

## HDR UK Reproducible Machine Learning Project

### Collaborators Meeting

**Date of meeting: 9.30am-12noon on Wednesday 31 March 2021**

**Purpose of this meeting:** To share progress and identify new opportunities around reproducible machine learning in health data science.

#### Meeting Schedule:

Time	Item	Speakers
9:30	Why we need reproducible machine learning to support trustworthy clinical insight	<b>Aiden Doherty</b> (Oxford)
9:50	How should researchers report machine learning in health data science?	<b>Gary Collins</b> (Oxford)
10:05	What are the opportunities and challenges around the generation of synthetic healthcare datasets?	<b>Allan Tucker</b> (Brunel)
10:20	How can one conduct reproducible clinically-relevant machine learning in restricted 'safe-haven' environments?	<b>James Liley</b> (Edinburgh)
10:35	<b>Break (15mins)</b>	
10:50	Tools and open datasets to support training activities around reproducible machine learning: an activity monitoring use-case.	<b>Shing Chan</b> (Oxford)
11:05	What opportunities are there to embed reproducible machine learning across national training programmes?	<b>Thanasis Tsanas</b> (Edinburgh)
11:20	How can we learn from the public voice to positively contribute to reproducible machine learning in healthcare?	<b>Sophie Staniszewska</b> (Warwick)
11:35	How might outputs from this project be more discoverable by researchers across HDR UK?	<b>Susheel Varma</b> (HDR UK)
11:50	AOB and Closing Remarks	<b>Aiden Doherty</b>
12:00	<b>Meeting Ends</b>	

**Attended:**

	<b>First Name</b>	<b>Surname</b>	<b>Institute</b>
1	<b>Aiden</b>	<b>Doherty (Chair)</b>	<b>University of Oxford</b>
2	<b>Bella</b>	<b>Pratt (Administration)</b>	<b>University of Oxford</b>
3	Louis	Aslett	Durham University / The Alan Turing Institute
4	Shing	Chan	University of Oxford
5	Peter	Charlton	University of Cambridge
6	Gary	Collins	University of Oxford
7	Emanuele	Di Angelantonio	University of Cambridge
8	Tim	Hubbard	King's College London
9	James	Liley	University of Edinburgh
10	Paolo	Missier	Newcastle University
11	Marcus	Munafo	University of Bristol
12	Hollydawn	Murray	HDR UK
13	Heidi	Seibold	Alan Turing Institute
14	Sophie	Staniszewska	University of Warwick
15	Darren	Treanor	University of Leeds
16	Thanasis	Tsanas	University of Edinburgh
17	Allan	Tucker	Brunel University
18	Catalina	Vallejos	Alan Turing Institute
19	Susheel	Varma	HDR UK
20	Sebastian	Vollmer	Alan Turing Institute

**The meeting slides will now follow in the order outlined in the programme above.**

# Why we need reproducible analytics to support trustworthy clinical insights

Aiden Doherty<sup>1</sup>, Chris Holmes<sup>2,3</sup>, Martin Landray<sup>1</sup>, Gary Collins<sup>4</sup>, Catalina Vallejos<sup>3,5</sup>, Sebastian Vollmer<sup>3</sup>, Emanuele Di Angelantonio<sup>6</sup>, Louis Aslett<sup>3</sup>, Alastair Denniston<sup>7</sup>, Verena Heise<sup>4</sup>, Bilal Mateen<sup>3</sup>, James Rudd<sup>6</sup>, Sophie Staniszewska<sup>7</sup>, Thanasis Tsanas<sup>5</sup>, Kirstie Whittaker<sup>3,6</sup>, Gabriella Rustici<sup>2</sup>, John Danesh<sup>6</sup>, Harry Hemingway<sup>8</sup>, Tim Hubbard<sup>8</sup>, Cathie Sudlow<sup>5</sup>

Shing Chan<sup>1</sup>,  
James Liley<sup>3,5</sup>,  
Arabella Pratt<sup>1</sup>

<sup>1</sup> Health Data Research UK Oxford <sup>2</sup> Health Data Research UK  
<sup>3</sup> Alan Turing Institute  
<sup>4</sup> University of Oxford  
<sup>5</sup> Health Data Research UK Scotland  
<sup>6</sup> Health Data Research UK Cambridge  
<sup>7</sup> Health Data Research UK Midlands  
<sup>8</sup> Health Data Research UK London

# Why machine learning could transform medicine

Data generation and data acquisition is becoming cheap

- driven by advances in digital measurement technologies
- genomes; images; health records; internet; wearables; EHRs
- BioBanks and longitudinal cohorts e.g. UK Biobank, 100,000 Genomes project

Coupled to increasing raw computing power (GPUs) that facilitate compute hungry algorithms

High level (governmental) recognition of AI and data as a resource

*Note: the actual methods, such as deep neural networks, are not too dissimilar to those used in the 1980s*



# Clinical areas where machine learning can add substantial value

Diagnosis

Risk prediction

Decision support

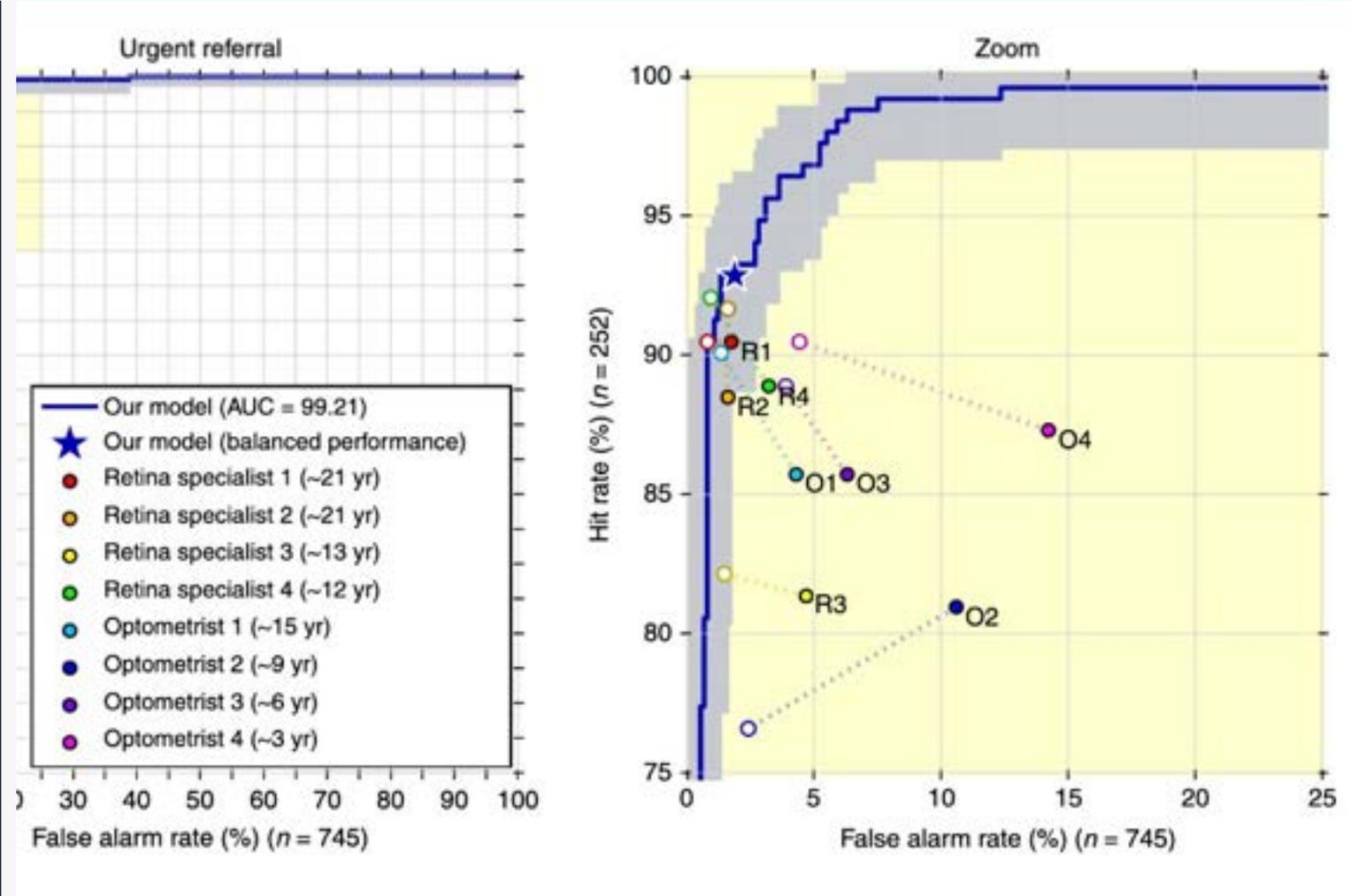
Typified by large quantities of rich, high-dimensional, and multi-modal data with clearly defined objectives. Areas where large data volumes make human decision analysis laborious.

Where we have cryptic “hard to define” features

- Imaging
- Electronic health records free text
- Wearables and consumer devices
- Other areas where handcrafted variables would be challenging and traditional statistical models difficult to elicit

# Machine learning offers promise to support diagnostic decisions

## Imaging - Diabetic retinopathy



# Machine learning offers promise for risk prediction

## Electronic health records - Acute kidney injury

Artificial intelligence + Add to myFT

### DeepMind creates algorithm to predict kidney damage in advance

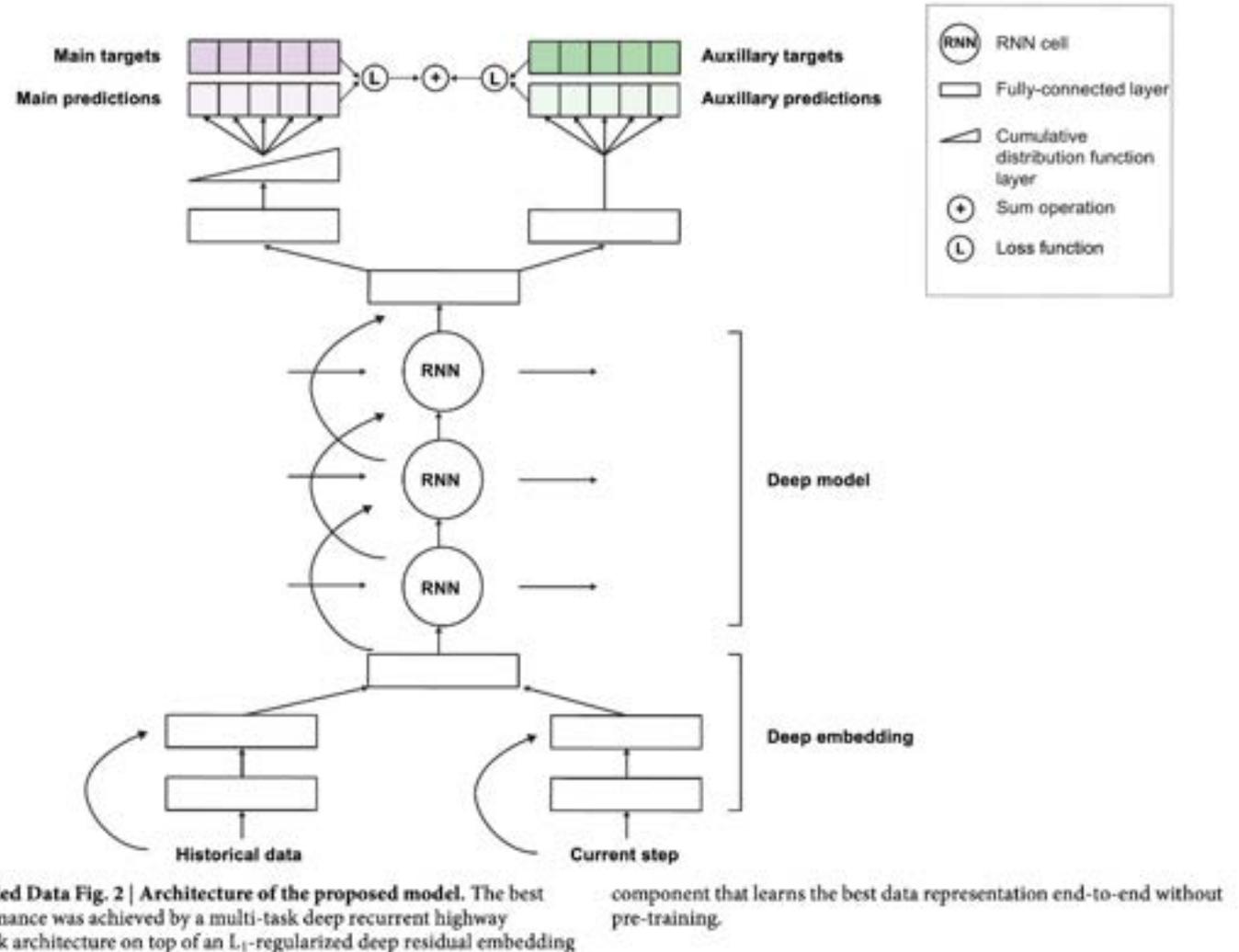
AI company's model can give 48 hours' warning of potentially fatal complications



The machine learning model was developed jointly by DeepMind Health, a division of Google, and the US Department of Veterans Affairs

Madhumita Murgia in London JULY 31 2019 28

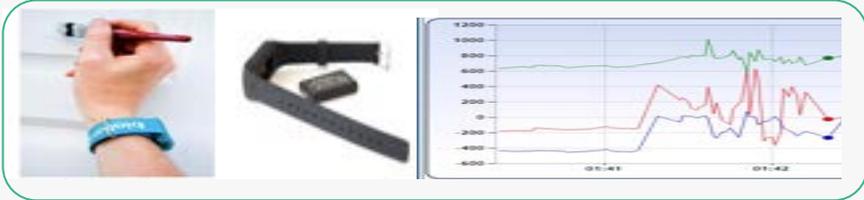
Artificial intelligence can now warn critical care doctors that their patients are at risk of developing severe kidney damage up to two days early, with the potential to save hundreds of thousands of lives every year.



# Machine learning can improve exposure measurement

## Physical activity & cardiovascular disease

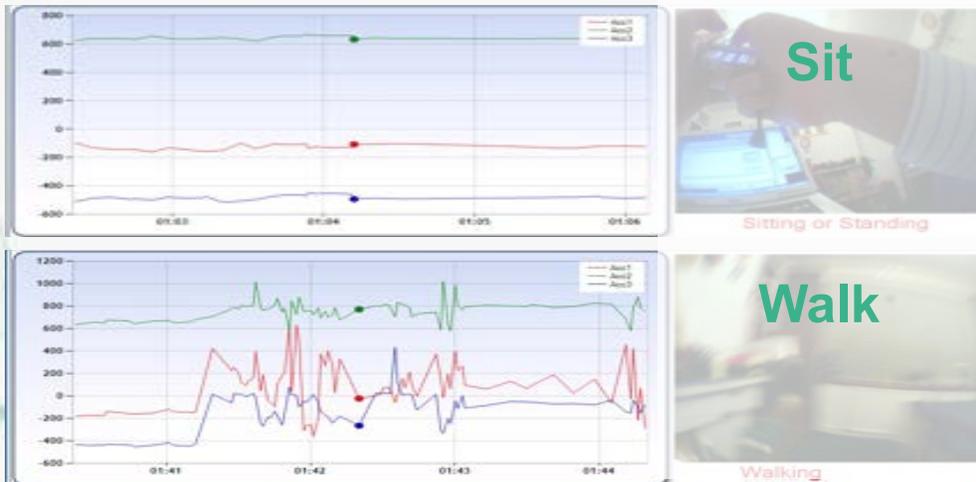
### Large health sensor datasets now available



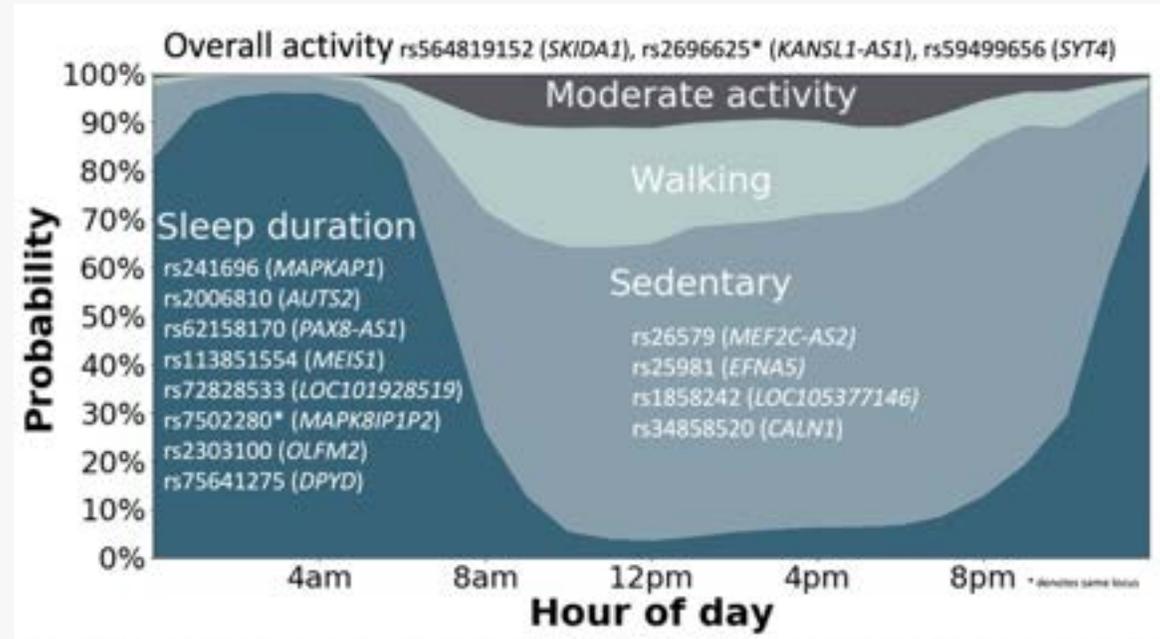
### UK Biobank - 100,000 people

- Range of health measurements
- Linkage to clinical outcomes

### Classify activities of daily living



### New insights on sleep & physical activity



New biological insight

Causal association with blood pressure & depression

Importance of activity has been underestimated

# Should I always use machine learning methods?



ELSEVIER



Journal of Clinical Epidemiology 110 (2019) 12–22

Journal of  
Clinical  
Epidemiology

HDRUK  
Health Data Research UK

## REVIEW

### A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models

Evangelia Christodoulou<sup>a</sup>, Jie Ma<sup>b</sup>, Gary S. Collins<sup>b,c</sup>, Ewout W. Steyerberg<sup>d</sup>,  
Jan Y. Verbakel<sup>a,e,f</sup>, Ben Van Calster<sup>a,d,\*</sup>

<sup>a</sup>Department of Development & Regeneration, KU Leuven, Herestraat 49 box 805, Leuven, 3000 Belgium

<sup>b</sup>Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, Botnar Research Centre, University of Oxford, Windmill Road, Oxford, OX3 7LD UK

<sup>c</sup>Oxford University Hospitals NHS Foundation Trust, Oxford, UK

<sup>d</sup>Department of Biomedical Data Sciences, Leiden University Medical Centre, Albinusdreef 2, Leiden, 2333 ZA The Netherlands

<sup>e</sup>Department of Public Health & Primary Care, KU Leuven, Kapucijnenvoer 33J box 7001, Leuven, 3000 Belgium

<sup>f</sup>Nuffield Department of Primary Care Health Sciences, University of Oxford, Woodstock Road, Oxford, OX2 6GG UK

Accepted 5 February 2019; Published online 11 February 2019

#### Abstract

**Objectives:** The objective of this study was to compare performance of logistic regression (LR) with machine learning (ML) for clinical prediction modeling in the literature.

