



An Introduction To Structured Missingness

Chris Harbron, Roche

25th April 2022 | For Public Use

Structured Missingness

What Is Structured Missingness?

Structured Missingness as an Emerging Research Area -
Grand Challenges

Characterising Structured Missingness

Acknowledgements

Alan Turing Institute

Robin Mitra
Ben MacArthur
Chris Holmes

Roche

Niels Hagenbuch
Sarah McGough
Ryan Copping

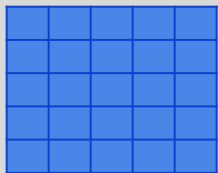
Structured Missingness Workshops Participants

Athena Sheppard	Ginestra Bianconi	Ruben Sanchez-Garcia
Alejandro Frangi	Julia Brettschneider	Ruby Chang
Alisha Davies	Luis Santos	Sara Johansson Fernstad
Aditi Shenvi	Marc De Kamps	Seppo Virtanen
Ana Basiri	Mark Gilthorpe	Sorina Maciuca
Anower Hossain	Michael Barnes	Stefanie Biedermann
Arun Sujenthiran	Maxine Mackintosh	Stefanie Bienert
Brieuc Lehmann	Musa Abdulkareem	Stephen Gardiner
David Leslie	Nina Deliu	Thomas Burnett
Deepak Parashar	Jack Noonan	Timothe Menard
Eleni-Rosalina Andrinopoulou	Nursen Aydin	Trevor Graham
Eda Ozyigit	Paolo Missier	Wenjuan Wang
Francisco Azuaje	Roy Ruddle	Xijin Chen
Ghita Berrada	Rebecca Ward	Xuan Vinh Doan

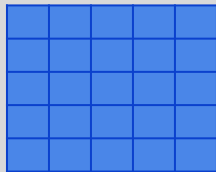
What Is Structured Missingness?

Inevitable Consequence of Combining Multiple Datasets At Scale

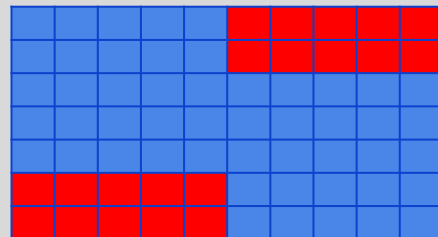
**Adding
Variables /
Different Data
Modalities**



