

Lecture 1: Medical Databases

CSCI6XXX/CHE6XXX/CSCI4XXX
(CSCI6093)

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

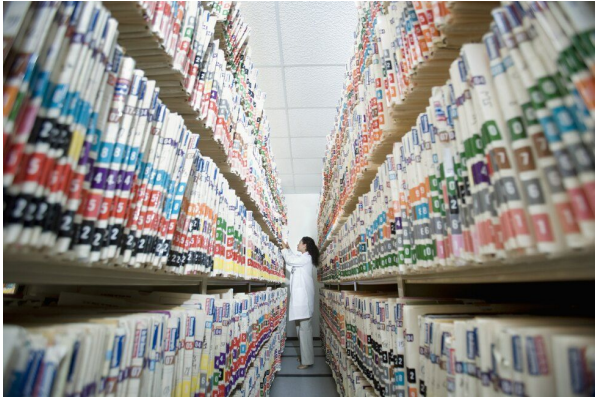
What is a database?

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 'PART A - PRESENT HEALTH HISTORY' and includes 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

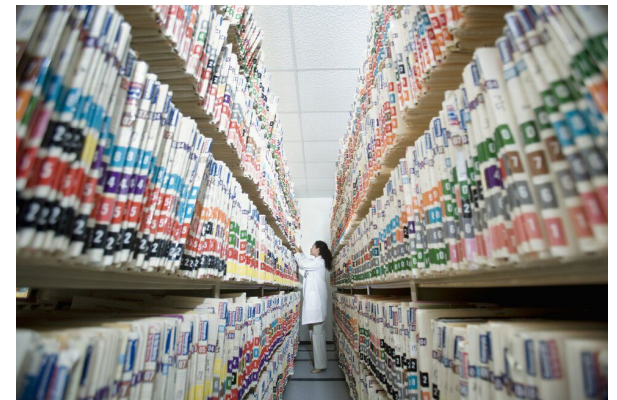
Examples include:

- Medical Charts
- Phone Book
- Dictionaries
- Spreadsheet

Ordering:

- Index
- Defined fields
- Standardisation

The image shows 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 'ANDRUS/CLINI-REC HEALTH HISTORY QUESTIONNAIRE' with sections for current medical problems, medications, allergies, and general health. The forms are marked 'CONFIDENTIAL' on the sides.



Databases (broadly) are ordered collections of data

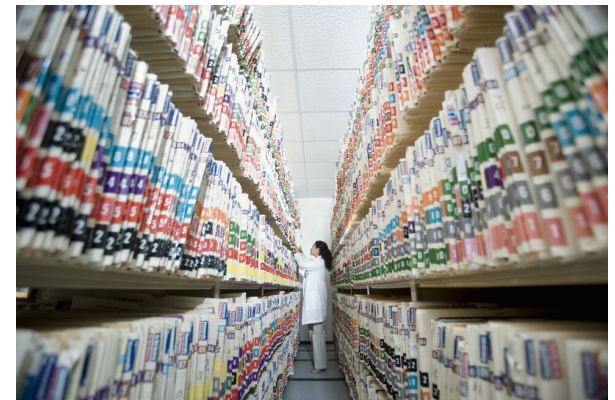
Examples include:

- 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.

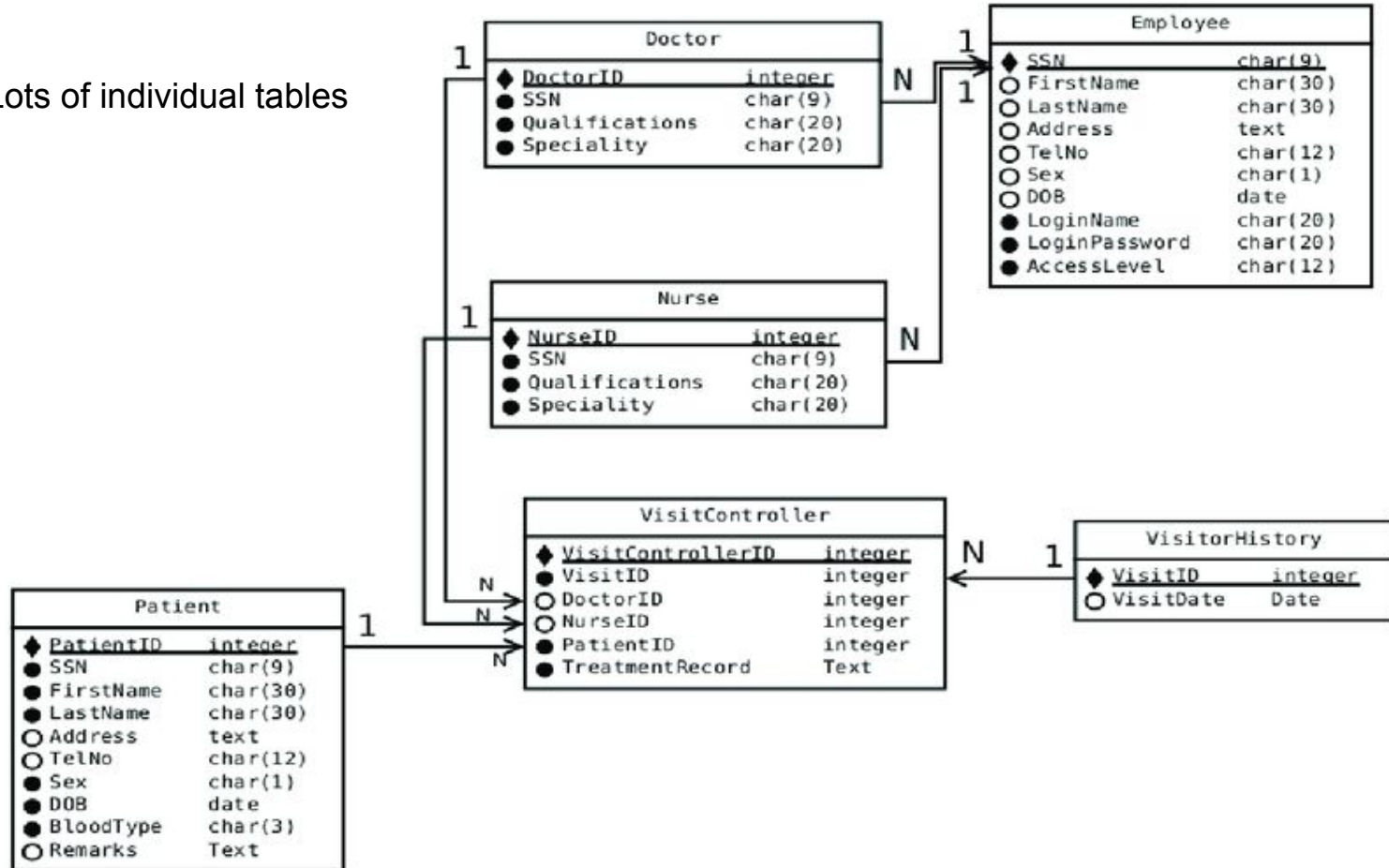


Organisation make some tasks easier/harder:

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

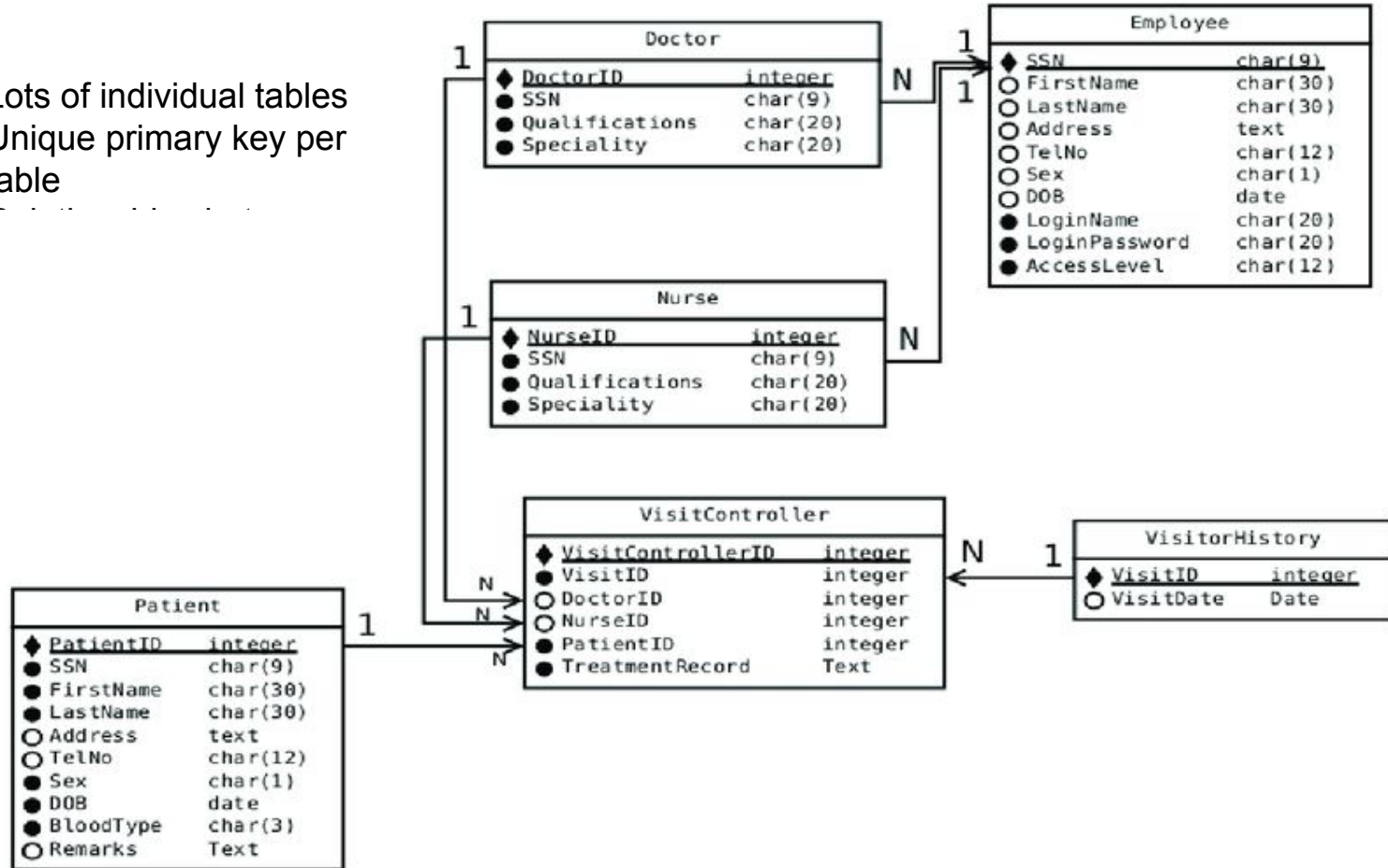
Most Common Type: Relational Databases

- Lots of individual tables



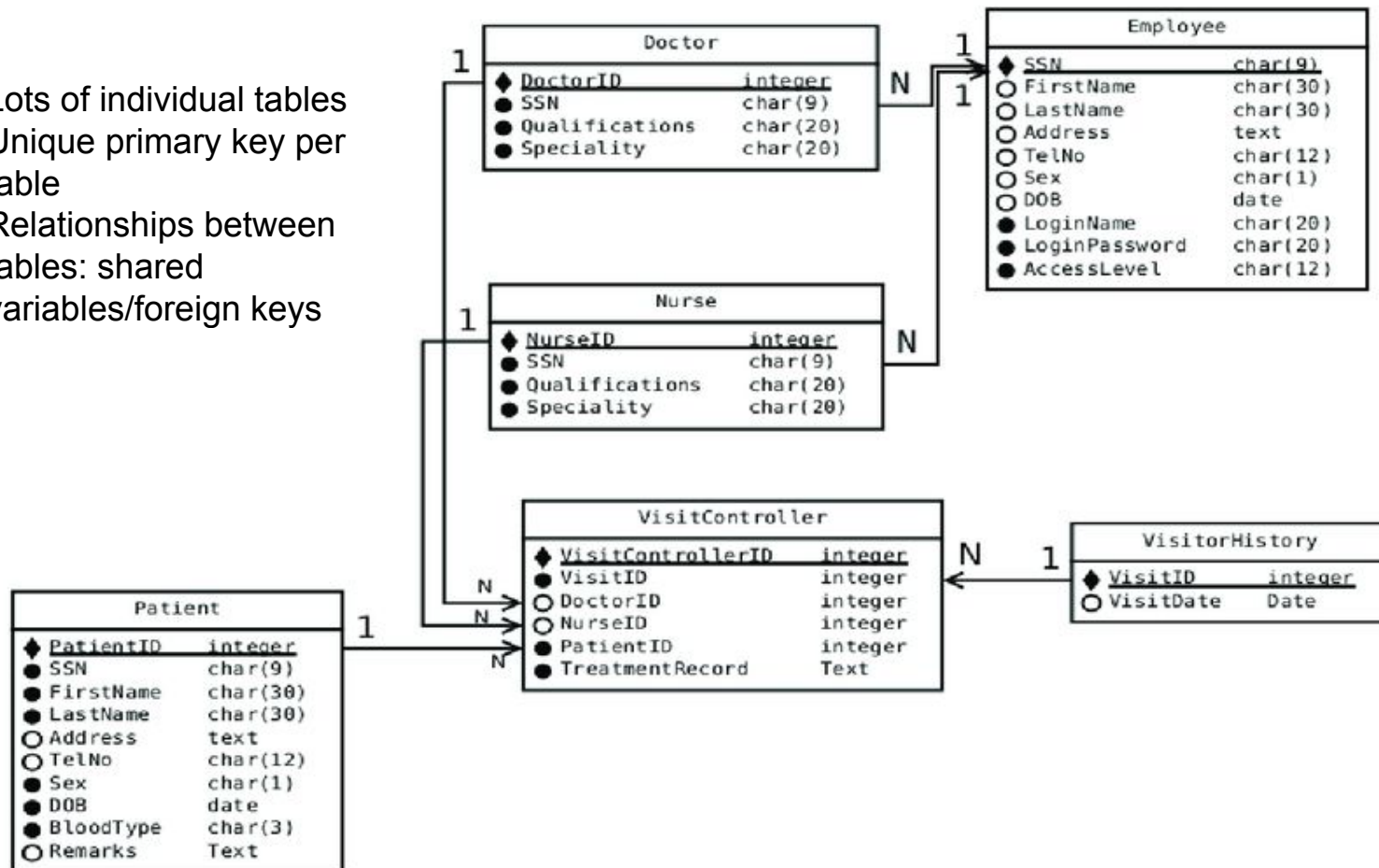
Most Common Type: Relational Databases

- Lots of individual tables
- Unique primary key per table



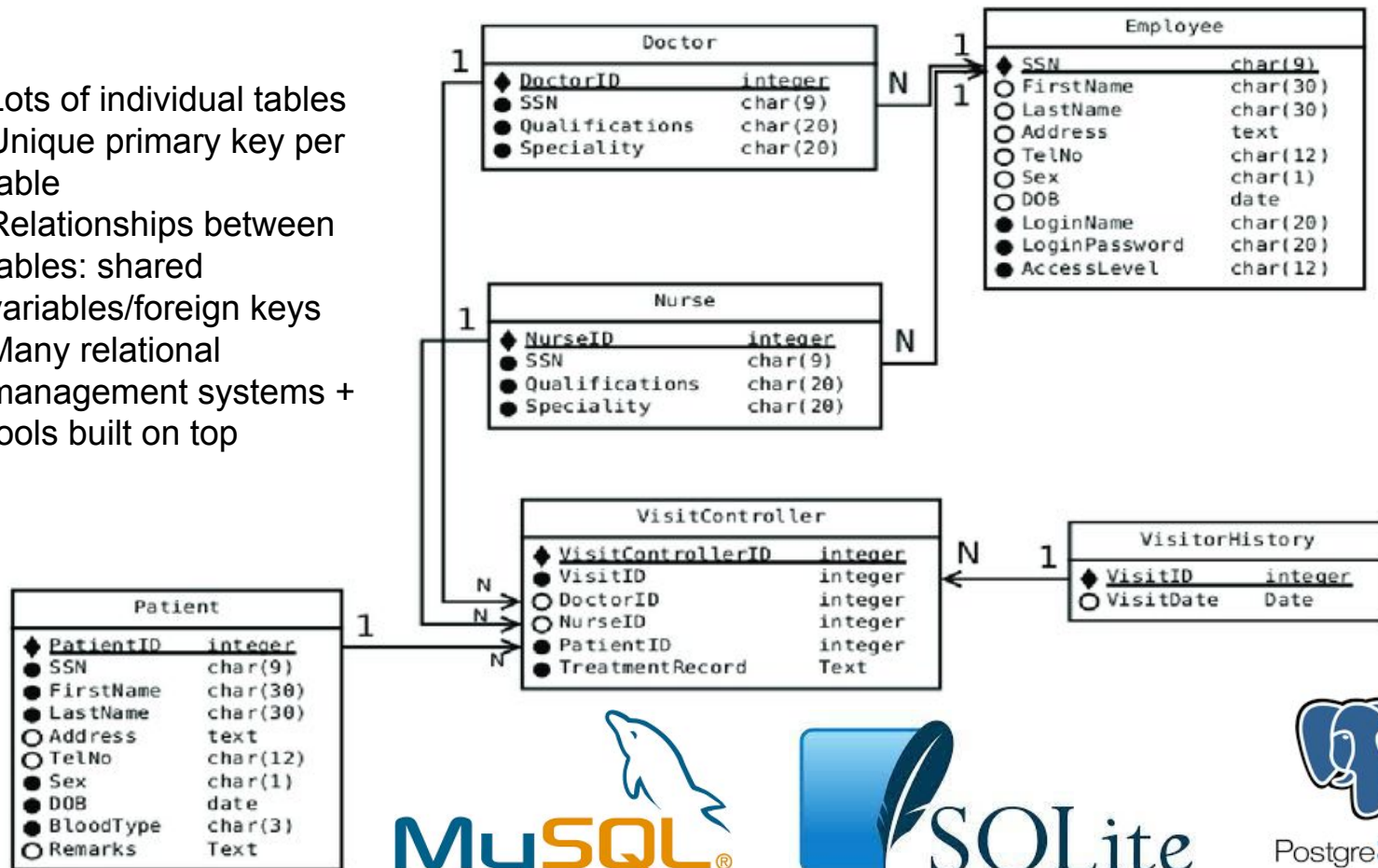
Most Common Type: Relational Databases

- 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



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)

