

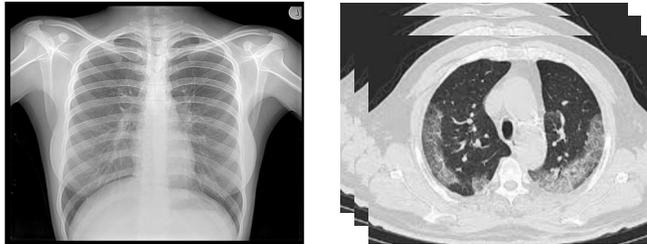
A tale of two crises: COVID-19 and ML reproducibility

Trustworthy AI for Medical and Health Research Workshop

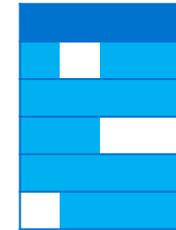
Michael Roberts
Department of Applied Mathematics and Theoretical Physics

A pandemic of reproducibility issues

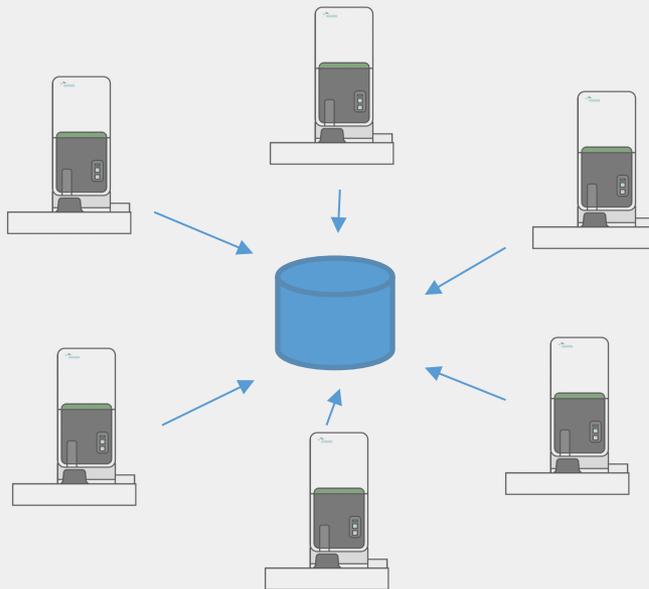
ML for
COVID-19



ML for
Incomplete
Data



ML at
scale



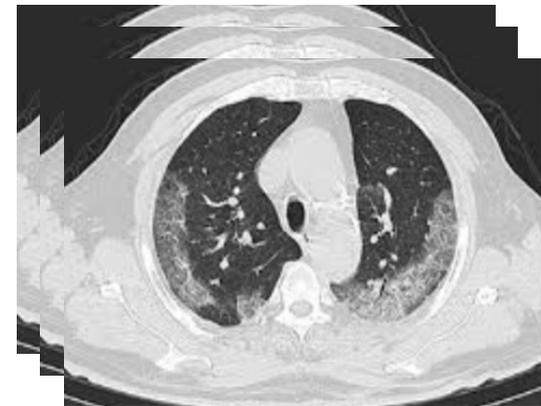
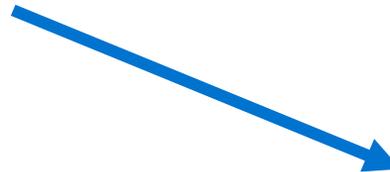
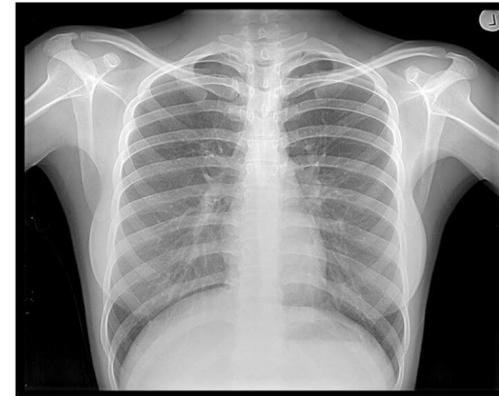
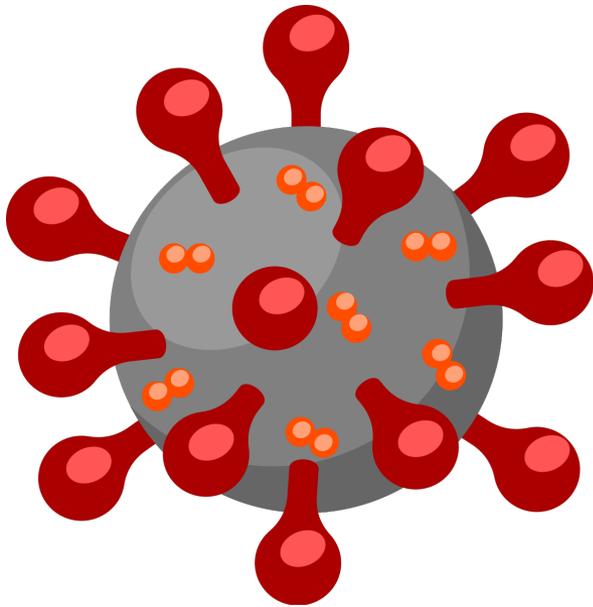
ML and
code