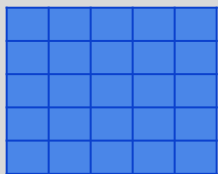
+



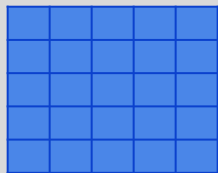
=



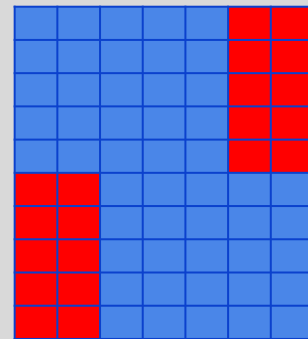
**Adding
Subjects**



+



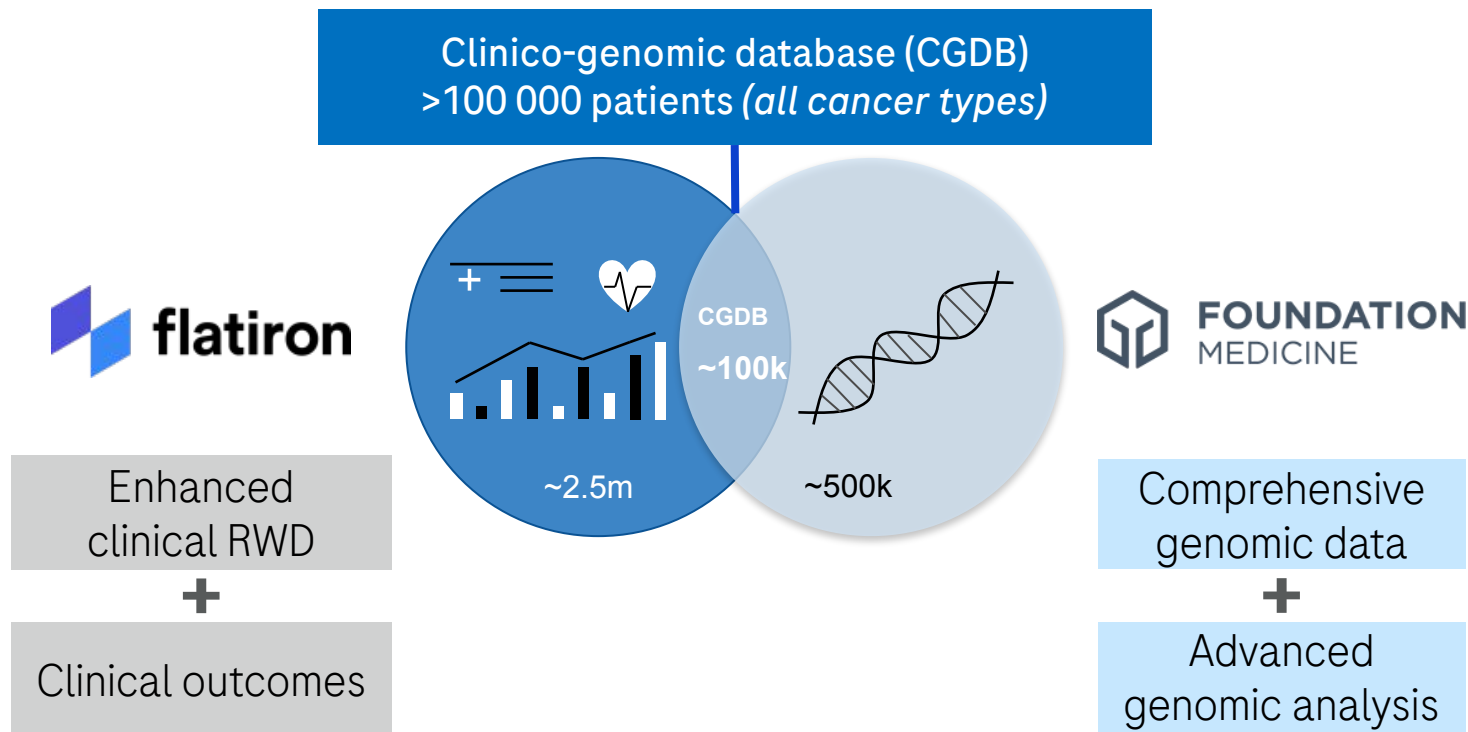
=



Many more subtle effects can also occur

An Example of Structured Missingness - CGDB

The Clinico-Genomic Database links Flatiron electronic health records with Foundation Medicine (FMI) comprehensive genomic profiling for tens of thousands of cancer patients in the U.S.



An Example of Structured Missingness - CGDB

Block Missingness From Measuring Different Genes

Each patient usually receives **1 test**.

Patients receive a variety of **different tests**

Tests are ordered to target specific treatment, prognosis, disease progression goals, and most importantly haematological vs solid tumours.

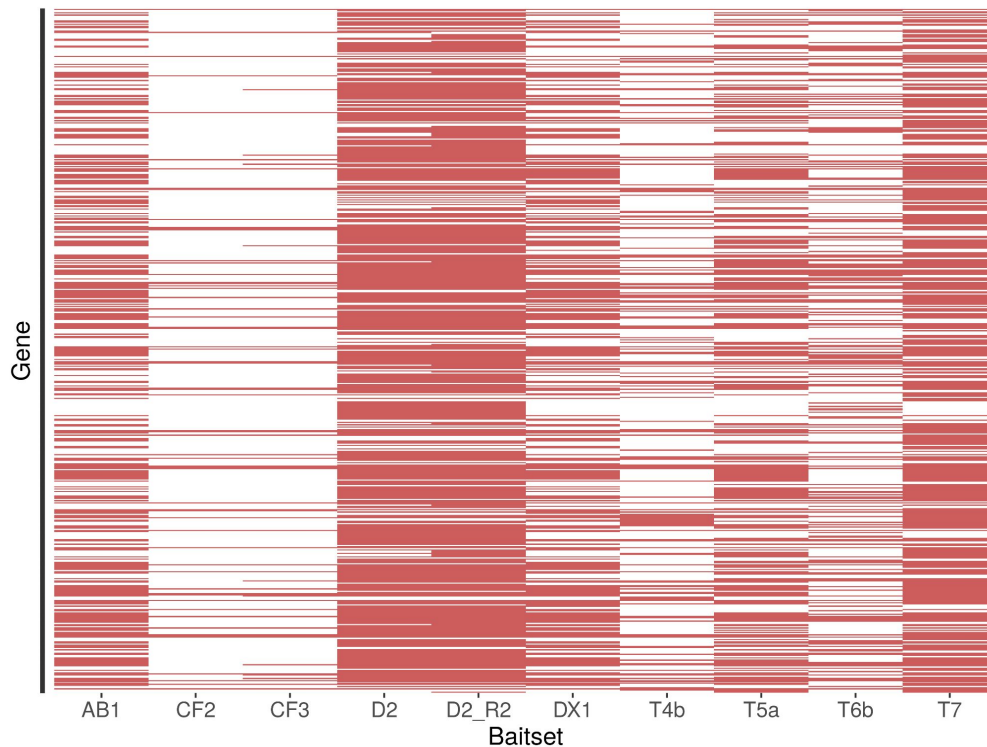
Tests also evolve over time

Tests can use different samples : solid tissue or liquid biopsy

Of **596 unique genes** measured in the CGDB, only **30** are measured across all tests.

Here, genes are **block missing** by test type.

Measured genes by test (baitset).