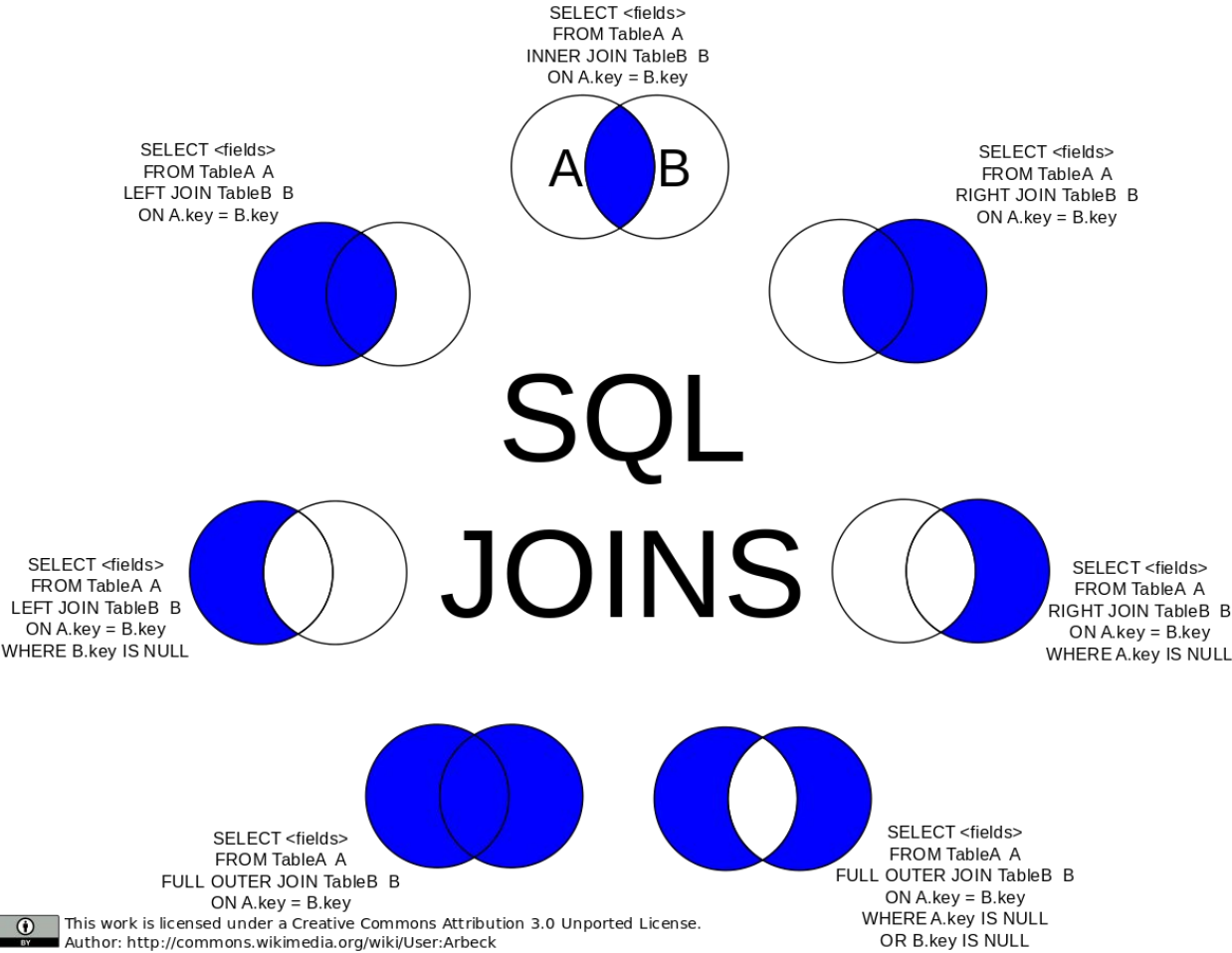
- Non-procedural Language
- Standardised/Powerful/Flexible
- Basis of many data tools
- Well-supported by dbplyr

```
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



Are all databases relational?

Non-Relational Databases

- Less common than relational in medicine

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



Column based



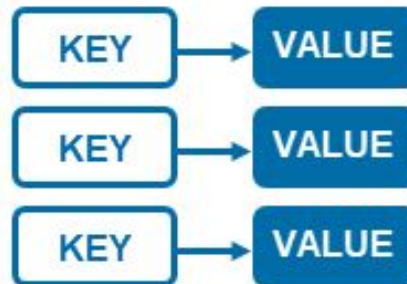
Non-Relational Databases

- Less common than relational in medicine

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



Column based



Key-value



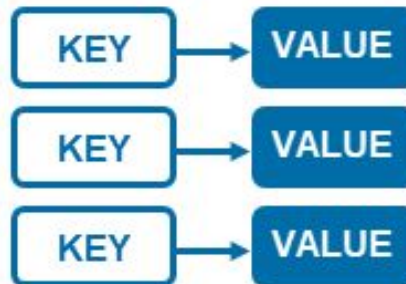
Non-Relational Databases

- Less common than relational in medicine

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



Column based



Key-value



Graph



Non-Relational Databases

- Less common than relational in medicine
- Querying can be... very easy or very complicated

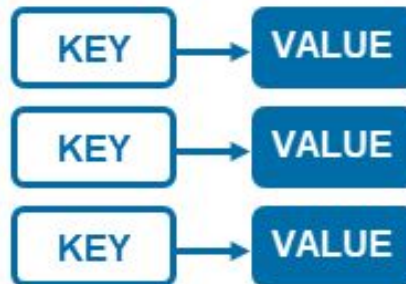
Find me the homepage of anyone known by Tim Berners-Lee.

```
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX card: <http://www.w3.org/People/Berners-Lee/card#>
SELECT ?homepage
FROM <http://www.w3.org/People/Berners-Lee/card>
WHERE {
    card:i foaf:knows ?known .
    ?known foaf:homepage ?homepage .
}
```

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



Column based



Key-value



Graph



Document



Non-Relational Databases

- Less common than relational in medicine
- Querying can be... very easy or very complicated
- Mostly for very large datasets:
 - User data / security audit data
 - Medical image data

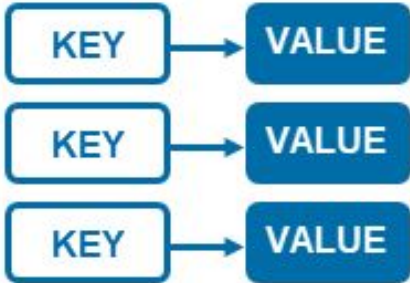
Find me the homepage of anyone known by Tim Berners-Lee.

```
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX card: <http://www.w3.org/People/Berners-Lee/card#>
SELECT ?homepage
FROM <http://www.w3.org/People/Berners-Lee/card>
WHERE {
    card:i foaf:knows ?known .
    ?known foaf:homepage ?homepage .
}
```

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



Column based



Key-value



Graph



Document



Non-Relational Databases

- Less common than relational in medicine
- Querying can be... very easy or very complicated
- Mostly for very large datasets:
 - User data / security audit data
 - Medical image data
- Or unusual data structures:
 - Contact tracing
 - Ontologies

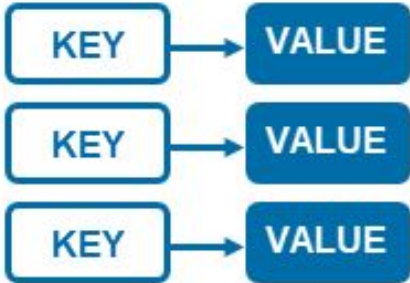
Find me the homepage of anyone known by Tim Berners-Lee.

```
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX card: <http://www.w3.org/People/Berners-Lee/card#>
SELECT ?homepage
FROM <http://www.w3.org/People/Berners-Lee/card>
WHERE {
    card:i foaf:knows ?known .
    ?known foaf:homepage ?homepage .
}
```

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



Column based



Key-value



Graph



Document



Non-Relational Databases

- Less common than relational in medicine
- Querying can be... very easy or very complicated
- Mostly for very large datasets:
 - User data / security audit data
 - Medical image data
- Or unusual data structures:
 - Contact tracing
 - Ontologies
- Or both:
 - Social media data

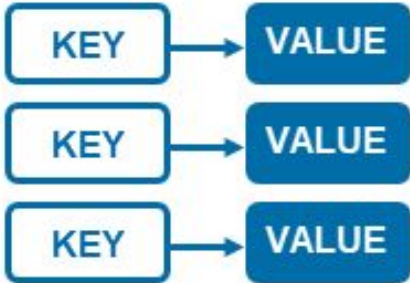
Find me the homepage of anyone known by Tim Berners-Lee.

