

# Lecture 1: Medical Databases

CSCI6410/EPAH6410/CSCI4148

Finlay Maguire (finlay.maguire@dal.ca)

# Learning Objectives

- Overview of the types of medical database
- Ways of maintaining data privacy with medical databases and some of their trade-offs
- How and why ontologies and survey weights are used in medical databases
  
- Key strategies/approaches for exploratory data analysis
- Different types of dimensionality reduction
- Basics of supervised learning
- Accessing feature importances
- Aggregating simple/weak models to improve performance: boosting and bagging

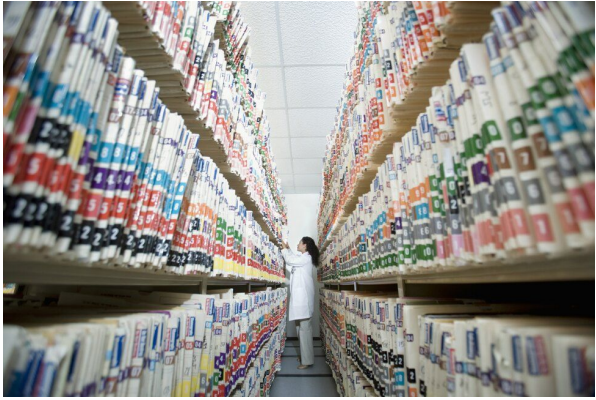
What is a database?

# Databases (broadly) are ordered collections of data

Examples include:

- Medical Charts

The image displays several overlapping medical history forms. The top-left form is titled 'PART A - PRESENT HEALTH HISTORY (continued)' and includes sections for 'IV. GENERAL HEALTH, ATTITUDE AND HABITS', 'I. FAMILY HEALTH', and 'II. HOSPITALIZATIONS, SURGERIES'. The top-right form is 'PART C - BODY SYSTEMS REVIEW'. The middle form is the 'ANDRUS/CLINI-REC HEALTH HISTORY QUESTIONNAIRE'. The bottom form is another 'PART A - PRESENT HEALTH HISTORY' with sections for 'I. CURRENT MEDICAL PROBLEMS', 'II. MEDICATIONS', 'III. ALLERGIES AND SENSITIVITIES', and 'IV. GENERAL HEALTH, ATTITUDE AND HABITS'. The word 'CONFIDENTIAL' is printed vertically on the left and right sides of the forms.

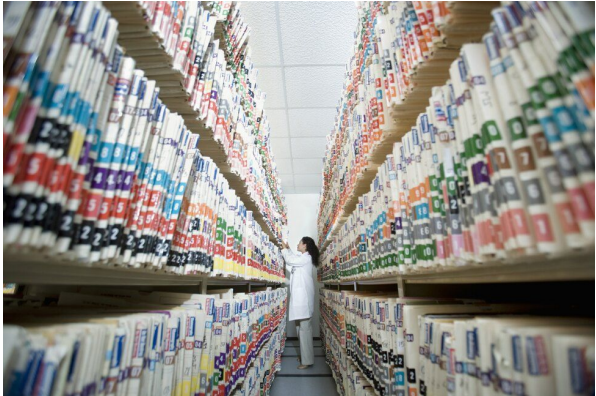


# Databases (broadly) are ordered collections of data

Examples include:

- Medical Charts
- Phone Book
- Dictionaries
- Spreadsheet

The image displays several overlapping medical forms. The top-left form is 'PART A - PRESENT HEALTH HISTORY (continued)' with sections for general health, family health, hospitalizations, and current medical problems. The top-right form is 'PART C - BODY SYSTEMS REVIEW' with a grid for various body systems. The middle form is the 'ANDRUS/CLINI-REC HEALTH HISTORY QUESTIONNAIRE' for a patient named ANDRUS/CLINI-REC. The bottom form is another 'PART A - PRESENT HEALTH HISTORY' with sections for current medical problems, medications, allergies, and general health. The forms are marked with 'CONFIDENTIAL' on the sides.



# Databases (broadly) are ordered collections of data

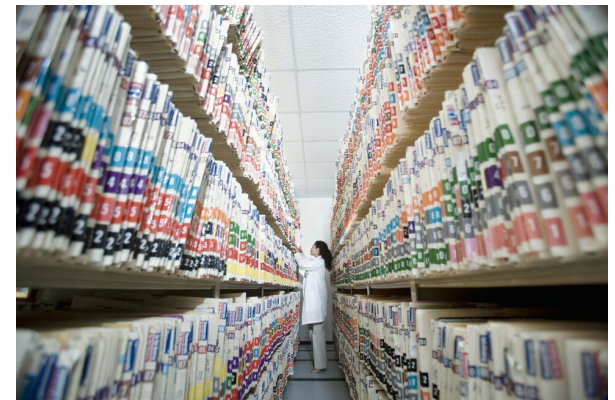
Examples include:

- Medical Charts
- Phone Book
- Dictionaries
- Spreadsheet

Ordering:

- Index
- Defined fields
- Standardisation

The image displays several overlapping medical forms. The top-left form is 'PART A - PRESENT HEALTH HISTORY (continued)' with sections for general health, family health, hospitalizations, and current medical problems. The top-right form is 'PART C - BODY SYSTEMS REVIEW' with a grid for various body systems. The middle form is the 'ANDRUS/CLINI-REC HEALTH HISTORY QUESTIONNAIRE' for a patient named ANDRUS/CLINI-REC. The bottom form is another 'PART A - PRESENT HEALTH HISTORY' with sections for current medical problems, medications, allergies, and general health. The forms are marked with 'CONFIDENTIAL' on the sides.



# Databases (broadly) are ordered collections of data

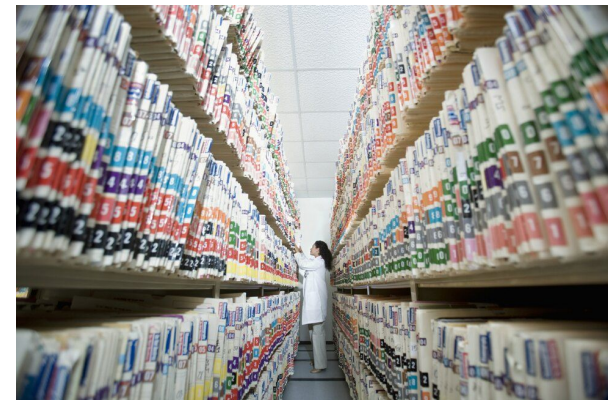
Examples include:

- Medical Charts
- Phone Book
- Dictionaries
- Spreadsheet

Ordering:

- Index
- Defined fields
- Standardisation

The image shows several overlapping medical forms. The top form is 'PART A - PRESENT HEALTH HISTORY (continued)' with sections for general health, family health, and hospitalizations. Below it is 'PART B - PAST HISTORY' and 'PART C - BODY SYSTEMS REVIEW'. In the center is the 'ANDRUS/CLINI-REC HEALTH HISTORY QUESTIONNAIRE' for a patient named ANDRUS/CLINI-REC. The forms are filled with various questions and checkboxes, and the word 'CONFIDENTIAL' is printed vertically on the left and right sides of the forms.

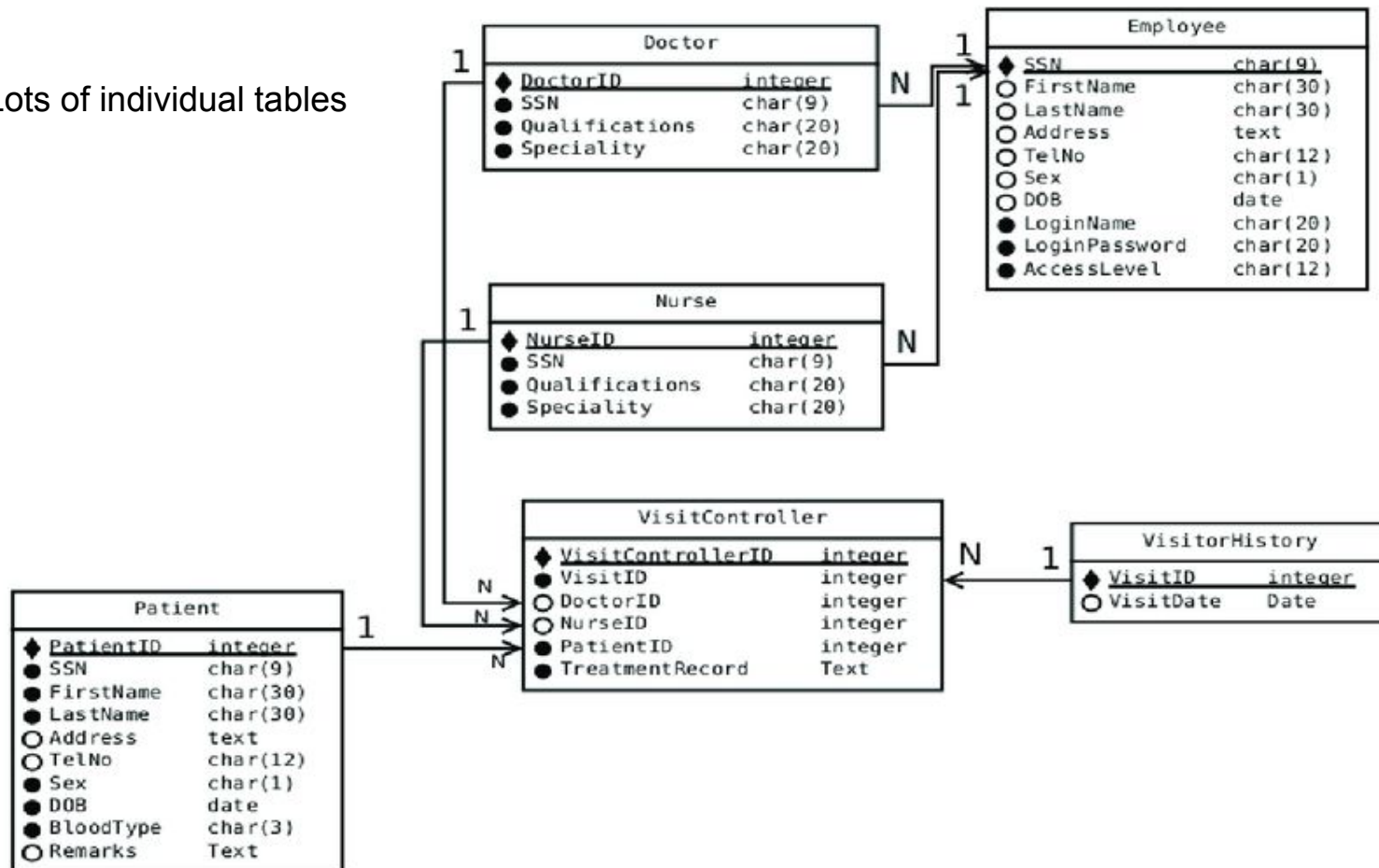


Organisation make some tasks easier/harder:

- Find all patients with the same condition
- Find the longest word in a dictionary
- Find an a number from an address in a phonebook

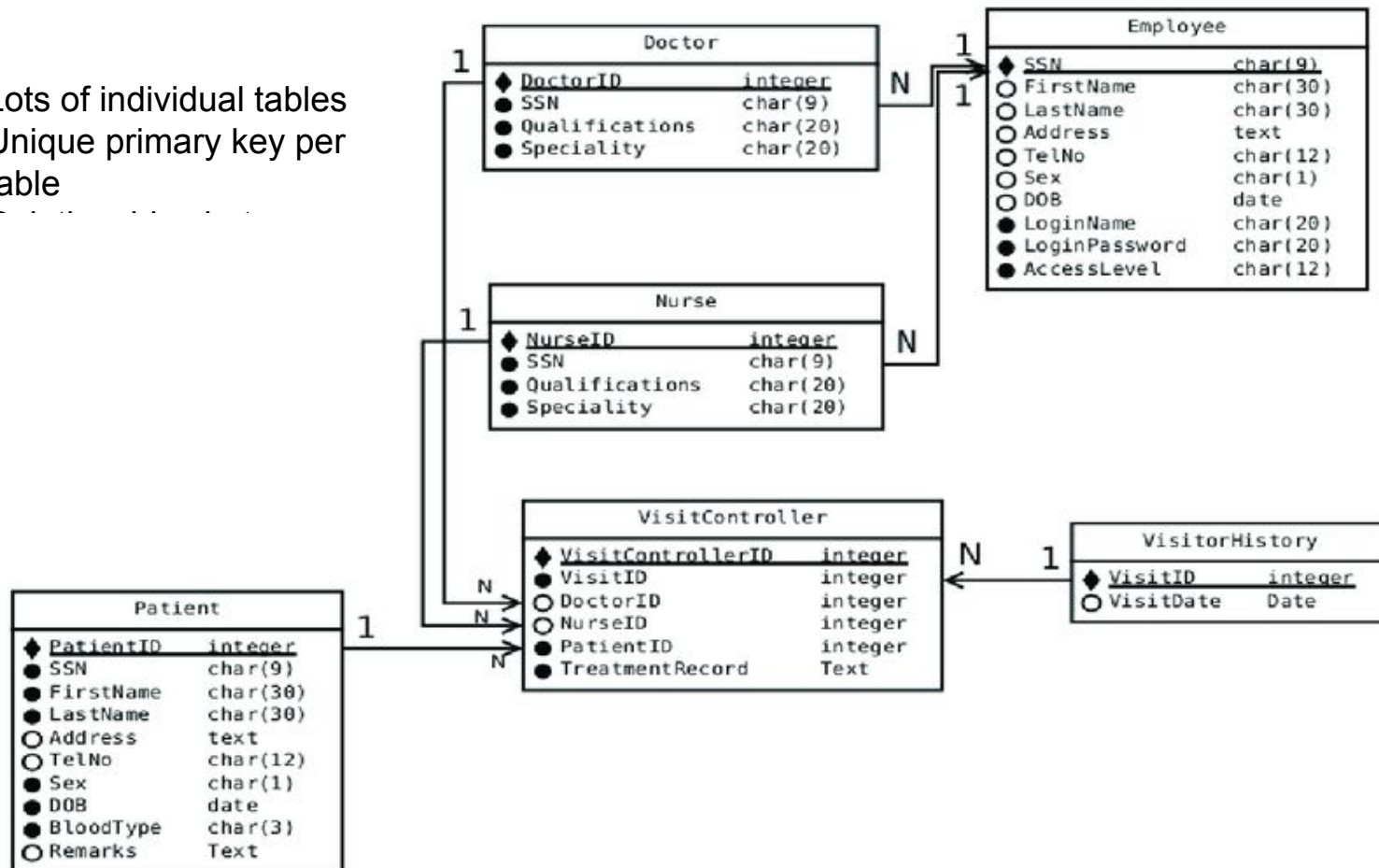
# Most Common Type: Relational Databases

- Lots of individual tables



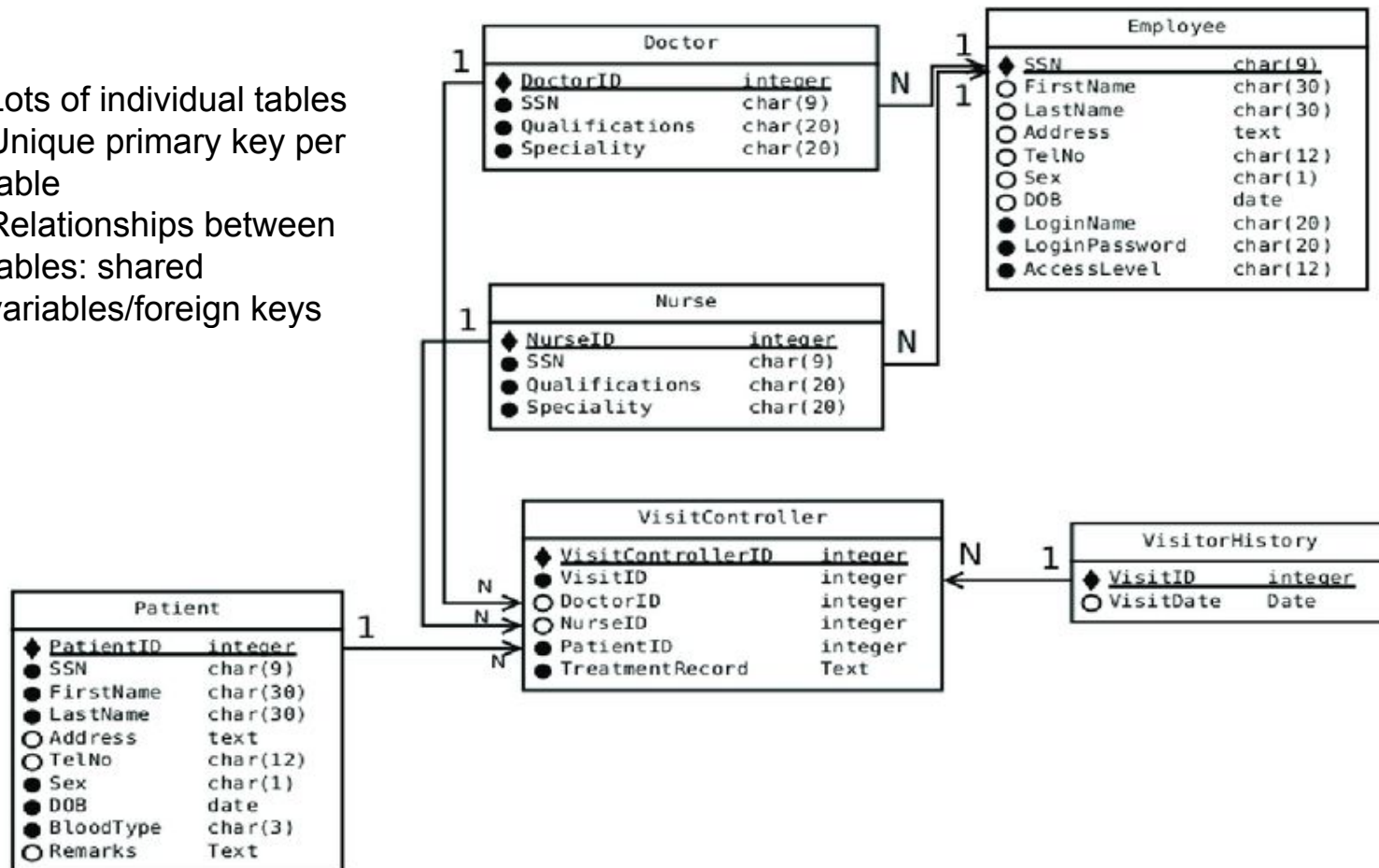
# Most Common Type: Relational Databases

- Lots of individual tables
- Unique primary key per table



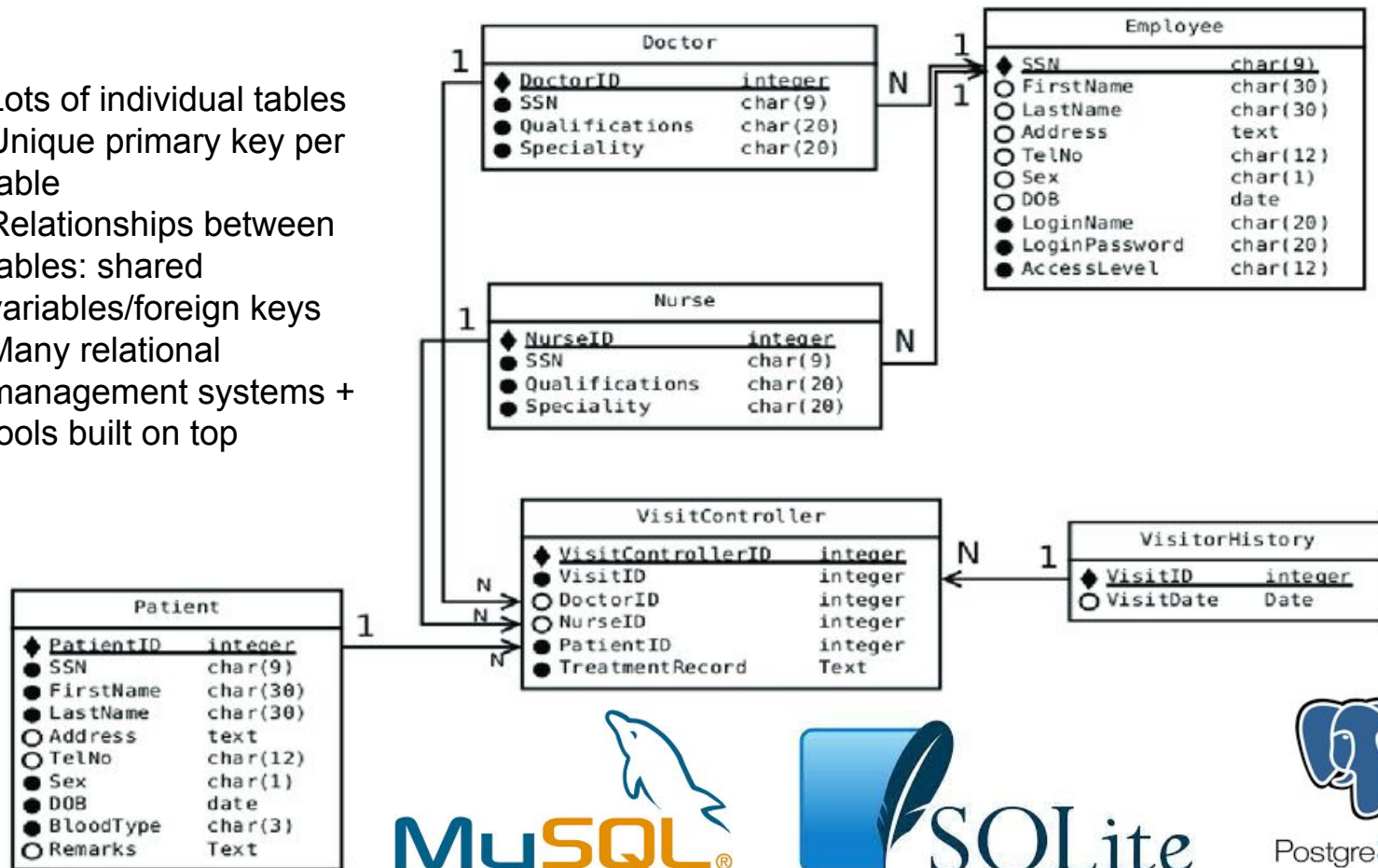
# Most Common Type: Relational Databases

- Lots of individual tables
- Unique primary key per table
- Relationships between tables: shared variables/foreign keys