**Study Design and Setting:** We conducted a Medline literature search (1/2016 to 8/2017) and extracted comparisons between LR and ML models for binary outcomes.

**Results:** We included 71 of 927 studies. The median sample size was 1,250 (range 72–3,994,872), with 19 predictors considered (range 5–563) and eight events per predictor (range 0.3–6,697). The most common ML methods were classification trees, random forests, artificial neural networks, and support vector machines. In 48 (68%) studies, we observed potential bias in the validation procedures. Sixty-four (90%) studies used the area under the receiver operating characteristic curve (AUC) to assess discrimination. Calibration was not addressed in 56 (79%) studies. We identified 282 comparisons between an LR and ML model (AUC range, 0.52–0.99). For 145 comparisons at low risk of bias, the difference in logit(AUC) between LR and ML was 0.00 (95% confidence interval, –0.18 to 0.18). For 137 comparisons at high risk of bias, logit(AUC) was 0.34 (0.20–0.47) higher for ML.

**Conclusion:** We found no evidence of superior performance of ML over LR. Improvements in methodology and reporting are needed for studies that compare modeling algorithms. © 2019 Elsevier Inc. All rights reserved.

**Keywords:** Clinical prediction models; Logistic regression; Machine learning; AUC; Calibration; Reporting

# Is the machine learning model robust to erroneous input?

*A small change to the input causes a large, highly confident, change to the prediction*

## The anatomy of an adversarial attack

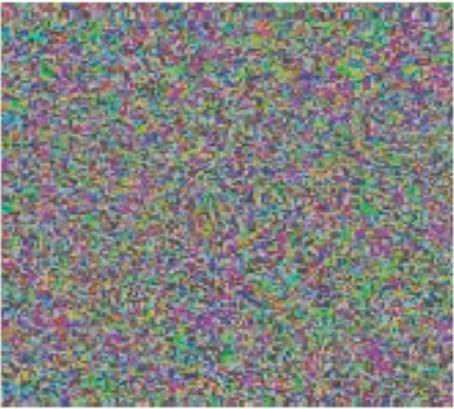
Demonstration of how adversarial attacks against various medical AI systems might be executed without requiring any overtly fraudulent misrepresentation of the data.

Original image



+ 0.04 ×

Adversarial noise

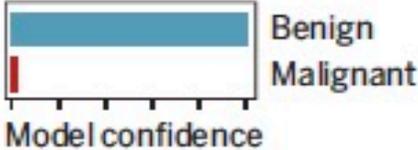


=

Adversarial example



Dermatoscopic image of a benign melanocytic nevus, along with the diagnostic probability computed by a deep neural network.



Perturbation computed by a common adversarial attack technique. See (7) for details.

Combined image of nevus and attack perturbation and the diagnostic probabilities from the same deep neural network.



# There is not an engrained culture of reproducible practices when using machine learning in medicine

- Reporting of machine learning methods is not standardised<sup>1</sup>
- Researchers using the open MIMIC dataset, for the same challenge, report vastly different sample sizes<sup>2</sup>
- While completing the same neuroimaging task, different analysis tools provide widely different results<sup>3</sup>
- Data sharing is rarely done and often restricted<sup>4</sup>

		Data	
		Same	Different
Analysis	Same	Reproducible	Replicable
	Different	Robust	Generalisable

New Online Views 3,710 Citations 0 Altmetric 131 Comments

Viewpoint

ONLINE FIRST FREE

January 6, 2020

## Challenges to the Reproducibility of Machine Learning Models in Health Care

Andrew L. Beam, PhD<sup>1,2</sup>; Arjun K. Manrai, PhD<sup>2,3</sup>; Marzyeh Ghassemi, PhD<sup>4,5</sup>

> Author Affiliations | Article Information

JAMA. Published online January 6, 2020. doi:<https://doi.org/10.1001/jama.2019.20866>

**R**eproducibility has been an important and intensely debated topic in science and medicine for the past few decades.<sup>1</sup> As the scientific enterprise has grown in scope and complexity, concerns regarding how well new findings can be reproduced and validated across different scientific teams and study populations have emerged. In some instances,<sup>2</sup> the failure to replicate numerous previous studies has added to the grow-

1 Rajkumar, A., Oren, E., Chen, K. et al. *npj Digital Med* 1, 18 (2018).

3 Bowring et al. *Hum Brain Mapp* 2019;40:3362-3384

5 Hutson *Science* 2018; 16 Feb 2018: Vol. 359, Issue 6377, pp. 725-726

2 Johnson, et al. *Machine Learning for Healthcare Conference*, 2017.

4 Miyakawa *Molecular Brain* (2020) 13:24

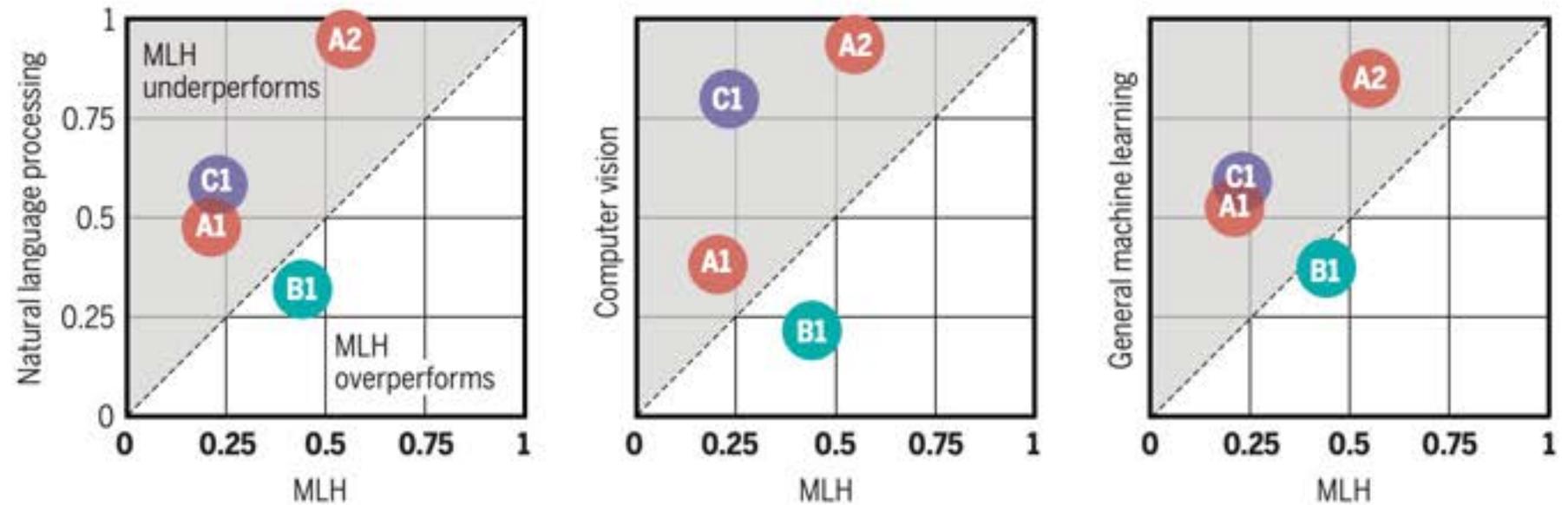
# Reproducible machine learning for health lags behind other fields

*An evaluation of 511 scientific papers across several machine learning subfields*

## SCIENCE TRANSLATIONAL MEDICINE | PERSPECTIVE

### Evaluation metrics

- A** Technical reproducibility
  - 1 Code available
  - 2 Public dataset
- B** Statistical reproducibility
  - 1 Variance reported
- C** Conceptual reproducibility (replicability)
  - 1 Multiple datasets



**Fig. 1. Reproducibility metrics for machine learning applications.** Shown are reproducibility metrics (A, B, and C) for evaluating scientific papers from four machine learning subspecialties: machine learning in health (MLH), natural language processing, computer vision, and general machine learning. Presented is the fraction of papers in a given subspecialty (y axis) versus those in MLH (x axis) that release their code (A1), release their data (A2), report their variance (B1), and leverage multiple datasets (C1). MLH consistently lags other subfields of machine learning on all measures of reproducibility apart from inclusion of proper statistical variance.

# Proposed solutions to improve reproducibility of ML in health



**Fig. 2. Improving reproducibility in MLH research.** Shown are eight recommendations for improving the reproducibility of machine learning in health research. These recommendations are subdivided according to which primary stakeholders directly drive these changes: the community of MLH researchers, health data providers (e.g., clinical care organizations), and journals and conferences (the primary publishers of machine learning research).

# National project – Reproducible machine learning

**Can reproducible machine learning be embedded in UK health data science to support trustworthy clinical insights?**

- 1) How should researchers report machine learning in health data science?
- 2) Can synthetic datasets be used to evaluate the stability of machine learning models in health data science?
- 3) What are the minimal requirements for reproducibility in restricted ‘safe-haven’ environments?
- 4) Can we initiate a culture of reproducible machine learning?

## Main achievements in first year

**Implementation project has provided flexibility to align with new opportunities around impactful driver projects**

WP1 reporting - TRIPOD-AI

WP2 synthetic data – wider HDR UK activities

WP3 safe haven environments - SPARRA

WP4 training - HDR UK Summer School & UK Biobank

PPIE – one of the first efforts in this space

Need to avoid temptation to ‘boil the ocean’. Machine learning and reproducibility are large fields in their own right.

*Meeting aim: To share progress and identify new opportunities around reproducible machine learning in health data science.*

<b>Time</b>	<b>Item</b>	<b>Speakers</b>
<b>9:30</b>	Why we need reproducible machine learning to support trustworthy clinical insight	<b>Aiden Doherty</b> (Oxford)
<b>9:50</b>	How should researchers report machine learning in health data science?	<b>Gary Collins</b> (Oxford)
<b>10:05</b>	What are the opportunities and challenges around the generation of synthetic healthcare datasets?	<b>Allan Tucker</b> (Brunel)
<b>10:20</b>	How can one conduct reproducible clinically-relevant machine learning in restricted 'safe-haven' environments?	<b>James Liley</b> (Edinburgh)
<b>10.35</b>	<b>Break (15mins)</b>	
<b>10:50</b>	Tools and open datasets to support training activities around reproducible machine learning: an activity monitoring use-case.	<b>Shing Chan</b> (Oxford)
<b>11.05</b>	What opportunities are there to embed reproducible machine learning across national training programmes?	<b>Thanasis Tsanas</b> (Edinburgh)
<b>11.20</b>	How can we learn from the public voice to positively contribute to reproducible machine learning in healthcare?	<b>Sophie Staniszewska</b> (Warwick)
<b>11.35</b>	How might outputs from this project be more discoverable by researchers across HDR UK?	<b>Susheel Varma</b> (HDR UK)
<b>11.50</b>	AOB and Closing Remarks	<b>Aiden Doherty</b>
<b>12:00</b>	<b>Meeting Ends</b>	

# Improving the reporting of clinical prediction models developed using machine learning (TRIPOD-AI)

Gary Collins

Professor of Medical Statistics  
Centre for Statistics in Medicine  
University of Oxford

# TRIPOD Statement (prediction models using regression)

- Published in Jan 2015, in 11 journals
- Focus on models developed using regression
  - However, guidance is relevant for ML
- **Explanation document (73 pages) focusses solely on examples from regression**
  - Touches on conduct
  - Opportunity to flag methodological issues
  - e.g., model presentation/availability
- **Needs to be tailored to the ML community**
  - e.g. examples, terminology, model presentation/availability
- **TRIPOD-AI launched [Collins & Moons, Lancet 2019]**

Annals of Internal Medicine RESEARCH AND REPORTING METHODS

## Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): The TRIPOD Statement

Gary S. Collins, PhD; Johannes B. Reitsma, MD, PhD; Douglas G. Altman, DSc; and Karel G.M. Moons, PhD

Prediction models are developed to aid health care providers in estimating the probability or risk that a specific disease or condition is present (diagnostic models) or that a specific event will occur in the future (prognostic models), to inform their decision making. However, the overwhelming evidence shows that the quality of reporting of prediction model studies is poor. Only with full and clear reporting of information on all aspects of a prediction model can risk of bias and potential usefulness of prediction models be adequately assessed. The Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD) Initiative developed a set of recommendations for the reporting of studies developing, validating, or updating a prediction model, whether for diagnostic or prognostic purposes. This article describes how the TRIPOD Statement was developed. An extensive list of items based on a review of the literature was created, which was reduced after a Web-based survey and revised during a 3-day meeting in June

2011 with methodologists, health care professionals, and journal editors. The list was refined during several meetings of the steering group and in e-mail discussions with the wider group of TRIPOD contributors. The resulting TRIPOD Statement is a checklist of 22 items, deemed essential for transparent reporting of a prediction model study. The TRIPOD Statement aims to improve the transparency of the reporting of a prediction model study regardless of the study methods used. The TRIPOD Statement is best used in conjunction with the TRIPOD explanation and elaboration document. To aid the editorial process and readers of prediction model studies, it is recommended that authors include a completed checklist in their submission (also available at [www.tripod-statement.org](http://www.tripod-statement.org)).

Ann Intern Med. 2015;162:55-63. doi:10.7326/M14-0697 [www.annals.org](http://www.annals.org)  
For author affiliations, see end of text.  
For contributors to the TRIPOD Statement, see the Appendix (available at [www.annals.org](http://www.annals.org)).

Annals of Internal Medicine RESEARCH AND REPORTING METHODS

## Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): Explanation and Elaboration

Karel G.M. Moons, PhD; Douglas G. Altman, DSc; Johannes B. Reitsma, MD, PhD; John P.A. Ioannidis, MD, DSc; Petra Macaskill, PhD; Ewout W. Steyerberg, PhD; Andrew J. Vickers, PhD; David F. Kanehl, MD; and Gary S. Collins, PhD

The TRIPOD (Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis) Statement includes a 22-item checklist, which aims to improve the reporting of studies developing, validating, or updating a prediction model, whether for diagnostic or prognostic purposes. The TRIPOD Statement aims to improve the transparency of the reporting of a prediction model study regardless of the study methods used. This explanation and elaboration document describes the rationale, clarifies the meaning of each item, and discusses why transparent reporting is important, with a view to assessing risk of bias and clinical usefulness of the prediction model. Each checklist item of the TRIPOD Statement is explained in detail and accom-

panied by published examples of good reporting. The document also provides a valuable reference of issues to consider when designing, conducting, and analyzing prediction model studies. To aid the editorial process and help peer reviewers and, ultimately, readers and systematic reviewers of prediction model studies, it is recommended that authors include a completed checklist in their submission. The TRIPOD checklist can also be downloaded from [www.tripod-statement.org](http://www.tripod-statement.org).

Ann Intern Med. 2015;162:W1-W73. doi:10.7326/M14-0699 [www.annals.org](http://www.annals.org)  
For author affiliations, see end of text.  
For members of the TRIPOD Group, see the Appendix.

# Importance of reporting

- Poor reporting limits (or prevents) the use of new research findings - avoidable waste
- Research reproducibility
  - ...impossible without complete reporting of all relevant aspects of scientific design, conduct, measurements, data and analysis
- Critical appraisal (risk of bias assessment)
  - Crucial for assessing study design and methods
- Readers need a clear understanding of exactly what was done
  - Clinicians, researchers, systematic reviewers, policy makers
  - Full and transparent reporting important for determining relevance and applicability
- Consequences of poor reporting include reducing or distorting the evidence base, usability of findings, and potential outcomes for clinical practice and ultimately to patients
- Huge number of reviews evaluating reporting of publications

# Why should we care?

## RESEARCH

 OPEN ACCESS

 Check for updates

 FAST TRACK

## Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal

Laure Wynants,<sup>1,2</sup> Ben Van Calster,<sup>2,3</sup> Gary S Collins,<sup>4,5</sup> Richard D Riley,<sup>6</sup> Georg Heinze,<sup>7</sup> Ewoud Schuit,<sup>8,9</sup> Marc M J Bonten,<sup>8,10</sup> Darren L Dahly,<sup>11,12</sup> Johanna A Damen,<sup>8,9</sup> Thomas P A Debray,<sup>8,9</sup> Valentijn M T de Jong,<sup>8,9</sup> Maarten De Vos,<sup>2,13</sup> Paula Dhiman,<sup>4,5</sup> Maria C Haller,<sup>7,14</sup> Michael O Harhay,<sup>15,16</sup> Liesbet Henckaerts,<sup>17,18</sup> Pauline Heus,<sup>8,9</sup> Michael Kammer,<sup>7,19</sup> Nina Kreuzberger,<sup>20</sup> Anna Lohmann,<sup>21</sup> Kim Luijken,<sup>21</sup> Jie Ma,<sup>5</sup> Glen P Martin,<sup>22</sup> David J McLernon,<sup>23</sup> Constanza L Andaur Navarro,<sup>8,9</sup> Johannes B Reitsma,<sup>8,9</sup> Jamie C Sergeant,<sup>24,25</sup> Chunhu Shi,<sup>26</sup> Nicole Skoetz,<sup>19</sup> Luc J M Smits,<sup>1</sup> Kym I E Snell,<sup>6</sup> Matthew Sperrin,<sup>27</sup> René Spijker,<sup>8,9,28</sup> Ewout W Steyerberg,<sup>3</sup> Toshihiko Takada,<sup>8</sup> Ioanna Tzoulaki,<sup>29,30</sup> Sander M J van Kuijk,<sup>31</sup> Bas C T van Bussel,<sup>1,32</sup> Iwan C C van der Horst,<sup>32</sup> Florian S van Royen,<sup>8</sup> Jan Y Verbakel,<sup>33,34</sup> Christine Wallisch,<sup>7,35,36</sup> Jack Wilkinson,<sup>22</sup> Robert Wolff,<sup>37</sup> Lotty Hooft,<sup>8,9</sup> Karel G M Moons,<sup>8,9</sup> Maarten van Smeden<sup>8</sup>

For numbered affiliations see end of the article.

