

Session 3: Data science for public health tools

Kirsty Lee Garson, Fatma Guerfali, Aquillah Kanzi, Amadou Diallo



UNIVERSITY OF CAPE TOWN
IYUNIVESITHI YASEKAPA • UNIVERSITEIT VAN KAAPSTAD



**wellcome
connecting
science**



**COVID-19
GENOMICS
GLOBAL TRAINING**

Session outline

- Key components of data science, the role of data science in public health
- Why is sharing data and ensuring that data is FAIR so crucial in public health?
 - The growth in magnitude and complexity of disease surveillance data
 - Generating findable, accessible, interoperable and reusable data
- Public repositories, community standards and other resources which facilitate shared progress
- Applying open science principles in training and capacity development
- Group activity: scenarios for discussion



wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

Session outcomes

At the end of this session, participants will be able to:

- Understand the principles of open science and its capacity to enable higher quality, more equitable knowledge generation
- Describe the elements of FAIR data and the importance of ensuring that public health data can be accessed and integrated timeously
- Recognise the value of applying and sharing best practices, established protocols, and code/workflows
- Describe some of the challenges and considerations involved in data sharing
- Identify ways to improve datasets by locating and applying appropriate standards for metadata

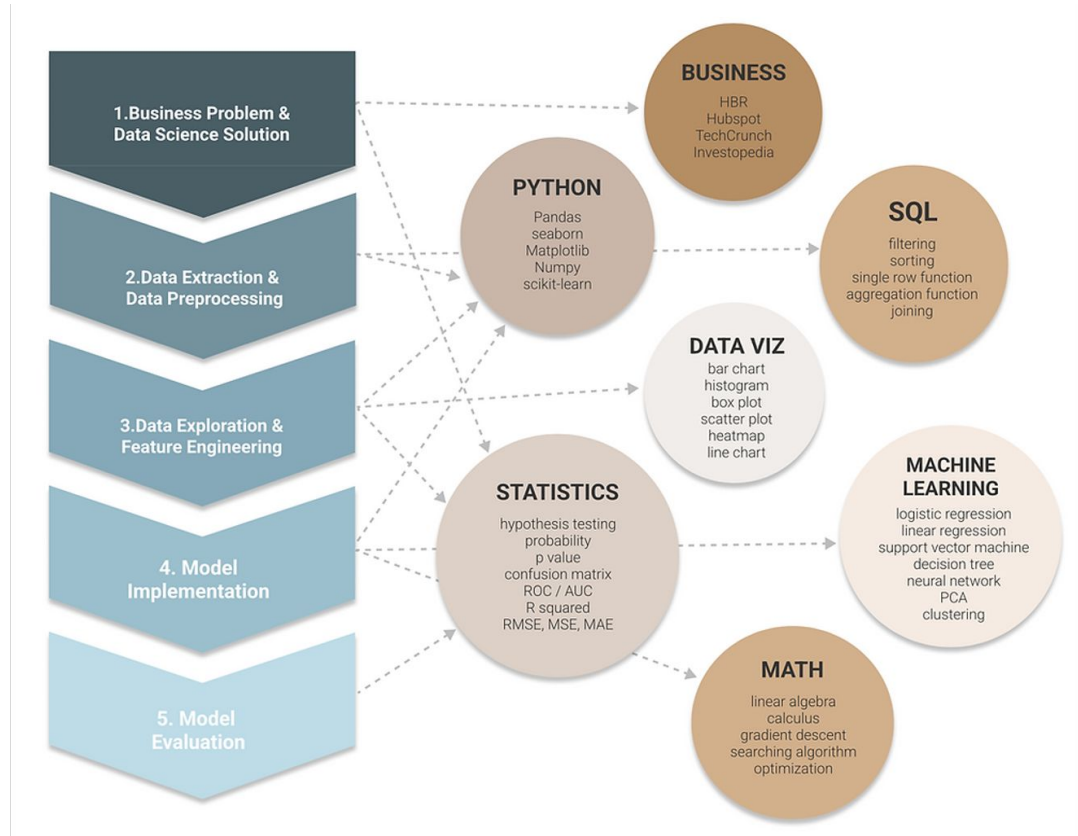


wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

Key components of Data Science

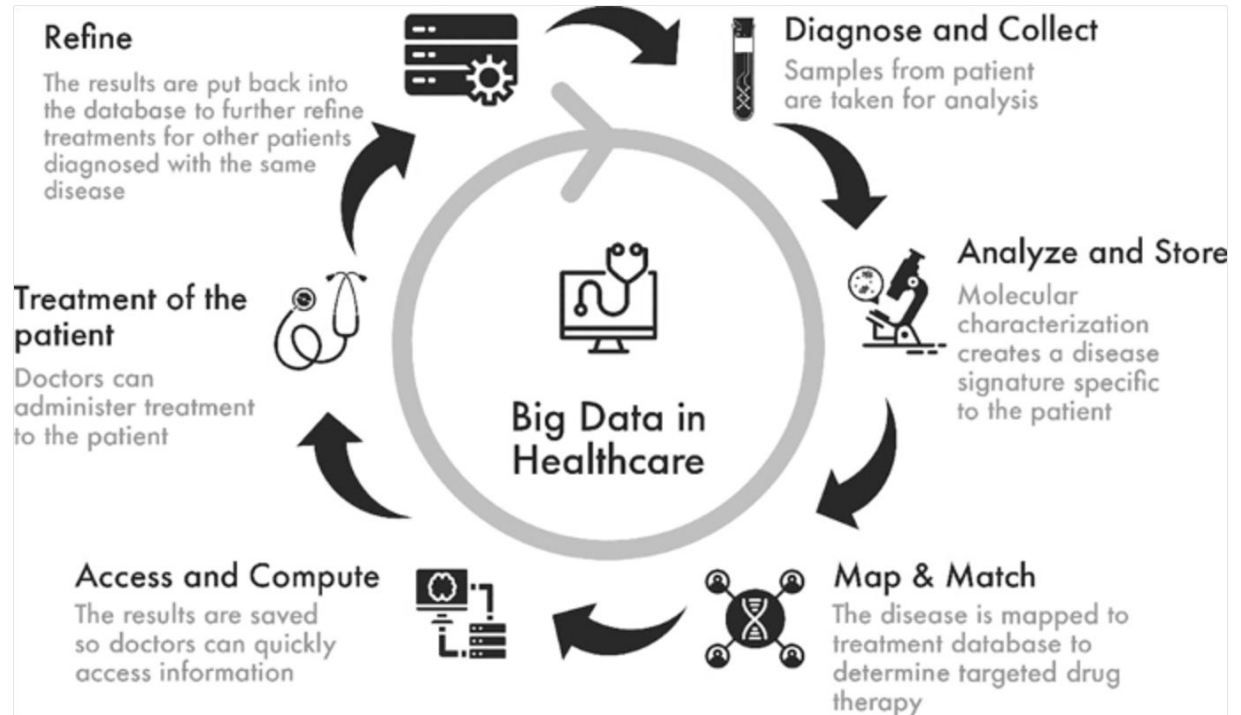


wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

Key components of Data Science in Public Health



Approaches to working with data in pathogen surveillance

In the past decade, the use of Next Generation Sequencing (NGS) has become increasingly applied in the areas of:

- Diagnostics
- Surveillance
- Research of infectious diseases

“Humans increasingly rely on computational support to deal with data as a result of the increase in volume, complexity, and creation speed of data.”