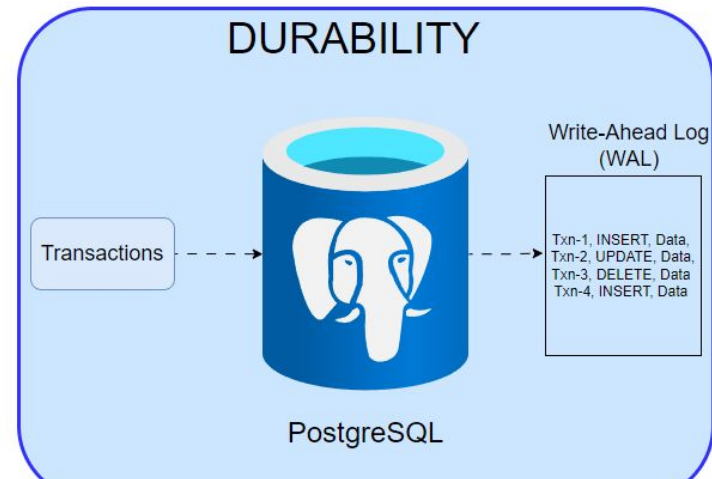
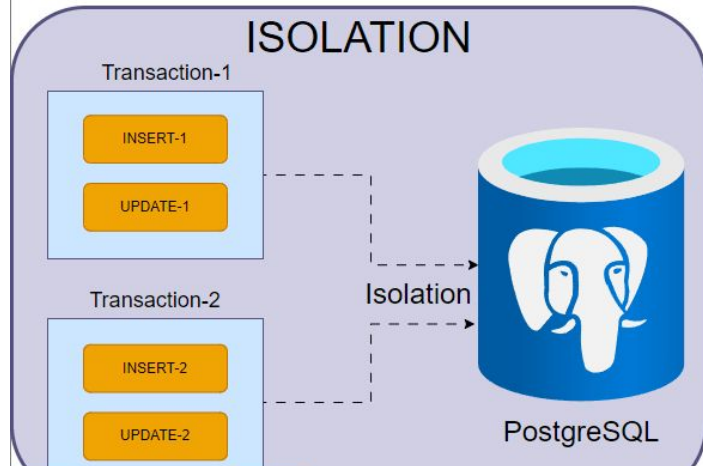
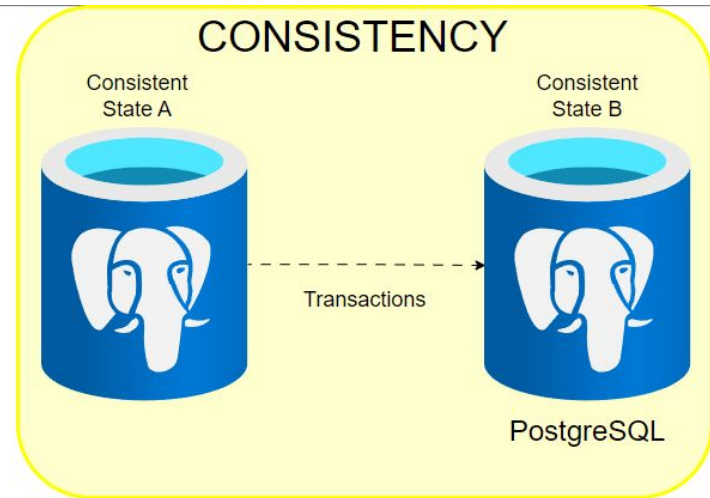
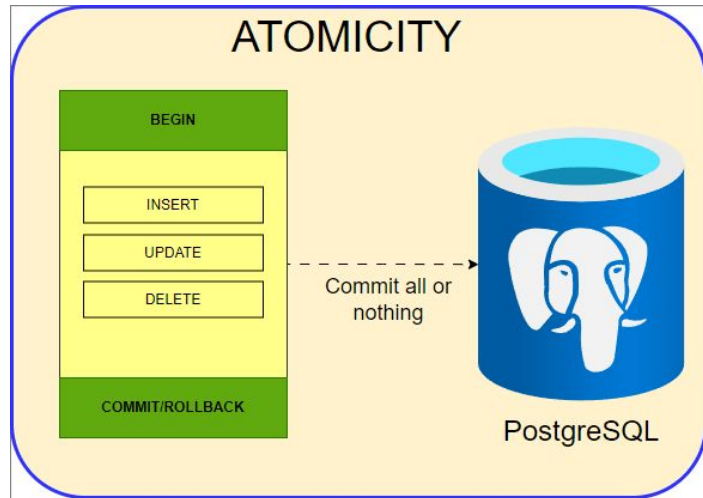


# Most Common Type: Relational Databases

- Lots of individual tables
- Unique primary key per table
- Relationships between tables: shared variables/foreign keys
- Many relational management systems + tools built on top



# Why bother with databases? ACID properties



# Queried using Structured Query Language (SQL)

- Non-procedural Language
- Standardised/powerful/flexible

# Queried using Structured Query Language (SQL)

- Non-procedural Language
- Standardised/powerful/flexible
- Basis of many data tools
- Well-supported by dbplyr

# Queried using Structured Query Language (SQL)

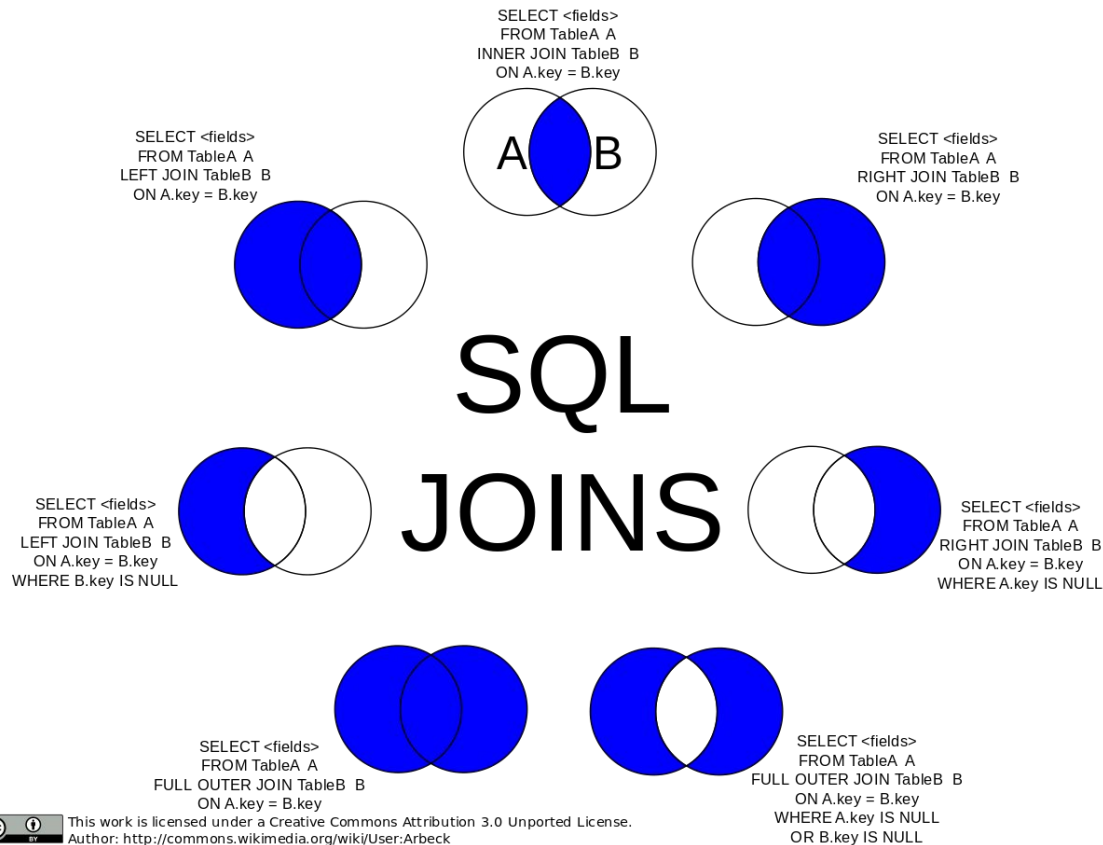
- Non-procedural Language
- Standardised/powerful/flexible
- Basis of many data tools
- Well-supported by dbplyr

```
flights %>%
  select(contains("delay")) %>%
  show_query()
#> <SQL>
#> SELECT `dep_delay`, `arr_delay`
#> FROM `nycflights13::flights`
```

```
flights %>%
  select(distance, air_time) %>%
  mutate(speed = distance / (air_time / 60)) %>%
  show_query()
#> <SQL>
#> SELECT `distance`, `air_time`, `distance` / (`air_time` / 60.0) AS `speed`
#> FROM (SELECT `distance`, `air_time`
#> FROM `nycflights13::flights`)
```

```
flights %>%
  group_by(month, day) %>%
  summarise(delay = mean(dep_delay)) %>%
  show_query()
#> Warning: Missing values are always removed in SQL.
#> Use `AVG(x, na.rm = TRUE)` to silence this warning
#> <SQL>
#> SELECT `month`, `day`, AVG(`dep_delay`) AS `delay`
#> FROM `nycflights13::flights`
#> GROUP BY `month`, `day`
```

# SQL enables complex joins/queries



**INNER JOIN** - only assigned patients and their doctors

**LEFT JOIN** - all doctors including those without patients

**RIGHT JOIN** - all patients including those without doctors

**OUTER JOIN** - all patients and doctors

This work is licensed under a Creative Commons Attribution 3.0 Unported License. Author: <http://commons.wikimedia.org/wiki/User:Arbeck>

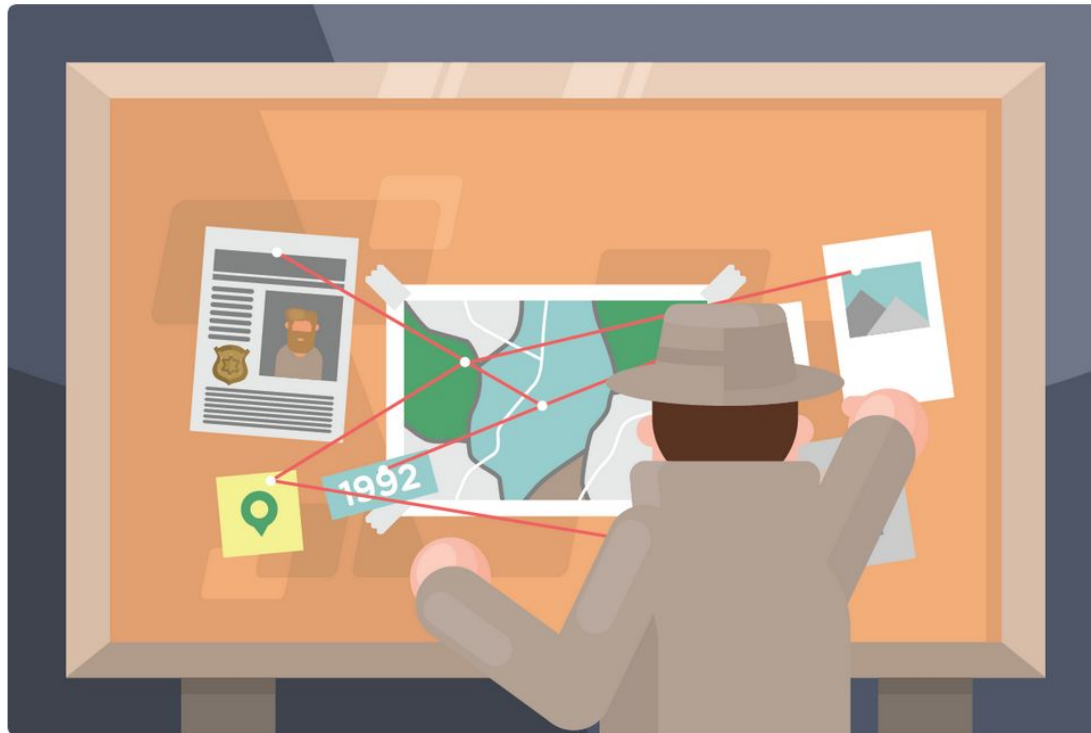
| doctors          |             |                  | patients          |             |                          |
|------------------|-------------|------------------|-------------------|-------------|--------------------------|
| <u>doctor_id</u> | <u>name</u> | <u>specialty</u> | <u>patient_id</u> | <u>name</u> | <u>primary_doctor_id</u> |
| 10               | Dr. Smith   | Cardiology       | 1                 | Alice       | 10                       |
| 20               | Dr. Jones   | Neurology        | 2                 | Bob         | 20                       |
| 30               | Dr. Lee     | Pediatrics       | 3                 | Carol       | NULL                     |
| 40               | Dr. Patel   | Oncology         | 4                 | Dave        | 10                       |

# Fun way to learn basic SQL

<https://mystery.knightlab.com/>

## *SQL Murder Mystery*

Can you find out whodunnit?



**Are all databases relational?**

# Non-Relational Databases AKA NoSQL

- Less common than relational in medicine
- General focus on flexibility & performance
- Mostly for very large/unusual datasets or high demand:
  - User data / security audit data
  - Medical image data
- Or unusual data structures:
  - Contact tracing
  - Ontologies
- Or both:
  - Social media data

<https://phoenixnap.com/kb/database-types>



**Column based**



Google  
Big Query

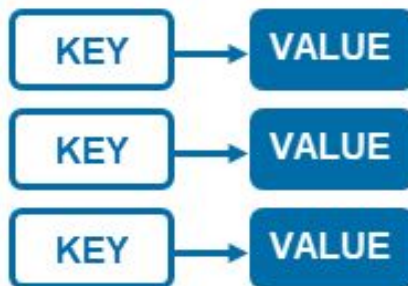
# Non-Relational Databases AKA NoSQL

- Less common than relational in medicine
- General focus on flexibility & performance
- Mostly for very large/unusual datasets or high demand:
  - User data / security audit data
  - Medical image data
- Or unusual data structures:
  - Contact tracing
  - Ontologies
- Or both:
  - Social media data

<https://phoenixnap.com/kb/database-types>



**Column based**



**Key-value**



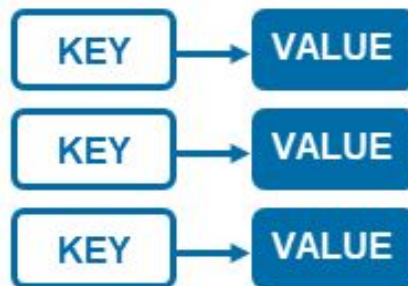
# Non-Relational Databases AKA NoSQL

- Less common than relational in medicine
- General focus on flexibility & performance
- Mostly for very large/unusual datasets or high demand:
  - User data / security audit data
  - Medical image data
- Or unusual data structures:
  - Contact tracing
  - Ontologies
- Or both:
  - Social media data

<https://phoenixnap.com/kb/database-types>



Column based



Key-value



Graph



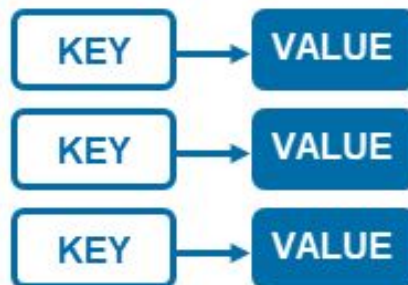
# Non-Relational Databases AKA NoSQL

- Less common than relational in medicine
- General focus on flexibility & performance
- Mostly for very large/unusual datasets or high demand:
  - User data / security audit data
  - Medical image data
- Or unusual data structures:
  - Contact tracing
  - Ontologies
- Or both:
  - Social media data

<https://phoenixnap.com/kb/database-types>



Column based



Key-value



Graph

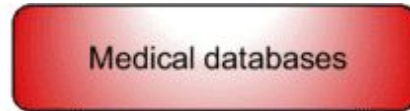


Document



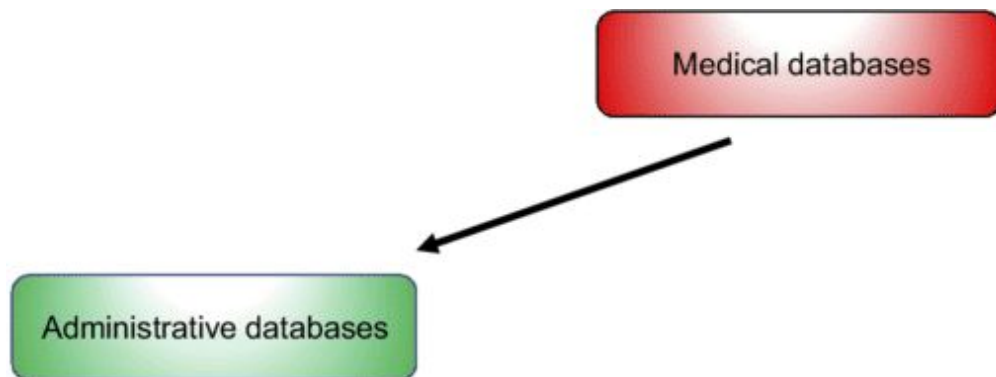
What are medical databases?

# Many types of database



■ All types of registries and databases that contain health-related data

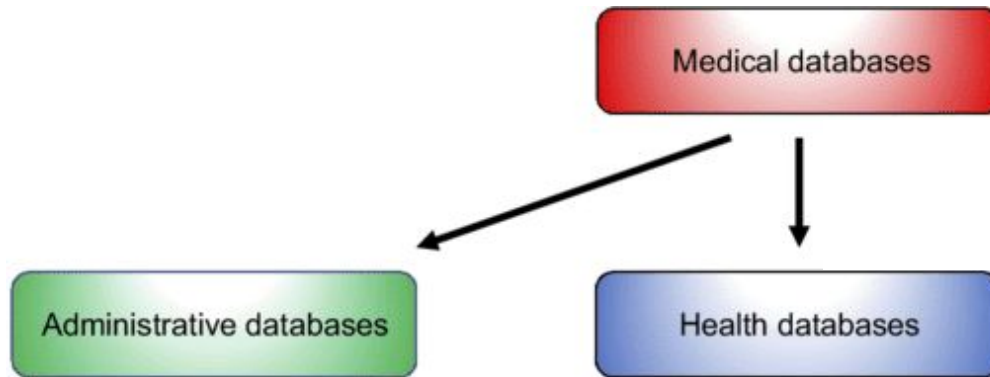
# Many types of database



■ All types of registries and databases that contain health-related data

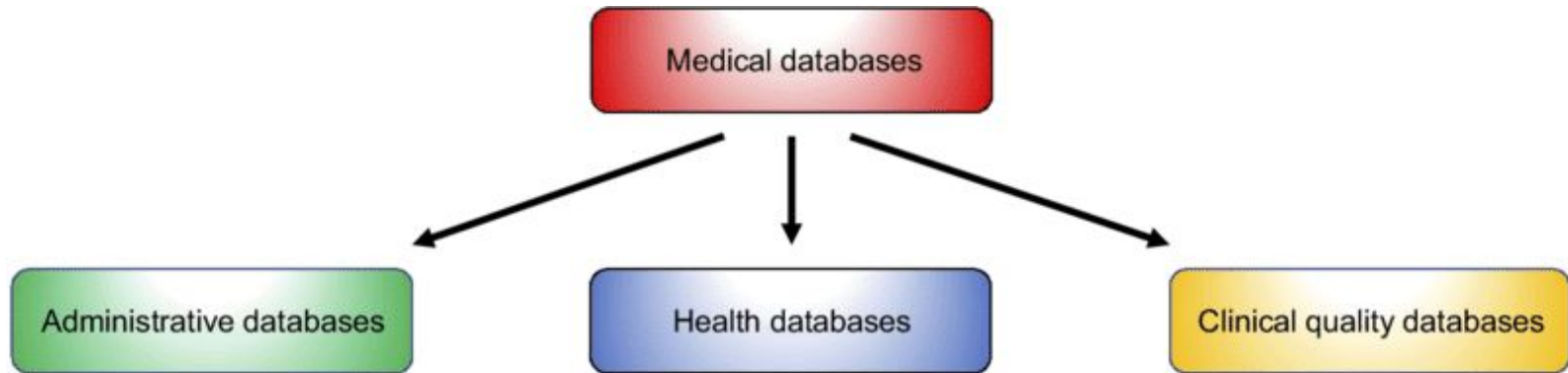
■ Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic

# Many types of database



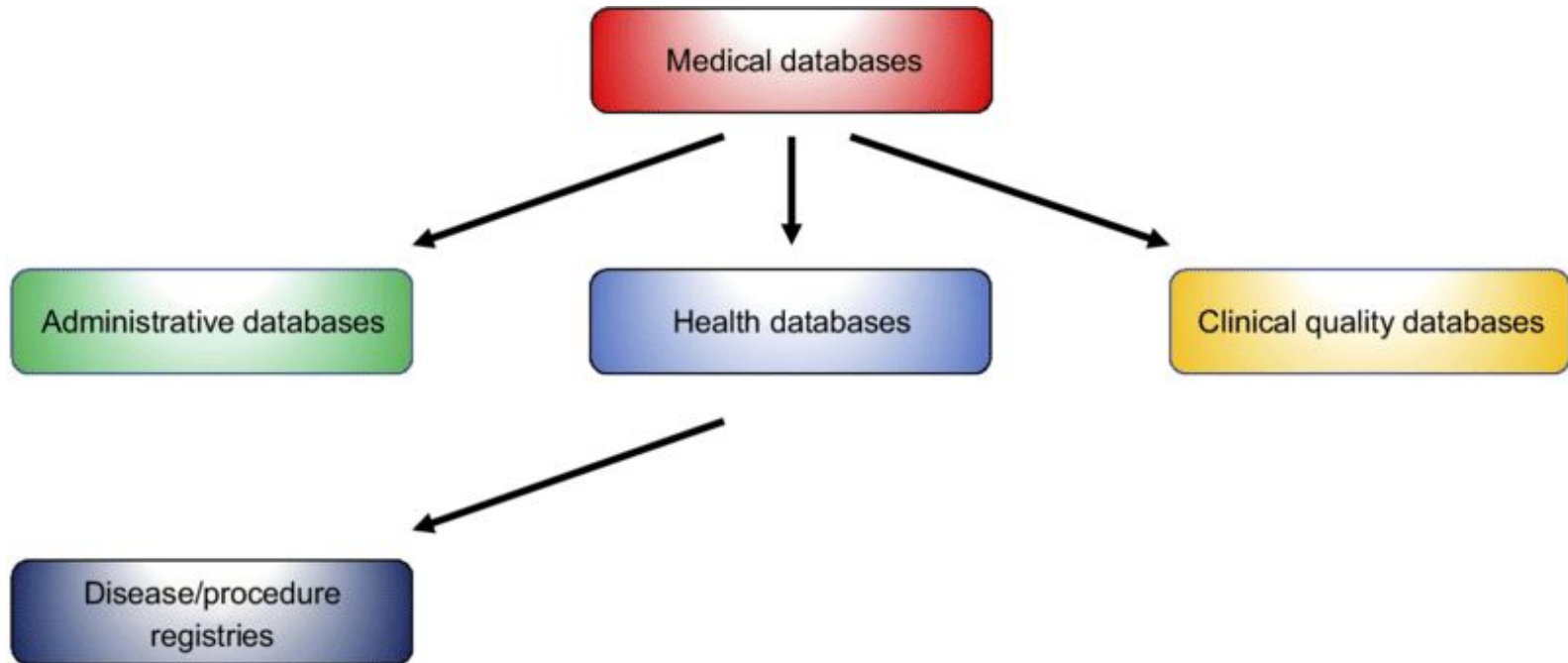
- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research