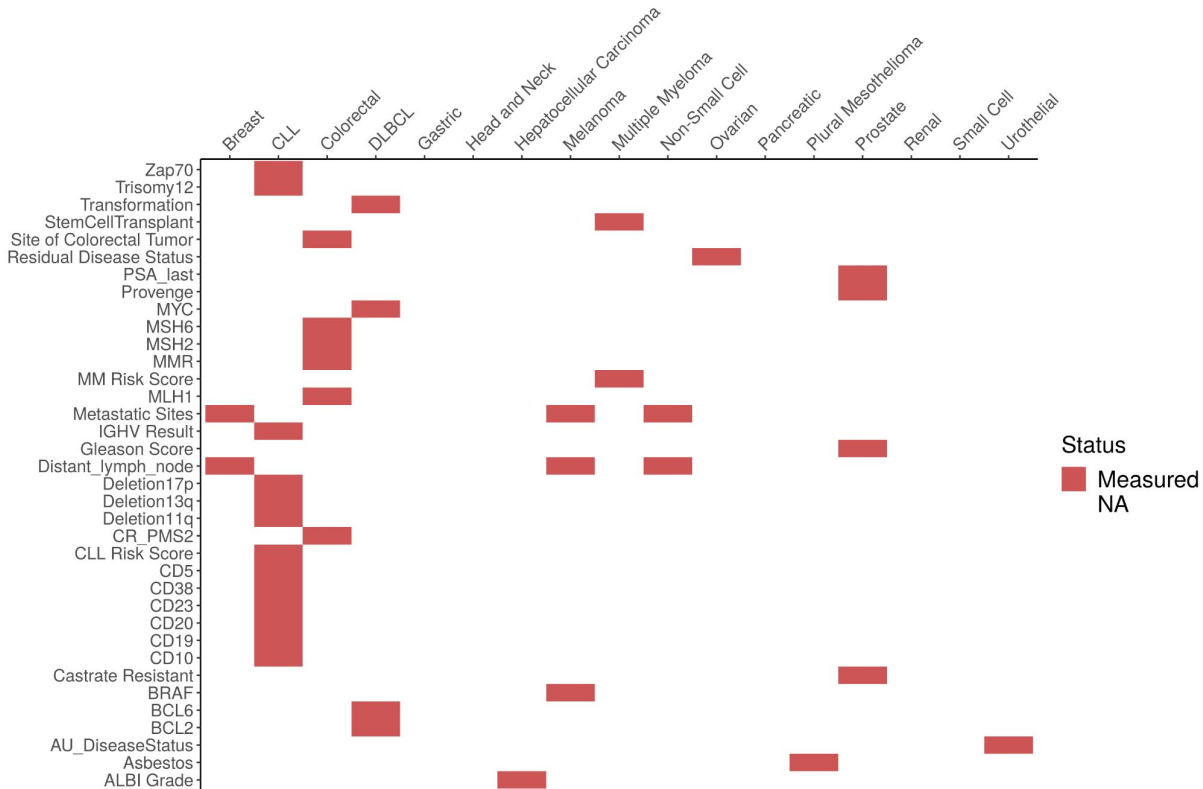
An Example of Structured Missingness - CGDB

Block missingness from combining cancer specific information across dozens of cancer types

Each cancer type collects cancer-specific information, such as the **Gleason Score** for Prostate Cancer patients or **Stem Cell Transplant** for DLBCL patients

Here, variables are “**block missing**” by cancer type.

Is imputation appropriate when the missing value **doesn't exist** or have any meaning?

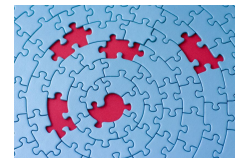


Why Is Structured Missingness Worth Considering?

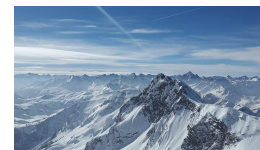
It's **inevitable** & **ubiquitous** if we are combining datasets at scale



Many analytics and machine learning methodologies require **complete data**



The structure may present additional **challenges** or additional **opportunities** compared to unstructured missing data



The missingness may contain information - i.e. be **informative**



The missingness may highlight **limitations** of the data
e.g. underrepresentation



Grand Challenges For Structured Missingness



Defining and Characterising
Structured Missingness



Exploring SM
Geometry and Visualisation



Prediction



Inference and Estimation



Causality



The Role of Imputation



Design Considerations



Benchmarking And Evaluation



Ethical Implications

Characterising Structured Missingness

Multiple Dimensions

1	Relationship of missingness patterns to values	<ul style="list-style-type: none">• MR : Missingness occurs independently• MO : Missingness related to values of other variables• MV : Missingness related to value of variable
2	Nature of relationship of missingness patterns to values	<ul style="list-style-type: none">• D : Deterministic• P : Probabilistic
3	Relationship of missingness patterns to missingness patterns in other variables	<ul style="list-style-type: none">• U : Unstructured• SS : Strong Structure - (Deterministic)• WS : Weak Structure - (Probabilistic)
4	Sub-characterisation by different patterns or structures of missingness	<ul style="list-style-type: none">• e.g.• (B) : Block Missing• (S) : Sequentially Missing
5	Does a missing value exists but is unobserved, or no value exists	<ul style="list-style-type: none">• E : Value exists but was not observed• N : Value doesn't exist for logical/biological reasons

