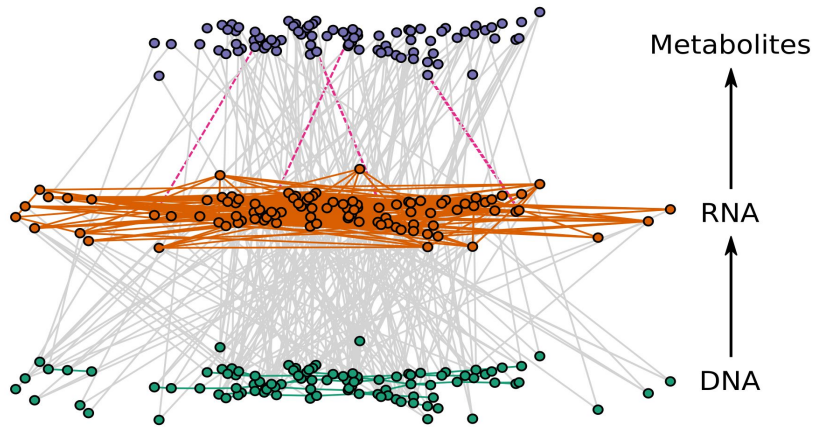


Lecture 0: Introduction to Applied Research in Health Data Science

CSCI6410/4148 & EPAH6410

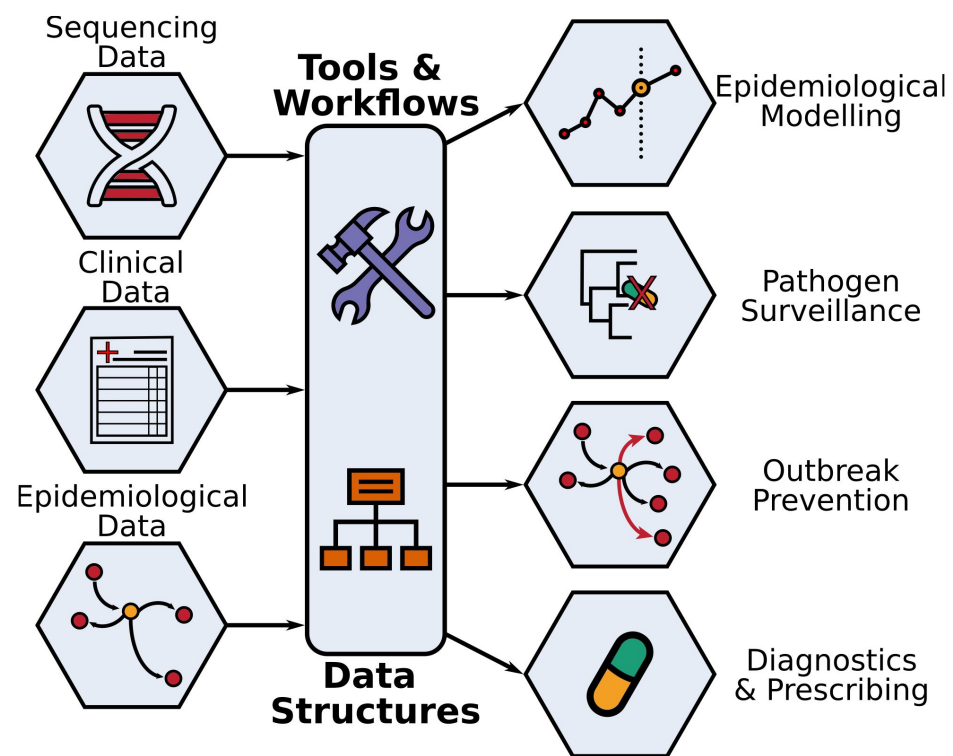
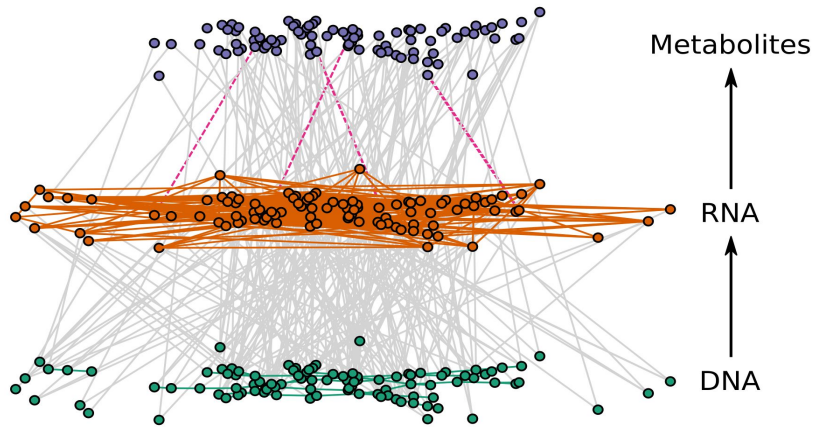
Finlay Maguire (finlay.maguire@dal.ca)
TA: Mehrana Calagari (m.calagari@dal.ca)

Why am I teaching this course?



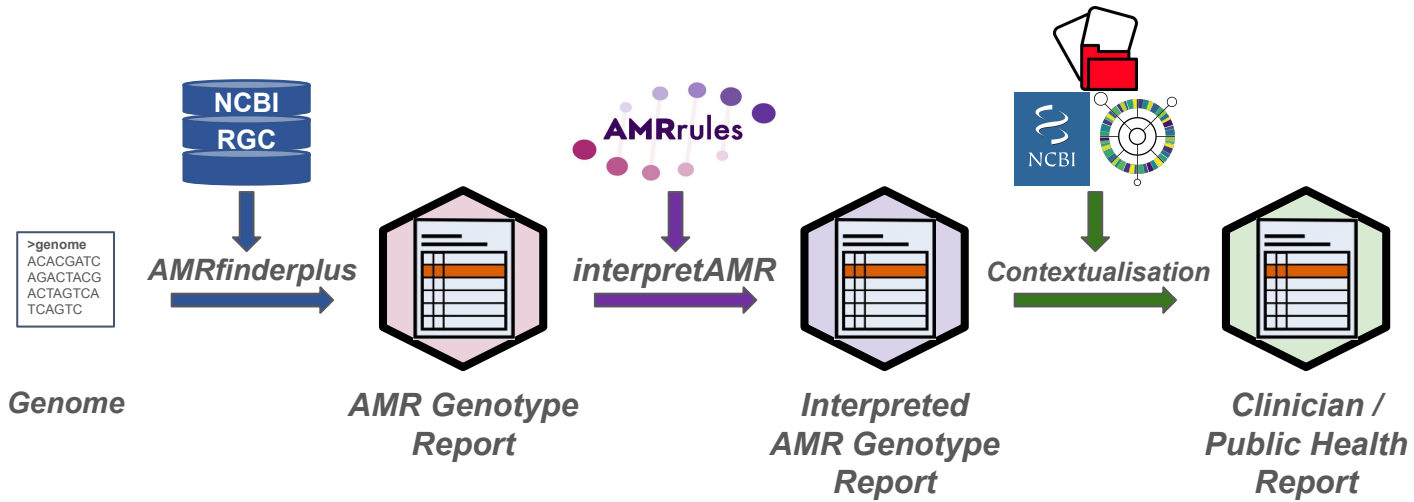
- **PhD (Bioinformatics):** using large noisy datasets to understand how microbial systems and mechanisms evolve.

Why am I teaching this course?

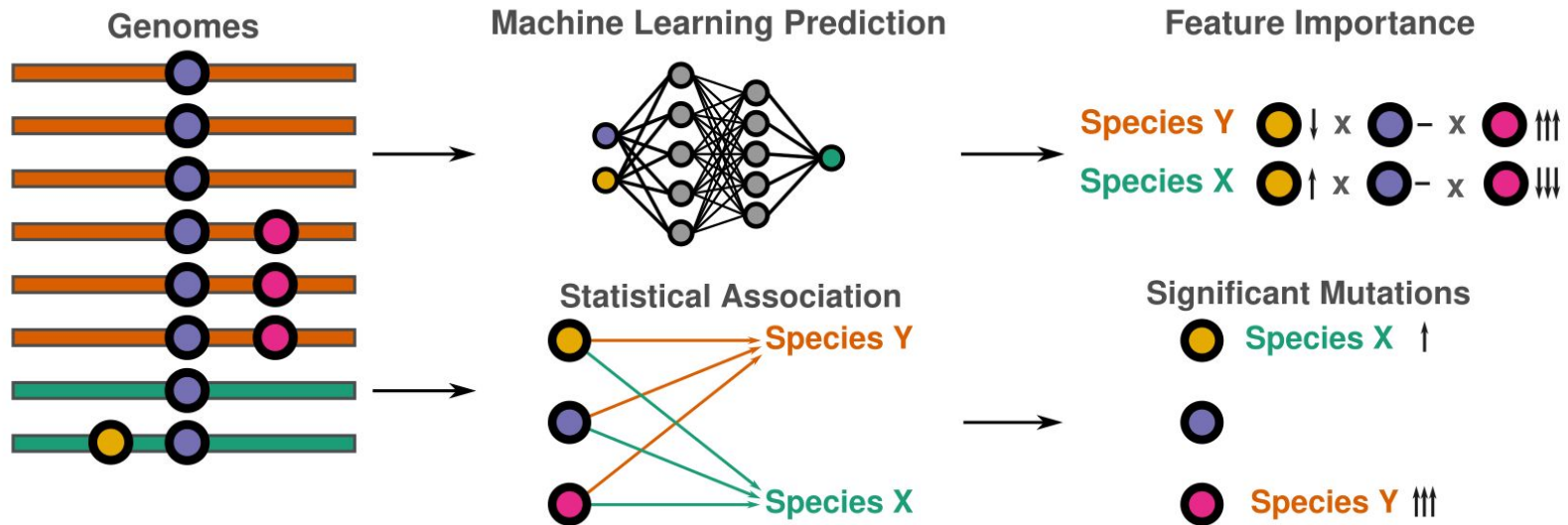
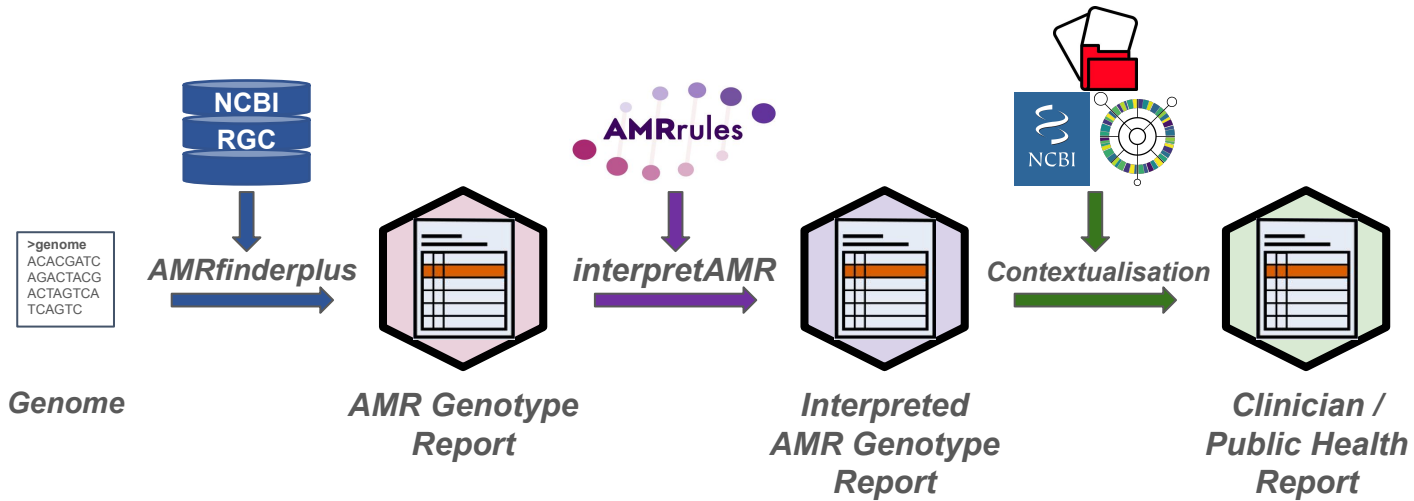


- **PhD (Bioinformatics)**: using large noisy datasets to understand how microbial systems and mechanisms evolve.
- **Research Group (Genomic Epidemiology)**: using large noisy datasets to better diagnose, track and predict infectious diseases.

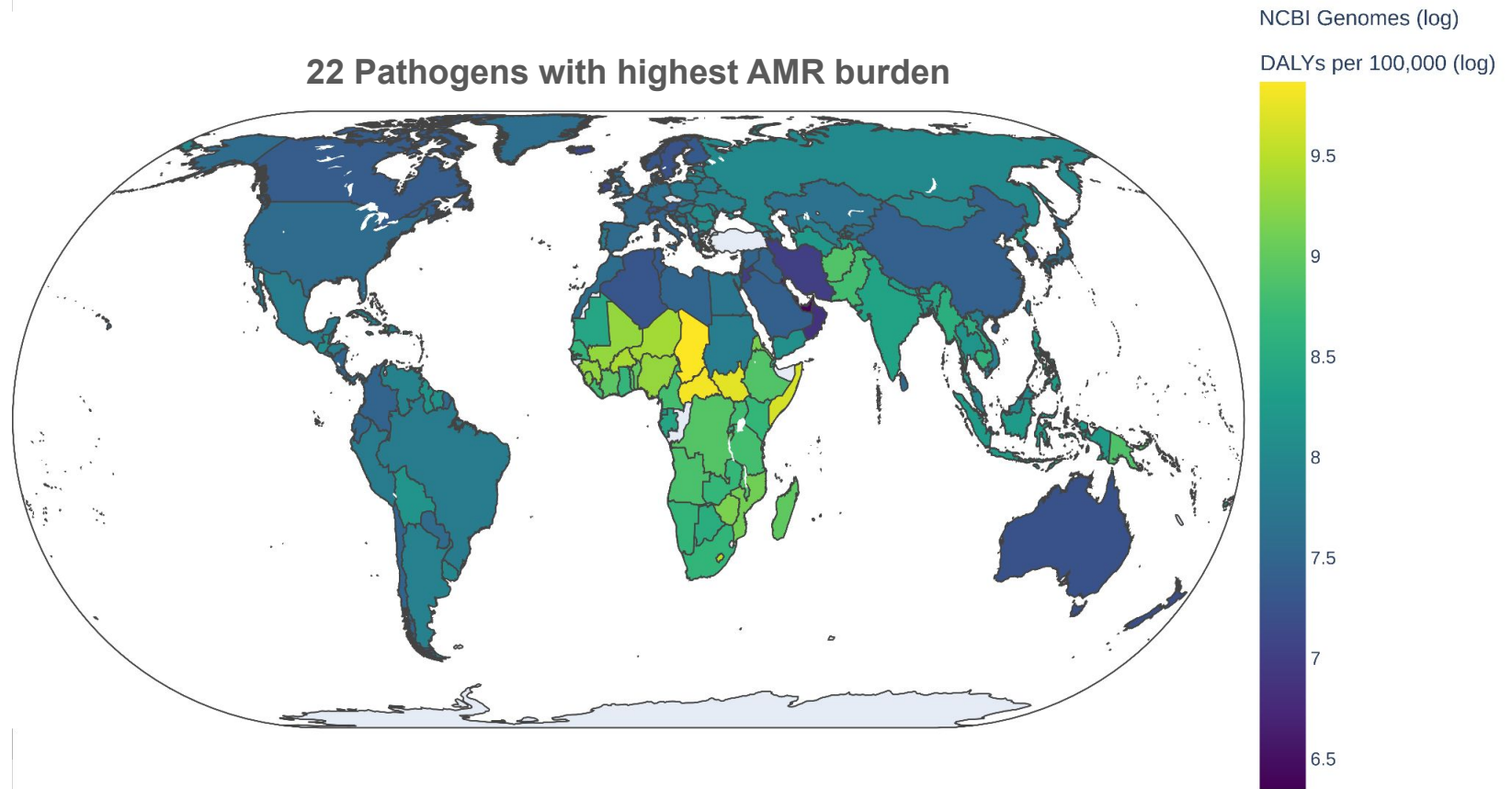
Why am I teaching this course?



Why am I teaching this course?

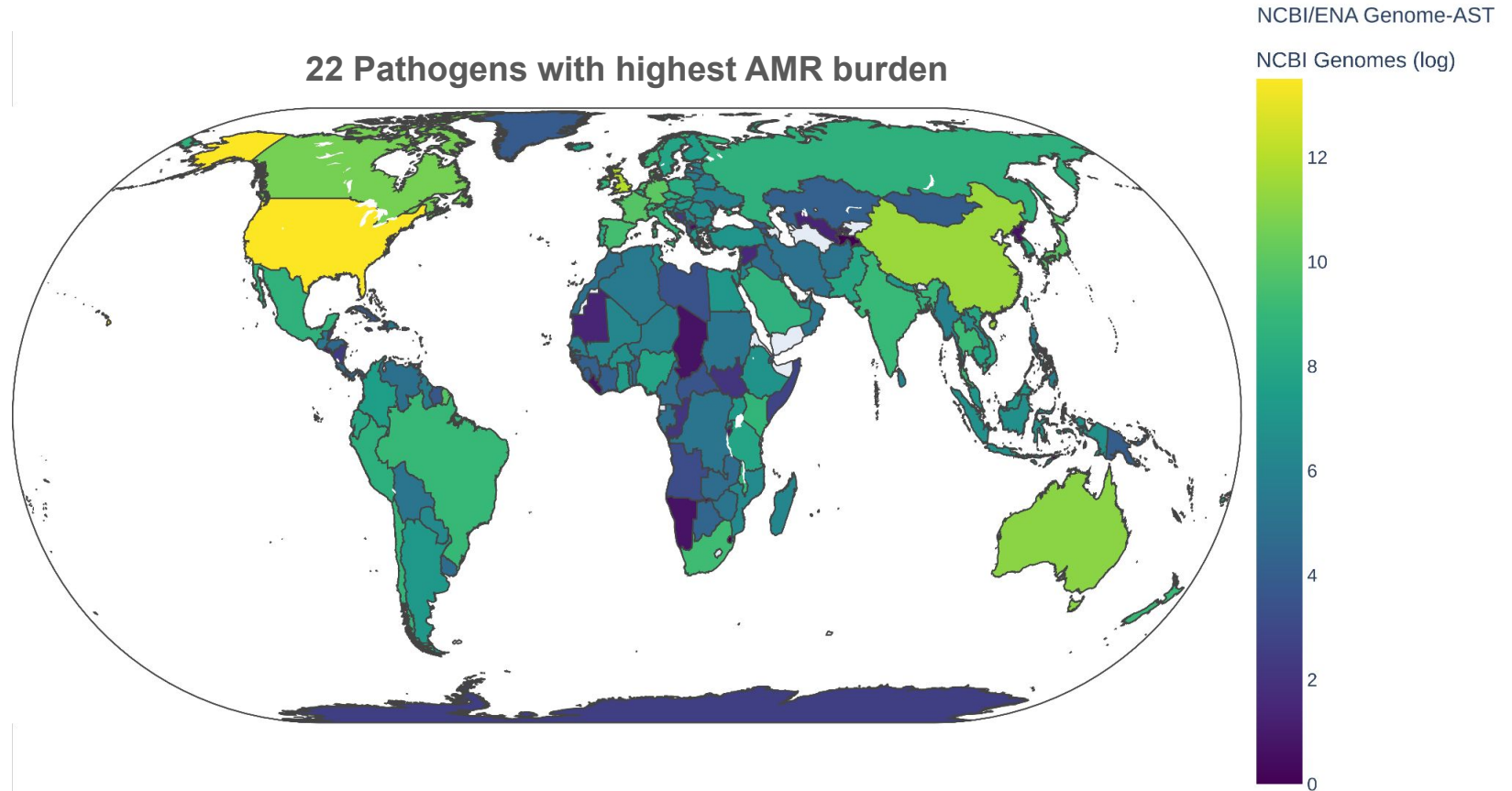


Why am I teaching this course?



