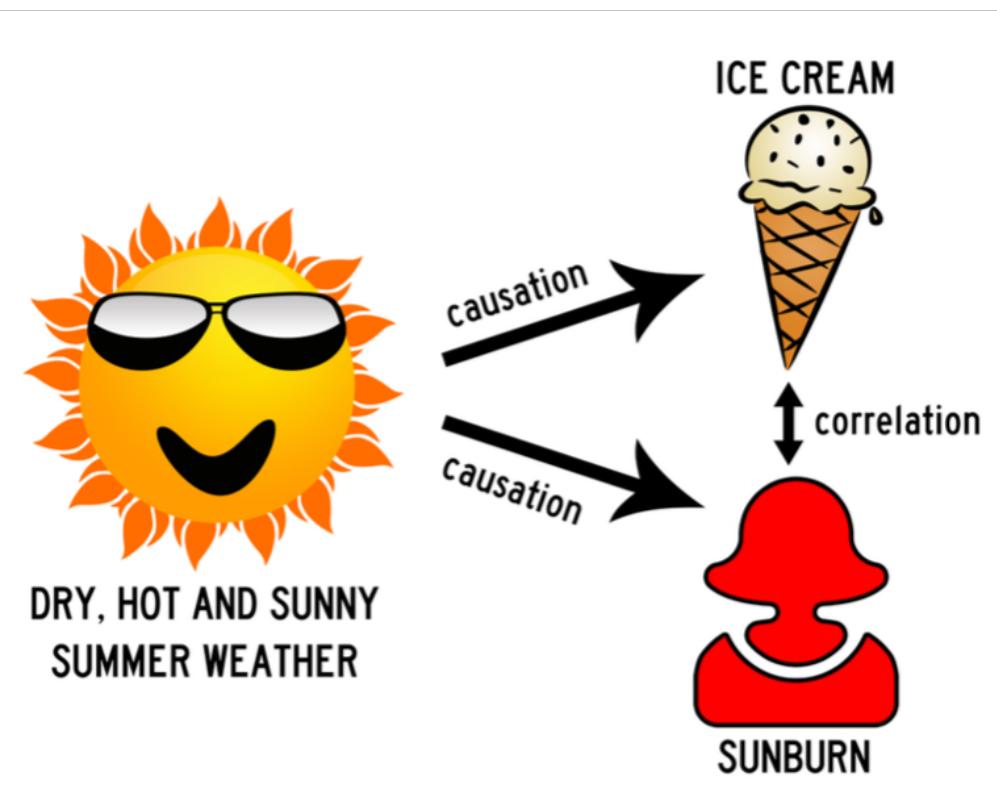


Figure 18: correlation vs. causation

exploratory Data Analysis (EDA)

- Task of analysing data using simple tools from statistics to simple plotting tools.
 - Discover patterns
 - Identify outliers and errors
 - Check assumptions
 - Identify promising variables for predictive modelling

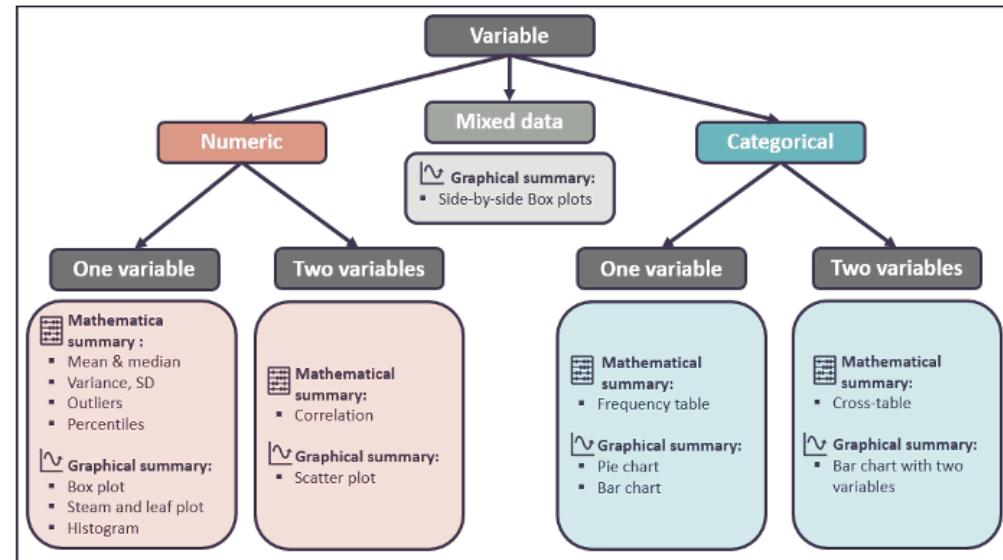


Figure 19: EDA

Exploratory Data Analysis (EDA)