```
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX card: <http://www.w3.org/People/Berners-Lee/card#>
SELECT ?homepage
FROM <http://www.w3.org/People/Berners-Lee/card>
WHERE {
    card:i foaf:knows ?known .
    ?known foaf:homepage ?homepage .
}
```

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



Column based



Key-value



Graph



Document



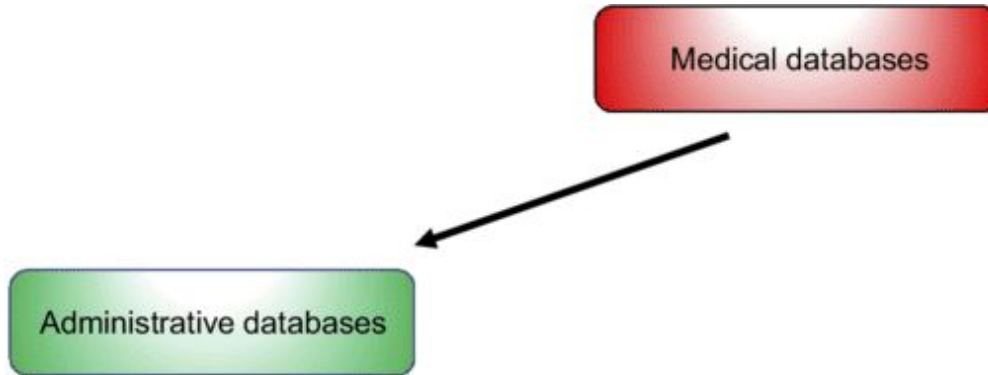
What are medical databases?

Many types of database



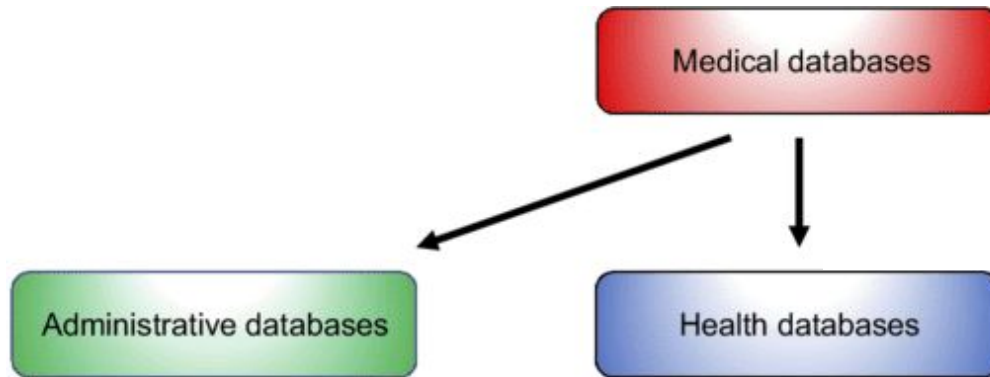
■ All types of registries and databases that contain health-related data

Many types of database



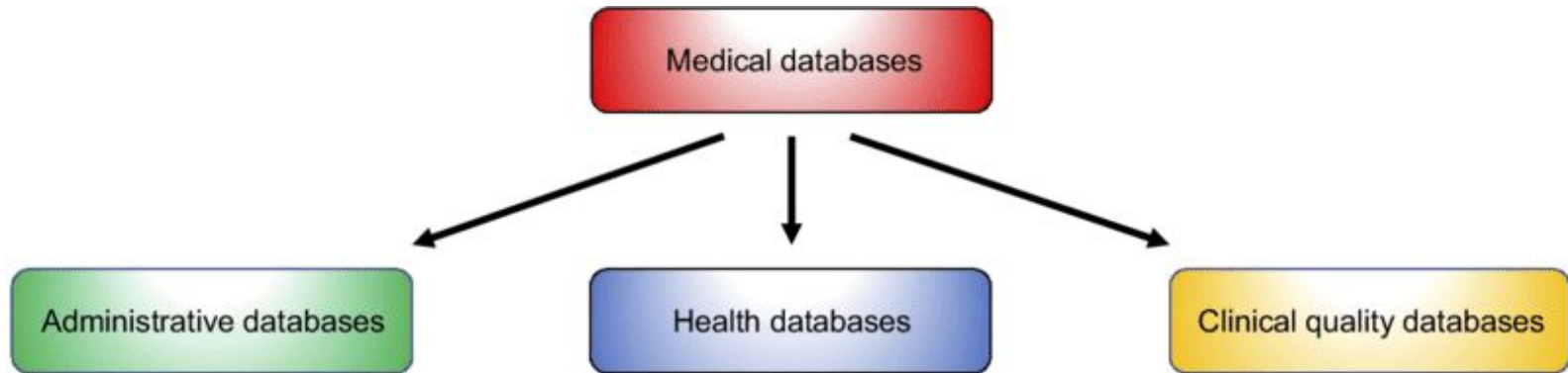
- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic

Many types of database



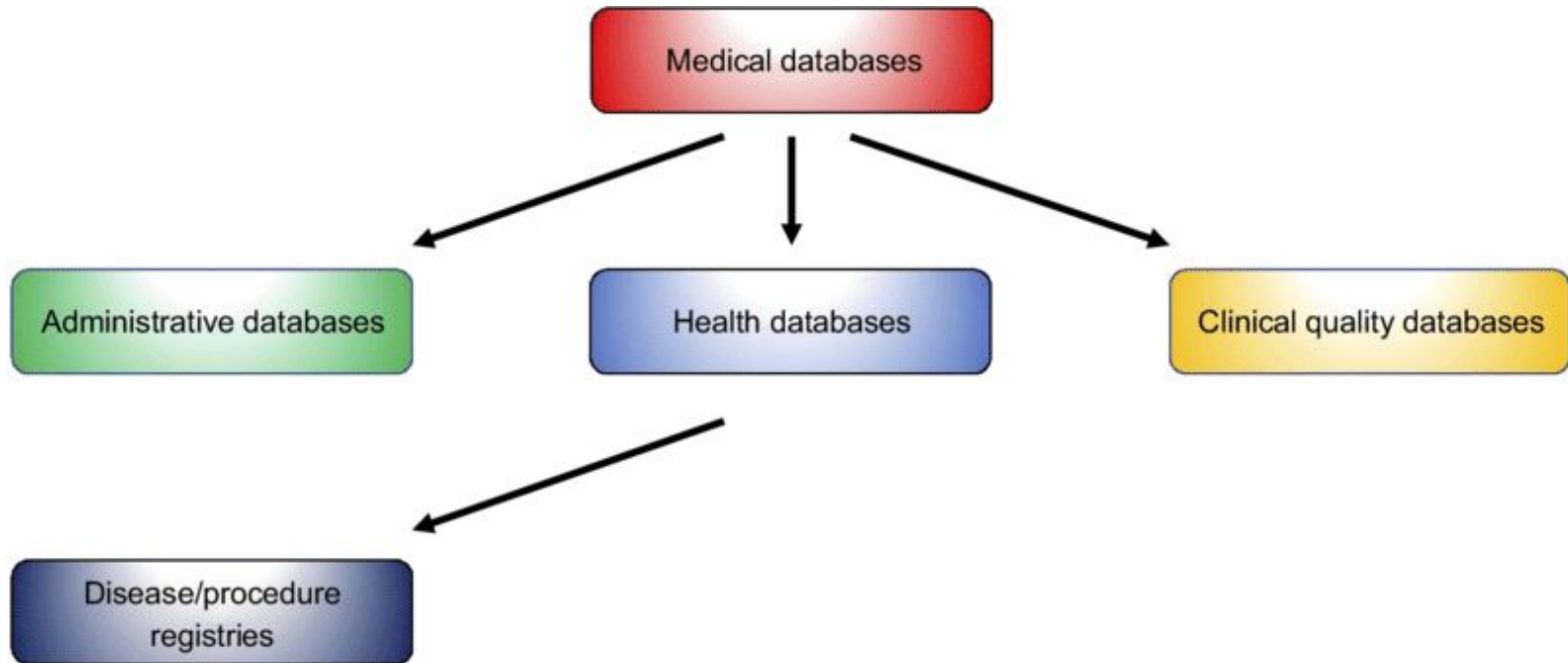
- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research

Many types of database



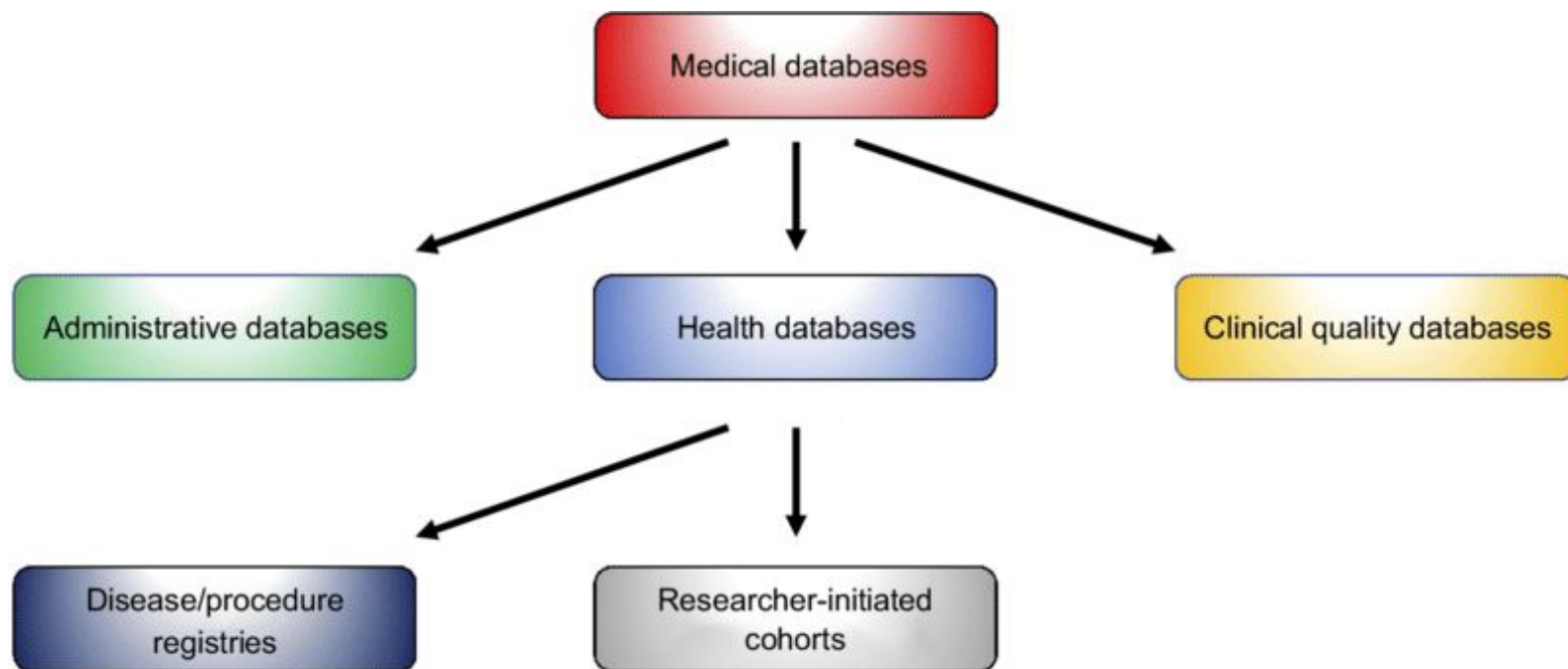
- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research
- Register detailed clinical data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

Many types of database



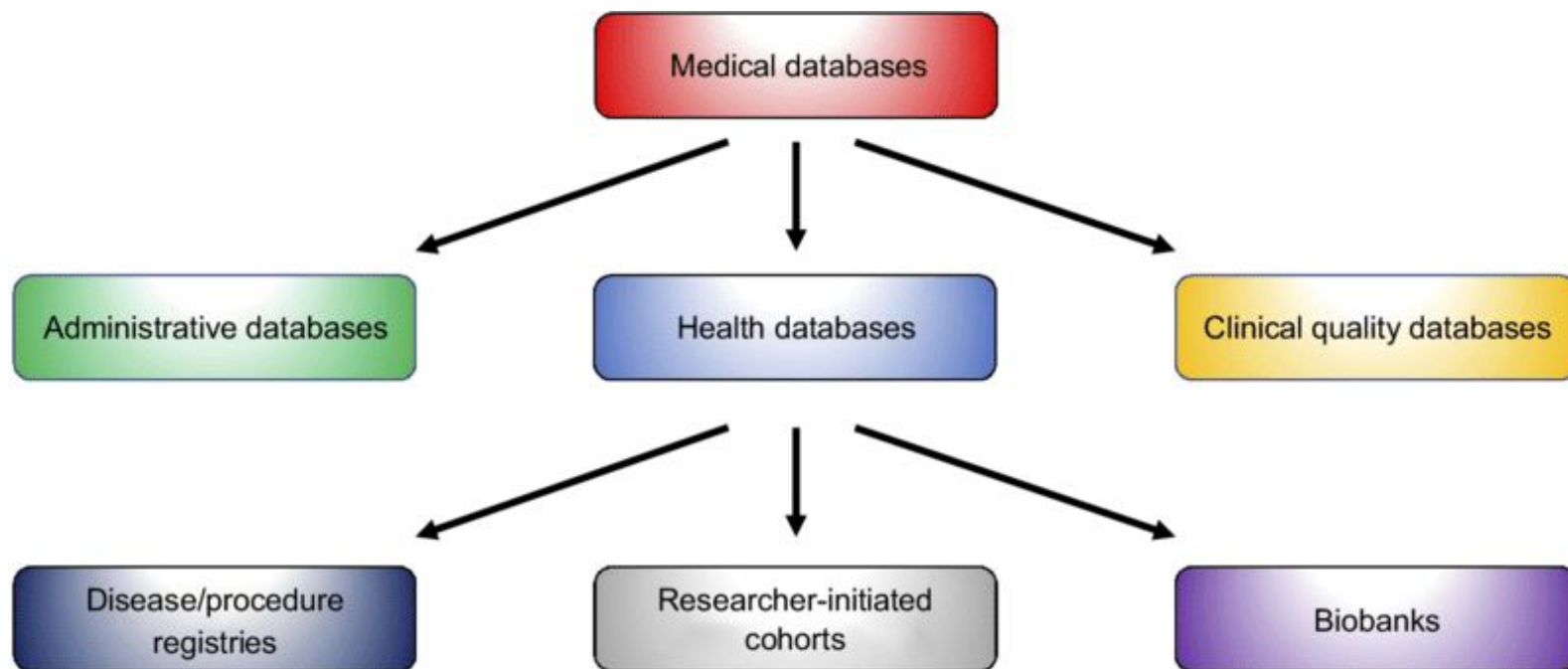
- All types of registries and databases that contain health-related data
- Register individuals according to geographic area, health insurance program, or attendance at a particular hospital or clinic
- Register health data for the purpose of surveillance and research
- Register detailed clinical data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

Consider 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

Sampling scope

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



Challenge of standardisation

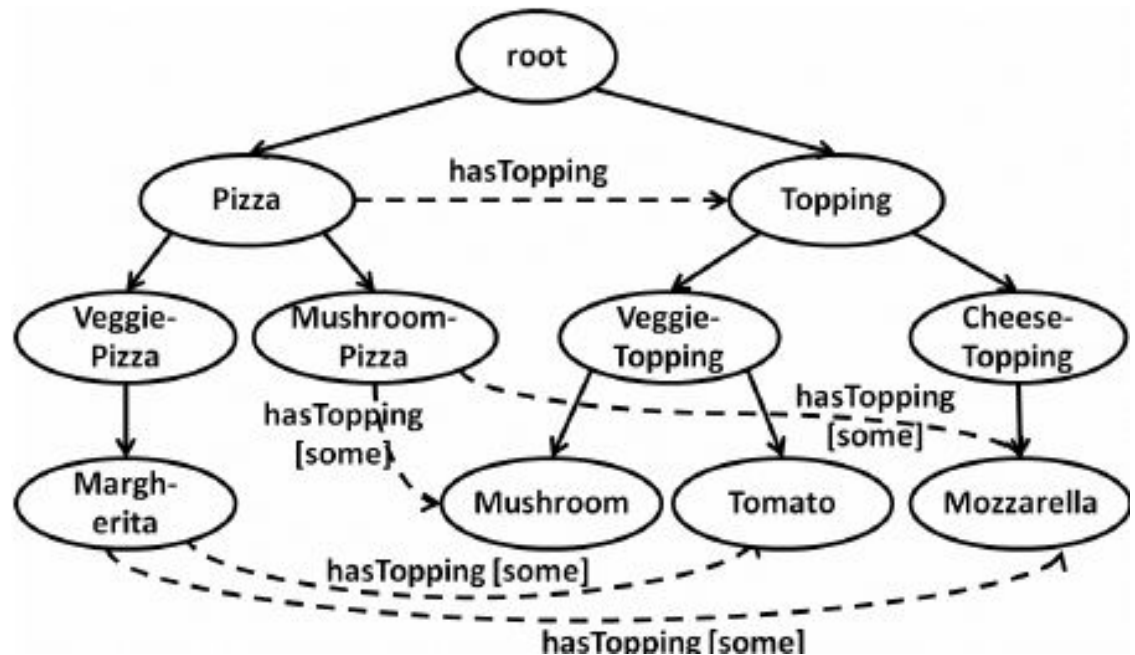


Generalisability

How do medical databases try to handle
standardisation?

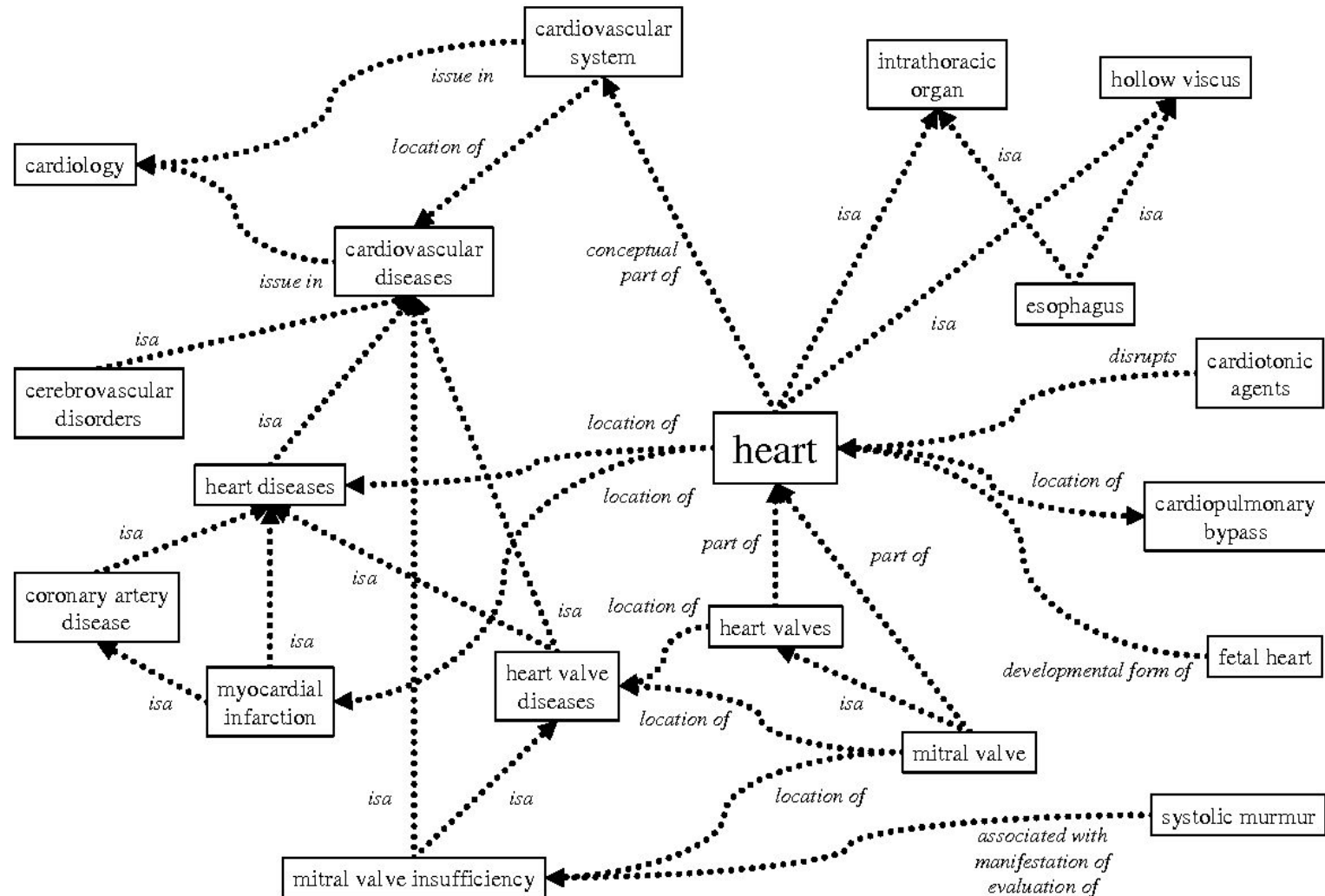
Ontologies for standardisation

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

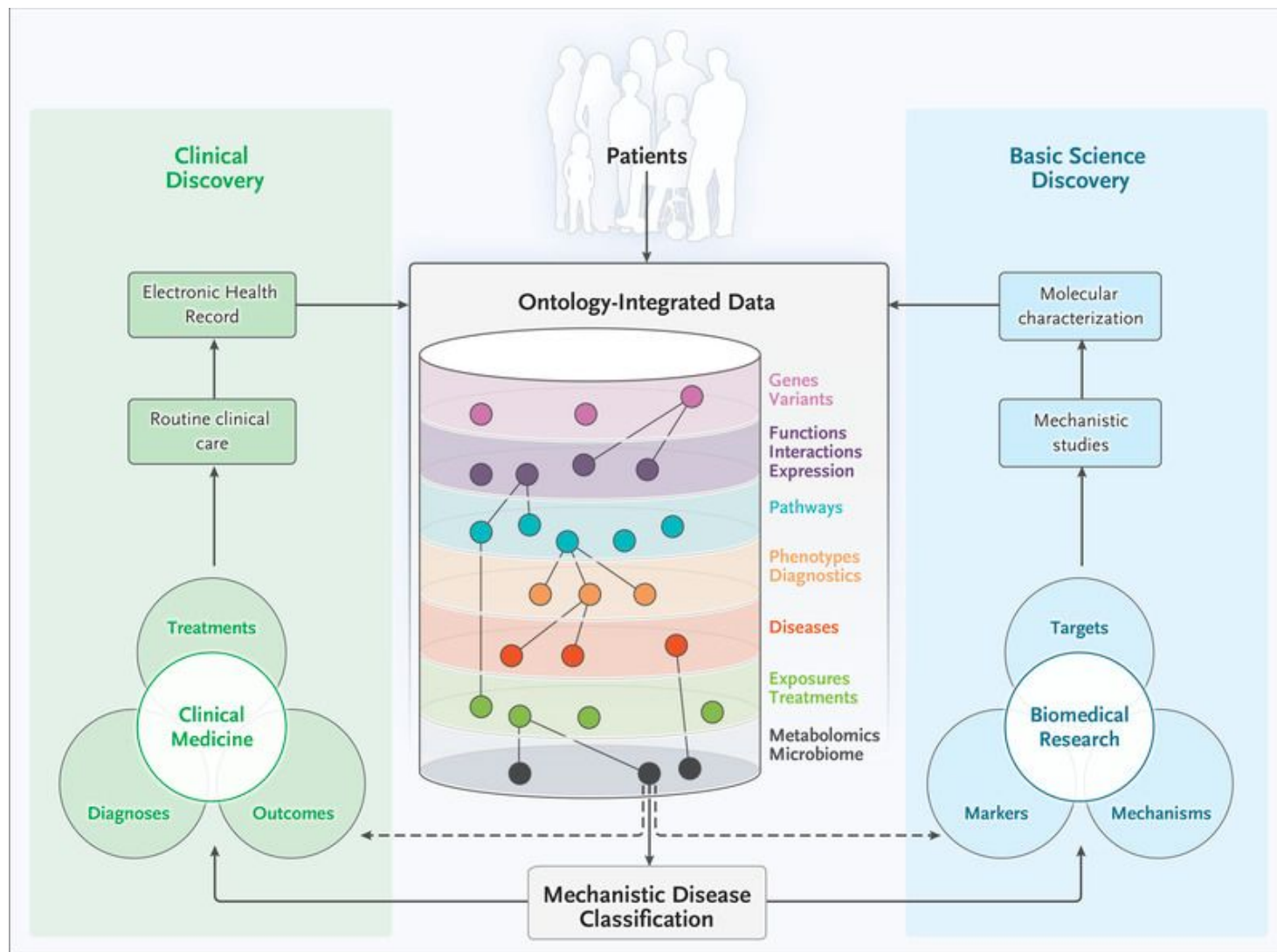


https://www.researchgate.net/figure/Example-pizza-ontology-represented-as-a-graph-G-a-and-a-changed-version-of-the-pizza_fig1_236842047

Medical Ontologies



Linking different ontologies



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

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

https://www.cdc.gov/nchs/icd/icd10cm_pcs_background.htm

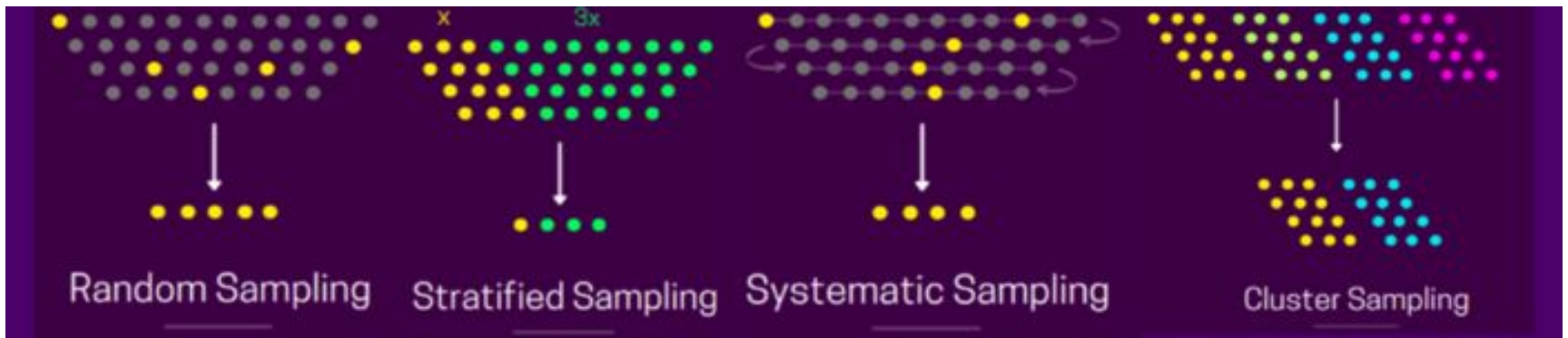
How do we sample for medical databases?

Sampling strategy

- Exhaustive isn't always exhaustive
- Numerous and often quite complex!
- Major source of bias so always carefully explore

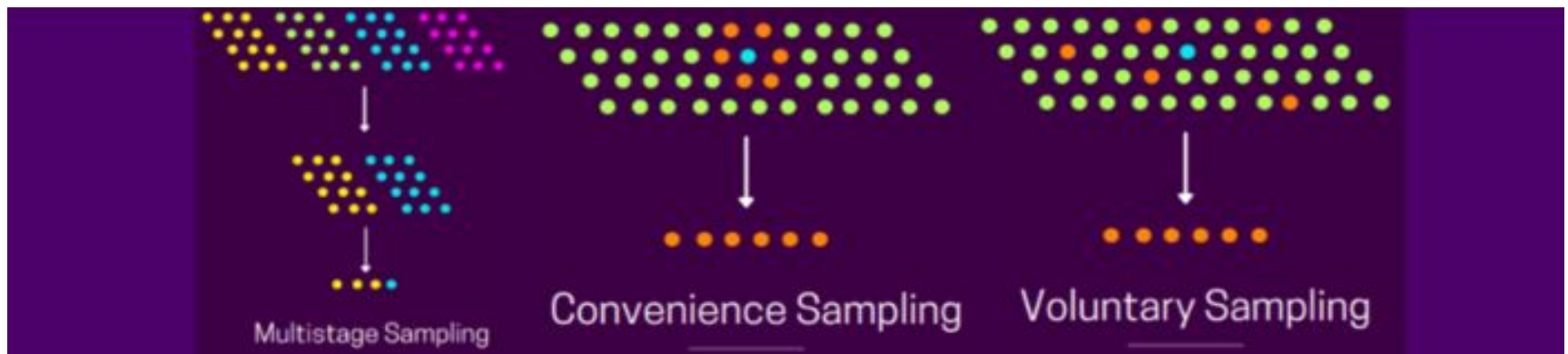
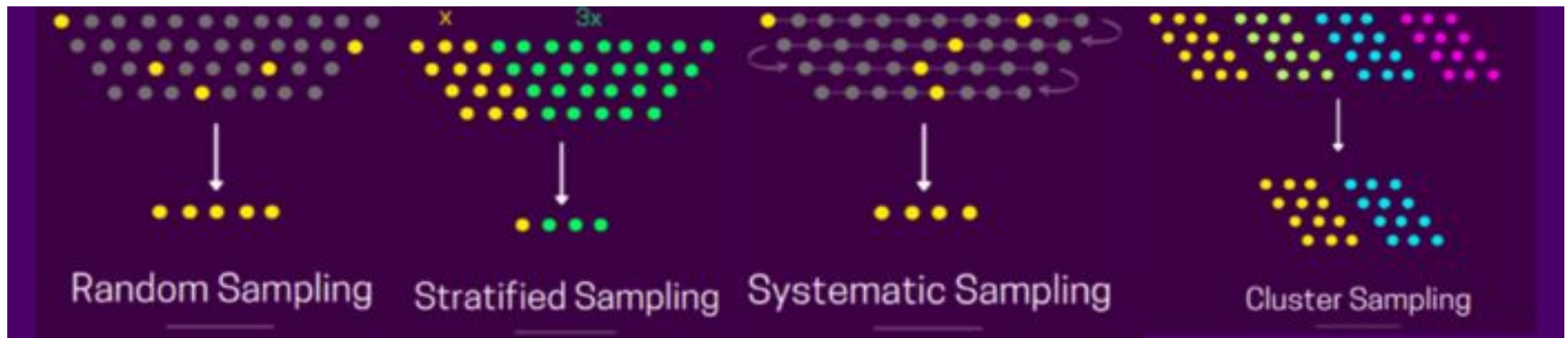
Sampling strategy

- Exhaustive isn't always exhaustive
- Numerous and often quite complex!
- Major source of bias so always carefully explore



Sampling strategy

- Exhaustive isn't always exhaustive
- Numerous and often quite complex!
- Major source of bias so always carefully explore

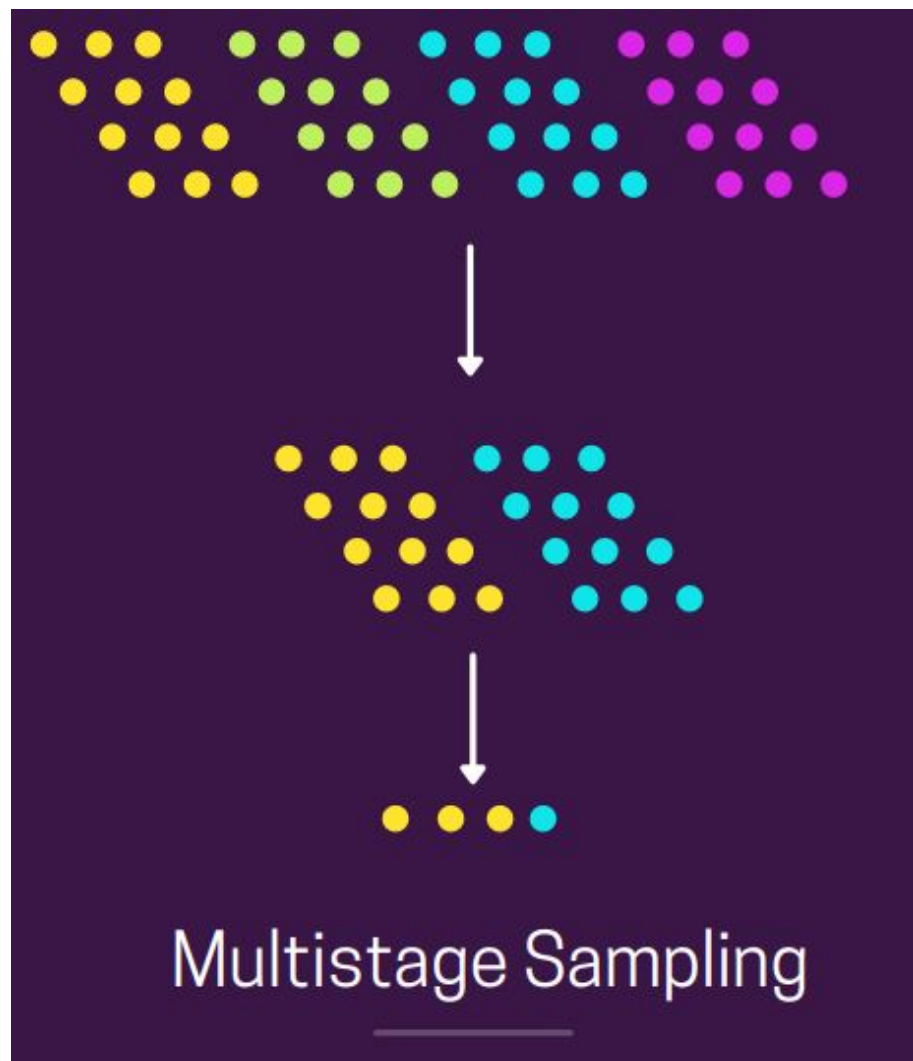