K. pneumoniae: 23,379,605 DALYs across 100 highest burden countries

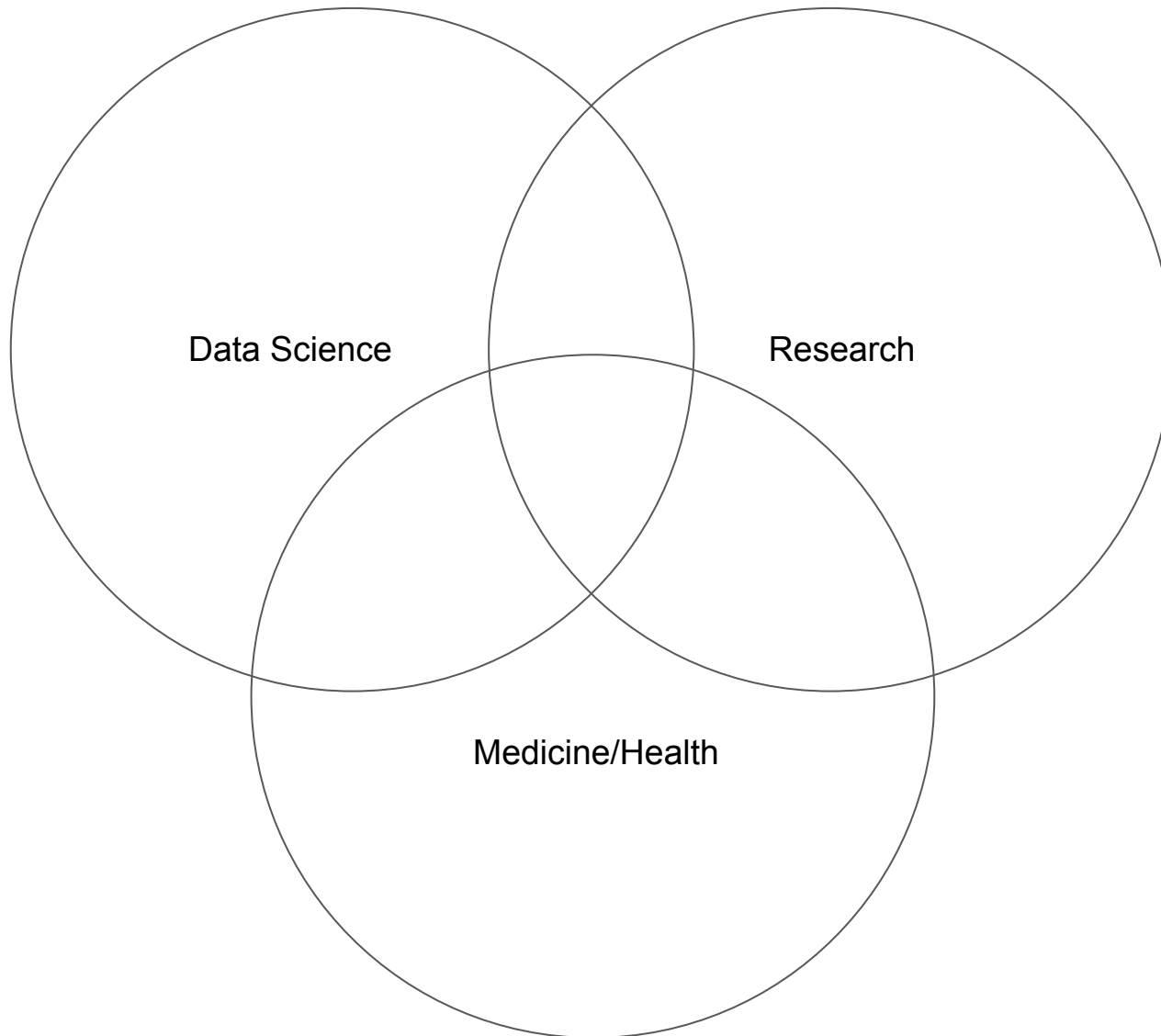
Why am I teaching this course?



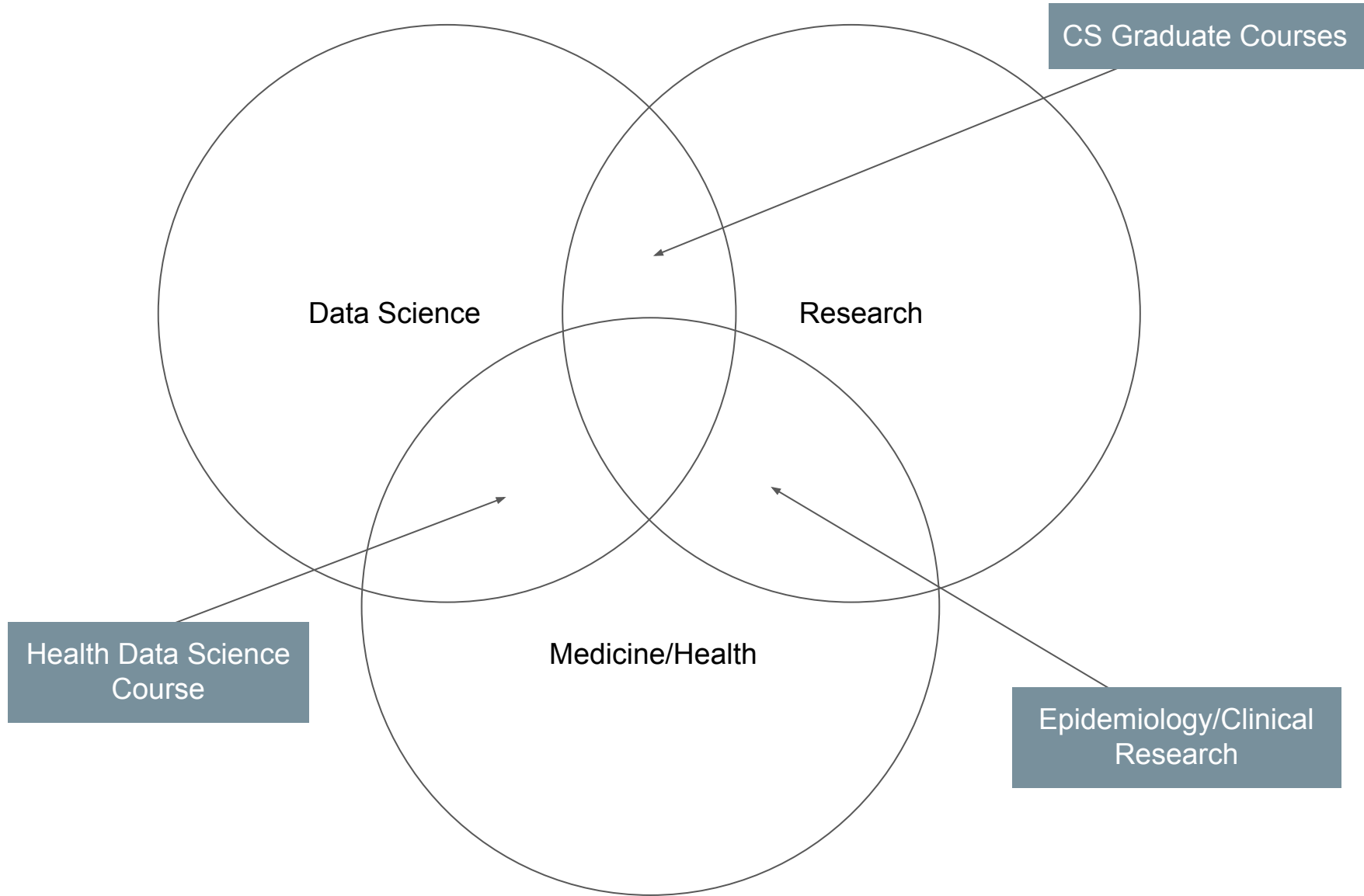
K. pneumoniae: 8,689/89,056 genomes from 100 highest burden countries

Overview of course

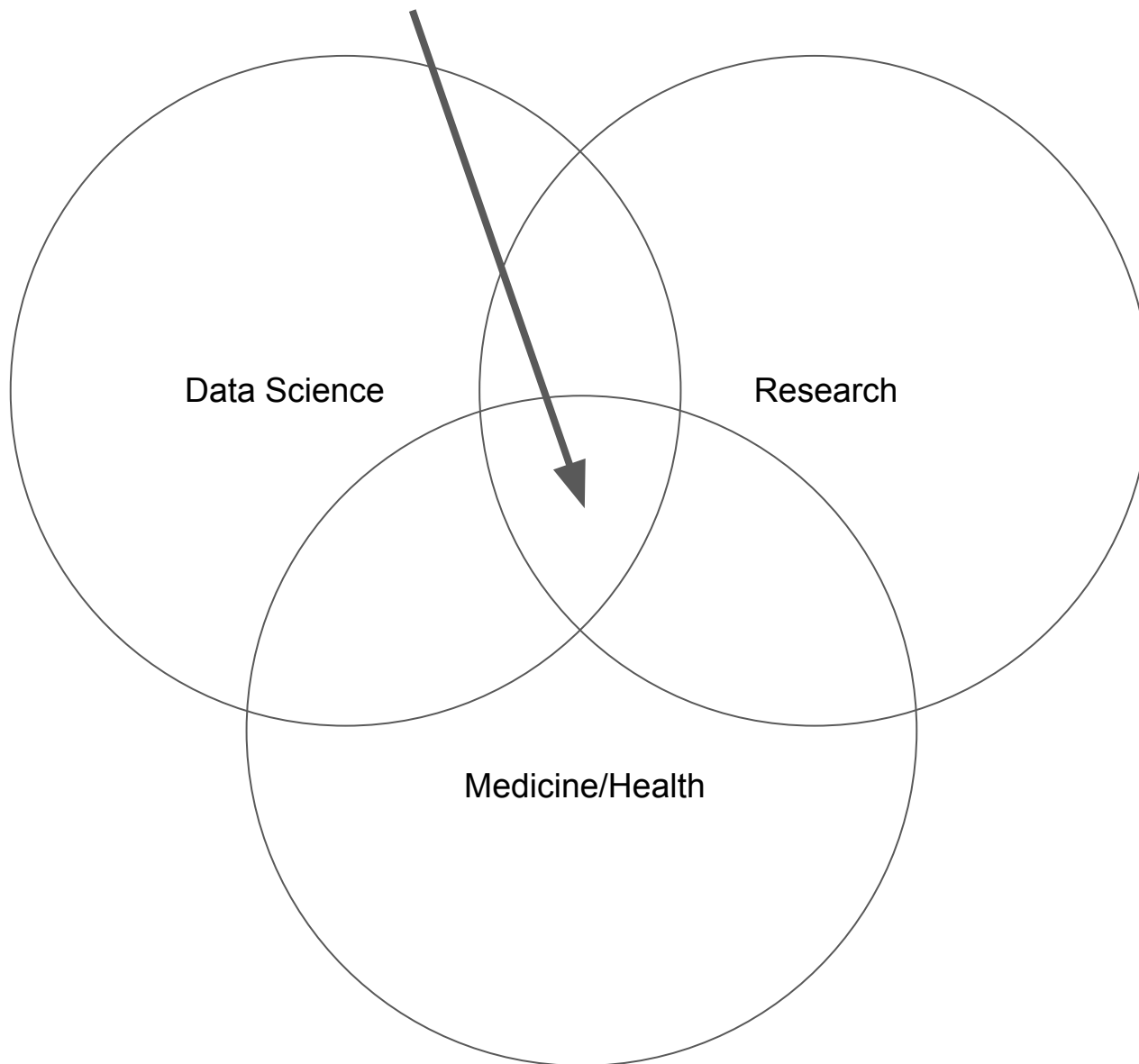
Applied Research in Health Data Science



Applied Research in Health Data Science



Applied Research in Health Data Science



Learning Outcomes

1. Understand the **4 principal sources and data types** of medical data:
 - a. longitudinal databases (tabular)
 - b. electronic medical records (structured, semi-structured, and unstructured text)
 - c. radiological imaging (image)
 - d. physiological (signal and time-series).

Learning Outcomes

1. Understand the **4 principal sources and data types** of medical data:
 - a. longitudinal databases (tabular)
 - b. electronic medical records (structured, semi-structured, and unstructured text)
 - c. radiological imaging (image)
 - d. physiological (signal and time-series).
2. Identify and apply **appropriate type of method** to the analysis of each data type

Learning Outcomes

1. Understand the **4 principal sources and data types** of medical data:
 - a. longitudinal databases (tabular)
 - b. electronic medical records (structured, semi-structured, and unstructured text)
 - c. radiological imaging (image)
 - d. physiological (signal and time-series).
2. Identify and apply **appropriate type of method** to the analysis of each data type
3. Gain the technical skills necessary for effective health data science research including **data management, reproducibility**, and version control.

Learning Outcomes

1. Understand the **4 principal sources and data types** of medical data:
 - a. longitudinal databases (tabular)
 - b. electronic medical records (structured, semi-structured, and unstructured text)
 - c. radiological imaging (image)
 - d. physiological (signal and time-series).
2. Identify and apply **appropriate type of method** to the analysis of each data type
3. Gain the technical skills necessary for effective health data science research including **data management, reproducibility**, and version control.
4. Understand the key **collaborative, legal, ethical, and knowledge translation** concepts required in interdisciplinary health data science research.

Learning Outcomes

1. Understand the **4 principal sources and data types** of medical data:
 - a. longitudinal databases (tabular)
 - b. electronic medical records (structured, semi-structured, and unstructured text)
 - c. radiological imaging (image)
 - d. physiological (signal and time-series).
2. Identify and apply **appropriate type of method** to the analysis of each data type
3. Gain the technical skills necessary for effective health data science research including **data management, reproducibility**, and version control.
4. Understand the key **collaborative, legal, ethical, and knowledge translation** concepts required in interdisciplinary health data science research.
5. Critically **appraise research literature** in health data science.

Learning Outcomes

1. Understand the **4 principal sources and data types** of medical data:
 - a. longitudinal databases (tabular)
 - b. electronic medical records (structured, semi-structured, and unstructured text)
 - c. radiological imaging (image)
 - d. physiological (signal and time-series).
2. Identify and apply **appropriate type of method** to the analysis of each data type
3. Gain the technical skills necessary for effective health data science research including **data management, reproducibility**, and version control.
4. Understand the key **collaborative, legal, ethical, and knowledge translation** concepts required in interdisciplinary health data science research.
5. Critically **appraise research literature** in health data science.
6. Combine these skills to develop high-quality collaborative health data science **research proposals**

What is not covered in this course

- **Breadth/depth** of each data science method: *each could be multiple graduate CS courses*

What is not covered in this course

- **Breadth/depth** of each data science method: *each could be multiple graduate CS courses*
- **Breadth/depth** of medical research: *again could be a whole PhD program*

What is not covered in this course

- **Breadth/depth** of each data science method: *each could be multiple graduate CS courses*
- **Breadth/depth** of medical research: *again could be a whole PhD program*
- True **messiness** of real data: *provide tools but experience is invaluable*

What is not covered in this course

- **Breadth/depth** of each data science method: *each could be multiple graduate CS courses*
- **Breadth/depth** of medical research: *again could be a whole PhD program*
- True **messiness** of real data: *provide tools but experience is invaluable*
- Some important forms of medical data (e.g., genomics): *see CSCI4181/6810, EPAH6052, come speak to me if interested in this specifically.*

Course Structure

Overview of data types & analysis methods:

- **Lectures** (Monday/Wednesday)

Course Structure

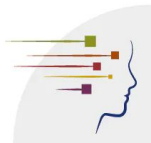
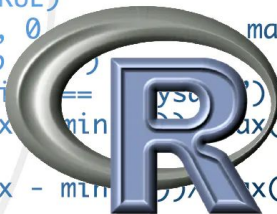
Overview of data types & analysis methods:

- **Lectures** (Monday/Wednesday)
- **Practical Exercises** (Friday/Monday)

Assessment: Submission of Practical Exercise Due the day before **following practical** (7.5% x 4)

(CSCI4148: drop lowest scoring assignment)

```
dens <- density(data, n = npts)
dx <- dens$x
dy <- dens$y
if(add == TRUE)
  plot(0., 0., main,
       ylab,
       if(orientati == "y")
         dx2 <- (dx - min(dx)) / dx(dx)
         x[1.]
         dy2 <- (dx - min(dx)) / dx(dy)
         y[1.]
         seqbelow <- rep(y[1.], length(dx))
         if(Fill == T)
           confshade(dx2, seqbelow, dy2
```



<https://www.coursera.org/learn/r-programming>

Course Structure

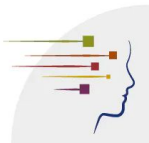
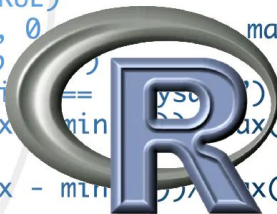
Overview of data types & analysis methods:

- **Lectures** (Monday/Wednesday)
- **Practical Exercises** (Friday/Monday)

Assessment: Submission of Practical Exercise Due the day before **following practical** (7.5% x 4)

(CSCI4148: drop lowest scoring assignment)

```
dens <- density(data, n = npts)
dx <- dens$x
dy <- dens$y
if(add == TRUE)
  plot(0., 0., main,
       ylab,
       if(orientati == "y")
         dx2 <- (dx - min(dx)) / dx(dx)
         x[1.]
         dy2 <- (dx - min(dx)) / dx(dy)
         y[1.]
         seqbelow <- rep(y[1.], length(dx))
         if(Fill == T)
           confshade(dx2, seqbelow, dy2
```



<https://www.coursera.org/learn/r-programming>

Research in health data science:

- **Journal Club** (Wednesday/Friday)

2 papers per week, randomly assigned rota for leading discussion of paper with rest of class.

Assessment:

Paper presentation (20%)

Participation in discussion (10%)

Course Structure

Overview of data types & analysis methods:

- **Lectures** (Monday/Wednesday)
- **Practical Exercises** (Friday/Monday)

Assessment: Submission of Practical Exercise Due the day before **following practical** (7.5% x 4)

(CSCI4148: drop lowest scoring assignment)

```
dens <- density(data, n = npts)
dx <- dens$x
dy <- dens$y
if(add == TRUE)
  plot(0., 0., main,
       ylab,
       if(orientati == 'y')
         dx2 <- (dx - min(dx)) / dx(dx)
         x[1.]
         dy2 <- (dy - min(dy)) / dy(dy)
         y[1.]
         seqbelow <- rep(y[1.], length(dx))
         if(Fill == T)
           confshade(dx2, seqbelow, dy2
```



<https://www.coursera.org/learn/r-programming>

Research in health data science:

- **Journal Club** (Wednesday/Friday)

2 papers per week, randomly assigned rota for leading discussion of paper with rest of class.

