

On the Stability and **Reproducibility** of Data Science Pipelines

Professor Gavin Brown
University of Manchester, UK



THE

EDITORIALS

POSTDOCS More
but fewer jobs
the way p.4

g
wer
p.441

Reality check

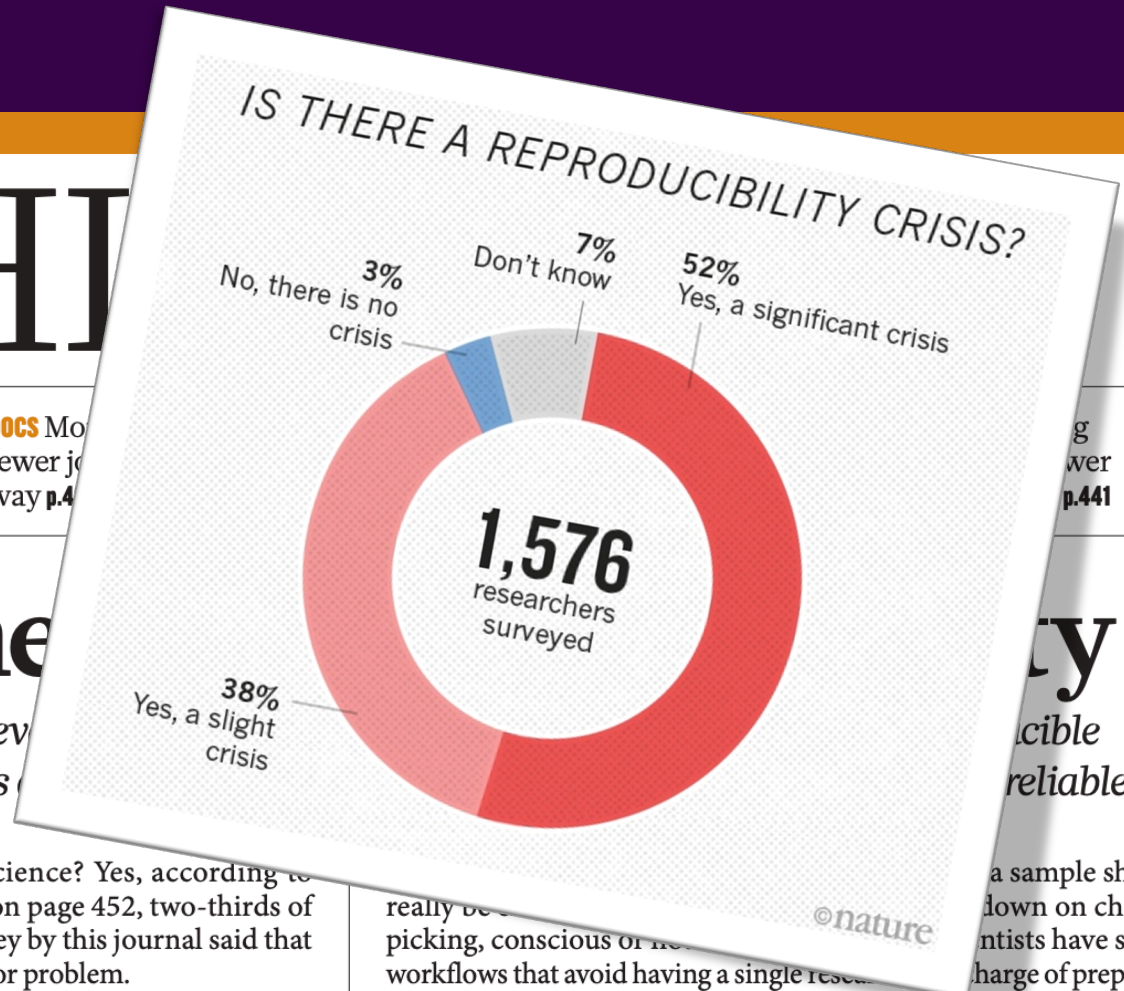
A survey of Nature readers reveals
results. Researchers, funders

ty
cible
reliable.