Issues with ...

Joint with Derek Driggs, Matthew Thorpe, Julian Gilbey, Angelica I. Aviles-Rivero, Cathal McCague, James Rudd, Evis Sala, Carola-Bibiane Schönlieb and many AIX-COVNET members

ML for COVID-19 imaging



Roberts, M., Driggs, D., Thorpe, M. et al.
Common pitfalls and recommendations for using machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans.
Nature Machine Intelligence 3, 199–217 (2021).
<https://doi.org/10.1038/s42256-021-00307-0>

nature
machine intelligence

email: mr808@cam.ac.uk

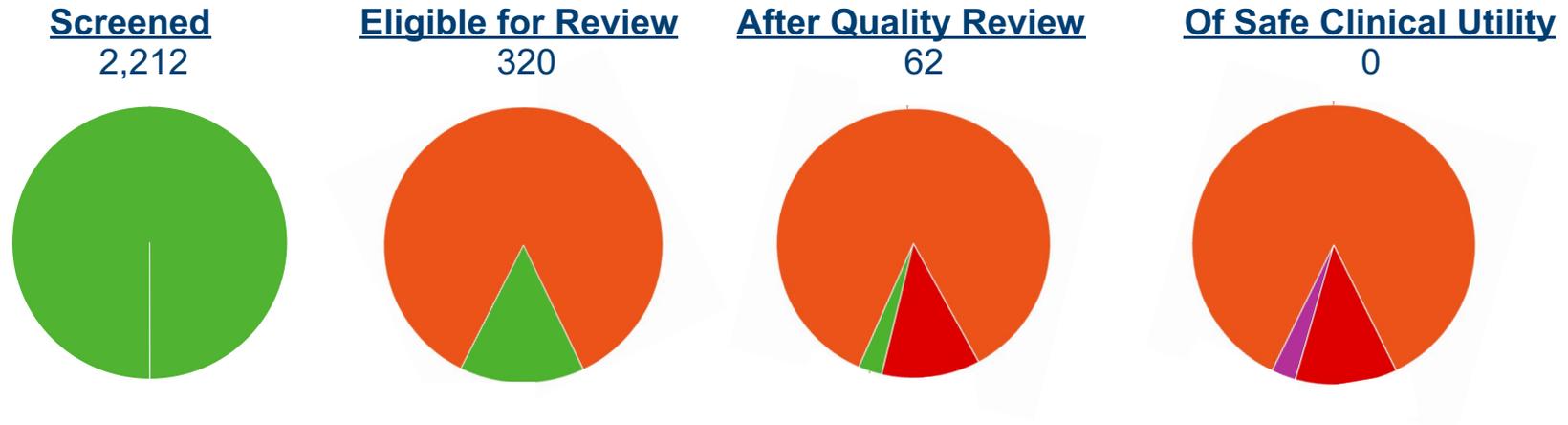
COVID-19: A perfect test case for image based ML

- Never before have we had access to data:
 - at such scale for a single disease
 - all collected around the same time period
 - collected on the same machines
 - for a disease with a short infection time

So why did image based ML fail to contribute significantly to the COVID-19 pandemic?

Systematic Review

- **Eligibility:** Any papers using ML and CXR/CT imaging for COVID-19 diagnosis or prognosis.



Basic pitfalls: Frankenstein datasets

- Know where your data comes from!