Characterising Structured Missingness

Relationship of Missingness to Other Missingness in Data

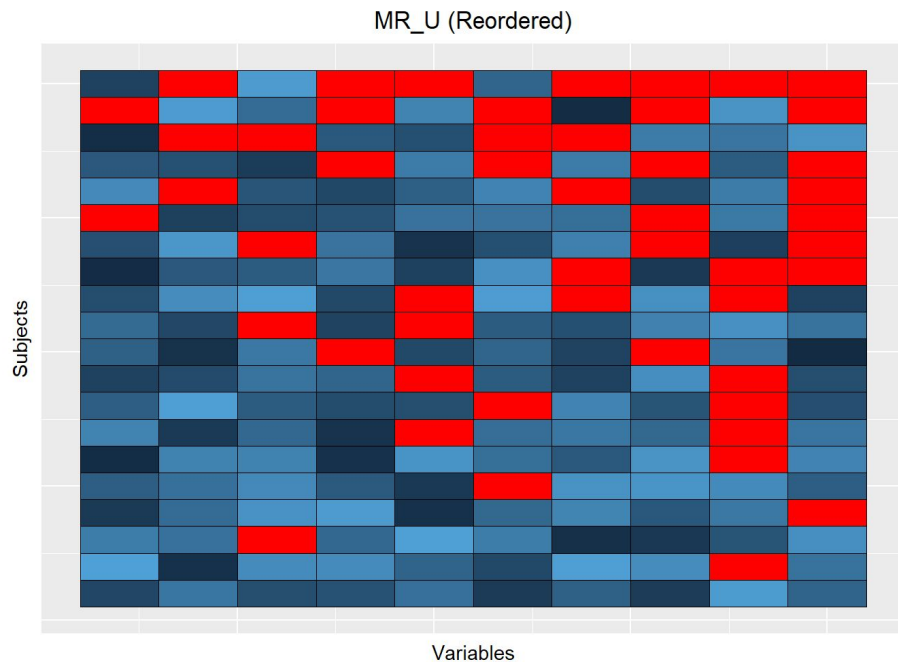
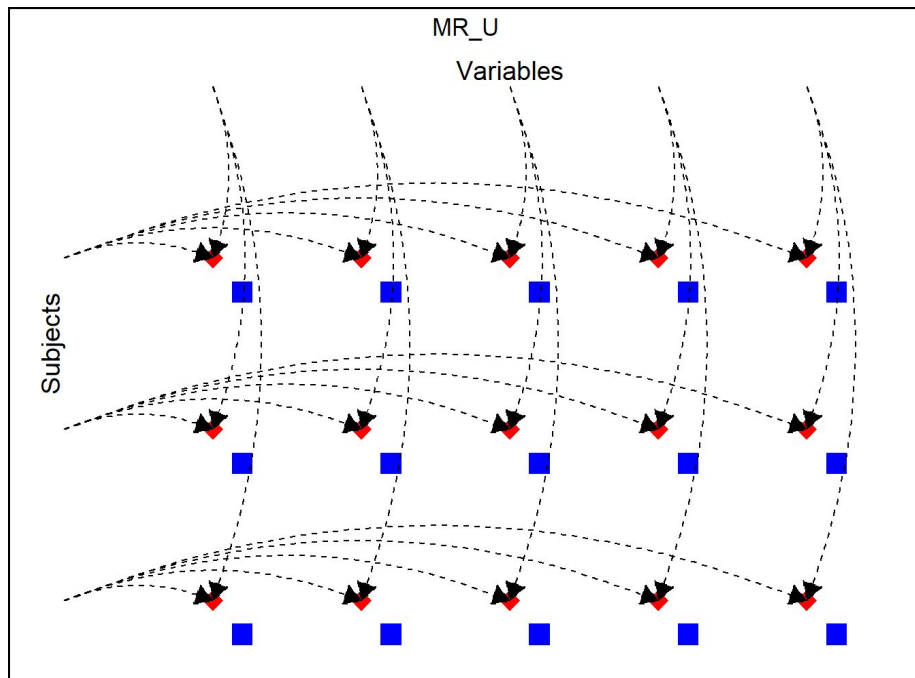
Relationship of Missingness To Values In Data	Unstructured	Structured			
		Strong		Weak	
	MR	MR-U		MR-SS	
				MR-SS(B)	MR-SS(S)
				MR-WS(B)	MR-WS(S)
Relationship of Missingness To Values In Data	MO - Prob	MOP-U		MOP-SS	
				MOP-SS(B)	MOP-SS(S)
				MOP-WS(B)	MOP-WS(S)
	MO - Det	MOD-U		MOD-SS	
				MOD-SS(B)	MOD-SS(S)
				MOD-WS(B)	MOD-WS(S)
Relationship of Missingness To Values In Data	MV - Prob	MVP-U		MVP-SS	
				MVP-SS(B)	MVP-SS(S)
				MVP-WS(B)	MVP-WS(S)
	MV - Det	MVD-U		MVD-SS	
				MVD-SS(B)	MVD-SS(S)
				MVD-WS(B)	MVD-WS(S)