Survey weights

- 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 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.
- Generally poorly supported by machine learning libraries.

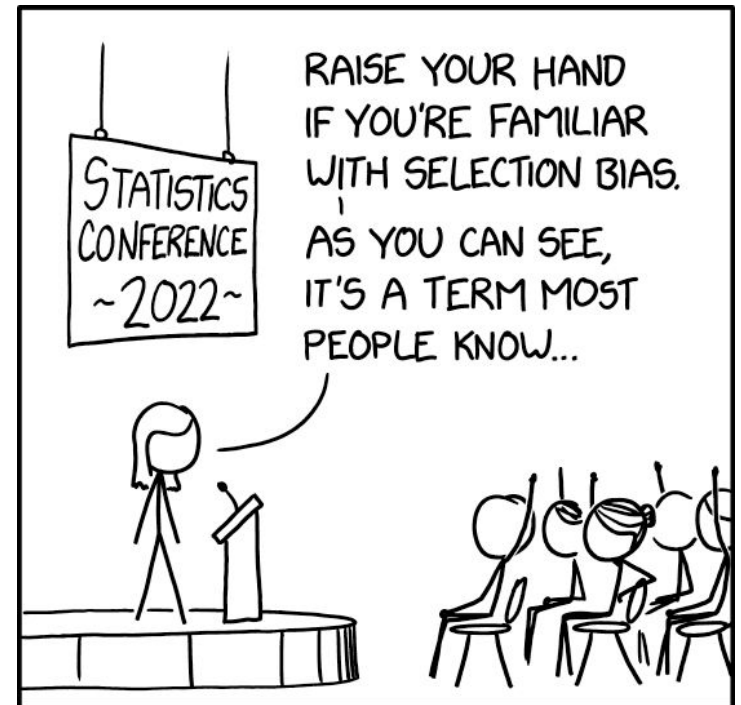


Types of weights

- 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

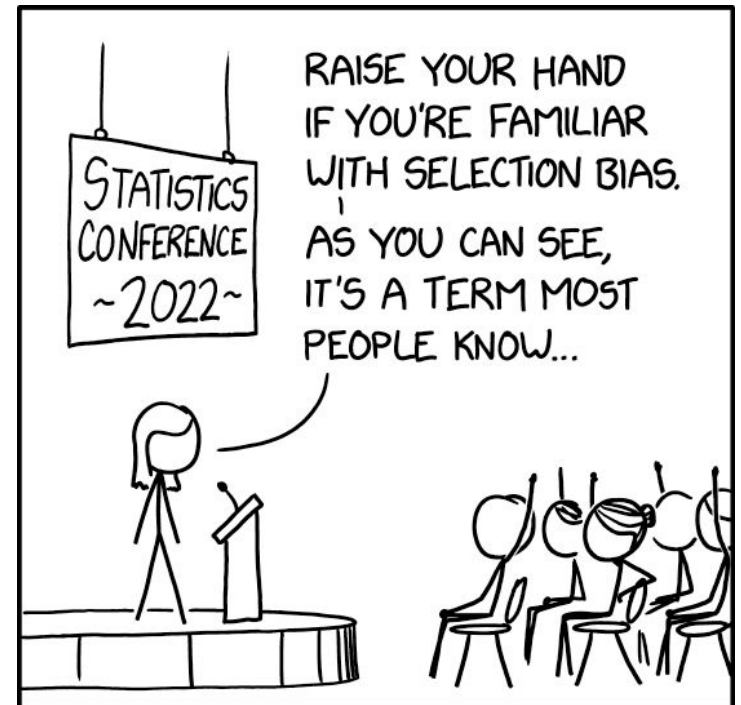
Types of weights

- 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



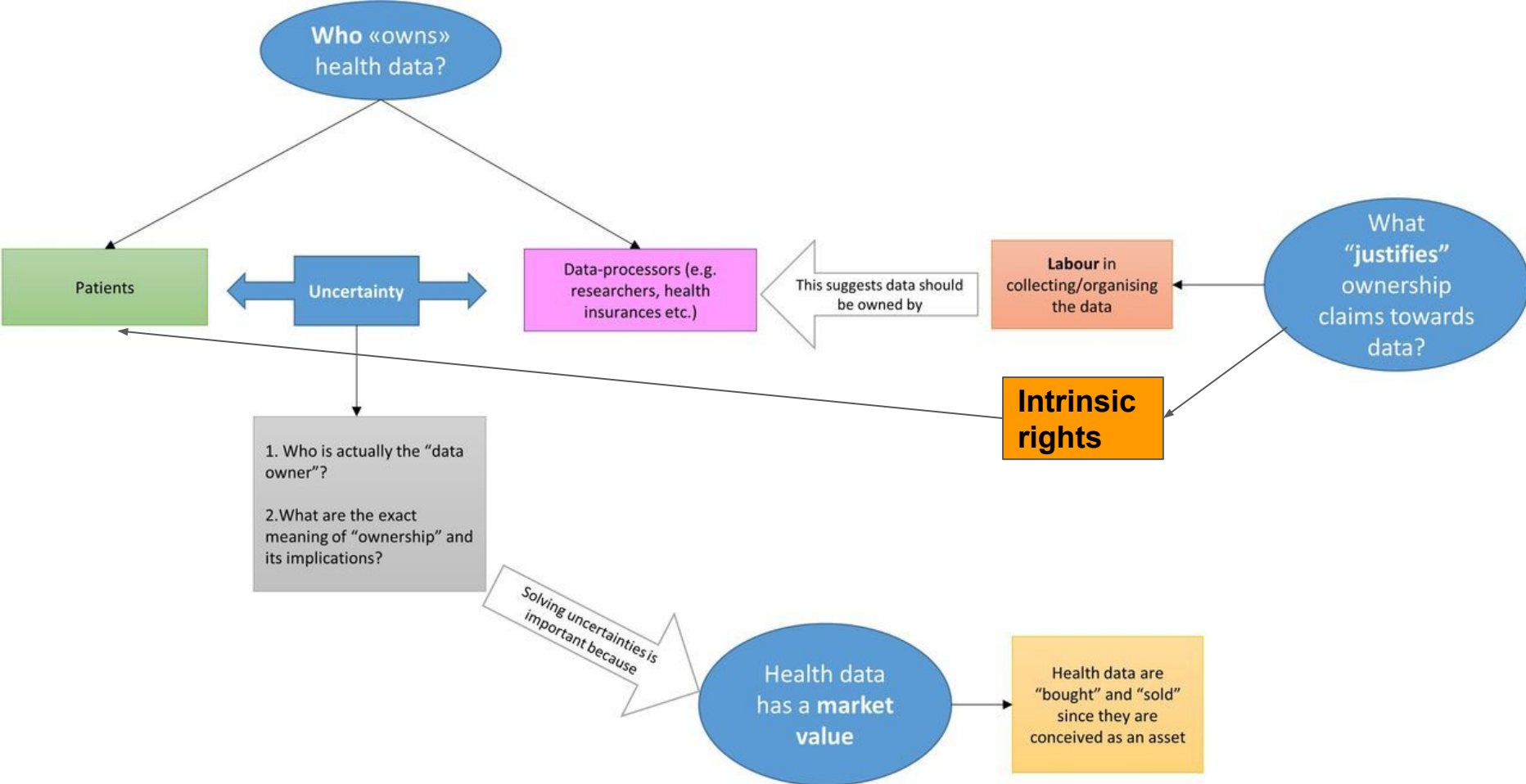
Types of weights

- 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



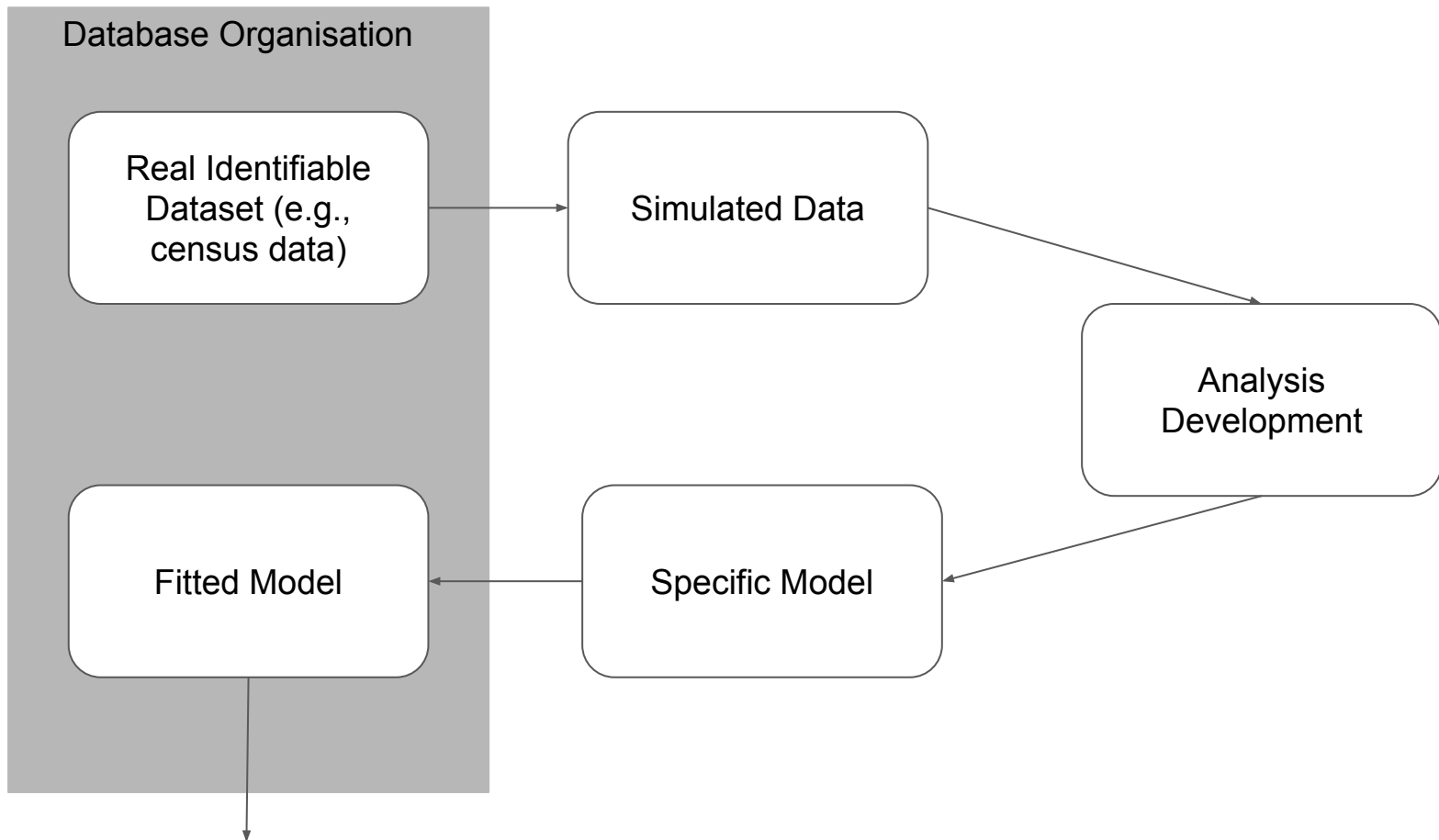
Who actually owns this data?

Data Ownership is Difficult



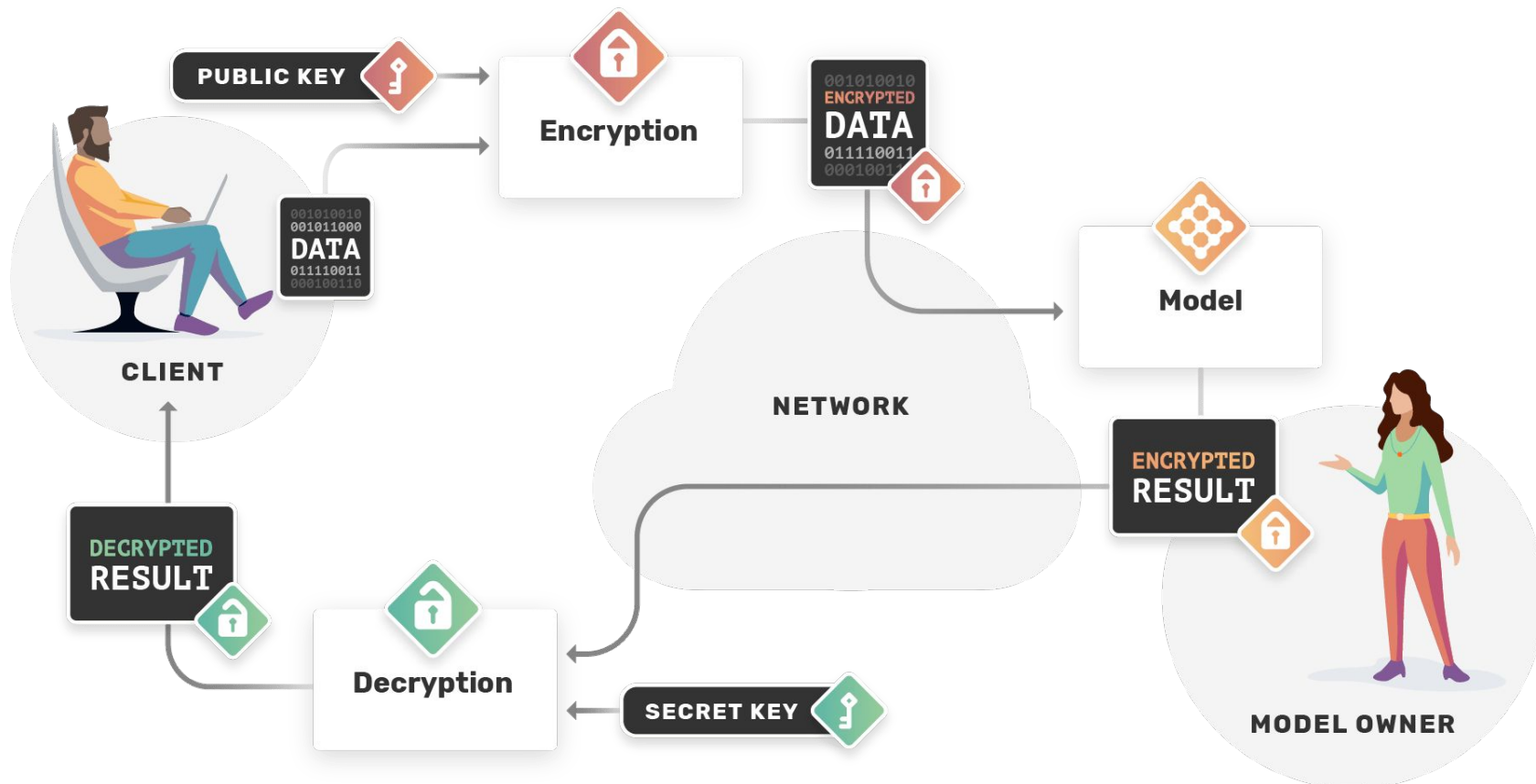
How do you protect privacy in these
databases?

No direct data access















Shared data but encrypted: homomorphic encryption

Partial to fully homomorphic encryption















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










Data privacy is a continuum

	EXPLICITLY PERSONAL	POTENTIALLY IDENTIFIABLE	NOT READILY IDENTIFIABLE
 <p>DIRECT IDENTIFIERS Data that identifies a person without additional information or by linking to information in the public domain (e.g., name, SSN)</p>	 <p>INTACT</p>	 <p>PARTIALLY MASKED</p>	 <p>PARTIALLY MASKED</p>
 <p>INDIRECT IDENTIFIERS Data that identifies an individual indirectly. Helps connect pieces of information until an individual can be singled out (e.g., DOB, gender)</p>	 <p>INTACT</p>	 <p>INTACT</p>	 <p>INTACT</p>
 <p>SAFEGUARDS and CONTROLS Technical, organizational and legal controls preventing employees, researchers or other third parties from re-identifying individuals</p>	 <p>NOT RELEVANT due to nature of data</p>	 <p>LIMITED or NONE IN PLACE</p>	 <p>CONTROLS IN PLACE</p>
<p>SELECTED EXAMPLES</p>	<p>Name, address, phone number, SSN, government-issued ID (e.g., Jane Smith, 123 Main Street, 555-555-5555)</p>	<p>Unique device ID, license plate, medical record number, cookie, IP address (e.g., MAC address 68:A8:6D:35:65:03)</p>	<p>Same as Potentially Identifiable except data are also protected by safeguards and controls (e.g., hashed MAC addresses & legal representations)</p>










Indirectly identifiable: Pseudonymous Data

	KEY CODED	PSEUDONYMOUS	PROTECTED PSEUDONYMOUS
 <p>DIRECT IDENTIFIERS Data that identifies a person without additional information or by linking to information in the public domain (e.g., name, SSN)</p>	 ELIMINATED or TRANSFORMED	 ELIMINATED or TRANSFORMED	 ELIMINATED or TRANSFORMED
 <p>INDIRECT IDENTIFIERS Data that identifies an individual indirectly. Helps connect pieces of information until an individual can be singled out (e.g., DOB, gender)</p>	 INTACT	 INTACT	 INTACT
 <p>SAFEGUARDS and CONTROLS Technical, organizational and legal controls preventing employees, researchers or other third parties from re-identifying individuals</p>	 CONTROLS IN PLACE	 LIMITED or NONE IN PLACE	 CONTROLS IN PLACE
