

The Impact of Training Data Shortfalls on Safety of AI-Based Clinical Decision Support Systems

Philippa Ryan, Berk Ozturk, Tom Lawton, Ibrahim Habli



Introduction

- Developing ML based diabetes comorbidity predictor
 - Provides "independent second opinion" on patient
- Training data is anonymized patient records
 - But hard to ensure balance, reduce bias
- What is safety impact? How do we mitigate this?
- Safety analysis
- Discussion and next steps

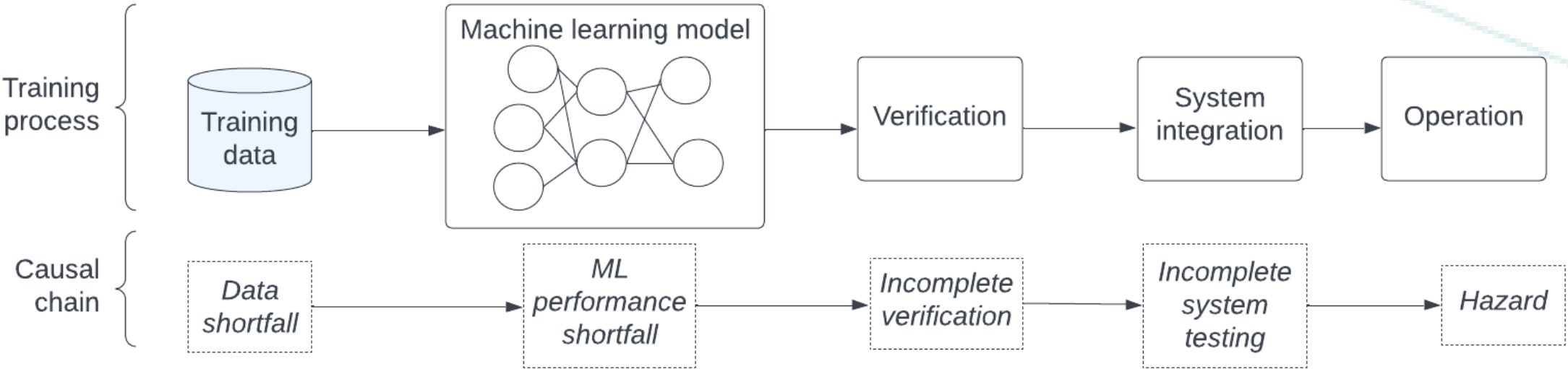
- Funded by
 - EPSRC Assuring Responsibility-Trustworthy Autonomous Systems
 - LRF Assuring Autonomy International Programme