Characterising Structured Missingness

MR_U : Unstructured Missing Randomly

Simplest Case : $P(M_{ij} = 1) = k \quad \forall i, j$

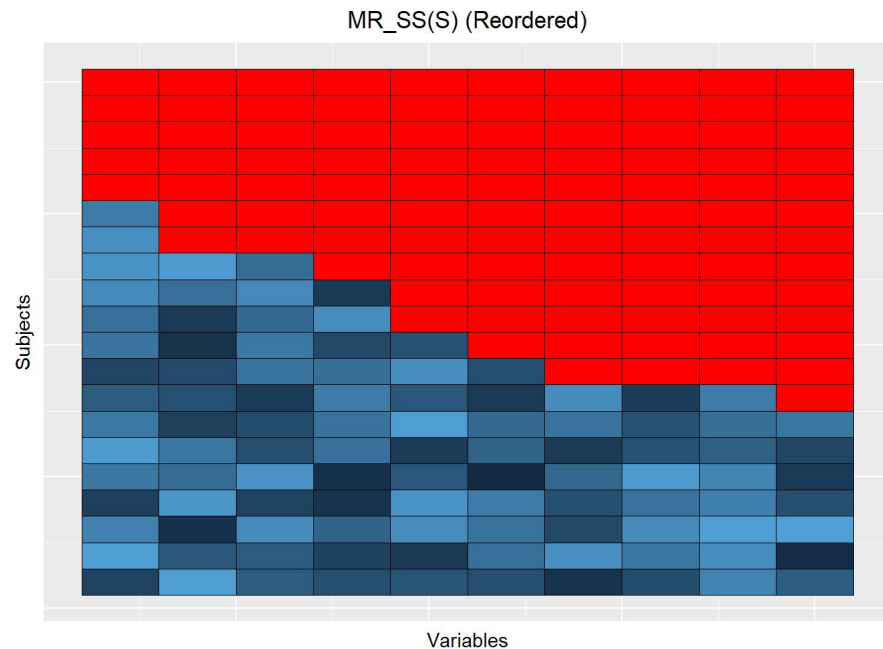
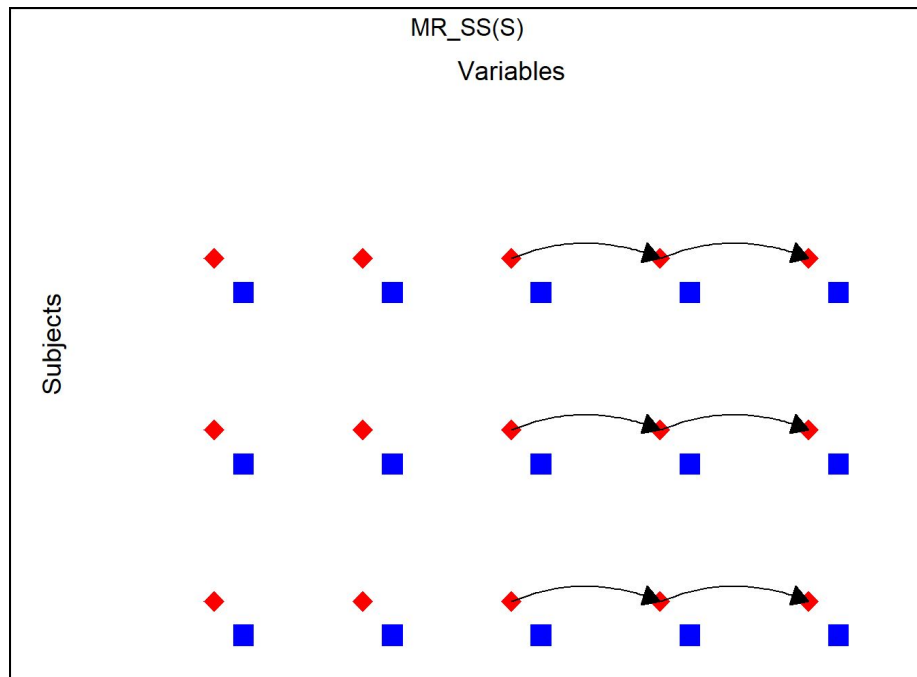
General Case : $P(M_{ij} = 1) = f(s(i), v(j))$



Characterising Structured Missingness

MR_SS(S) : Sequential Strong Structure Missing Randomly

$$M_{ij} = \max(M_{ik:k < j}, P_{ij})$$

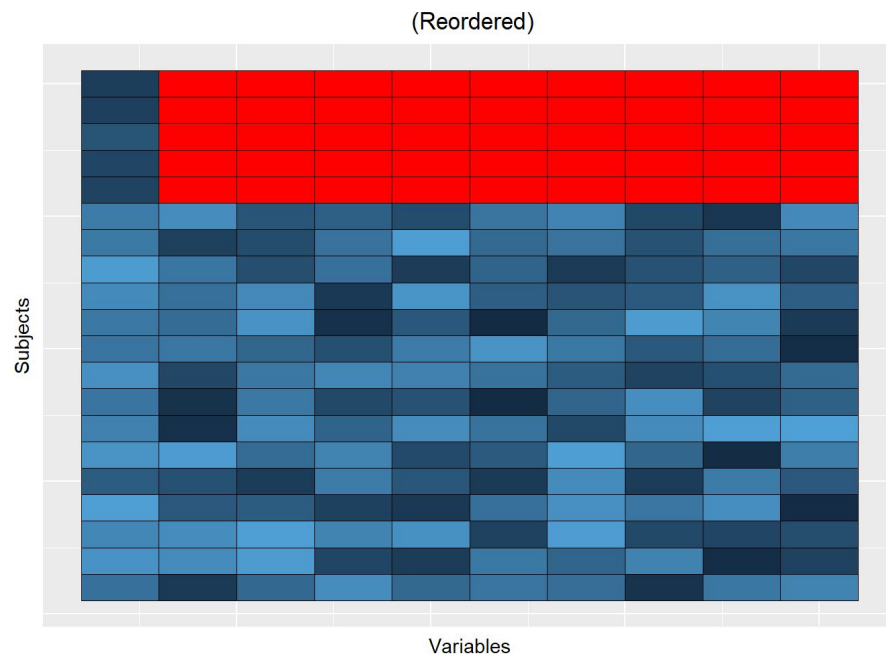
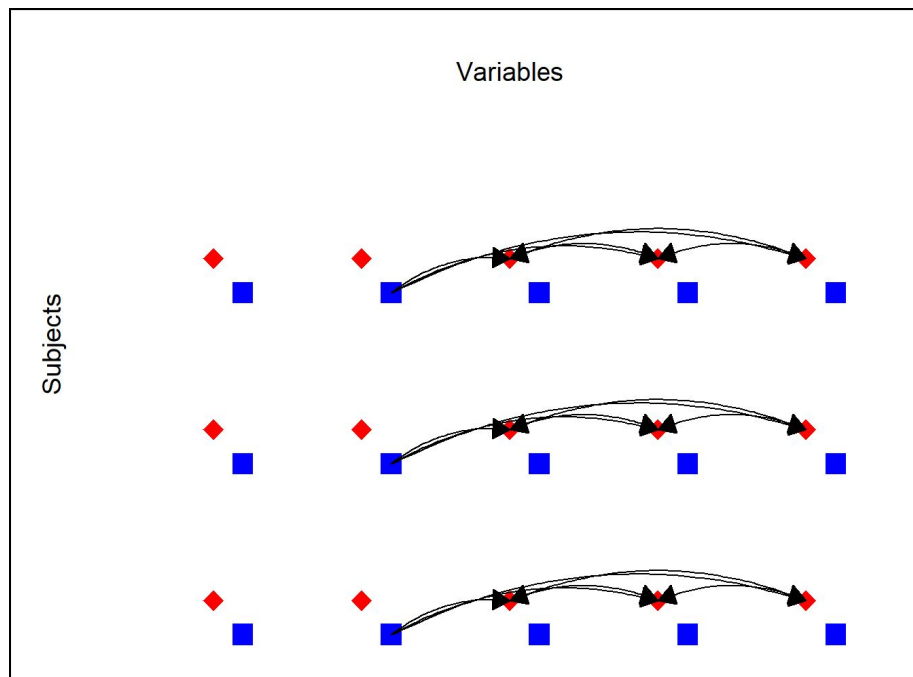