Correspondence to: L Wynants  
laure.wynants@maastrichtuniversity.nl  
(ORCID 0000-0002-3037-122X)

Additional material is published online only. To view please visit the journal online.

Cite this as: *BMJ* 2020;369:m1328  
<http://dx.doi.org/10.1136/bmj.m1328>

### ABSTRACT OBJECTIVE

To review and appraise the validity and usefulness of published and preprint reports of prediction models for diagnosing coronavirus disease 2019 (covid-19) in patients with suspected infection, for prognosis of patients with covid-19, and for detecting people in the general population at increased risk of covid-19 infection or being admitted to hospital with the

### DATA SOURCES

PubMed and Embase through Ovid, up to 1 July 2020, supplemented with arXiv, medRxiv, and bioRxiv up to 5 May 2020.

### STUDY SELECTION

Studies that developed or validated a multivariable covid-19 related prediction model.

### DATA EXTRACTION

At least two authors independently extracted data

# Update 3 (1-July-2020)

- 169 studies describing 232 prediction models
  - 7 risk scores, 118 diagnostic; 107 prognostic
  - Mixture of modelling procedures
- Bottom line: 226 at high risk of bias; 6 at unclear risk of bias
- “This review indicates that almost all published prediction models are **poorly reported**, and at **high risk of bias** such that their reported predictive **performance is probably optimistic**. “

# Reporting of machine learning models\*

**The completeness of reporting and adherence to the TRIPOD Statement of clinical prediction models using machine learning methods in oncology: a systematic review**

Paula Dhiman<sup>1,2</sup>, Jie Ma<sup>1</sup>, Constanza Andaur Navarro<sup>3</sup>, Beni Speich<sup>1,4</sup>, Garrett Bullock<sup>5</sup>, Shona Kirtley<sup>1</sup>, Richard D Riley<sup>6</sup>, Ben Van Calster<sup>7</sup>, Karel GM Moons<sup>3</sup>, Gary S Collins<sup>1,2</sup>.

<sup>1</sup> Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, University of Oxford, Oxford, OX3 7LD, UK

<sup>2</sup> NIHR Oxford Biomedical Research Centre, Oxford University Hospitals NHS Foundation Trust, Oxford, United Kingdom

<sup>3</sup> Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht, The Netherlands

<sup>4</sup> Department of Clinical Research, Basel Institute for Clinical Epidemiology and Biostatistics, University Hospital Basel, University of Basel, Basel, Switzerland

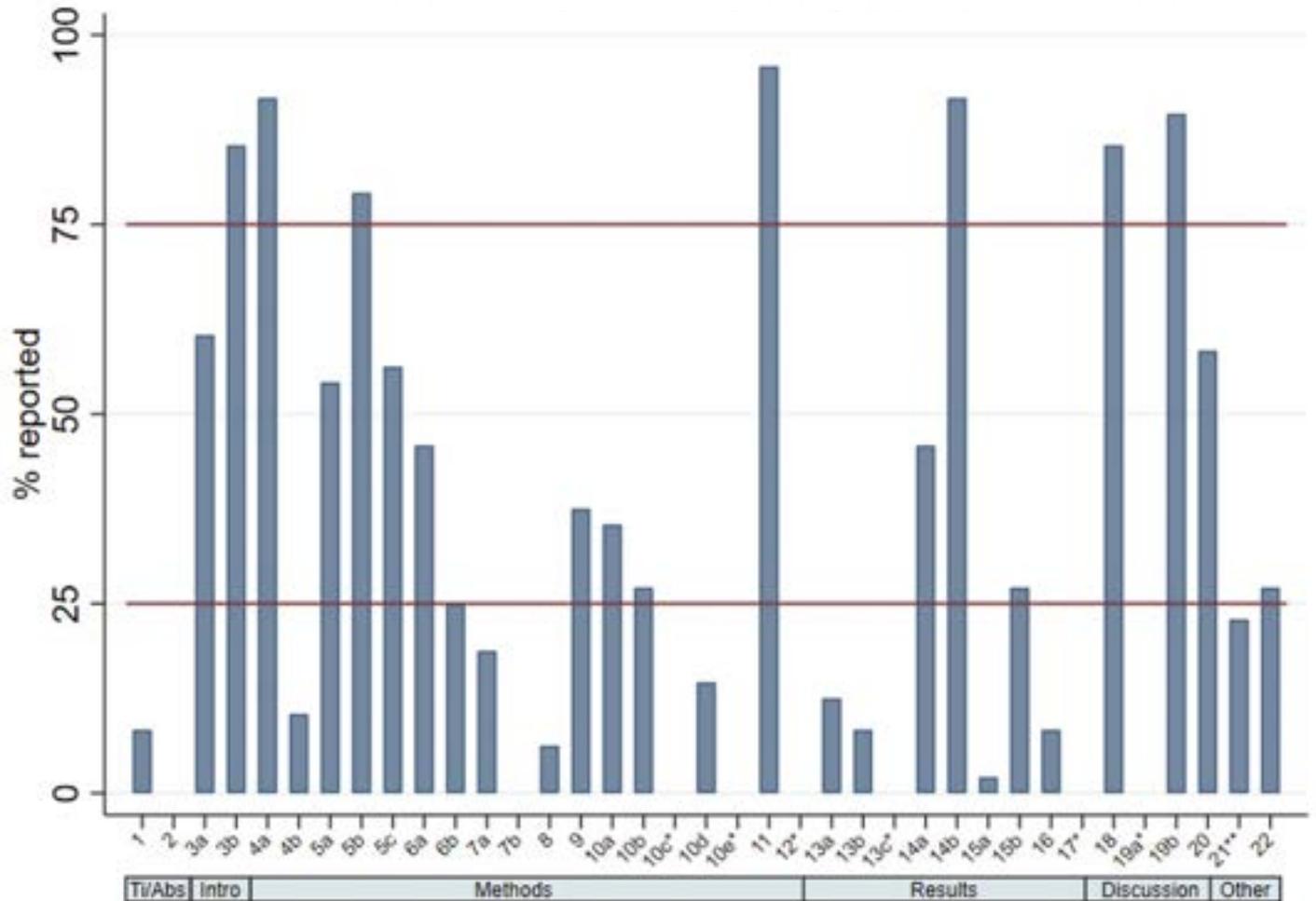
<sup>5</sup> Nuffield Department of Orthopaedics, Rheumatology, and Musculoskeletal Sciences, University of Oxford, Oxford, UK

<sup>6</sup> Centre for Prognosis Research, School of Primary, Community and Social Care, Keele University, Staffordshire, UK

\* Under review

# Reporting deficiencies

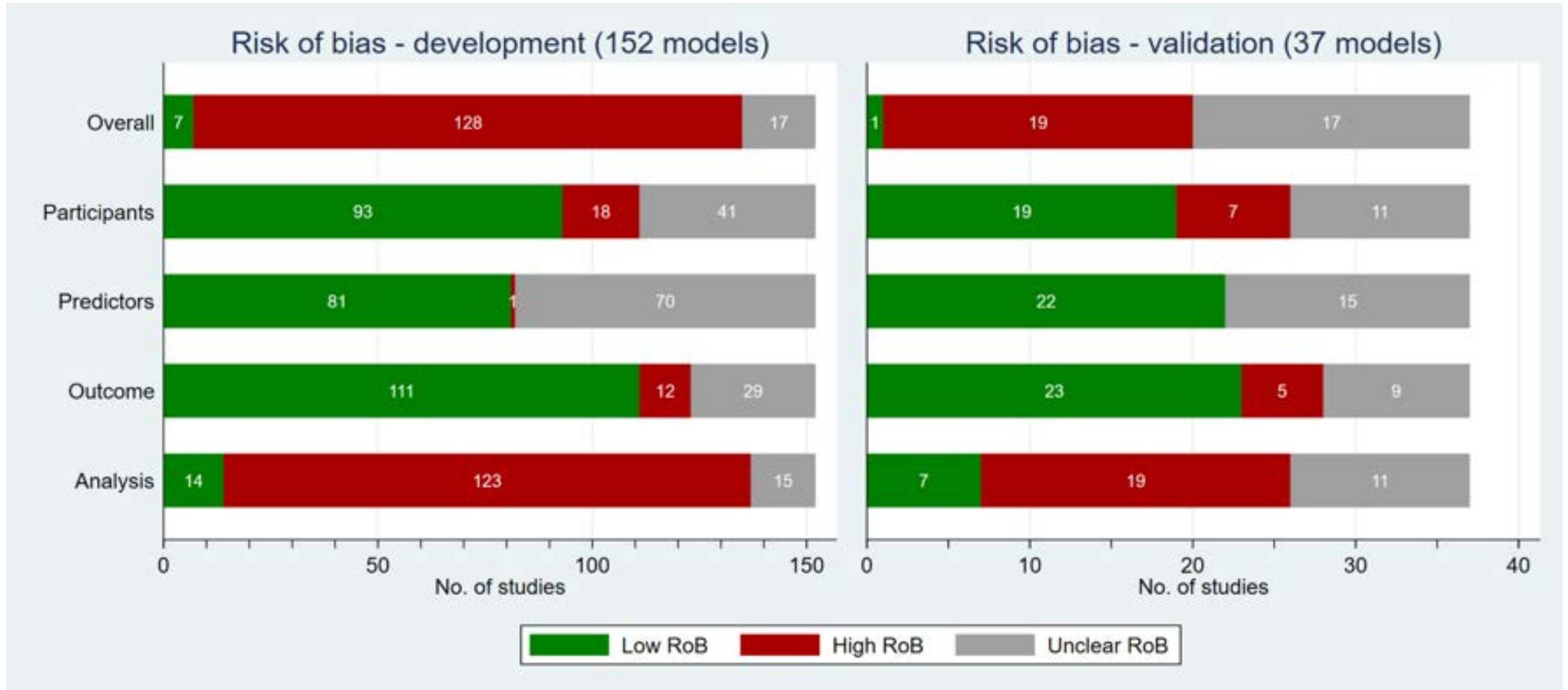
- Item 4b - study dates
- Item 8 - Sample size
- Item 10b - model building/internal validation
- Item 13b - characteristics of participants
- Item 15a - model availability
- Item 16 - performance measures with CIs



# TRIPOD adherence

Lower adherence ( $\leq 15\%$ )	Higher adherence ( $\geq 70\%$ )
<ul style="list-style-type: none"><li>x <b>Title</b></li><li>x <b>Abstract</b></li><li>x <b>Predictor assessment</b></li><li>x <b>Sample size</b></li><li>x <b>Description/comparison of between development and validation data</b></li><li>x <b>Participant flow and characteristics</b></li><li>x <b>Presentation and interpretation of full model</b></li><li>x <b>Model performance</b></li></ul>	<ul style="list-style-type: none"><li>✓ <b>Objectives</b></li><li>✓ <b>Source of data</b></li><li>✓ <b>Eligibility criteria</b></li><li>✓ <b>Reporting of risk groups</b></li><li>✓ <b>Unadjusted associations</b></li><li>✓ <b>Limitations</b></li><li>✓ <b>Overall interpretation of results</b></li></ul>

# Risk of bias – PROBAST



# TRIPOD for machine learning/AI

THE LANCET

Access provided by University of Oxford

COMMENT | VOLUME 393, ISSUE 10181, P1577-1579, APRIL 20, 2019

## Reporting of artificial intelligence prediction models

Gary S Collins  + Karel G M Moons

Published: April 20, 2019 · DOI: [https://doi.org/10.1016/S0140-6736\(19\)30037-6](https://doi.org/10.1016/S0140-6736(19)30037-6)  Check for updates

References

Article Info

Figures

Data-driven technologies that form the basis of the digital health-care revolution provide potentially important opportunities to deliver improvements in individual care and to advance innovation in medical research. Digital health technologies include mobile devices and health apps (m-health), e-health technology, and intelligent monitoring. Behind the digital health revolution are also methodological advancements using artificial intelligence and machine learning techniques. Artificial intelligence, which encompasses machine learning, is the scientific discipline that uses computer algorithms to learn from data, to help identify patterns in data, and make predictions. A key feature underpinning the excitement behind artificial intelligence and machine learning is their potential to analyse large and complex data structures to create prediction models that personalise and improve diagnosis, prognosis, monitoring, and administration of treatments, with the aim of improving individual health outcomes. Prediction models to support clinical decision making have existed for decades, and these include well known tools such as the Framingham Risk Score,<sup>1</sup> QRISK3,<sup>2</sup> Model for End-stage Liver Disease,<sup>3</sup> ABCD<sup>2</sup> score,<sup>4</sup> and the Nottingham Prognostic Index.<sup>5</sup> Health-care professionals, medical researchers, policy makers, guideline developers, patients, and members of the general public are all

# **A protocol for development of a reporting guideline (TRIPOD-AI) and risk of bias tool (PROBAST-AI) for diagnostic and prognostic prediction studies based on artificial intelligence**

Gary S. Collins<sup>1,2 \*</sup>, Paula Dhiman<sup>1,2</sup>, Constanza L. Andaur Navarro<sup>3,4</sup>, Jie Ma<sup>1</sup>, Lotty Hooft<sup>3,4</sup>, Johannes B. Reitsma<sup>3</sup>, Patricia Logullo<sup>1</sup>, Andrew L. Beam<sup>5,6</sup>, Lily Peng<sup>7</sup>, Ben Van Calster<sup>8,9,10</sup>, Maarten van Smeden<sup>3</sup>, Richard D. Riley<sup>11</sup>, Karel G.M. Moons<sup>3,4 \*</sup>

<sup>1</sup> Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology & Musculoskeletal Sciences, University of Oxford, Oxford, OX3 7LD, United Kingdom

<sup>2</sup> NIHR Oxford Biomedical Research Centre, John Radcliffe Hospital, Oxford, United Kingdom NIHR Oxford Biomedical Research Centre, John Radcliffe Hospital, Oxford, United Kingdom

<sup>3</sup> Julius Center for Health Sciences & Primary Care, and Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands

<sup>4</sup> Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, The Netherlands

<sup>5</sup> Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, Massachusetts, United States

<sup>6</sup> Department of Biomedical Informatics, Harvard Medical School, Boston, Massachusetts, United States

<sup>7</sup> Google Health, 3400 Hillview Ave, Palo Alto, CA 94304, United States

<sup>8</sup> KU Leuven, Department of Development and Regeneration, Leuven, Belgium

<sup>9</sup> Department of Biomedical Data Sciences, Leiden University Medical Centre (LUMC), Leiden, the Netherlands

<sup>10</sup> EPI-centre, KU Leuven, Leuven, Belgium

<sup>11</sup> Centre for Prognosis Research, School of Medicine, Keele University, Staffordshire, ST5 5BG, United Kingdom.

# Consensus statement

- Delphi about to be launched to get consensus on items to include in the checklist
  - Interested in participating in the Delphi then contact me
- Anticipate TRIPOD-AI to be not too dissimilar to the original TRIPOD
- Biggest difference will be in the terminology, examples, and methods guidance
- Getting agreement on how to present/make the model available

## Participant Information Sheet

### Developing the Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis for Artificial Intelligence (TRIPOD-AI)

#### What is the purpose of the survey?

The TRIPOD (Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis) statements aims to improve the reporting of clinical prediction model studies, however, there are distinct challenges in artificial intelligence studies (such as transparency and description of the model, methodology, availability of data, description and replicability), which are not addressed by the current guidance. Given the rapid expansion of research in this area, there is an urgent need for guidance on the conduct and reporting of clinical prediction model studies involving artificial intelligence and machine learning. We are proposing to develop an extension to the TRIPOD statement to provide guidance for authors of studies involving artificial intelligence and machine learning (TRIPOD-AI).

#### Who is conducting the survey?

This project is led by the TRIPOD-AI group, which supported by a grant from HDR UK, Cancer Research (C49297/A27294), the NIHR Biomedical Research Centre, Oxford, and the Netherlands Organisation for Scientific Research.

The TRIPOD-AI group is led by Professor Gary Collins, (Professor of Medical Statistics, University of Oxford, United Kingdom) and Professor Carl Moons, (Professor Clinical Epidemiology, UMC Utrecht, the Netherlands).

# Harmonisation of two languages

<b>Statistics</b>	<b>Machine learning</b>	<b>Statistics</b>	<b>Machine learning</b>
Covariates	Features	Prediction	Supervised learning
Outcome variable	Target	Latent variable modeling	Unsupervised learning
Model	Network, graphs	Fitting	Learning
Parameters	Weights	Prediction error	Error
Model for discrete var.	Classifier	Sensitivity	Recall
Model for continuous var.	Regression	Positive predictive value	Precision
Log-likelihood	Loss	Contingency table	Confusion matrix
Multinomial regression	Softmax	Measurement error model	Noise-aware ML
Measurement error	Noise	Structural equation model	Gaussian Bayesian network
Subject/observation	Sample/instance	Gold standard	Ground truth
Dummy coding	One-hot encoding	Derivation-validation	Training-test
Measurement invariance	Concept drift	Experiment	A/B test

# *Artificial intelligence faces reproducibility crisis*

Unpublished code and sensitivity to training conditions make many claims hard to verify

The most basic problem is that researchers often don't share their source code. At the AAAI meeting, Odd Erik Gundersen, a computer scientist at the Norwegian University of Science and Technology in Trondheim, reported the results of a survey of 400 algorithms presented in papers at two top AI conferences in the past few years. He found that only 6% of the presenters shared the algorithm's code. Only a third shared the data they tested their algorithms on, and just half shared "pseudocode"—a limited summary of an algorithm. (In many cases, code is also absent from AI papers published in journals, including *Science* and *Nature*.)

Researchers say there are many reasons for the missing details: The code might be a work in progress, owned by a company, or held tightly by a researcher eager to stay ahead of the competition. It might be dependent on other code, itself unpublished. Or it might be that the code is simply lost, on a crashed disk or stolen laptop—what Rougier calls the “my dog ate my program” problem.