Is there a reproducibility crisis in science? Yes, according to the readers of *Nature*. As we report on page 452, two-thirds of researchers who responded to a survey by this journal said that current levels of reproducibility are a major problem.

The ability to reproduce experiments is at the heart of science, yet failure to do so is a routine part of research. Some amount of irreproducibility is inevitable: profound insights can start as fragile signals,

really be... down on cherry-picking, conscious or not... workflows that avoid having a single researcher in charge of preparing images or collecting results. Dozens of respondents reported steps to make better use of statistics, randomization or blinding. One described an institution-level initiative to teach scien-



This is **your** data science pipeline.

Observation	Departed	HofReedung	HofReedung	Industrial	QC1
SPR191381	1642	36	36	83304	45.702
SPR191384	15483	36	36	424004	45.483
SPR191386	17620	36	36	628522	50.604
SPR191389	21836	36	36	788106	
SPR191393	22817	36	36	83670	
SPR191398	14821	36	36	5342	
SPR191423	11023	36	36	28624	
SPR191430	19574	36	36	542064	
SPR191377	20829	36	36	371	
SPR191457	20812	36	36	106	
SPR191462	20836	36	36	107036	
SPR191467	17706	36	36	832094	
SPR191464	43473	36	36	1048022	
SPR191411	17495	36	36	626878	
SPR191403	28232	36	36	1018153	
SPR191465	384202	36	36	14183432	
SPR191415	34824	36	36	1572424	
SPR191419	8371	36	36	288056	
SPR191417	8371	36	36	288056	

Customer	ICL_max	ICL_min	ICL_max	ICL_min
1	55	7	44	10CL45
2	57	8	46	
3	53	10	40	
4	49	14	41	45SHCL <5>
5	54	4	49	
6	52	9	50	
7	44	10	52	
8	45	9	50	
9	50	13	53	50SHCL <5>
10	52	3	53	
11	55	5	54	
12	57	1	55	
13	53	11	55	55SHCL <6>
14	47	2	57	
15	55	10	58	
16	48	11	64	65SHCL <5>
17	64	15	65	65SHCL
18	52	15	67	

Tiny change!

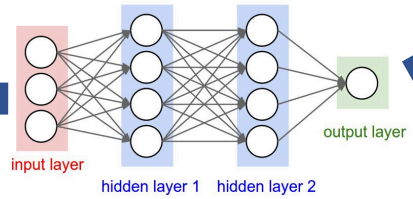
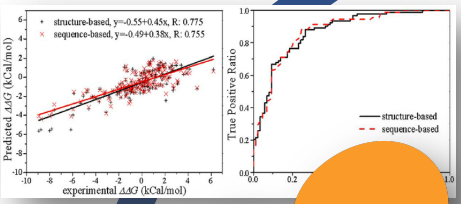
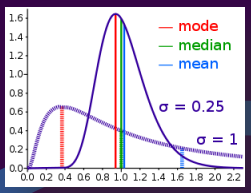


Predictions
Conclusions
Investments

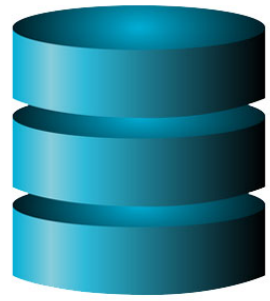


“How much do you **trust** your data pipeline?”

“How **reproducible** is your result?”



Predictions
 Conclusions
 Investments



Is this trustworthy?