Characterising Structured Missingness

MOD_SS(B):

Strong Block Structure Missing Deterministically Based On Other Variables

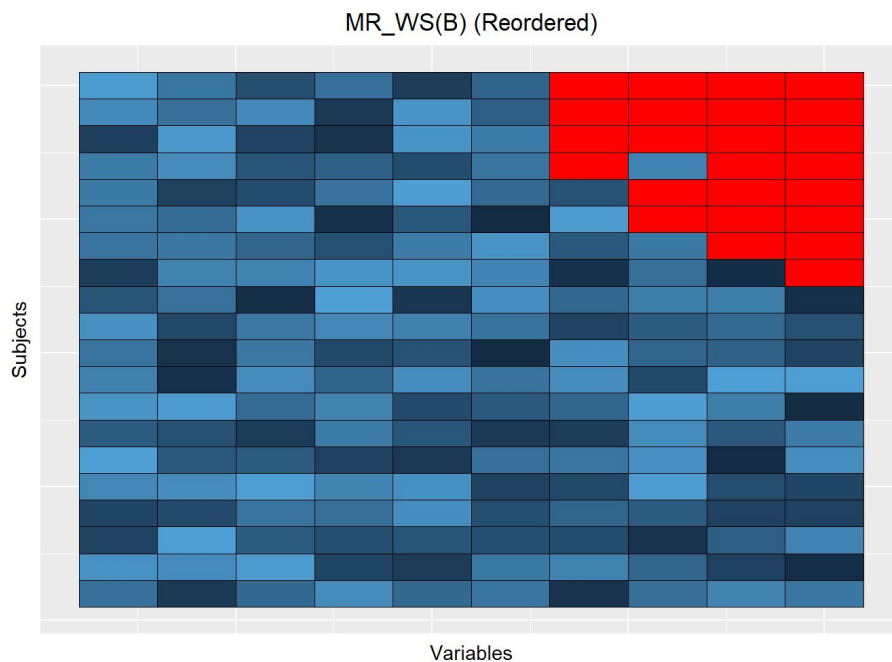
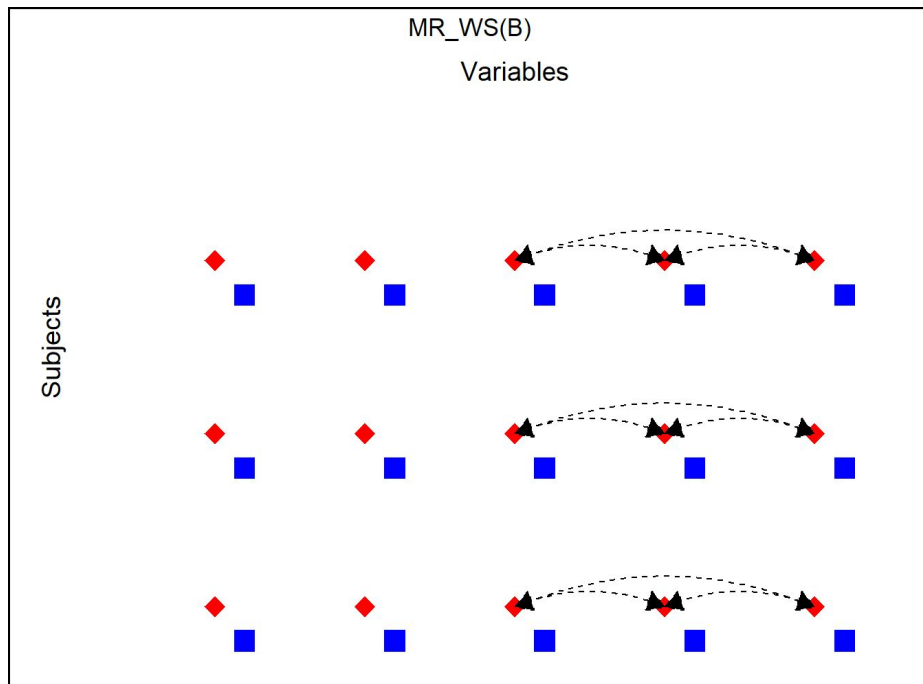
$$M_{ij} = f(X_{ik}) \quad \forall j \in S, k \notin S$$