# Many types of database



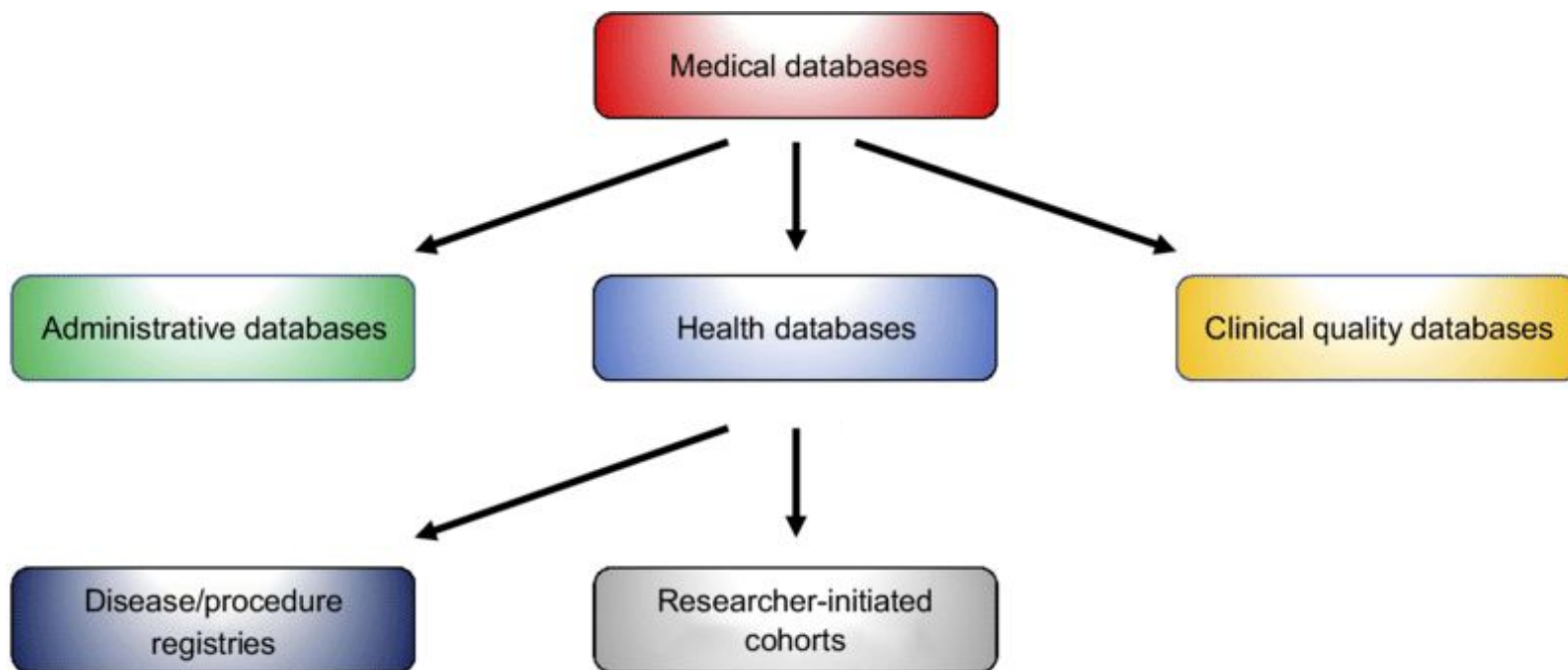
- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research
- Register detailed clinical data for clinical quality control

# Many types of database



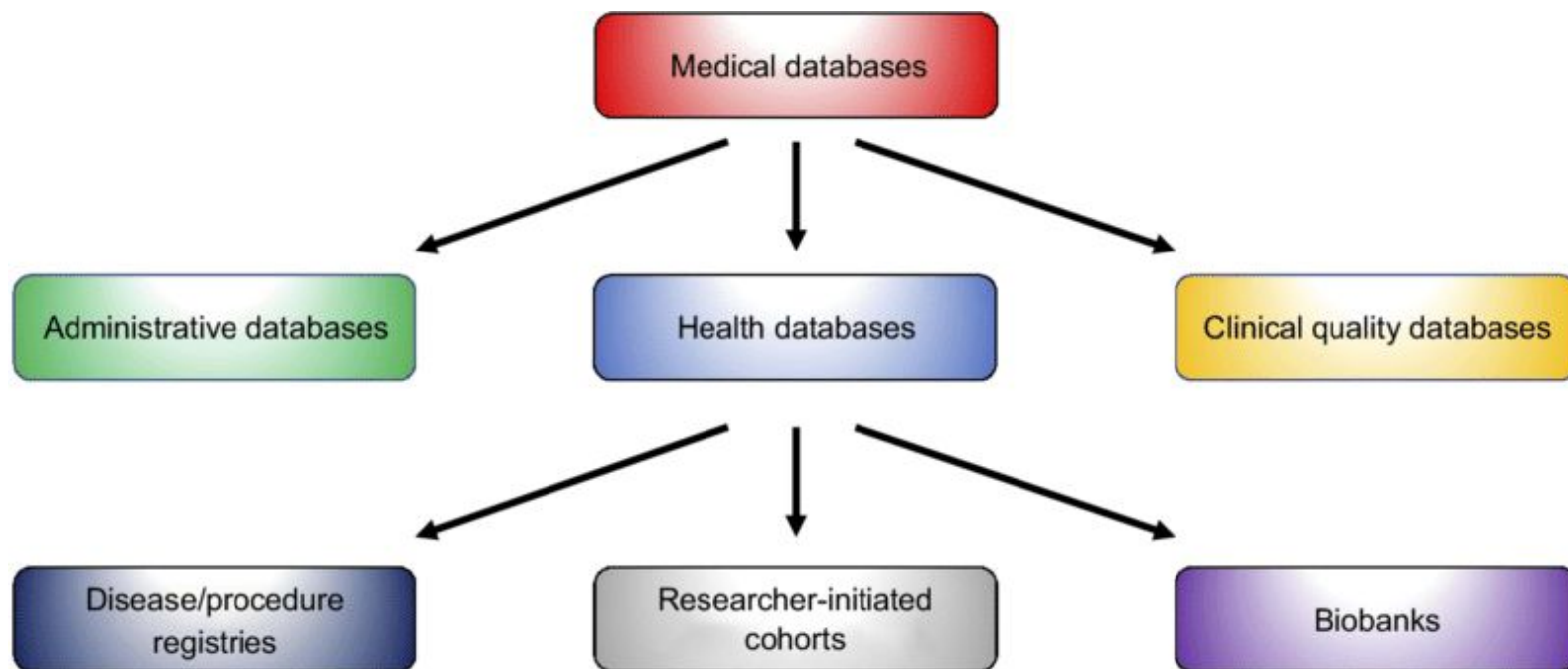
- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research
- Register detailed clinical data for clinical quality control
- Register patients according to diagnosis or procedure

# Many types of database



- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research
- Register detailed clinical data for clinical quality control
- Register patients according to diagnosis or procedure
- Register individuals according to prespecified criteria (eg, area of residency, age, sex, conscription, adoption, pregnancy, or survey participation)

# Many types of database



- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research
- Register detailed clinical data for clinical quality control
- Register patients according to diagnosis or procedure
- Register individuals according to prespecified criteria (eg, area of residency, age, sex, conscription, adoption, pregnancy, or survey participation)
- Store biological samples (eg, blood and tissue)

# Analysis based on primary record type & sampling scope

## Primary Record Type:

- Individual procedures e.g., arthroplasty
- Prescriptions e.g., colistin
- Disease/Illness e.g., ovarian cancer
- Hospital Admission/Discharge
- Individual health interactions
- Patient
- Person
- Population

# Analysis based on primary record type & sampling scope

## Primary Record Type:

- Individual procedures e.g., arthroplasty
- Prescriptions e.g., colistin
- Disease/Illness e.g., ovarian cancer
- Hospital Admission/Discharge
- Individual health interactions
- Patient
- Person
- Population

## Sampling Scope:

- Single physician
- Group of physicians
- Hospital
- Health Authority
- Province
- National
- International



Generalisability

# Analysis based on primary record type & sampling scope

## Primary Record Type:

- Individual procedures e.g., arthroplasty
- Prescriptions e.g., colistin
- Disease/Illness e.g., ovarian cancer
- Hospital Admission/Discharge
- Individual health interactions
- Patient
- Person
- Population

## Sampling Scope:

- Single physician
- Group of physicians
- Hospital
- Health Authority
- Province
- National
- International



Challenge of standardisation

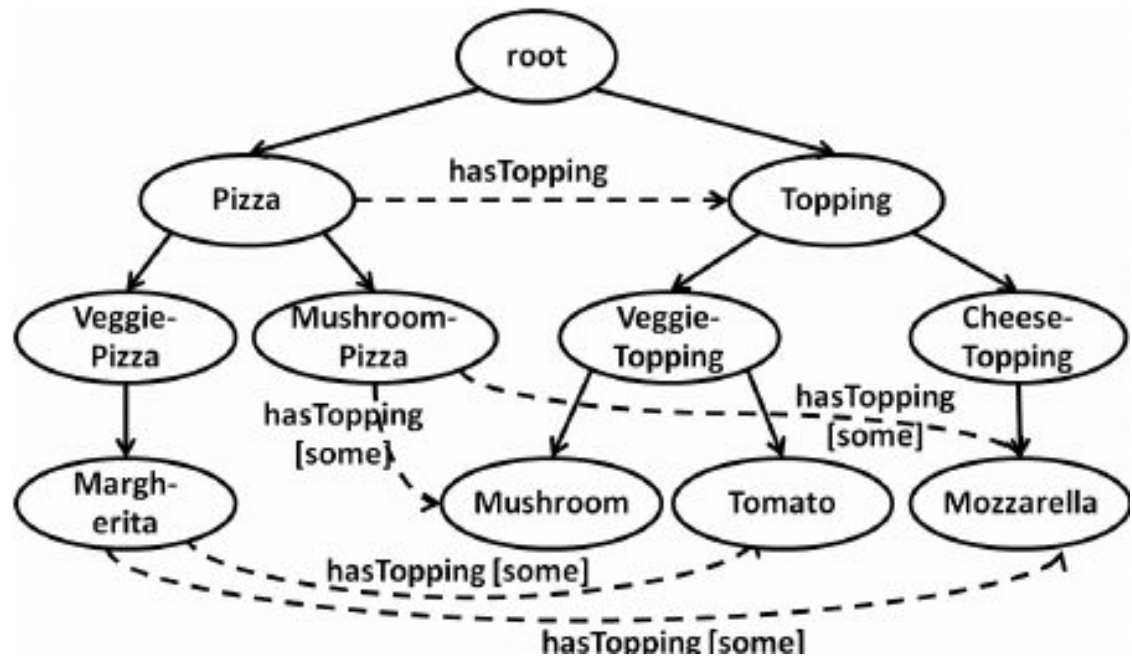


Generalisability

How do medical databases try to handle standardisation?

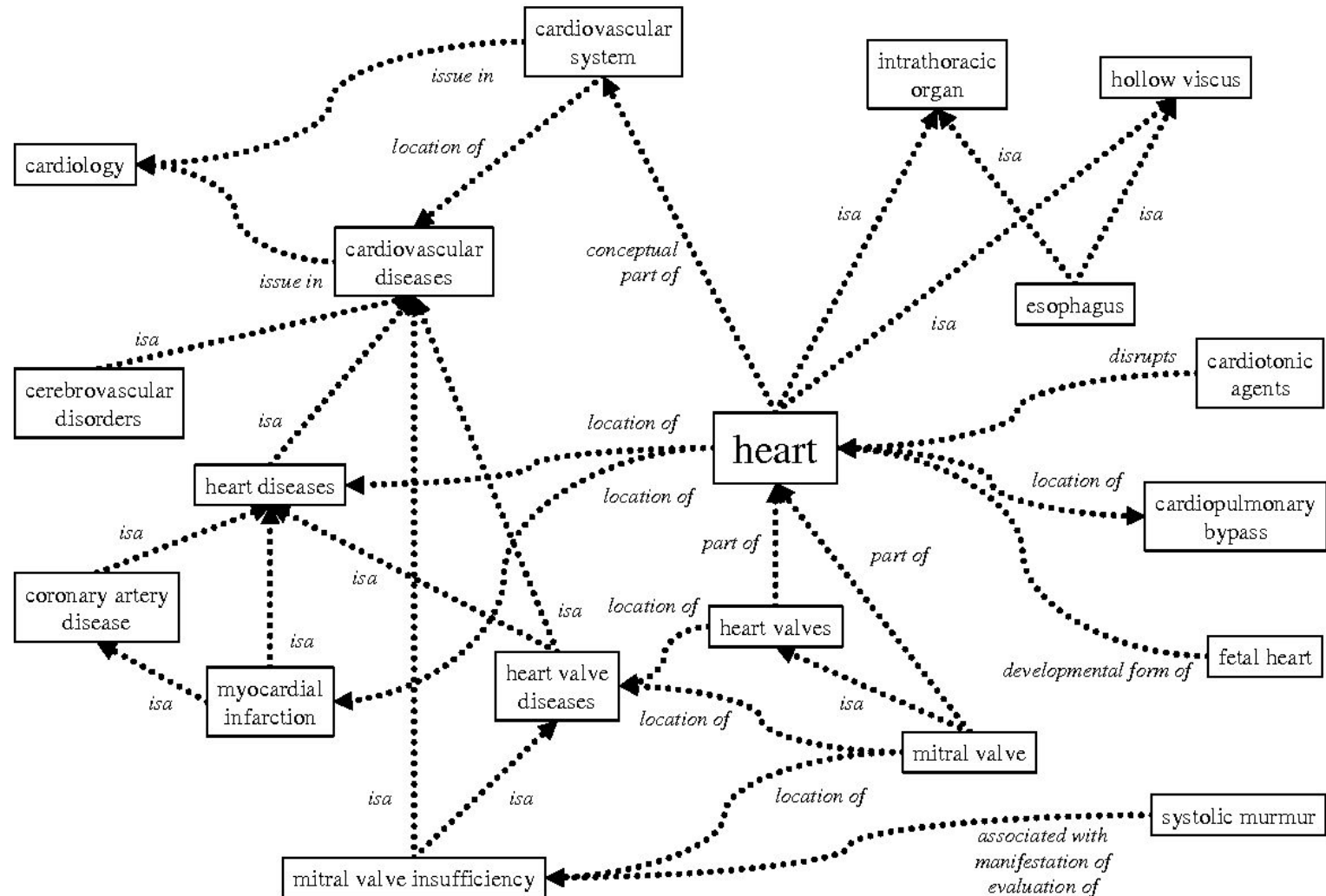
# Ontologies for standardisation

- Standardised terms e.g., Pizza, Tomato, Mozzarella
- Standardised types of relationships between terms
- Acyclic links between terms
- Manual curation
- Automated curation

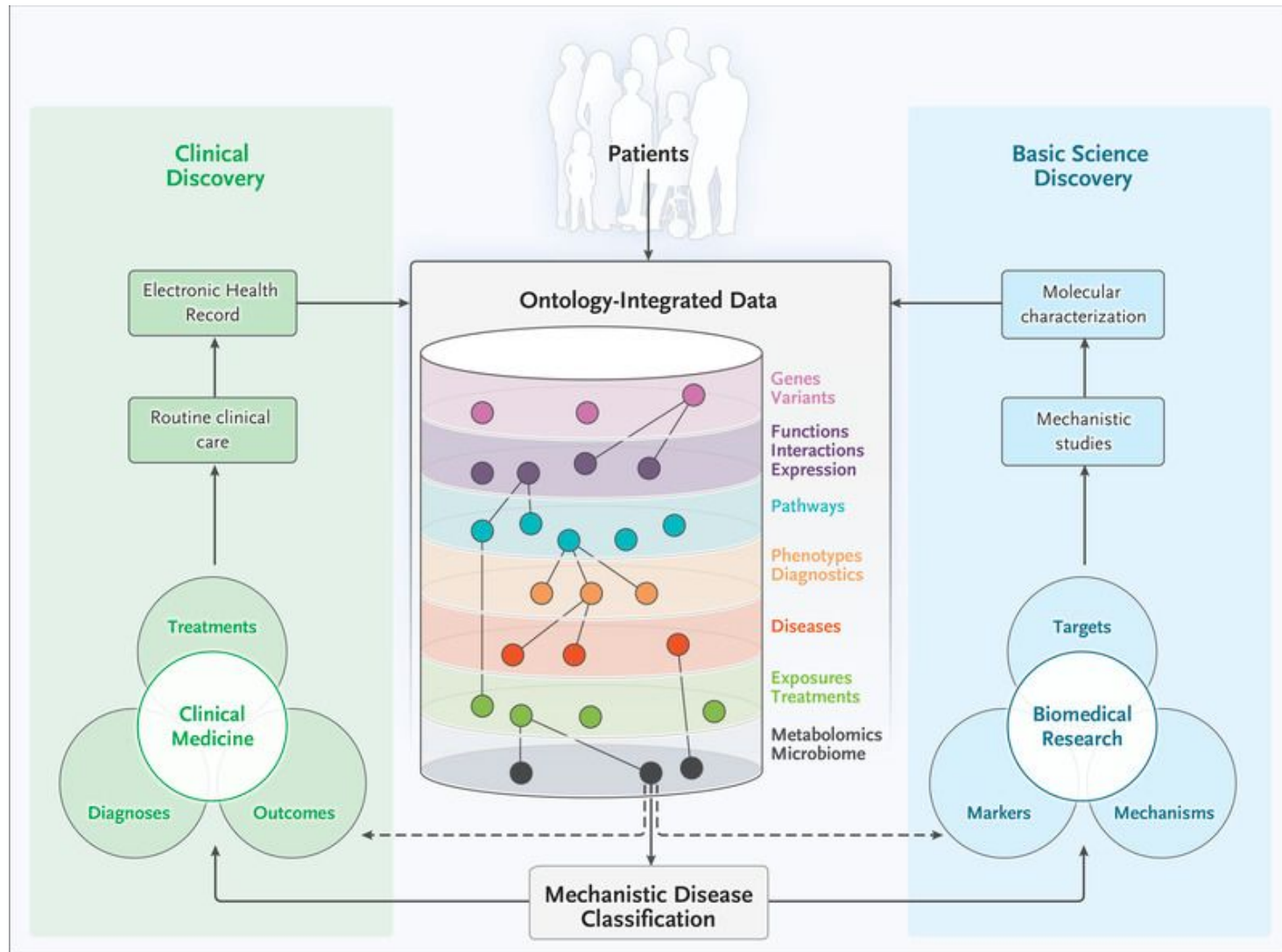


[https://www.researchgate.net/figure/Example-pizza-ontology-represented-as-a-graph-G-a-and-a-changed-variant-of-the-pizza\\_fig1\\_236842047](https://www.researchgate.net/figure/Example-pizza-ontology-represented-as-a-graph-G-a-and-a-changed-variant-of-the-pizza_fig1_236842047)

# Medical Ontologies



# Ontologies for linking diverse types of data



# International Statistical Classification of Diseases and Related Health Problems (ICD-9, ICD-10)

- 2 ontologies
  - ICD-X-CM (medical diagnoses)
  - ICD-X-PCS (procedure coding)

# International Statistical Classification of Diseases and Related Health Problems (ICD-9, ICD-10)

- 2 ontologies
  - ICD-X-CM (medical diagnoses)
  - ICD-X-PCS (procedure coding)
- ICD-9 -> ICD-10 (2015)

| Differences Between ICD-9-CM and ICD-10 Code Sets |              |                  |
|---|--------------|------------------|
|   | ICD-9-CM     | ICD-10 code sets |
| Procedure   | 3,824 codes  | 71,924 codes     |
| Diagnosis   | 14,025 codes | 69,823 codes     |

| ICD-10 Code Structure Changes (selected details) |   |  |
|--|---|--|
|  | Old   | New  |
| Diagnosis Structure                              | ICD-9-CM <ul style="list-style-type: none"> <li>• 3-5 characters</li> <li>• First character is numeric or alpha</li> <li>• Characters 2-5 are numeric</li> </ul>  | ICD-10-CM <ul style="list-style-type: none"> <li>• 3-7 characters</li> <li>• Character 1 is alpha</li> <li>• Character 2 is numeric</li> <li>• Characters 3 – 7 can be alpha or numeric</li> </ul> |
| Procedure Structure                              | ICD-9-CM <ul style="list-style-type: none"> <li>• 3-4 characters</li> <li>• All characters are numeric</li> <li>• All codes have at least 3 characters</li> </ul> | ICD-10-PCS <ul style="list-style-type: none"> <li>• ICD-10-PCS has 7 characters</li> <li>• Each can be either alpha or numeric</li> <li>• Numbers 0-9; letters A-H, J-N, P-Z</li> </ul>            |

[https://www.cdc.gov/nchs/icd/icd10cm\\_pcs\\_background.htm](https://www.cdc.gov/nchs/icd/icd10cm_pcs_background.htm)

# International Statistical Classification of Diseases and Related Health Problems (ICD-9, ICD-10)

- 2 ontologies
  - ICD-X-CM (medical diagnoses)
  - ICD-X-PCS (procedure coding)
- ICD-9 -> ICD-10 (2015)
- “V97.33XD: Sucked into jet engine, subsequent encounter.”
- “Y93.D: V91.07XD: Burn due to water-skis on fire, subsequent encounter.”
- “Z63.1: Problems in relationship with in-laws.”
- “W22.02XD: V95.43XS: Spacecraft collision injuring occupant, sequela.”

| Differences Between ICD-9-CM and ICD-10 Code Sets |              |                  |
|---|--------------|------------------|
|   | ICD-9-CM     | ICD-10 code sets |
| Procedure   | 3,824 codes  | 71,924 codes     |
| Diagnosis   | 14,025 codes | 69,823 codes     |

| ICD-10 Code Structure Changes (selected details) |   |  |
|--|---|--|
|  | Old   | New  |
| Diagnosis Structure                              | ICD-9-CM <ul style="list-style-type: none"> <li>• 3-5 characters</li> <li>• First character is numeric or alpha</li> <li>• Characters 2-5 are numeric</li> </ul>  | ICD-10-CM <ul style="list-style-type: none"> <li>• 3-7 characters</li> <li>• Character 1 is alpha</li> <li>• Character 2 is numeric</li> <li>• Characters 3 – 7 can be alpha or numeric</li> </ul> |
| Procedure Structure                              | ICD-9-CM <ul style="list-style-type: none"> <li>• 3-4 characters</li> <li>• All characters are numeric</li> <li>• All codes have at least 3 characters</li> </ul> | ICD-10-PCS <ul style="list-style-type: none"> <li>• ICD-10-PCS has 7 characters</li> <li>• Each can be either alpha or numeric</li> <li>• Numbers 0-9; letters A-H, J-N, P-Z</li> </ul>            |

[https://www.cdc.gov/nchs/icd/icd10cm\\_pcs\\_background.htm](https://www.cdc.gov/nchs/icd/icd10cm_pcs_background.htm)

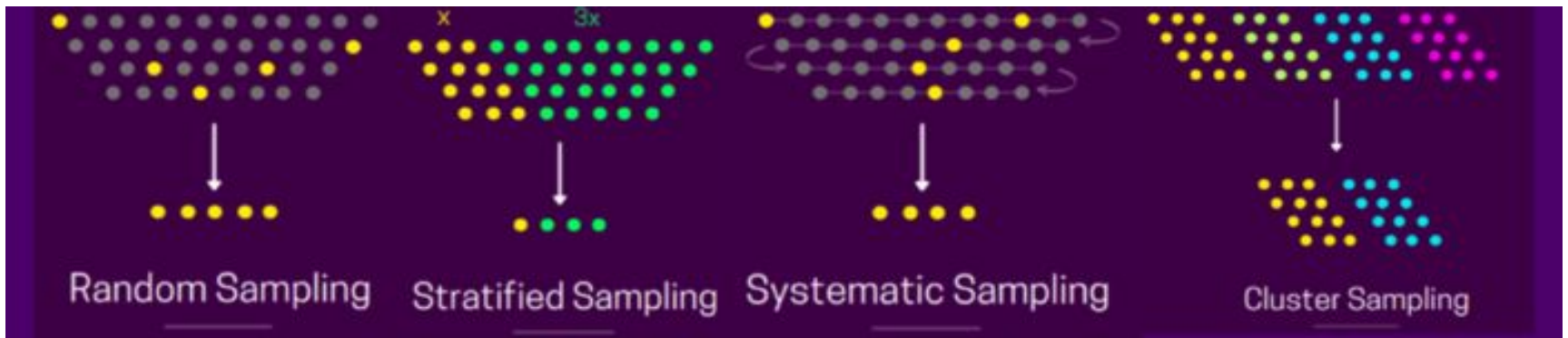
How do we sample from medical databases?