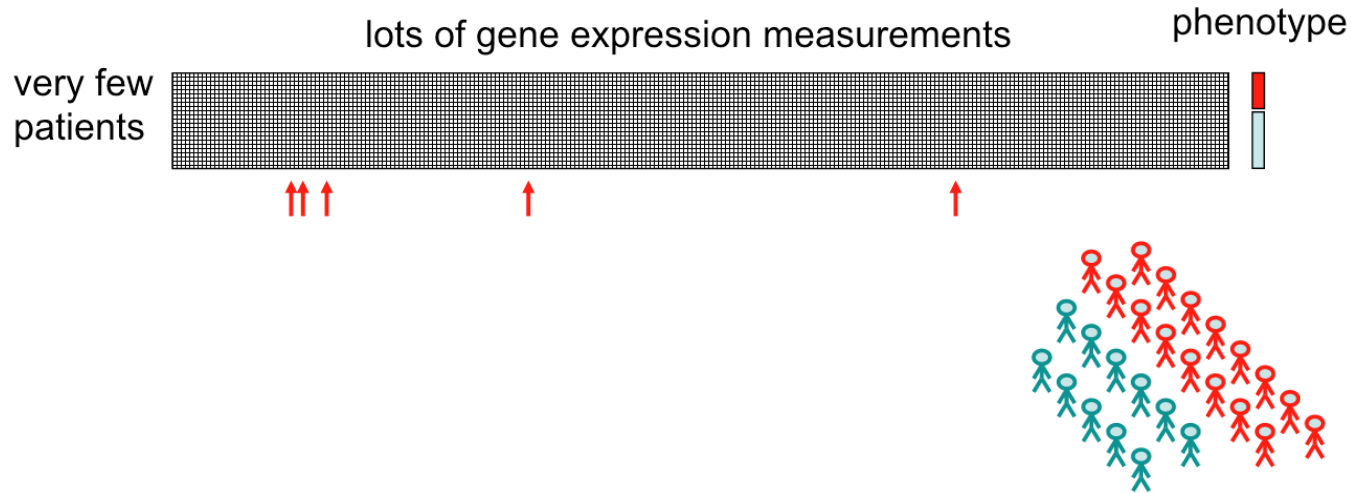
Reproducibility = Trust

Reproducibility is **not a “yes/no”** question.

Conjecture: ***We can measure reproducibility.***

Let's take something specific.

Data-driven biomarker selection



Only a subset of features actually influence the phenotype.

How much do you **trust your choices**?

$x_1, x_2, x_3, x_4, x_5, x_6, x_7, \dots, \text{etc}, \dots, x_{499}, x_{500}$



Your
Data Science
Pipeline



$x_1, x_3, x_5, x_6, x_{493}$

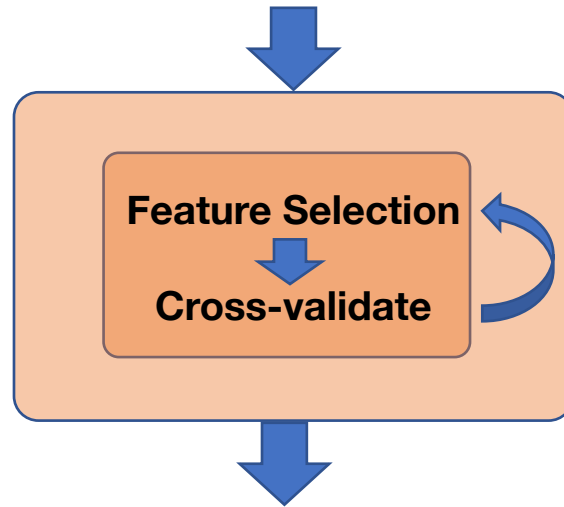


How much do you **trust your choices**?

$x_1, x_2, x_3, x_4, x_5, x_6, x_7, \dots, \text{etc}, \dots, x_{499}, x_{500}$



Your
Data Science
Pipeline



$x_1, x_3, x_5, x_6, x_{493}$

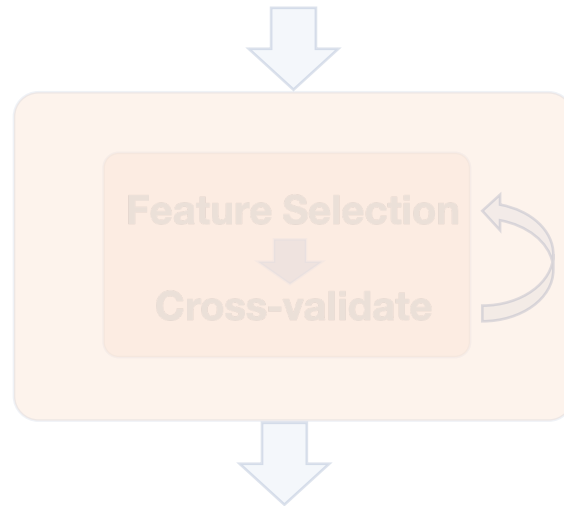


How much do you **trust your choices**?