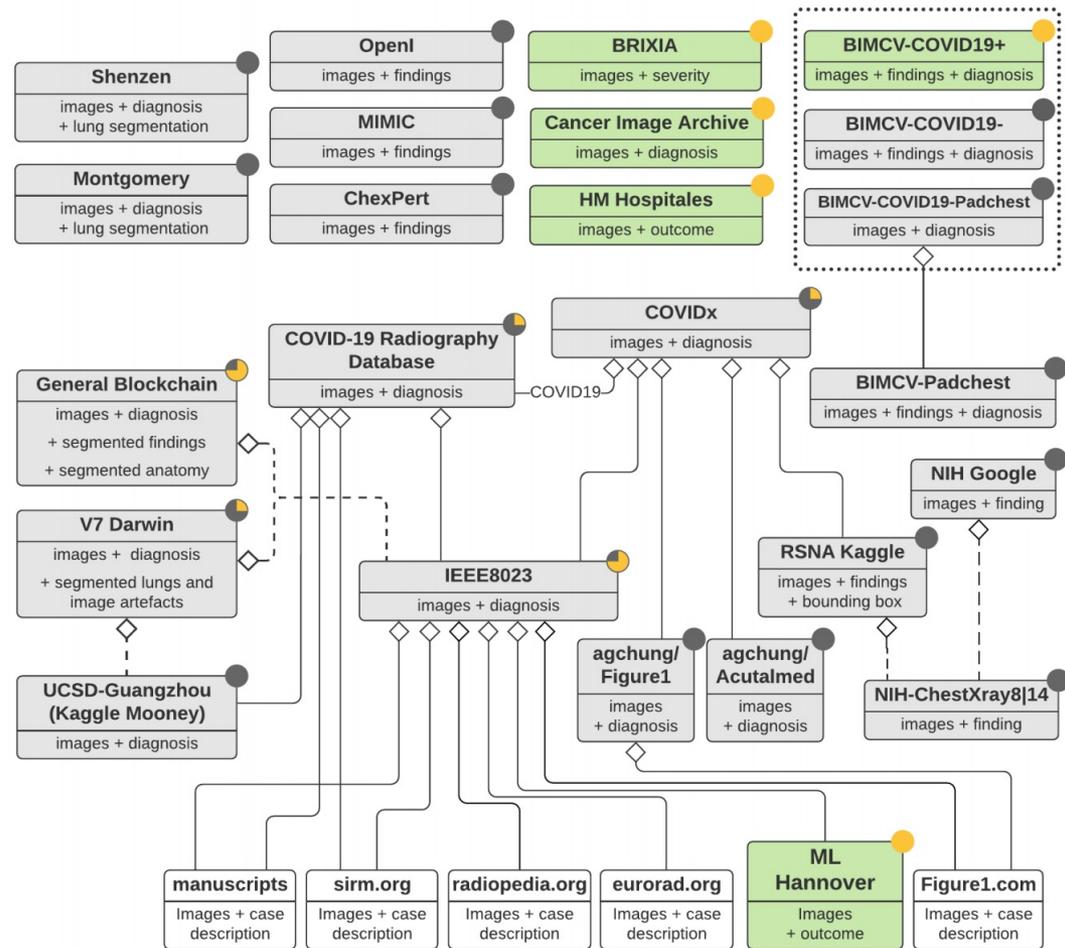


Figure from Cruz et al. (2021)

Basic pitfalls: biases in images

- Know where your data comes from!
- Appreciate the biases in your data.

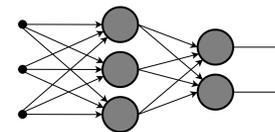


Basic pitfalls: biases in labels

- Know where your data comes from!
- Appreciate the biases in your data.
- Ground truth assigned based on images.