Training Data for ML



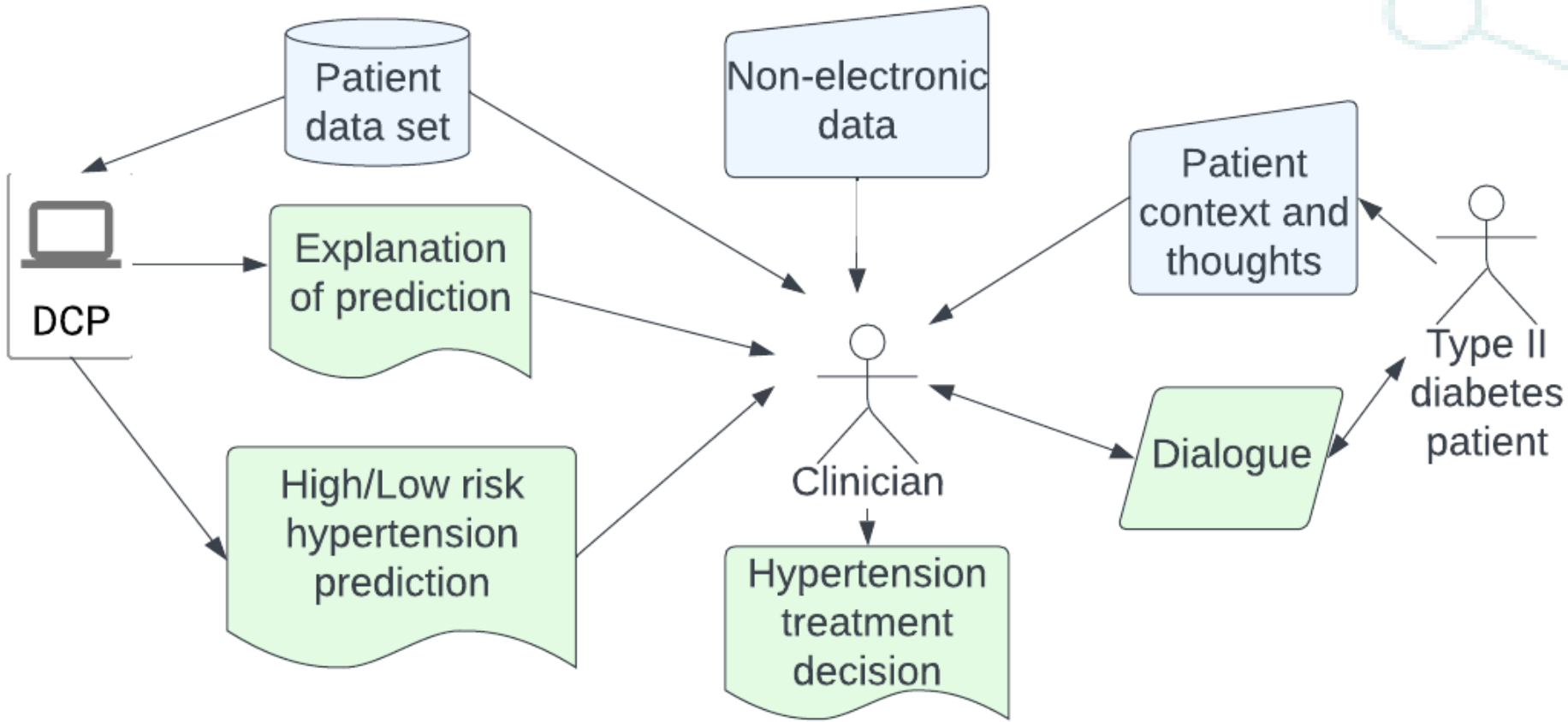
- All Machine Learning needs good quality training data
 - Data embodies the functionality you want it to learn
 - Synthetic user generated data
 - Issues with validity (values, representative of reality)
 - Better for coverage (generate cases)
 - Real world datasets
 - Fewer issues with validity for individual data points
 - Harder to argue future coverage and distribution
- Any problems with training data reflected in final ML