	<p>Clinical or research datasets where only curator retains key (e.g., Jane Smith, diabetes, HgB 15.1 g/dl = Csrk123)</p>	<p>Unique, artificial pseudonyms replace direct identifiers (e.g., HIPAA Limited Datasets, John Doe = 5L7T LX619Z) (unique sequence not used anywhere else)</p>	<p>Same as Pseudonymous, except data are also protected by safeguards and controls</p>

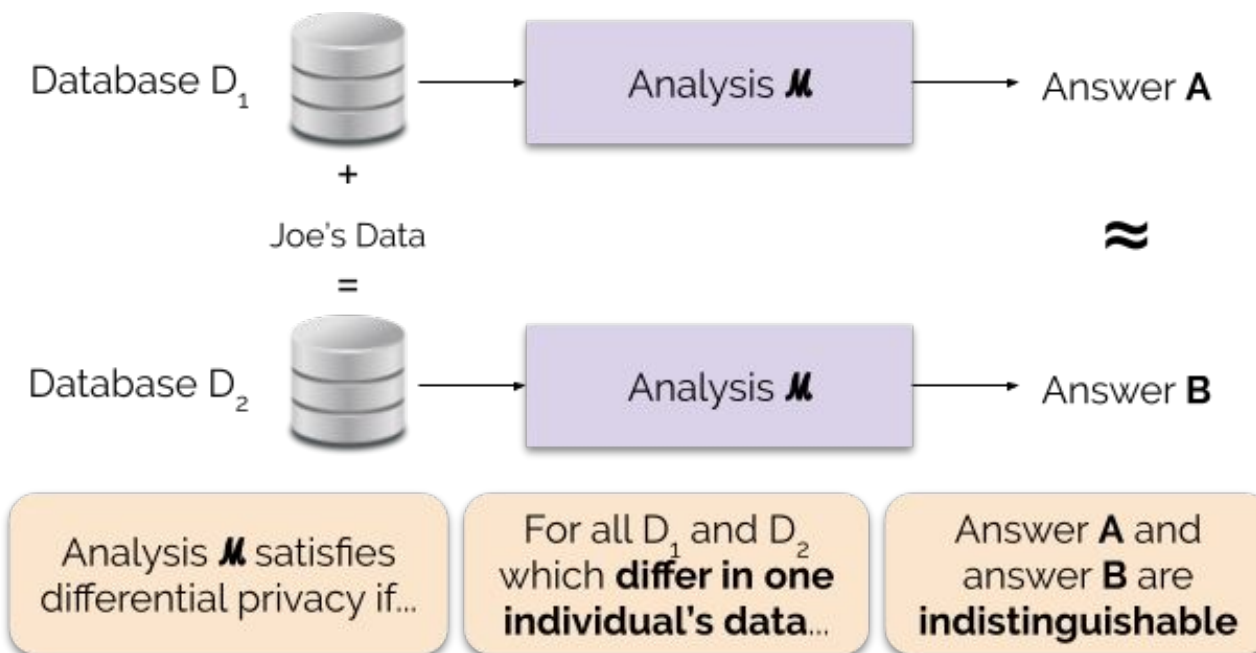
Identifiers removed/broken: De-Identified Data

	DE-IDENTIFIED	PROTECTED DE-IDENTIFIED
 <p>DIRECT IDENTIFIERS Data that identifies a person without additional information or by linking to information in the public domain (e.g., name, SSN)</p>	 <p>ELIMINATED or TRANSFORMED</p>	 <p>ELIMINATED or TRANSFORMED</p>
 <p>INDIRECT IDENTIFIERS Data that identifies an individual indirectly. Helps connect pieces of information until an individual can be singled out (e.g., DOB, gender)</p>	 <p>ELIMINATED or TRANSFORMED</p>	 <p>ELIMINATED or TRANSFORMED</p>
 <p>SAFEGUARDS and CONTROLS Technical, organizational and legal controls preventing employees, researchers or other third parties from re-identifying individuals</p>	 <p>LIMITED or NONE IN PLACE</p>	 <p>CONTROLS IN PLACE</p>
	Data are suppressed, generalized, perturbed, swapped, etc. (e.g., GPA: 3.2 = 3.0-3.5, gender: female = gender: male)	Same as De-Identified, except data are also protected by safeguards and controls

Non-identifiability Guarantee: Anonymous Data

	ANONYMOUS	AGGREGATED ANONYMOUS
 <p>DIRECT IDENTIFIERS Data that identifies a person without additional information or by linking to information in the public domain (e.g., name, SSN)</p>	 <p>ELIMINATED or TRANSFORMED</p>	 <p>ELIMINATED or TRANSFORMED</p>
 <p>INDIRECT IDENTIFIERS Data that identifies an individual indirectly. Helps connect pieces of information until an individual can be singled out (e.g., DOB, gender)</p>	 <p>ELIMINATED or TRANSFORMED</p>	 <p>ELIMINATED or TRANSFORMED</p>
 <p>SAFEGUARDS and CONTROLS Technical, organizational and legal controls preventing employees, researchers or other third parties from re-identifying individuals</p>	 <p>NOT RELEVANT due to nature of data</p>	 <p>NOT RELEVANT due to high degree of data aggregation</p>
	<p>For example, noise is calibrated to a data set to hide whether an individual is present or not (differential privacy)</p>	<p>Very highly aggregated data (e.g., statistical data, census data, or population data that 52.6% of Washington, DC residents are women)</p>

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 → $\Pr[\mathcal{M}(D_1) \in O]$

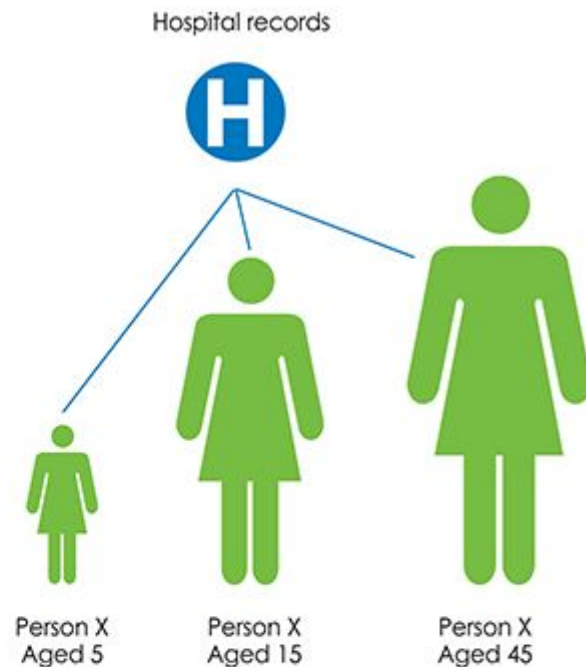
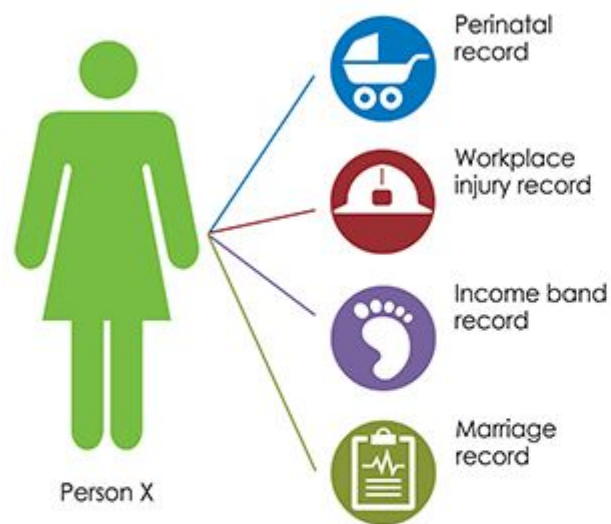
Probability of seeing output O on input D_2 → $\Pr[\mathcal{M}(D_2) \in O]$

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

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



Many different data access processes

- Buy access and get processed data
- Apply for individual fields and justify why
- Full pre-registration of analysis