$x_1, x_2, x_3, x_4, x_5, x_6, x_7, \dots, \text{etc}, \dots, x_{499}, x_{500}$



Your
Data Science
Pipeline



Drop a **random 1%**
of examples

Will not make a difference
...or will it?

“Stability”

x2, x_3, x_5, x_6 , **x491**



“Stability”

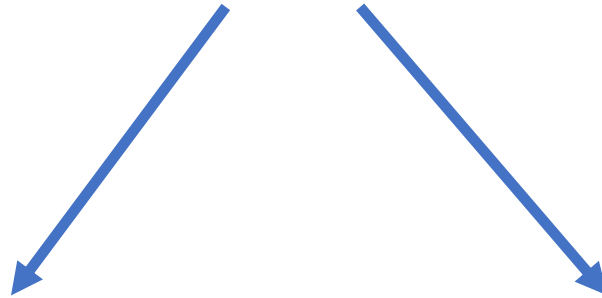
A specific instance of **reproducibility**.

For the task of data-driven **biomarker selection**.

But ... how to measure it?

Estimating **Stability**

[**x1**, x2, **x3**, x4, **x5**, **x6** ..., **x493**, ..., x499, x500]



[x1, x3, x5, x6, x493]

In

[x2, x4, x7, ..., x492, x494... x500]

Out

Estimating **Stability**

Set intersection? (i.e. features in common)

$$\phi(s_i, s_j) = 3$$

My selected biomarkers.

When using **ALL** data.

[x1, x3, x5, x6, x493]

[x2, x3, x5, x6, x491]

Small change if I drop a random **1%** of data

Estimating **Stability**... which set measure?

Dunne et al. (2002)	Hamming	$1 - \frac{ s_i \setminus s_j + s_i \setminus s_j }{d}$
Kalousis et al. (2005)	Jaccard	$\frac{ s_i \cap s_j }{ s_i \cup s_j }$
Yu et al. (2008)	Dice-Sørensen	$\frac{2 s_i \cap s_j }{ s_i + s_j }$
Goh and Wong (2016)	Ochiai	$\frac{ s_i \cap s_j }{\sqrt{ s_i s_j }}$
Shi et al (2006)	POG	$\frac{ s_i \cap s_j }{ s_i }$
Kuncheva (2007)	Consistency	$\frac{r_{i,j} - \frac{k^2}{d}}{k - \frac{k^2}{d}}$
Lustgarten et al. (2009)	Lustgarten	$\frac{r_{i,j} - \frac{k_i k_j}{d}}{\min(k_i, k_j) - \max(0, k_i + k_j - d)}$
Wald et al. (2013)	Wald	$\frac{r_{i,j} - \frac{k_i k_j}{d}}{\min(k_i, k_j) - \frac{k_i k_j}{d}}$