# Sampling strategy

- Exhaustive in a database isn't usually exhaustive in true population
- Numerous and often quite complex!
- Major source of bias so always carefully explore

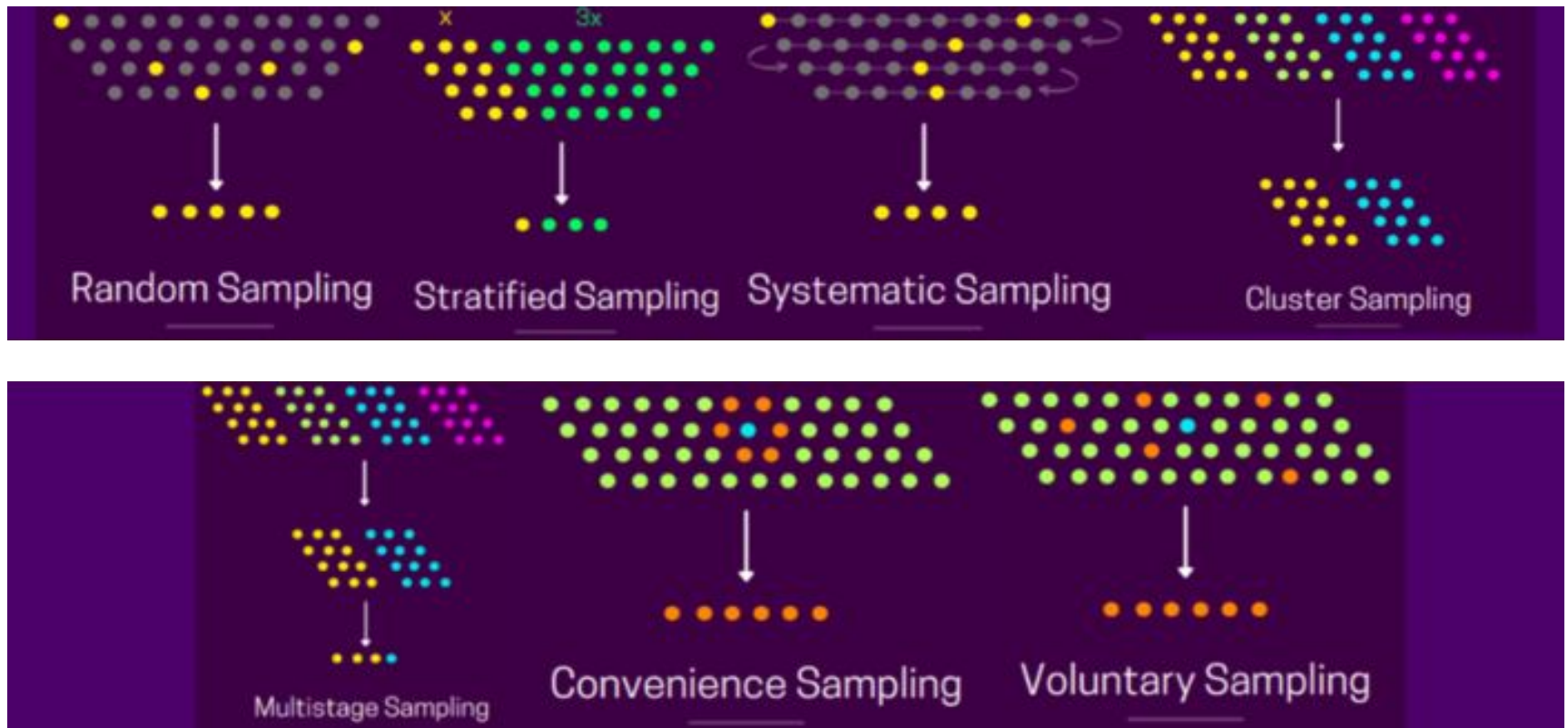
# Sampling strategy

- Exhaustive in a database isn't usually exhaustive in true population
- Numerous and often quite complex!
- Major source of bias so always carefully explore



# Sampling strategy

- Exhaustive in a database isn't usually exhaustive in true population
- Numerous and often quite complex!
- Major source of bias so always carefully explore

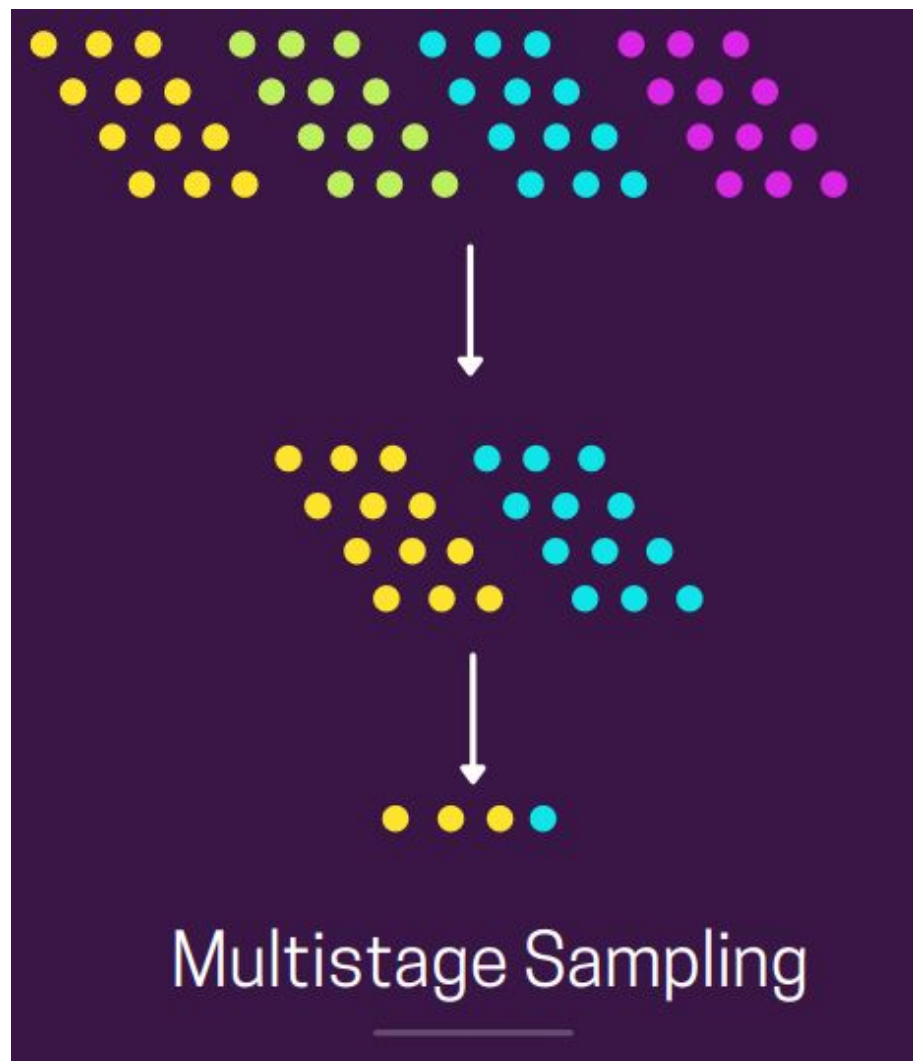


# Survey/Sample weights

- Value/weight assigned to each record
- Make statistics calculated from database more representative of population
  - Weight=0.5 underweight this case
  - Weight=1
  - Weight=2 overweight the contribution of this case

# Survey/Sample weights

- Value/weight assigned to each record
- Make statistics calculated from database more representative of population
  - Weight=0.5 underweight this case
  - Weight=1
  - Weight=2 overweight the contribution of this case
- Complex sampling strategies (e.g., deliberate oversampling of some populations, biasing recruitment) mean weights **MUST** be used.
- Not directly supported in all machine learning libraries (`sample_weights` implemented for some models)

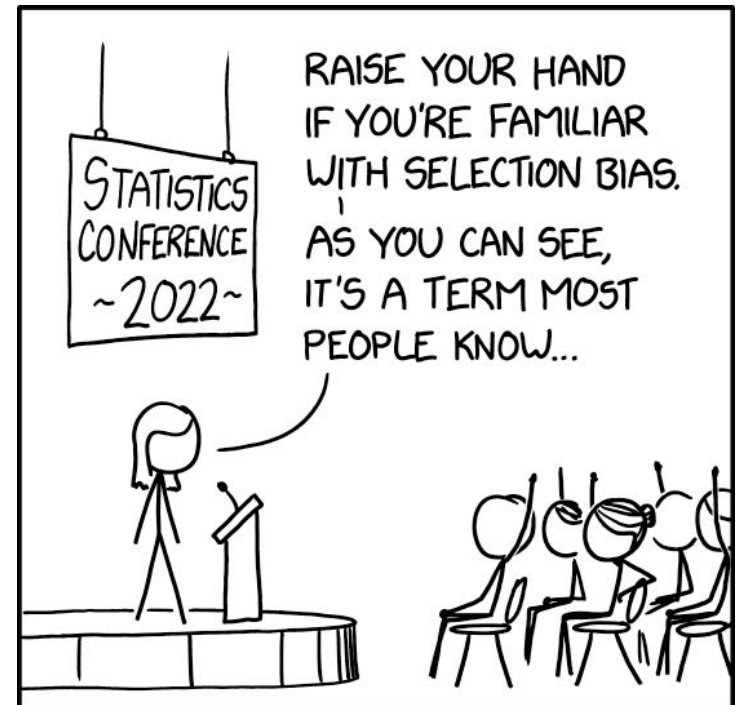


# Types of weights

- Design Weights
  - Based on sampling strategy i.e., “design” of survey/database/data collection
  - Common to over-sample under-represented or rare groups
  - Need to correct for this or will overestimate statistics e.g., lower weight of over-sampled groups

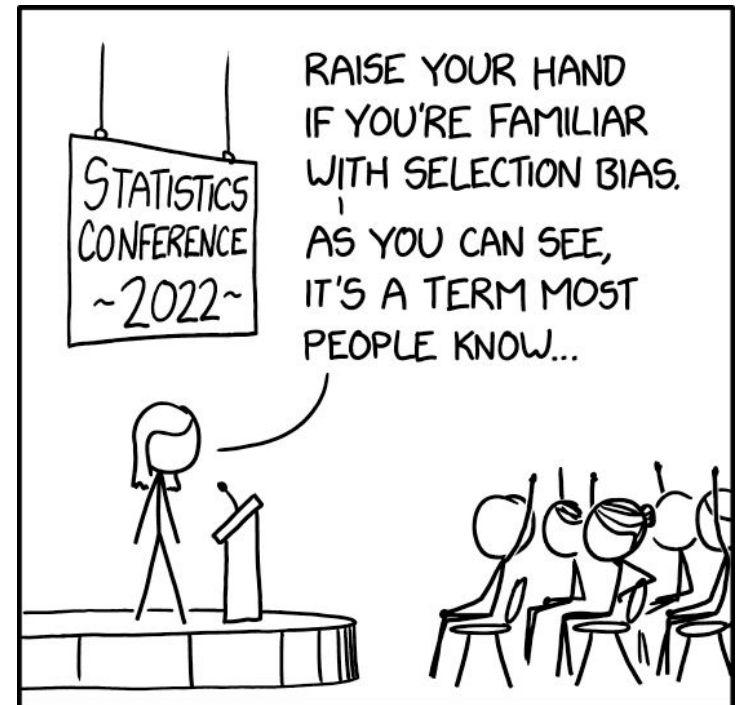
# Types of weights

- Design Weights
  - Based on sampling strategy i.e., “design” of survey/database/data collection
  - Common to over-sample under-represented or rare groups
  - Need to correct for this or will overestimate statistics e.g., lower weight of over-sampled groups
- Post-stratification / Non-response weights
  - Based on collected data
  - Typically biases in whose data is collected
  - Over-represented groups need to be under-weighted



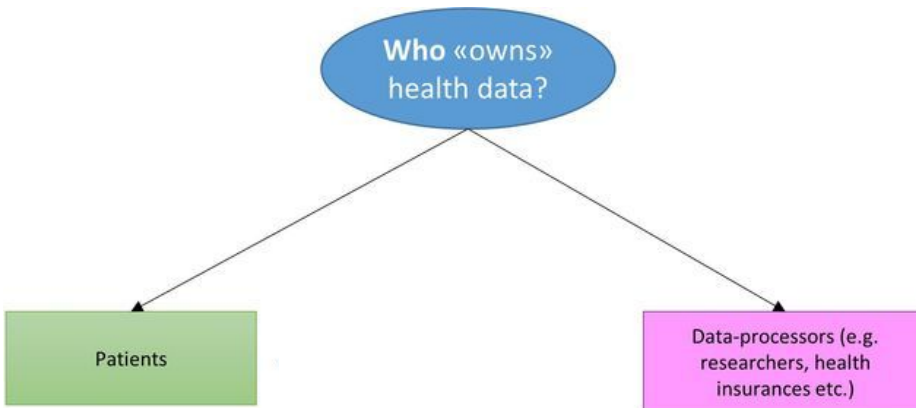
# Types of weights

- Design Weights
  - Based on sampling strategy i.e., “design” of survey/database/data collection
  - Common to over-sample under-represented or rare groups
  - Need to correct for this or will overestimate statistics e.g., lower weight of over-sampled groups
- Post-stratification / Non-response weights
  - Based on collected data
  - Typically biases in whose data is collected
  - Over-represented groups need to be under-weighted
- Often many different weights are combined

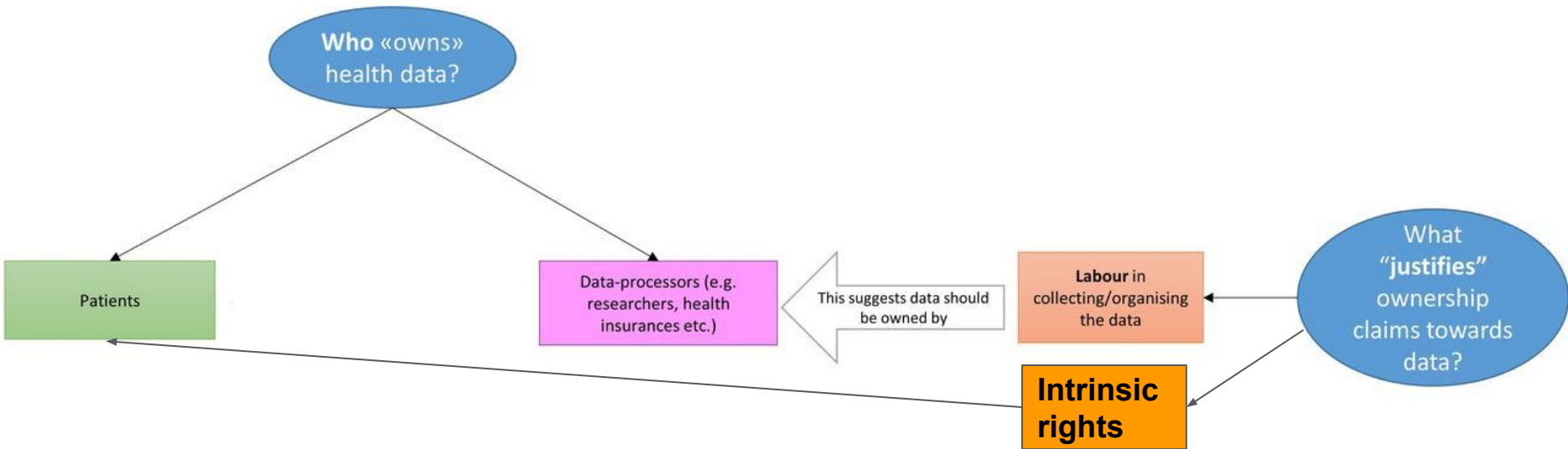


Who actually owns this data?