Let's take a short break!

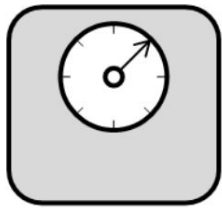
So, you've got access to a database, what
now?

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

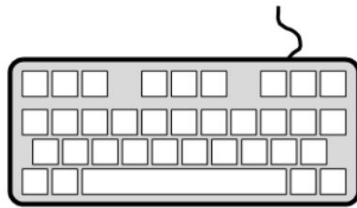
Data Quality



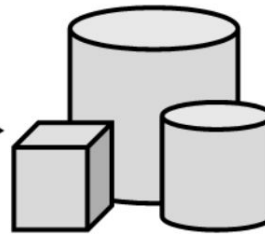
Actual value:
200.6 lbs.



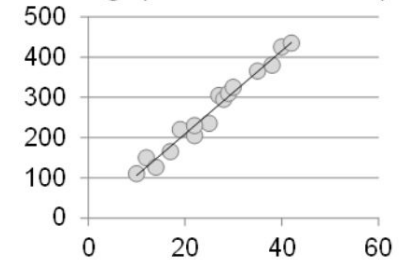
Recorded value:
200 lbs.



Data warehouse value:
200 kg



Analytic value:
100 kg (mean 200 & 0)



Measured (same day)

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

Measured (diff. days)

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

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

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

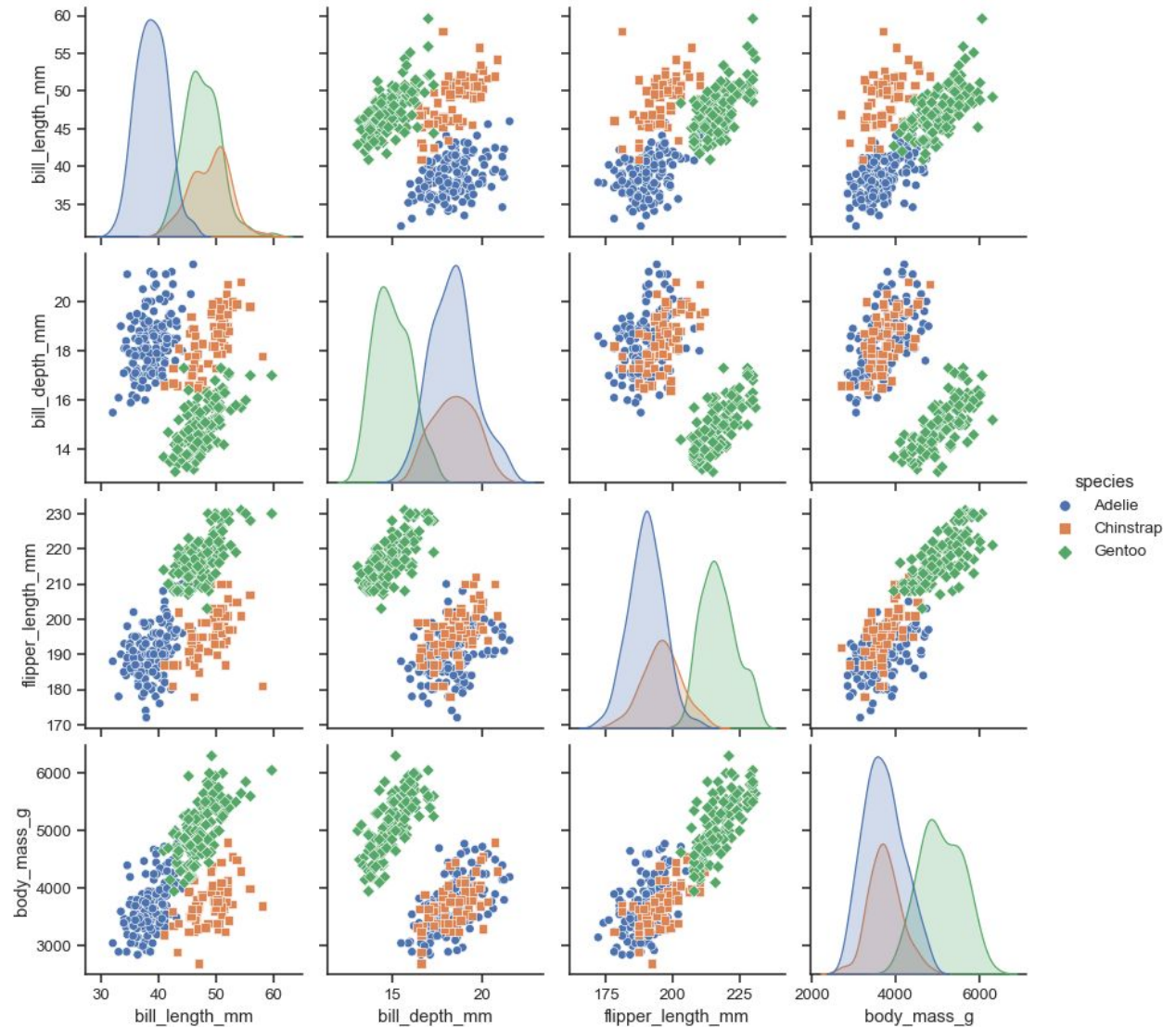
Under review

Slide from Dr. Hadi Kharrazi

9 months & >25 rules to clean weight

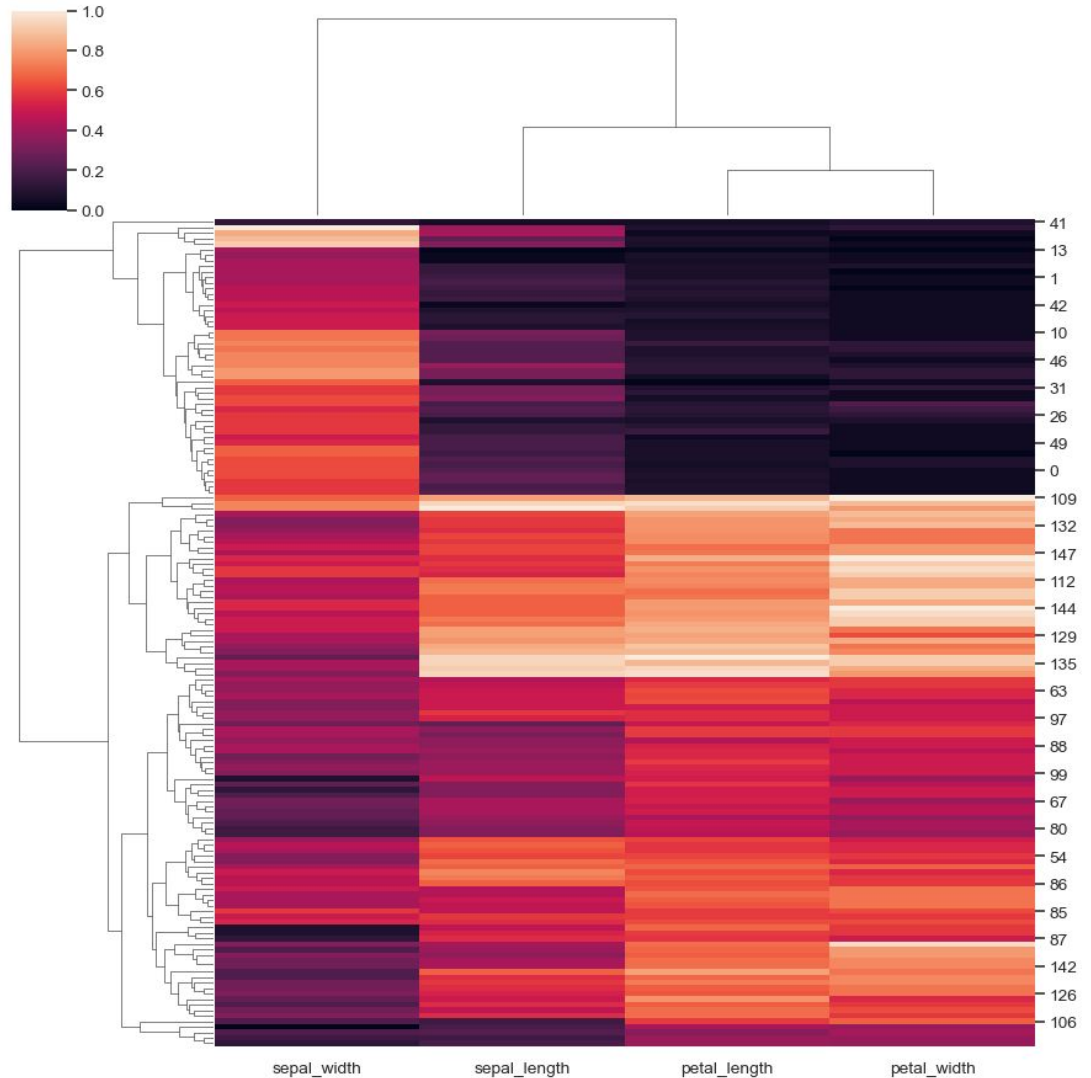
Exploratory Data Analysis

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



Exploratory Data Analysis

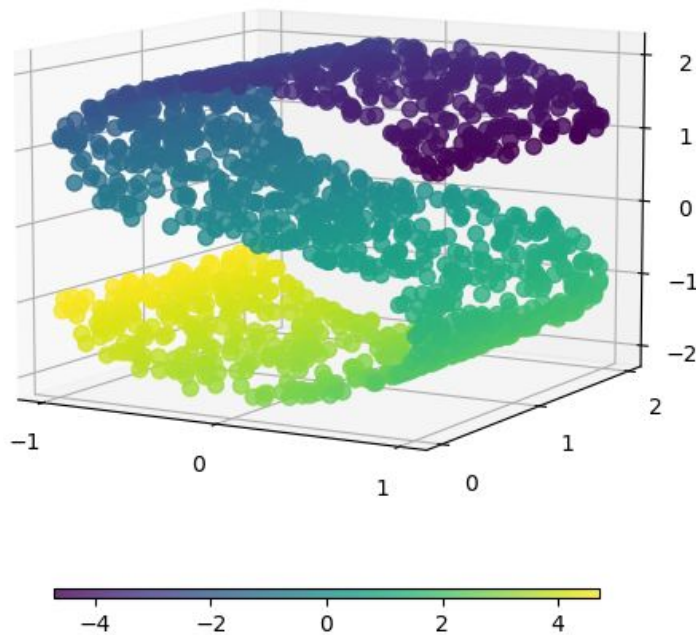
- Individual variable distributions
- Pairwise variable distributions
- Distributions relative to variable(s) of interest
- Hierarchical clustering of variables
- Point analysis of extreme values



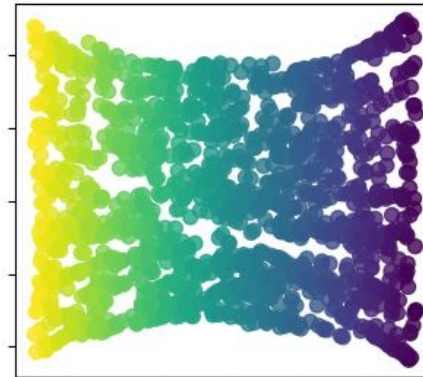
How do I look at all the data together?

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

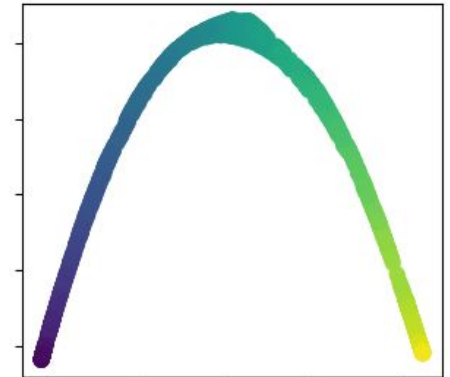
Original S-curve samples



Isomap Embedding



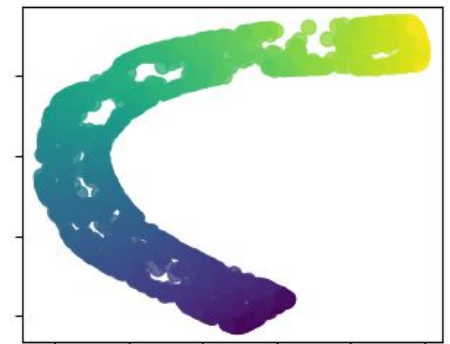
Spectral Embedding



Multidimensional scaling

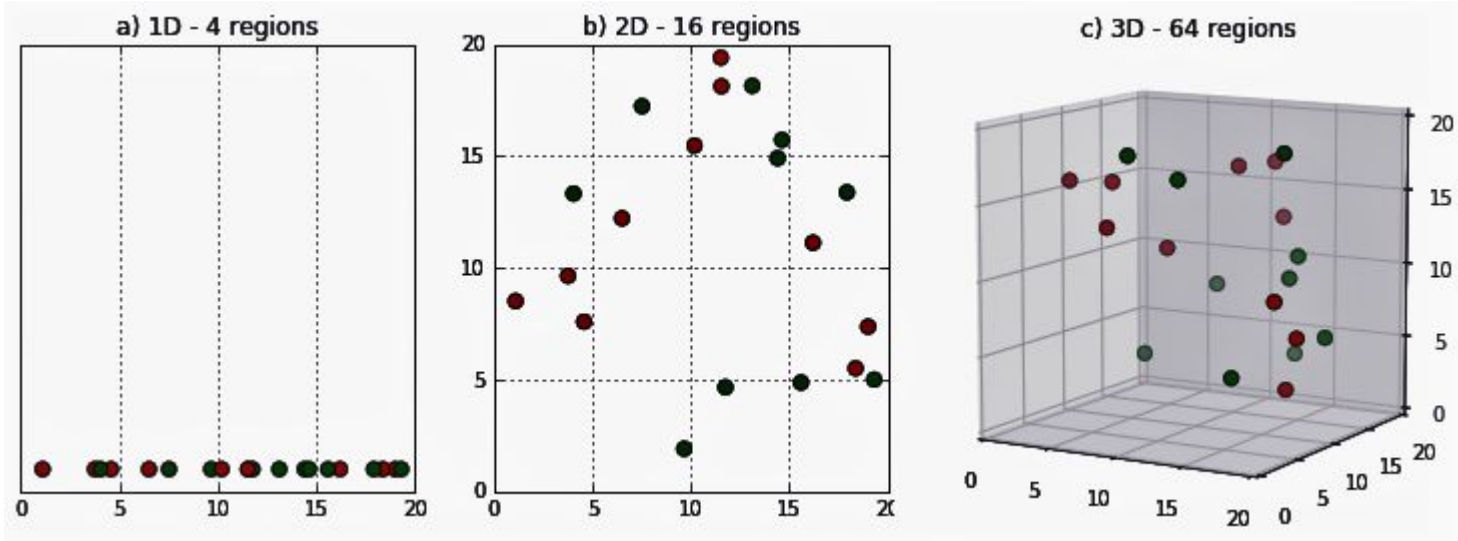


T-distributed Stochastic Neighbor Embedding



Why is this hard?

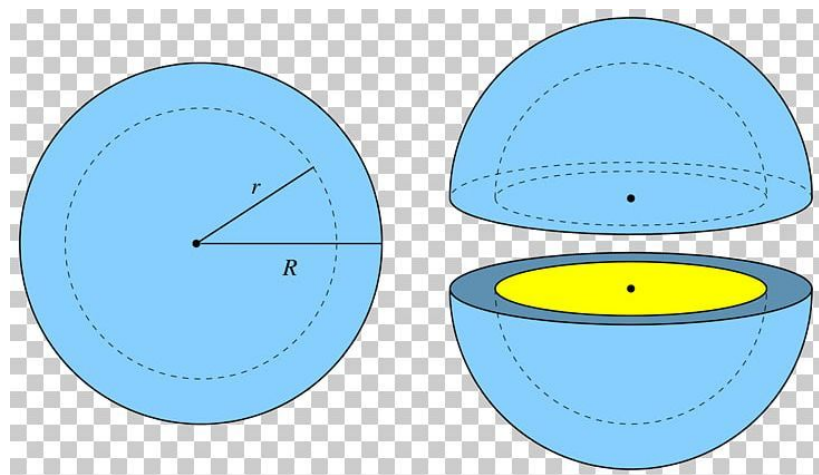
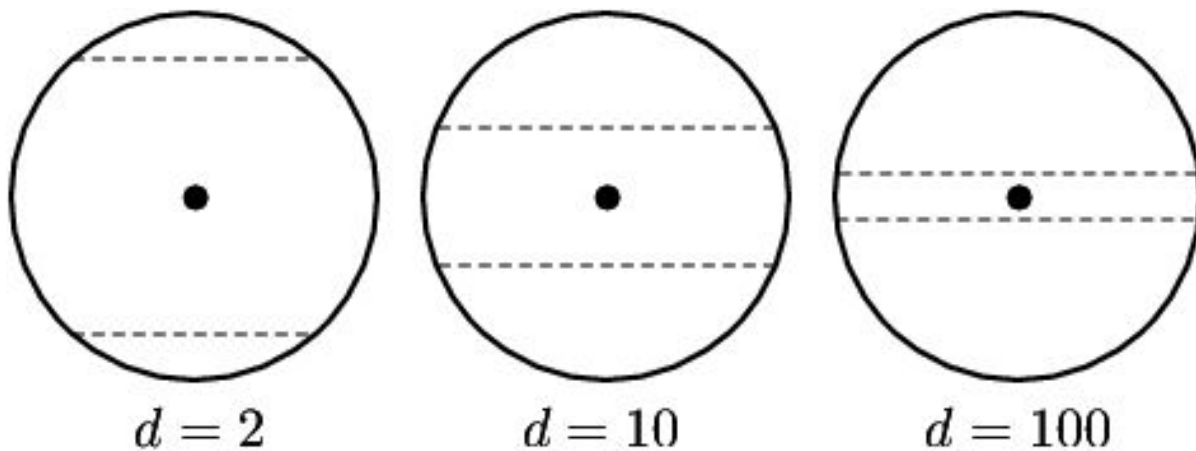
High dimensional data is sparse



<https://medium.com/analytics-vidhya/the-curse-of-dimensionality-and-its-cure-f9891ab72e5c>

High dimensional space is counterintuitive

Orthogonality -> Band-size to capture 99% of the volume of a sphere:



Mass becomes increasingly “shell-like”

No representation is perfect

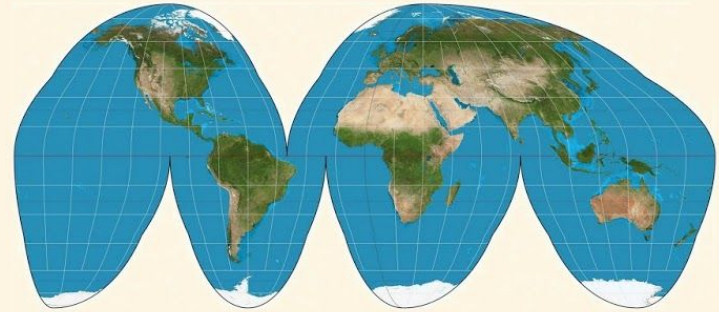
MERCATOR



GALL-PETERS



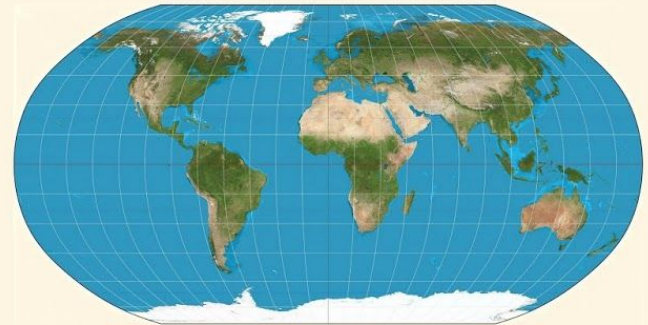
GOODE-HOMOLOGINE



WATERMELON



ALBERS

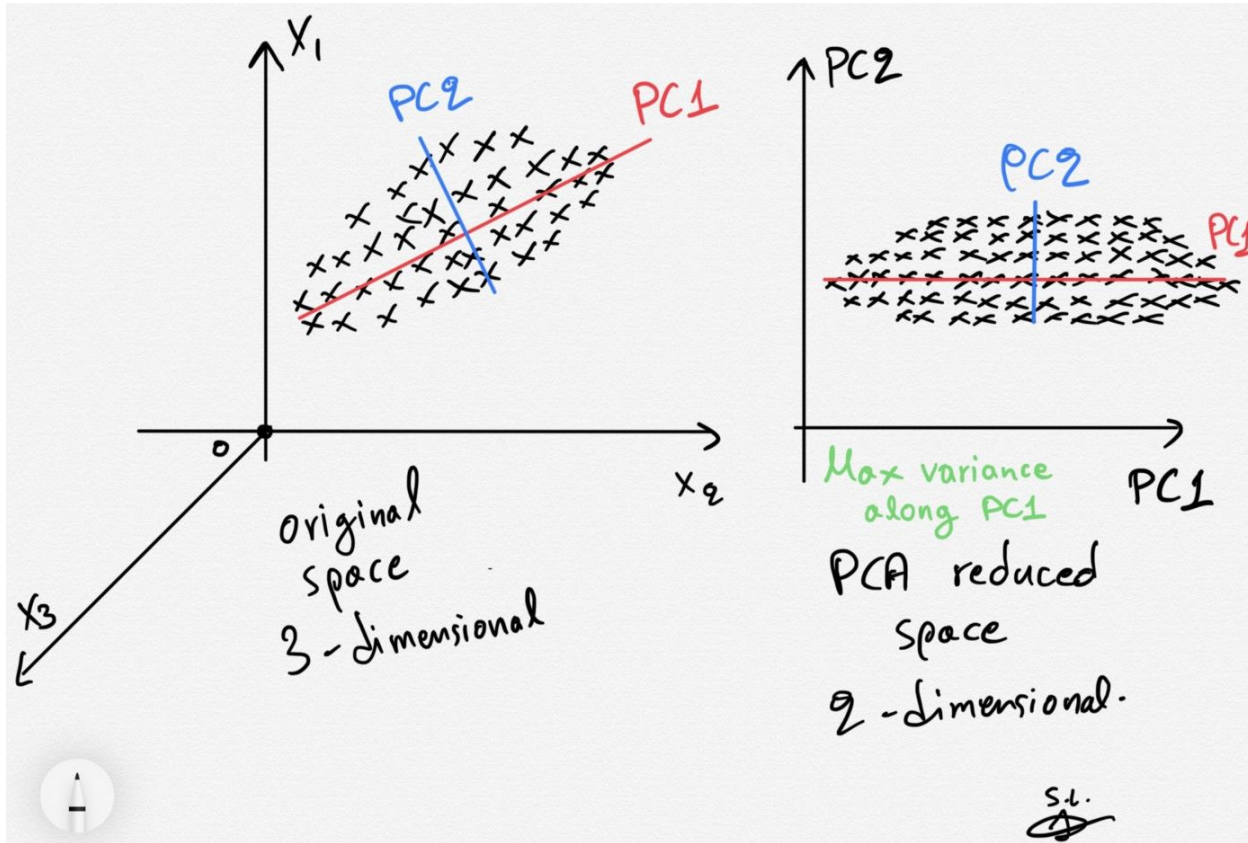


ROBINSON

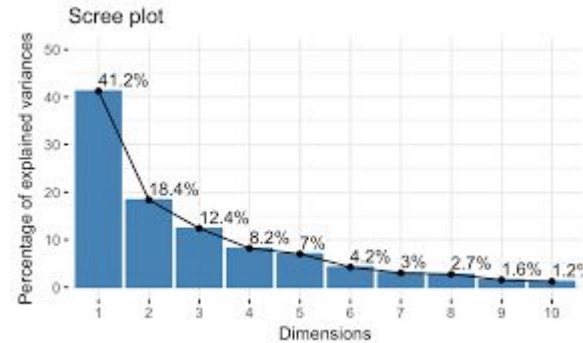
So, how can we do it?

Principal Component Analysis (PCA): Variance

Mean center data -> Generate Covariance Matrix ->
Eigendecomposition -> Sort Eigenvalues



- How many components?
Screen/elbow plot



- What variables contribute most to PCs? BiPlot

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.

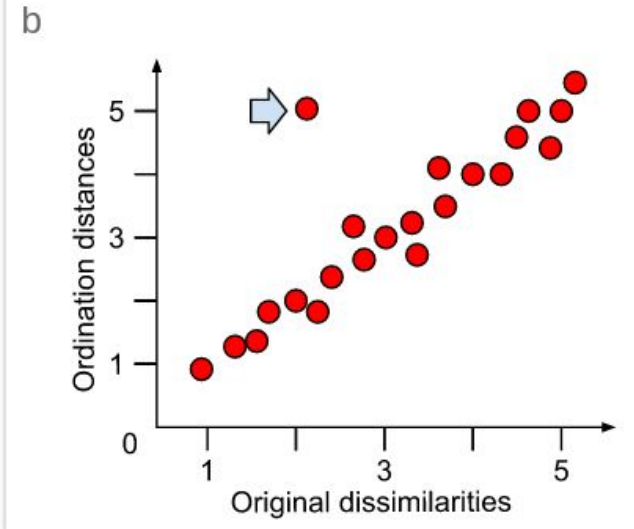
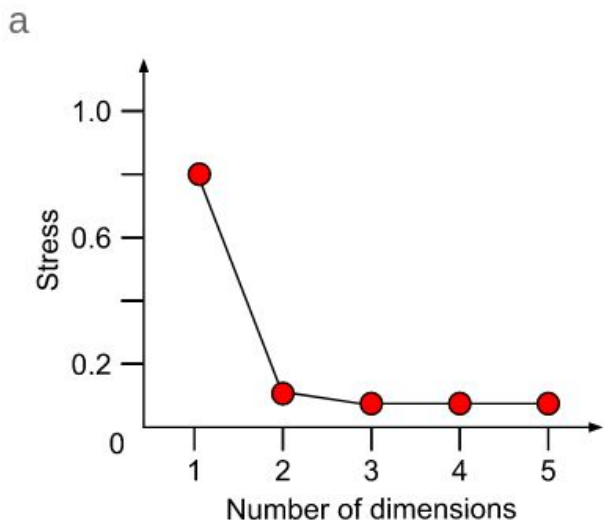
Where x_1, \dots, x_N are data points with their new set of coordinates in lower dimensional space.

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.

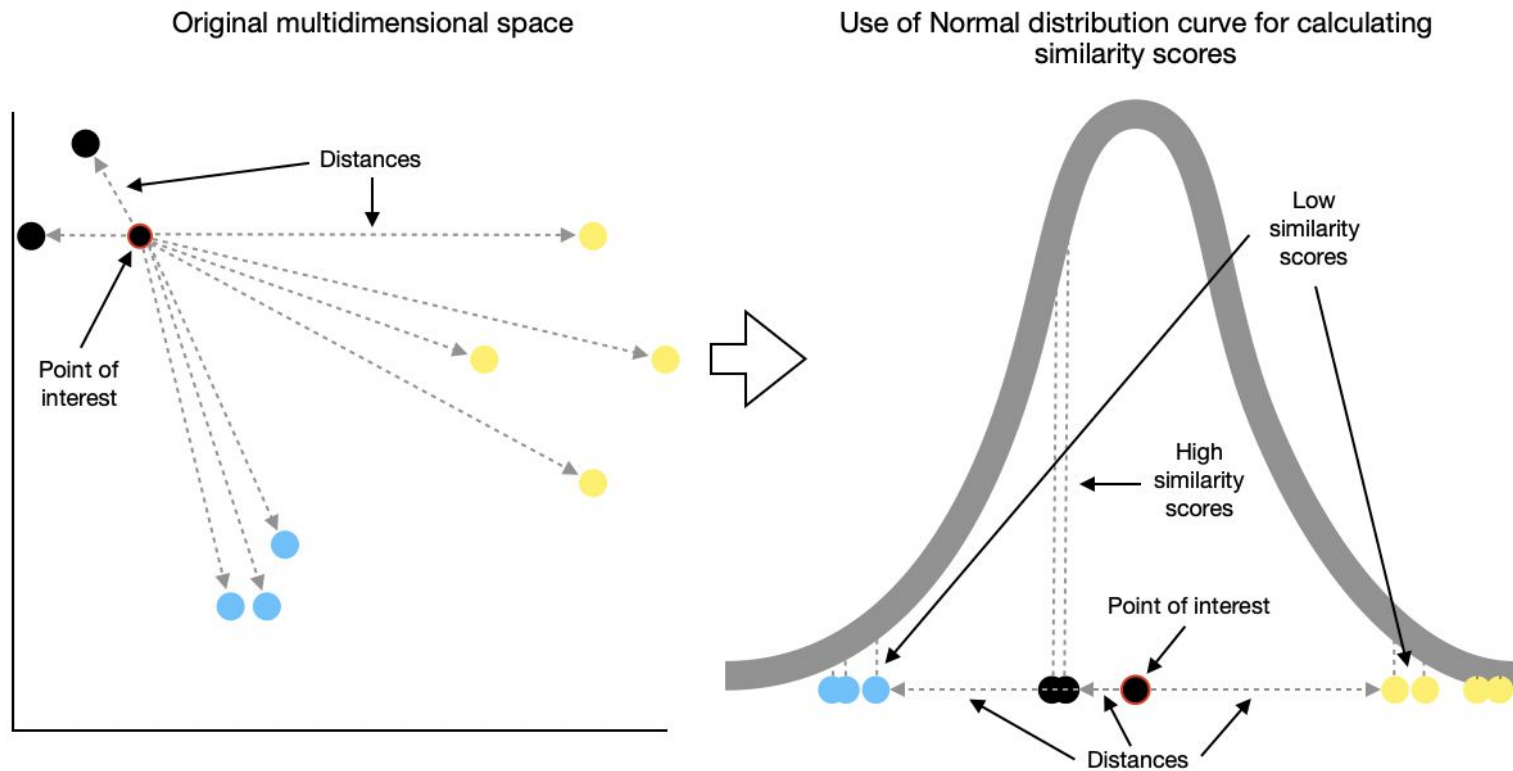
The closer the value of $||x_i - x_j||$ is to d_{ij} the