- Before going from analytic to predictive
 - Treating ML as a black box can be dangerous
- Rational selection and investigation of potential features
 - Leads to more interpretable models
 - Less complex (minimal set of features) and more generalisable predictive models (Occam's Razor principle)

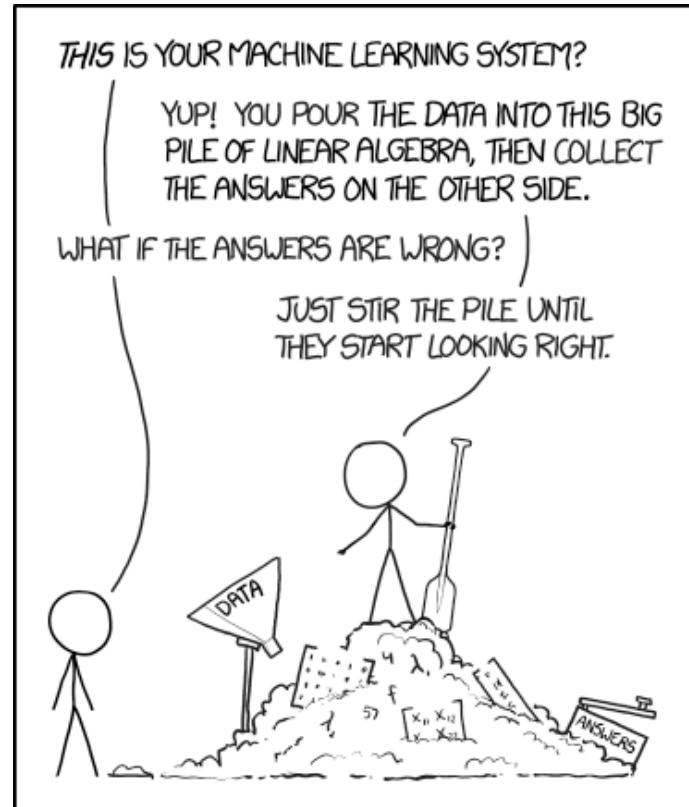


Figure 20: <https://xkcd.com/1838/>

Exploratory Data Analysis via Visualization

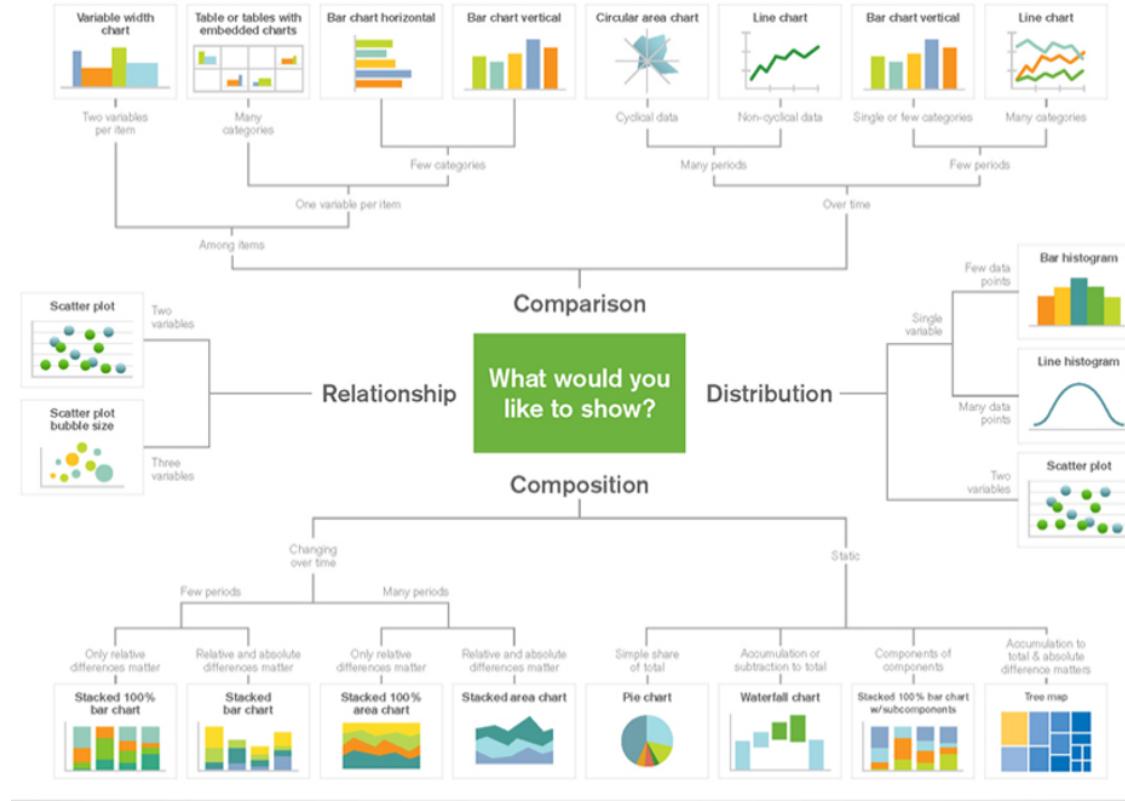


Figure 21: EDA visualization

Reproducibility

- Data wrangling has lots of steps, lots of humanness. How do we make it reproducible?
- This is a hard problem!