Assessment:

Paper presentation (20%)

Participation in discussion (10%)

Development of a research proposal:

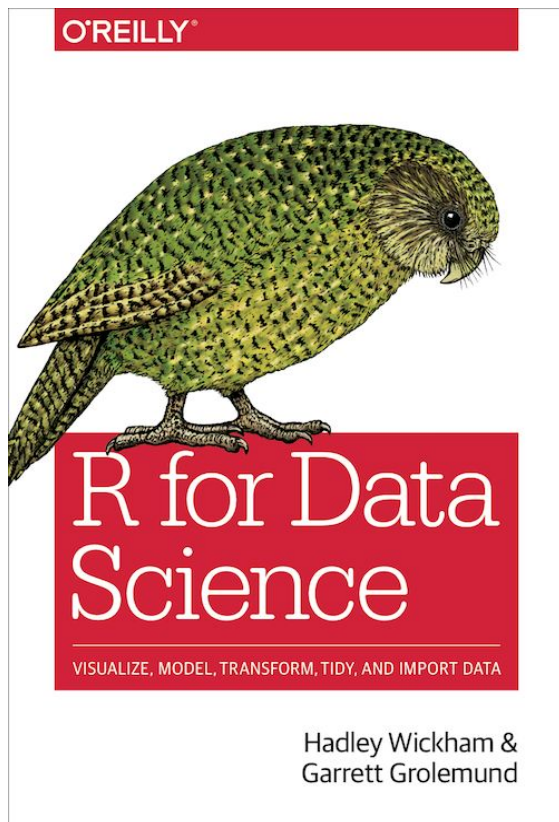
- **Class** (Wednesday/Friday)

Assessment:

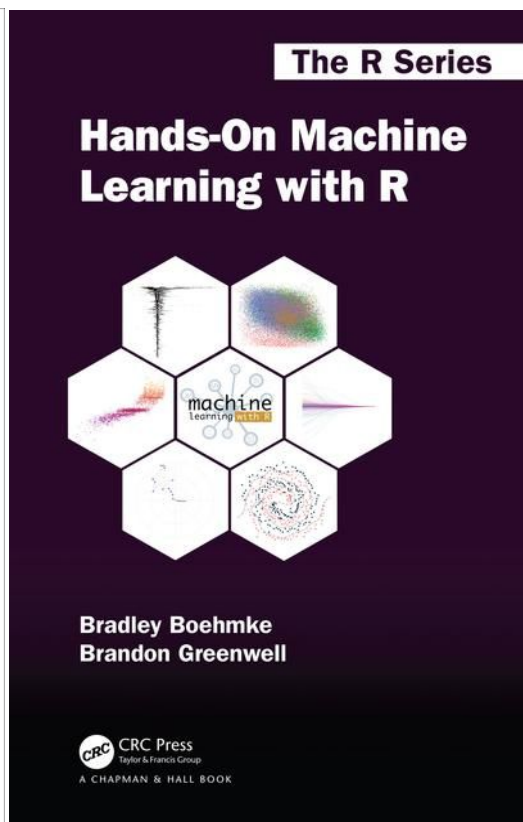
Presentation **last full week of class** (20%)

Submitted **final day of class** (20%)

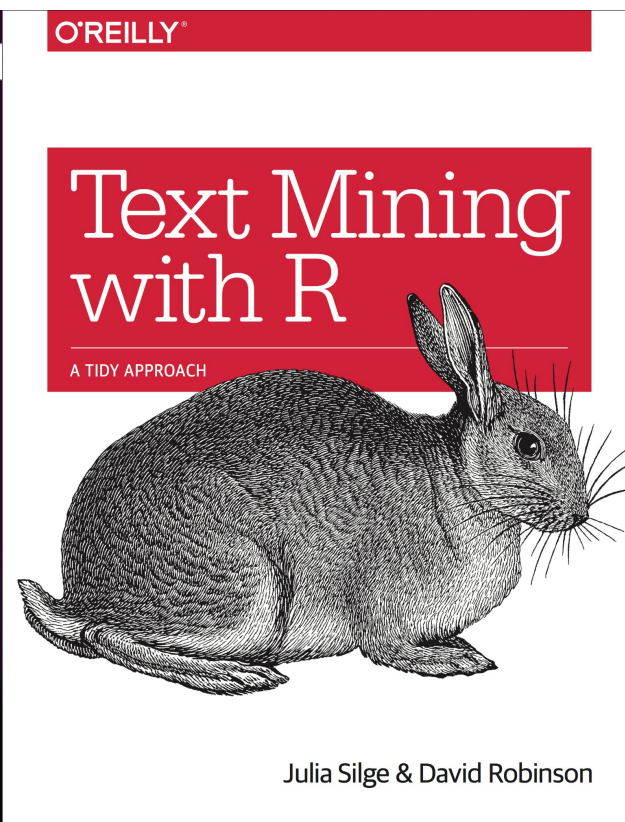
Course Materials



<https://r4ds.had.co.nz/>



<https://bradleyboehmke.github.io/HOML/>



<https://www.tidytextmining.com/>

Course Website




Dalhousie University


CSCI6410/CSCI4148/EPAH6410: Applied Research in Health Data Science

Summer 2025-2026

 HOME

 SCHEDULE

 LECTURES

 PRACTICALS

 PROPOSAL

 LITERATURE

CSCI6410/CSCI4148/EPAH6410: Applied Research in Health Data Science / Summer
2025-2026

https://maguire-lab.github.io/health_data_science_research_2026/

Course Website

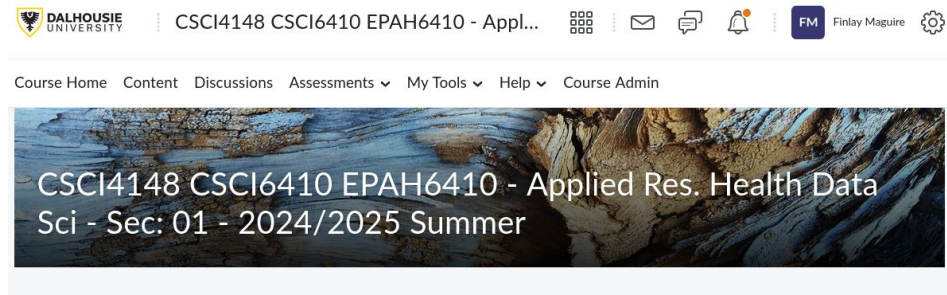








Dalhousie University
CSCI6410/CSCI4148/EPAH6410: Applied Research in Health Data Science
Summer 2025-2026

[HOME](#) [SCHEDULE](#) [LECTURES](#) [PRACTICALS](#) [PROPOSAL](#) [LITERATURE](#)

CSCI6410/CSCI4148/EPAH6410: Applied Research in Health Data Science / Summer 2025-2026

https://maguire-lab.github.io/health_data_science_research_2026/



DALHOUSIE UNIVERSITY | CSCI4148 CSCI6410 EPAH6410 - Appl...      Finlay Maguire 

Course Home Content Discussions Assessments ▾ My Tools ▾ Help ▾ Course Admin

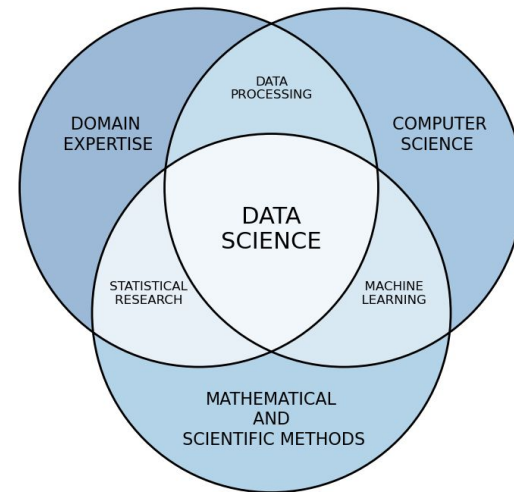
CSCI4148 CSCI6410 EPAH6410 - Applied Res. Health Data Sci - Sec: 01 - 2024/2025 Summer

Grades/Submissions:

<https://dal.brightspace.com/d2l/home/385844>

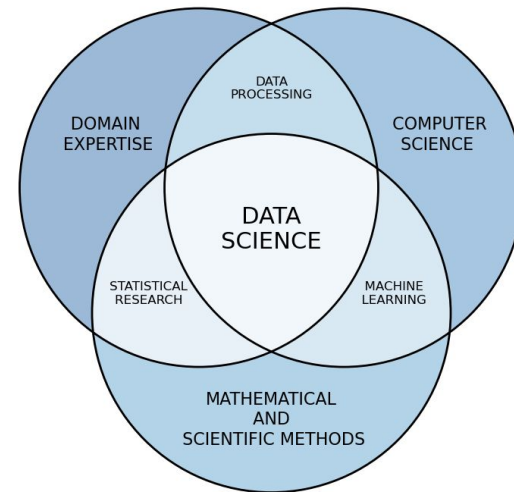
What is ~~health~~ data science?

Data Science: *Data-intensive interdisciplinary approaches to understand and predict with secondary/live data*



Data Science: *Data-intensive interdisciplinary approaches to understand and predict with secondary/live data*

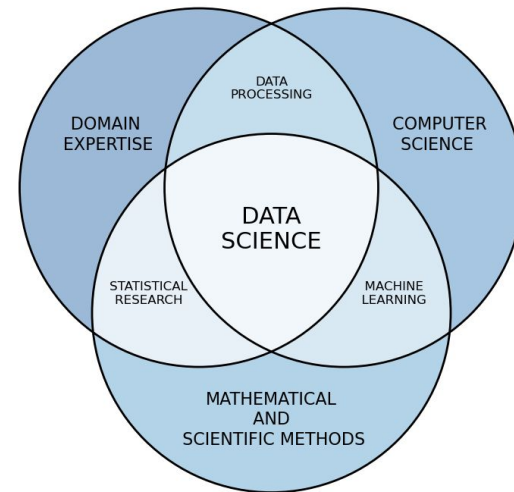
A range of partial and totally overlapping terms:



Data Science: *Data-intensive interdisciplinary approaches to understand and predict with secondary/live data*

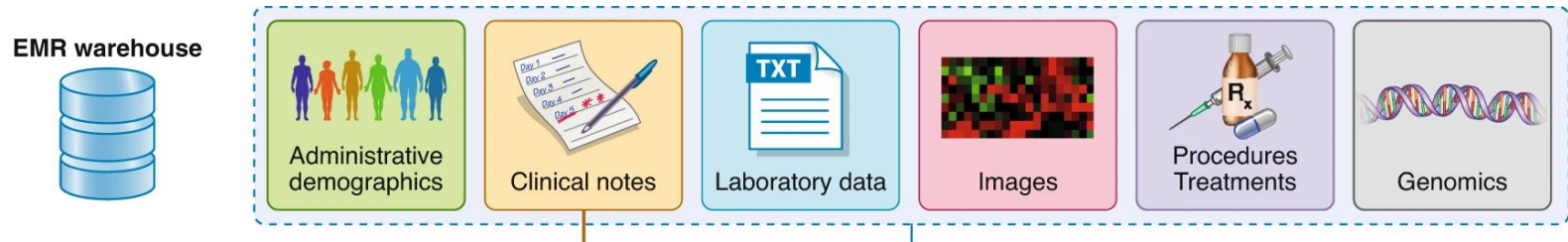
A range of partial and totally overlapping terms:

- Data Analytics
- Data Engineering
- Data Mining
- {Health,Bio,Medical}Informatics
- Database Analysis
- Business Intelligence
- Epidemiology
- Statistics
- **Machine Learning**
- Pattern Recognition
- Predictive Analytics
- Quantitative Researcher
- Scientist
- Analyst
- Algorithmic Modeling



OK, what is **Health** Data Science?

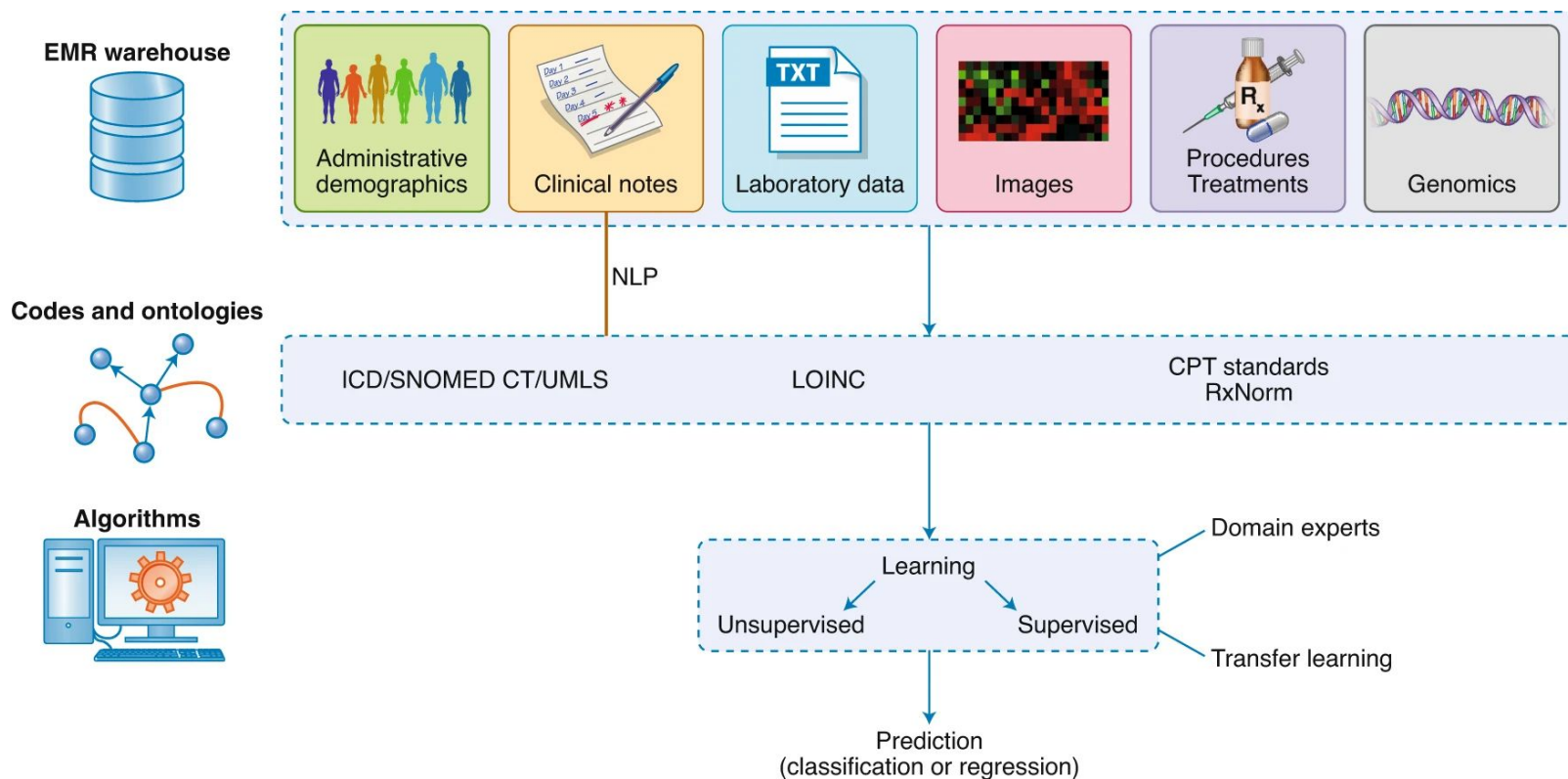
Data Science applied to Health Data



Why “health data” instead of “medical data”:

medical (being defined as MD-related) \subset health - **contentious**

Data Science applied to Health Data



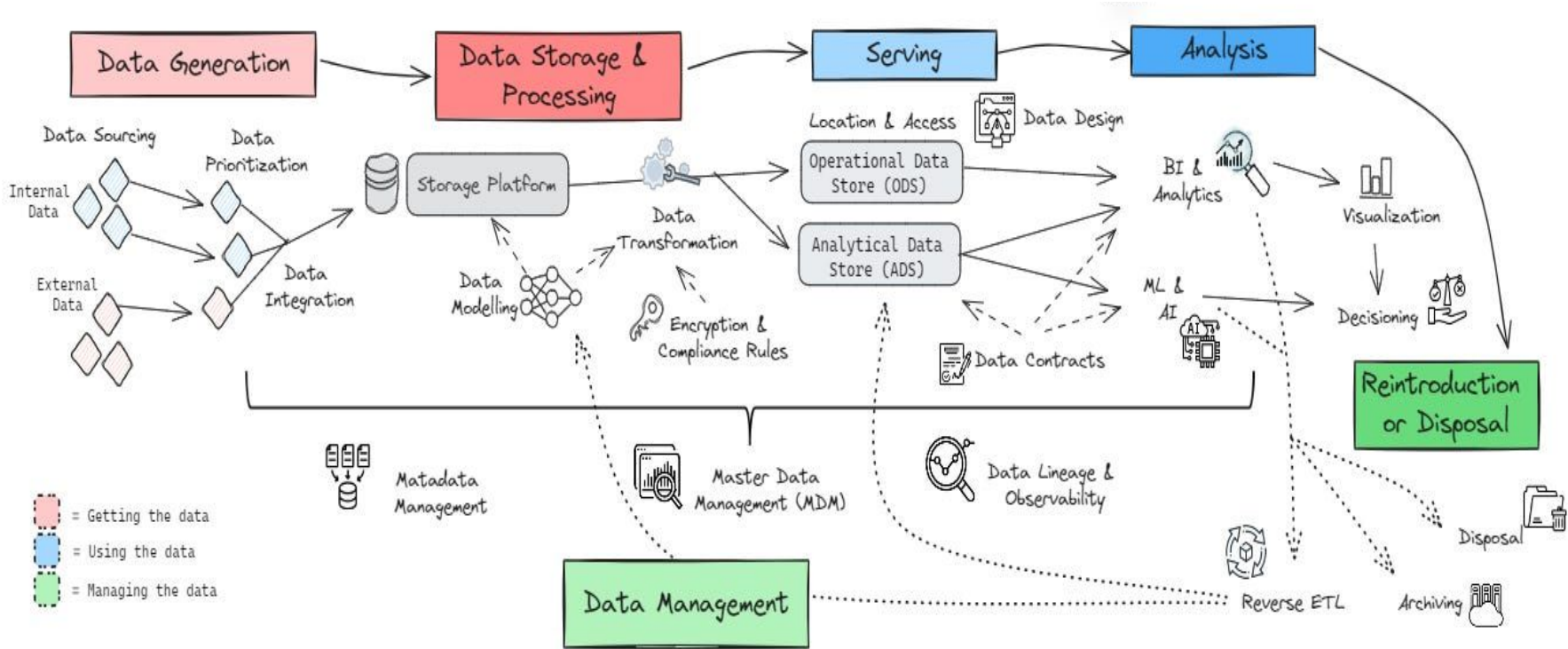
<https://www.nature.com/articles/s41588-020-0698-y/figures/2>

Why “health data” instead of “medical data”:

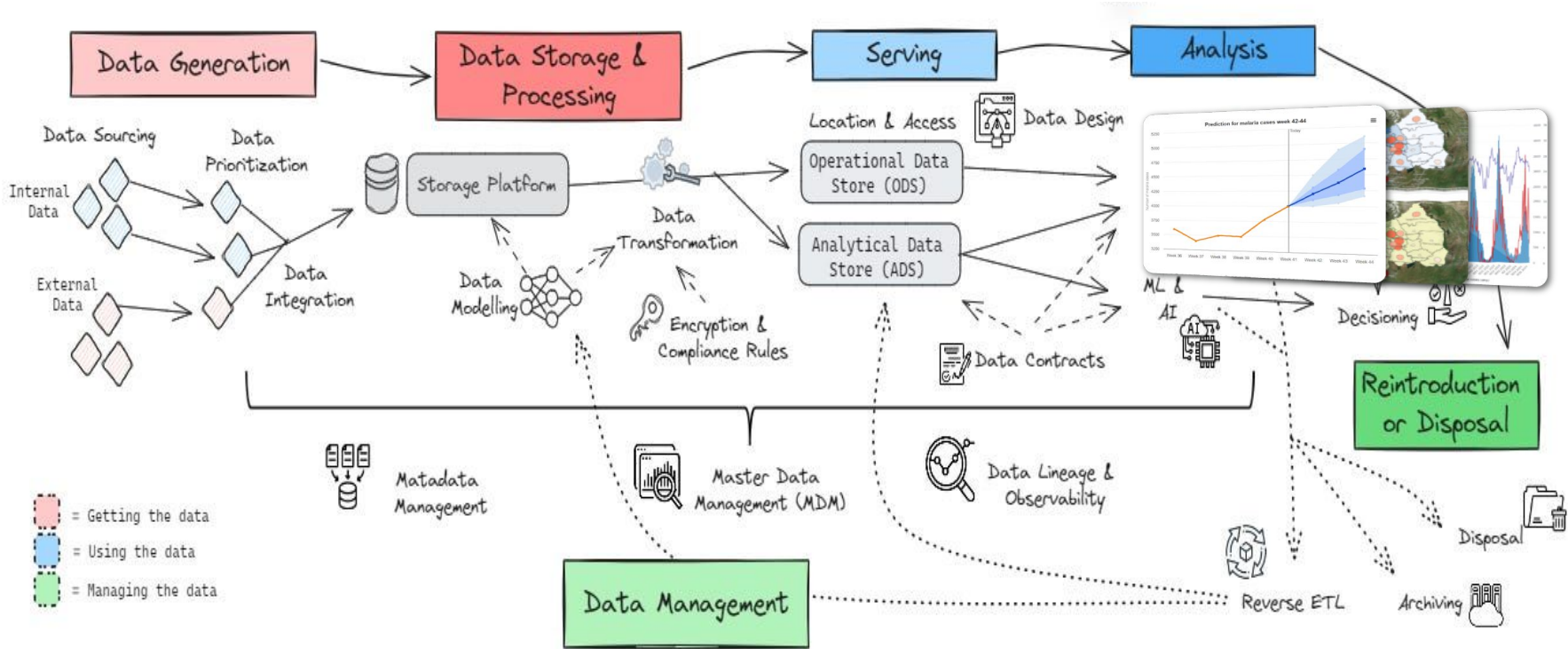
medical (being defined as MD-related”) \subset health - **contentious**

How does (health) data science (generally) differ from traditional epidemiology using secondary data?

