



THE UNIVERSITY OF
MELBOURNE

Introduction to Data Science

Daniel Capurro, MD, PhD
Computing and Information Systems
Melbourne School of Engineering
Centre for Digital Transformation of Health - MDHS





Short introduction

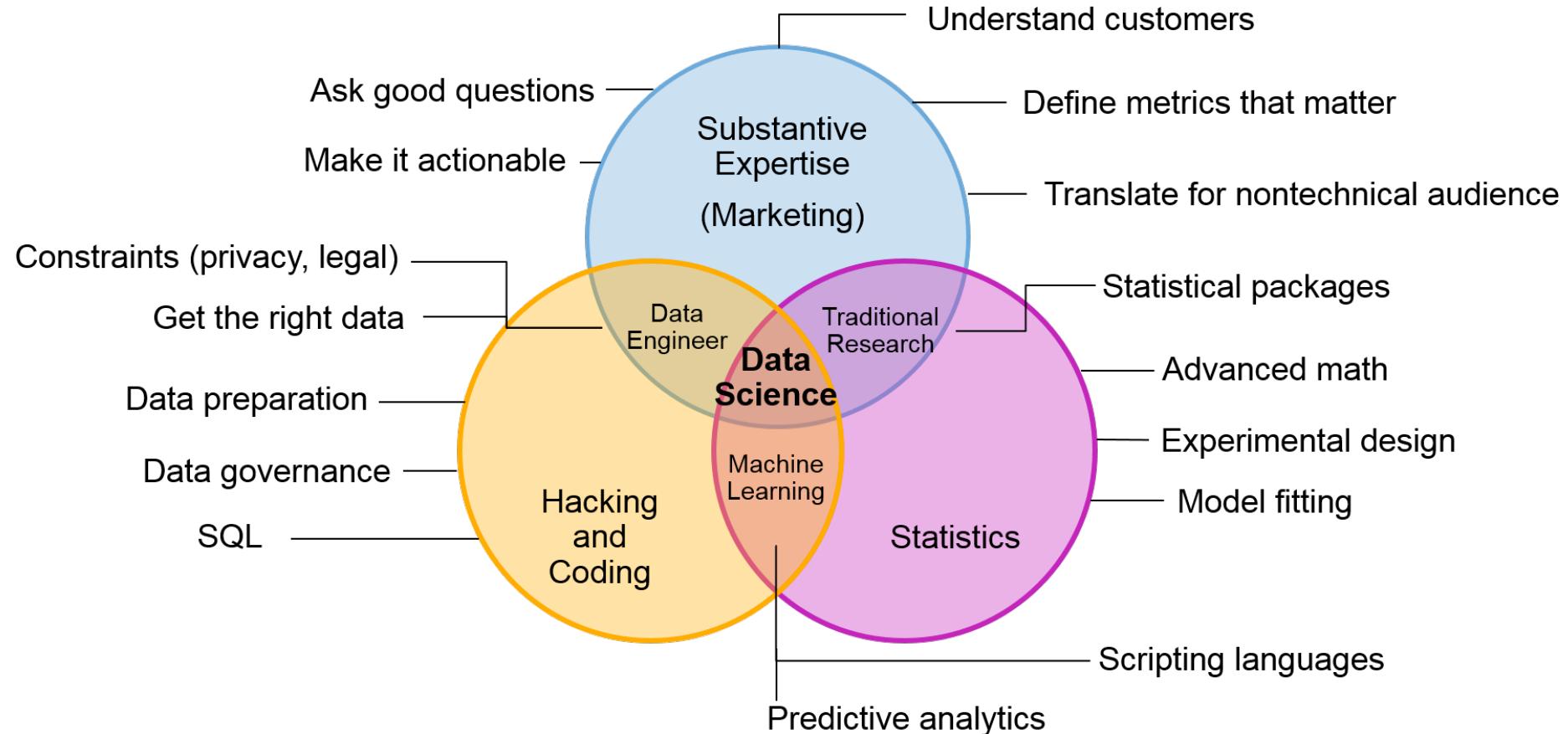
- Medical Doctor, practiced Internal Medicine up to a year ago
- PhD in Biomedical Informatics
- Senior Lecturer in Digital Health, School of Computing and Information Systems
- Centre for the Digital Transformation of Health





Definition...