Non-Metric: Ranks



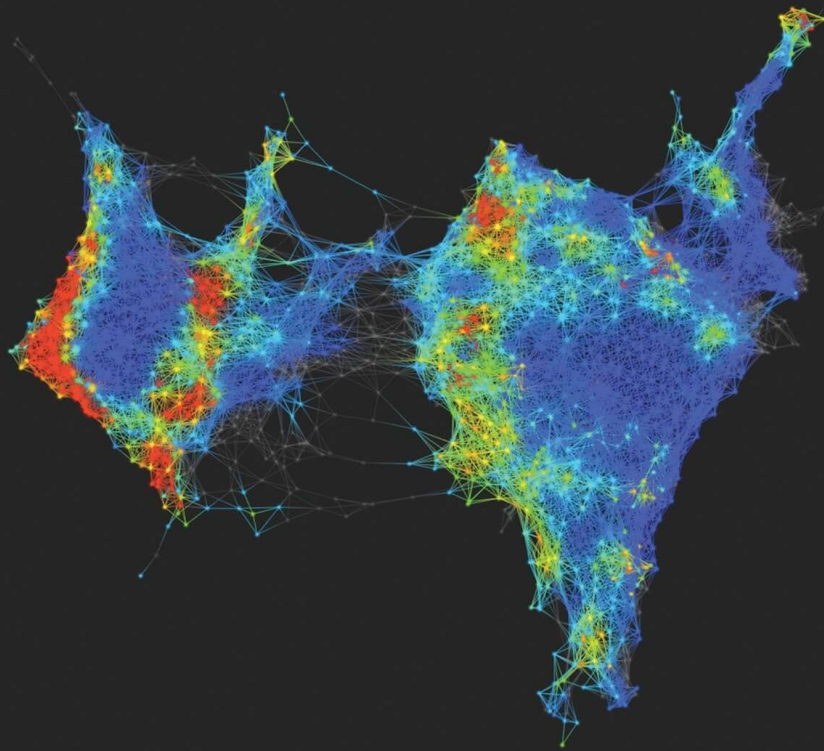
t-SNE/UMAP: Probabilities

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



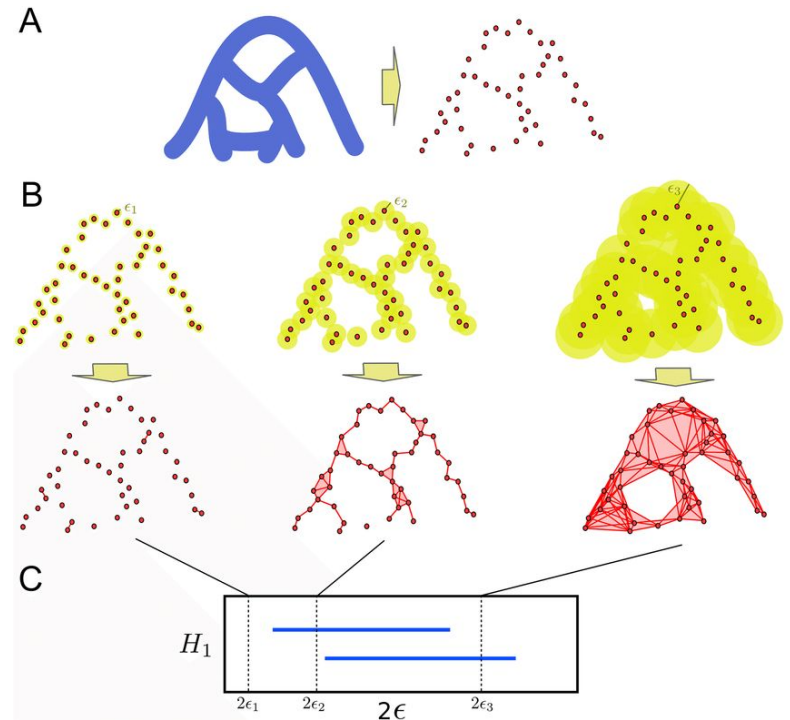
Topological Data Analysis

Biobank Database: Patient Stratification **AYASDI**

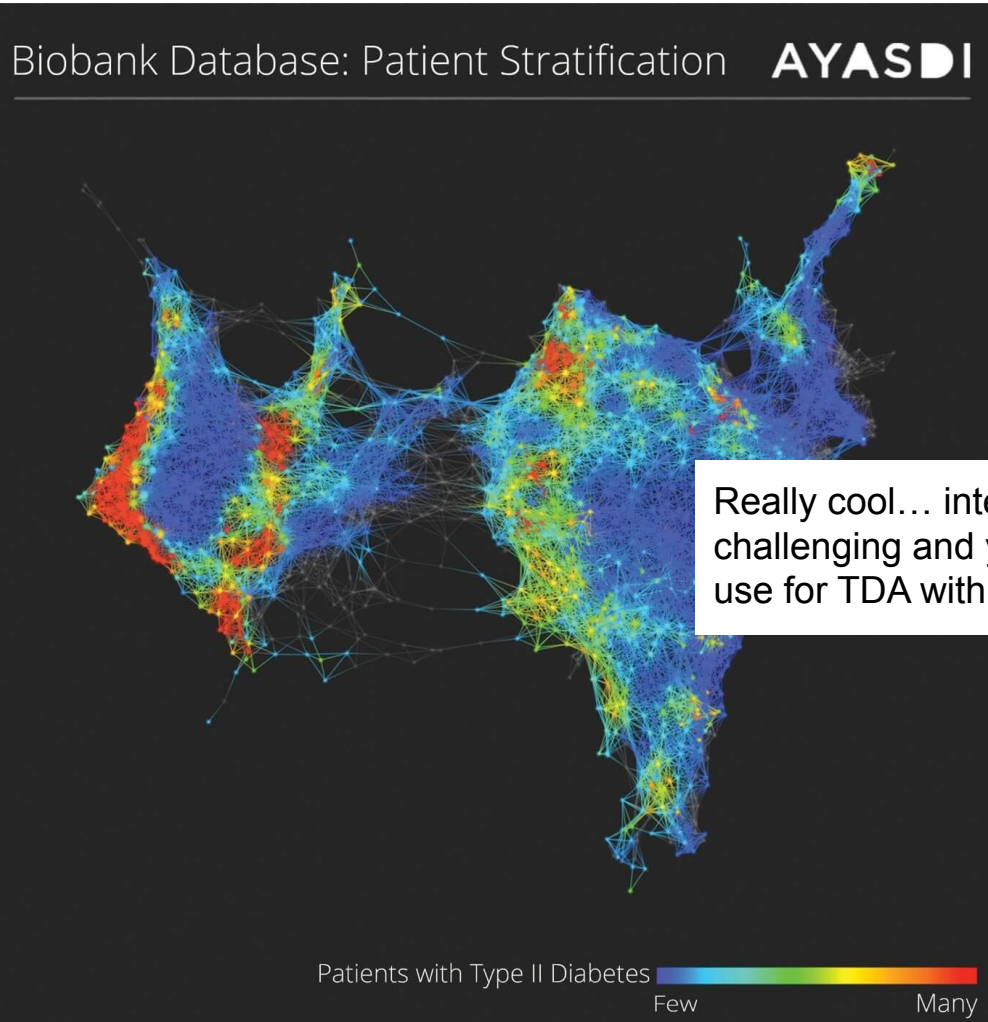


Patients with Type II Diabetes 
Few Many

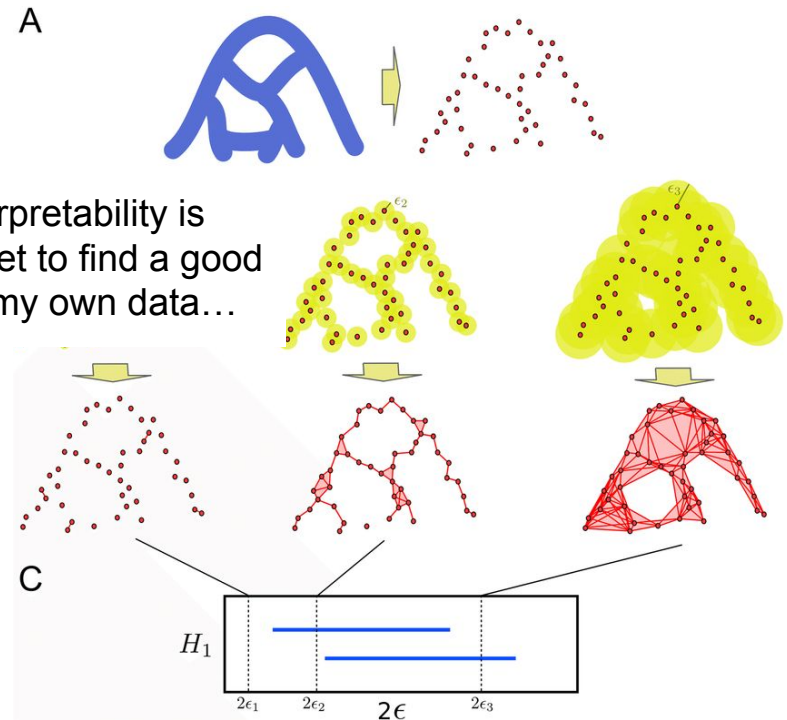
- Point clouds \rightarrow increase radius \rightarrow simplicial complexes \rightarrow topological characteristics



Topological Data Analysis

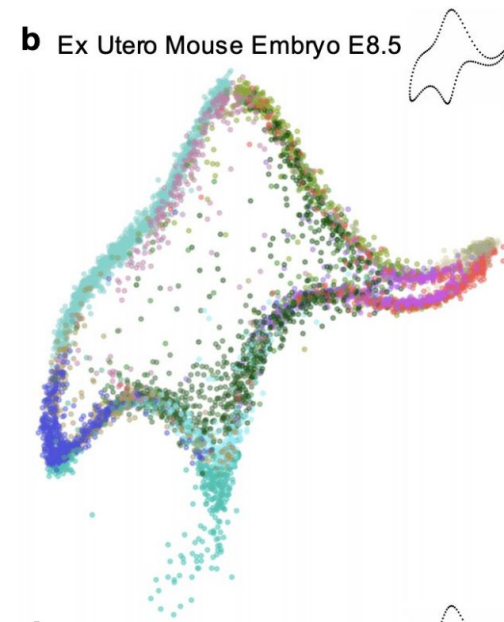
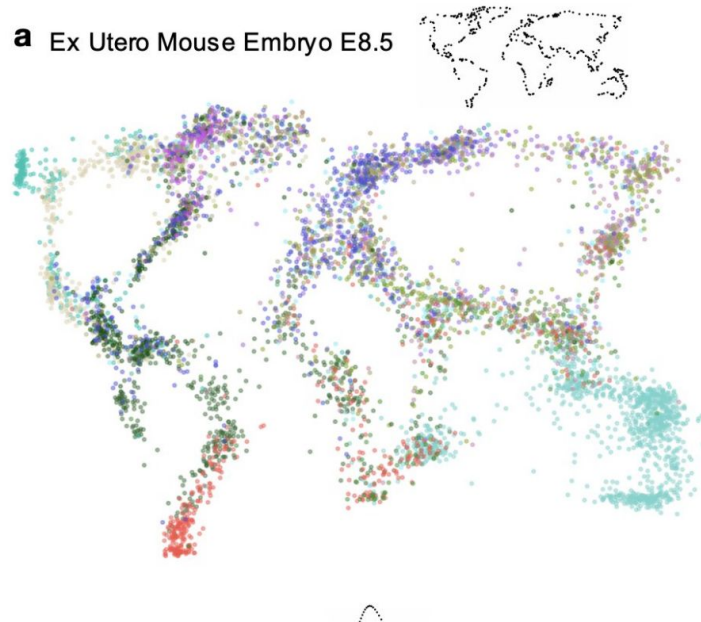


- Point clouds \rightarrow increase radius \rightarrow simplicial complexes \rightarrow topological characteristics



Avoid over-interpreting single plots

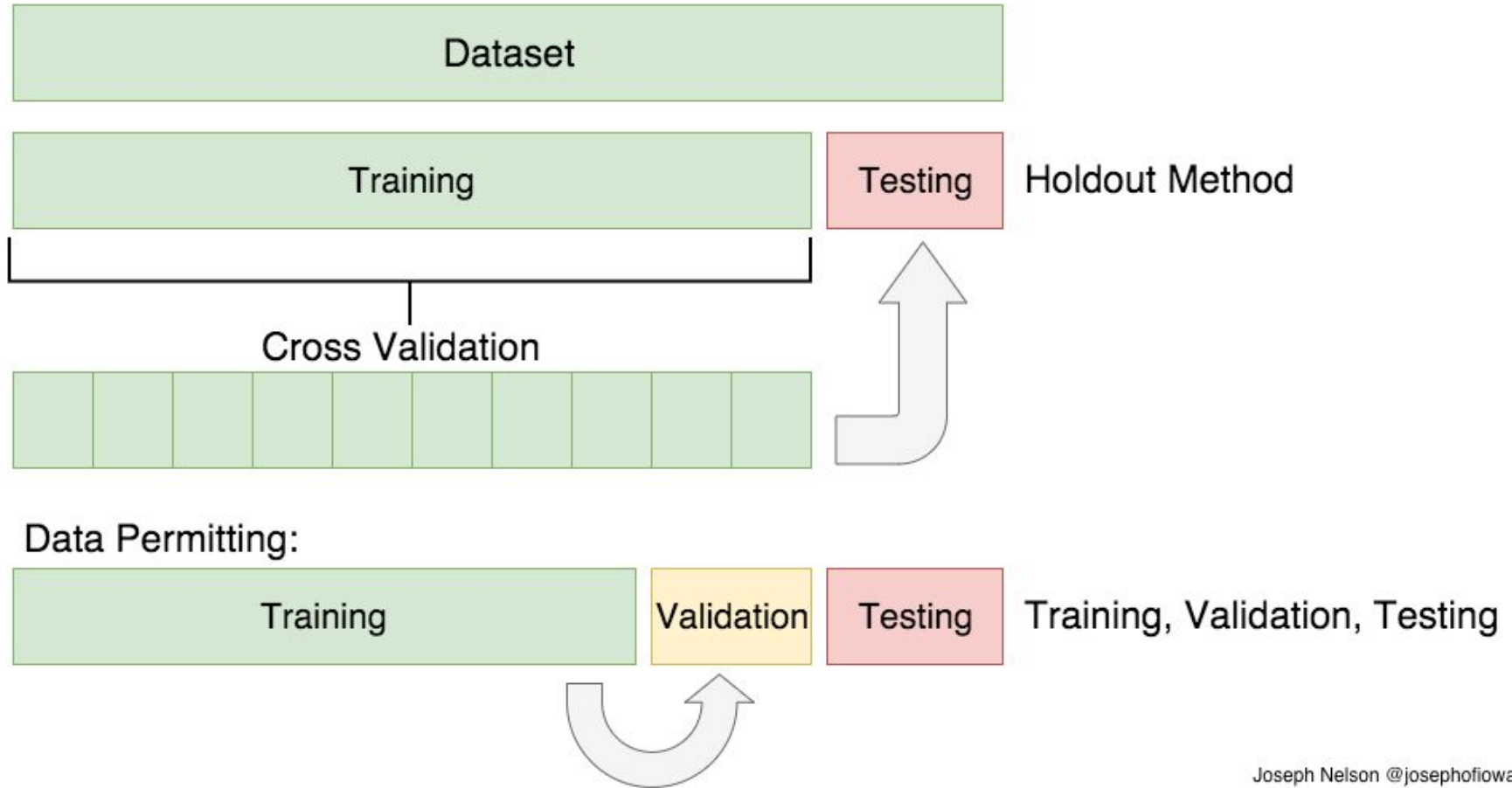
- Sensitive to hyperparameters
- Beware analysing these non-linear projections
- Can contribute to confirmation bias



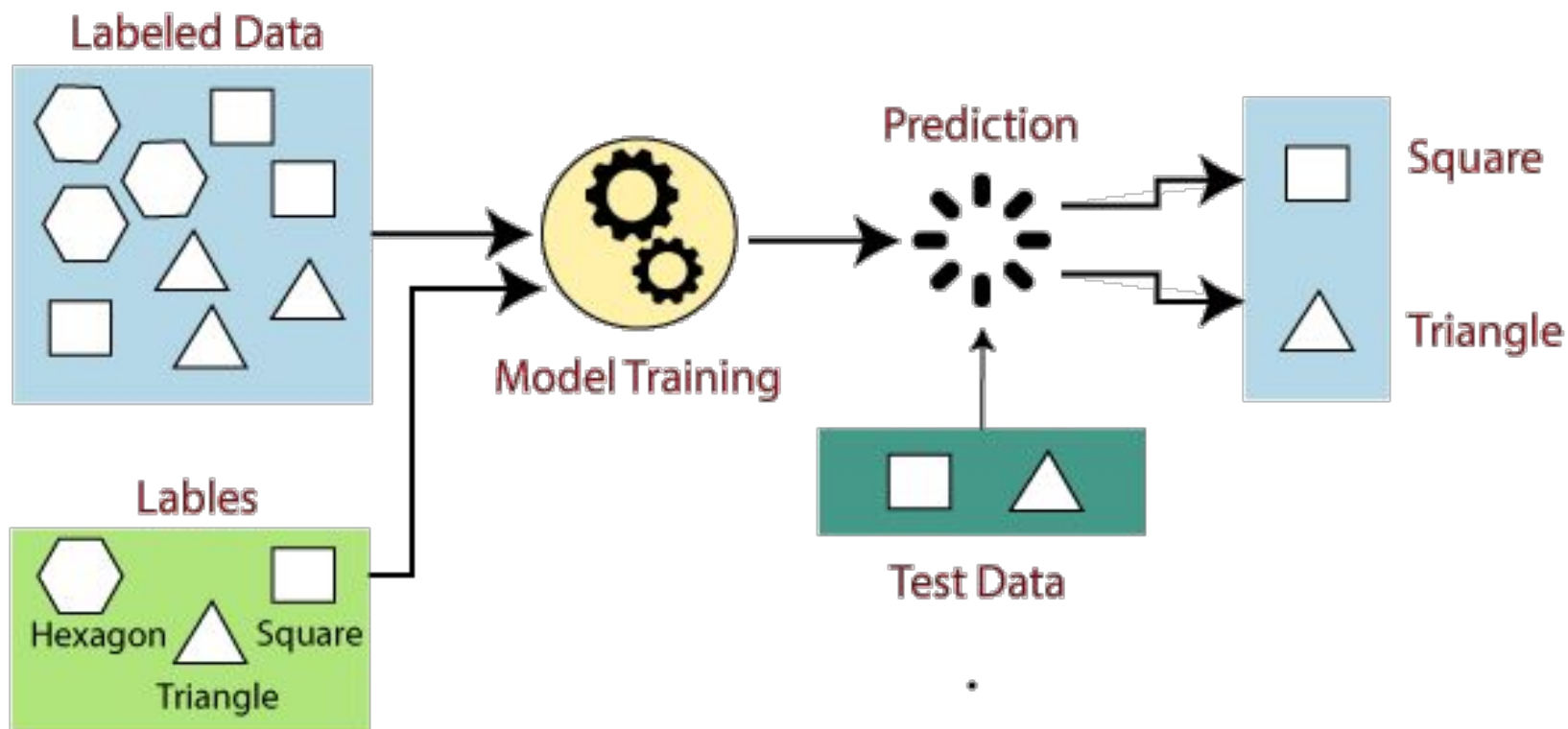
<https://www.biorxiv.org/content/10.1101/2021.08.25.457696v3>

Predicting using tabular data

Overfitting 101: Test-Train Split

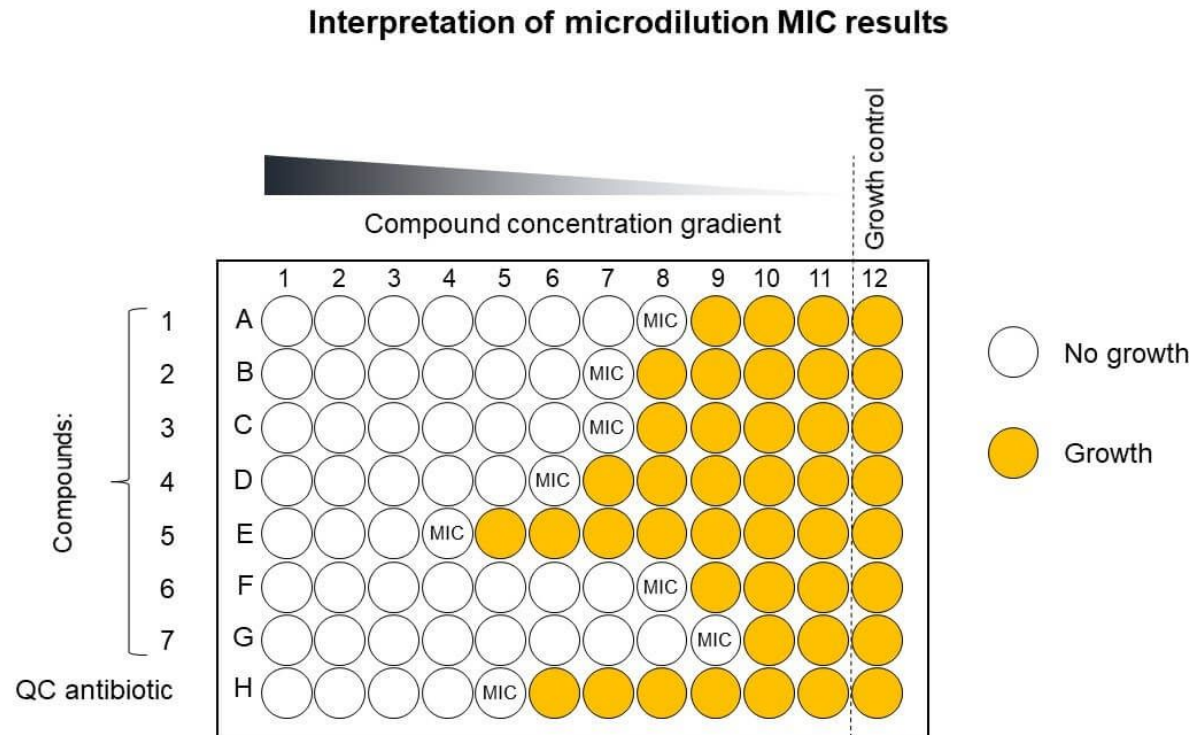


Predicting Labels or Values



Values can be complex: interval prediction

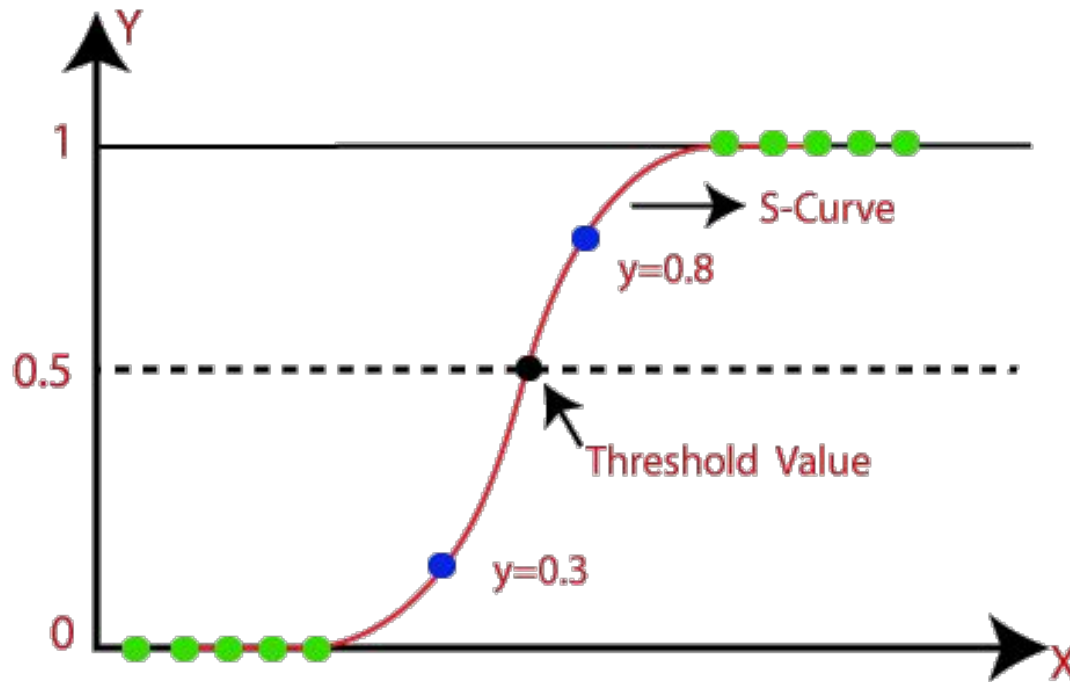
- MIC > highest concentration = right-censored
- MIC < lowest concentration = left-censored
- Serial Dilutions: MIC of x actually $[x/2, 2x]$ = unequal error



Start simple: linear regression

	Common name	Built-in function in R	Equivalent linear model in R	Exact?	The linear model in words	Icon
Simple regression: $\text{lm}(y \sim 1 + x)$	y is independent of x P: One-sample t-test N: Wilcoxon signed-rank	t.test(y) wilcox.test(y)	$\text{lm}(y \sim 1)$ $\text{lm}(\text{signed_rank}(y) \sim 1)$	✓ for N > 14	One number (intercept, i.e., the mean) predicts y . - (Same, but it predicts the <i>signed rank</i> of y .)	
	P: Paired-sample t-test N: Wilcoxon matched pairs	t.test(y1, y2, paired=TRUE) wilcox.test(y1, y2, paired=TRUE)	$\text{lm}(y_2 - y_1 \sim 1)$ $\text{lm}(\text{signed_rank}(y_2 - y_1) \sim 1)$	✓ for N > 14	One intercept predicts the pairwise y₂-y₁ differences. - (Same, but it predicts the <i>signed rank</i> of y₂-y₁ .)	
	y ~ continuous x P: Pearson correlation N: Spearman correlation	cor.test(x, y, method="Pearson") cor.test(x, y, method="Spearman")	$\text{lm}(y \sim 1 + x)$ $\text{lm}(\text{rank}(y) \sim 1 + \text{rank}(x))$	✓ for N > 10	One intercept plus x multiplied by a number (slope) predicts y . - (Same, but with <i>ranked x</i> and y)	
	y ~ discrete x P: Two-sample t-test P: Welch's t-test N: Mann-Whitney U	t.test(y1, y2, var.equal=TRUE) t.test(y1, y2, var.equal=FALSE) wilcox.test(y1, y2)	$\text{lm}(y \sim 1 + G_2)^4$ $\text{gls}(y \sim 1 + G_2, \text{weights}=\dots^4)$ $\text{lm}(\text{signed_rank}(y) \sim 1 + G_2)^4$	✓ ✓ for N > 11	An intercept for group 1 (plus a difference if group 2) predicts y . - (Same, but with one variance <i>per group</i> instead of one common.) - (Same, but it predicts the <i>signed rank</i> of y .)	