Data science integrates within the wider data ecosystem



Data science integrates within the wider data ecosystem



Important as datafication pervades medicine (and more)

Medical Notes:

- 1-50MB

Laboratory Values:

- 10-2000MB

Physiological Sensors:

- 0.1-200GB

Genomics:

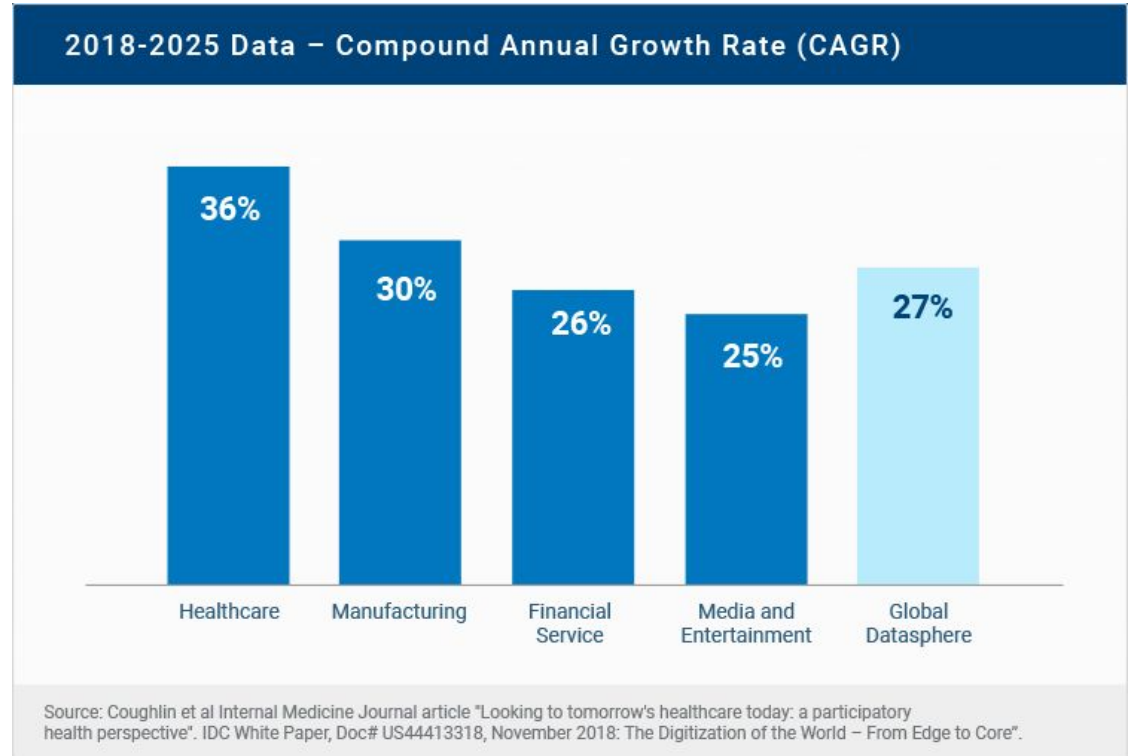
- 25-250GB

Medical Imaging:

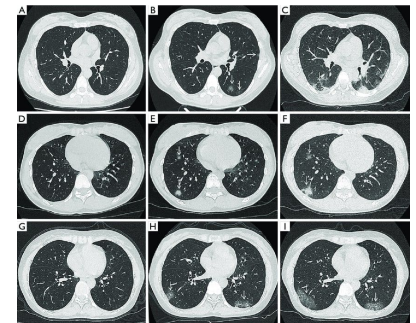
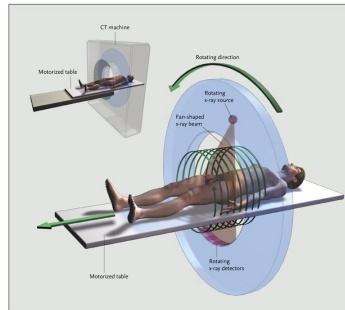
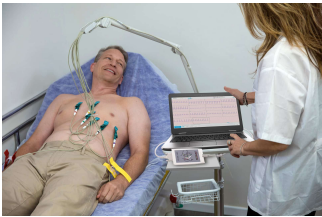
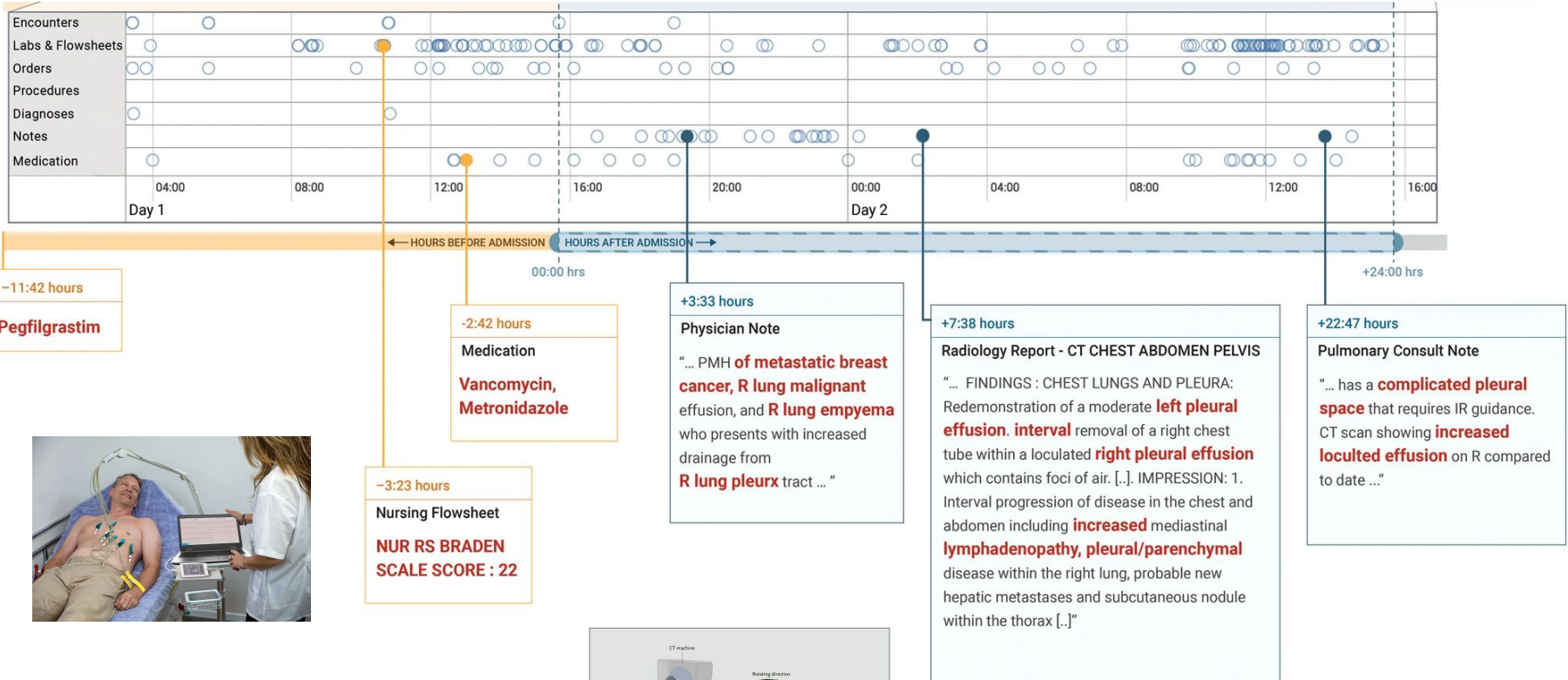
- 0.1-10TB

Hospital:

- 0.5-50PB per day



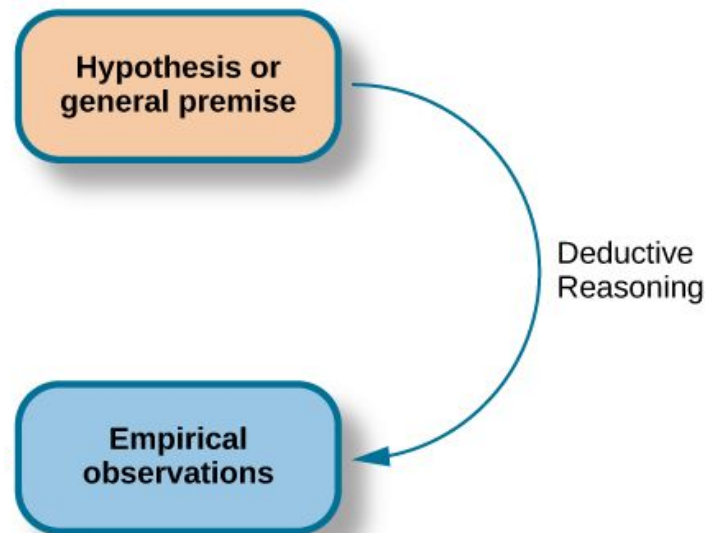
Data Science handles unstructured and multi-modal data



Data science supports inductive approaches

Deductive:

- “Condition X, causes Y”
- Collect data
- Perform (typically) frequentist statistical tests
- Reject or confirm null hypothesis



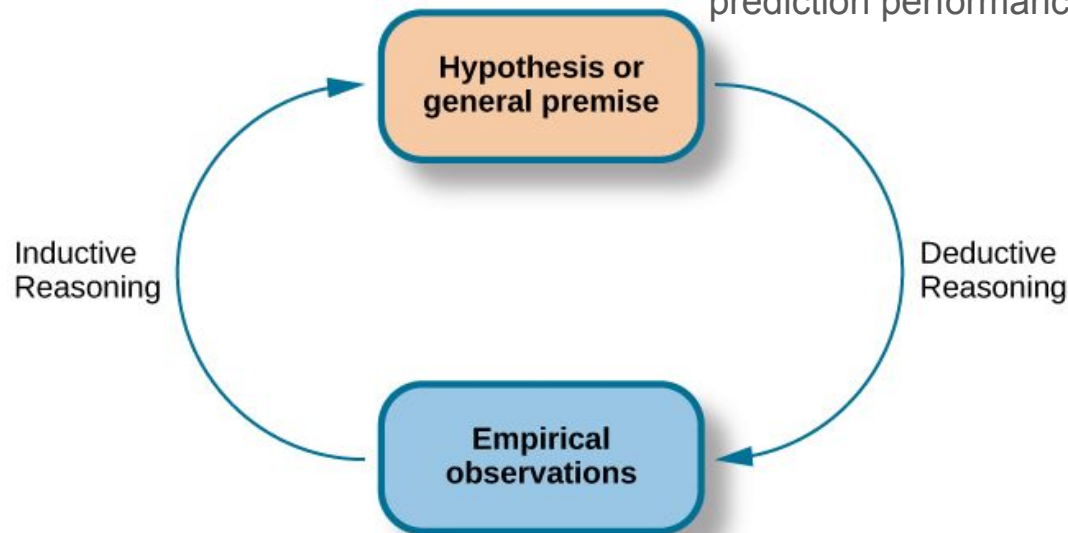
Data science supports inductive approaches

Deductive:

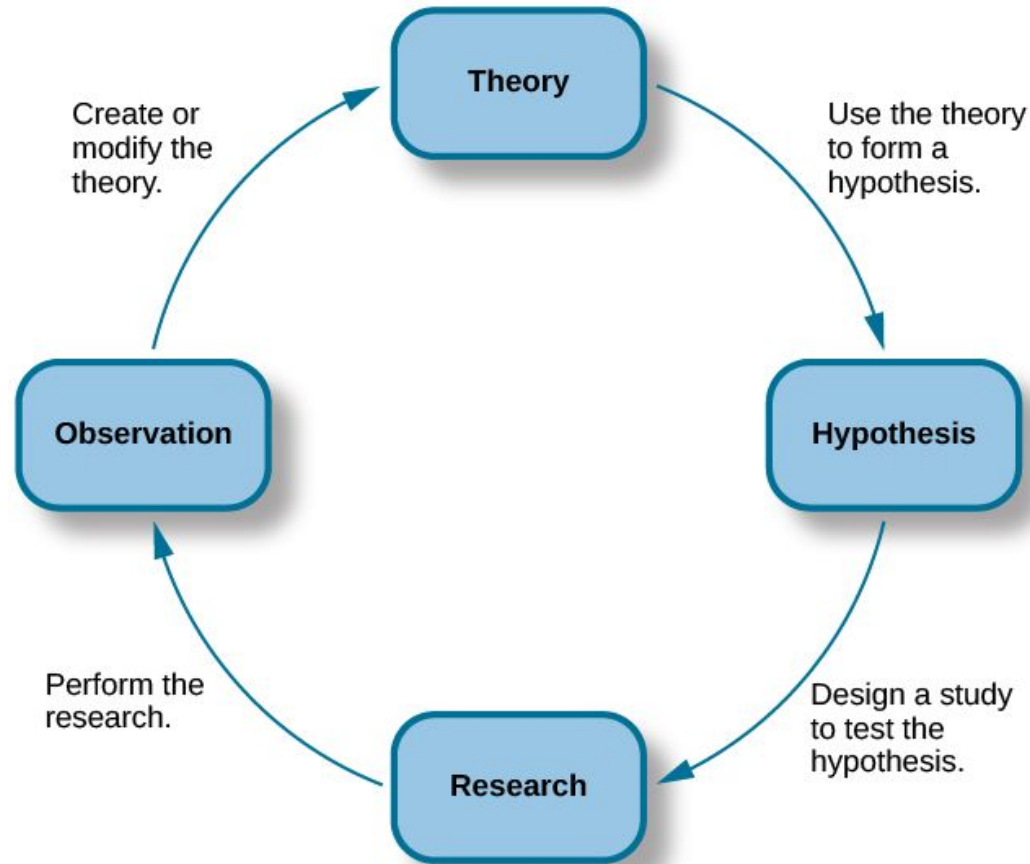
- “Condition X, causes Y”
- Collect data
- Perform (typically) frequentist statistical tests
- Reject or confirm null hypothesis

Inductive:

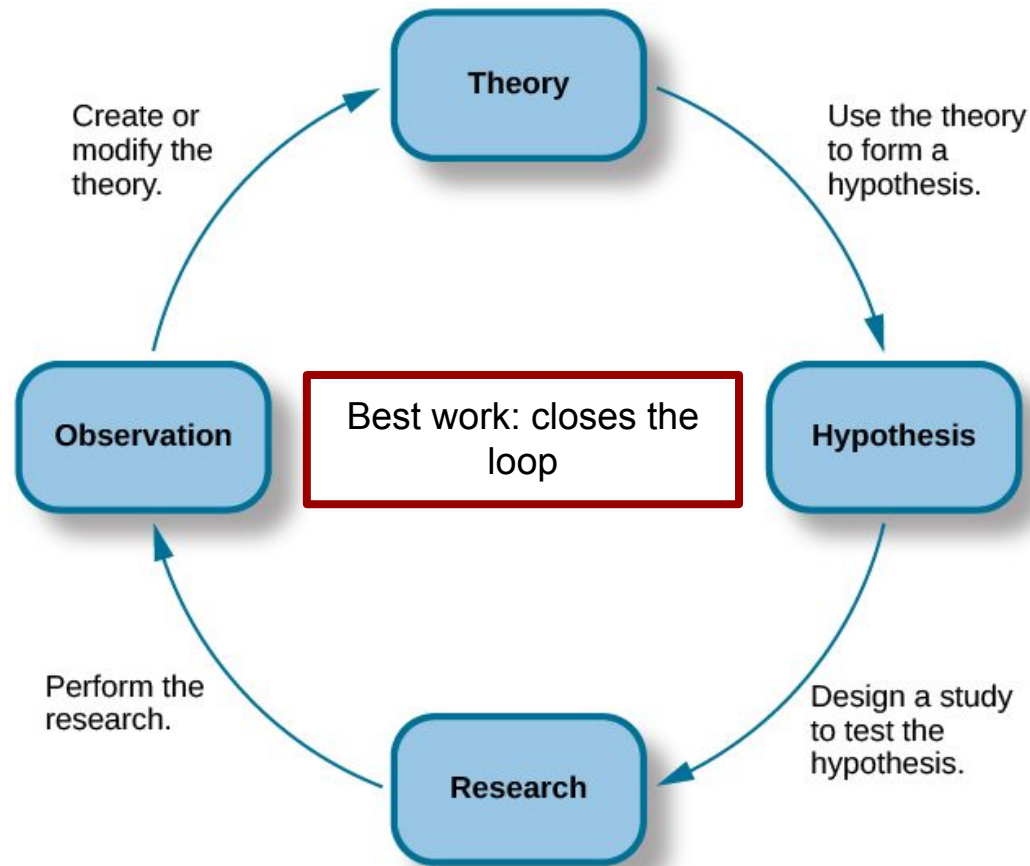
- Collect data
- Identify patterns in the data
- Observe X and Y seem connected somehow
- Quantify strength of association e.g., prediction performance