Latent failures



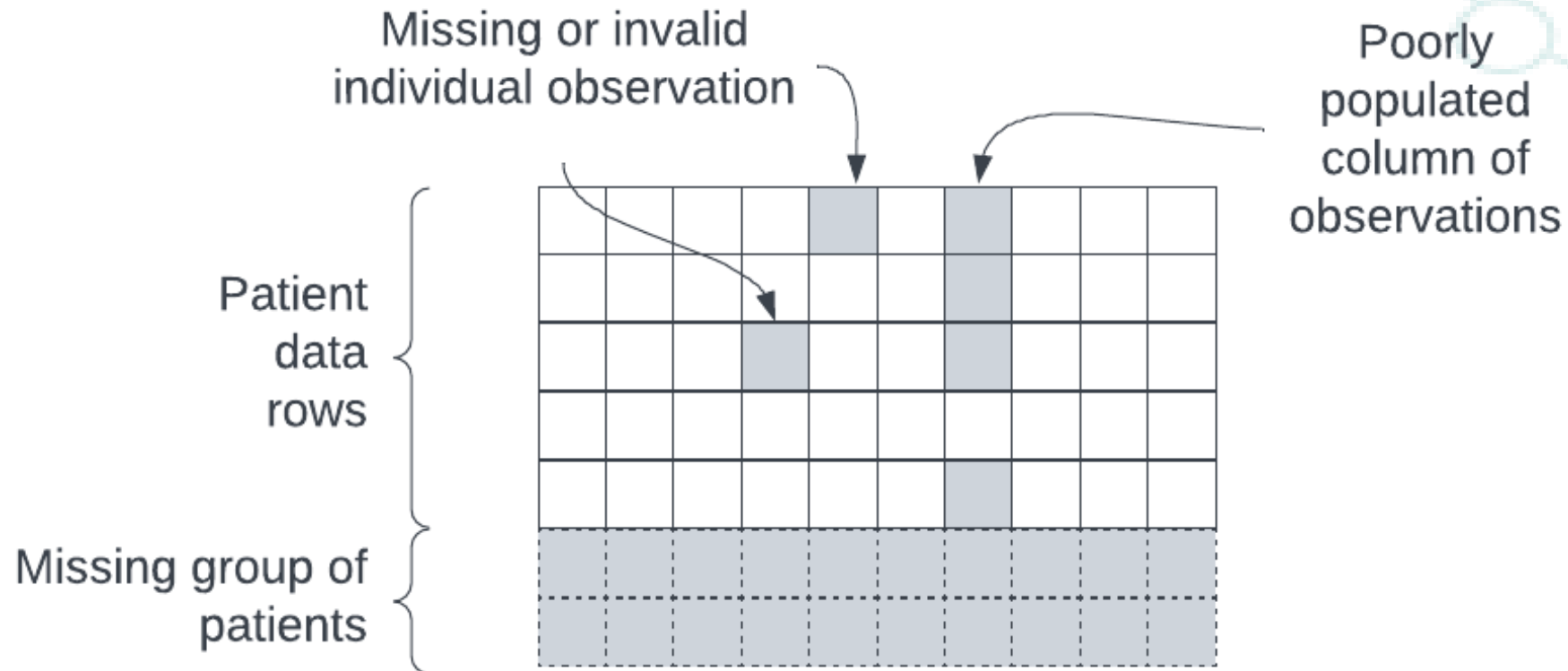
DCP use case

Hypertension version



The training data

Connecting Bradford - database



- 43,000+ data training rows used of Type II Diabetes patients
- Reduced feature space (14,000+) to 20 FOI
 - Reviewed by clinician for validation

What can we do?

Pre-process and synthesize data

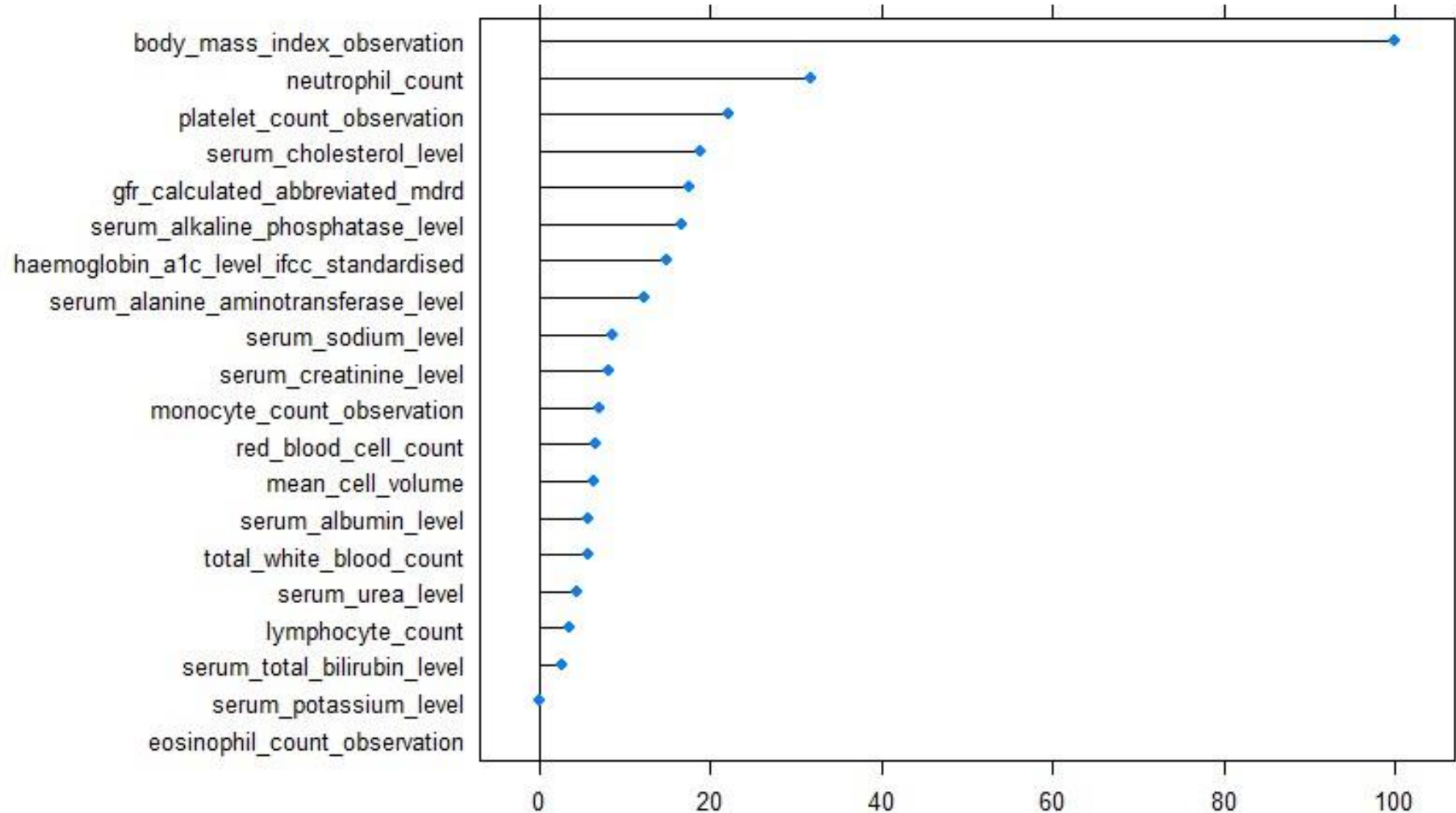
- Missing values common problem with medical diagnosis ML
- Can compensate => data imputation
 - Lots of methods e.g., average, median
 - Bag imputation
 - Uses ML to predict likely values for missing cases
- But can introduce bias
- Concern is understanding *system* risk
- Not just maximise metrics
- Bias considerations must consider system failures