- Not really a consensus definition...
- Intersection between
 - Statistics
 - Computer Science (Machine Learning, Big Data)
 - Business Analytics (in our case, healthcare analytics)
- Definition will probably evolve over time





Data Science a paradigm shift

Traditional research:

1. Hypothesis
2. Study design to answer that hypothesis
3. Collect data
4. Analyze data
5. Reject or accept your hypothesis, generate new hypothesis, etc.

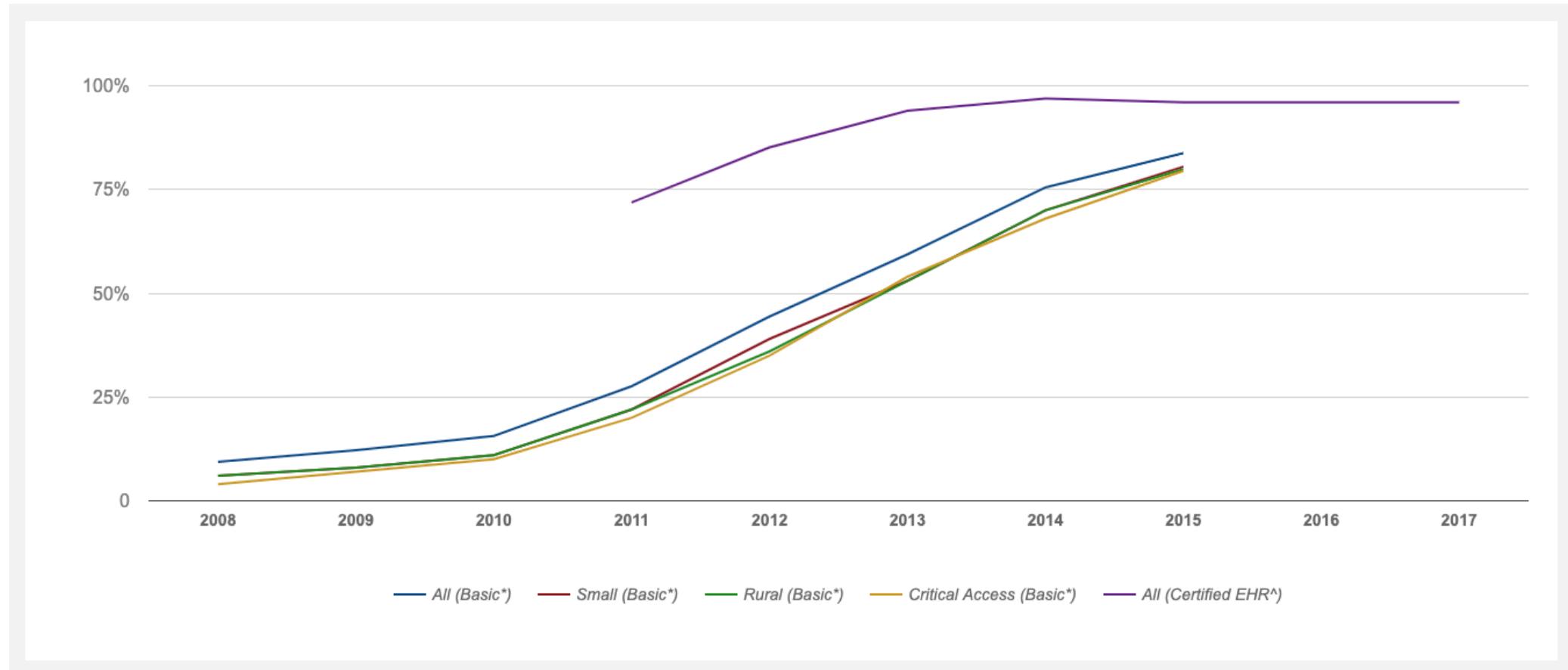


What enables the paradigm shift?

- Data is being routinely collected in information systems
- We continuously interact with ‘digital environments’

In healthcare this (mostly) means EHR adoption

US hospital adoption of EHRs





But also other sources of health/behaviour data

AJPH RESEARCH

Twitter as a Tool for Health Research: A Systematic Review

Lauren Sinnenberg, BA, Alison M. Buttenheim, PhD, MBA, Kevin Padrez, MD, Christina Mancheno, BA, Lyle Ungar, PhD, and Raina M. Merchant, MD, MSHP

Smartphone sensors





Special Collection: Smart Wellness Services for Lifestyles of Health and Sustainability (LOHAS)