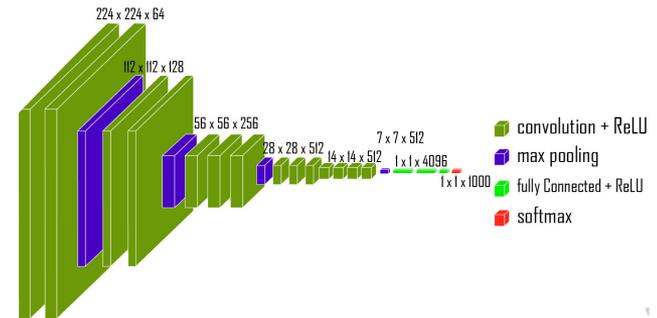
Positive
Negative



Positive
Negative

Basic pitfalls: biases in models

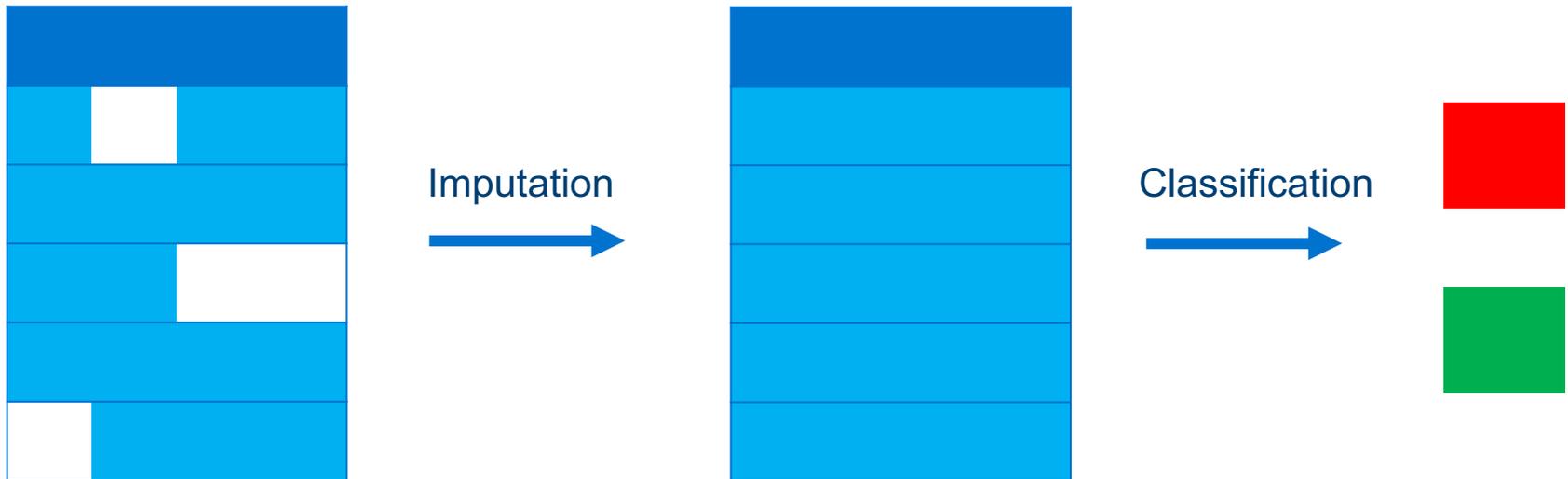
- Know where your data comes from!
- Appreciate the biases in your data.
- Ground truth assigned based on images.
- Resolution driven by pretrained networks



Issues with ...

Joint with Tolou Shadbahr, Julian Gilbey, Jan Stanczuk, Philip Teare, Sören Dittmer, John Aston, Carola-Bibiane Schönlieb and many AIX-COVNET members

ML for Incomplete Data



Why is imputation important?

- QRISK is a model for predicting your risk of a cardiovascular disease.

[1] Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, May M, Brindle P et al. Derivation and validation of QRISK, a new cardiovascular disease risk score for the United Kingdom: prospective open cohort study *BMJ* 2007; 335 :136 doi:10.1136/bmj.39261.471806.55

[2] <https://www.bmj.com/rapid-response/2011/11/01/multiple-imputation-needs-be-used-care-and-reported-detail>

Why is imputation important?

- QRISK is a model for predicting your risk of a cardiovascular disease.
- Initially, the published paper [1] found no link between cholesterol and outcome.

[1] Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, May M, Brindle P et al. Derivation and validation of QRISK, a new cardiovascular disease risk score for the United Kingdom: prospective open cohort study *BMJ* 2007; 335 :136 doi:10.1136/bmj.39261.471806.55

[2] <https://www.bmj.com/rapid-response/2011/11/01/multiple-imputation-needs-be-used-care-and-reported-detail>

Why is imputation important?

- QRISK is a model for predicting your risk of a cardiovascular disease.
- Initially, the published paper [1] found no link between cholesterol and outcome.
- Other researchers [2] found that when only the complete data was considered, the link to cholesterol was found.

[1] Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, May M, Brindle P et al. Derivation and validation of QRISK, a new cardiovascular disease risk score for the United Kingdom: prospective open cohort study *BMJ* 2007; 335 :136 doi:10.1136/bmj.39261.471806.55

[2] <https://www.bmj.com/rapid-response/2011/11/01/multiple-imputation-needs-be-used-care-and-reported-detail>

Why is imputation important?