Training process



- Data selection
 - 42,000+ data training rows used of Type II Diabetes patients
 - Reduced feature space (14,000+) to 20 FOI
 - Removed duplicate records
 - Normalised values
 - Compensated missing values using bag imputation
- Trained multiple ML models
 - Naïve Bayes, NN, random forest, SVM
- Ensemble gave best results
 - Accuracy and Kappa values
- NICE guidelines used

Feature Importance Levels



Hazards

Summary

- DCP output could influence decision
- False positive
 - Patient categorised high risk when they are not
 - Provided with medication they don't need with side-effects (severe)
- False negative
 - Patient categorised low risk when they are not
 - Risk of heart attack/stroke (catastrophic)
- Likelihood of incorrect diagnosis from DCP hard to predict
 - Varies per patient

Safety analysis

HAZOP type

- “Flow” – training data into the training process
- Guideword examples:
 - More - indicates a bias in the data, e.g., over representation of particular patient group in the dataset
 - No or Not - FOI or set of FOIs are missing
 - Less - fewer examples of FOI than are desirable for good performance are present
 - Early/Before - indicates that a FOI may be present but out of date with respect to the co-morbidity presenting itself
 - Reverse – opposite diagnosis included

Guideword	Deviation	Cause	Effect	Mitigation
No or not	Samples for ethnic group not included in training data (TD)	No/limited patients of ethnic group were patients	ML not trained or verified adequately for ethnic group with higher genetic risk of hypertension	Manual review of DB by expert, show clinician prototypical examples, patient discussion
Part of	Partially missing BMI in TD samples	BMI not consistently recorded	ML performance biased based on the data imputation method used, leads to poor performance for high or low BMI patients	Use bag imputation for TD records to reduce bias, recommend collection of BMI for future TD, show clinician prototype examples, patient discussion
More	Over representation in TD of high BMI patients	Most patients examined had high BMI	Prediction biased towards patients with high BMI, meaning patients with low BMI have less accurate predictions	Manual review of DB by expert, training samples picked across all ranges, show clinician prototype examples, patient discussion
More	Over representation in TD of certain ethnic group	Over diagnosis by trained ML for patients of other ethnic groups	TD dominated by ethnic group with genetic disposition to hypertension	Manual review of DB by expert, show clinician prototype examples, patient discussion
Early/ Before and More	BMI data is out of date and training patients have changed BMI by time of diagnosis	DB not kept up to date, TD sampled from wrong part of patient history	ML underestimates likelihood of hypertension	TD selected from samples near to hypertension diagnosis, manual review of DB by expert, patient discussion
Instead	BMI value no longer highest FOI for some FOI distribution	Performance outlier from ML	Wrong prediction for hypertension	Show clinician FOI from training and for each prediction at point of use, patient discussion

Discussion

- Prototypical examples
 - Issue of patient confidentiality
 - Would need to obfuscate these further
- Limited to 20 FOI during training may miss data patterns
 - Some FOI result of hypertension not cause
- Missing data can be significant
 - Patient too unwell for tests
 - Long term trend in their health
 - Or could just be poor record keeping!
 - How do we incorporate in ML process?
- Scalability
 - How to perform manual review of such a large set of data?

Summary

- Issues with training data lead to latent ML faults
 - Subtle and varied
- Need to understand risk not just maximise metrics
- *System focused* hazard analysis
 - Can help identify risk from bias with more clarity
 - We can put *targeted* mitigations in place
- May be complex trade-offs
- Next steps – developing DCP for myocardial infarction (heart attacks)



UNIVERSITY
of York

**ASSURING
AUTONOMY**
INTERNATIONAL PROGRAMME