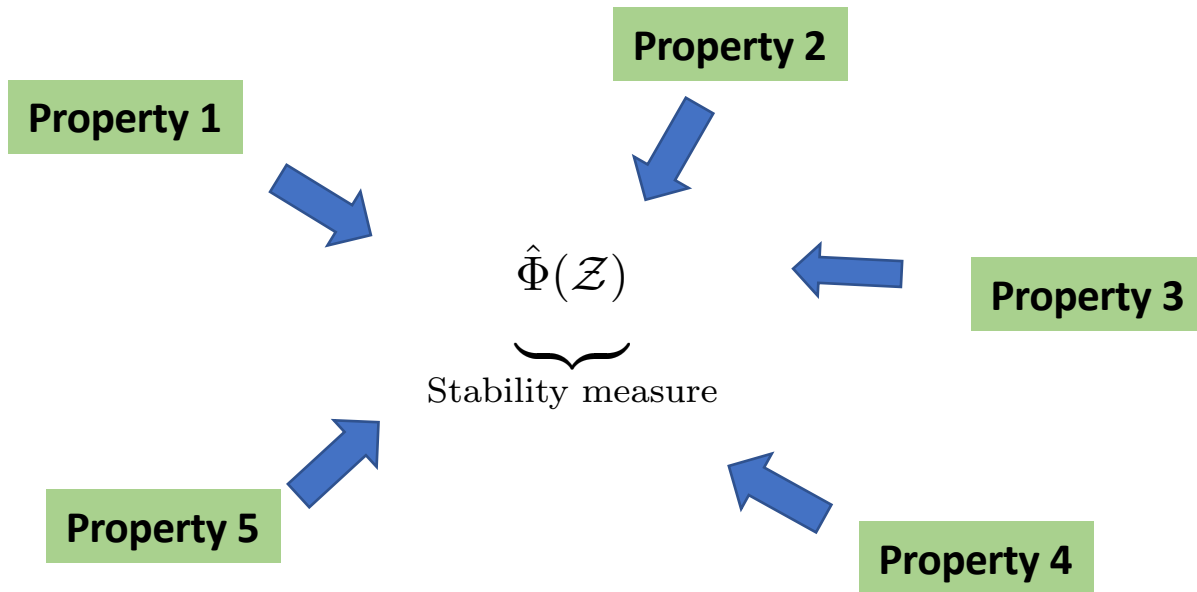
Set theoretic measures
1957-2017.

**Many definitions.
(about 20)**

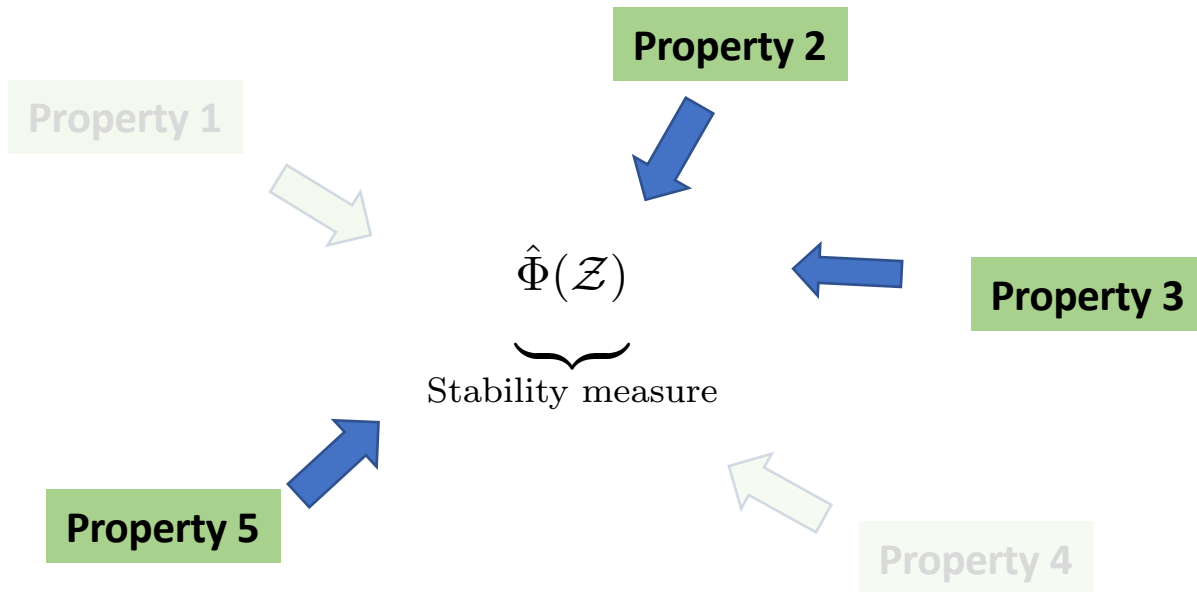
Mostly heuristic.
Conflicting opinions.

**No principled way
to choose
between them.**

What properties do we want?