Multiple regression: $\text{lm}(y \sim 1 + x_1 + x_2 + \dots)$	P: One-way ANOVA N: Kruskal-Wallis	aov(y ~ group) kruskal.test(y ~ group)	$\text{lm}(y \sim 1 + G_2 + G_3 + \dots + G_N)^4$ $\text{lm}(\text{rank}(y) \sim 1 + G_2 + G_3 + \dots + G_N)^4$	✓ for N > 11	An intercept for group 1 (plus a difference if group ≠ 1) predicts y . - (Same, but it predicts the <i>rank</i> of y .)	
	P: One-way ANCOVA	aov(y ~ group + x)	$\text{lm}(y \sim 1 + G_2 + G_3 + \dots + G_N + x)^4$	✓	- (Same, but plus a slope on x .) <i>Note: this is discrete AND continuous. ANCOVAs are ANOVAs with a continuous x.</i>	
	P: Two-way ANOVA	aov(y ~ group * sex)	$\text{lm}(y \sim 1 + G_2 + G_3 + \dots + G_N + S_2 + S_3 + \dots + S_K + G_2^*S_2 + G_3^*S_3 + \dots + G_N^*S_K)$	✓	Interaction term: changing sex changes the y ~ group parameters. <i>Note: G_{2..K} is an indicator (0 or 1) for each non-intercept levels of the group variable. Similarly for S_{2..K} for sex. The first line (with G) is main effect of group, the second (with S_j) for sex and the third is the group * sex interaction. For two levels (e.g. male/female), line 2 would just be "S₂" and line 3 would be S₂ multiplied with each G_i.</i>	[Coming]
	Counts ~ discrete x N: Chi-square test	chisq.test(groupXsex_table)	Equivalent log-linear model $\text{glm}(y \sim 1 + G_2 + G_3 + \dots + G_N + S_2 + S_3 + \dots + S_K + G_2^*S_2 + G_3^*S_3 + \dots + G_N^*S_K, \text{family}=\dots)^4$	✓	Interaction term: (Same as Two-way ANOVA.) <i>Note: Run glm using the following arguments: glm(model, family=poisson())</i> As linear-model, the Chi-square test is $\log(y_j) = \log(N) + \log(\alpha_j) + \log(\beta_j) + \log(\alpha_j\beta_j)$ where α_j and β_j are proportions. See more info in the accompanying notebook .	Same as Two-way ANOVA
	N: Goodness of fit	chisq.test(y)	$\text{glm}(y \sim 1 + G_2 + G_3 + \dots + G_N, \text{family}=\dots)^4$	✓	(Same as One-way ANOVA and see Chi-Square note.)	1W-ANOVA

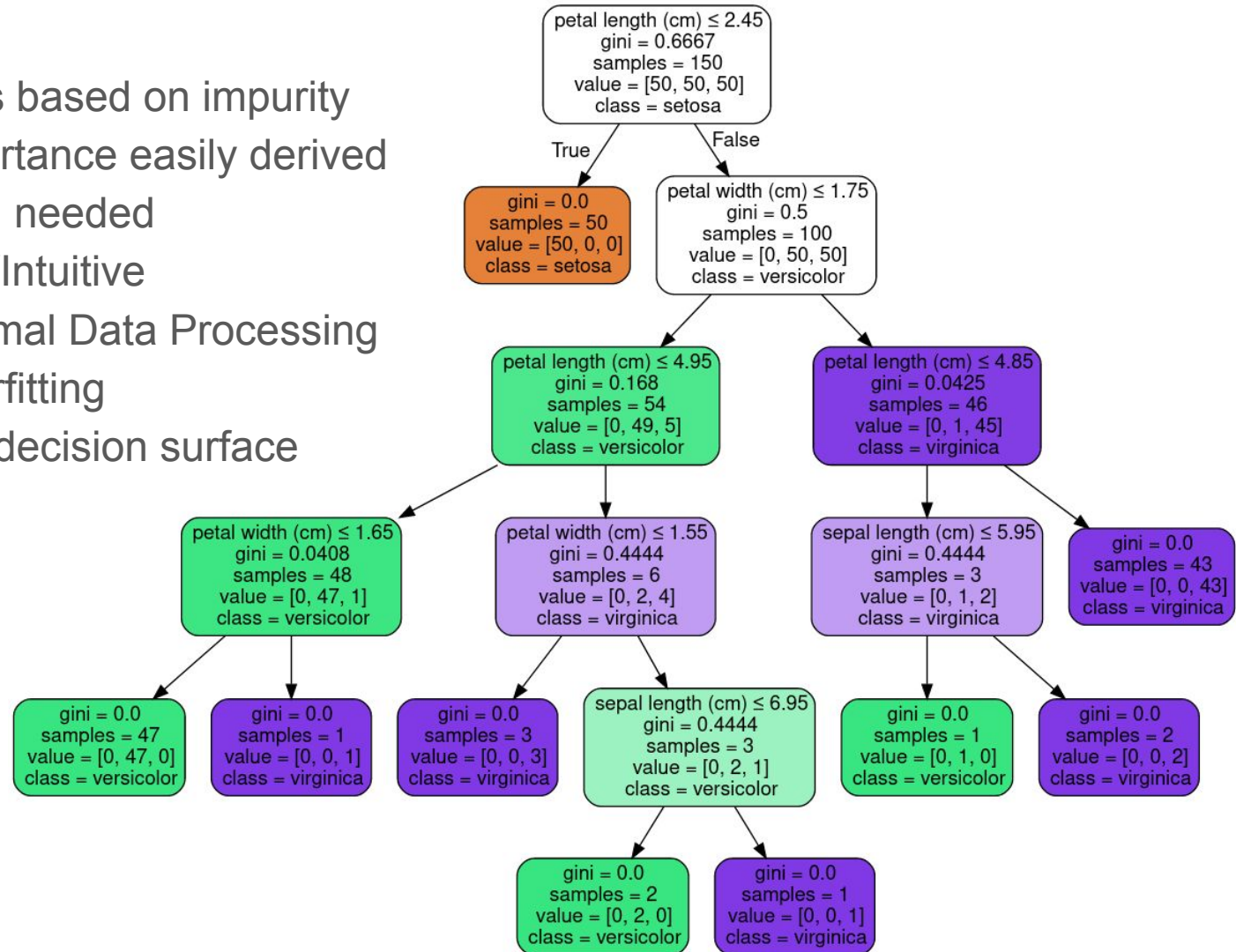
Add a sigmoid for classification



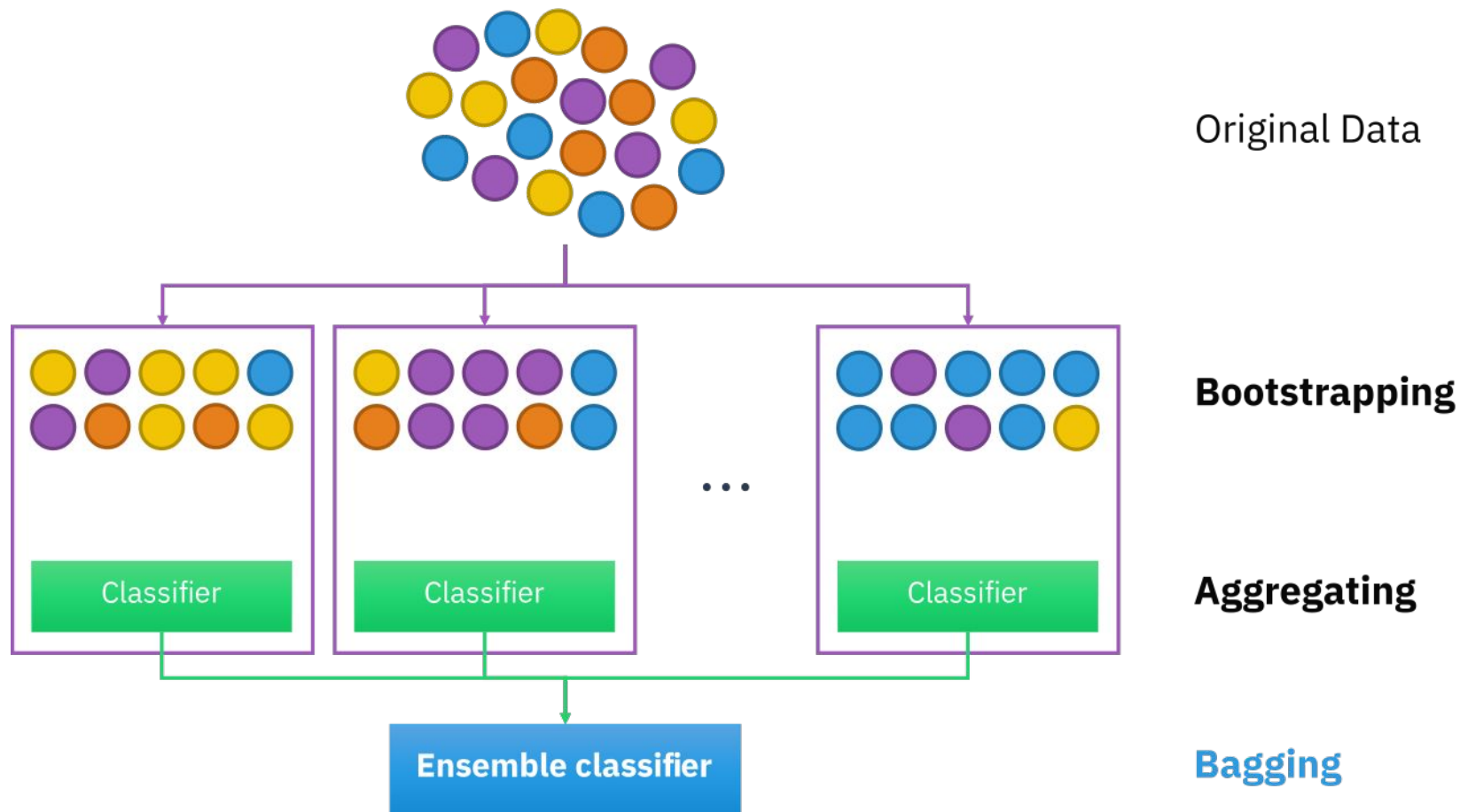
- Crude measure of feature importance (model coefficients)
- Specific feature selection can be a good idea
- Support for regularisation (Lasso/L1 -> sparsity vs Ridge/L2 -> minimal vs ElasticNet -> balance)
- Statistics has developed much better practices for treatment/interpretation of logistic regression

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

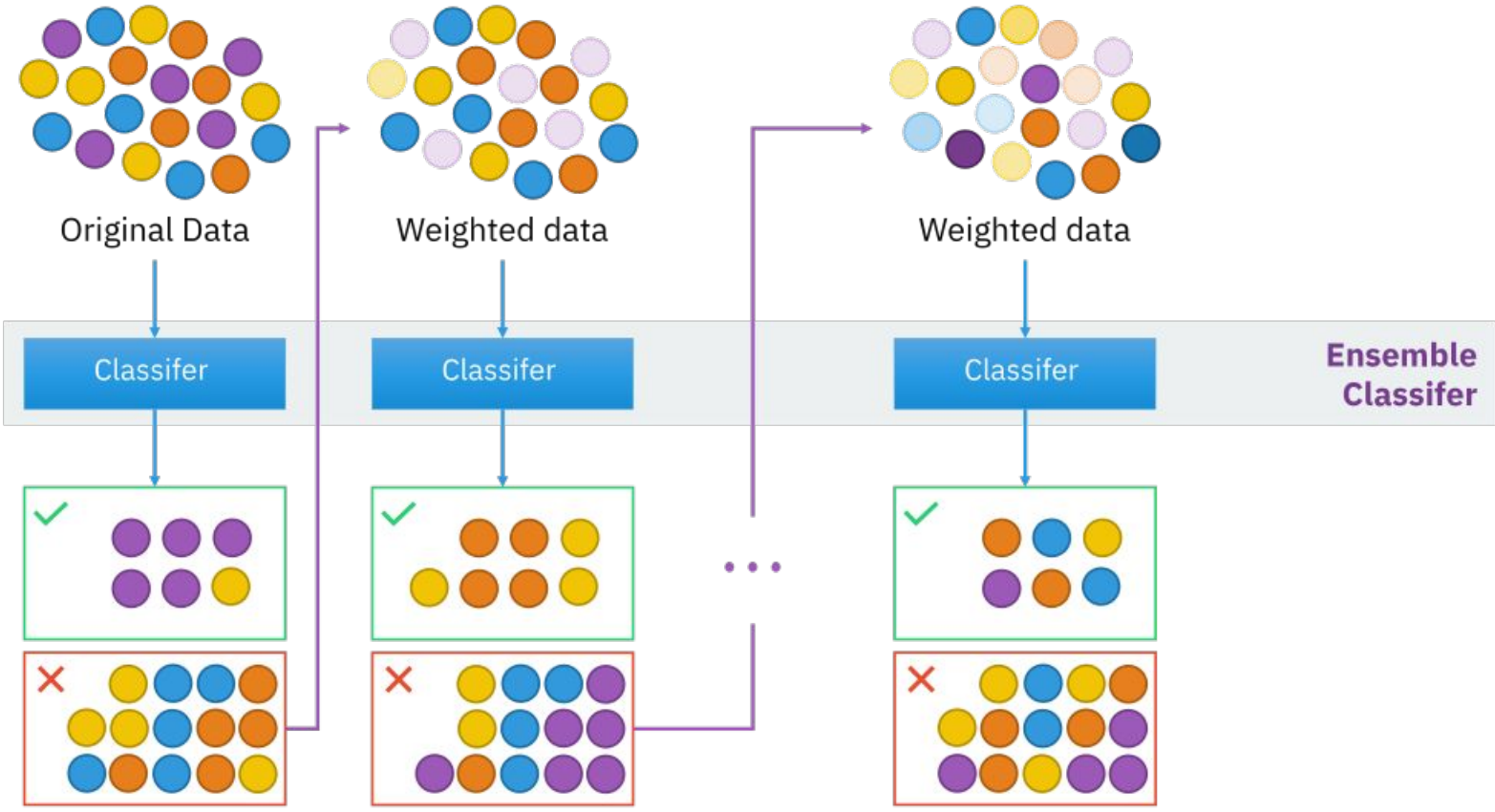


Many Decision Trees: Bagging



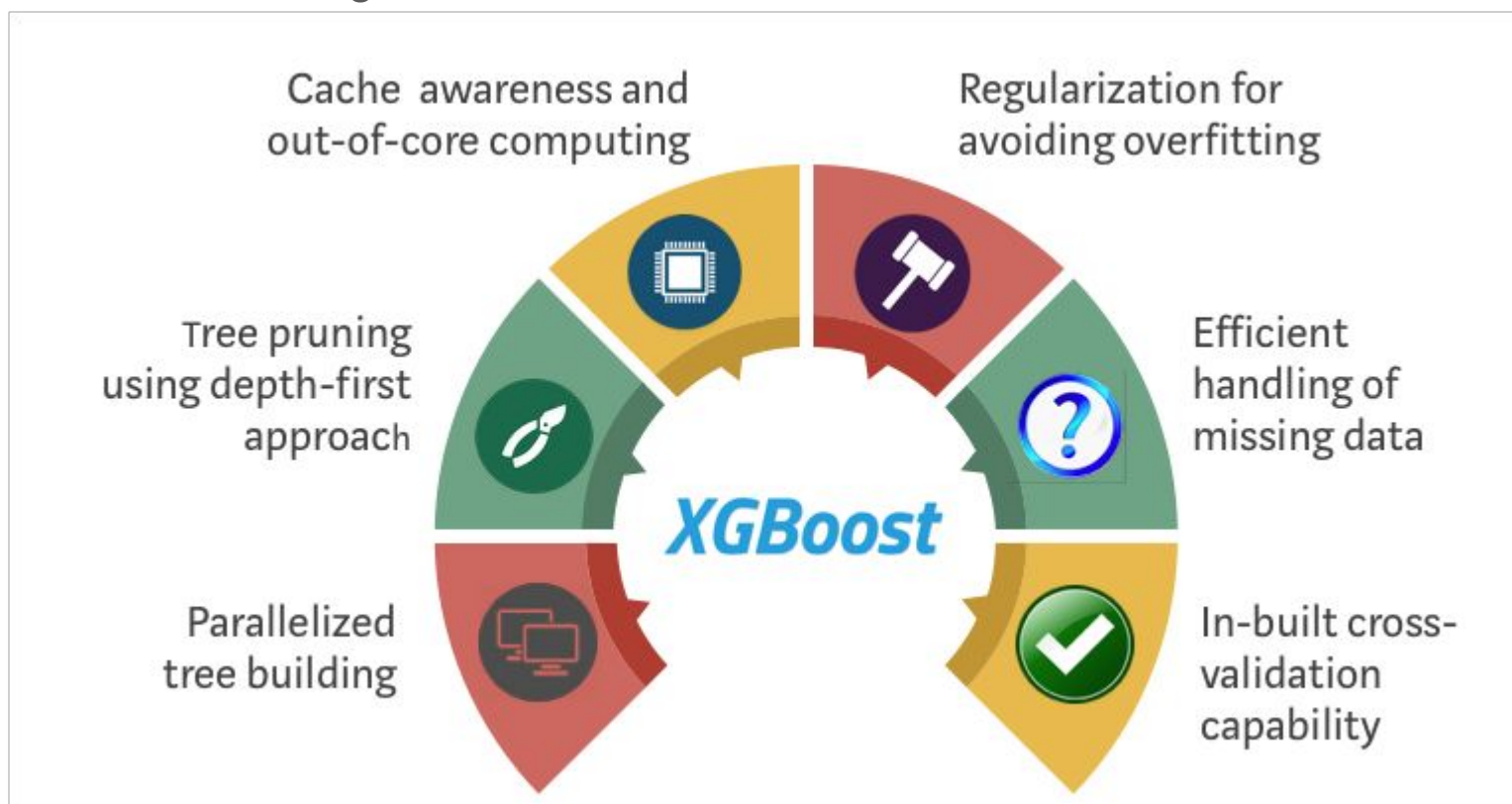
Random Forest: Bagging + Random Subset Per Split
Feature Importance: Average impurity decrease

Boosting: AdaBoost



Gradient Boosting: XGBoost

- Fix on pseudo-residuals instead of weights
- Use stochastic gradient descent



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