Characterising Structured Missingness

MR_WS(B) : Weak Block Structure Missing Randomly

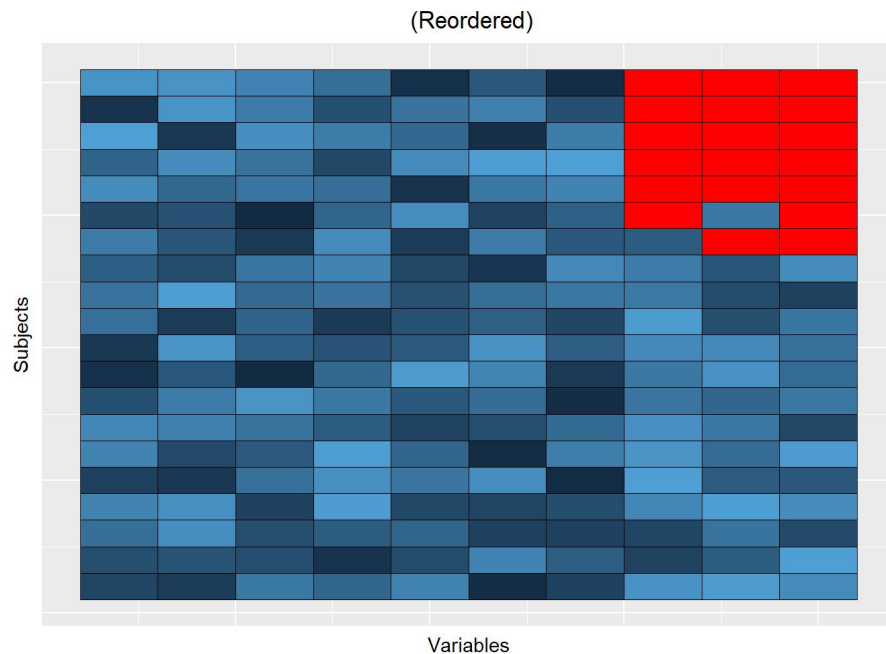
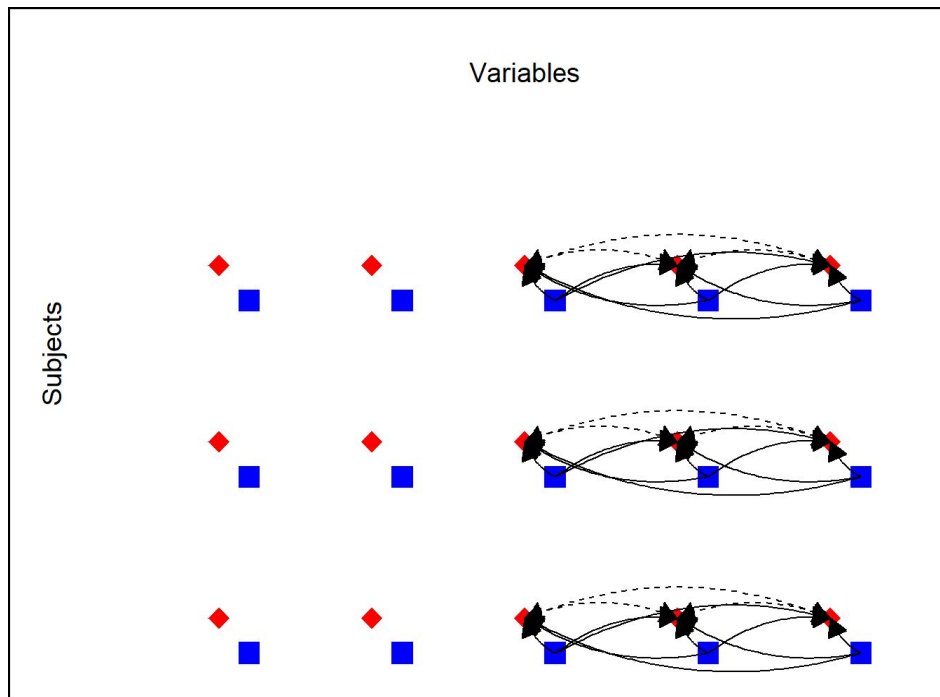
$$P(M_{ij} = 1) = f(M_{i,-j})$$



Characterising Structured Missingness

MVD_WS(B) : Weak Block Structure Deterministically On Variables' Values

$$P(M_{ij} = 1) = f(T_{ij}) \quad j \in S$$



Where Next?