- QRISK is a model for predicting your risk of a cardiovascular disease.
- Initially, the published paper [1] found no link between cholesterol and outcome.
- Other researchers [2] found that when only the complete data was considered, the link to cholesterol was found.
- After improving the imputation method, the link to cholesterol was recovered.

[1] Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, May M, Brindle P et al. Derivation and validation of QRISK, a new cardiovascular disease risk score for the United Kingdom: prospective open cohort study *BMJ* 2007; 335 :136 doi:10.1136/bmj.39261.471806.55

[2] <https://www.bmj.com/rapid-response/2011/11/01/multiple-imputation-needs-be-used-care-and-reported-detail>

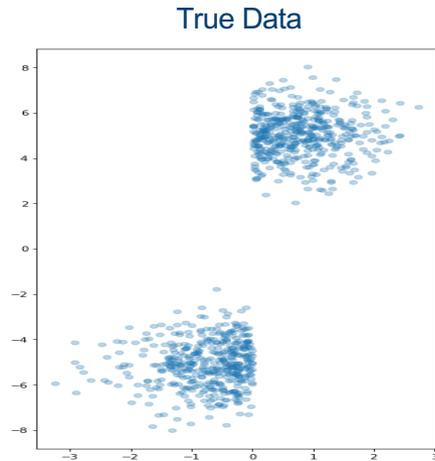
Why is imputation important?

- QRISK is a model for predicting your risk of a cardiovascular disease.
- Initially, the published paper [1] found no link between cholesterol and outcome.
- Other researchers [2] found that when only the complete data was considered, the link to cholesterol was found.
- After improving the imputation method, the link to cholesterol was recovered.
- An improved algorithm is now a standard used in the UK NHS.

[1] Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, May M, Brindle P et al. Derivation and validation of QRISK, a new cardiovascular disease risk score for the United Kingdom: prospective open cohort study *BMJ* 2007; 335 :136 doi:10.1136/bmj.39261.471806.55

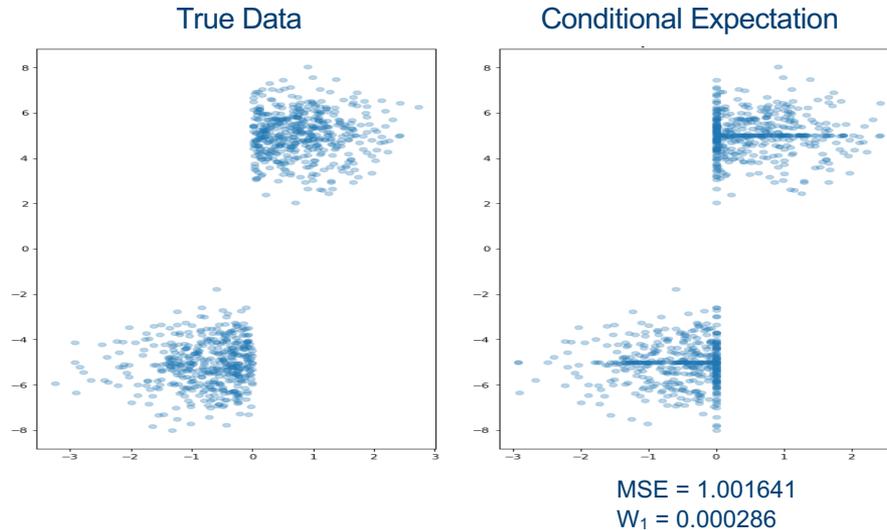
[2] <https://www.bmj.com/rapid-response/2011/11/01/multiple-imputation-needs-be-used-care-and-reported-detail>

How should we measure quality?



Mean / root mean square error is a common metric for measuring imputation quality.