What properties do we want?



Desirable property 2: *Strict Monotonicity*

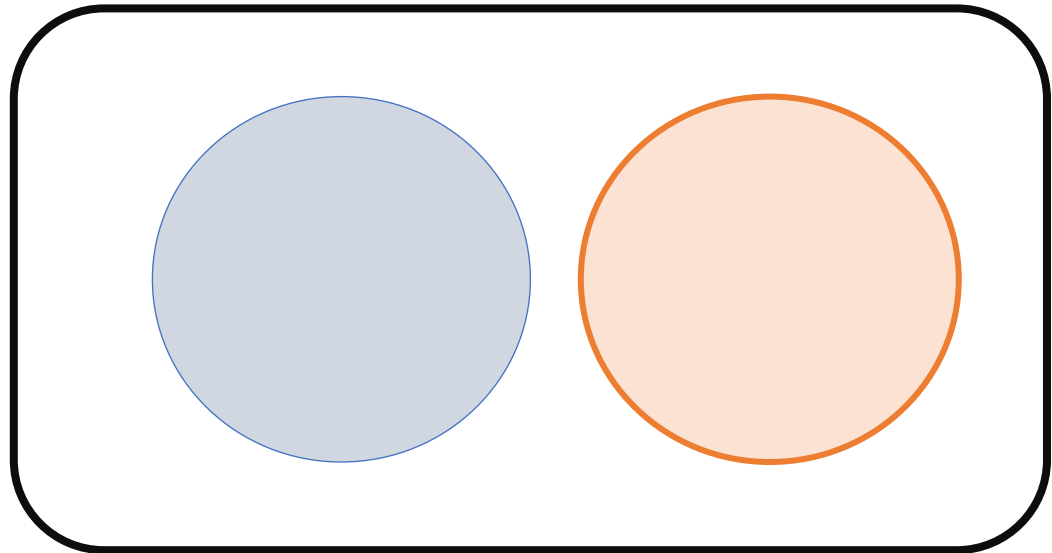
Property 2

...as the sets **overlap** more, the measure should **increase**.

maximal

$\hat{\Phi}(\mathcal{Z})$

minimal



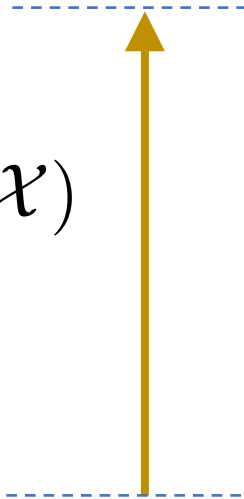
Desirable property 3: *Known upper/lower Bounds*

Property 3

Maximum
stability

$$\hat{\Phi}(\mathcal{X})$$

Minimum
stability

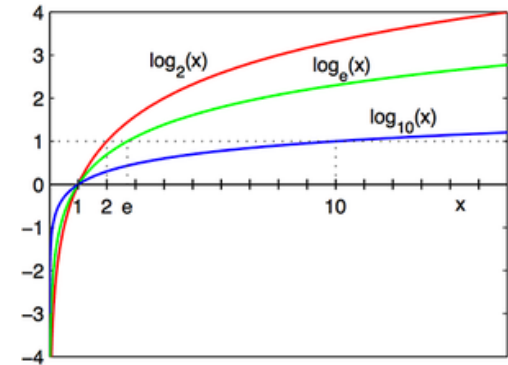


For **interpretability** and **comparison**
across problems/algorithms,

... it should have

known, finite upper/lower bounds,

No logarithms!



Desirable property 5: *Correction for chance*

Property 5

Selecting 2 features from 200...

Trial 1 0000000000**1**00**1**00000000...
Trial 2 0000000000000000**1**000000**1**...

Is very different to selecting 2 from 5...

Trial 1 **11**000
Trial 2 0**11**00

High chance of
intersection,
even if random!