# Reporting, Model availability + independent evaluation

e.g.,

- Make it available on a repository (e.g., GitHub)
- Grant access to get predictions for your data set
- Gain access to the code by setting-up non-disclosure agreements

## Minimum information about clinical artificial intelligence modeling: the MI-CLAIM checklist

Here we present the MI-CLAIM checklist, a tool intended to improve transparent reporting of AI algorithms in medicine.

Beau Norgeot, Giorgio Quer, Brett K. Beaulieu-Jones, Ali Torkamani, Raquel Dias, Milena Gianfrancesco, Rima Arnaout, Isaac S. Kohane, Suchi Saria, Eric Topol, Ziad Obermeyer, Bin Yu and Atul J. Butte

The application of artificial intelligence (AI) in medicine is an old idea<sup>1</sup>, but methods for this in the past involved programming computers with patterns or rules ascertained from human experts, which resulted in deterministic, rule-based systems. The study of AI in medicine has grown tremendously in the past few years

due to increasingly available datasets from medical practice, including clinical images, genomics, and electronic health records, as well as the maturity of methods that use data to train computers<sup>2</sup>. The use of data labeled by clinical experts to train machine, probabilistic, and statistical models is called supervised machine learning. Successful

uses of these new machine-learning approaches include targeted real-time early-warning systems for adverse events<sup>3</sup>, the detection of diabetic retinopathy<sup>4</sup>, the classification of pathology and other images, the prediction of the near-term future state of patients with rheumatoid arthritis<sup>5</sup>, patient-discharge disposition<sup>6</sup>, and more.

NATURE MEDICINE | VOL 26 | OCTOBER 2020 | 129-130 | www.nature.com/naturemedicine

### Reproducibility (Part 6): choose appropriate tier of transparency

- Tier 1: complete sharing of the code
- Tier 2: allow a third party to evaluate the code for accuracy/fairness; share the results of this evaluation
- Tier 3: release of a virtual machine (binary) for running the code on new data without sharing its details
- Tier 4: no sharing

### Matters arising

## Transparency and reproducibility in artificial intelligence

<https://doi.org/10.1038/s41586-020-2176-y>

Received: 1 February 2020

Accepted: 10 August 2020

Published online: 14 October 2020

Check for updates

Benjamin Heller-Kalish<sup>1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000</sup>

Armin Haller, S. W. McElreath et al. Nature <https://doi.org/10.1038/s41586-020-2176-y> (2020)

Table 2 | Frameworks to share code, software dependencies and deep-learning models

Resource	URL
<b>Code</b>	
BitBucket	<a href="https://bitbucket.org">https://bitbucket.org</a>
GitHub	<a href="https://github.com">https://github.com</a>
GitLab	<a href="https://about.gitlab.com">https://about.gitlab.com</a>
<b>Software dependencies</b>	
Conda	<a href="https://conda.io">https://conda.io</a>
Code Ocean	<a href="https://codeocean.com">https://codeocean.com</a>
Gigantum	<a href="https://gigantum.com">https://gigantum.com</a>
Colaboratory	<a href="https://colab.research.google.com">https://colab.research.google.com</a>
<b>Deep-learning models</b>	
TensorFlow Hub	<a href="https://www.tensorflow.org/hub">https://www.tensorflow.org/hub</a>
ModelHub	<a href="http://modelhub.ai">http://modelhub.ai</a>
ModelDepot	<a href="https://modeldepot.io">https://modeldepot.io</a>
Model Zoo	<a href="https://modelzoo.co">https://modelzoo.co</a>
<b>Deep-learning frameworks</b>	
TensorFlow	<a href="https://www.tensorflow.org/">https://www.tensorflow.org/</a>
Caffe	<a href="https://caffe.berkeleyvision.org/">https://caffe.berkeleyvision.org/</a>
PyTorch	<a href="https://pytorch.org/">https://pytorch.org/</a>

# (anticipated) TRIPOD-AI Timelines

- 1<sup>st</sup> round of the Delphi survey to be launched shortly
- 2<sup>nd</sup> round of the Delphi survey to be launched in May
- Post Delphi checklist by June/July
- (online) Consensus meeting in the Autumn

\*more information can be found on the OSF (<https://osf.io/zyacb/>)

# Complementary initiatives

- PROBAST-AI (1<sup>st</sup> round of Delphi survey about to be launched)
  - Risk of bias tool for prediction models using AI
- STARD-AI (1st round of Delphi survey closed)
  - Reporting of diagnostic test accuracy studies using AI
- QUADAS-AI (Delphi in preparation)
  - Risk of bias tool for diagnostic test accuracy studies using AI
- DECIDE-AI (1st round of Delphi survey will shortly be closed)
  - Reporting of studies looking at early clinical evaluation, human factors evaluation, preparatory steps towards large-scale clinical trials
- SPIRIT-AI/CONSORT-AI (already published in BMJ/Nature Med,...)
  - Protocols and results of clinical trials of AI interventions

# Opportunities and challenges around the generation of synthetic healthcare data

Dr Allan Tucker  
Intelligent Data Analytics Group,  
Department of Computer Science, Brunel University London.





# NEWS

Home | Coronavirus | US Election | UK | World | Business | Politics | Tech | Science

## Health

# Care.data: How did it go so wrong?



Nick Trigg  
Health correspondent

19 February 2014



digitalhealth  
news + networks + intelligence

Search

News Features Covid-19 Jobs CCIO Network Health CIO Intelligence Events Networks  
AI and Analytics Clinical Software Cyber Security Digital Patient Infrastructure Interoperability

### NEWS

Data sharing | Great North Care Record | Joe McDonald | LHCRES | Matthew Gould | NHS Digital | NHS England | NHSX | Patient data | Sengyne Health | Simon Stevens | System C

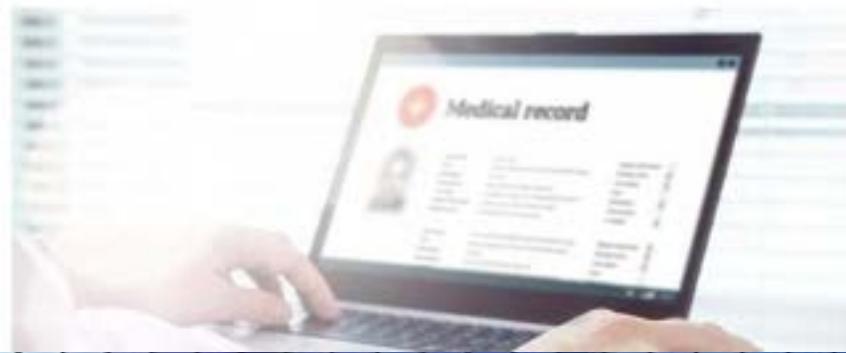
Andrea Downey

19 December 2019

Share this...



## NHS 'risks repeat of care.data in talks to commercialise medical records'



# Background

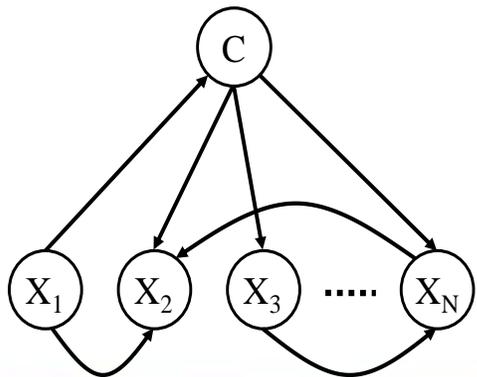
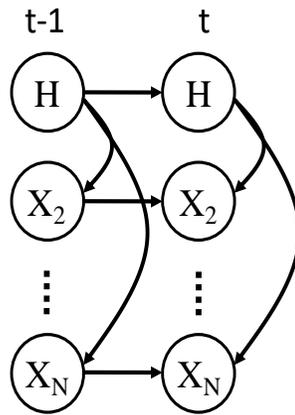
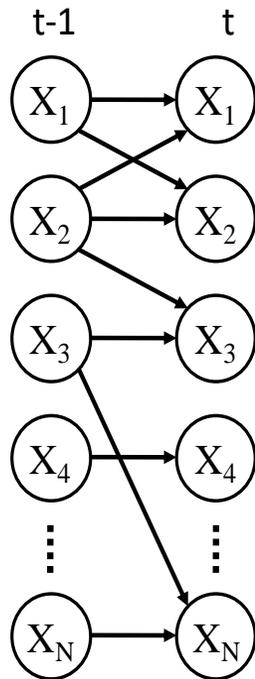
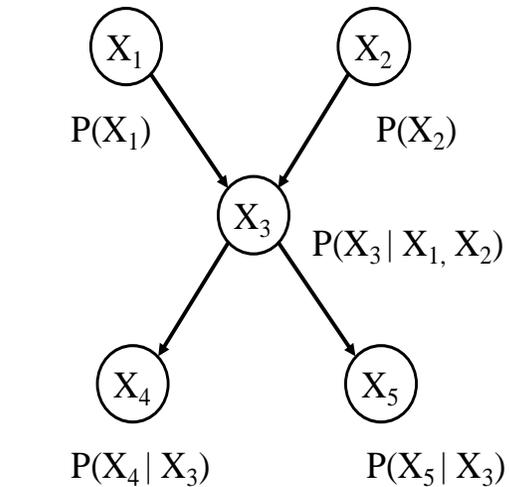
- Adding Noise / Jitter
- K-anonymisation: removing or generalising data
  - but what is an acceptable K?
  - not good on high dimensions
- Differential Privacy: removing an individual has no effect on analysis
  - but repeated requests for aggregate data enabling id of individuals
  - Impact of rare cases (elderly individuals / rare disease / personalised medicine)



# Bias and Transparency

- Impacts of “Black Box” models
- Is the model / data / society biased?
  - Inductive bias of generative model learning algorithms hasn’t attracted enough attention
  - Data analysis is always “Secondary”, collected for different reasons
  - Covid is highlighting the biases we have in diagnosing / treating different groups
- What dependencies are represented in the model?
- Can we confirm known biological relationships?

# Probabilistic Graphical Networks (flexible conditional data generation)



## PrivBayes: Private Data Release via Bayesian Networks

Authors: Jan Zhang, Graham Cormode, Cecilia M. Procopio, Divyesh Srivastava, Xiaokui Xiao

npj | digital medicine

<https://doi.org/10.1145>

Explore our content | Journal information

nature > npj digital medicine > articles > article

Article | Open Access | Published: 09 November 2020

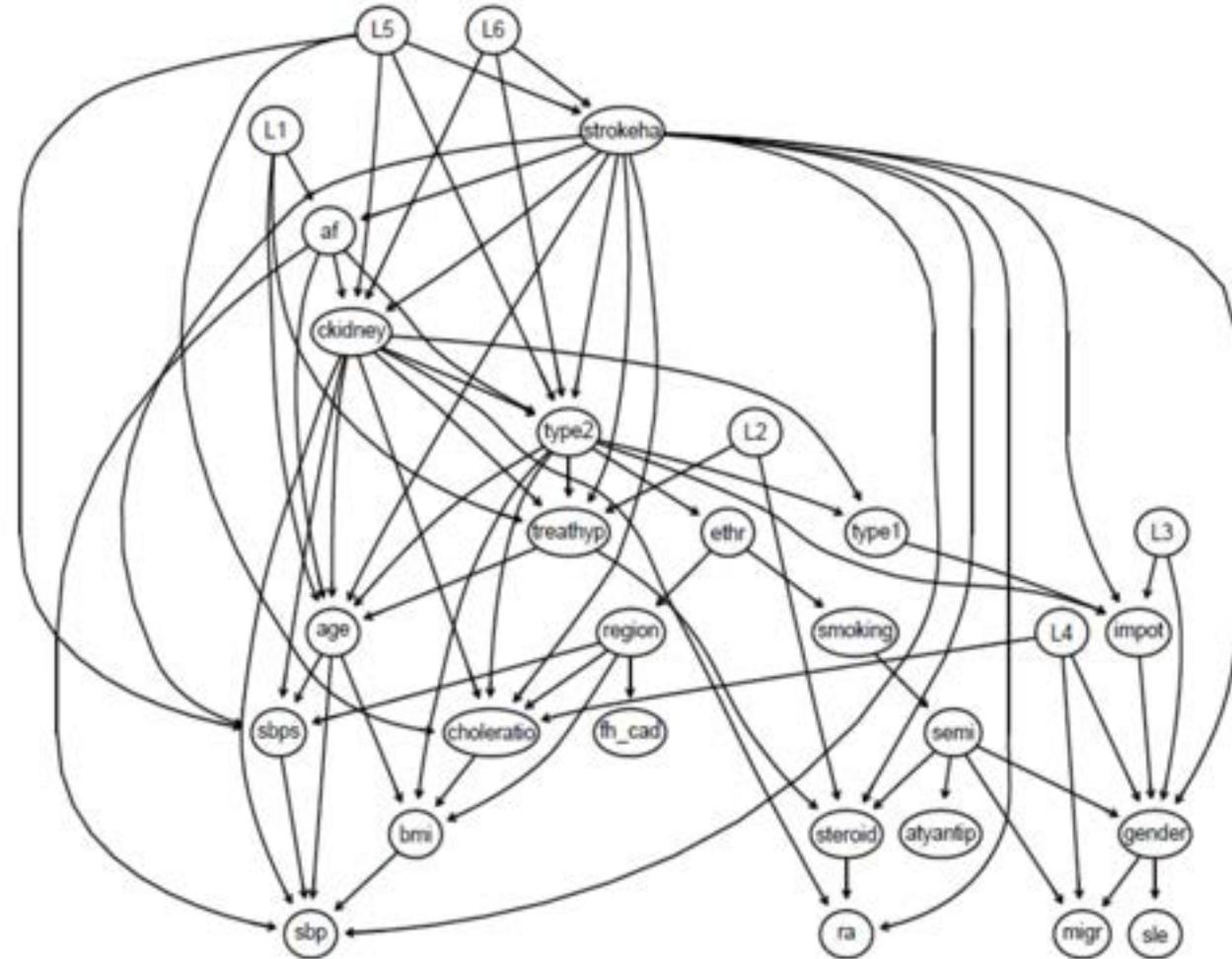
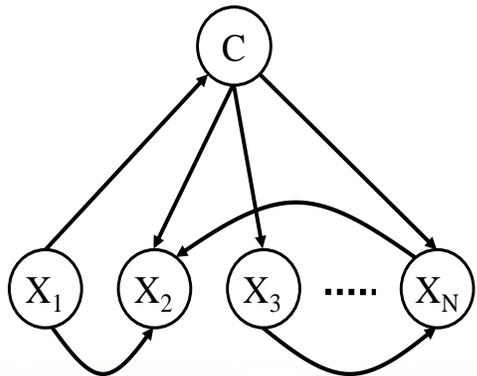
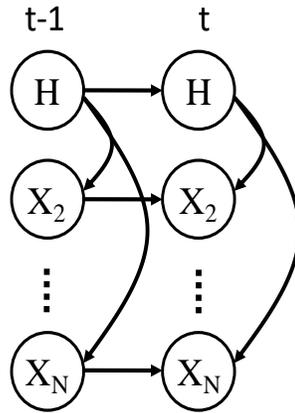
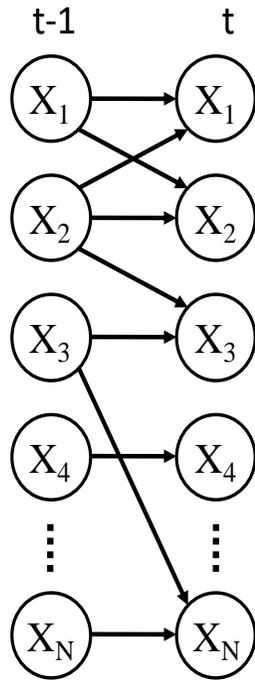
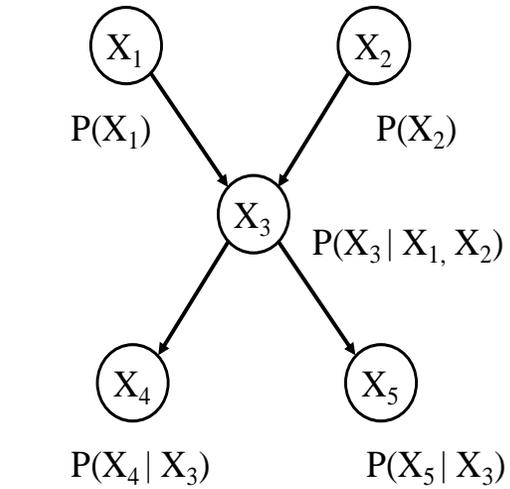
## Generating high-fidelity synthetic patient data for assessing machine learning healthcare software