Data science aligns with knowledge cycle

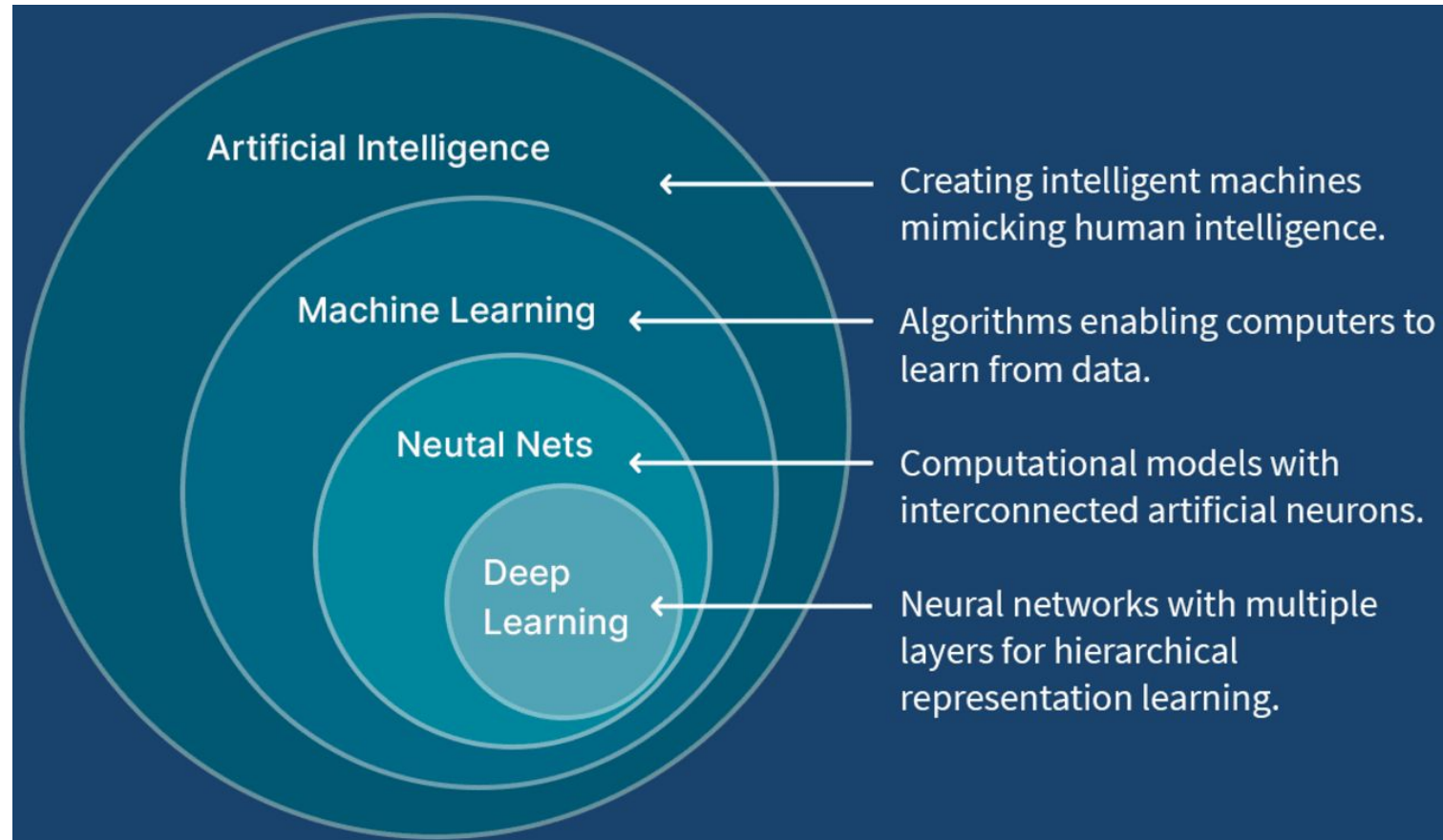


Data science aligns with knowledge cycle



Most obvious difference: Data Science often
is based in Machine Learning

Machine Learning is a subset of Artificial Intelligence



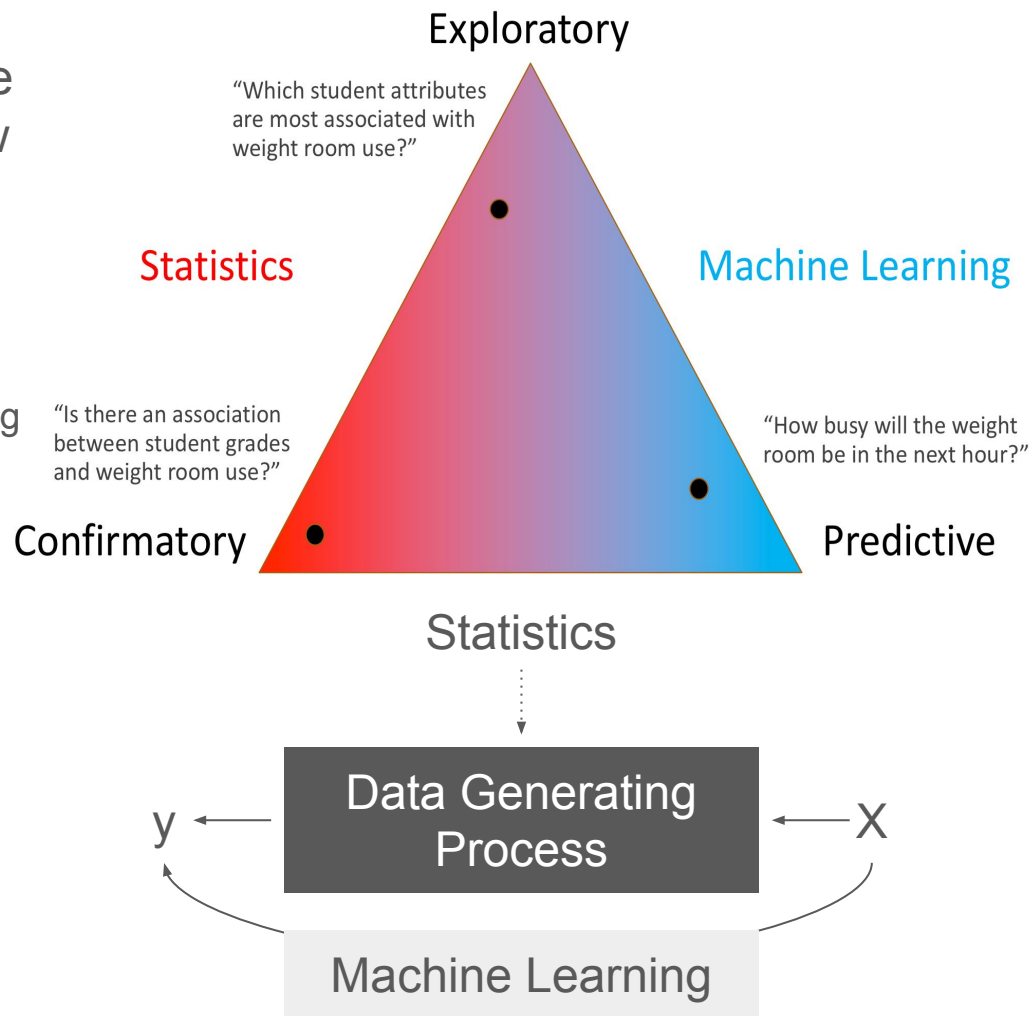
Is Machine Learning just a CS-flavoured
rebrand of statistics?

Large overlap but difference in people and priorities

- Many shared methods
- Difference in focus/priorities/culture

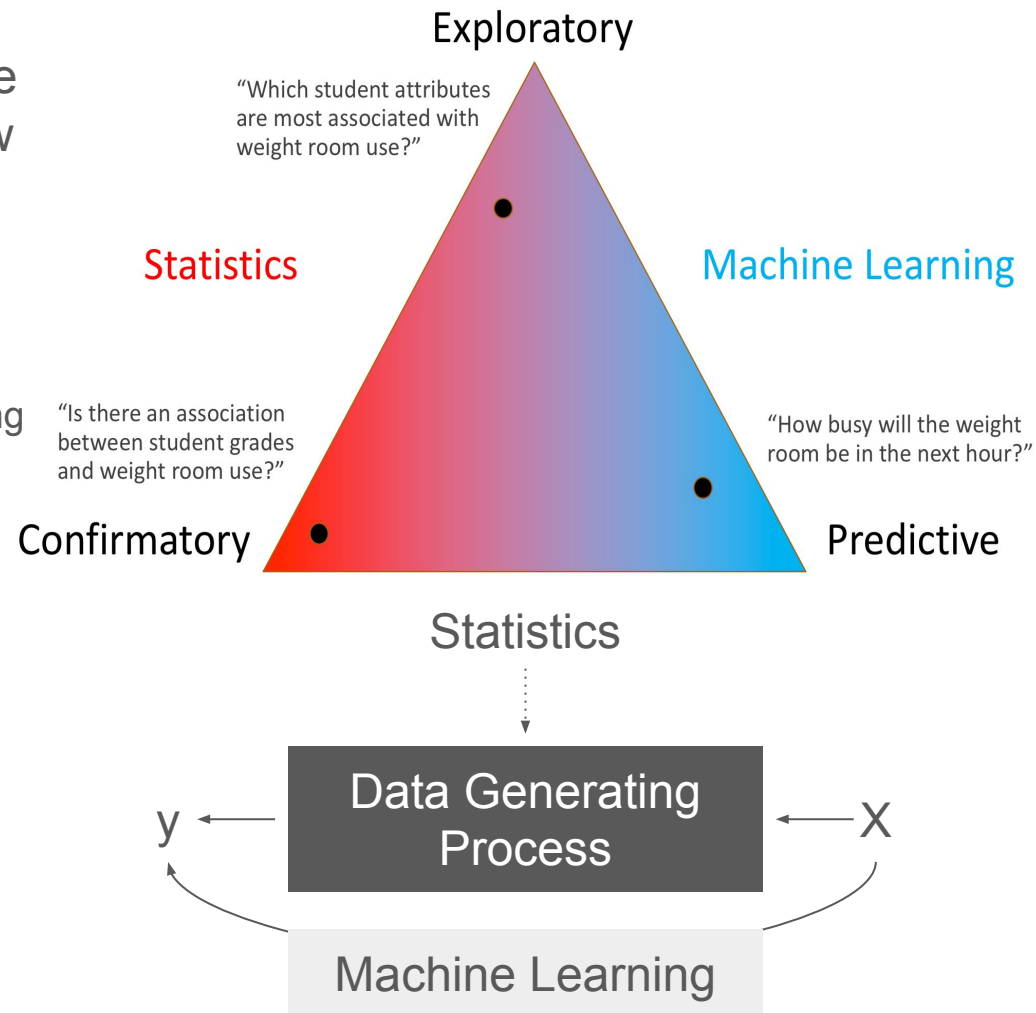
Large overlap but difference in people and priorities

- Many shared methods
- Difference in focus/priorities/culture
- Statistics ~ tries to understand how outcome was generated by data
- ML infers/learns a process for linking data to outcome
- Alternative framing:
 - Data Modelling vs Algorithmic Modelling



Large overlap but difference in people and priorities

- Many shared methods
- Difference in focus/priorities/culture
- Statistics ~ tries to understand how outcome was generated by data
- ML infers/learns a process for linking data to outcome
- Alternative framing:
 - Data Modelling vs Algorithmic Modelling
- DS/ML Pitfalls (can be):
 - Less rigorous/principled
 - Prone to reinventing the wheel
- DS/ML Benefits (can be):
 - More flexible
 - Less prescriptive/intimidating



Why learn about Health Data Science
research?

Opportunity of Health Data Science

Benefits (and pitfalls!) of data science in general combined with:

- Huge amounts of health data

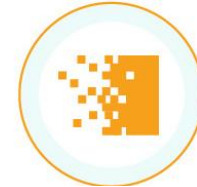
Growth in healthcare data

1 exabyte = 1 billion gigabytes

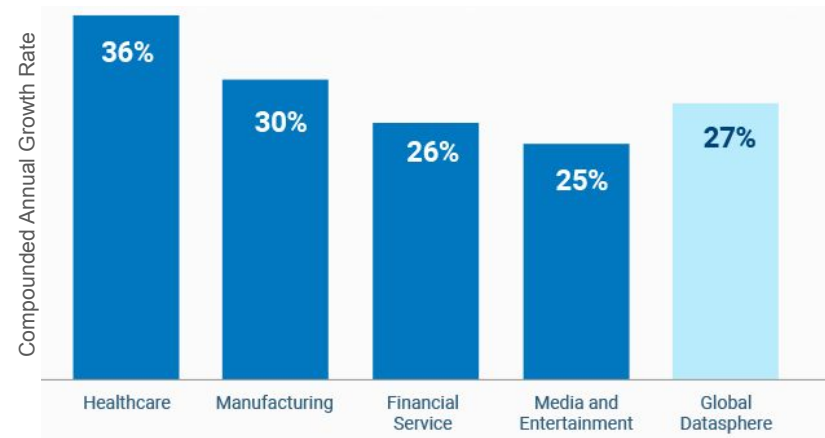
Source: Stanford Medicine 2017, IDC 2014



2013
153
EXABYTES



2020
2,314
EXABYTES



Opportunity of Health Data Science

Benefits (and pitfalls!) of data science in general combined with:

- Huge amounts of health data
- Many **interesting** and **important problems**

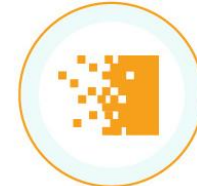
Growth in healthcare data

1 exabyte = 1 billion gigabytes

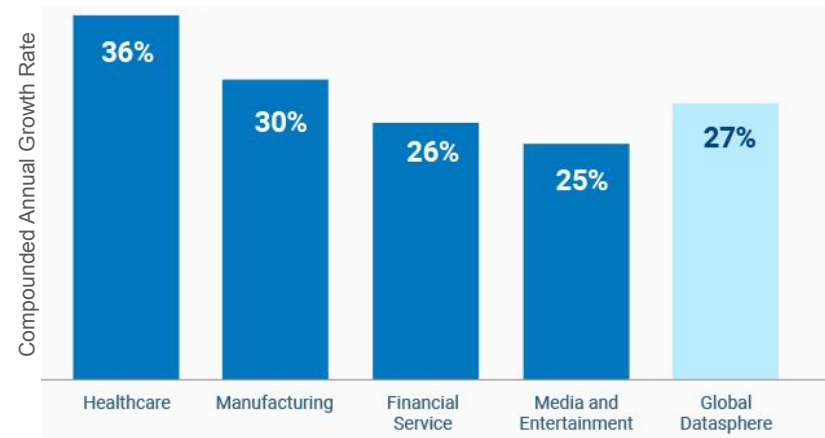
Source: Stanford Medicine 2017, IDC 2014



2013
153
EXABYTES



2020
2,314
EXABYTES



Opportunity of Health Data Science

Benefits (and pitfalls!) of data science in general combined with:

- Huge amounts of health data
- Many **interesting** and **important problems**
- Many domain experts desperate for data-related help with these problems

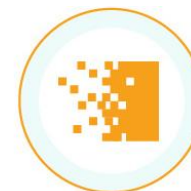
Growth in healthcare data

1 exabyte = 1 billion gigabytes

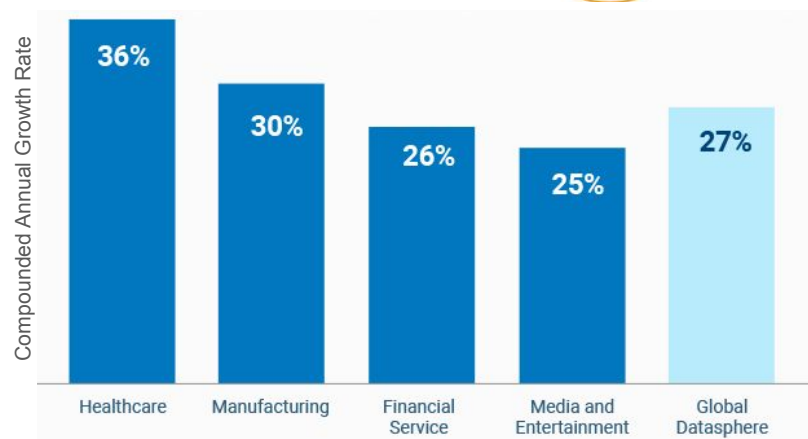
Source: Stanford Medicine 2017, IDC 2014



2013
153
EXABYTES



2020
2,314
EXABYTES



Opportunity of Health Data Science

Benefits (and pitfalls!) of data science in general combined with:

- Huge amounts of health data
- Many **interesting** and **important problems**
- Many domain experts desperate for data-related help with these problems
- Relative few skilled data science practitioners

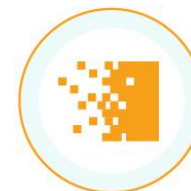
Growth in healthcare data

1 exabyte = 1 billion gigabytes

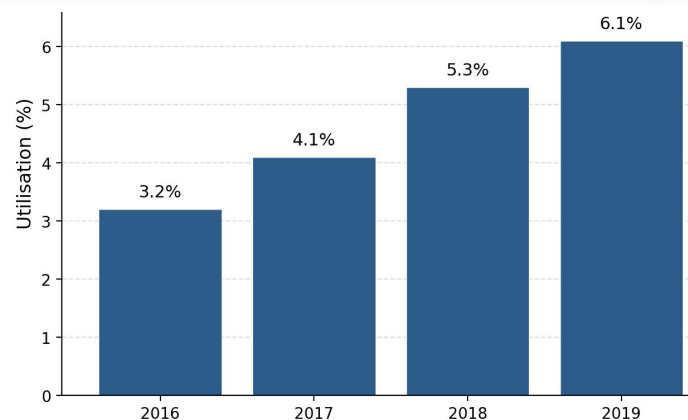
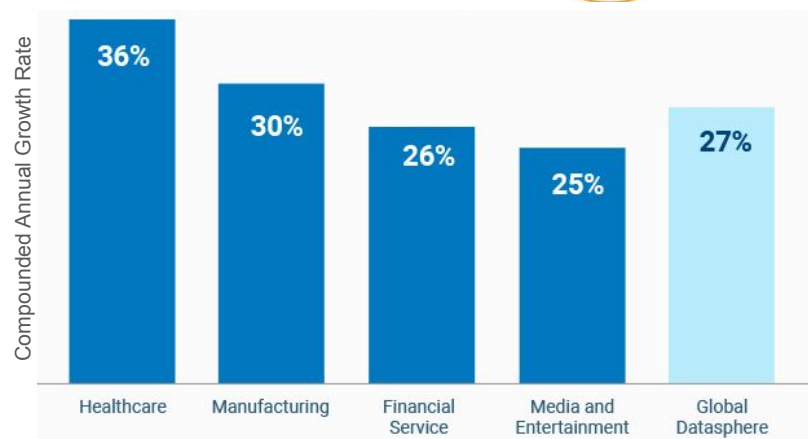
Source: Stanford Medicine 2017, IDC 2014



2013
153
EXABYTES



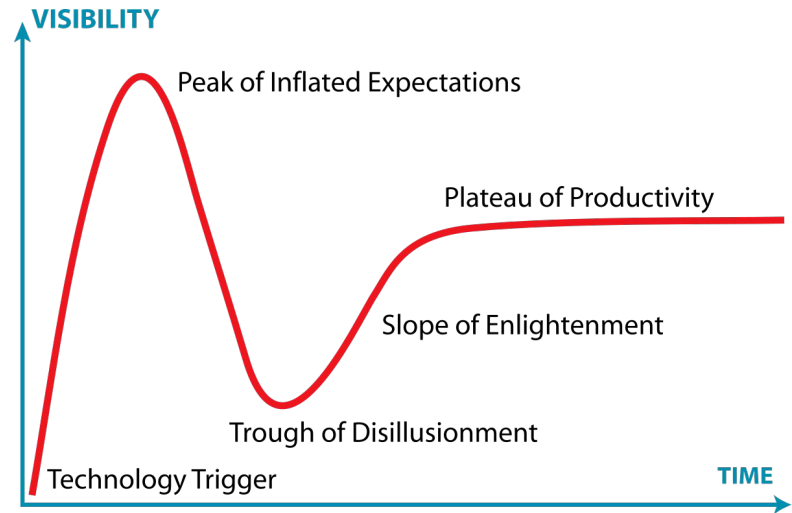
2020
2,314
EXABYTES



Xuanwu Hospital: % of staff requesting EMR data for research (Li et al., JMIR Med Inform 2021)

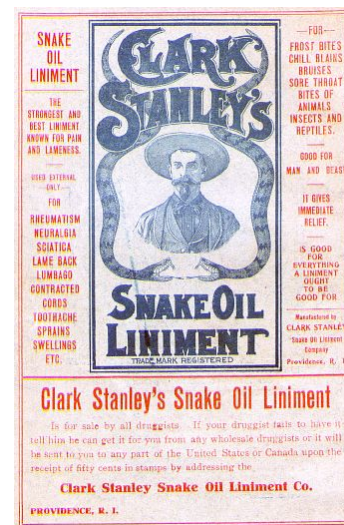
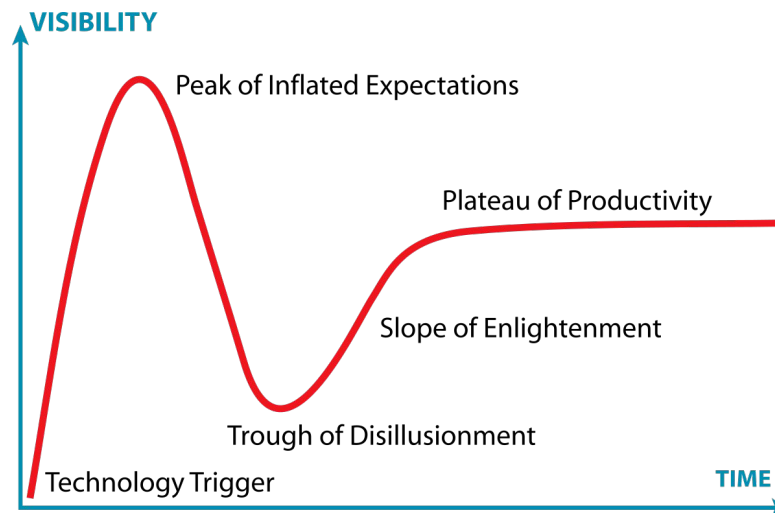
(Some) Challenges of Health Data Science