Reproducibility: Specifying Anaconda environments

environment.yml files

```
name: dwd
```

```
channels:
```

- pyviz
- conda-forge

```
dependencies:
```

- python=3.11
- pandas=2.0.3
- hvplot=0.8.4
- jupyterlab=4.0.3

Reproducibility: Docker (Podman, Singularity)

```
FROM jupyter/base-notebook:5cb1a915c4bf
MAINTAINER chapmanbe <brian.chapman@utah.edu>
USER root

RUN apt-get update && apt-get upgrade -y && apt-get install -y \
    locales-all \
&& rm -rf /var/lib/apt/lists/*

RUN conda update conda -y && conda install -c conda-forge -c pyviz -y \
    pandas=2.0.3 \
    hvplot=0.8.4

WORKDIR /home/jovyan

COPY Resources/ Resources/
ADD Pandas-DataWrangling.ipynb .

# RUN nbstripout --install
CMD ["start-notebook.sh"]
```

Reproducibility

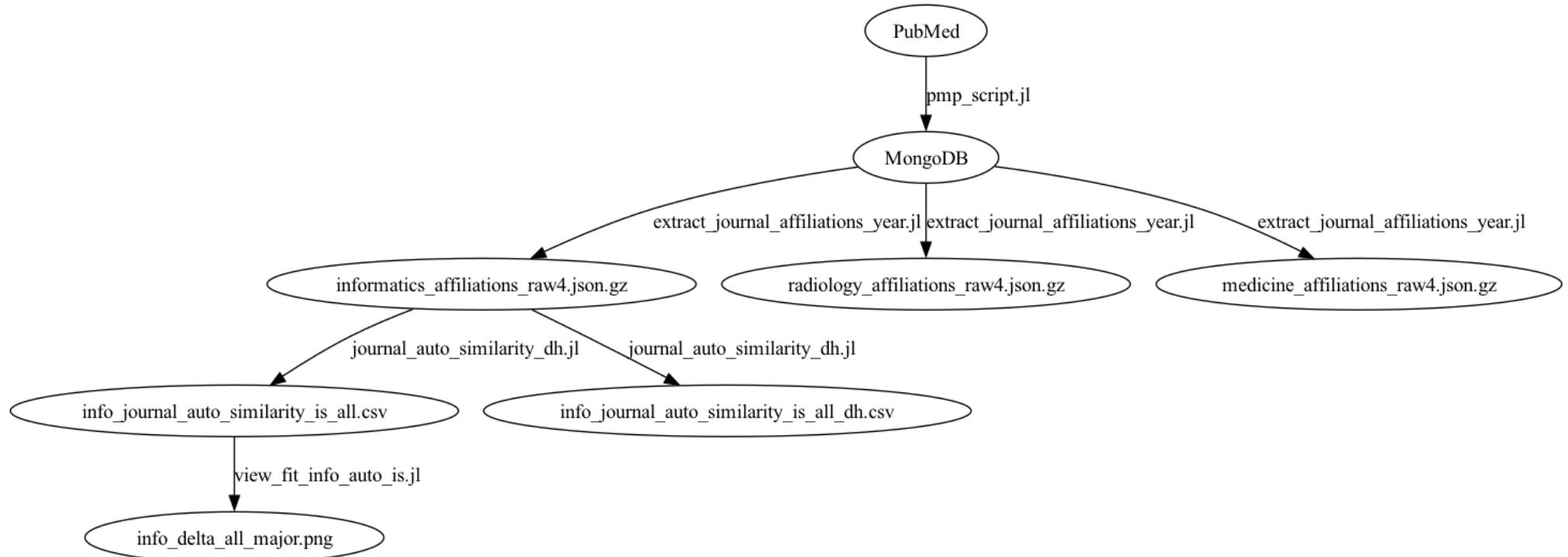


Figure 22: my data flow

Reproducibility

- Reproducible Data Science with Python
- Reproducible Data Analysis in Jupyter

References I

- Grimvall, G. 2011. *Quantify!: A Crash Course in Smart Thinking*. Johns Hopkins University Press. <https://books.google.com.au/books?id=5PqTZbgl-7QC>.
- Hand, D. J. 2020. *Dark Data: Why What You Don't Know Matters*. Princeton University Press. <https://books.google.com.au/books?id=FL2mDwAAQBAJ>.
- Kang, H. 2013. “The prevention and handling of the missing data.” *Korean J Anesthesiol* 64 (5): 402–6.