Allan Tucker, Zhenchen Wang, Ylenia Rotalinti & Puja Myles

npj Digital Medicine 3, Article number: 147 (2020) | Cite this article

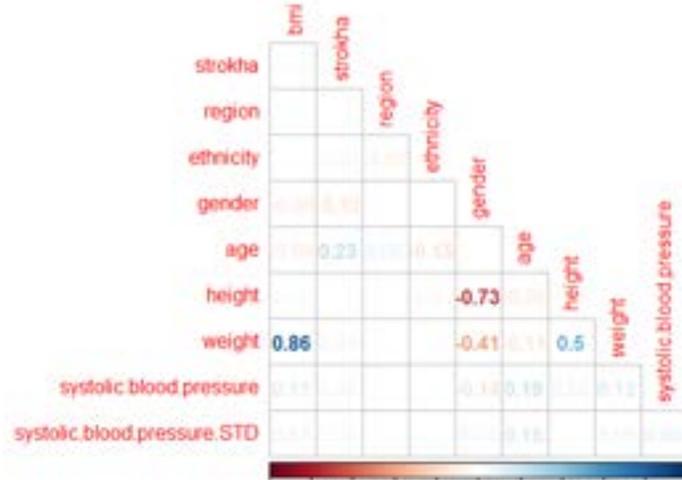
630 Accesses | 5 Altmetric | Metrics

# Probabilistic Graphical Networks (flexible conditional data generation)

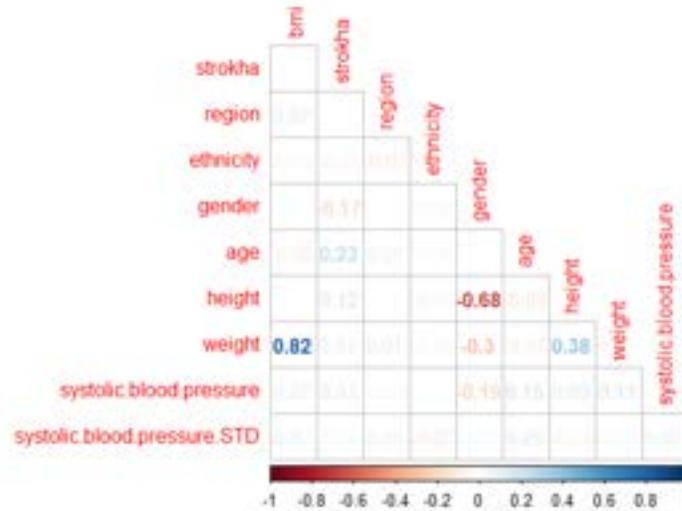


# Fidelity – CPRD

ground truth

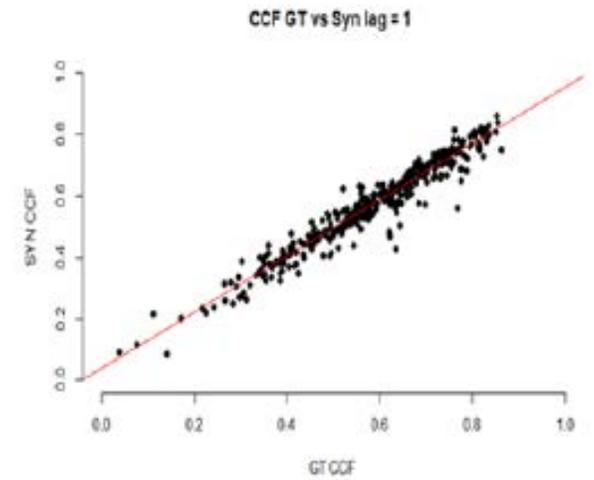
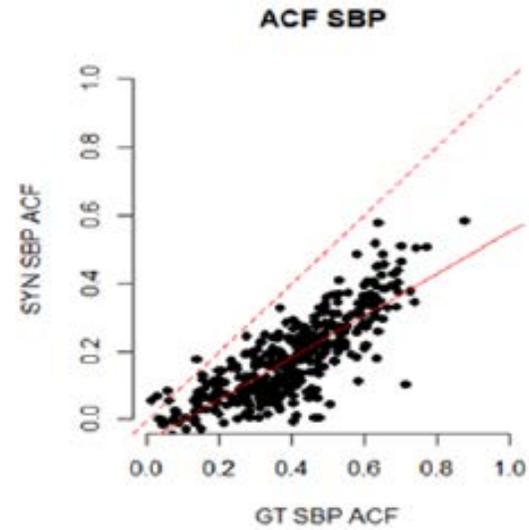
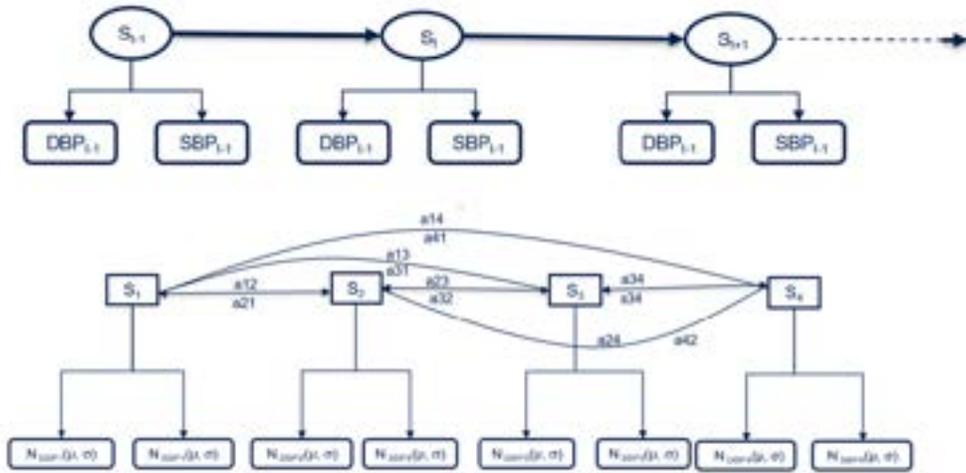


synthetic data

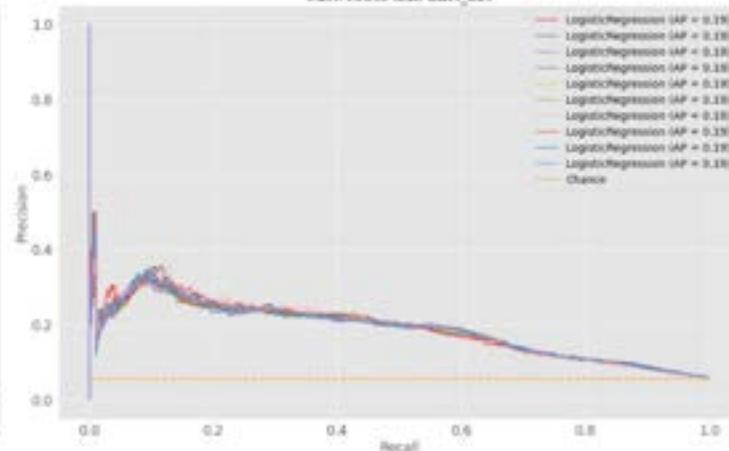
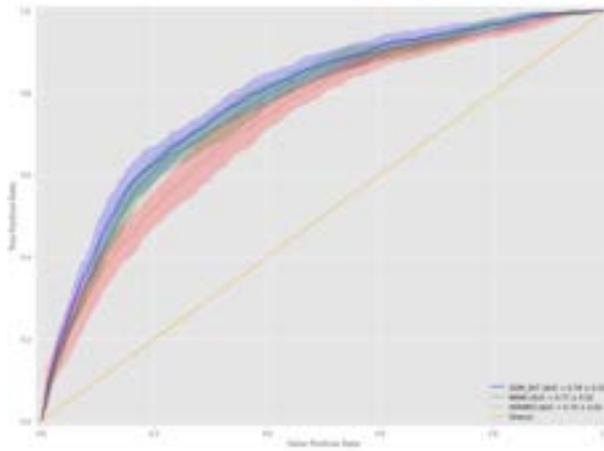
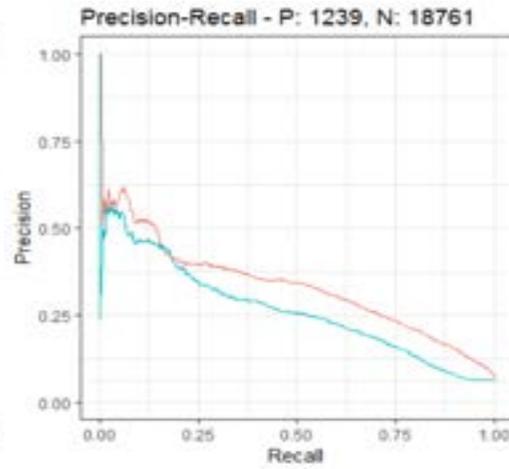
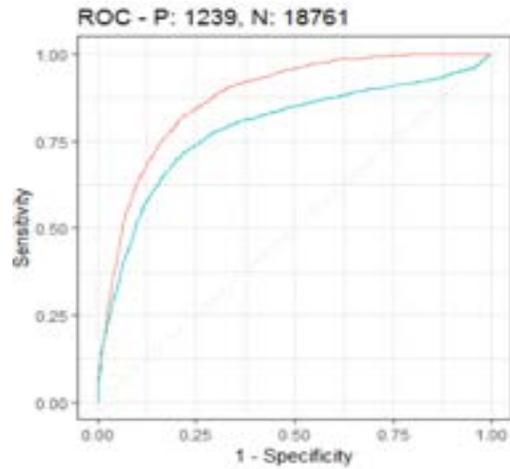


	a) Ground Truth – including missing data		b) Synthetic Data - Modelling GT in a) with missing nodes		c) Synthetic Data - Modelling GT in a) with latent variables	
Factor	Distinct / Missing	Distribution	Distinct / Missing	Distribution	Distinct / Missing	Distribution
age	84 / 0		84 / 0		84 / 0	
strokeha	2 / 0	F:93.6% T:6.4%	2 / 0	F:93.8% T:6.2%	2 / 0	F:93.6% T:6.4%
bmi	580 / 15.73		437 / 7.08		437 / 0.01	
choleratio	573 / 88.59		80 / 41.93		80 / 0.08	
smoking	5 / 0	1: (76.7%) 1: (14.2%) 2: (3.4%) 3: (3.9%) 4: (1.8%)	5 / 0	1: (76.8%) 1: (14.1%) 2: (3.3%) 3: (4.0%) 4: (1.8%)	5 / 0	1: (76.8%) 1: (14.1%) 2: (3.3%) 3: (4.0%) 4: (1.8%)

# Fidelity – MIMIC



# Validation – MIMIC / Sensyne Health / CPRD



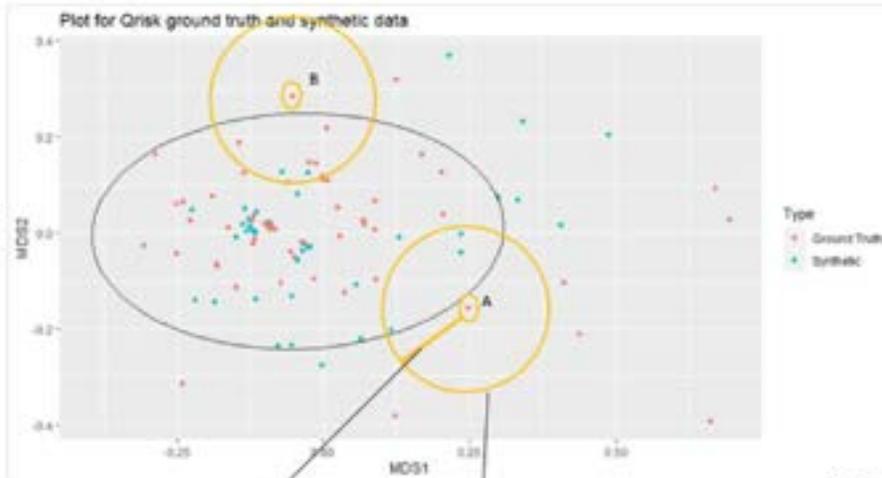
Iteration	ROC & PR curves (P=Positive cases, N = Negative cases)	Granger causality p value ( $\alpha=0.05$ ): PR, ROC	AUC GT: PR, ROC	AUC SYN: PR, ROC
1		<0.001, <0.001	0.293, 0.812	0.346, 0.881
2		<0.001, <0.001	0.292, 0.818	0.341, 0.881
3		<0.001, <0.001	0.274, 0.797	0.349, 0.880

# Privacy Tests

- Explore effect of sample size and number of variables on risk of “clones” / “inliers” / “outliers”.
- Results are based on 10 iterations of resampling without replacement.
- Risk of inliers / outliers (single real datapoints close to a synthetic datapoint) always very small.
- Actual number of outliers generated stays relatively stable (between 10 and 70).

GT population size	$R_{clone}$	$R_{in}, Pr = 0.001$	$R_{out}, Pt = 0.999$
100,000	0.016	462 (0.4620%)	25 (0.0250%)
200,000	0.013	770 (0.3850%)	34 (0.0170%)
300,000	0.014	613 (0.2043%)	24 (0.0080%)
400,000	0.012	553 (0.1383%)	53 (0.0133%)
500,000	0.016	529 (0.1058%)	19 (0.0038%)
600,000	0.009	254 (0.0423%)	45 (0.0075%)
700,000	0.008	534 (0.0763%)	24 (0.0034%)
800,000	0.011	581 (0.0726%)	13 (0.0016%)
900,000	0.012	518 (0.0576%)	33 (0.0037%)
1,000,000	0.012	30 (0.0030%)	45 (0.0045%)
2,000,000	0.01	78 (0.0039%)	29 (0.0015%)

# Simulating attack - Matching to Real Patients

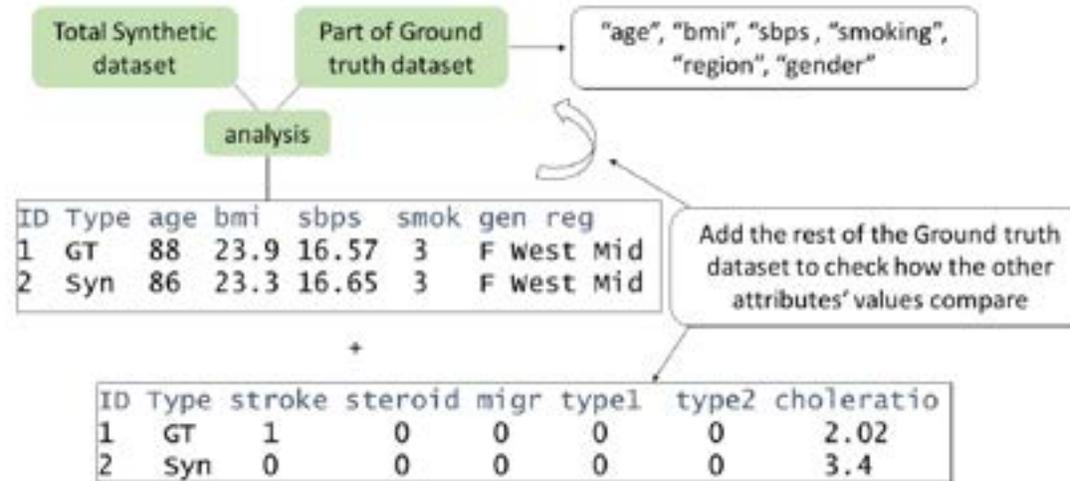


Distance threshold

3 Synthetic NN are close to the GT patient A and 2 Synthetic patients are close to GT patient B within the distance radius.

GT Attributes	Num Of GT Outliers	Num Of GT Patients in 10000 Sample	Num of Single Synth NN	Proportion of GT Outliers with Single Synth NN
age, smoking, region, gender, ethnicity, ckidney	396.3	6248	113.8	1.82%
age, smoking, region, gender, ethnicity, bmi, sbps, ckidney	377	6250	77.7	1.24%
age, smoking, region, gender, ethnicity, bmi, choleraio, sbp, sbps, ckidney	363.1	6289	48.22	0.76%
age, smoking, region, gender, ethnicity, bmi, sbps, ckidney, sle, atyantip, type1, steroid	363.4	6180	43.5	0.70%

Example output:



Press release

# New synthetic datasets to assist COVID-19 and cardiovascular research

Ground-breaking innovation to support medical technologies

Published 29 July 2020

From: [Medicines and Healthcare products Regulatory Agency](#)



### Related content

[Drug Safety Update: monthly PDF newsletter](#)

[Safe use of emollient skin creams to treat dry skin conditions](#)



Computational Intelligence  
AN INTERNATIONAL JOURNAL



# Data Privacy vs Utility

- Synthetic data appears useful in validating models
- Classes of generative models offer promise
  - Large scale datasets take up resources
  - Flexible in handling data types (continuous, categorical, time etc.)
- But always a trade-off between utility and privacy:
  - Aggregated data vs personalised interventions
  - Rare disease
  - Marginalising over many features
- Dealing with bias is important ...

# Thanks

- Barbara Draghi
- Namir Oues
- Juan de Benedetti
- Zhenchen Wang
- Ylenia Rotalinti
- Puja Myles



IDA  
Research



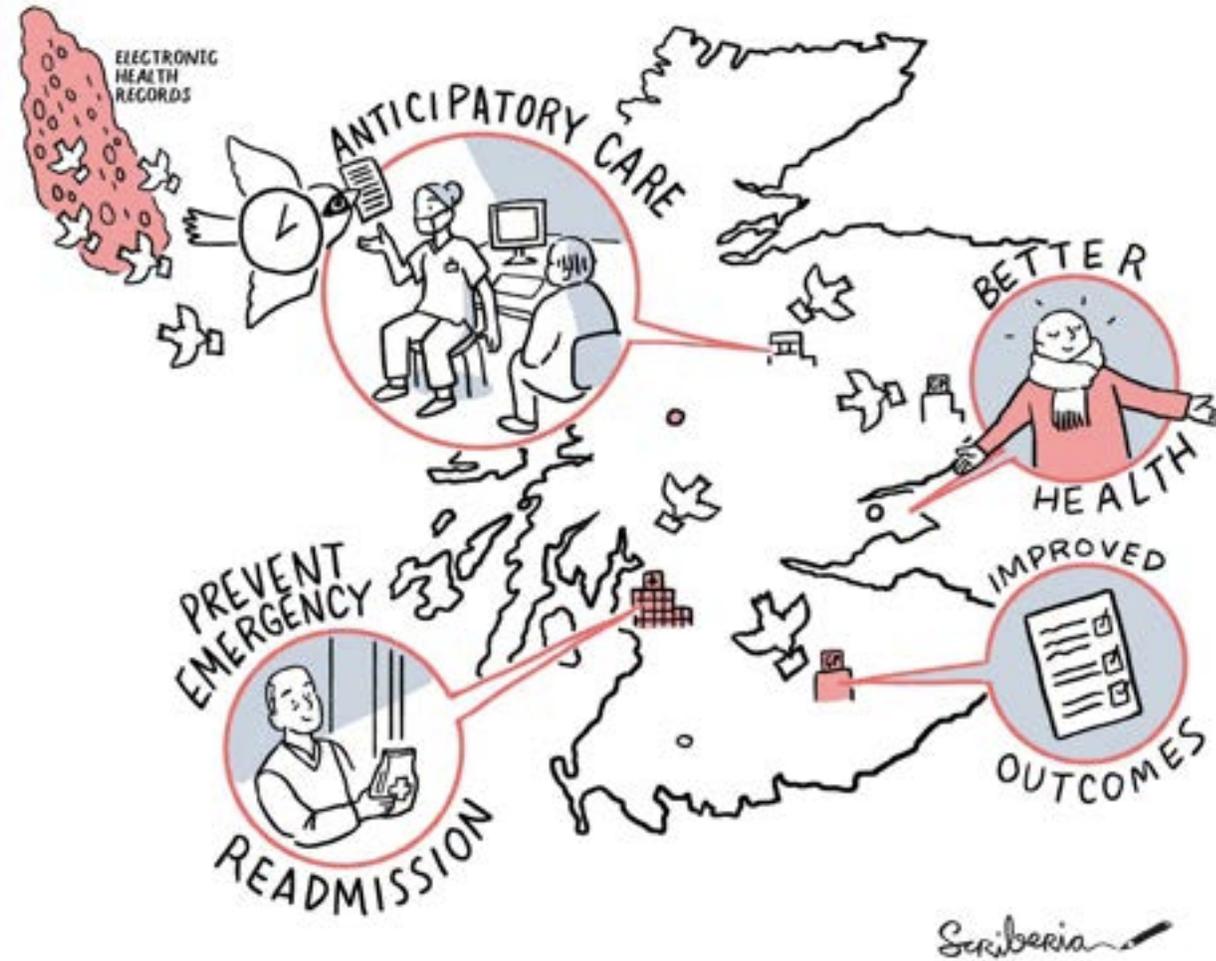
Medicines &  
Healthcare products  
Regulatory Agency