INTERNATIONAL JOURNAL OF
**ENGINEERING BUSINESS
MANAGEMENT**

A real-time fall detection system based on the acceleration sensor of smartphone

*International Journal of Engineering
Business Management*

Volume 10: 1–8

© The Author(s) 2018

DOI: 10.1177/1847979017750669

journals.sagepub.com/home/enb



Youngmin Lee¹, Hongjin Yeh², Ki-Hyung Kim², and Okkyung Choi³ 



New paradigm: 2 options

1. Hypothesis-driven:

Hypothesis

Study design using previously collected data

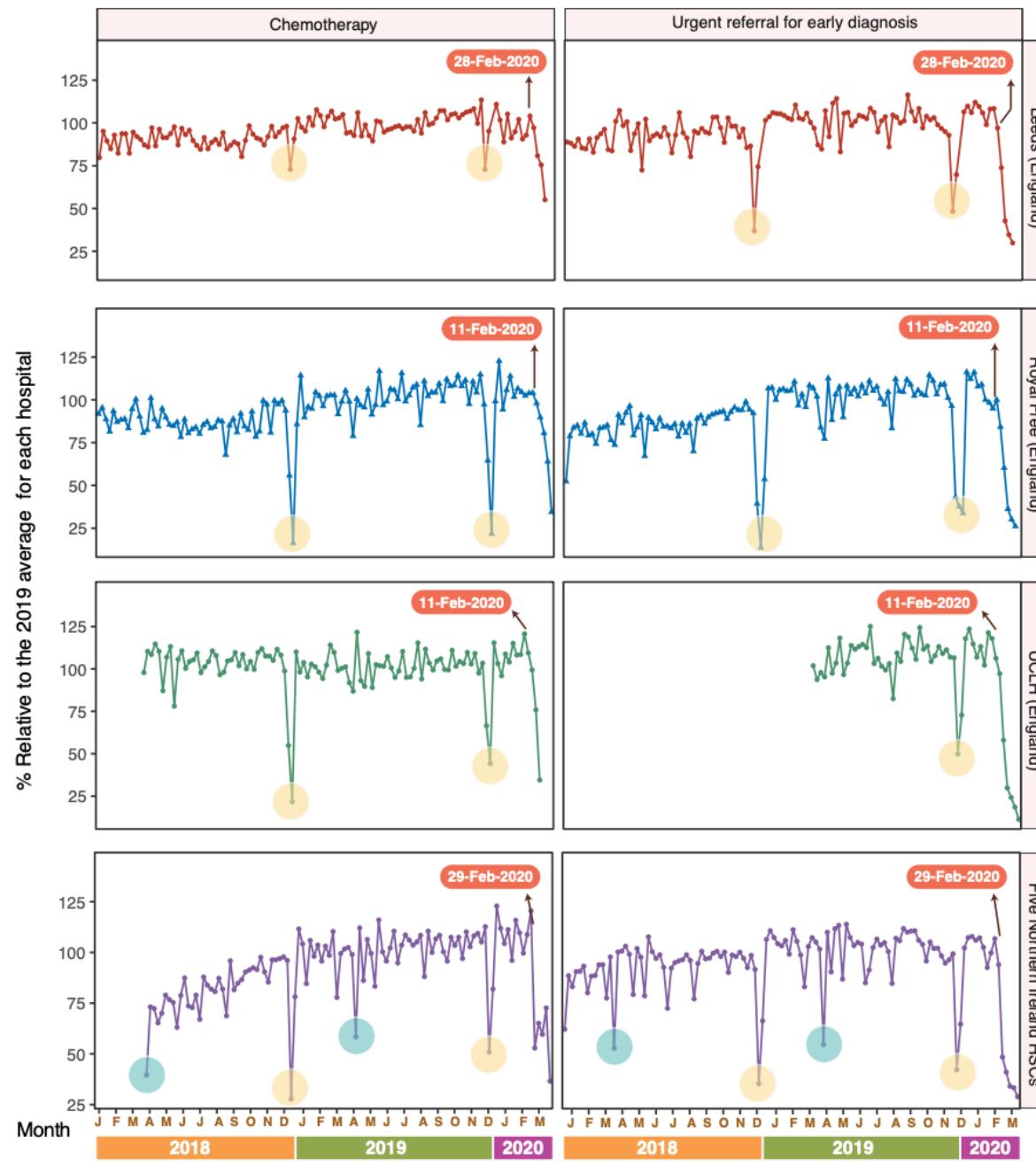
Data linkage, data preparation, data cleaning

Analyze data

Reject or accept your hypothesis, generate new hypothesis, etc.