Dunne et al. (2002)	Hamming	$1 - \frac{ s_i \setminus s_j + s_i \setminus s_j }{d}$
Kalousis et al. (2005)	Jaccard	$\frac{ s_i \cap s_j }{ s_i \cup s_j }$
Yu et al. (2008)	Dice-Sørensen	$\frac{2 s_i \cap s_j }{ s_i + s_j }$
Goh and Wong (2016)	Ochiai	$\frac{ s_i \cap s_j }{\sqrt{ s_i s_j }}$
Shi et al (2006)	POG	$\frac{ s_i \cap s_j }{ s_i }$
Kuncheva (2007)	Consistency	$\frac{r_{i,j} - \frac{k^2}{d}}{k - \frac{k^2}{d}}$
Lustgarten et al. (2009)	Lustgarten	$\frac{r_{i,j} - \frac{k_i k_j}{d}}{\min(k_i, k_j) - \max(0, k_i + k_j - d)}$
Wald et al. (2013)	Wald	$\frac{r_{i,j} - \frac{k_i k_j}{d}}{\min(k_i, k_j) - \frac{k_i k_j}{d}}$

Remember this?

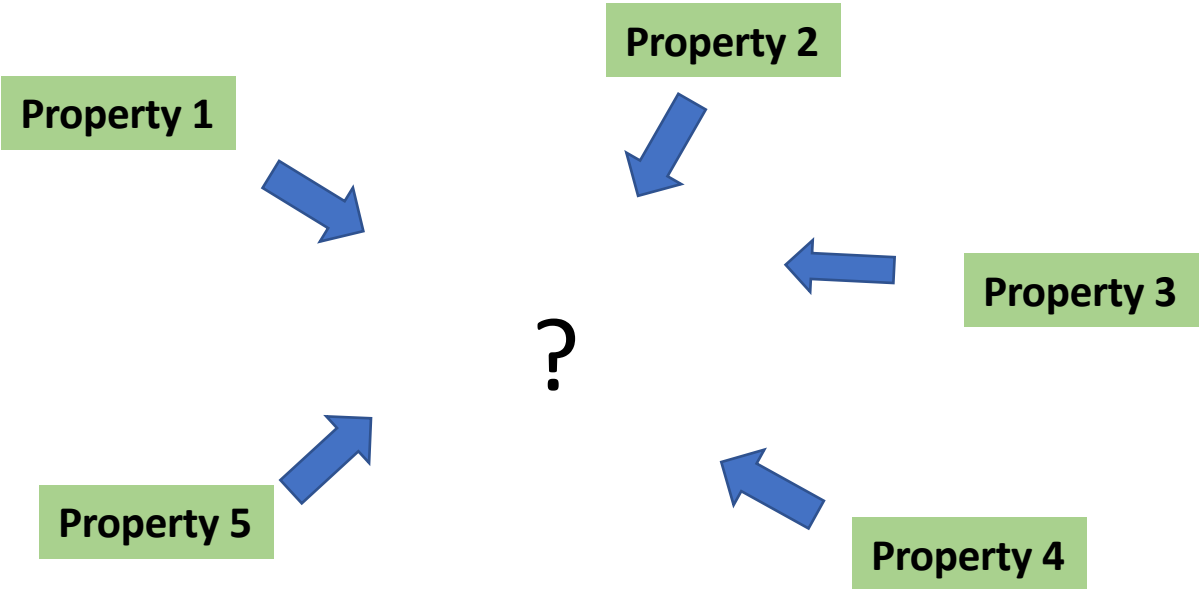
Results...

- 1
- 2
- 3
- 4
- 5

**See paper for the
75 proofs!**

Name	Fully defined	Monotonicity	Bounds	Maximum	Correction
Hamming					
Jaccard					
Dice					
Ochiai					
POG					
Kuncheva					
Lustgarten					
Wald			?		
nPOG					
Goh					
Davis					
Krízek					
Guzmán					
CW_{rel}					
<i>Lausser</i>					

So where do these properties point?