Brunel  
University  
London

# Reproducibility in population-wide emergency risk prediction

Scottish Patients At Risk of Re-admission and Admission  
'SPARRA', version 4



James Liley  
31/03/2021

# Overview

## 1. The SPARRA model

- a). Motivation and aim
- b). Model overview
- c). Findings

## 2. Data safe havens and challenges

- a) 'Hard' obstacles to reproducibility
- b) 'Soft' obstacles to reproducibility

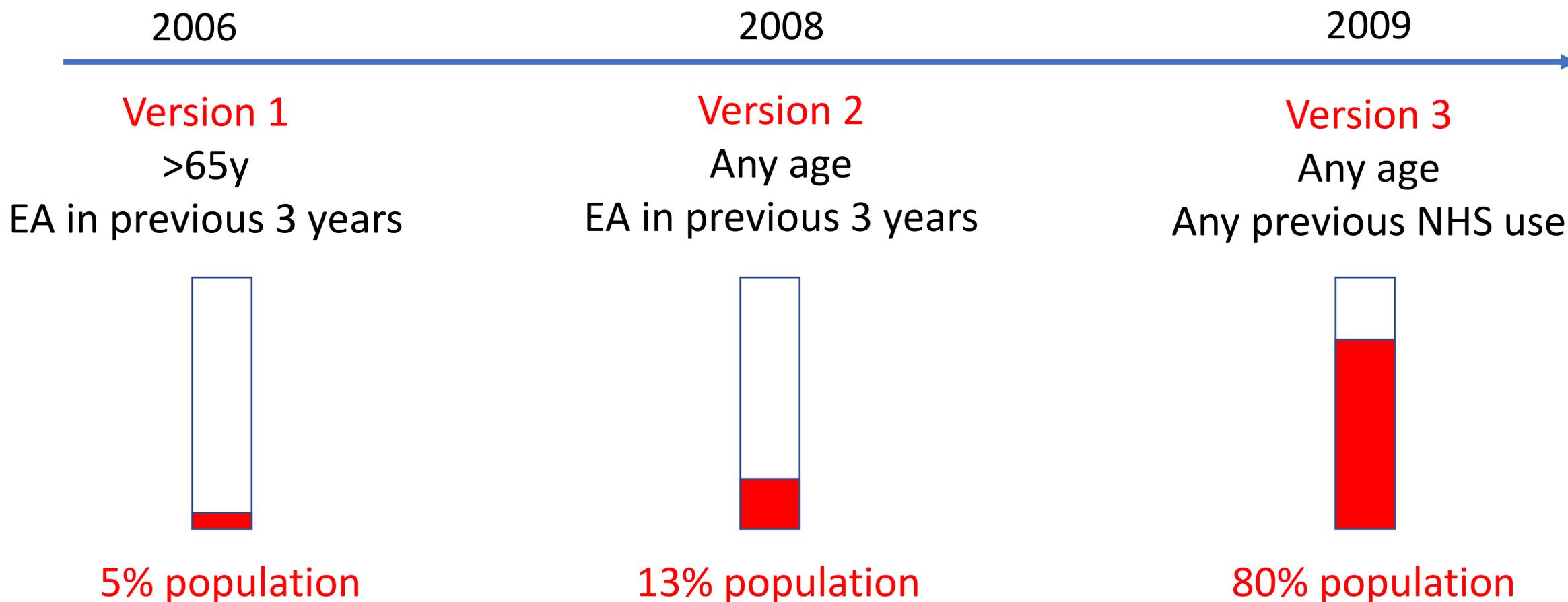
## 3. Practices to improve reproducibility



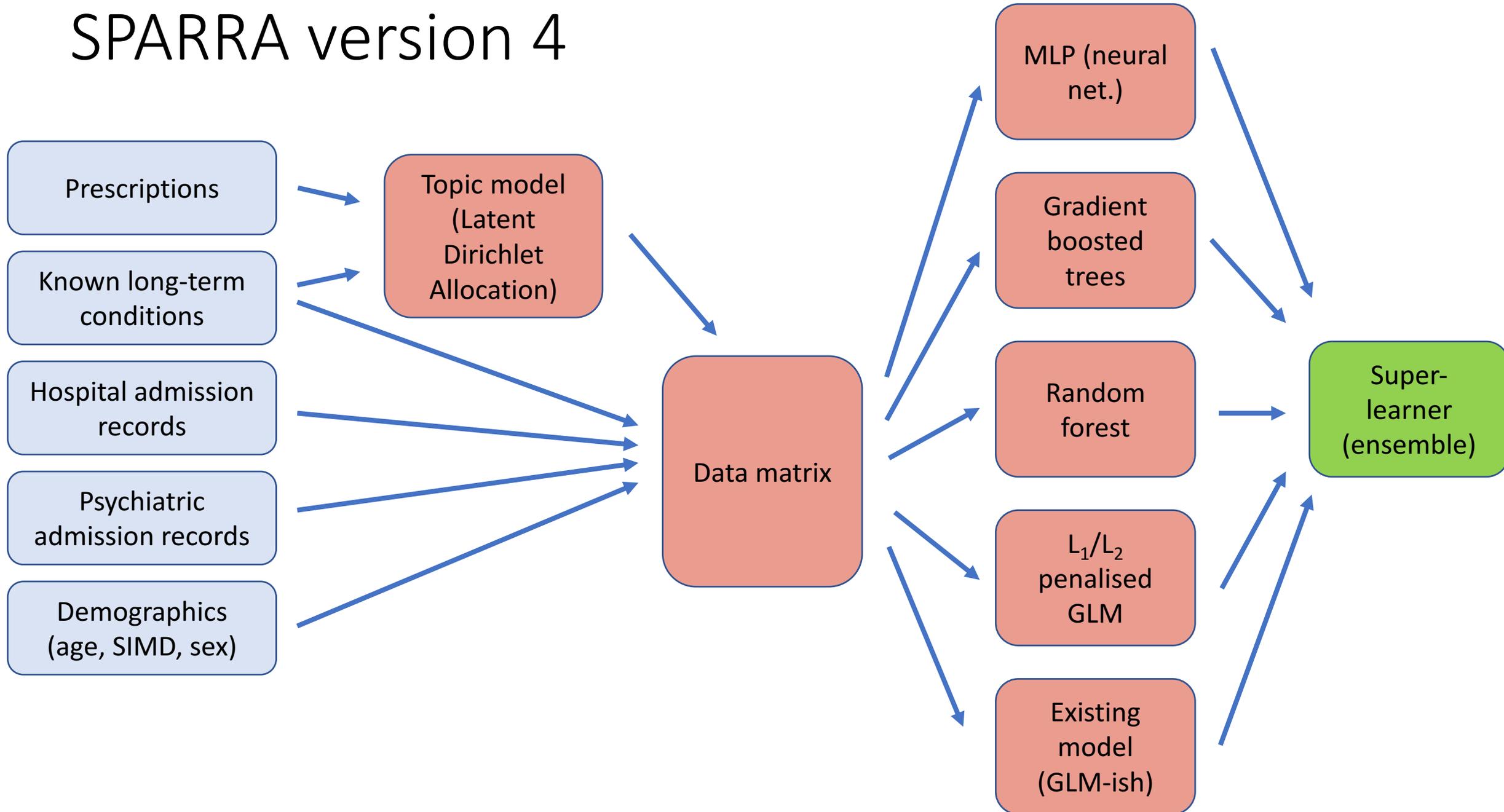
'The NHS should work with other public services and with patients and carers to provide continuous, **anticipatory** care to ensure that, as far as possible, **health care crises are prevented from happening**'

A National Framework for Service Change in the NHS in Scotland,  
NHS Scotland, 2005

SPARRA: estimate probability that a patient will have an emergency hospital admission in the coming year on the basis of records of previous interactions with the healthcare system.



# SPARRA version 4



# Updating considerations

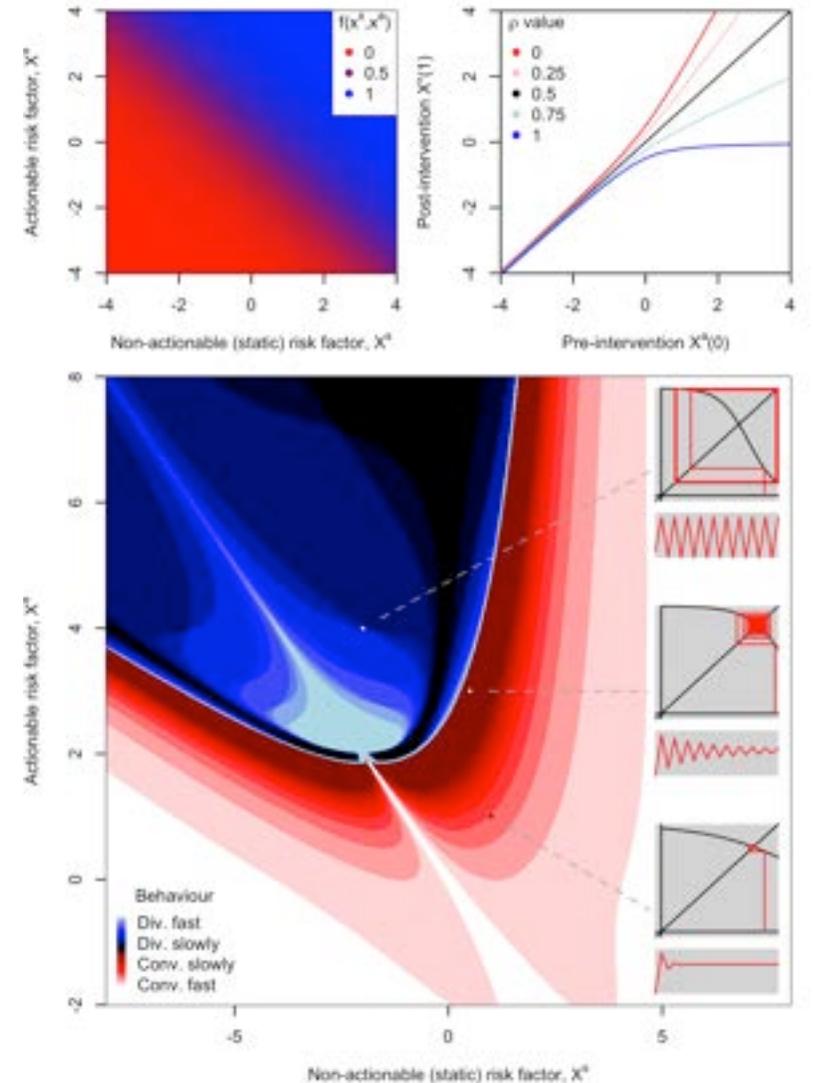
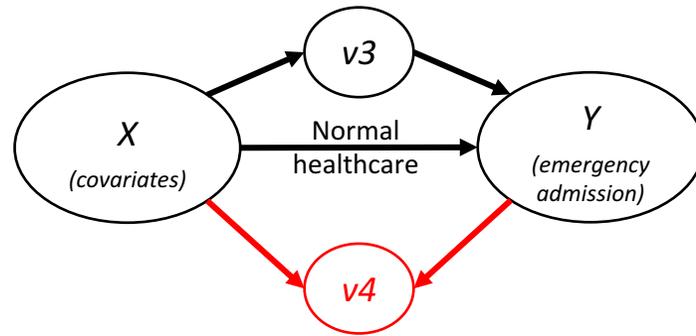
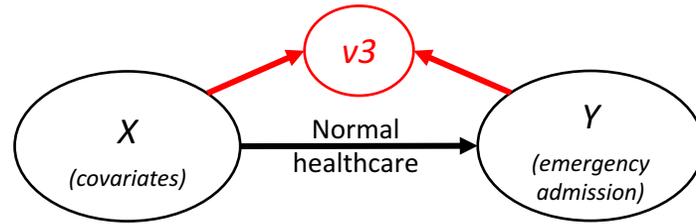
Does an existing model overestimate risk because the model is out-of-date, or because it is driving intervention?

Replacing an existing score means the new score is used in a different system to the one in which it was fitted.

'Naive' updating leads to problems

Proposed solution: act on  $v_3$ , then act on  $v_4$

Approximate by deploying maximum of  $v_3$  and  $v_4$  scores

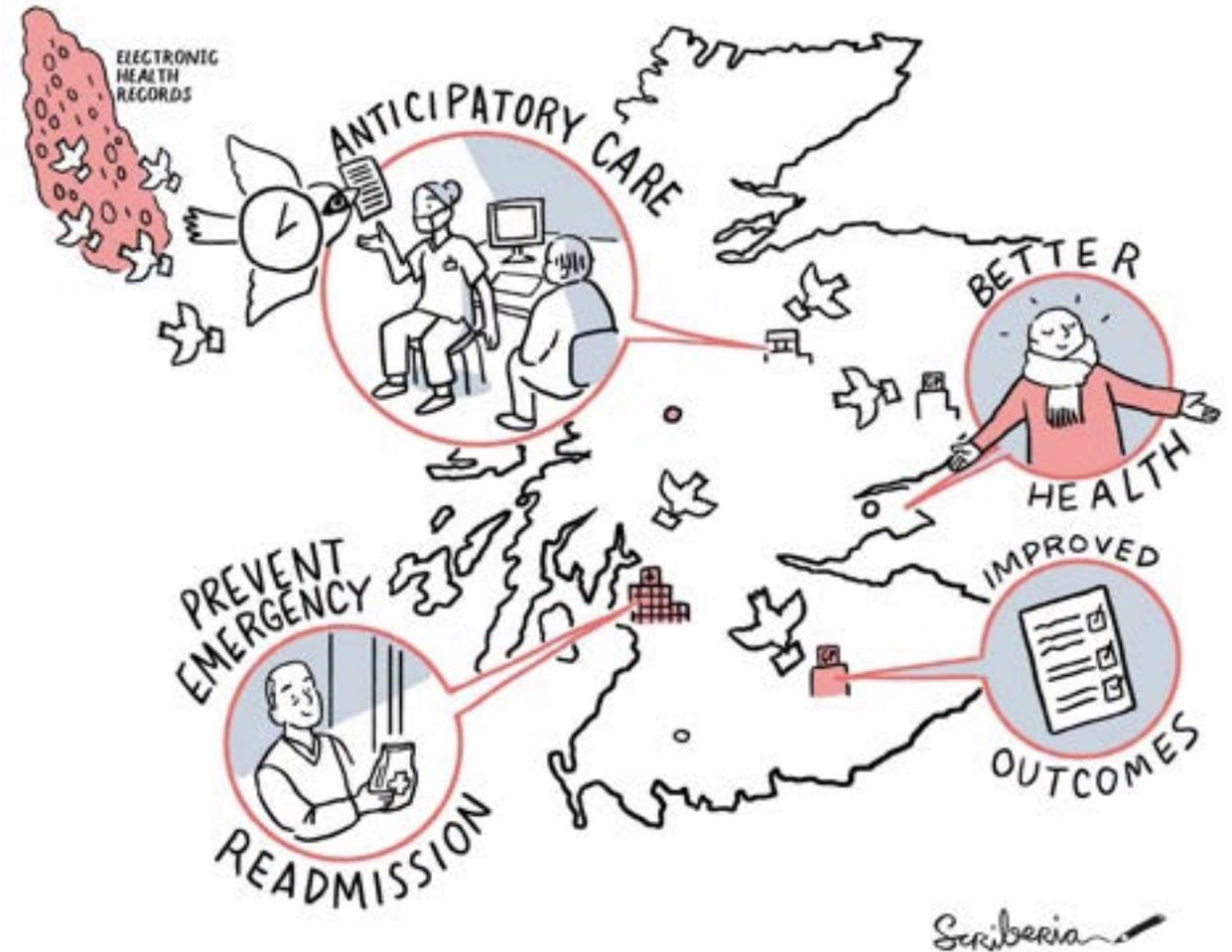


# Deployment

SPARRA scores are deployed to GPs Scotland-wide, and may be used to guide intervention or public health action.

We need to make sure that the score going to GPs is the **same** as the score we are writing about in the paper.

This project was paid for by public funds and influences Scotland's healthcare strategies. It **should be as open as possible**.



# Difficulties with reproducibility

## Data safe havens (DSH) – ‘Five Safes’

1. Restriction the **personnel** accessing the DSH
2. Clear **project** delineations
3. No communications to/from DSH
4. Ensure individuals in data are **not identifiable**.
5. **Disclosure control** on data/code exports

## Other considerations

1. Population-wide scale
2. Complex data to results pipeline
3. Complex legal protections

## Implications (DSH)

### **Raw data cannot be removed from DSH (4)**

Other researchers cannot access DSH to perform checks (1)

No internet access on DSH (3)

Highly dependent on inflexible machine architecture

All code must be readable if it is to be exported (3,5)

**Work progresses slowly (1-5)**

**Data export is slow (3,5)**

## Implications (other)

Generally not feasible to re-run entire pipeline whenever an error is discovered (1)

Many people to keep ‘in the loop’ (2)

**SPARRA v3 algorithm/coefficients embargoed (3)**

**Difficult to attain permission to publish code (3)**

# Hard obstacle 1: full reproducibility is impossible

		Data	
		Same	Different
Analysis	Same	Reproducible	Replicable
	Different	Robust	Generalisable

**Gold standard reproducibility: verify that algorithm gives the same results given the same input**

**Gold standard impractical as**

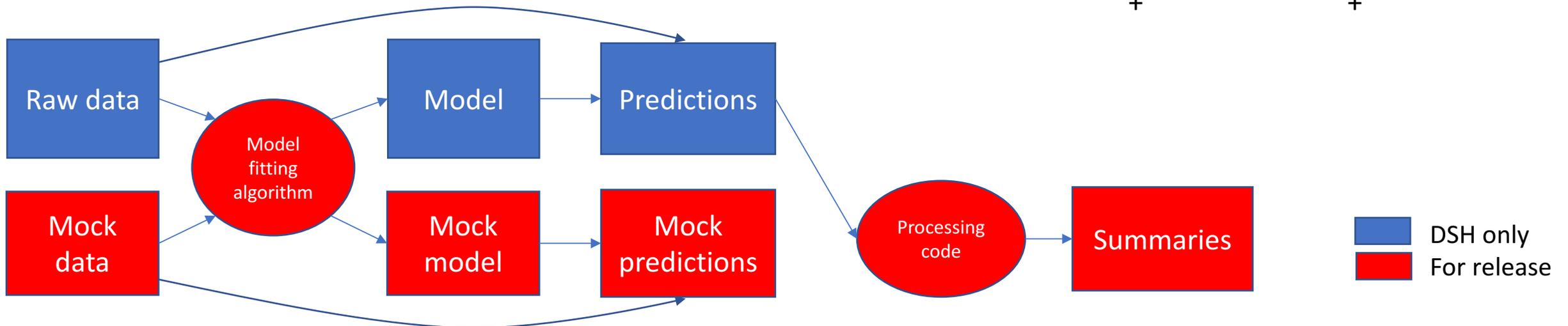
- raw data cannot be exported
- model parameters cannot (generally) be exported
- verifying researchers cannot access DSH.