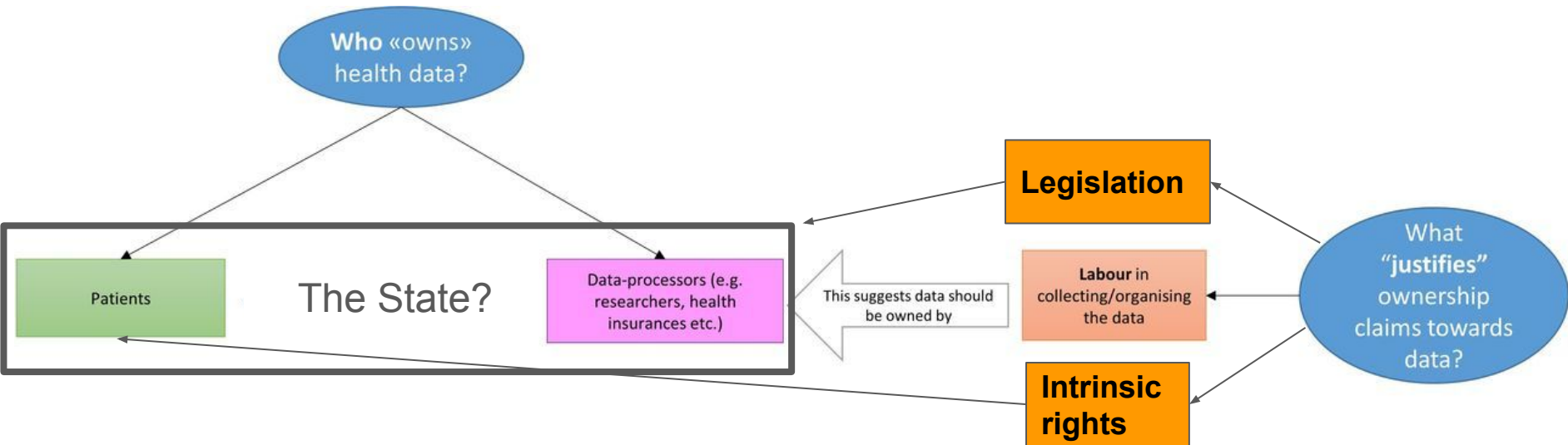
# Data ownership is complicated



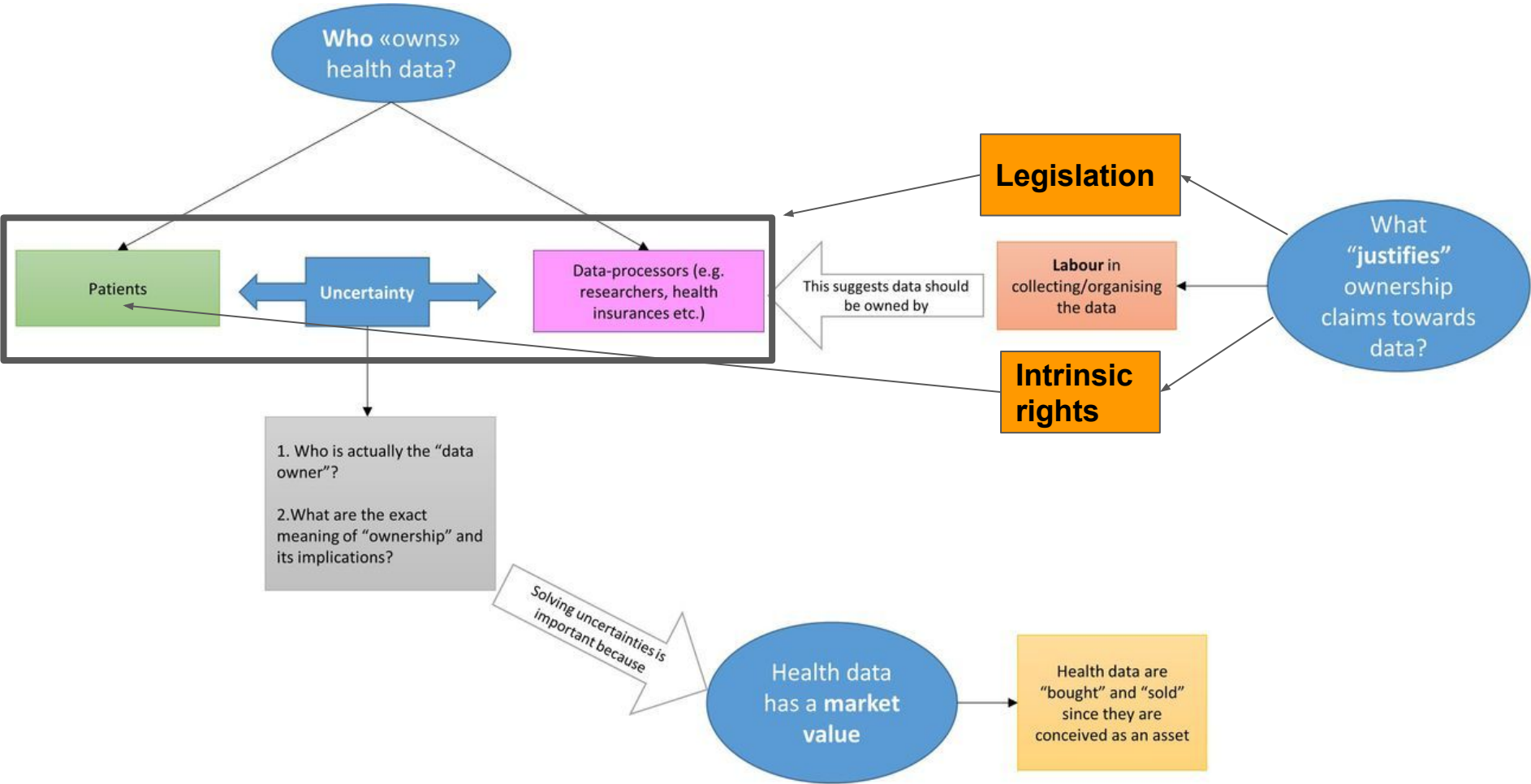
# Data ownership is complicated



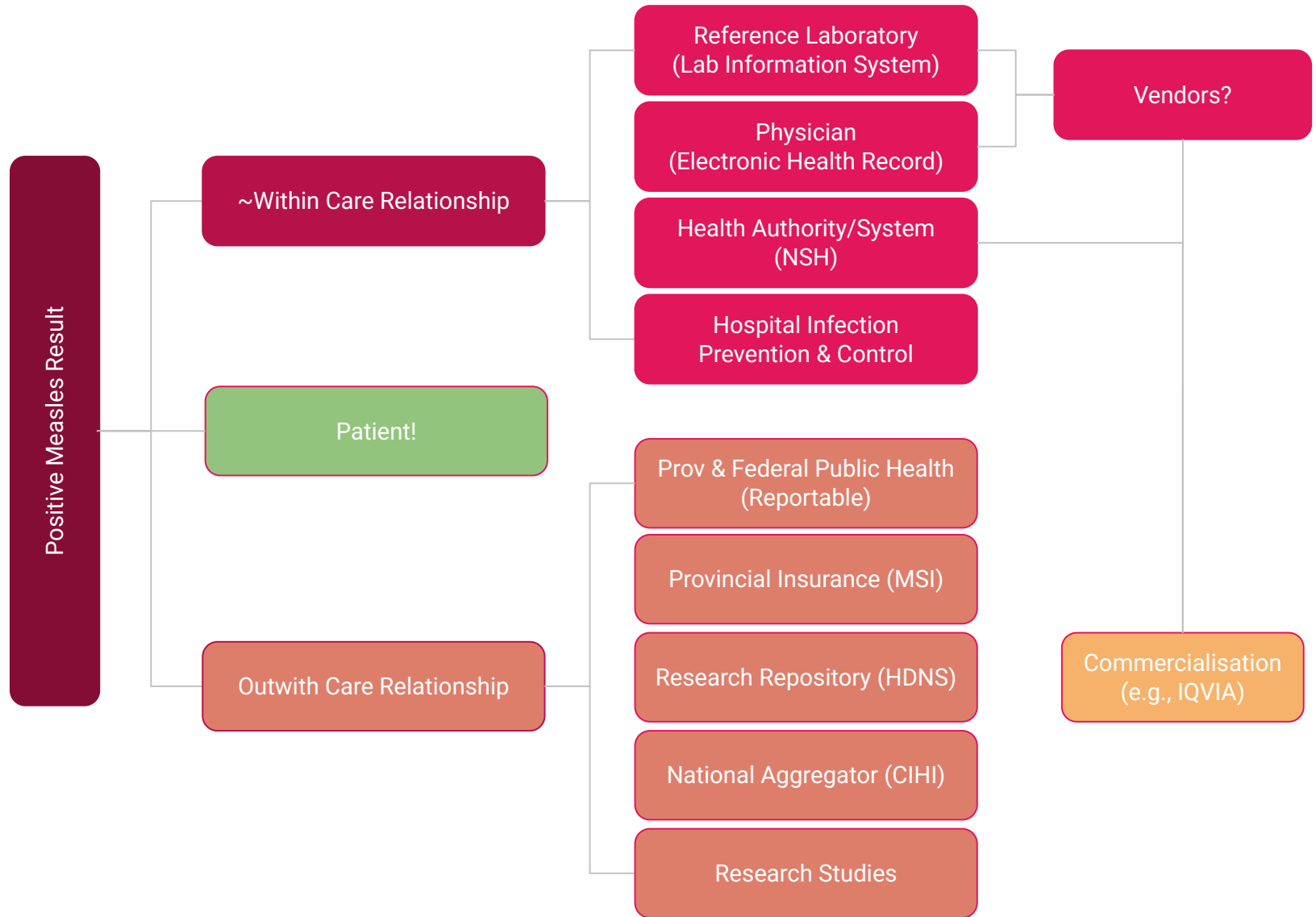
# Data ownership is complicated



# Data ownership is complicated

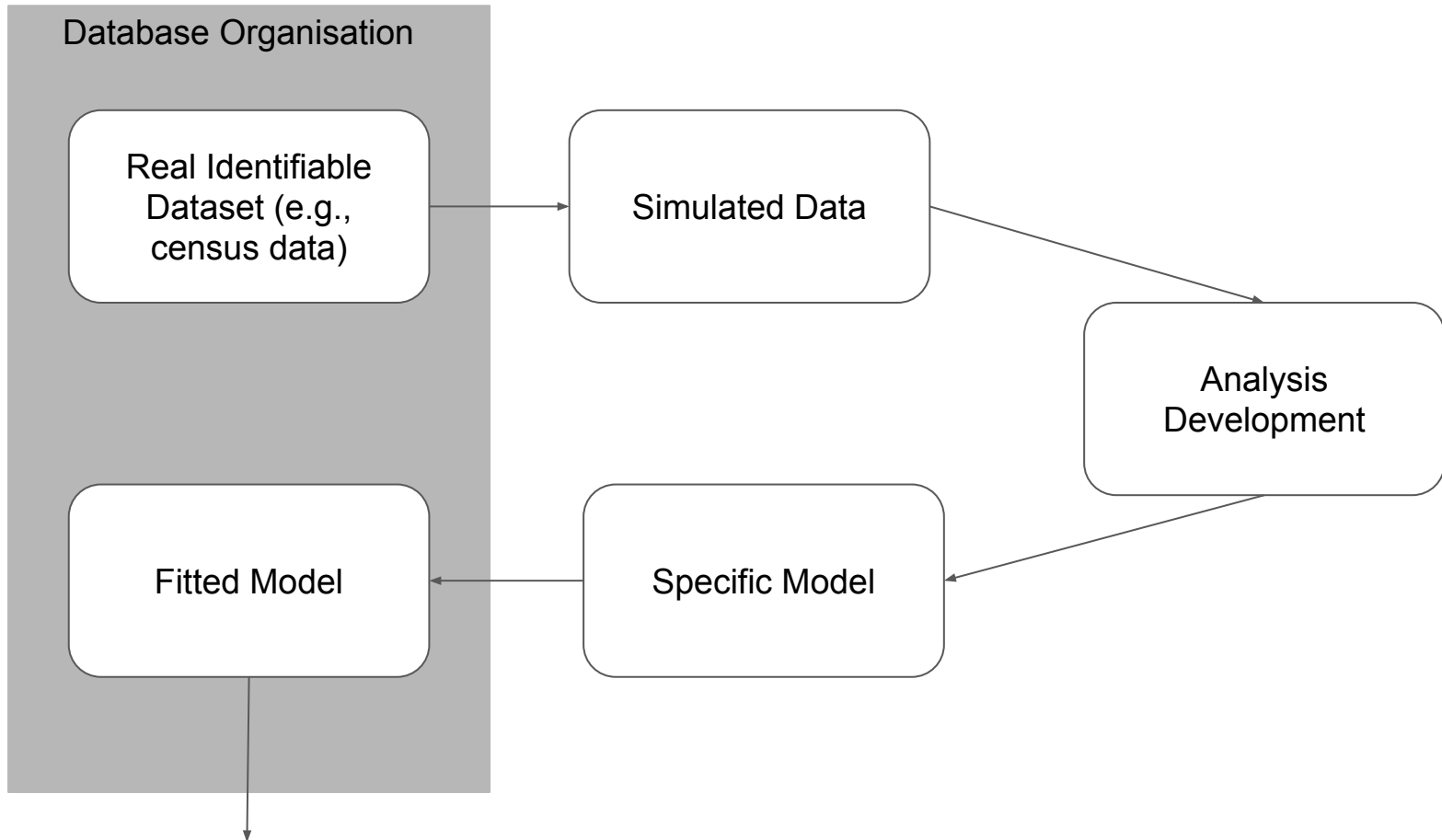


# Ownership over which copy of same data?

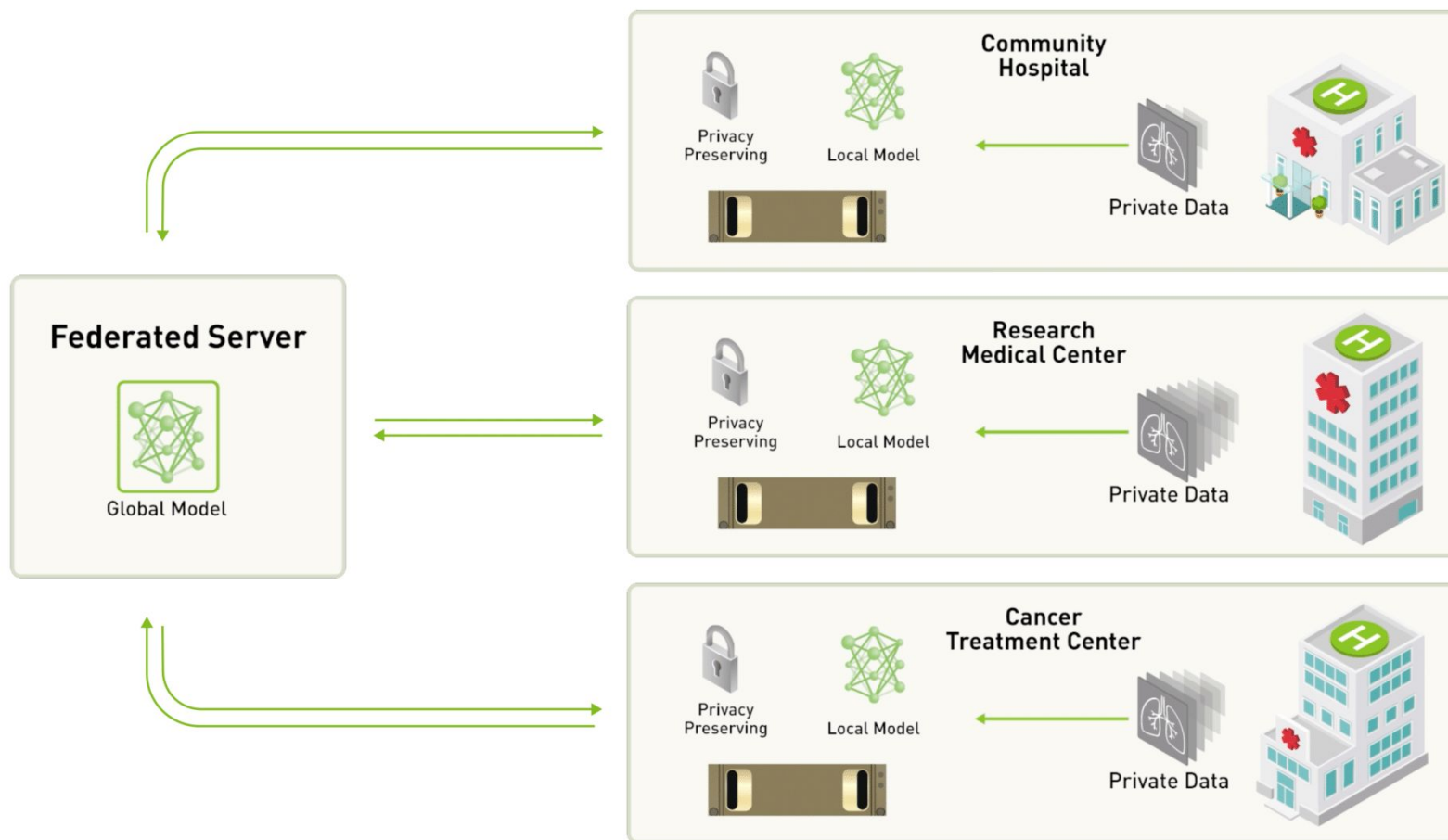


How do you protect privacy in these  
databases?

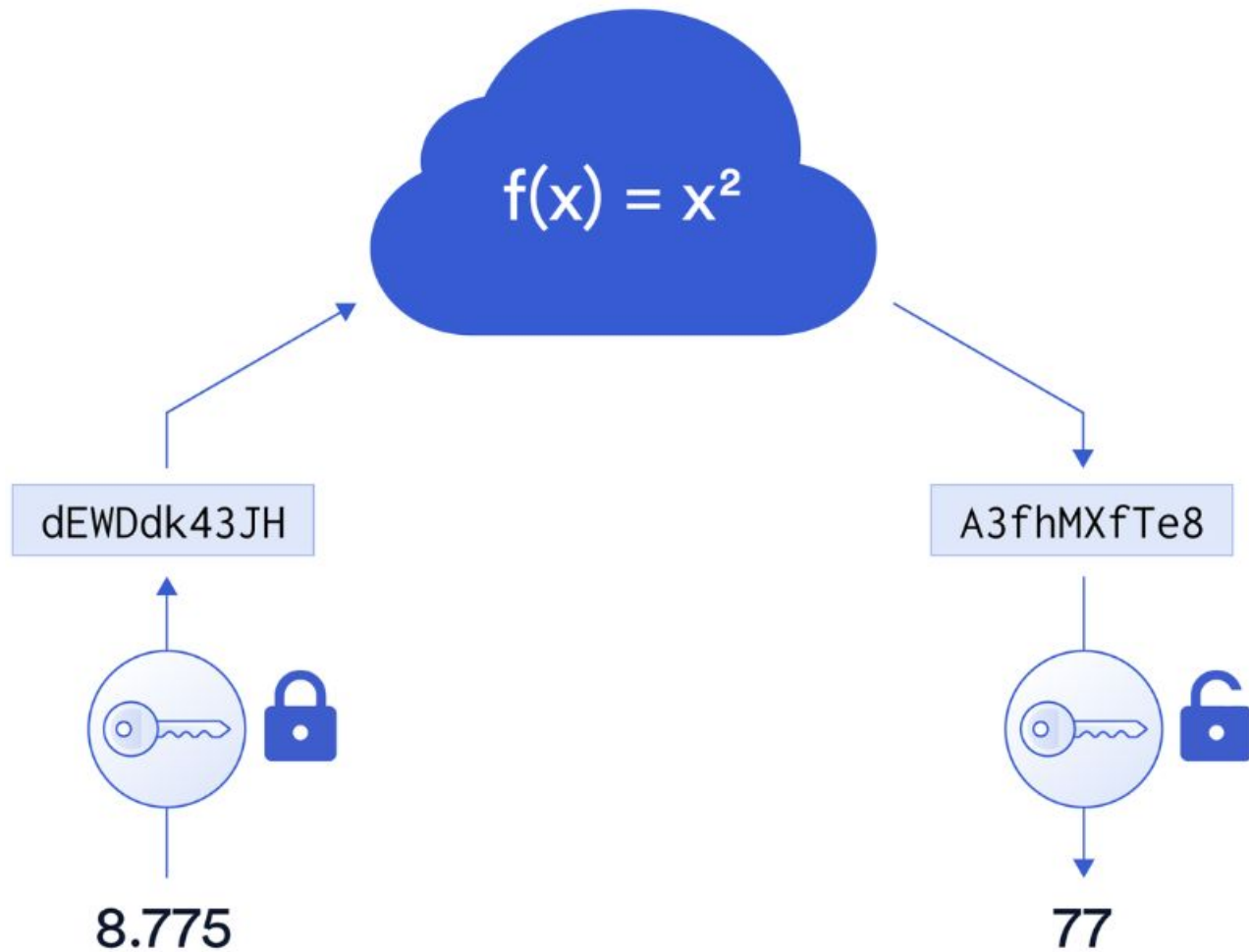
# No direct data access



# No external data access: federated training



# Shared data but encrypted: homomorphic encryption



These are difficult and limited... so how can we share data directly but safely?

# Data privacy is a continuum

## **Direct Identifiers:**

- Data that identifies a person without additional information or by linking to information in the public domain (e.g., name, SIN, MSI Number)

## **Indirect Identifiers:**

- Data that identifies an individual indirectly but can so in aggregate (e.g., DoB, gender, IP address, license plate)

## **Pseudonymous Data:**

- Direct identifiers eliminated or transformed but indirect intact (e.g., names replace with unique identifier)

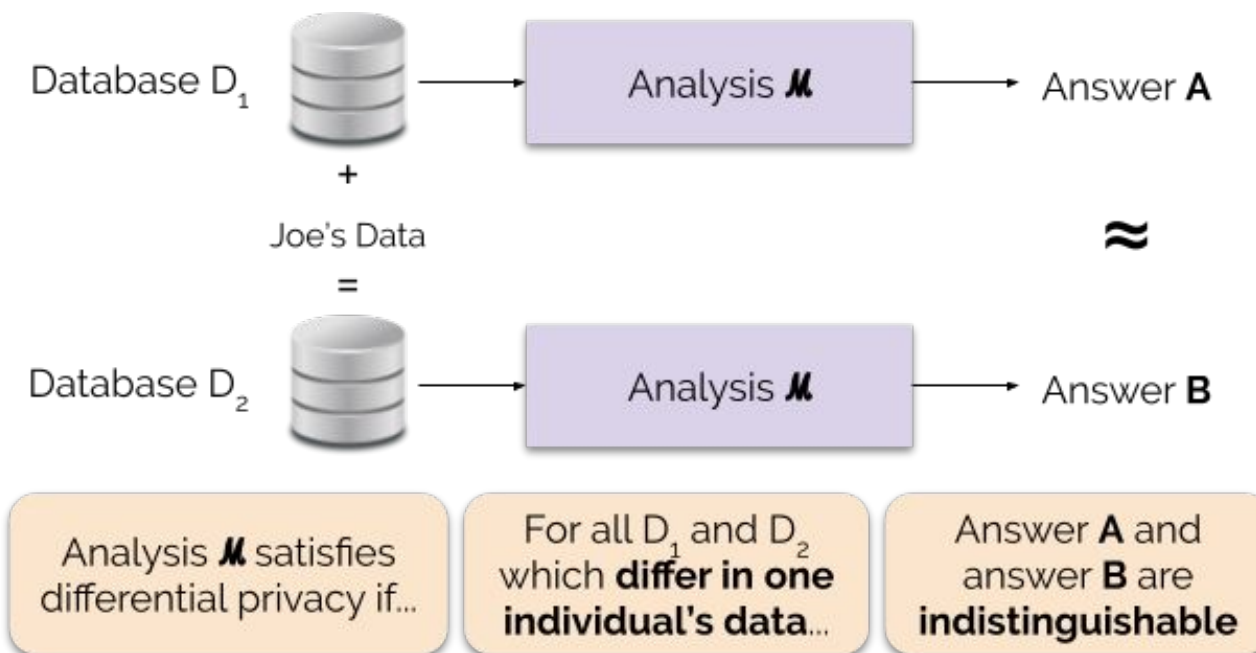
## **De-identified Data:**

- Direct and indirect identifiers removed or modified to break linkage to real-world identities (e.g., data suppressed, generalized, perturbed, swapped - ranges applied)

## **Anonymous Data:**

- Direct and indirect identifiers removed or modified to ensure mathematical guarantees against re-identification (e.g., aggregated census data, noise added to dataset)

# Differential privacy: no singling out individuals



# Differential privacy: no singling out individuals



Analysis  $\mathcal{M}$  satisfies differential privacy if...

For all  $D_1$  and  $D_2$  which **differ in one individual's data...**

Answer **A** and answer **B** are **indistinguishable**

Probability of seeing output  $O$  on input  $D_1$  →  $\Pr[\mathcal{M}(D_1) \in O]$

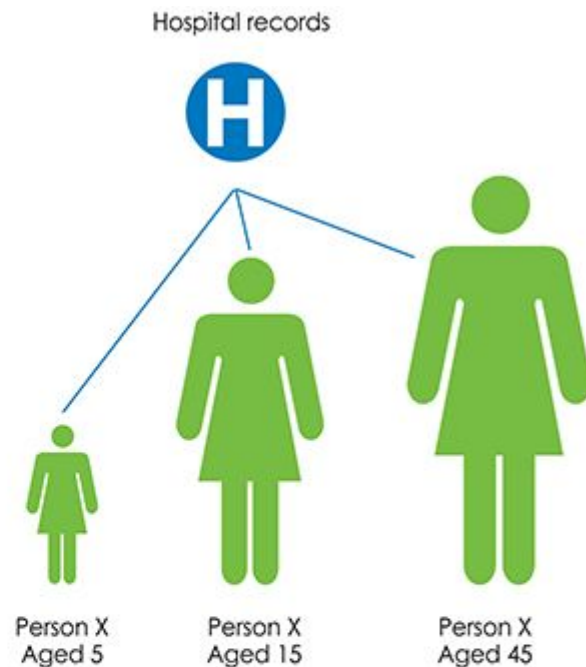
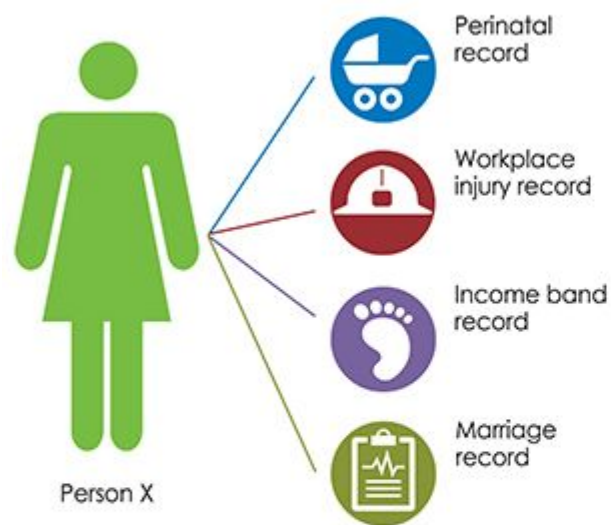
Probability of seeing output  $O$  on input  $D_2$  →  $\Pr[\mathcal{M}(D_2) \in O]$

$$\frac{\Pr[\mathcal{M}(D_1) \in O]}{\Pr[\mathcal{M}(D_2) \in O]} \leq e^\epsilon$$

Indistinguishability: bounded ratio of probabilities

# Data linkage is powerful but dangerous

- Linking between databases and resources -> identifiability
- Can be done probabilistically
- Often needs additional ethics/applications
- Can break a lot of data privacy operations



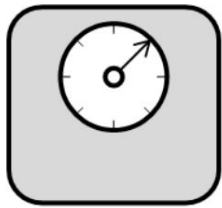
So, you've got access to a database, what  
now?

# Data Cleaning: even “simple” fields can be a nightmare

## Data Quality



Actual value:  
200.6 lbs.

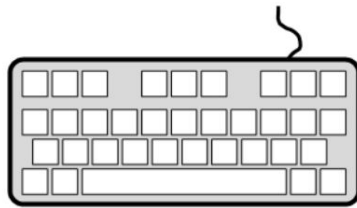


### Measured (same day)

- Validity challenge  
198.9 | 198.9 | 198.9 lbs.
- Reliability challenge  
200.6 | 198.9 | 202.2 lbs.

### Measured (diff. days)

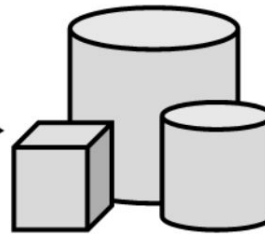
Recorded value:  
200 lbs.