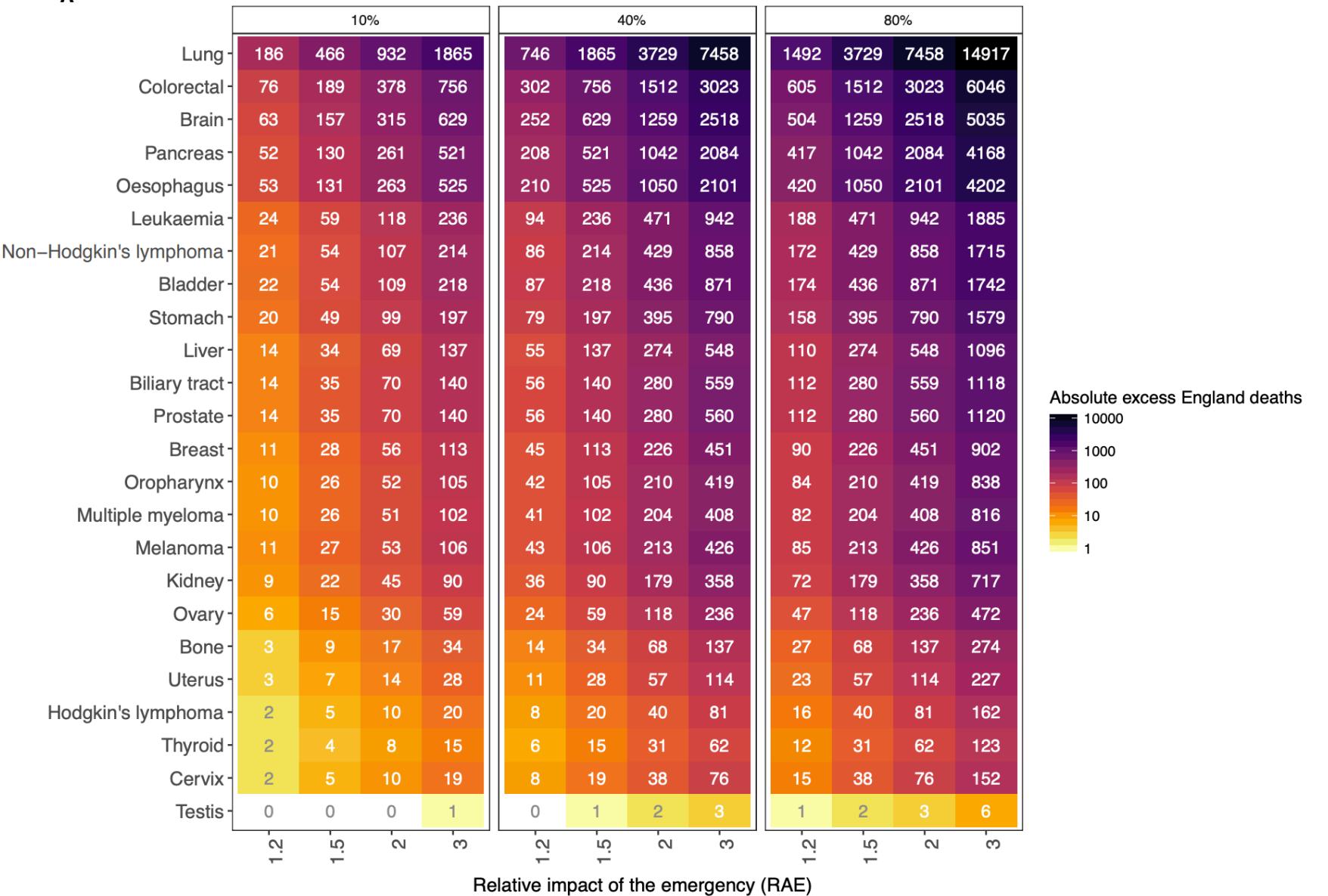
Estimating excess mortality in people with cancer and multimorbidity in the COVID-19 emergency