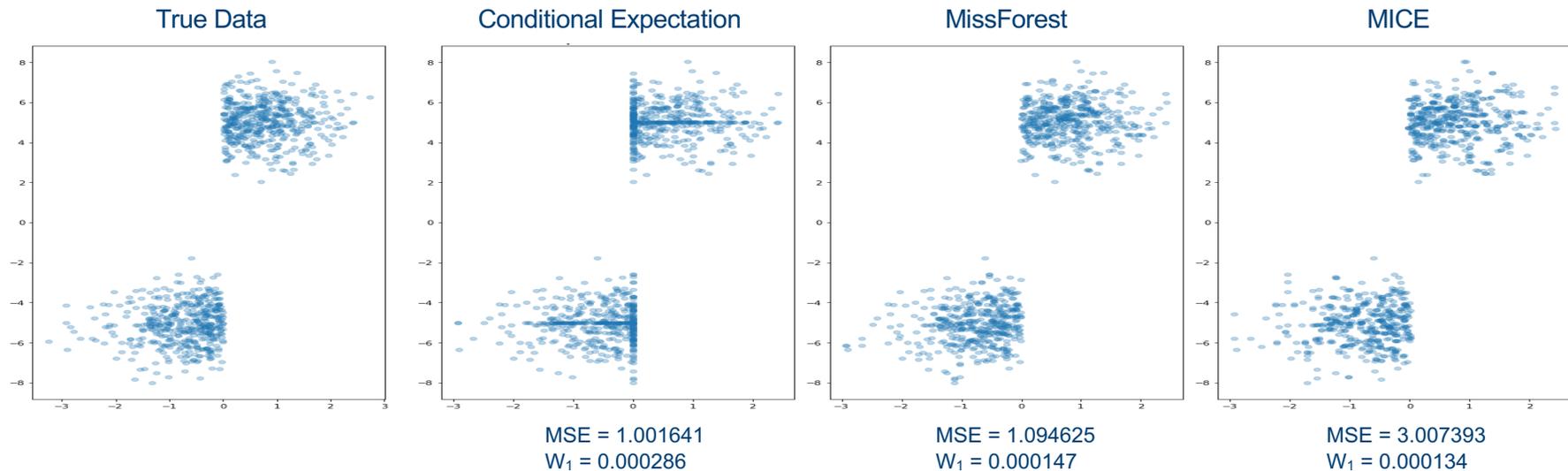
How should we measure quality?



Mean / root mean square error is a common metric for measuring imputation quality.

W_1 = Wasserstein distance between imputed and real value distributions

How should we measure quality?



Mean / root mean square error is a common metric for measuring imputation quality.

W_1 = Wasserstein distance between imputed and real value distributions

What are the issues?

- We find that many new methods fail both to:
 - recreate data distributions
 - give stable imputations
- This compromises model interpretability
- Missingness may be informative for the models

Issues with ...

ML and code



Marauding as software engineers

10 years ago, this was a software engineer....



Marauding as software engineers

10 years ago, this was a software engineer....

Now it is also a:

- Data scientist
- Mathematician
- Statistician
- Machine learning engineer



Issues with the new world

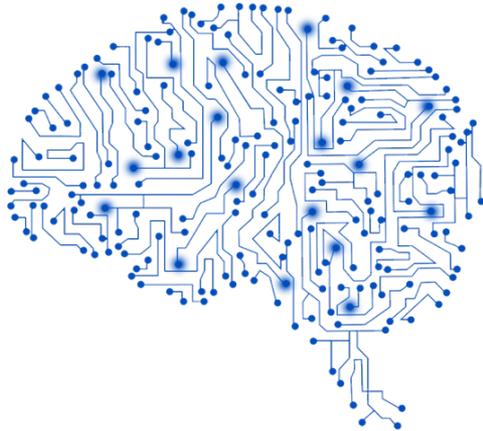
An explosion of data and sources

We must often structure and preprocess this data ourselves

Without deep domain knowledge, we have lost control of the biases in the data

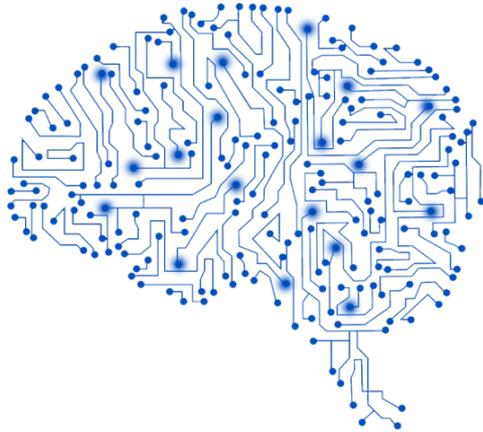


Issues with the new world

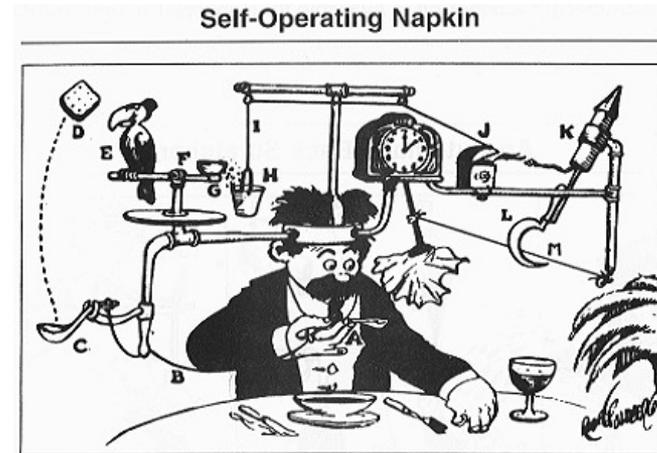


**Idealistic illustration
of machine learning**

Issues with the new world

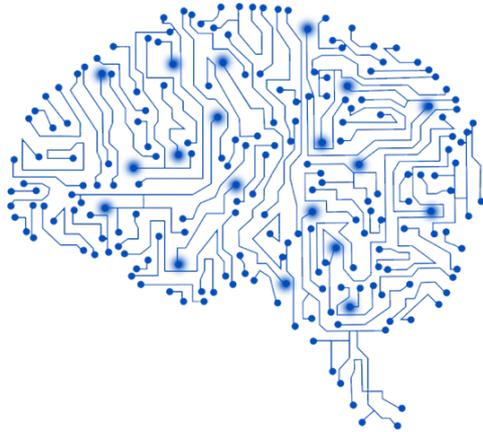


**Idealistic illustration
of machine learning**



The reality ...

Issues with the new world

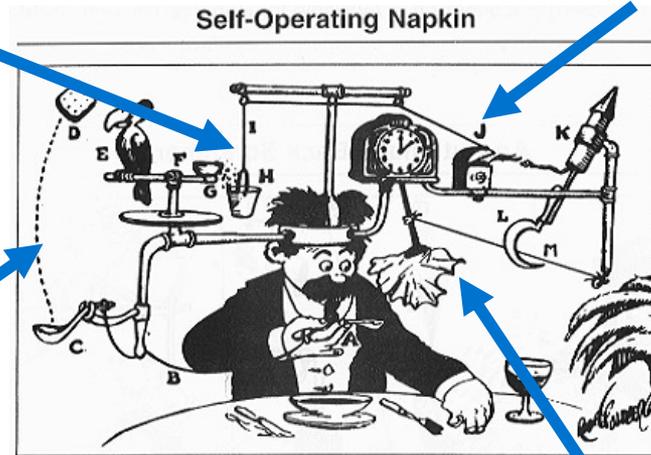


**Idealistic illustration
of machine learning**

**Data
pre-processing**

Model training

**Data
loading**



The reality ...

Evaluation

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.
- Clinical trials are slow, and in phases, for a reason...

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.
- Clinical trials are slow, and in phases, for a reason...
- Create a study plan in advance addressing imbalance, power, etc

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.
- Clinical trials are slow, and in phases, for a reason...
- Create a study plan in advance addressing imbalance, power, etc
- The literature needs purifying.

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.
- Clinical trials are slow, and in phases, for a reason...
- Create a study plan in advance addressing imbalance, power, etc
- The literature needs purifying.
- We all need to use checklists, journals need to enforce their use

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.
- Clinical trials are slow, and in phases, for a reason...
- Create a study plan in advance addressing imbalance, power, etc
- The literature needs purifying.
- We all need to use checklists, journals need to enforce their use
- Complex datasets + complex systems = a nightmare for reproducibility

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.
- Clinical trials are slow, and in phases, for a reason...
- Create a study plan in advance addressing imbalance, power, etc
- The literature needs purifying.
- We all need to use checklists, journals need to enforce their use
- Complex datasets + complex systems = a nightmare for reproducibility
- Practitioners need to rethink coding methodologies, e.g. AGILE

General recommendations / thoughts

- We need to rethink the incentive structure for machine learning based research.
- Clinical trials are slow, and in phases, for a reason...
- Create a study plan in advance addressing imbalance, power, etc
- The literature needs purifying.
- We all need to use checklists, journals need to enforce their use
- Complex datasets + complex systems = a nightmare for reproducibility
- Practitioners need to rethink coding methodologies, e.g. AGILE
- If you don't enforce your assumptions in the code, someone will break it

Thank you