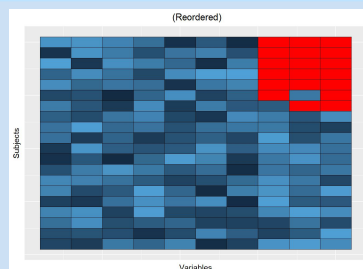
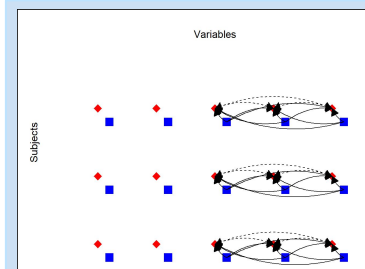
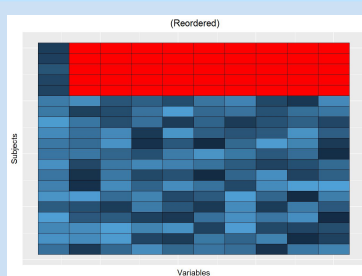
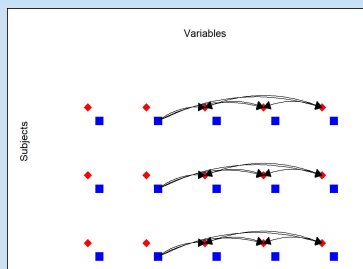
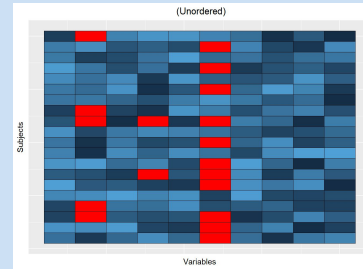
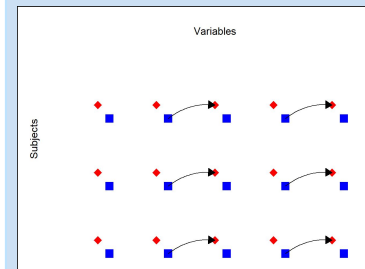
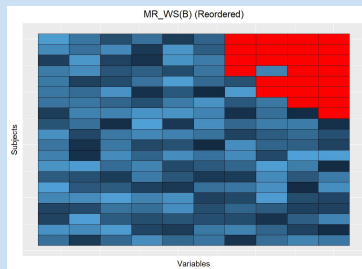
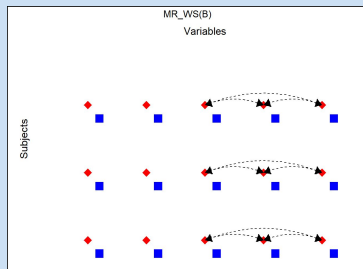
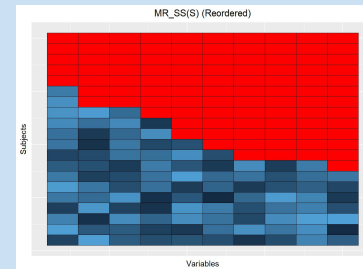
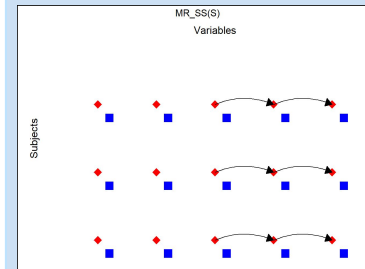
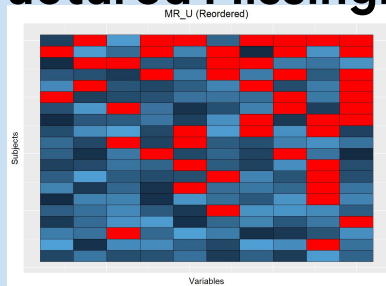
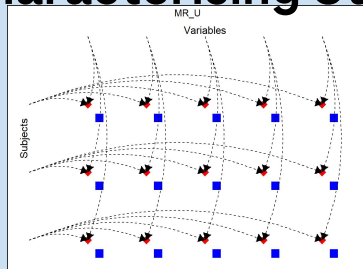
Roche-Turing Partnership Projects in SM

Publications on Grand Challenges in & Characterising SM

Continue to Build an SM Community
Slack Channels & Future Events

Doing now what patients need next

Characterising Structured Missingness



Initial Project Theme

Structured Missingness using CGDB as Motivation

Missing Data is a ubiquitous challenge across healthcare data, which compromises our ability to learn from data. This issue is exacerbated by **structure in the missing values**.

To make the most of data resources we need new methods to handle structured missingness, tailored to the particular challenges of healthcare data.

Clinico-genomic database (CGDB)
~ 80 000 patients
(all cancer types)



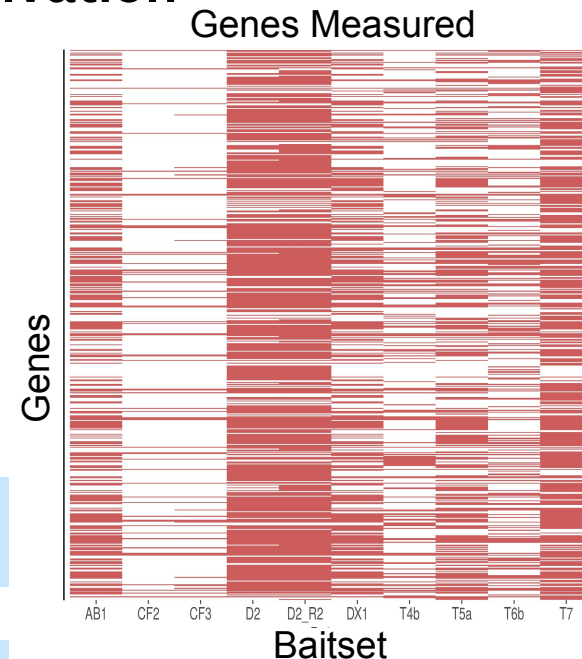
Data Model

Enhanced clinical RWD

Comprehensive genomic data

Clinical outcomes

Advanced genomic analysis



596 genes measured
30 measured across all baitsets

This is an essential building block which will be required across many different AA healthcare applications 20