Alvina G. Lai, Ph.D.^{1,2,✉}, Laura Pasea, Ph.D.^{1,2*}, Amitava Banerjee, DPhil^{1,2,3*}, Spiros Denaxas, Ph.D.^{1,2,6,7}, Michail Katsoulis, Ph.D.^{1,2}, Wai Hoong Chang, MSc^{1,2}, Bryan Williams, Ph.D.^{4,5,6}, Deenan Pillay, Ph.D.⁸, Mahdad Noursadeghi, Ph.D.⁸, David Linch, FMedSci^{6,9}, Derralynn Hughes, FRCPath^{10,11}, Martin D. Forster, Ph.D.^{4,10}, Clare Turnbull, Ph.D.¹², Natalie K. Fitzpatrick, MSc^{1,2}, Kathryn Boyd, MD¹³, Graham R. Foster, Ph.D.¹⁴, DATA-CAN¹⁵, Matt Cooper, Ph.D.¹⁵, Monica Jones, PGDip¹⁵, Kathy Pritchard-Jones, FMedSci^{15,16,17,18}, Richard Sullivan, Ph.D.¹⁹, Geoff Hall, Ph.D.^{15,20,21}, Charlie Davie, FRCP^{11,15,16}, Mark Lawler, Ph.D.^{15,22}, and Harry Hemingway, FMedSci^{1,2,6,✉}



A

Proportion of the population who are affected (PAE) by the COVID-19 emergency





New paradigm: 2 options

2. Data-driven

(Usually) No hypothesis

Data linkage, data preparation, data cleaning

Identify patterns (knowledge) in databases

Supervised and unsupervised methods

New hypothesis that needs to be tested

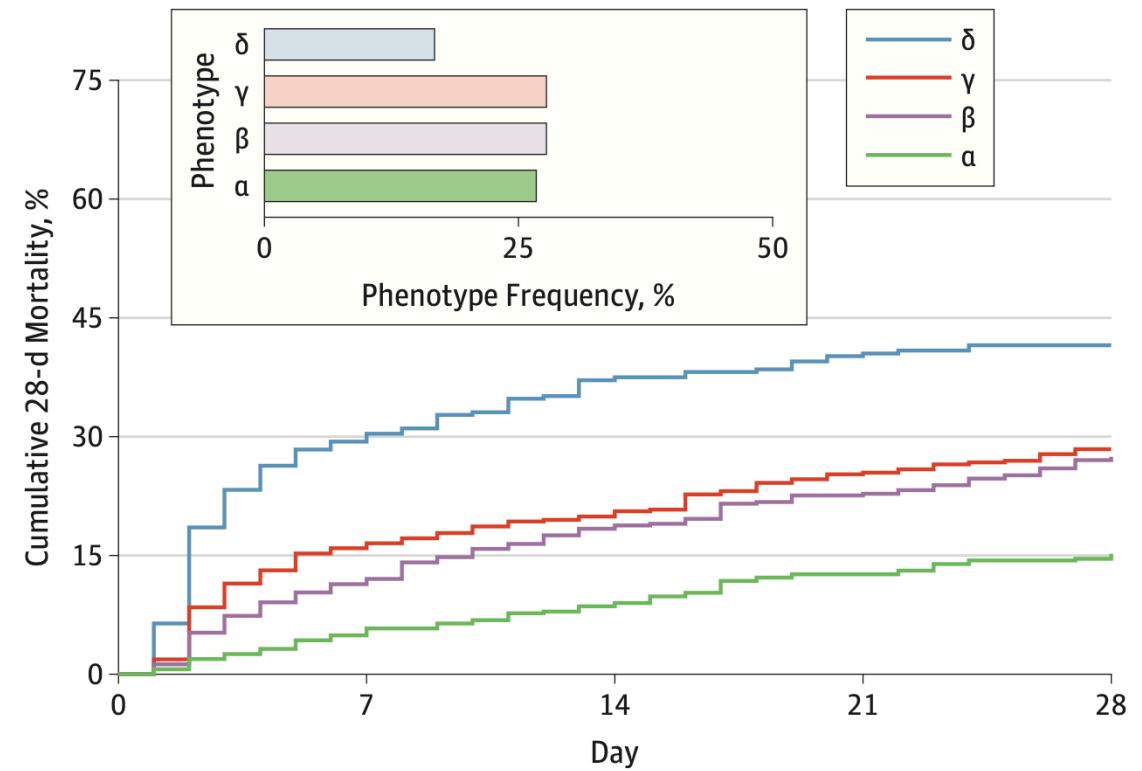
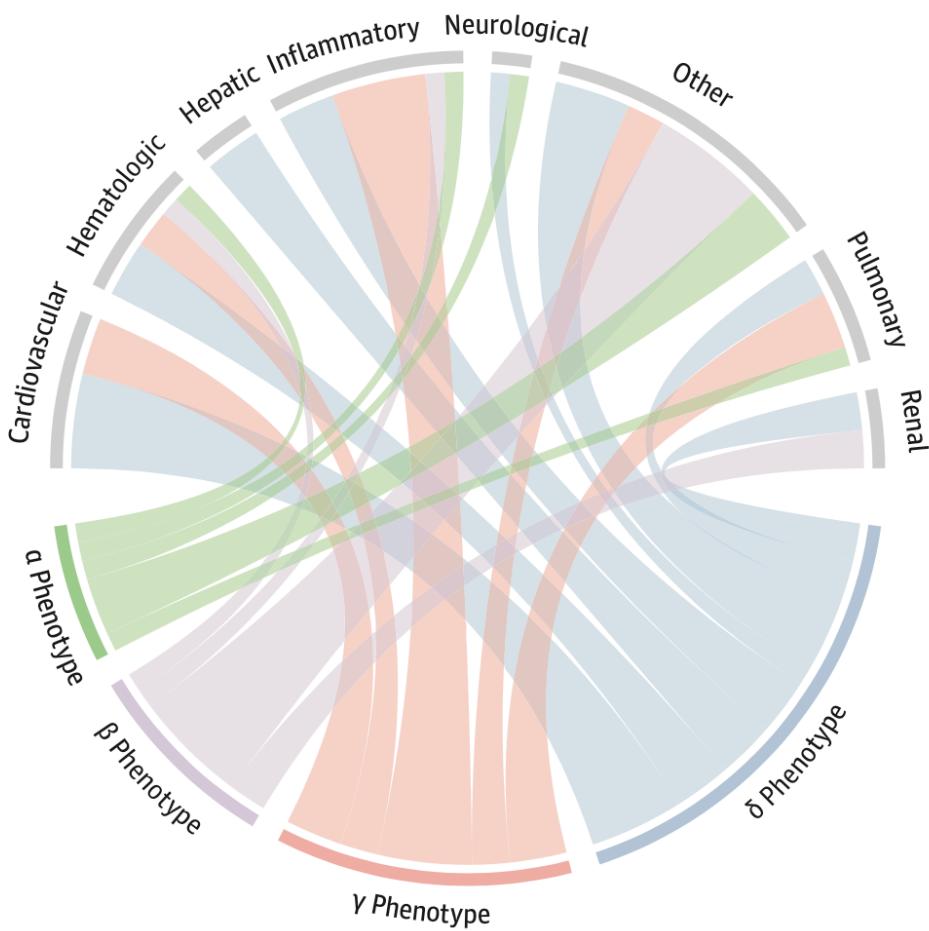