### User Typed (one entry)

- Typos  
200.6 lbs. → 20.06, 2006
- Mismatching units  
200.6 lbs. → 200.6 kg
- Assumptions/truncations  
200.6 lbs. → 200 lbs.  
NULL → 0
- Free-text additions  
200.6 lbs. → 200.6 pounds

Data warehouse value:  
200 kg

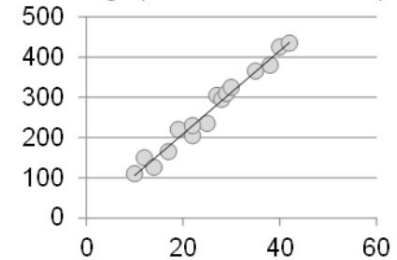


### DB Operations (one entry)

- Truncations/Rounding  
200.6 → 200.0
- Error conversions  
200.6 pounds → NULL  
200.6 lbs. → 200.6 kg
- Cleaning  
200+ lbs. → 200.0

Analytic value:

100 kg (mean 200 & 0)



### Analytics (data points)

- Aggregation of data points  
200 | 0 → mean of 100
- Selecting a representative  
190 | 200 | 210 → 210 (first)  
190 | 200 | 210 → 200 (mean)  
190 | 200 | 210 → 210 (last)
- Removing outliers  
200 | 200 | 350 → 200 | 200 | NULL

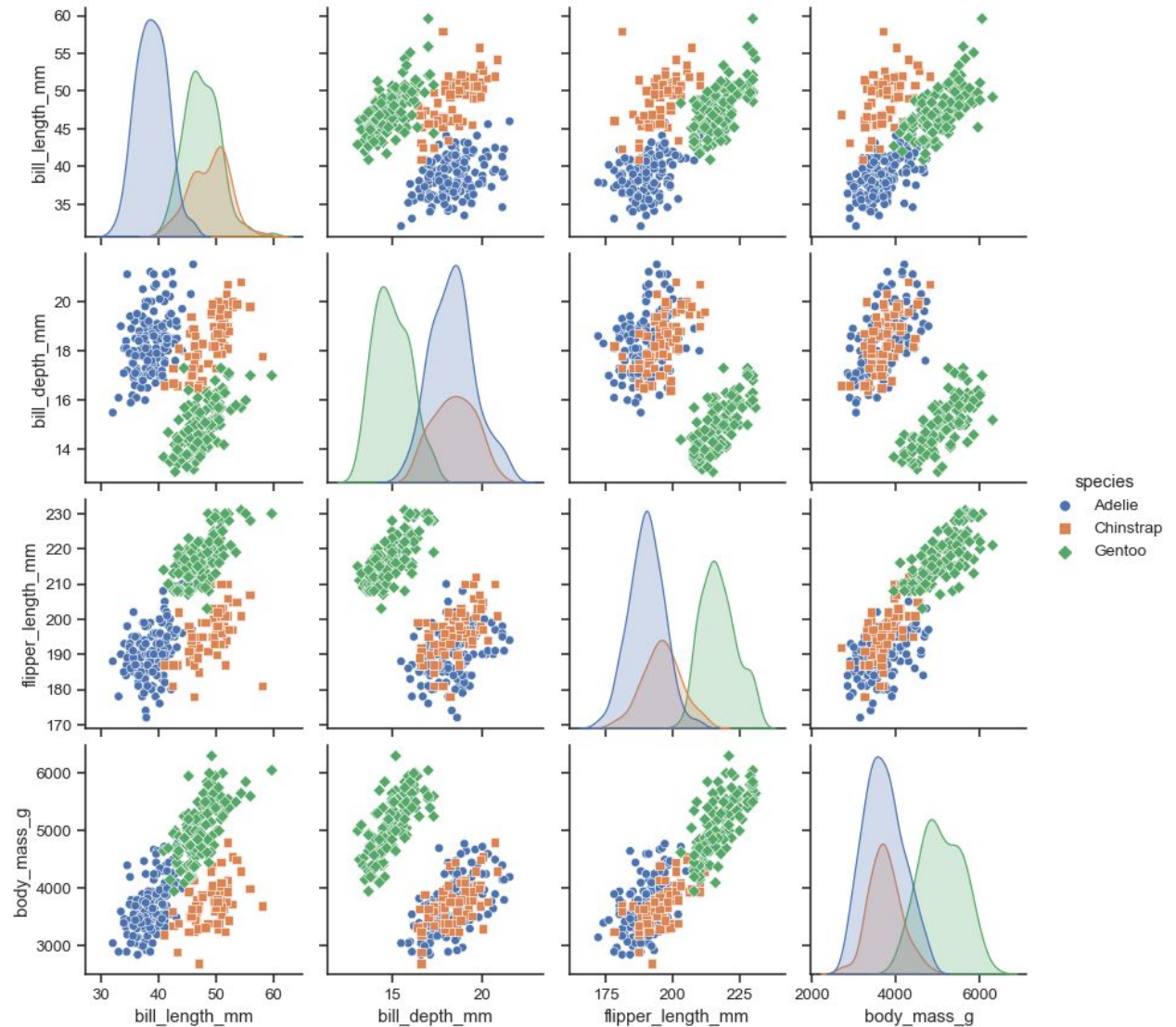
*Under review*

Slide from Dr. Hadi Kharrazi

# 9 months & >25 rules to clean weight

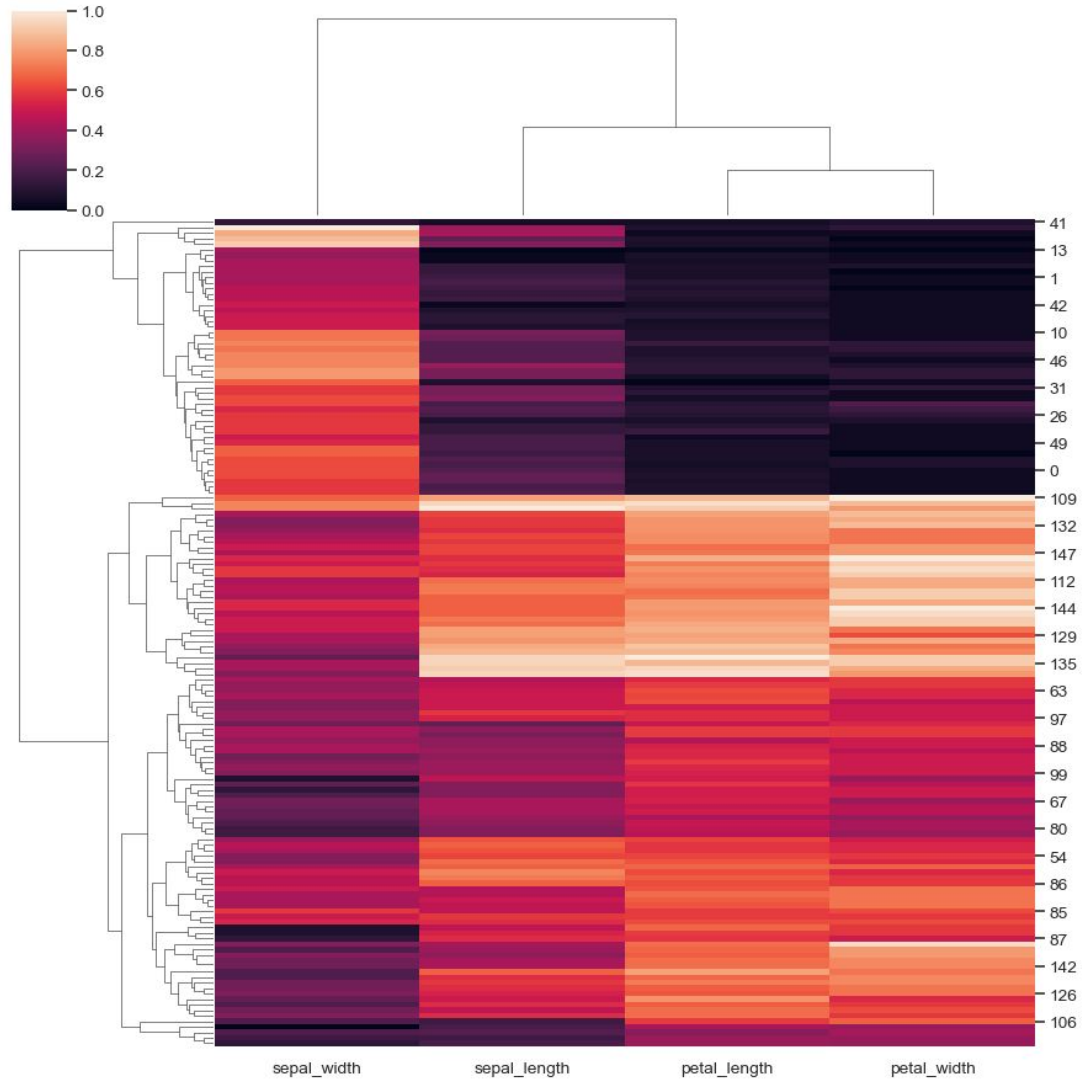
# Exploratory Data Analysis

- Individual variable distributions
- Pairwise variable distributions
- Distributions relative to variable(s) of interest
- Point analysis of extreme values



# Exploratory Data Analysis

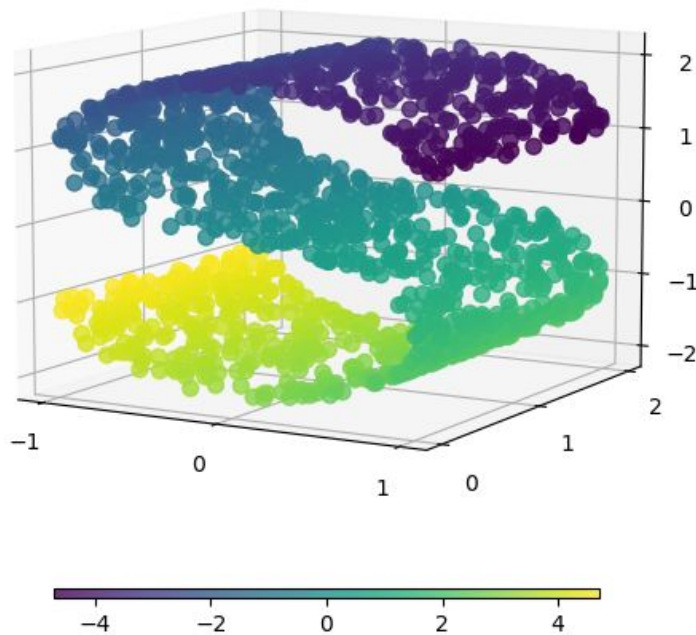
- Individual variable distributions
- Pairwise variable distributions
- Distributions relative to variable(s) of interest
- Hierarchical clustering of variables
- Point analysis of extreme values



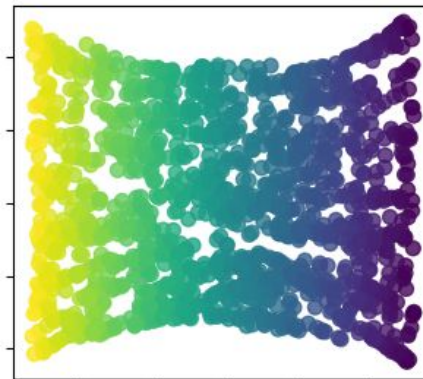
How do I look at all the data together?

# Many dimensions to few: Manifold learning, Ordination, Decomposition, Dimensionality reduction

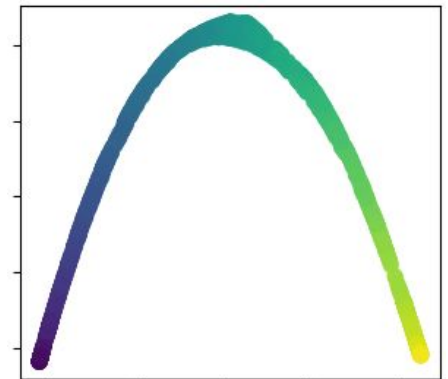
Original S-curve samples



Isomap Embedding



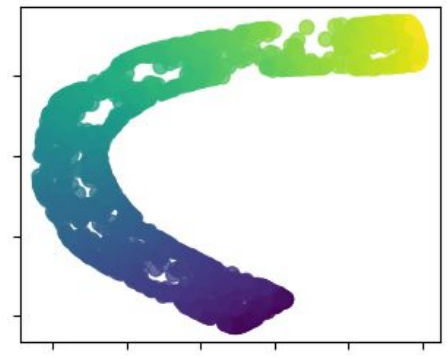
Spectral Embedding



Multidimensional scaling

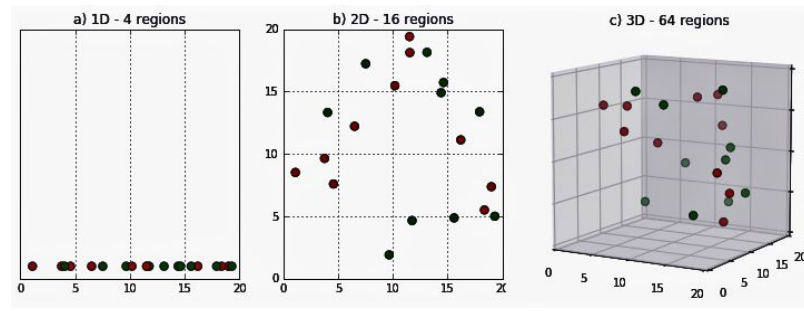


T-distributed Stochastic Neighbor Embedding

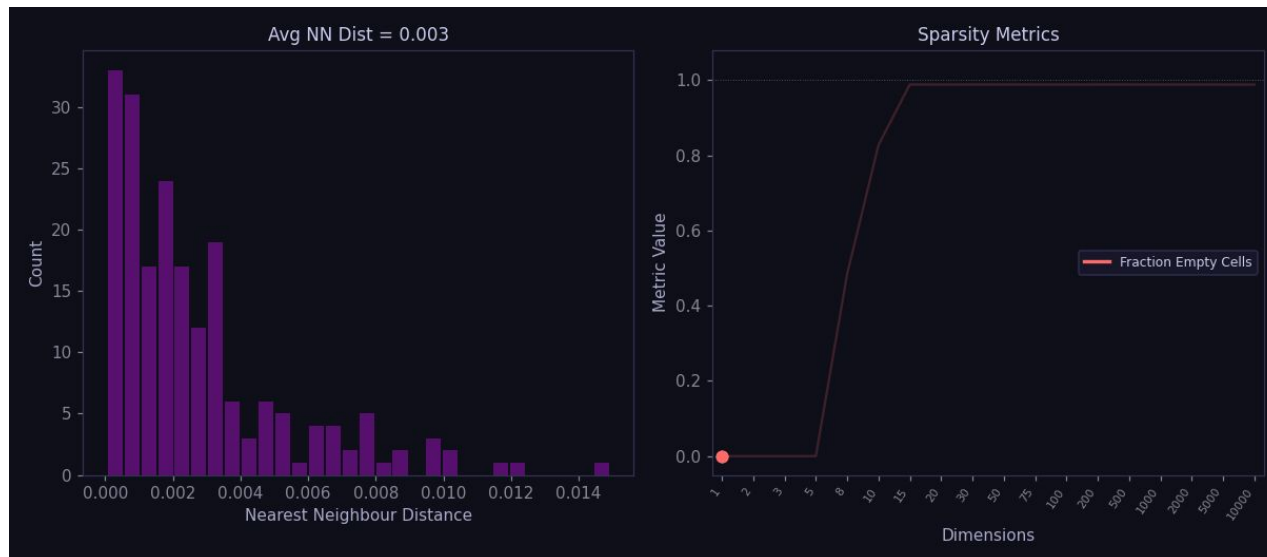
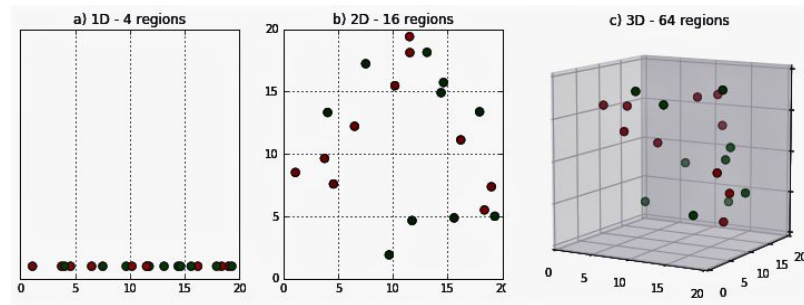


Why is this hard? High dimensions are counter-intuitive!

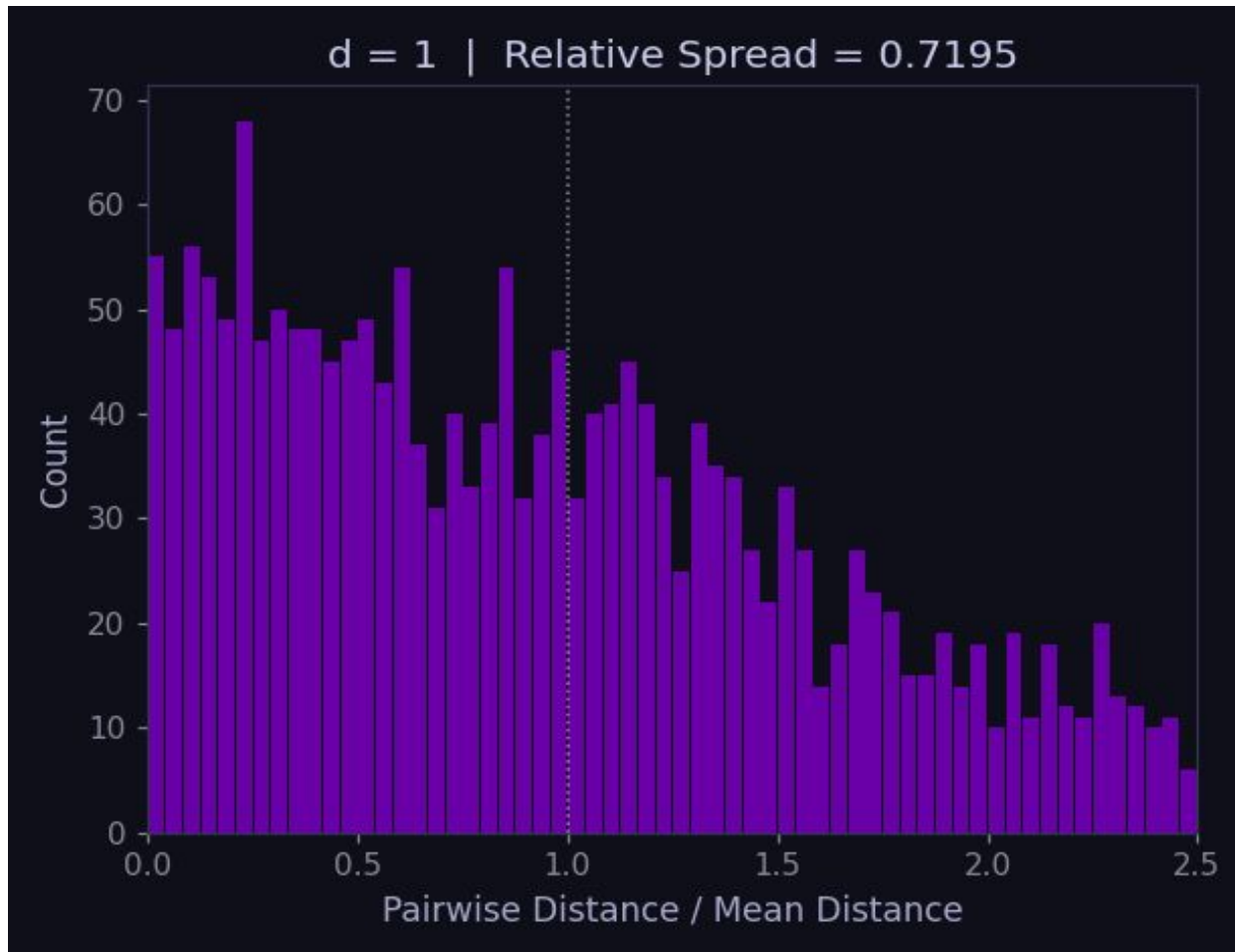
# Data becomes increasingly sparse in high dimensions



# Data becomes increasingly sparse in high dimensions



All points start to look equally far apart



No lower dimension representation will ever be perfect

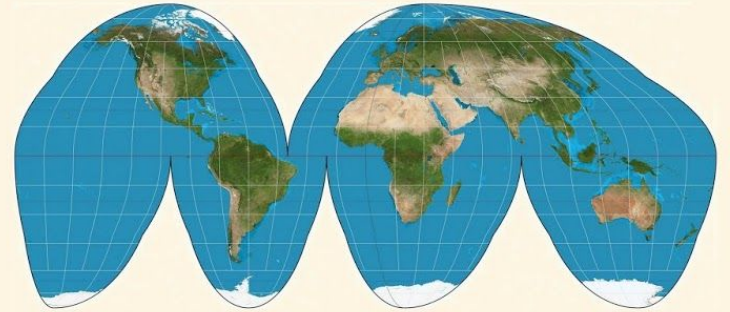
**MERCATOR**



**GALL-PETERS**



**GOODE-HOMOLOGINE**



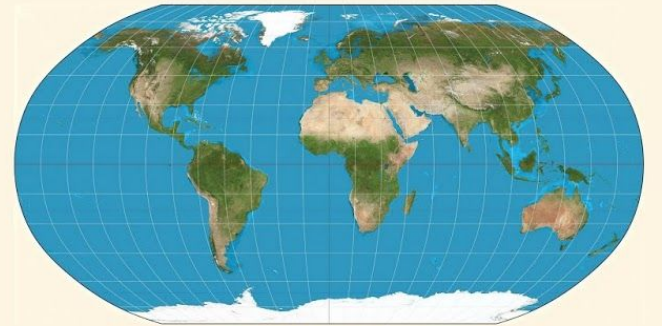
**WATERMELON**



**ALBERS**



**ROBINSON**



So, how can we do it?

# Principal Component Analysis - Simplest Method

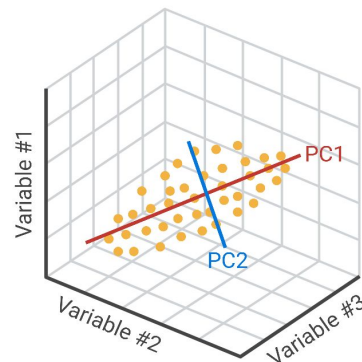
Reorient the data in the direction of maximal variance

1. Center the data
2. Calculate the covariance matrix
3. Perform eigendecomposition
4. Sort and select n principal components
5. Project the data onto the reduced space

$$\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$$

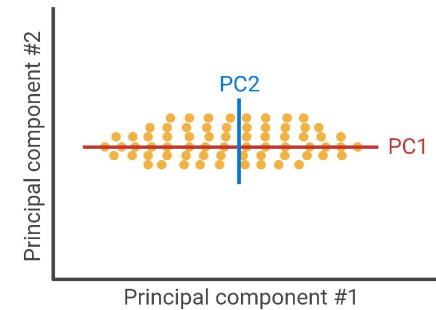
$\mathbf{Q}$ : Eigen vectors of  $\mathbf{A}$  ( $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$ )  
 $\mathbf{\Lambda}$ : Eigen values of  $\mathbf{A}$  ( $\lambda_1, \lambda_2, \lambda_3$ )  
 $\mathbf{Q}^{-1}$ : Eigen vectors of  $\mathbf{A}$  ( $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$ )

Original data  
(high-dimensions)



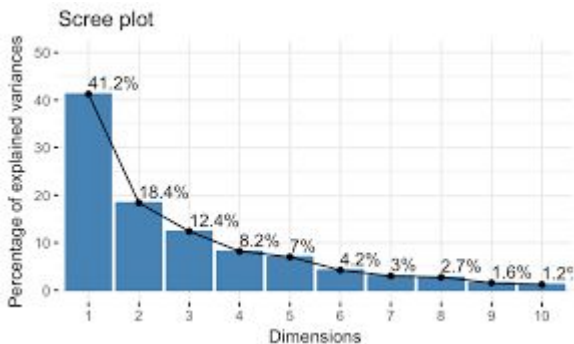
PCA dimensionality reduction

Lower-dimensional embedding



- Maximize variance along **PC1**
- Minimize residuals along **PC2**

How many components? Scree/elbow plot



# MultiDimensional Scaling (MDS): Distances

$$Stress_D(x_1, x_2, \dots, x_N) = \sqrt{\sum_{i \neq j=1, \dots, N} (d_{ij} - ||x_i - x_j||)^2}$$

The goal of the algorithm is to minimize the value of stress.