Technically can be reproduced on other safe havens;  
this is not usually practical.

Most important: **works in deployment machine**

**Management**

- Mockup dataset with similar global properties to real data
- Export of full model fitting algorithm, verified on mock data
- Export of all processing code
- Complete unit test for algorithm
- Extensive export of summaries: figure + data + code



# Hard obstacle 2: computation

## Software considerations

1. No internet access on DSH; no github/stackExchange
2. Software updates difficult/slow.

## Hardware considerations

1. Limited computational capacity
2. Rigid machine structure
3. No control over machine

## Management (software)

1. Code often must be written from scratch. Protocol:
  - all code commented and written in a consistent style
  - all code independently checked and ran by another team member
2. Containerised deployment (Docker, MS Azure)

## Management (hardware)

1. Careful algorithmic design; garbage collection, deleting data after use
2. Develop code to run on one machine or two
3. Multiple save checkpoints.

# Soft obstacle 1: increased personnel

## Personnel difficulties

1. Different teams of people generating and curating raw data (Public Health Scotland), analysing data and fitting models (us), and deploying the model (PHS)
2. No control over computation: crashes, data removal, etc. Pipelines cannot usually be ran end-to-end.

## Management

1. Close liaison with PHS
  - Weekly/biweekly meetings
  - 'Knowledge transfer'
  - Multidisciplinary teams
2. We cannot fix crashes/memory issues ourselves; so
  - pre-empt and avoid issues
  - pipelines split into independent sub-pipelines.

# Soft obstacle 2: everything is more difficult!

## General difficulties

1. Getting data on and off the DSH is very slow
2. Fixing any problem with computation is slow
3. Sharing preliminary data with colleagues is difficult
4. Checking raw data is difficult
5. Beholden to internet connectivity
6. Beholden to general workload of Public Health Scotland
7. Cannot copy-paste into data safe haven
8. No '#' key on data safe haven
9. Need to test new code extensively
10. No root access on data safe haven
11. Cannot delete temporary R files, so accidentally crashed sessions irreversibly use up memory

## Management

**Put reproducibility concerns at the forefront of workflow**

# Summary

Data safe haven work at population scale leads to difficulties in ensuring reproducibility, some of which are unavoidable.

Mock-up data is useful for verification of algorithms

For real-world applications, many people of different backgrounds need to be able to reproduce the results and use the methods. Ongoing 'knowledge transfer' and communication are important.

Working in data safe havens is slower and more difficult than ML researchers may be accustomed to. Additional time has to be allowed for completion of projects so that reproducibility can be maintained.

# Acknowledgments

University of Edinburgh/Alan Turing Institute (ATI)

**Catalina Vallejos**

Vallejos and Sanguinetti groups

Turing Institute

**Gergo Bohner**

Nathan Cunningham

Franz Kiraly

Ioanna Manolopoulou

**Bilal Mateen**



Katrina Payne

**Sebastian Vollmer**

University of Durham/Turing Institute

**Louis Aslett**

**Sam Emerson**

Public Health Scotland

Katie Borland

**David Carr**

Scott Heald

Sam Oduro

**Jill Ireland**

Keith Moffat

**Rachel Porteous**

Stephen Riddell

**Beth Bruce (eDRIS)**

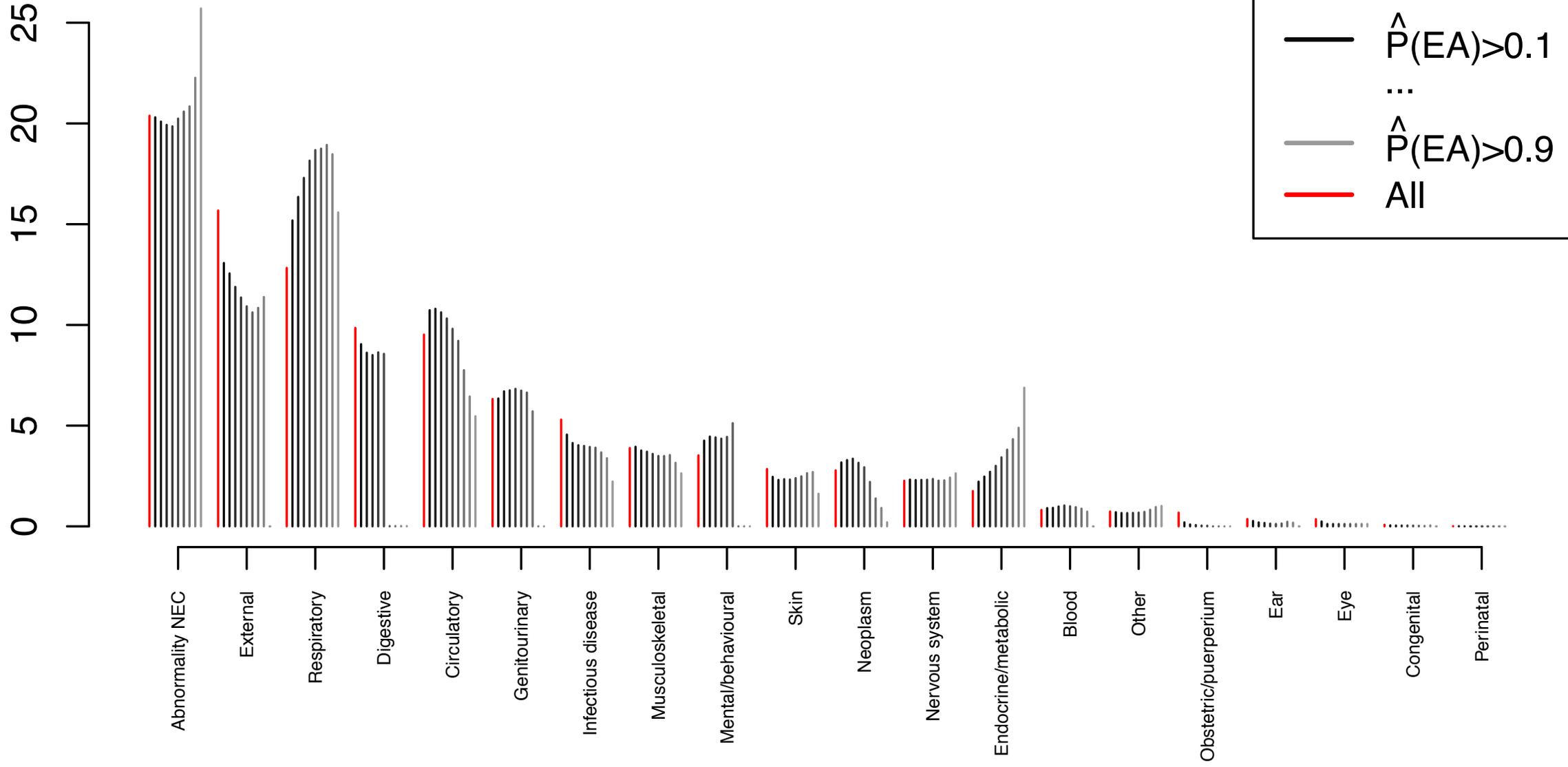
**Lizzie Nicholson (eDRIS)**



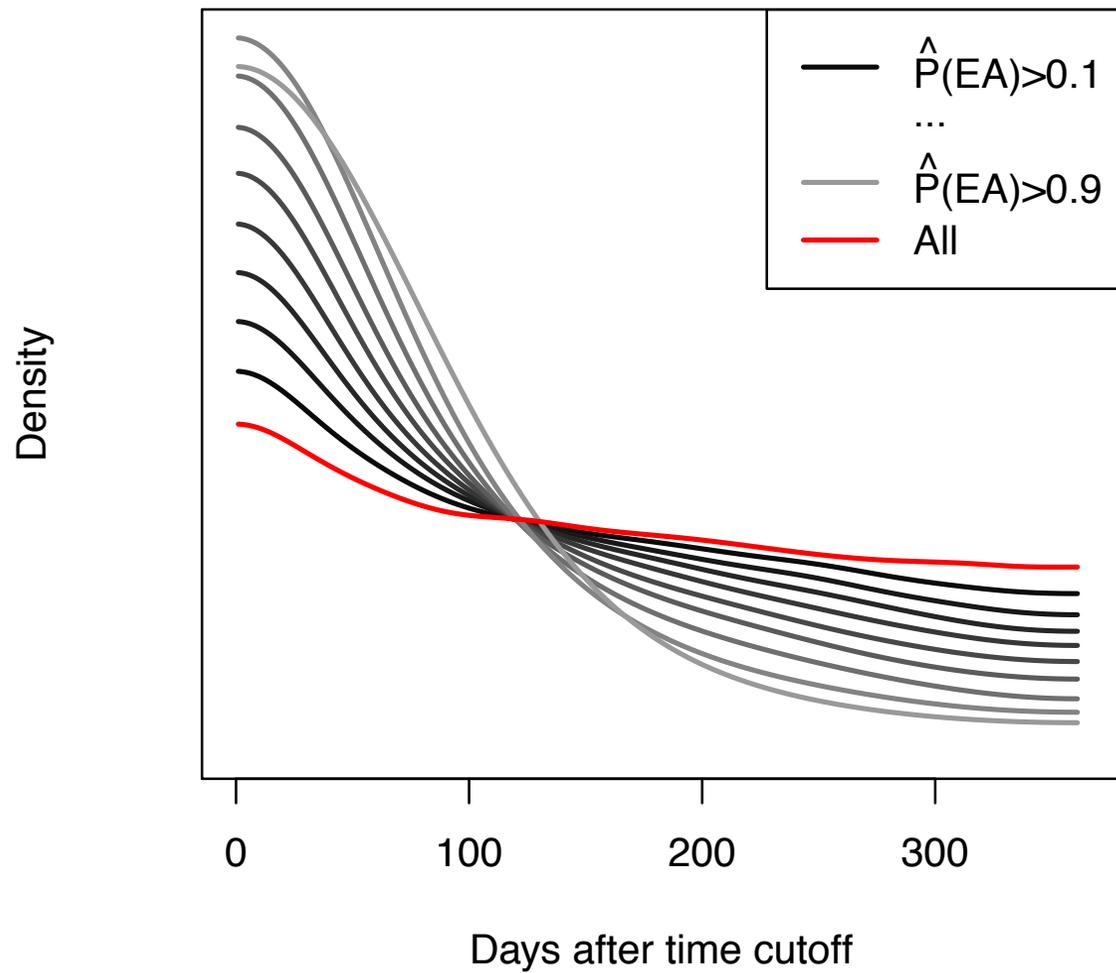
**The  
Alan Turing  
Institute**



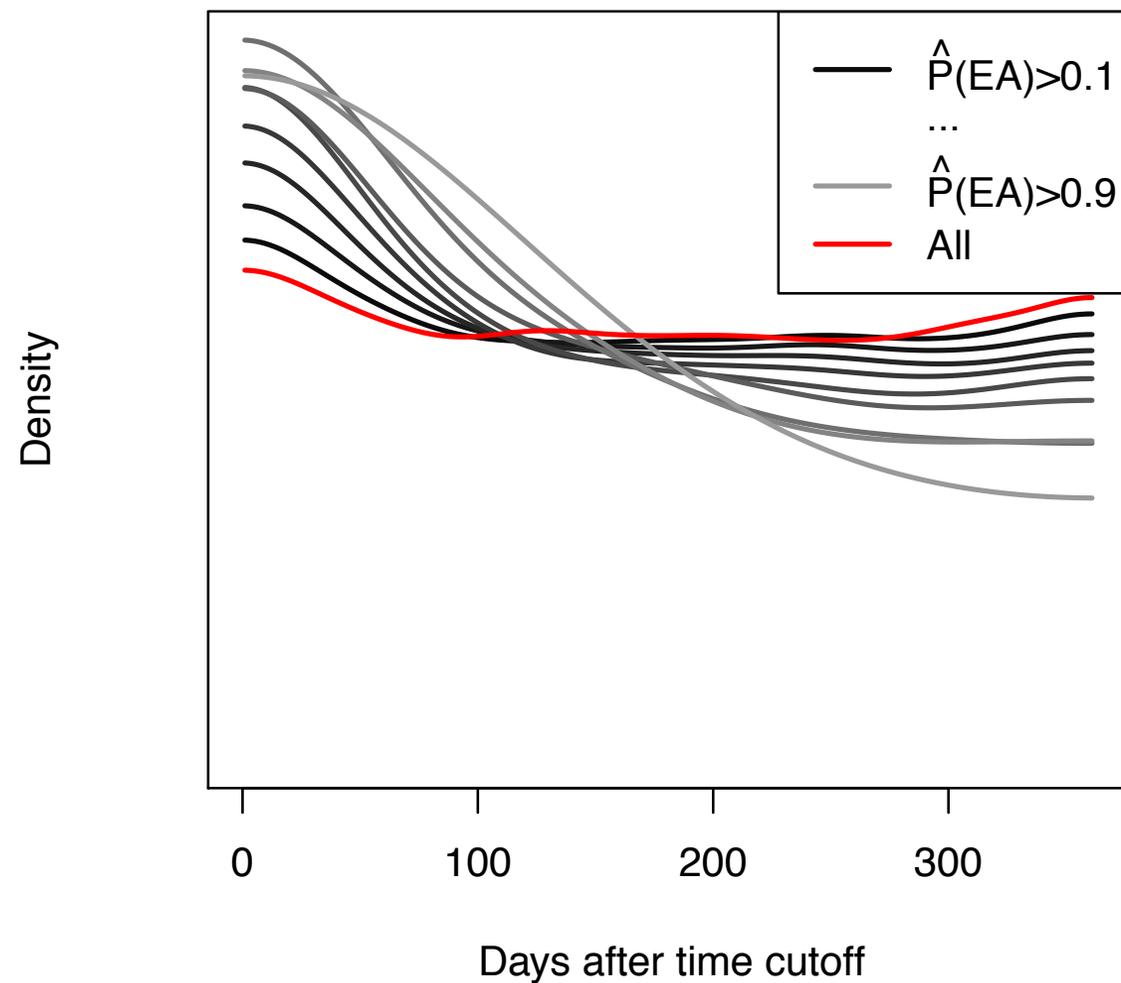
% patients



Density of time-to-first-admission



Density of time-to-only-admission



Tools and open datasets to support training  
activities around reproducible machine learning:  
An activity monitoring use-case  
Work Package II

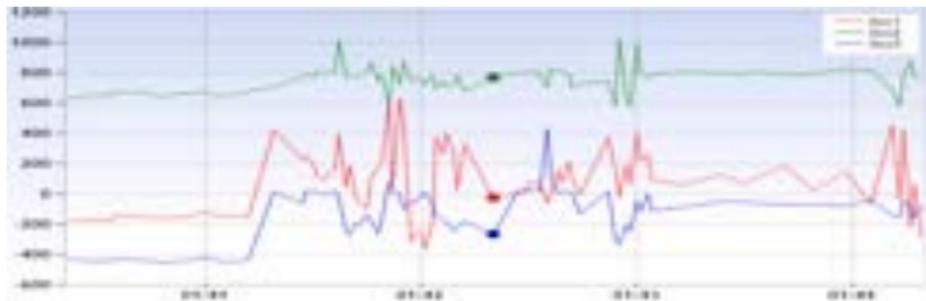
Shing Chan and Aiden Doherty

Big Data Institute & Nuffield Department of Population Health  
University of Oxford

March 31, 2021

# Accelerometers

Measures movement (more specifically, acceleration)



Applications

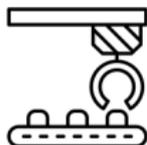


Image Credits: Made, Made x Made Icons, The Noun Project

# Wearable technology



Image Credits: Macrovector, Shutterstock; Business Insider; Xiaomi

# New opportunities for health care

Activity information remains under-exploited in health research

## Measuring activity is challenging

- ▶ How to measure



- ▶ What to measure

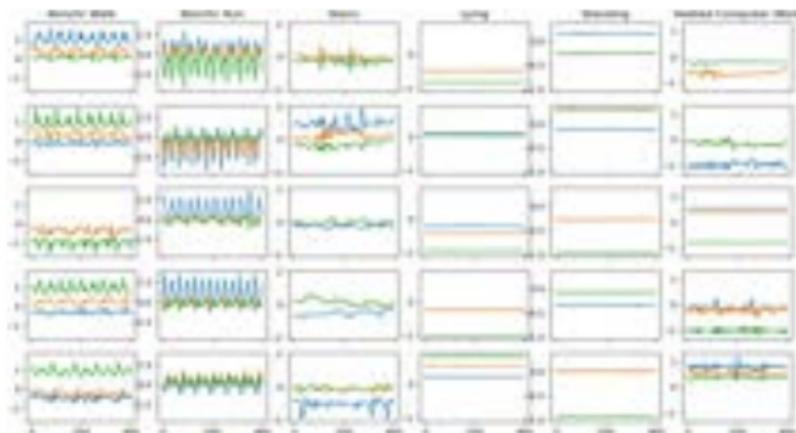


## New directions: Wearables

- ▶ Objective, hi-res, free-living
- ▶ Efficient/scalable (e.g. UKBiobank)