JAMA | Original Investigation | CARING FOR THE CRITICALLY ILL PATIENT

Derivation, Validation, and Potential Treatment Implications of Novel Clinical Phenotypes for Sepsis

Christopher W. Seymour, MD, MSc; Jason N. Kennedy, MS; Shu Wang, MS; Chung-Chou H. Chang, PhD; Corrine F. Elliott, MS; Zhongying Xu, MS; Scott Berry, PhD; Gilles Clermont, MD, MSc; Gregory Cooper, MD, PhD; Hernando Gomez, MD, MPH; David T. Huang, MD, MPH; John A. Kellum, MD, FACP, MCCM; Qi Mi, PhD; Steven M. Opal, MD; Victor Talisa, MS; Tom van der Poll, MD, PhD; Shyam Visweswaran, MD, PhD; Yoram Vodovotz, PhD; Jeremy C. Weiss, MD, PhD; Donald M. Yealy, MD, FACEP; Sachin Yende, MD, MS; Derek C. Angus, MD, MPH

A All phenotypes combined





THE UNIVERSITY OF
MELBOURNE

How would that work?

(one) Data Science Pipeline

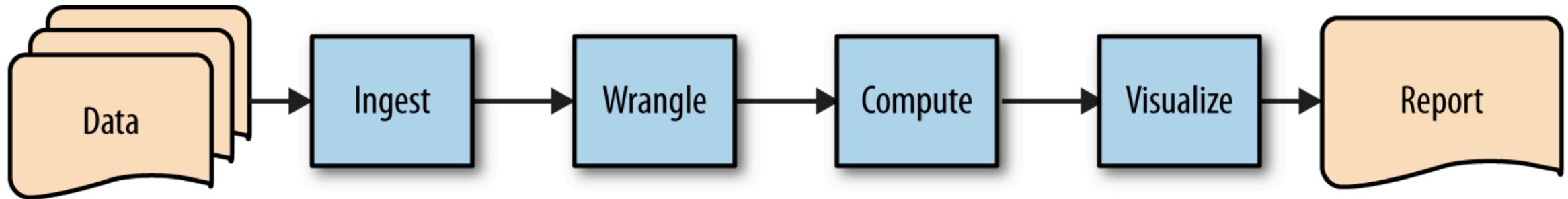
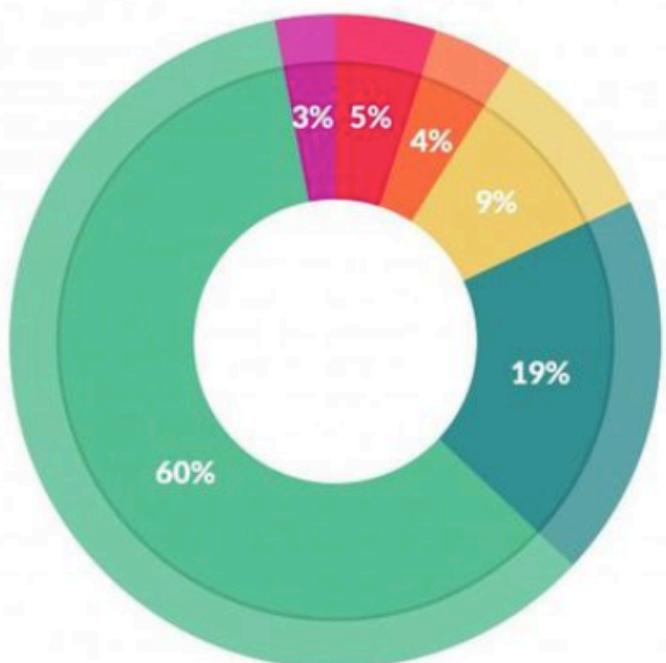


Figure 1-1. The data science pipeline

Data preparation

Data preparation accounts for about 80% of the work of data scientists



What data scientists spend the most time doing

- *Building training sets:* 3%
- *Cleaning and organizing data:* 60%
- *Collecting data sets:* 19%
- *Mining data for patterns:* 9%
- *Refining algorithms:* 4%
- *Other:* 5%



Some steps involved in data preparation

- Record linkage
- Variable Standardization - Labelling
- Out-of-range values
- Logical inconsistencies
- Variable transformation
- Discretization
- Missing Data
- Data reduction



Labelling clinical data is complex

Pt is 87 yo woman, highschool teacher with past medical history that includes

- status post cardiac catheterization in April 2019.

She presents today with palpitations and chest pressure.

HPI : Sleeping trouble on present dosage of Clonidine. Severe Rash on face and leg, slightly itchy

Meds : Vyvanse 50 mgs po at breakfast daily,

Clonidine 0.2 mgs -- 1 and 1 / 2 tabs po qhs

HEENT : Boggy inferior turbinates, No oropharyngeal lesion

Lungs : clear

Heart : Regular rhythm

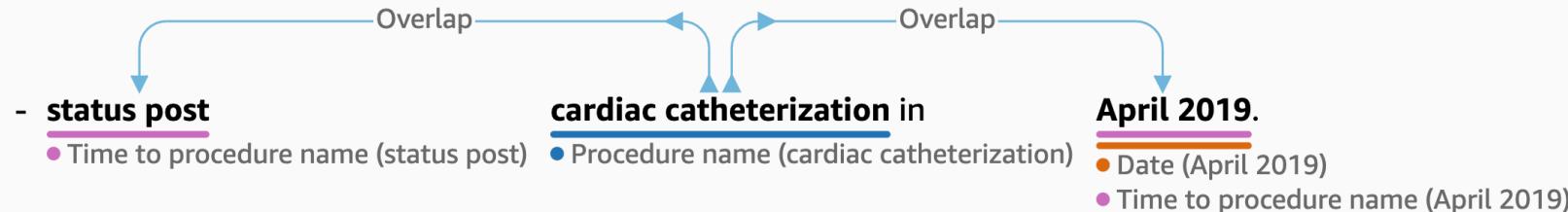
Skin : Mild erythematous eruption to hairline

Follow-up as scheduled

Clinical Natural Language Processing

Pt is **87** yo woman, **highschool teacher** with past medical history that includes

- Age (87)
- Profession (highschool teacher)



She presents **today** with

- Time to dx name (today)

palpitations and

- Dx name (palpitations)

chest pressure.