- Lots of hype



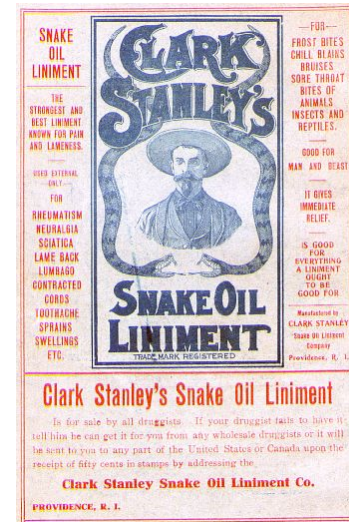
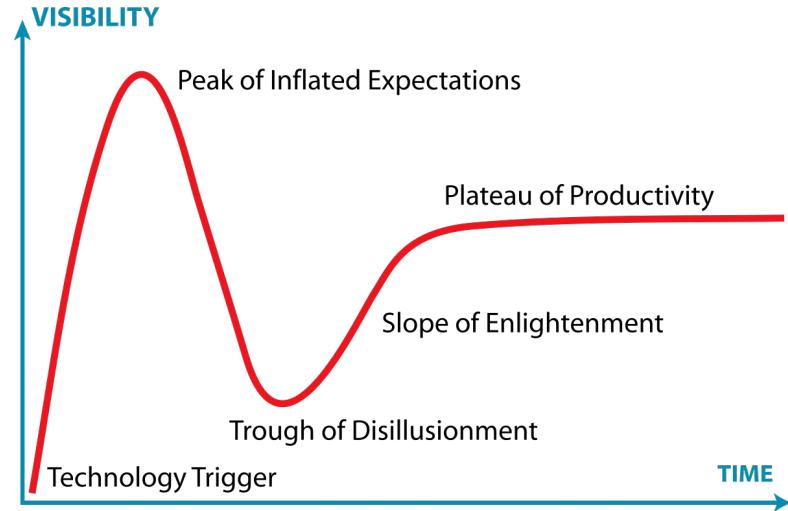
(Some) Challenges of Health Data Science

- Lots of hype
- Lots of grifters



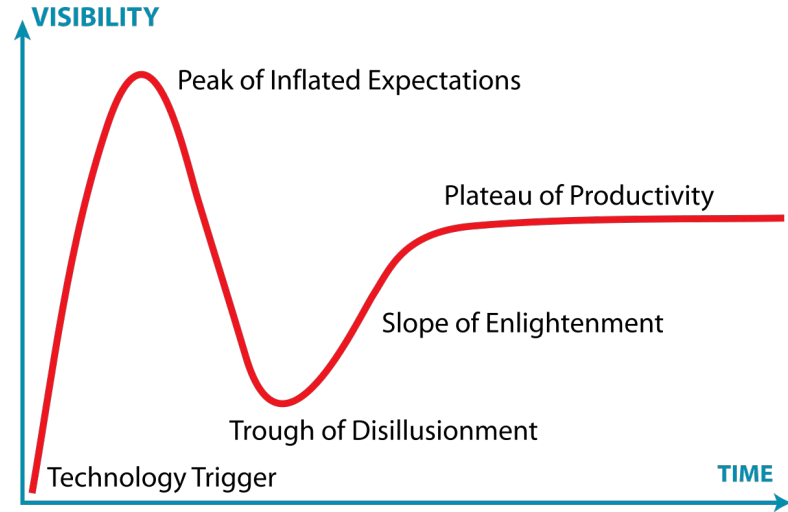
(Some) Challenges of Health Data Science

- Lots of hype
- Lots of grifters
- Data quality issues
- Contextual/Metadata quality issues

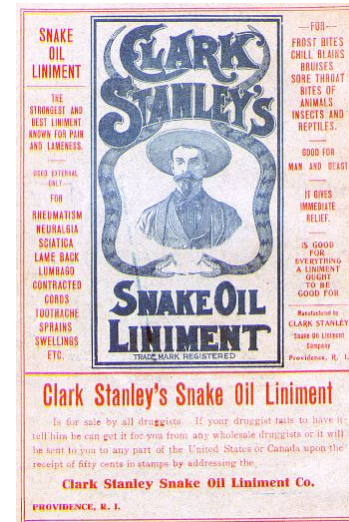


(Some) Challenges of Health Data Science

- Lots of hype
- Lots of grifters
- Data quality issues
- Contextual/Metadata quality issues
- Regulatory challenges
- Influence of US health system
- Ethical pitfalls
- Treatment to the mean
- Knowledge Translation and Operations: **Hard**

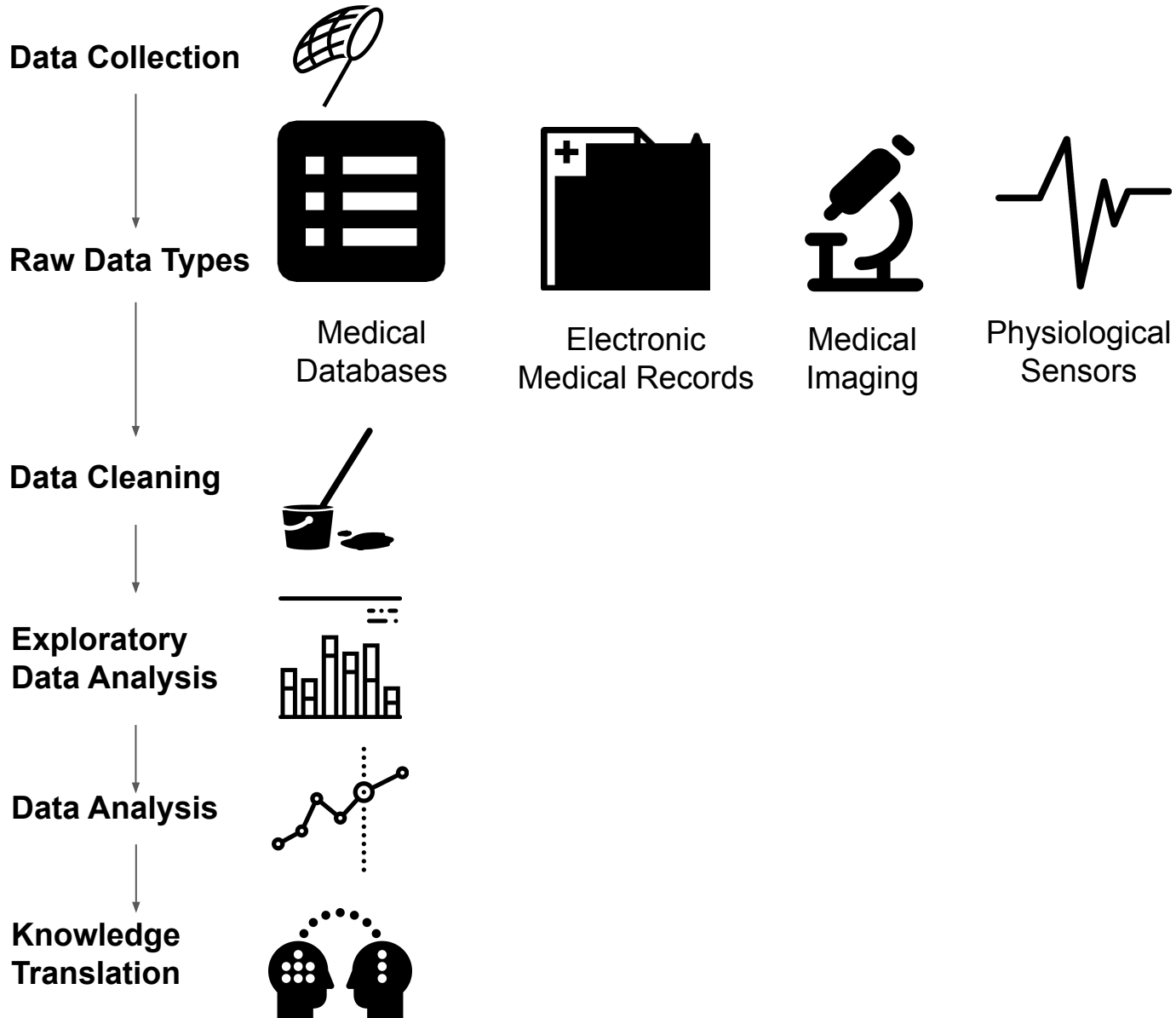


<https://www.r-bloggers.com/2019/08/new-course-learn-advanced-data-cleaning-in-r/>

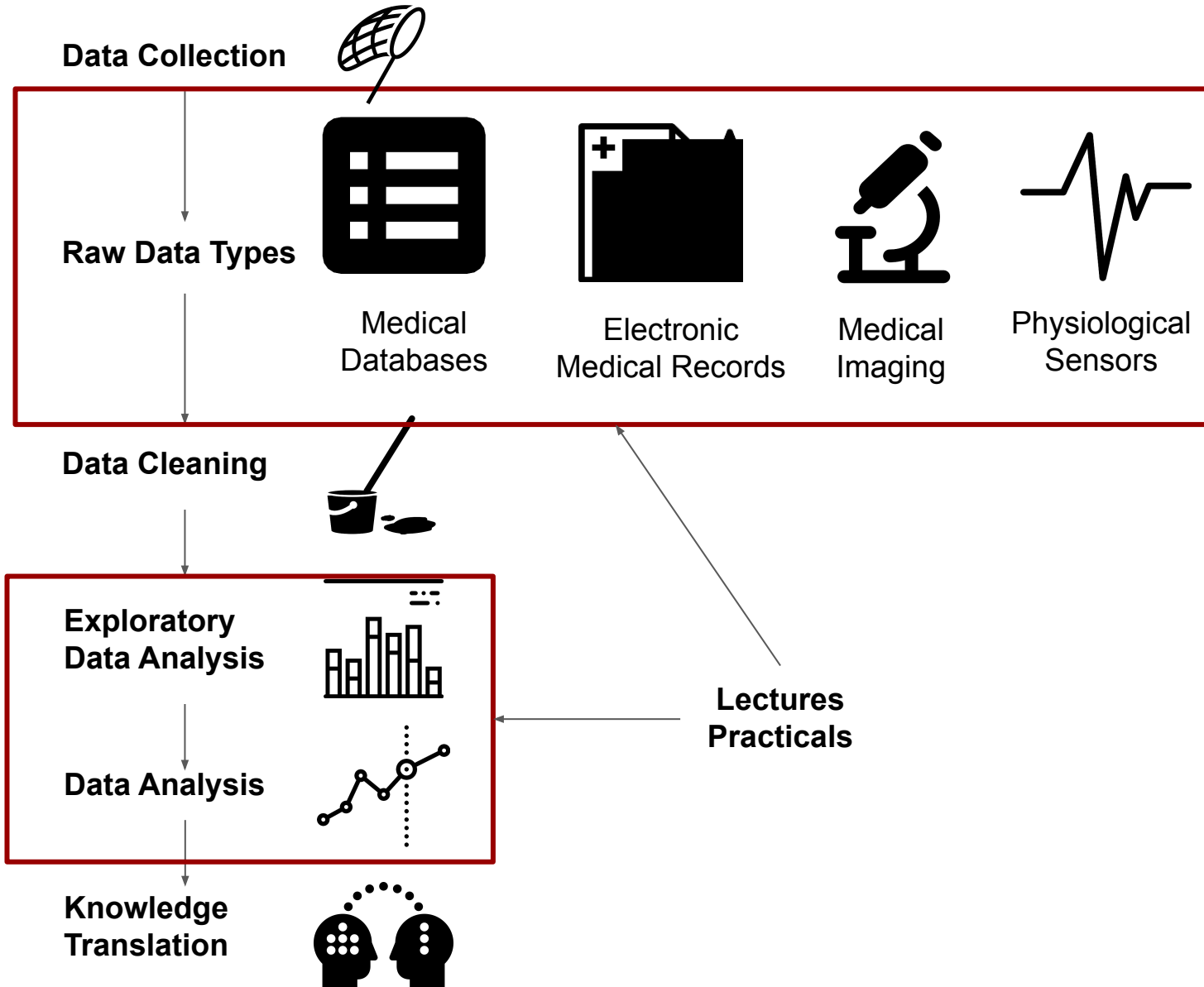


What parts of health data science will this course cover?

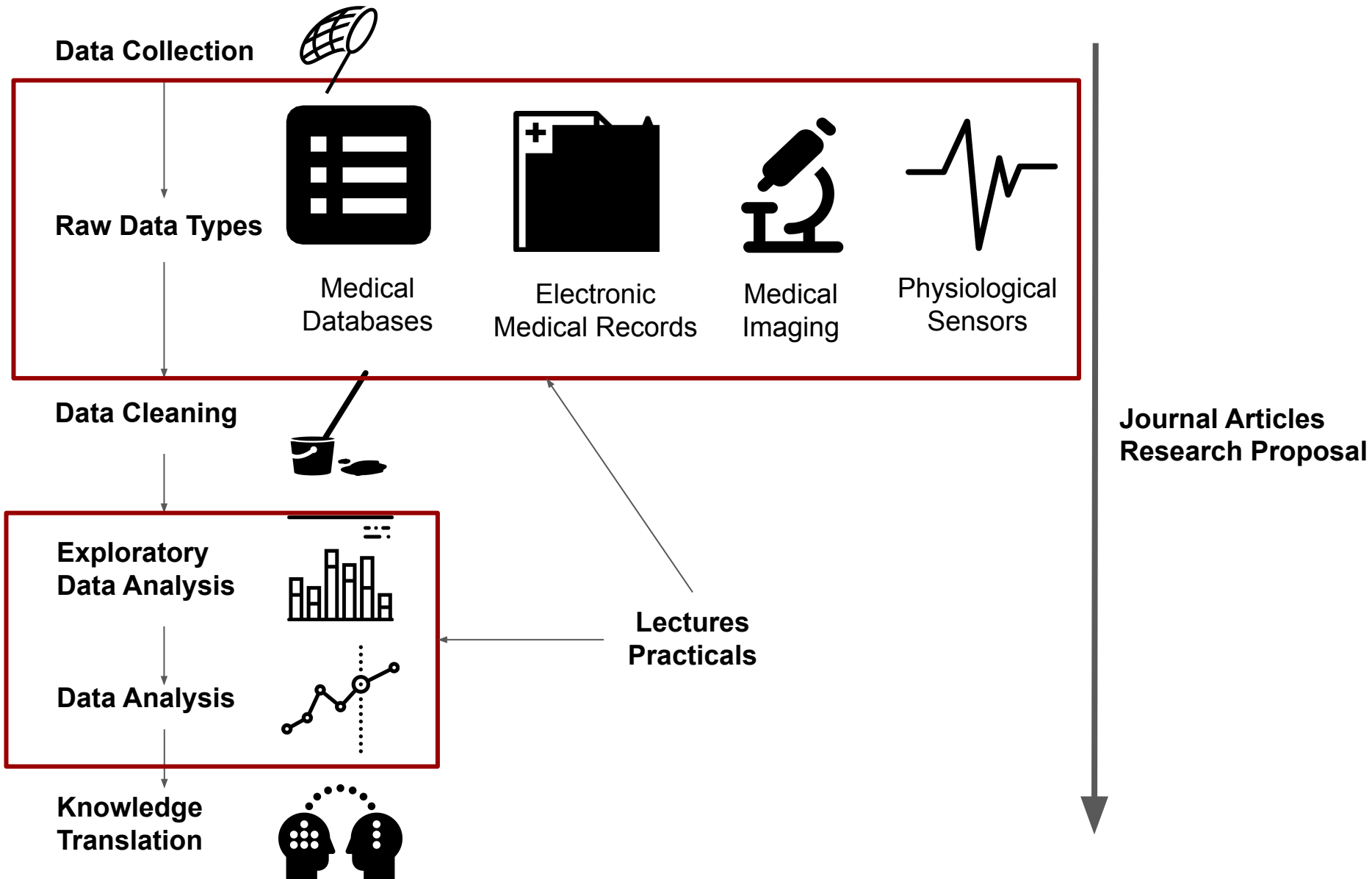
What parts of health data science will this course cover?



What parts of health data science will this course cover?



What parts of health data science will this course cover?



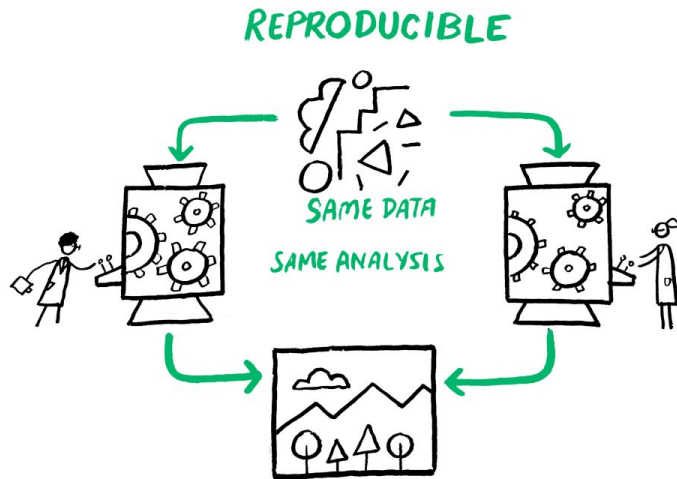
Let's take a 5 minute break!

Tools for Reproducible Health Data Science

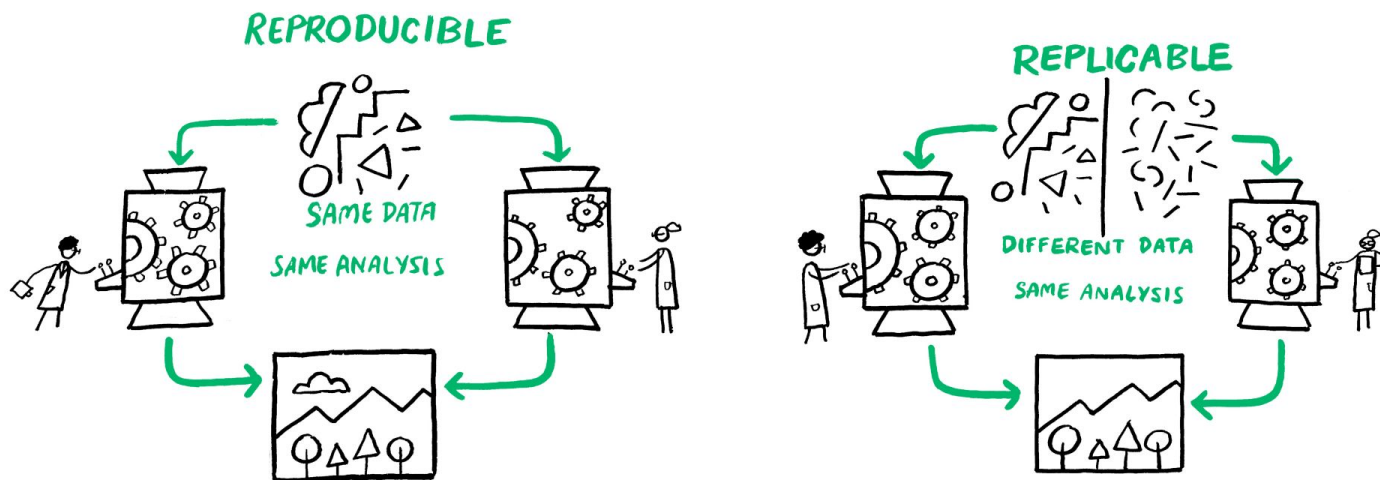
Rstudio, Rmarkdown, Git

Why do we care about reproducibility?