# Activity recognition



Large and representative (i.e. free-living) data is key

Image Credits: Made x Made, ProSymbols, The Noun Project

# Existing openly available wrist-worn accelerometer datasets

Most are small and collected in artificial settings

Dataset	Size	Annotations	Free-living
PAMAP2	9 ppl $\times$ 1 hrs	18 activities (walk, run, sit, ...)	No
MHEALTH	10 ppl $\times$ 15 mins	12 activities (walk, run, sit, ...)	No
ADL	16 ppl $\times$ 1 hrs	14 activities (brush teeth, comb hair, drink glass, ...)	No
UC Berkley WARD	20 ppl $\times$ 1 hrs	13 activities (walk, sit, stand, jog, stairs, ...)	No
CMU-MMAC	43 ppl $\times$ 5 mins	kitchen/cooking activities (pizza, salad, brownie, ...)	No
Opportunity	12 ppl $\times$ 1 hrs	morning activities (prepare breakfast, cleanup, ...)	Semi
<b>Capture-24</b>	<b>158 ppl <math>\times</math> 24 hrs</b>	<b>&gt;200 free-living annotations</b>	<b>Yes</b>

## Capture-24: An activity monitoring dataset

- ▶ Wrist-worn accelerometer (Axivity3, triaxial, 100Hz)
- ▶ 158 subjects  $\times$  24 hours
- ▶ Oxford area, 2015
- ▶ Free-living
- ▶ More than 200 different annotations
- ▶ Compatible with the UKBiobank

## Capture-24: Collection



Accelerometer watch



Body camera



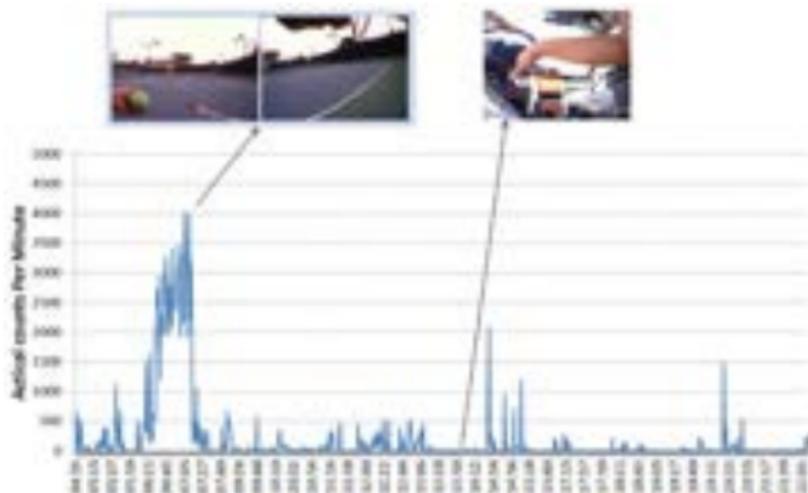
Sleep diary



Accelerometer + Body camera

## Capture-24: Annotation

Accelerometer timestamps marked using camera images



## Capture-24: Annotations

### More than 200 different annotations

1. home activity;miscellaneous;sitting;7010 lying and watching television with TV on as the primary activity
2. leisure;recreation;outdoor;5175 walking/running playing with child(ren)
3. home activity;miscellaneous;sitting;7021 sitting without observable activities
4. leisure;miscellaneous;walking;17133 walking upstairs
5. leisure;miscellaneous;walking;17070 descending stairs
6. home activity;miscellaneous;sitting;9060 sitting reading or using a mobile phone/smartphone/tablet or talking on the phone/computer (skype chatting)
7. home activity;miscellaneous;standing;5146 standing packing/unpacking household items occasional lifting
8. leisure;miscellaneous;21070 (generic) walking/standing combination indoor
9. leisure;miscellaneous;21016 sitting child care only active periods
10. leisure;miscellaneous;21010 sitting non-desk work (with or without eating at the same time)
11. leisure;miscellaneous;17031 loading /unloading a car implied walking

...

## Capture-24: Data Cleaning and Anonymisation



- ▶ Camera's low temporal resolution
- ▶ Sleep diary checked against accelerometry



- ▶ Rare annotations omitted or modified
- ▶ Study time perturbed (random offset), dates omitted
- ▶ Ages converted into age bands

Image credits: QualityIcons, The Noun Project

# Capture-24: Successful use-cases



Model trained on  
Capture-24

**biobank**<sup>uk</sup>



>100,000 subjects  
7 days wear

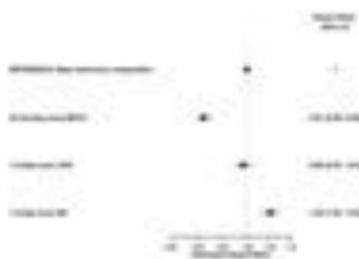
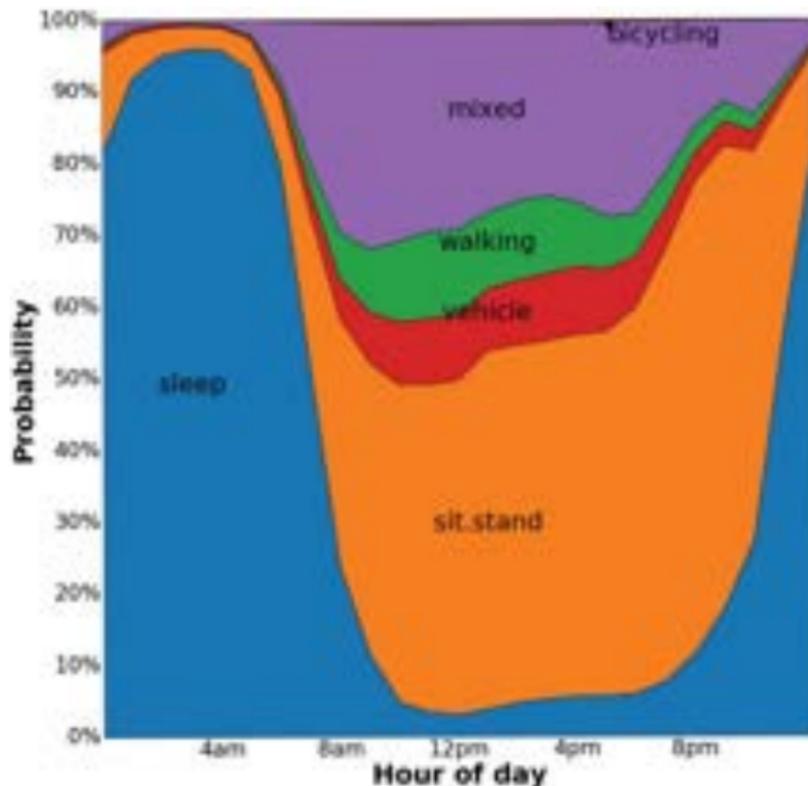


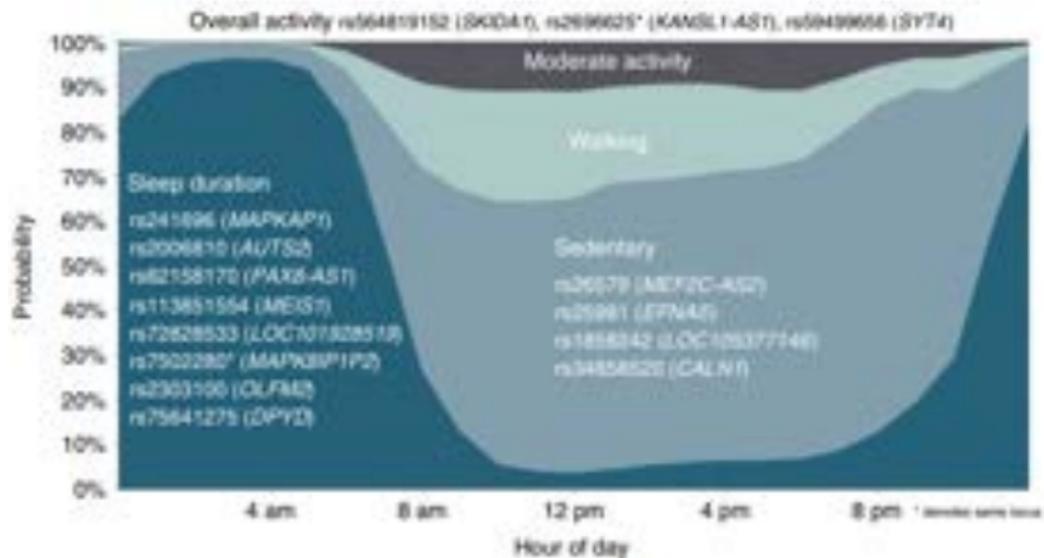
Image credits: Made x Made, ProSymbols, QualityIcons, The Noun Project

## Capture-24: Successful use-cases



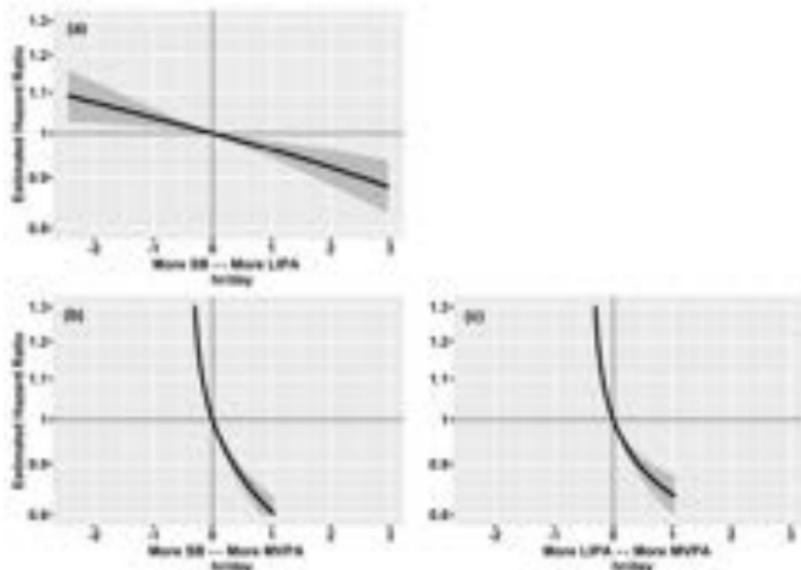
Statistical machine learning of sleep and physical activity phenotypes from sensor data in 96,220 UK Biobank participants – Willetts et al. 2018

## Capture-24: Successful use-cases



GWAS identifies 14 loci for device-measured physical activity and sleep duration – Doherty et al. 2018

## Capture-24: Successful use-cases



Reallocating time from device-measured sleep, sedentary behaviour or light physical activity to moderate-to-vigorous physical activity is associated with lower cardiovascular disease risk – Walmsley et al. 2020

# Capture-24: Successful use-cases

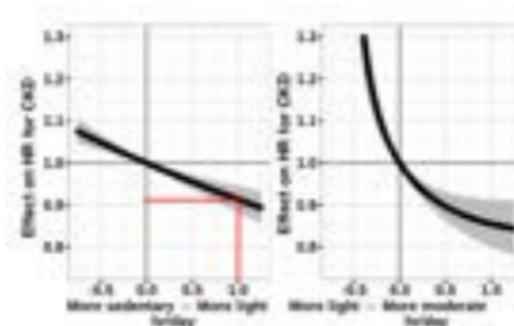
## Oxford Intensive 2-weeks Workshop Results

Depression			
Exposure	OR (95% CI)	p-value	Model pseudo R <sup>2</sup>
WHO 2 quartile - reference	1.00 (1.00, 1.00)	< 0.001	0.000
WHO 3 quartile	1.26 (1.00, 1.60)	0.047	0.000
WHO 4 quartile	1.84 (1.30, 2.60)	0.001	0.000
WHO 5 quartile	2.60 (1.80, 3.80)	0.001	0.000

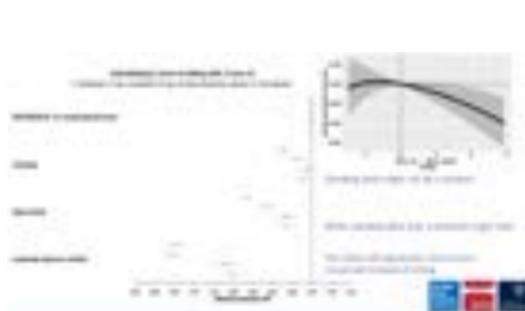
  

Depression			
Exposure	OR (95% CI)	p-value	Model pseudo R <sup>2</sup>
WHO 2 quartile - reference	1.00 (1.00, 1.00)	0.001	0.000
WHO 3 quartile	1.40 (1.00, 1.90)	0.047	0.000
WHO 4 quartile	1.80 (1.30, 2.50)	0.001	0.000
WHO 5 quartile	2.60 (1.80, 3.80)	0.001	0.000

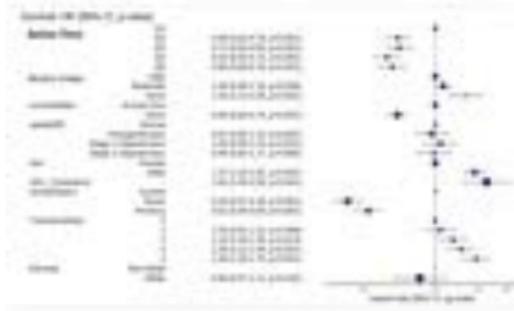
Mental disorders



Chronic kidney disease



Type II Diabetes



All-cause mortality

# Capture-24: Resources & Links

## Resources

- ▶ Accelerometer processing & analysis tool:  
<https://github.com/activityMonitoring/biobankAccelerometerAnalysis>
- ▶ Workshop materials: Activity recognition on the Capture-24:  
[https://github.com/activityMonitoring/week1\\_cdt\\_data\\_challenge](https://github.com/activityMonitoring/week1_cdt_data_challenge)

## Contact Us

aiden.doherty@ndph.ox.ac.uk

shing.chan@ndph.ox.ac.uk

# NHS Digital Academy

Opportunities to embed reproducible machine learning across national training programmes

Dr Athanasios Tsanas ('Thanasis')  
Associate Professor in Data Science  
Usher Institute, Medical School  
University of Edinburgh

[atsanas@ed.ac.uk](mailto:atsanas@ed.ac.uk)



DARTH



THE UNIVERSITY  
of EDINBURGH

 usher  
institute

# NHS Digital Academy



Health Education England  
Digital Transformation

[Home](#)

[About the NHS Digital Academy](#)

[Our programmes](#)

A large group of diverse people, mostly wearing blue t-shirts with the NHS Digital Academy logo, posing for a group photo. The photo is taken from a slightly elevated angle, showing the tops of their heads and shoulders. The background is dark, making the blue shirts stand out.

**NHS Digital  
Academy**

# Background in a nutshell

- Wachter review for UK healthcare
- 2017 Secretary of State commitment
- Train 300 aspiring NHS leaders
- Phase1: 4/2018 – 9/2021
- Programme run by Imperial + Edinburgh



# NHS DA programme

- Leadership
- 6 Modules
- Module 5:  
Actionable Data  
Analytics and  
Clinical Decision  
Support



# Latest developments

- NHS DA renewal (5 years)?
- HDRUK new partner
- Module 5 material widely accessible through HDRUK
- HDRUK teaching platform





Professor Sophie Staniszewska

# The contribution of Public Voice in Machine Learning in Health Care



# Patient and public involvement in research

By public involvement we mean research being carried out ‘with’ or ‘by’ members of the public rather than ‘to’, ‘about’ or ‘for’ them as defined by NIHR INVOLVE

The impact of public involvement in NIHR health and social care research is defined as:

“The changes, benefits and learning gained from the insights and experiences of patients, carers and the public when working in partnership with researchers and others involved in NIHR initiatives”

(NIHR INVOLVE 2019)

# Why involve the public in research?

Makes research more relevant, focused on questions of importance to patients and the public

Enhances quality of research eg. ensuring a trial measures the right outcomes

A moral/ethical imperative

“Nothing about me without me”

Democratic accountability to the taxpayers

ML: Fairness, accountability and transparency



# Co-production

Sharing power

Including all perspectives and skills

Respecting and valuing the knowledge of all

Reciprocity

Build and maintain relationships

Joint understanding and consensus and clarity over roles and responsibilities

WARWICK

MEDICAL SCHOOL



# Progress so far for our project

Advertised for public contributors with responses from 50 people

Interviewed potential participants in March 2021 with academic lead from Warwick, PPI lead from Oxford and public contributor from Warwick

Appointed a panel of 12 people with a range of experience and perspectives who are keen to work with the research team

Strategic decision to recruit public contributors to maximise the diversity of public voice

## Co-production in action:

# The MEMVIE Study – An example of the potential of co-production in a complex area

Staniszewska, S., Hill, E.M., Grant, R. *et al.* Developing a Framework for Public Involvement in Mathematical and Economic Modelling: Bringing New Dynamism to Vaccination Policy Recommendations. *Patient* (2021).  
<https://doi.org/10.1007/s40271-020-00476-x>



# What did it involve?

- 21 meetings over 5 years
- Each lasted 2-3 h. Email contact in-between with the group commenting on documents
- Deliberative knowledge space and Think Aloud techniques encouraged ideas and thoughts to emerge
- Public contributors were able to challenge the data, the basis for the collection of data and the interpretation of that data, thinking outside of the box in a safe space where modellers could rework their thinking
- The meetings enabled thematic development over time as the Reference Group contributors worked with the academic contributors on continuous iterations of the emerging framework



# Academic contributor

*“When I joined midway through the duration of the MEMVIE project, I had not had any previous exposure to public involvement as part of the research process. I found it extremely beneficial to have an additional forum to describe our modelling process, discuss model assumptions and examine data. **From my perspective, being given the opportunity to convey the work to public members through reasoned discourse, ensured justification of modelling aspects, aiding model integrity and validity.** In addition, public involvement generated broader discussion surrounding data curation and data collection (such as questionnaire content), producing recommendations that can be used to inform future developments.”*

# Conclusion

Public voice should be an integral part of how we think about Machine Learning

Consider moving from the passive (is this ok with the public?) to co-producing thinking with the public (which concepts to guide us, how do we create new knowledge together, can we develop frameworks for voice embedded in our methods)

Co-production holds great promise for facilitating the exploration of public voice in machine learning through deliberative knowledge spaces where ideas and concepts can be explored

Capture and publish your involvement to contribute to the evidence base to guide practice

# Contact



[Sophie.Staniszevska@warwick.ac.uk](mailto:Sophie.Staniszevska@warwick.ac.uk)

Professor Sophie Staniszevska

Division of Health Sciences

Warwick Medical School

University of Warwick