Where  $x_1, \dots, x_N$  are data points with their new set of coordinates in lower dimensional space.

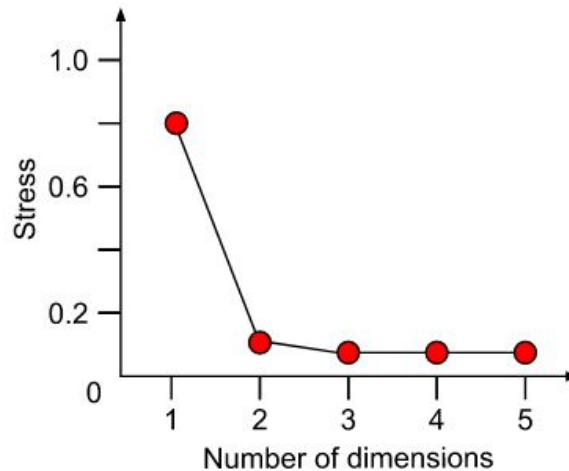
$d_{ij}$  is the actual distance we have calculated between the two corresponding data points in their original dimensional space.

$||x_i - x_j||$  is the distance between the two corresponding data points in their lower dimensional space.

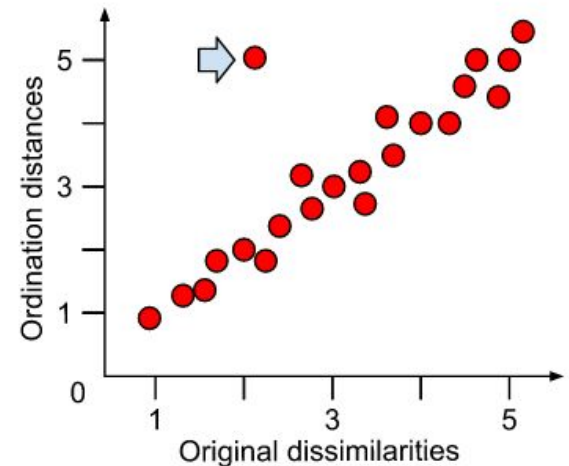
The closer the value of  $||x_i - x_j||$  is to  $d_{ij}$  the smaller will be the value of stress.

Non-Metric: Ranks

a

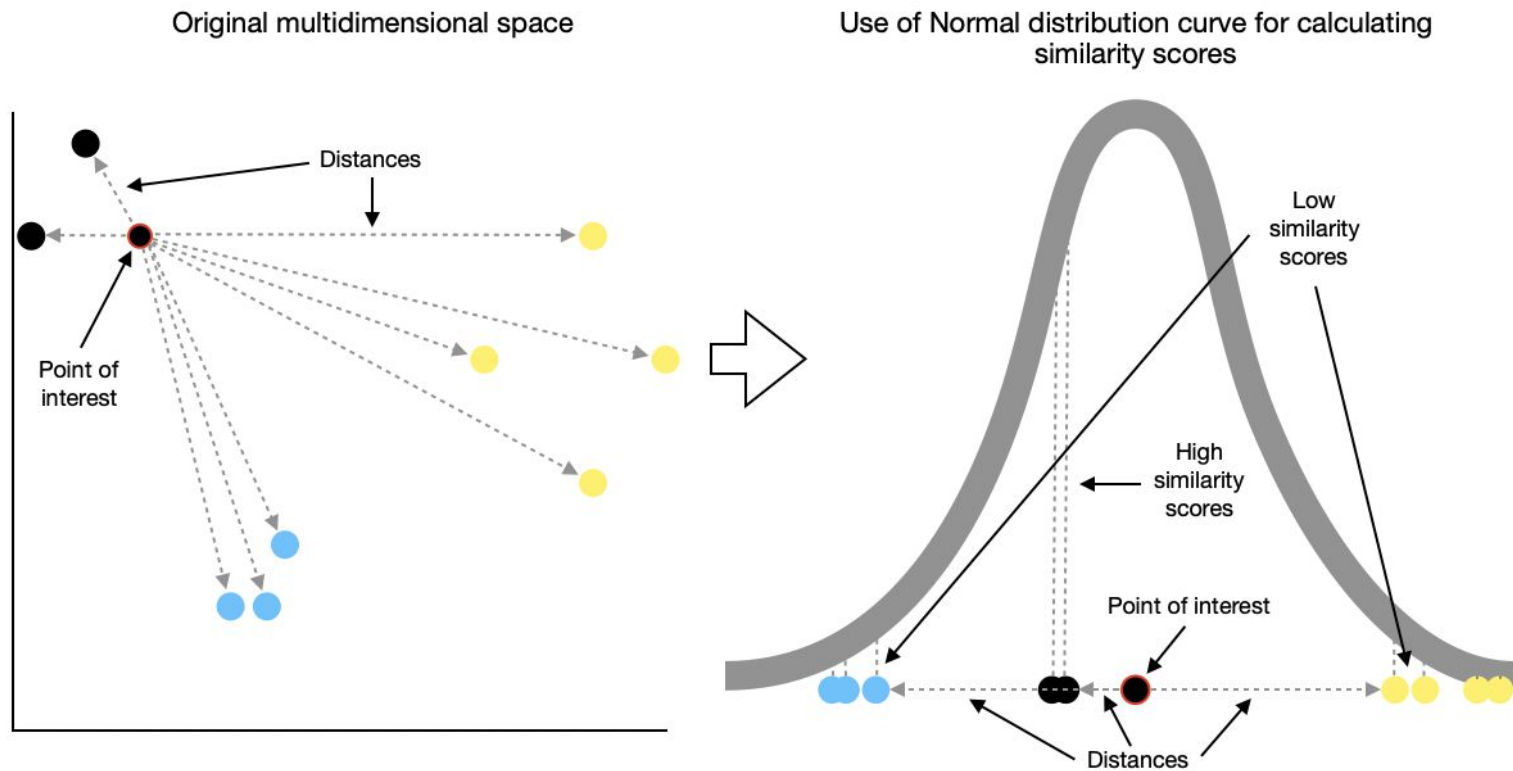


b



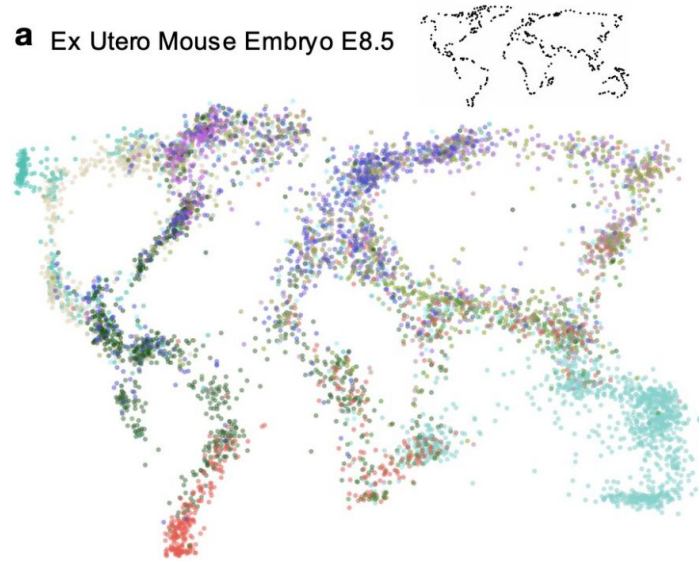
# t-SNE/UMAP: Probabilities

- Pairwise probability distribution in all dimensions
- Pairwise probability distribution in few dimensions
- Stochastic minimisation of KL divergence between distributions

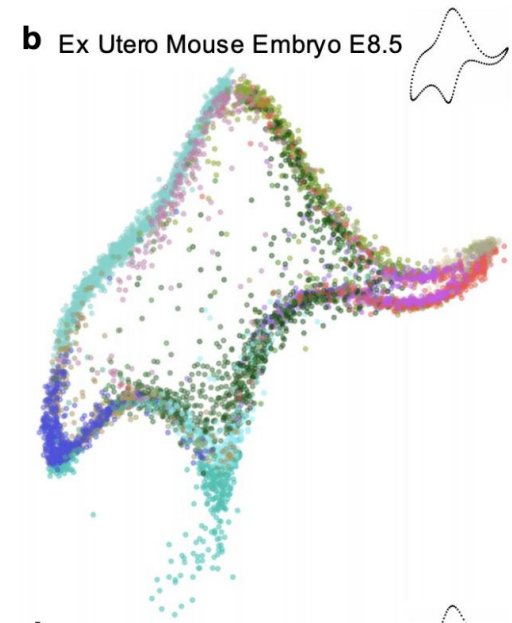


# Avoid over-interpreting single plots

- Sensitive to hyperparameters
- Beware analysing these non-linear projections
- Can contribute to confirmation bias



<https://www.biorxiv.org/content/10.1101/2021.06.23.457690v3>

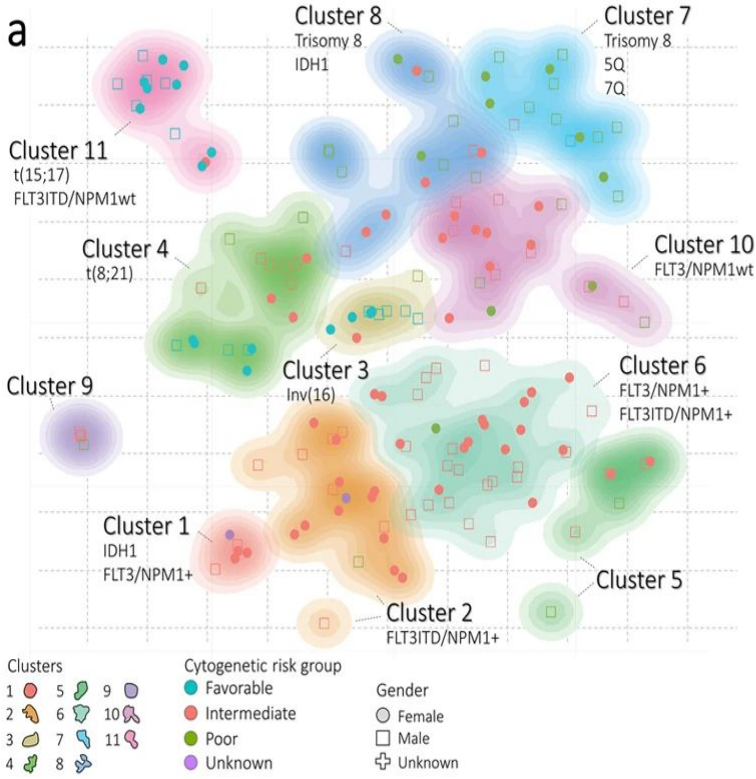
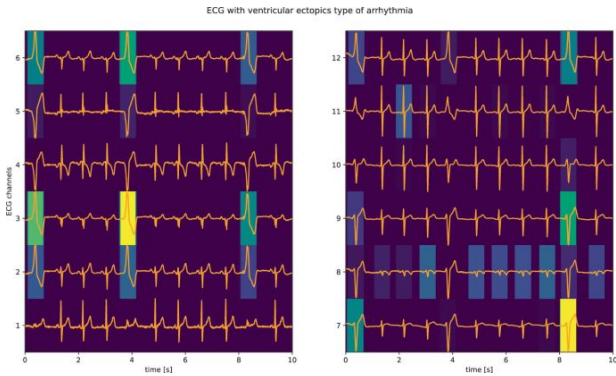
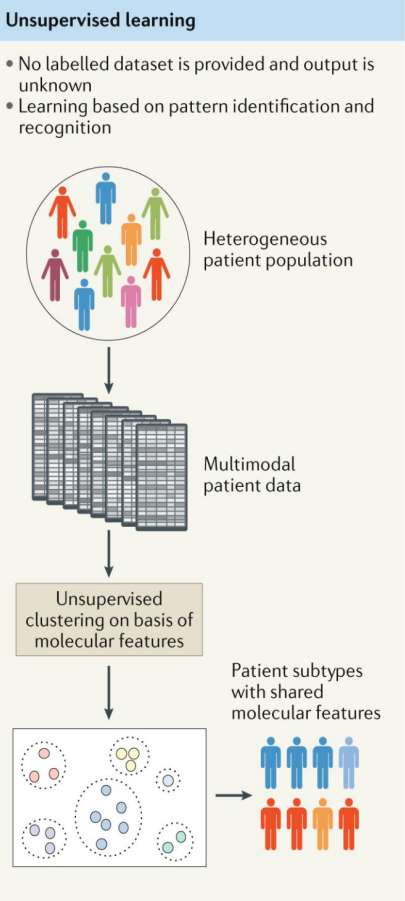


*"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk." - Von Neumann*

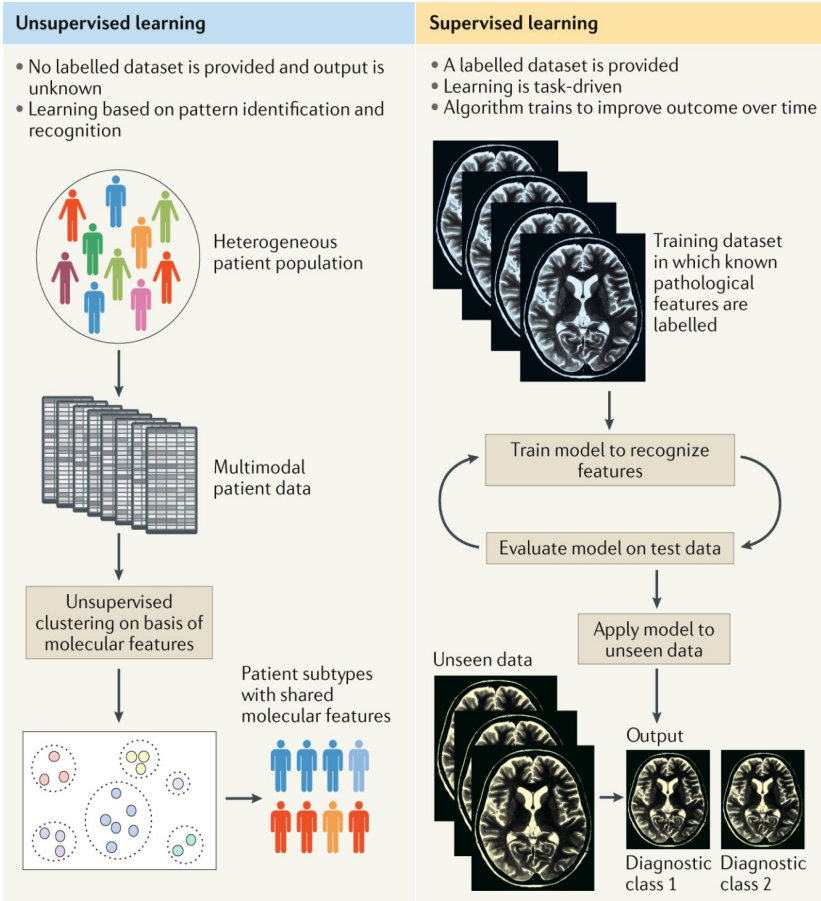
# Finding structure in data: unsupervised learning

## Anomaly Detection - ECG

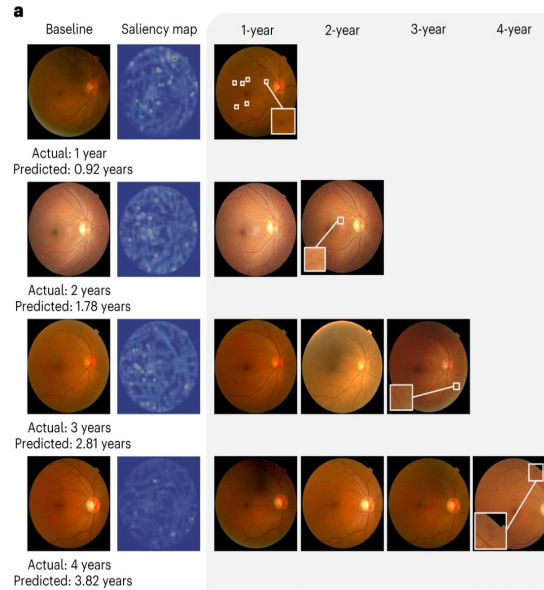
## Clustering - Acute Myeloid Leukemia Subtypes



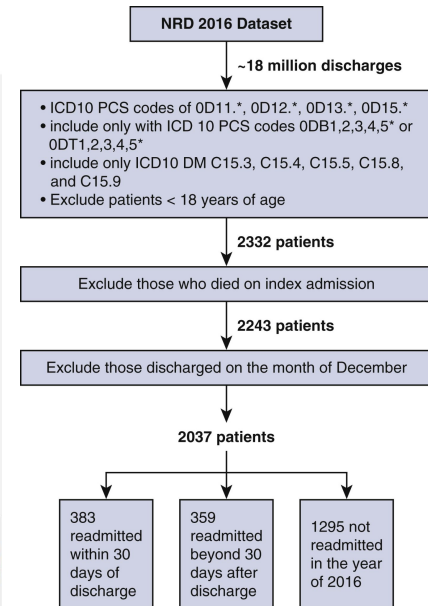
# Predicting outcomes from data: supervised learning



## Prediction of Diabetic Retinopathy Progression



## Prediction of 30-Day Readmission



Dai, L., Sheng, B., Chen, T. et al. A deep learning system for predicting time to progression of diabetic retinopathy. *Nat Med* 30, 584–594 (2024). <https://doi.org/10.1038/s41591-023-02702-z>

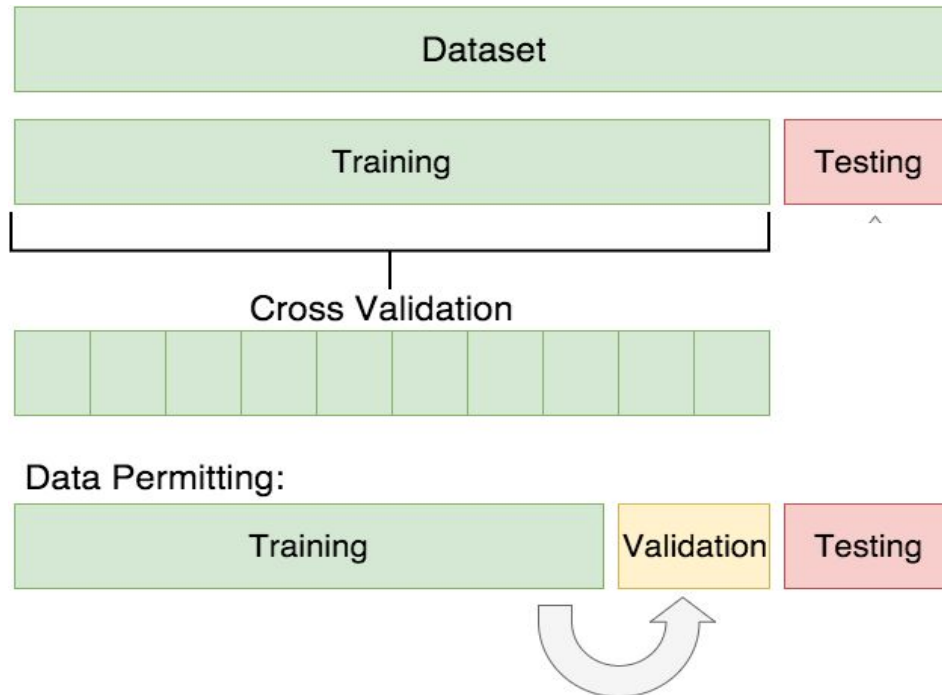
Bolourani, Siavash, et al. "Using machine learning to predict early readmission following esophagectomy." *The Journal of Thoracic and Cardiovascular Surgery* 161.6 (2021): 1926-1939.

# Supervised Learning: Classification and Regression

| Patient ID | Age | Sex | BMI  | Systolic BP (mmHg) | Diastolic BP (mmHg) | Fasting Glucose (mmol/L) | HbA1c (%) | Cholesterol (mmol/L) | Smoker | Family History Diabetes | Diagnosis (Classification) | 10-yr CVD Risk (Regression) |
|------------|-----|-----|------|--------------------|---------------------|--------------------------|-----------|----------------------|--------|-------------------------|----------------------------|-----------------------------|
| P001       | 54  | M   | 28.3 | 138                | 88                  | 6.2                      | 6.8       | 5.1                  | Yes    | Yes                     | Diabetic                   | 18.4%                       |
| P002       | 41  | F   | 22.1 | 118                | 75                  | 4.9                      | 5.2       | 4.7                  | No     | No                      | Healthy                    | 4.1%                        |
| P003       | 67  | M   | 31.7 | 155                | 95                  | 7.8                      | 7.4       | 6.3                  | Yes    | Yes                     | Diabetic                   | 31.2%                       |
| P004       | 35  | F   | 25.6 | 122                | 80                  | 5.3                      | 5.5       | 4.2                  | No     | Yes                     | Healthy                    | 5.8%                        |
| P005       | 58  | M   | 29.9 | 145                | 91                  | 6.9                      | 6.5       | 5.8                  | No     | No                      | Diabetic                   | 14.7%                       |
| P006       | 72  | F   | 27.4 | 160                | 97                  | 8.4                      | 8.1       | 6.9                  | Yes    | Yes                     | Diabetic                   | 38.5%                       |
| P007       | 29  | M   | 23.8 | 115                | 73                  | 4.7                      | 5         | 4                    | No     | No                      | Healthy                    | 2.3%                        |
| P008       | 63  | F   | 33.2 | 148                | 93                  | 7.1                      | 7         | 5.5                  | No     | Yes                     | Diabetic                   | 25.6%                       |

# Supervised Learning: Classification and Regression

| Patient ID | Age | Sex | BMI  | Systolic BP (mmHg) | Diastolic BP (mmHg) | Fasting Glucose (mmol/L) | HbA1c (%) | Cholesterol (mmol/L) | Smoker | Family History Diabetes | Diagnosis (Classification) | 10-yr CVD Risk (Regression) |
|------------|-----|-----|------|--------------------|---------------------|--------------------------|-----------|----------------------|--------|-------------------------|----------------------------|-----------------------------|
| P001       | 54  | M   | 28.3 | 138                | 88                  | 6.2                      | 6.8       | 5.1                  | Yes    | Yes                     | Diabetic                   | 18.4%                       |
| P002       | 41  | F   | 22.1 | 118                | 75                  | 4.9                      | 5.2       | 4.7                  | No     | No                      | Healthy                    | 4.1%                        |
| P003       | 67  | M   | 31.7 | 155                | 95                  | 7.8                      | 7.4       | 6.3                  | Yes    | Yes                     | Diabetic                   | 31.2%                       |
| P004       | 33  | F   | 23.8 | 122                | 80                  | 5.3                      | 5.9       | 4.2                  | No     | Yes                     | Healthy                    | 9.8%                        |
| P005       | 58  | M   | 29.9 | 145                | 91                  | 6.9                      | 6.5       | 5.8                  | No     | No                      | Diabetic                   | 14.7%                       |
| P006       | 72  | F   | 27.4 | 160                | 97                  | 8.4                      | 8.1       | 6.9                  | Yes    | Yes                     | Diabetic                   | 28.5%                       |
| P007       | 29  | M   | 23.8 | 115                | 73                  | 4.7                      | 5         | 4                    | No     | No                      | Healthy                    | 2.3%                        |
| P008       | 63  | F   | 33.2 | 148                | 93                  | 7.1                      | 7         | 5.5                  | No     | Yes                     | Diabetic                   | 25.8%                       |



What does it mean to fit a supervised model?