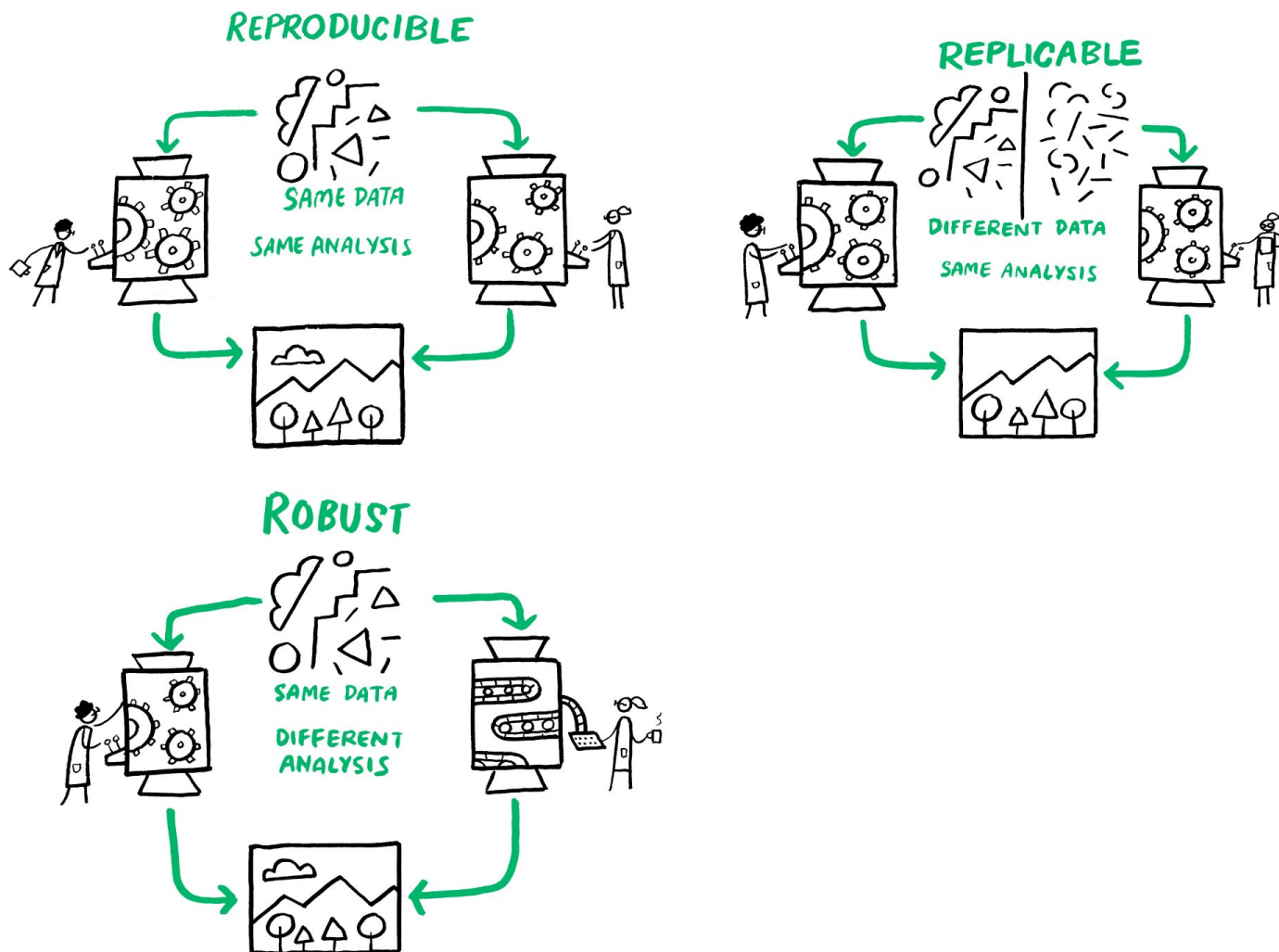
Reproducibility should be the bare minimum



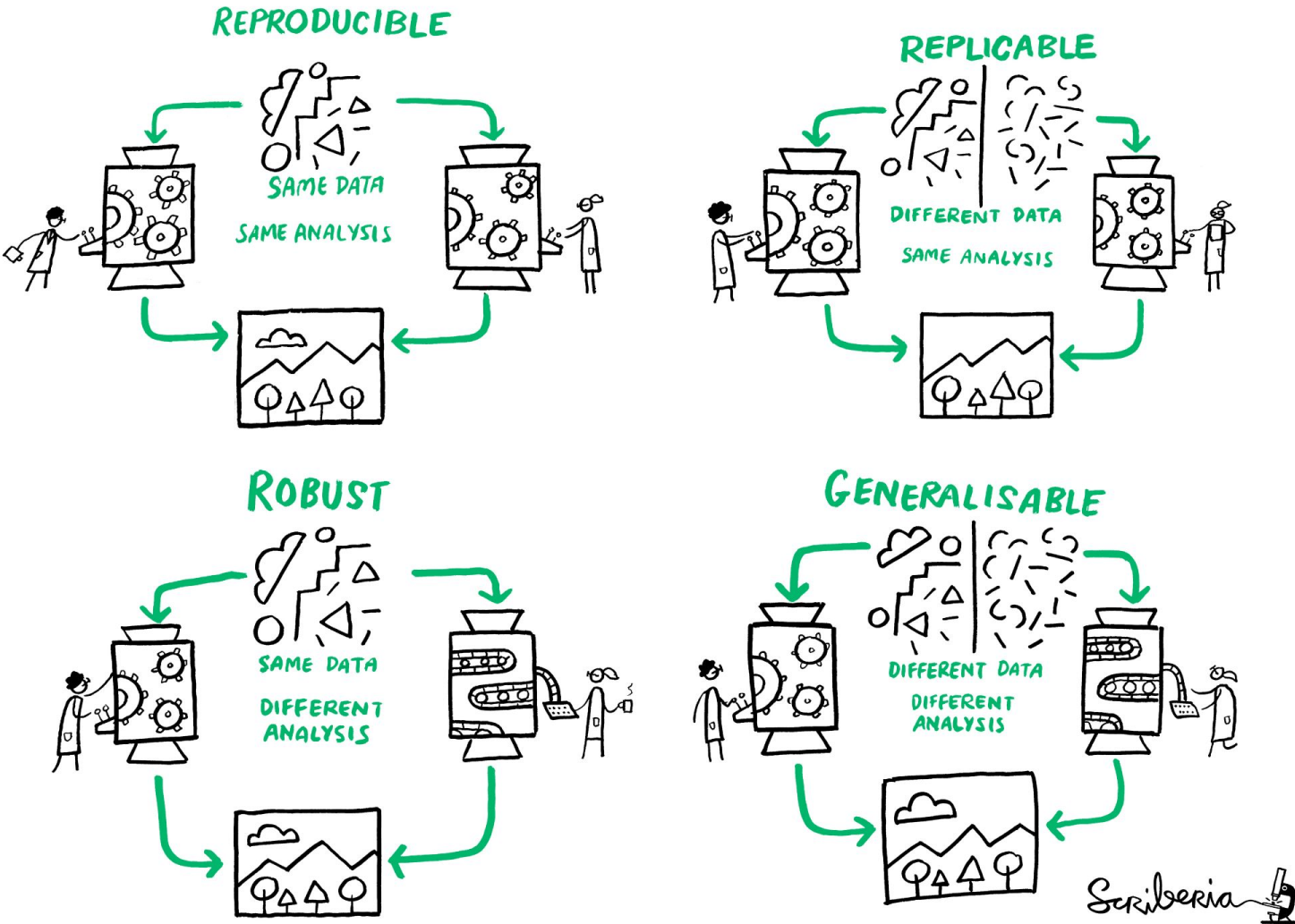
Reproducibility should be the bare minimum



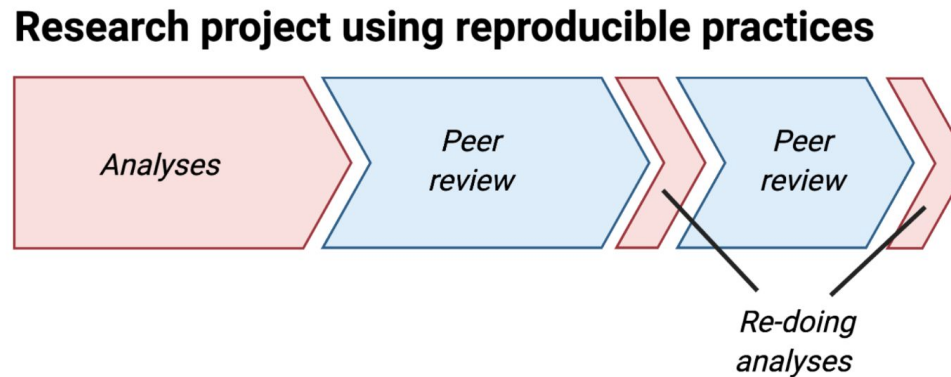
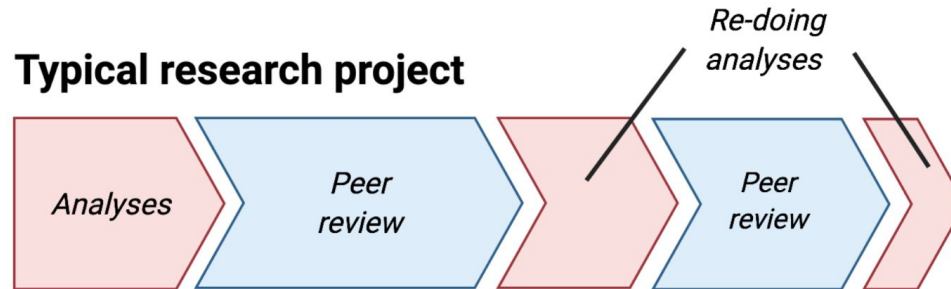
Reproducibility should be the bare minimum



Reproducibility should be the bare minimum



Makes your own life easier



 @dsquintana

oliviorgimenez.github.io/reproducible-science-workshop

What do we need to do to have reproducible research?

Reproducibility checklist

- Don't do anything by hand (even "one-off" tasks)

Reproducibility checklist

- Don't do anything by hand (even "one-off" tasks)
- Script every interaction with data:
 - Data collection
 - Moving data on your computer
 - Formatting datasets
 - Cleaning data
 - Exploratory data analysis
 - Main analyses
 - Report generation

Reproducibility checklist

- Don't do anything by hand (even "one-off" tasks)
- Script every interaction with data:
 - Data collection
 - Moving data on your computer
 - Formatting datasets
 - Cleaning data
 - Exploratory data analysis
 - Main analyses
 - Report generation
- Minimise interactivity/point and click interactions

Reproducibility checklist

- Don't do anything by hand (even "one-off" tasks)
- Script every interaction with data:
 - Data collection
 - Moving data on your computer
 - Formatting datasets
 - Cleaning data
 - Exploratory data analysis
 - Main analyses
 - Report generation
- Minimise any untracked interactivity/point and click interactions
- Version control all data, code, and documentation

Reproducibility checklist

- Don't do anything by hand (even "one-off" tasks)
- Script every interaction with data:
 - Data collection
 - Moving data on your computer
 - Formatting datasets
 - Cleaning data
 - Exploratory data analysis
 - Main analyses
 - Report generation
- Minimise any untracked interactivity/point and click interactions
- Version control all data, code, and documentation
- Use a random seed

Reproducibility checklist

- Don't do anything by hand (even "one-off" tasks)
- Script every interaction with data:
 - Data collection
 - Moving data on your computer
 - Formatting datasets
 - Cleaning data
 - Exploratory data analysis
 - Main analyses
 - Report generation
- Minimise any untracked interactivity/point and click interactions
- Version control all data, code, and documentation
- Use a random seed
- Keep track of the exact version of every library/program you use
- Remove hard-coded/system specific paths

How do we actually do these things?

Choose a language that makes it easy to do most/all of your analysis

Choose a language that makes it easy to do most/all of your analysis



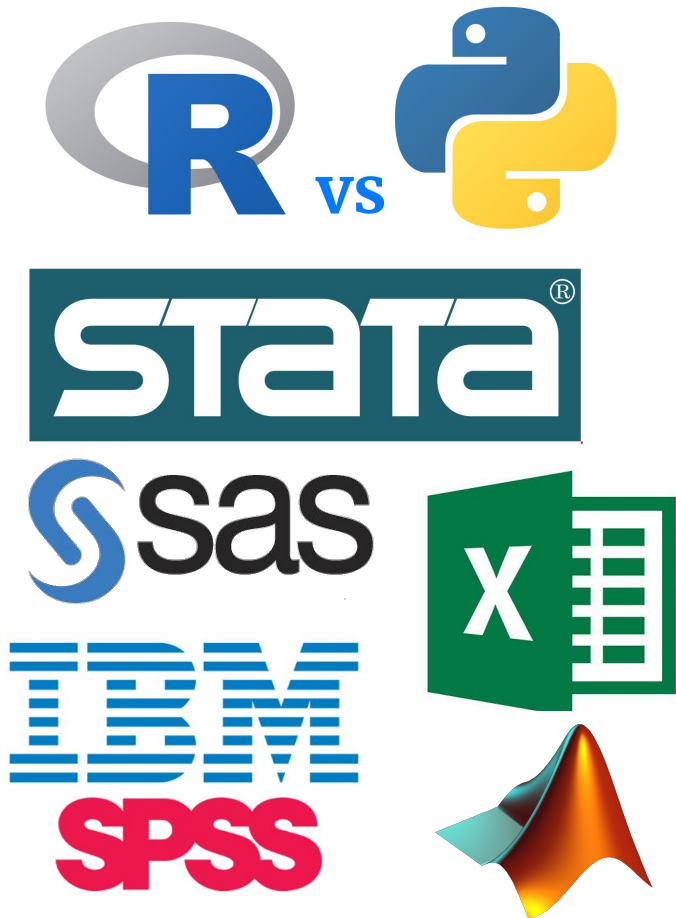
May 2026	May 2025	Change	Programming Language	Ratings	Change
1	1		Python	19.98%	-5.37%
2	3	▲	C	11.55%	+1.84%
3	4	▲	Java	7.94%	-1.37%
4	2	▼	C++	7.92%	-2.02%
5	5		C#	5.41%	+1.19%
6	6		JavaScript	3.08%	-0.60%
7	8	▲	Visual Basic	2.90%	+0.28%
8	12	▲	R	1.77%	+0.31%
9	10	▲	SQL	1.57%	-0.33%
10	9	▼	Delphi/Object Pascal	1.44%	-0.85%

<https://www.tiobe.com/tiobe-index/>

2026 consolidation toward R for statistical:

- MATLAB #20
- SAS #28
- Wolfram/Mathematica: #50-100
- S #102
- SPSS #112
- Stata #124

Choose a language that makes it easy to do most/all of your analysis



May 2026	May 2025	Change	Programming Language	Ratings	Change
1	1		Python	19.98%	-5.37%
2	3	▲	C	11.55%	+1.84%
3	4	▲	Java	7.94%	-1.37%
4	2	▼	C++	7.92%	-2.02%
5	5		C#	5.41%	+1.19%
6	6		JavaScript	3.08%	-0.60%
7	8	▲	Visual Basic	2.90%	+0.28%
8	12	▲▲	R	1.77%	+0.31%
9	10	▲	SQL	1.57%	-0.33%
10	9	▼	Delphi/Object Pascal	1.44%	-0.85%

<https://www.tiobe.com/tiobe-index/>

2026 consolidation toward R for statistical:

- MATLAB #20
- SAS #28
- Wolfram/Mathematica: #50-100
- S #102
- SPSS #112
- Stata #124

Use a data science focused IDE: Rstudio

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for loading packages, creating a 'daily' dataset, and plotting the number of flights per weekday.
- Environment:** Shows the 'daily' object with 365 observations and 3 variables.
- Console:** Shows the execution of the code, including the output of the 'daily' dataset and the first three rows of the boxplot.
- Plots:** Displays a boxplot titled 'Number of 2013 New York Flights Each Weekday'.

```
1 library(nycflights13) ## package containing flights dataset
2 library(lubridate)
3 library(dplyr)
4 library(ggplot2)
5
6 head(flights, n = 3)
7 daily <- flights %>%
8   mutate(date = make_date(year, month, day)) %>%
9   count(date) %>%
10  mutate(wday = wday(date, label = TRUE))
11 head(daily, n = 3)
12 ggplot(daily, aes(wday, n)) +
13   geom_boxplot(outlier.colour = "hotpink") +
14   labs(x = "Weekday", y = "Flights",
15        subtitle = "Number of 2013 New York Flights Each Weekday")
16
```

Console Output:

```
# A tibble: 3 x 19
  year month day dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier
  <int> <int> <int> <int> <int> <int> <dbl> <int> <int> <dbl> <chr>
1 2013 1 1 517 515 2 830 819 11 UA
2 2013 1 1 533 529 4 850 830 20 UA
3 2013 1 1 542 540 2 923 850 33 AA
# ... with 9 more variables: flight <int>, tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
# distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
> daily <- flights %>%
+   mutate(date = make_date(year, month, day)) %>%
+   count(date) %>%
+   mutate(wday = wday(date, label = TRUE))
> head(daily, n = 3)
# A tibble: 3 x 3
  date           n wday
  <date> <int> <ord>
1 2013-01-01 842 Tue
2 2013-01-02 943 Wed
3 2013-01-03 914 Thu
> ggplot(daily, aes(wday, n)) +
+   geom_boxplot(outlier.colour = "hotpink") +
+   labs(x = "Weekday", y = "Flights",
+        subtitle = "Number of 2013 New York Flights Each Weekday")
>
```

Boxplot Data Summary:

Weekday	Min	Q1	Median	Q3	Max
Sun	720	810	840	910	990
Mon	910	940	960	990	1000
Tue	760	840	870	900	950
Wed	720	840	870	900	950
Thu	740	840	870	900	950
Fri	820	840	870	900	950
Sat	680	710	750	780	860

.Rproj file
here()

set.seed()
sessionInfo()

Use notebooks to document analyses: Rmarkdown/Quarto

The screenshot displays the RStudio interface with an R Markdown notebook open. The notebook content is as follows:

```
1 ---
2 title: "Viridis Notebook"
3 output: html_notebook
4 ---
5
6 ```{r include = FALSE}
7 library(viridis)
8 ```
9
10 The code below demonstrates two color palettes in the
11 [viridis](https://github.com/sjmgarnier/viridis) package. Each
12 plot displays a contour map of the Maunga Whau volcano in
13 Auckland, New Zealand.
14
15 ## Viridis colors
16
17 ```{r}
18 image(volcano, col = viridis(200))
19 ```
```

The notebook shows two contour plots of the Maunga Whau volcano. The first plot, titled "Viridis colors", uses the viridis color palette. The second plot, titled "Magma colors", uses the magma color palette. Both plots show a contour map of the volcano with axes ranging from 0.0 to 1.0. The viridis plot shows a color gradient from dark purple to bright yellow, while the magma plot shows a color gradient from dark purple to bright red.

The RStudio interface includes a console at the bottom, a file browser on the left, and a viewer on the right. The viewer displays the rendered HTML output of the notebook, including the title "Viridis Notebook", the introductory text, the section heading "Viridis colors", the R code chunk, and the resulting contour plot. The viewer also includes a "Code" dropdown menu with options to "Show All Code", "Hide All Code", and "Download Rmd", and a "Hide" button.

Use notebooks to document analyses: Rmarkdown/Quarto

settings). Therefore, from this time onward, case counts are likely underestimated and the sequenced virus diversity is not necessarily representative of the virus circulating in the overall population.