Definition 2 (Effective Stability for Pairwise Feature Redundancy).
 Given a matrix \mathbf{C} specifying feature relationships, the effective stability is

$$\hat{\Phi}_{\mathbf{C}}(\mathcal{Z}) = 1 - \frac{\sum_{f=1}^d \sum_{f'=1}^d c_{f,f'} \widehat{\text{cov}}(Z_f, Z_{f'})}{\sum_{f=1}^d \sum_{f'=1}^d c_{f,f'} \text{cov}(Z_f, Z_{f'} | H_0)} = 1 - \frac{\text{tr}(\mathbf{C}\mathbf{S})}{\text{tr}(\mathbf{C}\mathbf{\Sigma}^0)}, \quad (5)$$

where \mathbf{S} is an unbiased estimator of the variance-covariance matrix of \mathcal{Z} , i.e. $\mathbf{S}_{f,f'} = \widehat{\text{Cov}}(Z_f, Z_{f'}) = \frac{M}{M-1} (\hat{p}_{f,f'} - \hat{p}_f \hat{p}_{f'})$, $\forall f, f' \in \{1 \dots d\}$, while $\mathbf{\Sigma}^0$ is the covariance matrix of \mathcal{Z} under the null model of feature selection

1. All 5 **desirable properties**, as discussed.
2. Clean **statistical** interpretation
 ...Confidence intervals and hypothesis tests come for free
3. Computable in **closed form**, as opposed to quadratic

On the Stability of Feature Selection Algorithms

Sarah Nogueira
Konstantinos Sechidis
Gavin Brown
School of Computer Science
University of Manchester
Manchester M13 9PL, UK

SARAH.NOGUEIRA@MANCHESTER.AC.UK
KONSTANTINOS.SECHIDIS@MANCHESTER.AC.UK
GAVIN.BROWN@MANCHESTER.AC.UK

Editor: Isabelle Guyon

Abstract

Feature Selection is central to modern data science, from exploratory data analysis to predictive model-building. The “stability” of a feature selection algorithm refers to the *robustness* of its feature preferences, with respect to data sampling and to its stochastic nature. An algorithm is ‘unstable’ if a *small* change in data leads to *large* changes in the chosen feature subset. Whilst the idea is simple, *quantifying* this has proven more challenging—we note numerous proposals in the literature, each with different motivation

Case Study: Non-Small Cell Lung Cancer

Efficacy of **gefitinib** vs **chemotherapy** for lung cancer.

2 competing biomarker sets. Which do we trust?

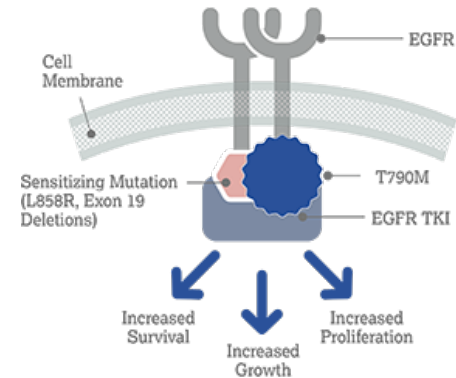
Rank	GBM	CMIM
1	EGFR expression (X_4)	EGFR mutation (X_2)
2	Disease stage (X_{10})	Serum ALP (X_{13})
3	WHO perform. status (X_1)	Blood leukocytes (X_{21})
4	Serum ALT (X_{12})	Serum ALT (X_{12})

Case Study: Non-Small Cell Lung Cancer

	GBM		CMIM
Stability $\hat{\Phi}(\mathcal{Z})$	0.87	>	0.68
- within Group A	0.96		0.45
- within Group B	0.82		0.80
- within Group C	0.14		0.43
Effective stability $\hat{\Phi}_C(\mathcal{Z})$	0.87	<	0.91

All **EGFR** gene mutations

(known to play a role in NSCLC)



Measure within-group stability

to see what's happening...

Changes our view

of the "best" algorithm to invest in.

On the **Reproducibility** of Data Science Pipelines

The **Take-Home** Message

Reproducibility is not a yes /no question.

Reproducibility = Trust

The industry needs methods to **quantify reproducibility.**



On the Stability and **Reproducibility** of Data Science Pipelines

Professor Gavin Brown
University of Manchester, UK