The value of any disease surveillance dataset can't be fully realised apart from **current and historical data for the regional, national and global context**

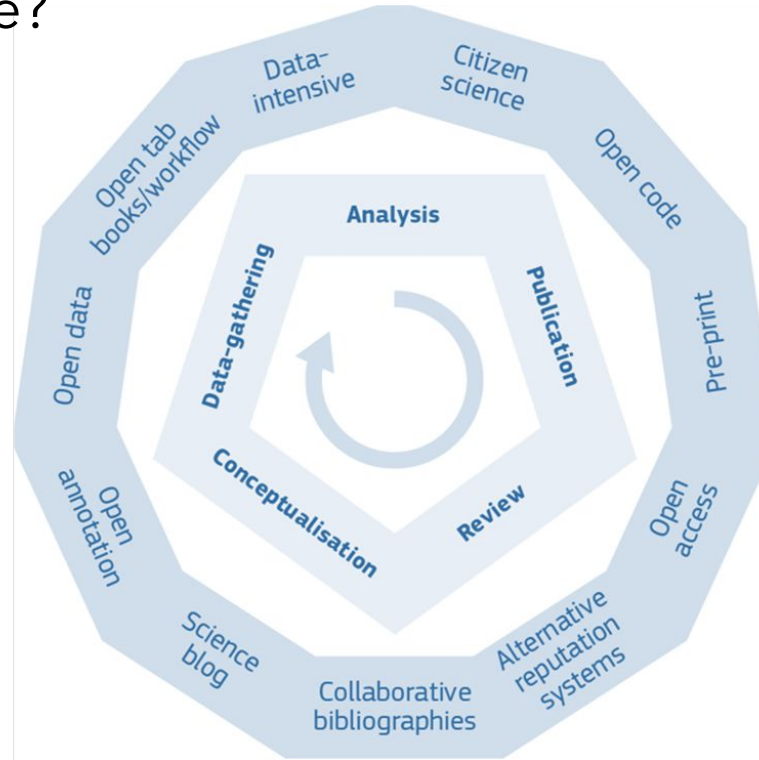
The ability to quickly access and re-use data is **vital**

Findable
Accessible
Interoperable
Reproducible



Why Open Science?

Quality
Efficiency
Reproducibility



Democratizing
access and
participation

<https://www.fosteropenscience.eu/learning/what-is-open-science/#/>



**wellcome
connecting
science**



**COVID-19
GENOMICS
GLOBAL TRAINING**

Developing consensus on best practices,
sharing protocols, making findings
comparable across environments



protocols.io

A secure platform for developing
and sharing reproducible methods

<https://www.protocols.io/>



**Public Health Alliance for
Genomic Epidemiology**

Open. Consensus. Practice. Tools.

**Improving Openness And
Interoperability In Public Health
Bioinformatics**

A Global Coalition.

An international consortium working to
develop shared standards in public health



**wellcome
connecting
science**



**COVID-19
GENOMICS
GLOBAL TRAINING**

The importance of informative and harmonizable metadata

PHA4GE SARS-CoV-2 Contextual Data Specification Package:

<https://github.com/pha4ge/SARS-CoV-2-Contextual-Data-Specification>

SARS-CoV-2 contextual data specification package

Spreadsheet-based (.xlsx) collection template

It contains the following items (tabs in the spreadsheet):

1. **a template for populating the complete set of contextual data;**

The collection template contains "required" (colour-coded yellow), "strongly recommended" (colour-coded purple) and "optional" (colour-coded white) fields.

2. **guidance for populating the template;**

The reference guide aims to facilitate the use of the collection template. It contains field definitions, further guidance/instructions, and examples of structured data.

3. **ontology-mapped controlled vocabulary for the picklists.**

Lists of controlled vocabulary, agreed upon by PHA4GE, are provided here for populating the template.



The importance of informative and harmonizable metadata

Example of metadata standards for public health:

H3ABioNet Case Report Form Standards:

<https://h3abionet.org/data-standards/phenotype-data-collection-standard>



- COVID-19: REDCap Project Template; Data Dictionary; User Guideline



wellcome
connecting
science



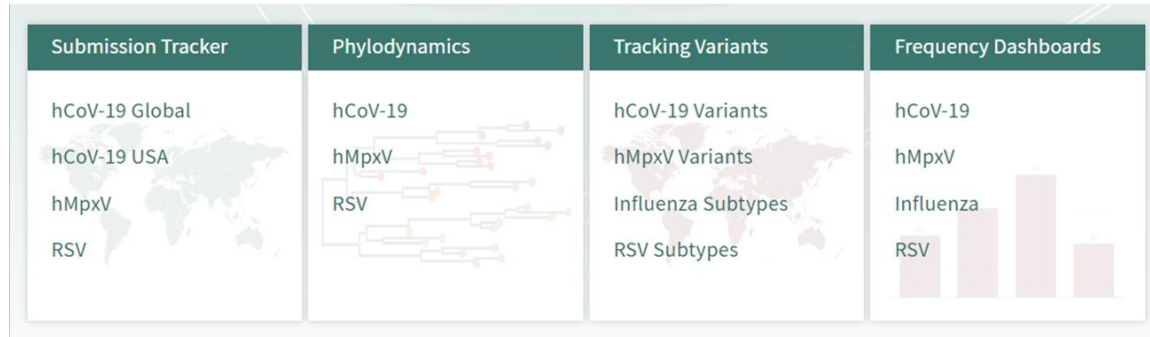
COVID-19
GENOMICS
GLOBAL TRAINING

Data repositories and tools for disease surveillance data



"Rapid sharing of data from all influenza viruses and the coronavirus causing COVID-19"

*Emphasis: "acknowledging the **Originating Laboratories** providing the specimens, and the **Submitting Laboratories** generating sequence and other metadata"*



<https://gisaid.org/>



wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

Data repositories and tools for disease surveillance data

International Nucleotide Sequence Database Collaboration (INSDC)

NCBI

ENA
European Nucleotide Archive

DDBJ
DNA Data Bank of Japan

Data type	DDBJ	EMBL-EBI	NCBI
Next Generation reads	Sequence Read Archive	European Nucleotide Archive (ENA)	Sequence Read Archive
Assembled Sequences	DDBJ		GenBank
Samples	BioSample		BioSample
Studies	BioProject		BioProject

Additional Resources for Data Science in Public Health

ViralAI, a global network for genomic surveillance and infectious disease research

<https://viral.ai/collections>

European COVID-19 Data Portal

<https://www.covid19dataportal.org/>

eLwazi Open Data Science Platform

<https://www.elwazi.org/>



wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

The value of training materials which are FAIR

NGS Academy for the Africa CDC Pathogen Genomics Initiative

Home About Courses Tools & Resources Contact

Other Pathogen Surveillance Courses

Show 10 entries Search: malaria

Organisation	Course Title	Course Level	Context	Audience	Primary focus areas	More
Wellcome Connecting Science	Malaria Experimental Genetics	Beginner	Research/academia	Bioinformaticians based at health institutes and experimental biologists at research institutions	Pathogen genomics for clinicians, Bioinformatics, Data analysis, data visualization, data management or a related topic	more

Home
About
Courses
Tools & Resources
Contact

Click [here](#) for the training survey

<https://uct-cbio.github.io/ngs-academy/courses>

<https://redcap.h3abionet.org/redcap/surveys/?s=FCA9MTKHPH>



wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

The value of training materials which are FAIR

 | **THE
CARPENTRIES**

**We teach foundational coding
and data science skills to
researchers worldwide.**

**DATA
CARPENTRY**

library
carpentry

software
carpentry

<https://carpentries.org/>



**wellcome
connecting
science**



**COVID-19
GENOMICS
GLOBAL TRAINING**

Africa PGI Data Management and Exchange Platform



Dr Gerald Mboowa

Africa Centres for Disease Control and Prevention

Christoffels, A *et al.* A pan-African pathogen genomics data sharing platform to support disease outbreaks. *Nat Med* (2023).

<https://doi.org/10.1038/s41591-023-02266-y>



wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

Group activity

Work with participants from your regional context to discuss the following questions



wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING

Group activity

Discussion:

- What concerns and considerations shape how, whether and where we share data? (consider the differences between sharing complete metadata vs. selected metadata variables vs. raw reads etc.)
- How can we work to facilitate more effective data sharing in our local contexts?
- How can we model, practice and negotiate for changes to facilitate open science practices in resource-limited environments?



wellcome
connecting
science



COVID-19
GENOMICS
GLOBAL TRAINING