- Dx name (chest pressure)
- System organ site (chest)

HPI : **Sleeping trouble** on present dosage of

- Dx name (Sleeping trouble)

Clonidine. Severe

- Generic name (Clonidine)

Rash on

- Dx name (Rash)

face and

- System organ site (face)

leg, slightly

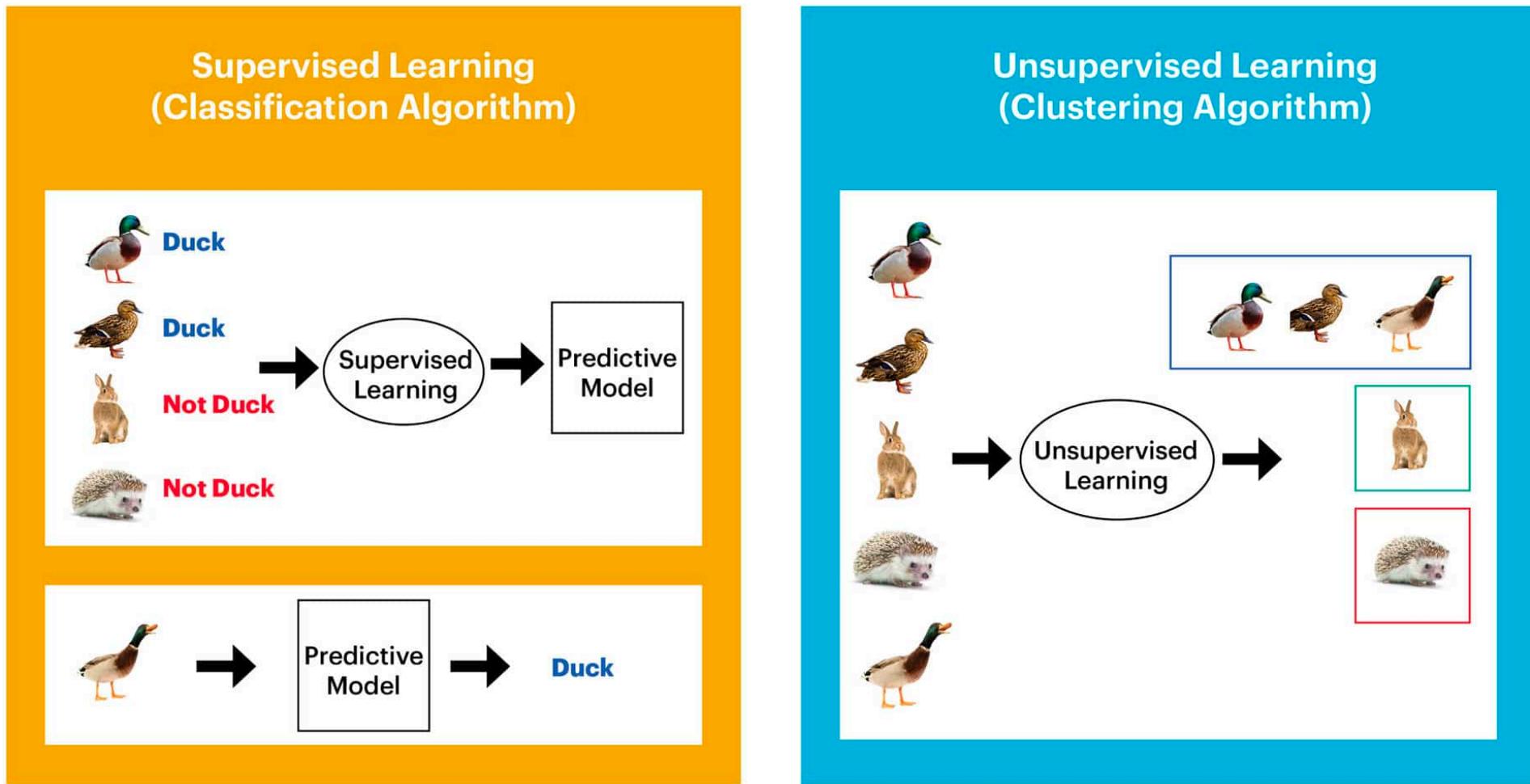
- System organ site (leg)

itchy

- Dx na

Symptom

Supervised vs. Unsupervised Machine Learning



Opportunities and Challenges

Some opportunities

- Accelerate discoveries
- Discovery of non-apparent patterns
- Rare diseases
- New data sources
- Real-time analysis
- Machine learning and Artificial Intelligence

Some challenges

- Technical
 - Data quality
 - Data standards
 - Big Data
 - Analytical methods
- Non-technical
 - Validation
 - Privacy
 - Bias
 - Unintended consequences