# Decision boundaries and loss functions

Find  $\beta$  so our function maps  $X$  to  $y$

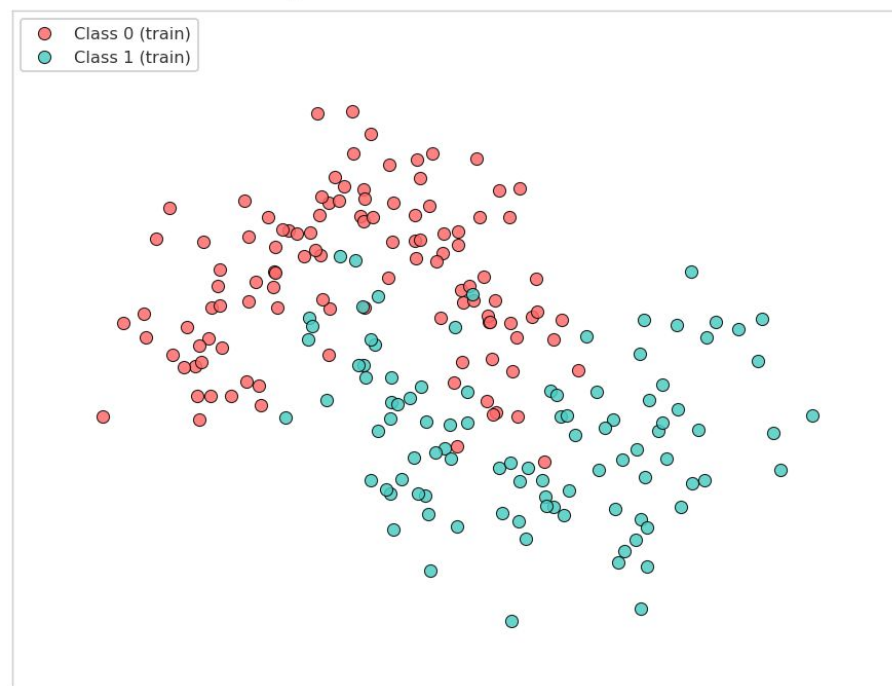
$$\hat{y} = f_{\beta}(x)$$

In other words: find  $\beta$  that defines boundary between labels

Minimising some loss function like binary cross-entropy (log-loss)

$$\mathcal{L}(\beta) = -\frac{1}{N} \sum_{i=1}^N [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

Binary Classification Dataset



train set

# Decision boundaries and loss functions

Find  $\beta$  so our function maps  $X$  to  $y$

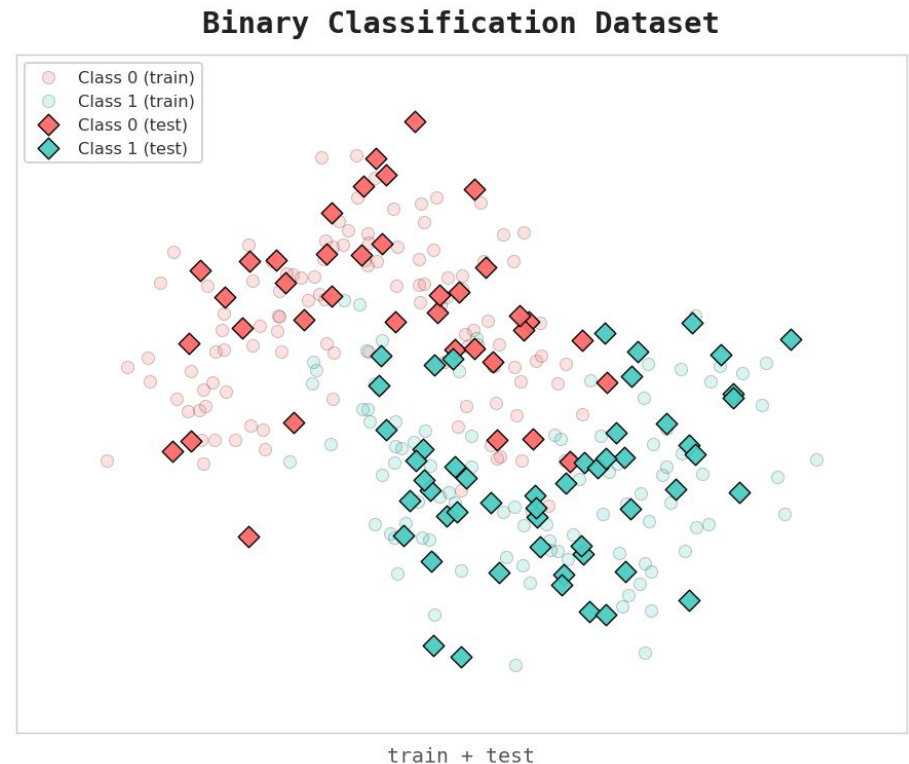
$$\hat{y} = f_{\beta}(x)$$

In other words: find  $\beta$  that defines boundary between labels

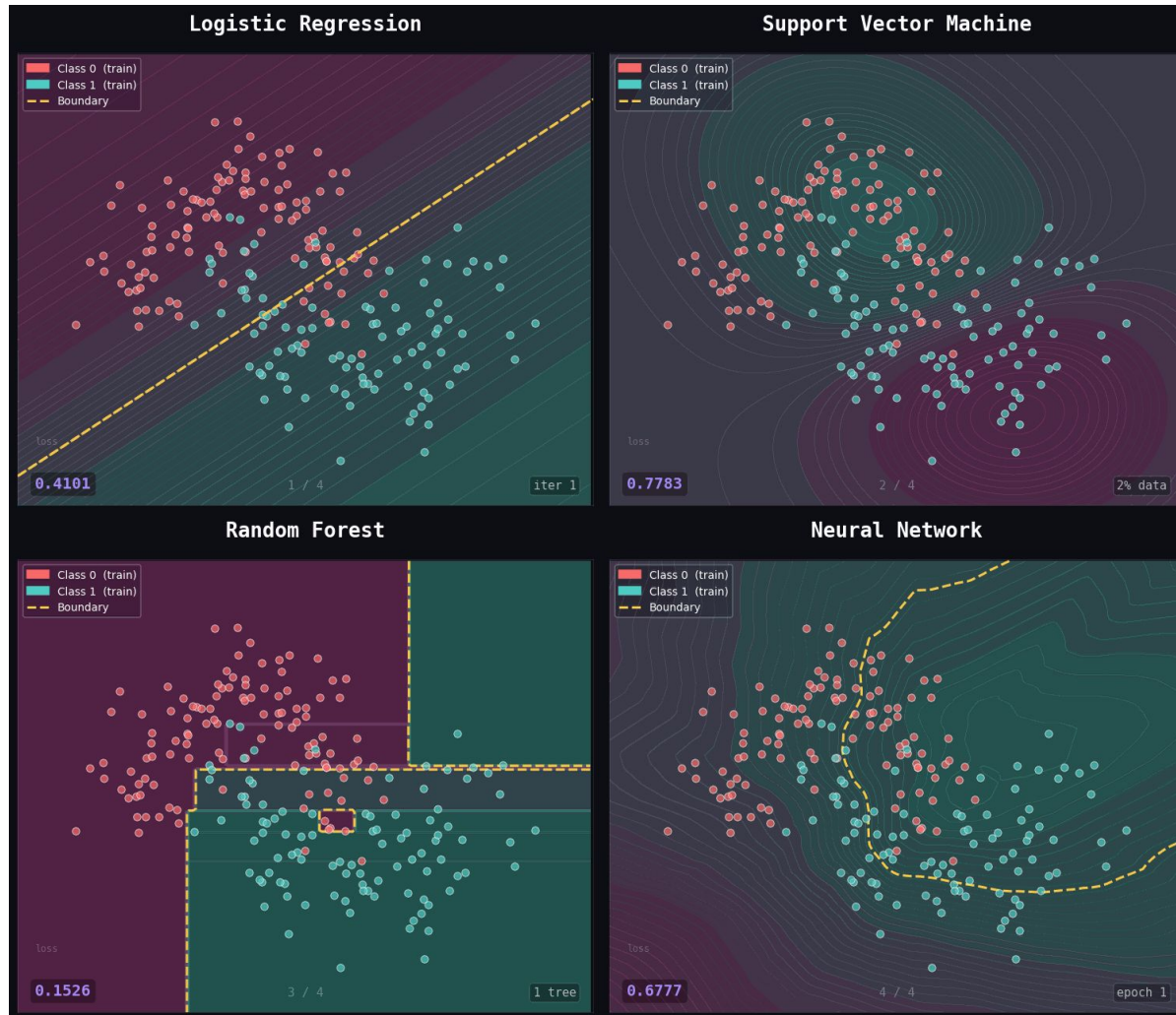
Minimising some loss function like binary cross-entropy (log-loss)

$$\mathcal{L}(\beta) = -\frac{1}{N} \sum_{i=1}^N [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

Ideally in a way that generalises to new data!



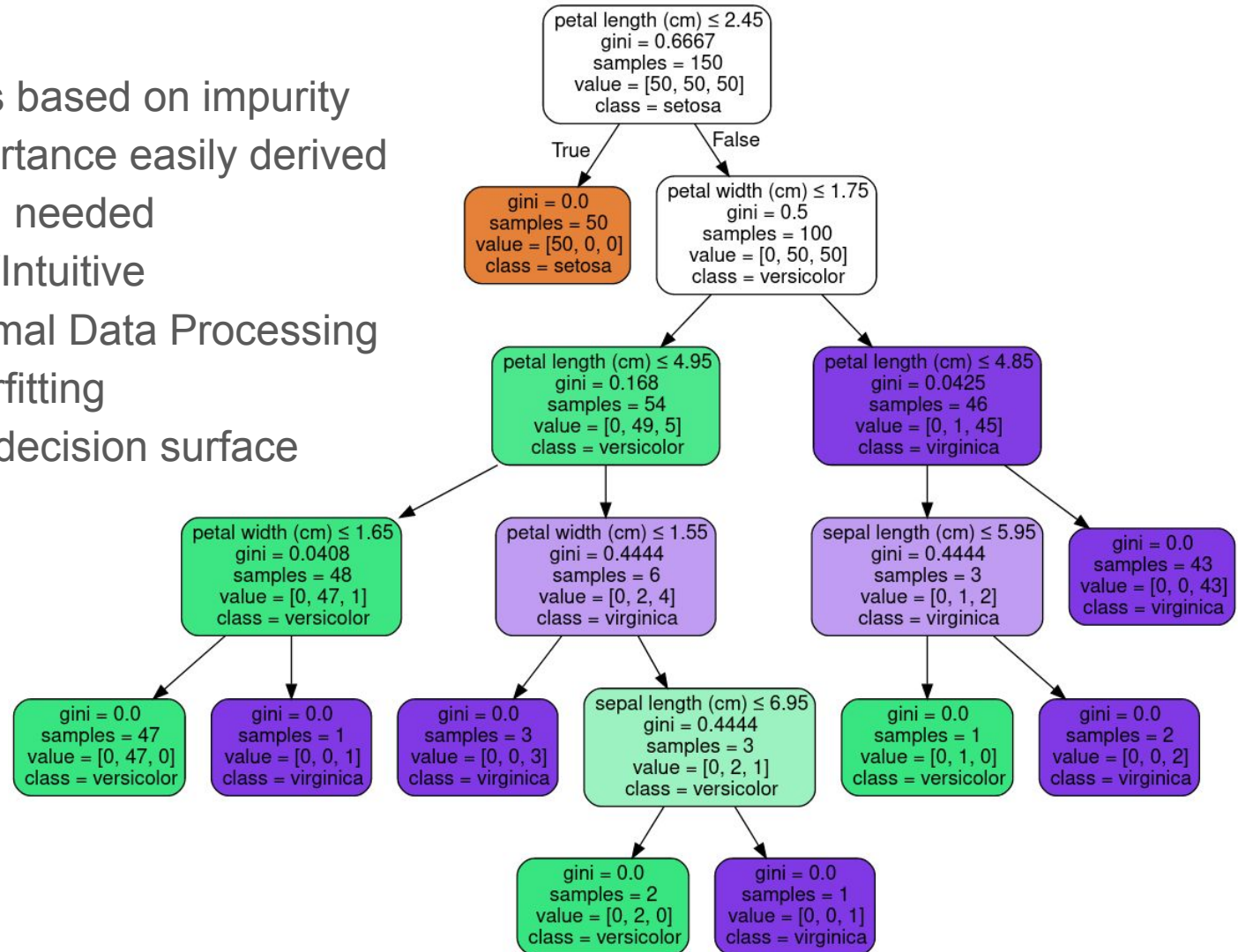
# Decision boundaries and loss functions



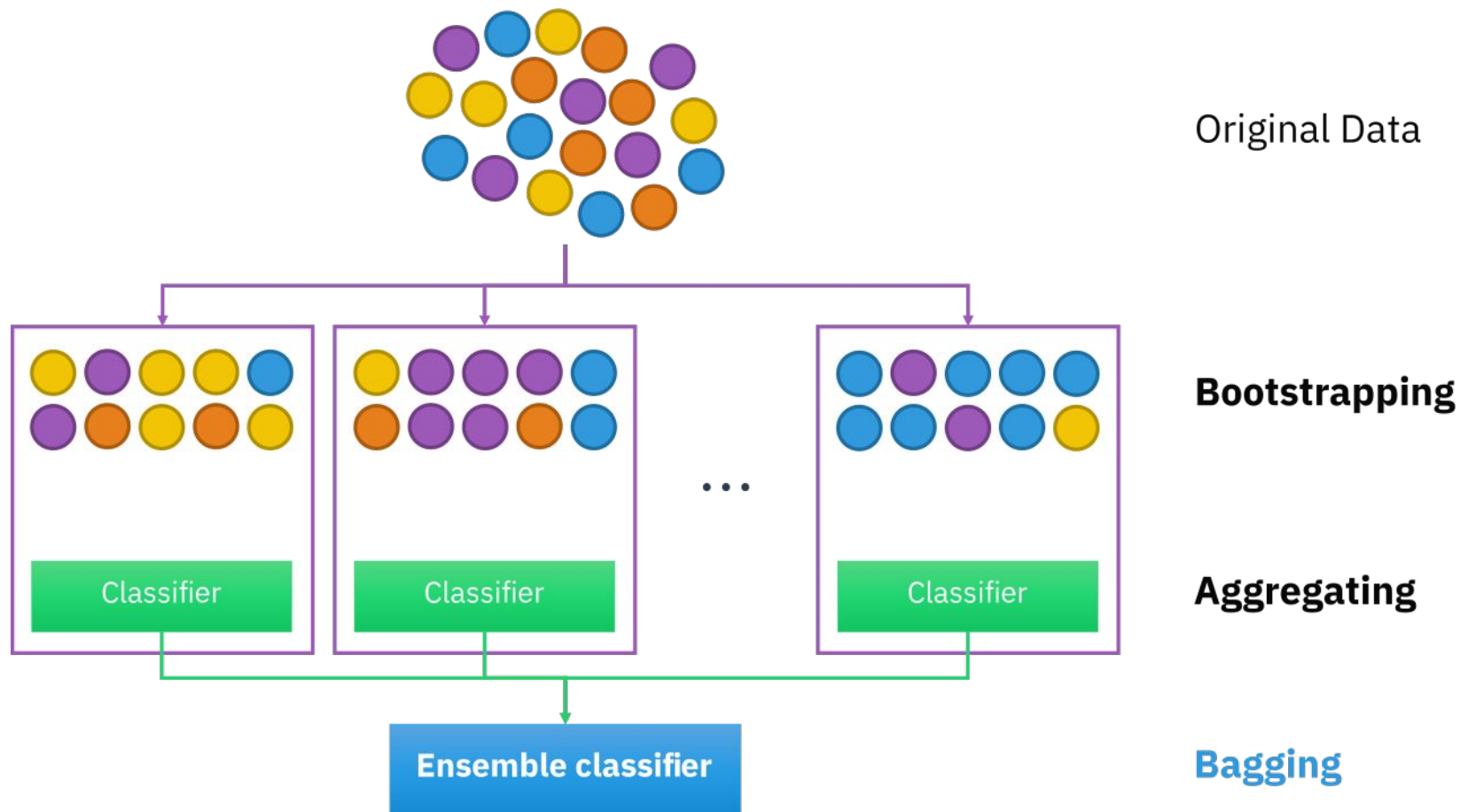
Can't get a well-performing model? Just  
combine lots of weak ones!

# Decision Trees

- Dataset splits based on impurity
- Feature importance easily derived
- Pruning often needed
- Interpretable/Intuitive
- Require Minimal Data Processing
- Prone to overfitting
- Non-smooth decision surface

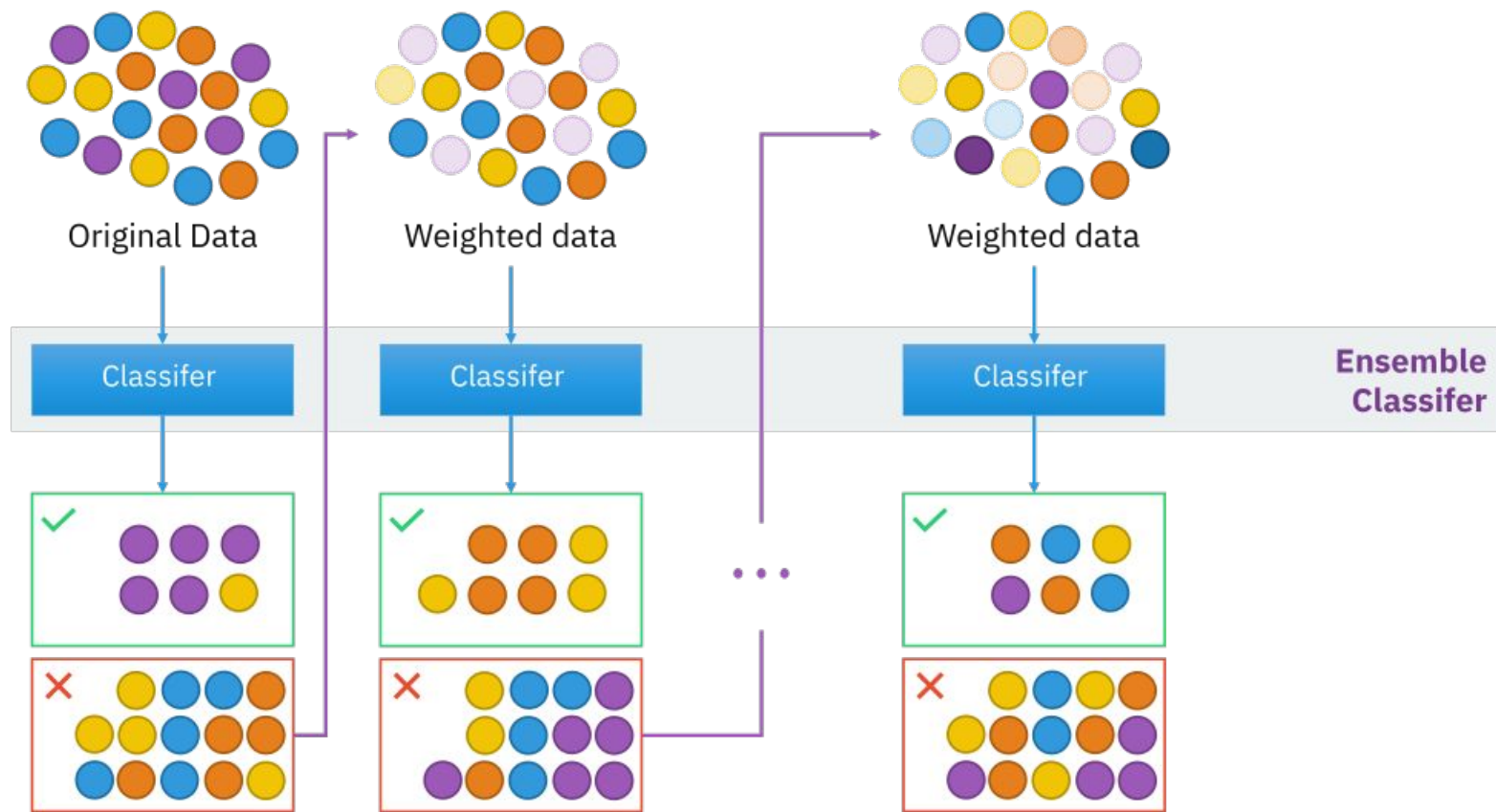


# Many Decision Trees: Bagging



Random Forest: Bagging + Random Subset Per Split  
Feature Importance: Average impurity decrease

# Boosting: AdaBoost, Gradient Boosting, XGBoost



AdaBoost, Gradient Boosting, XGBoost

# Decision Trees methods regularly outperform deep learning on tabular data

Tree-based methods deal well with common features of tabular data (even compared to well-tuned neural networks):

- Heterogeneous data
- Ignoring uninformative data
- Non-smooth decision boundaries
- Moderate size & dimensionality
- Skewed or heavy-tailed feature distributions and other forms of dataset
- Rotational invariance (column/row order is not informative)

**But:** difference is often negligible (except in computational efficiency!)

---

**Why do tree-based models still outperform deep learning on typical tabular data?**

Léo Grinsztajn  
Soda, Inria Saclay  
leo.grinsztajn@inria.fr

Edouard Oyallon  
MLIA, Sorbonne University

Gaël Varoquaux  
Soda, Inria Saclay

---

**When Do Neural Nets Outperform Boosted Trees on Tabular Data?**

Duncan McElfresh<sup>\*1,2</sup>, Sujay Khandagale<sup>3</sup>, Jonathan Valverde<sup>4</sup>, Vishak Prasad C<sup>5</sup>,  
Ganesh Ramakrishnan<sup>5</sup>, Micah Goldblum<sup>6</sup>, Colin White<sup>1,7</sup>

<sup>1</sup> Abacus.AI, <sup>2</sup> Stanford, <sup>3</sup> Pinterest, <sup>4</sup> University of Maryland,

<sup>5</sup> IIT Bombay, <sup>6</sup> New York University, <sup>7</sup> Caltech

# Overview

- Medical databases are usually relational and are defined by their origin, primary record type, scope, and sampling strategy
- Standardisation is important and ontologies support that in medical databases
- Survey weights are key to compensate for complex sampling
- There is a continuum of approaches to retain data privacy (and data ownership is a complex issue)
- Individual and joint distributions are key EDA tools
- Dimensionality reduction (PCA, MDS, t-SNE) is very useful but can be challenging/misleading
- Start with simple classifiers e.g., logistic regression/decision tree
- Combine weak classifiers via bagging (bootstrapping data: Random Forest special form) or boosting (sequential training model on errors: AdaBoost/XGBoost) to improve performance.
- XGBoost gold-standard but requires more tuning than AdaBoost