BC AB SK MB ON QC NS NB NL

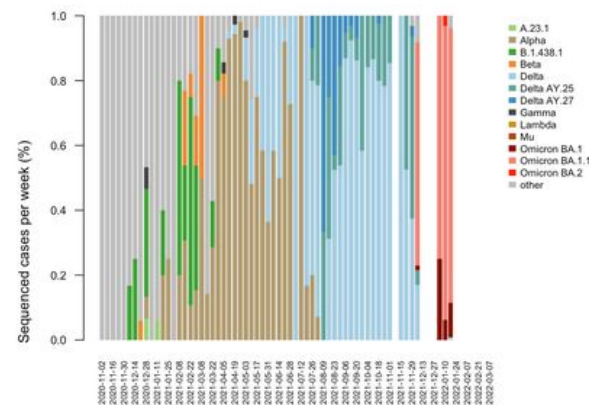
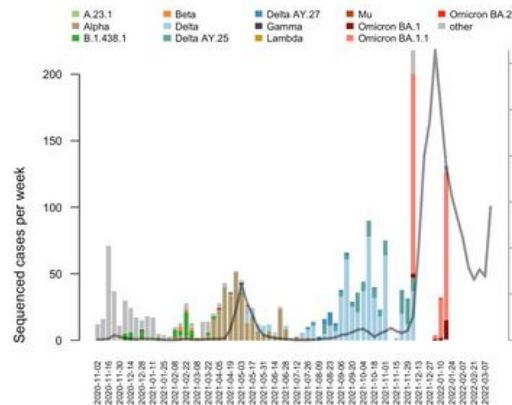
Nova Scotia

Additional up-to-date COVID data for this province can be found here:

<https://experience.arcgis.com/experience/204d6ed723244dfbb763ca3f913c5cad>

Hide

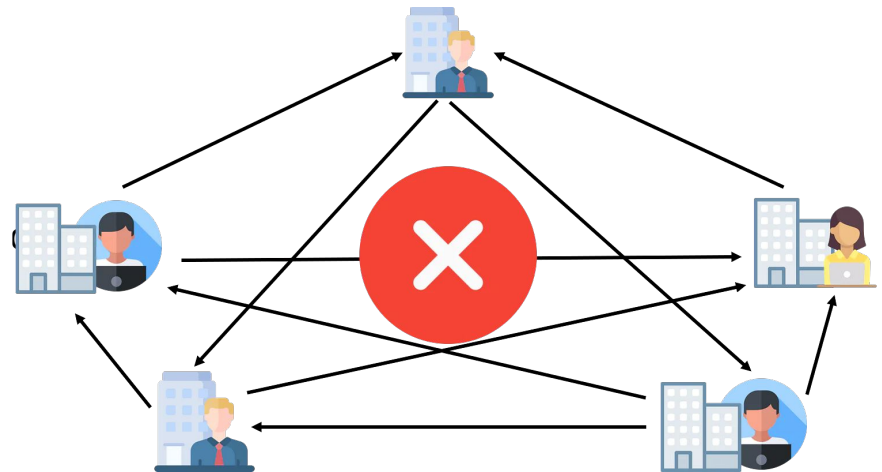
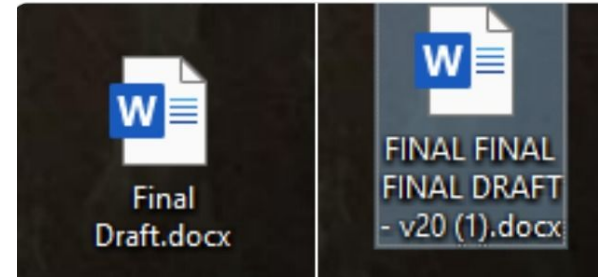
```
plot.variants(region='Nova Scotia')
plot.variants(region='Nova Scotia', scaled=T)
```



<https://covarr-net.github.io/duotang/duotang.html#>

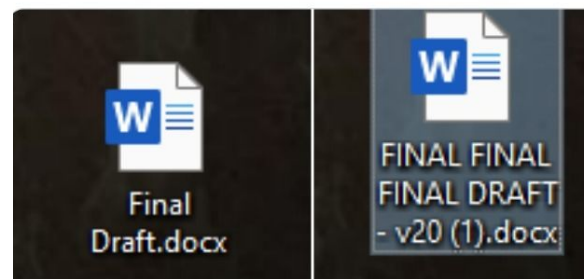
Use standard version control systems

- Ever had a nightmare of versioning even when just you?
- Add more people and the chaos grows exponentially!



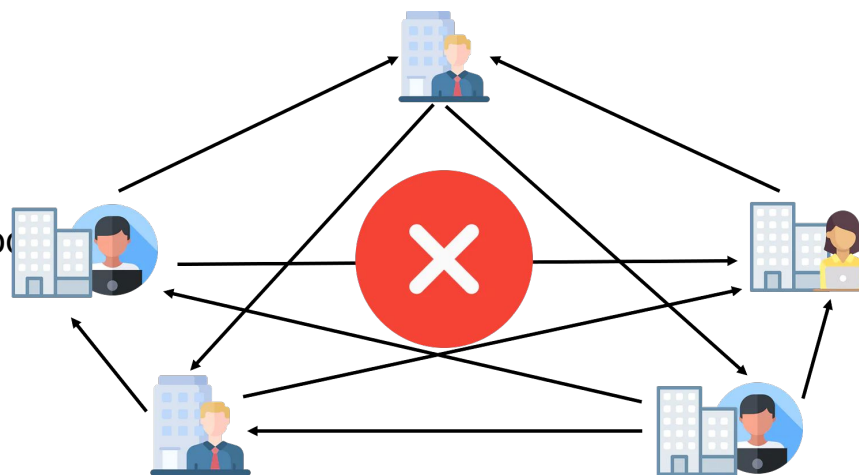
Use standard version control systems

- Ever had a nightmare of versioning even when just you?
- Add more people and the chaos grows exponentially!

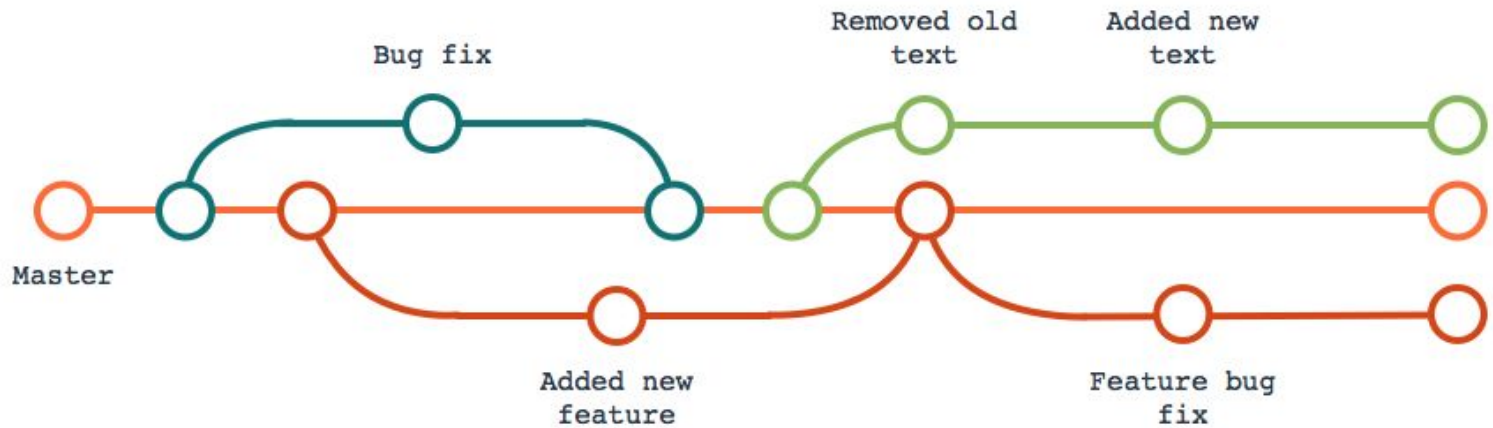


Version control let's you:

- Revert mistakes
- Acts as a comprehensive backup
- Let's you maintain multiple versions of your analysis
- Let's you compare different versions of your code
- Track down the who/what broke the analysis
- Work out why you did something in the past
- Build on someone else's work
- Share your own work
- Experiment without risk

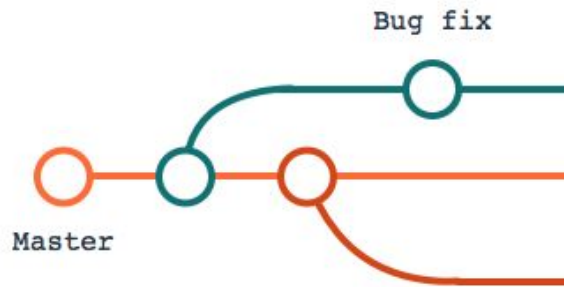


Git Version Control



- Most popular
- Decentralised
- Designed for
- GitLab/GitHub Services

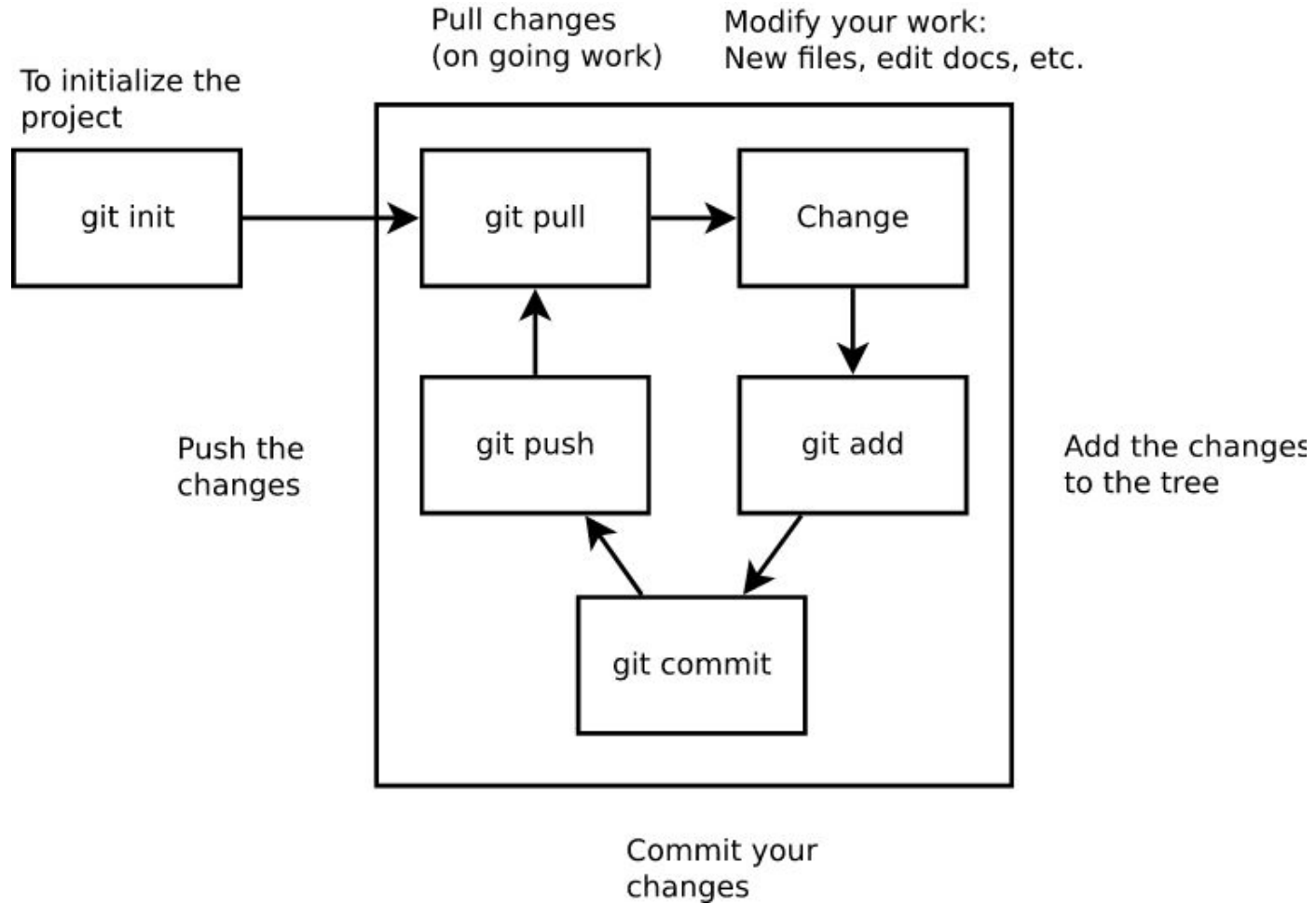
Git Version Control



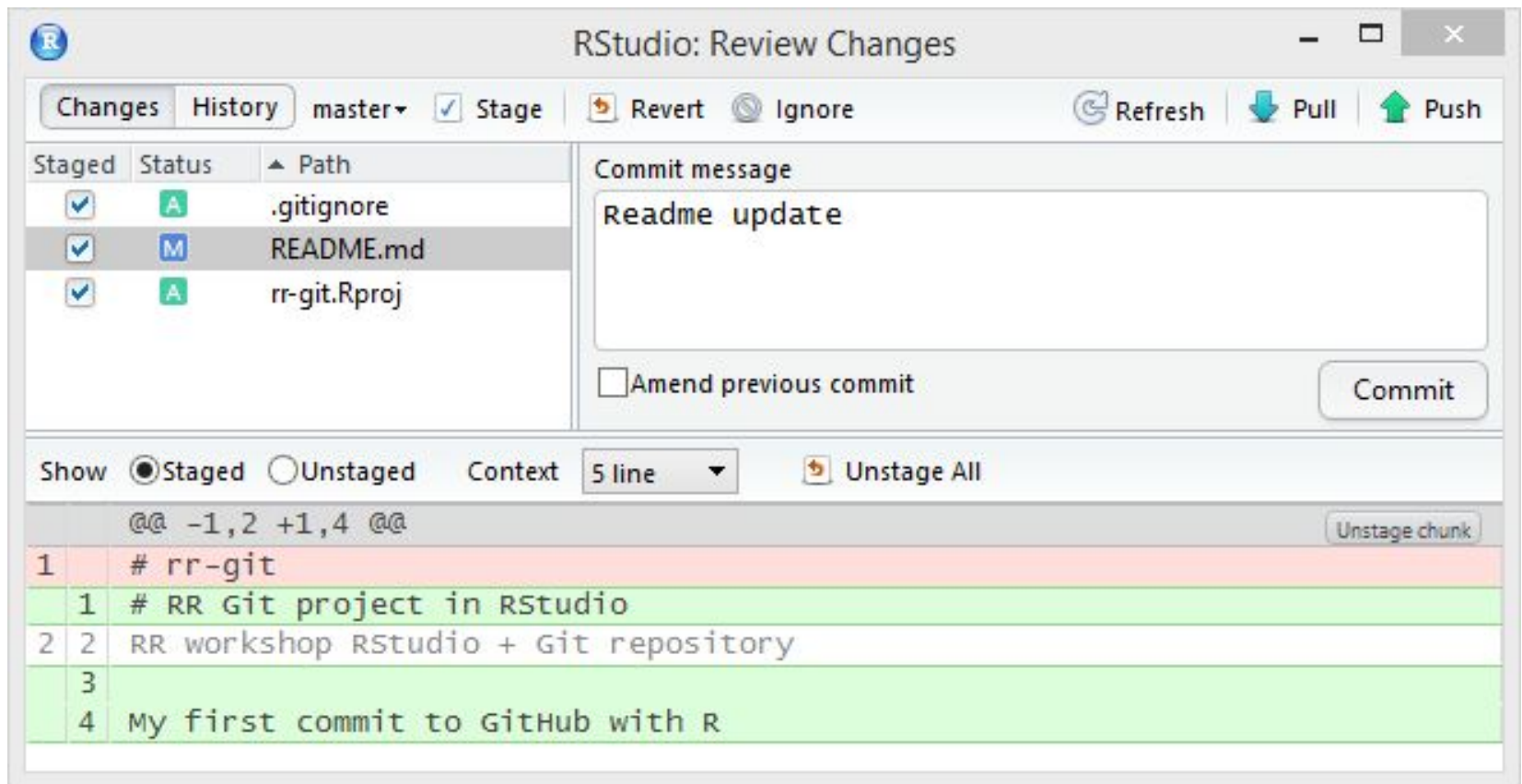
- Most popular
- Decentralised
- Designed for
- GitLab/GitHub Services



Git Workflow



Git is integrated into Rstudio!



The screenshot shows the RStudio 'Review Changes' window. At the top, there are tabs for 'Changes' and 'History', and a dropdown menu set to 'master'. Below these are buttons for 'Stage', 'Revert', and 'Ignore'. On the right side, there are buttons for 'Refresh', 'Pull', and 'Push'. The main area is divided into two panes. The left pane shows a table of staged files:

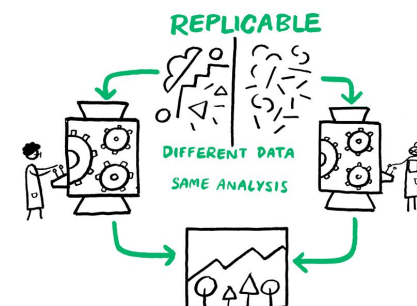
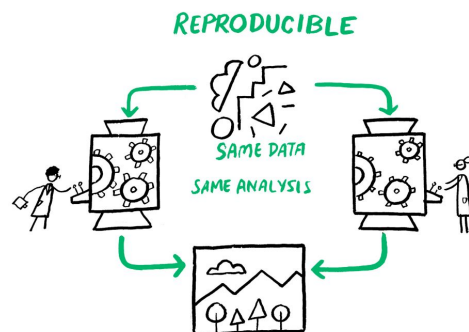
Staged	Status	Path
<input checked="" type="checkbox"/>	A	.gitignore
<input checked="" type="checkbox"/>	M	README.md
<input checked="" type="checkbox"/>	A	rr-git.Rproj

The right pane is for the commit message, with a text box containing 'Readme update' and a 'Commit' button. Below the commit message pane, there is a checkbox for 'Amend previous commit'. At the bottom of the window, there are controls for 'Show' (radio buttons for 'Staged' and 'Unstaged'), 'Context' (a dropdown menu set to '5 line'), and 'Unstage All'. Below these controls, there is a diff view showing changes to the README.md file:

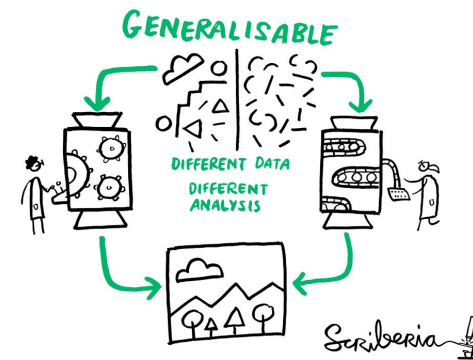
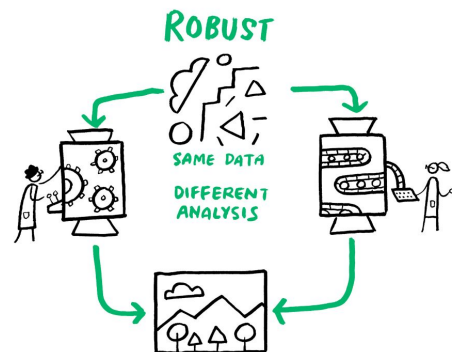
```
@@ -1,2 +1,4 @@
1 # rr-git
1 # RR Git project in RStudio
2 2 RR workshop RStudio + Git repository
3
4 My first commit to GitHub with R
```

Combine Git+Rmd Notebooks for Reproducibility

1. Add analysis to notebook
2. Add changes to git
3. Find out you made a mistake
4. Revert changes



1. Share notebook with collaborator
2. They make changes
3. You make changes
4. Merge changes into single analysis



Summary

- Overview of course: Database/EMR/Imaging/Signal
- Main assessments: practicals, journal article presentations, research proposal
- Data science is statistics with an EDA/Inductive/Data-focused Spin
- Health Data Science is a massive and growing area with lots of opportunity and challenges
- R is a powerful and useful tool for health data science
- Reproducibility is vital to good ~~health-data~~ science
- Rstudio, Rmarkdown notebooks and Git based version control facilitate that reproducibility

Friday's Practical

- Will go over the practical use of R, Rstudio, Rmd Notebooks, Git
- Try and install rstudio, git, and rmarkdown beforehand.
- 1st practical will not contribute to your course grade

Wednesday's Journal Articles

A Beginner's Guide to Conducting Reproducible Research

Jesse M. Alston, Jessica A. Rick First published: 15 January 2021 <https://doi.org/10.1002/bes2.1801>

Challenges to the Reproducibility of Machine Learning Models in Health Care

Andrew L. Beam, Arjun K. Manrai, Marzyeh Ghassemi <https://doi.org/10.1001/jama.2019.20866>

JAMA • 28 January 2020 • Volume 323, Number 4

Reproducibility in machine learning for health research: Still a ways to go

[Matthew B. A. McDermott](#) [Shirly Wang](#) [Nikki Marinsek](#) [Rajesh Ranganath](#) [Luca Foschini](#) [Marzyeh Ghassemi](#)

Science Translational Medicine • 24 Mar 2021 • Vol 13, Issue 586 • [DOI: 10.1126/scitranslmed.abb1655](https://doi.org/10.1126/scitranslmed.abb1655)