

1 LETTER TO THE EDITOR

2 **Principal component analysis-based latent-space**
3 **dimensionality under-estimation, with uncorrelated latent**
4 **variables**

5 Thomas M. H. Hope,^{1,2} Ajay Halai,³ Jenny Crinion,⁴ Paola Castelli,² Cathy J. Price¹ and
6 Howard Bowman⁵

7 **Author affiliations:**

8 1 Wellcome Centre for Human Neuroimaging, Department of Imaging Neuroscience, Institute
9 of Neurology, University College London, London, WC1N 3AR, UK

10 2 Department of Psychology, John Cabot University, 00165, Rome, Italy

11 3 MRC Cognition and Brain Sciences Unit, University of Cambridge, Cambridge CB2 7EF,
12 UK

13 4 Institute of Cognitive Science, Department of Experimental Psychology, University College
14 London, London, WC1N 3AR, UK

15 5 School of Psychology, University of Birmingham, Edgbaston, Birmingham B15 2TT, UK

16
17 Correspondence to: Thomas Hope

18 Wellcome Centre for Human Neuroimaging, Department of Imaging Neuroscience, Institute
19 of Neurology, University College London, 12 Queen Square, London, WC1N 3AR, UK

20 E-mail: t.hope@ucl.ac.uk

21
22 In many scientific disciplines, features of interest cannot be observed directly, so must instead
23 be inferred from observed behaviour. In the study of the damaged brain, those ‘features of
24 interest’ might be the function or disruption of dissociable cognitive sub-systems, and the
25 ‘observed behaviour’ might be accuracies and / or reaction times recorded in standardised,
26 behavioural tasks. This inverse inference from observed data to features of interest is
27 increasingly approached using latent variable analyses^{1–4}. One of the simplest and most popular

1 of these methods, is Principal Components Analysis (PCA). During the last decade, stroke
2 outcomes research, using PCA, has yielded a surprising result: latent spaces appear lower-
3 dimensional than expected. These analyses typically find no more than 5 latent variables, and
4 sometimes just 1¹⁻⁴, even when applied to scores from wide-ranging batteries of tasks, which
5 could potentially capture impairments to many more dissociable sensory, motor and cognitive
6 sub-systems.

7 Recently, this apparent ‘dimensionality under-estimation problem’ has been explained as
8 potentially arising from spatial correlations in natural stroke lesion distributions⁵. The authors
9 used simulated data derived from real stroke induced lesions, in which impairment severity
10 scores were assigned based on the extent of damage to non-overlapping brain regions. Since
11 the regions were independent, the impairments should have been independent: i.e., the latent-
12 space dimensionality should have been the same as the number of simulated scores. But
13 instead, the authors observed that PCA typically found lower-dimensional latent spaces,
14 because natural stroke-induced lesions tend to damage neighbouring (non-overlapping) brain
15 regions together, causing the impairments to be correlated in practice even though they need
16 not have been correlated in theory⁵. The implication is that PCA-based analyses of stroke
17 outcomes data might tell us as much about lesion distributions as they ever can about the
18 fundamental organisation of cognition.

19 Here, we show that dimensionality under-estimation can occur entirely regardless of lesion
20 distributions – even when post-stroke impairments are independent by construction. We show
21 that this effect is partly a function of task impurity, the extent to which behavioural performance
22 in individual tasks is thought to emerge from the interaction of many different cognitive skills.
23 And we show that dimensionality under-estimation can be ameliorated by employing more
24 multivariate behavioural data (i.e., more tasks).

25

26 **Materials and methods**

27 We use PCA to analyse synthetic, multivariate behavioural data, which are linear mixtures of
28 known latent variable values. No lesion data were included. Following the approach employed
29 by Sperber and colleagues, we count the components derived by PCA as those whose
30 eigenvalues surpass a threshold and employ two different thresholds: the Kaiser criterion
31 (threshold eigenvalue = 1)⁶, and the Jolliffe criterion (threshold eigenvalue = 0.7)⁷. The more
32 conservative Kaiser criterion is more popular, in our experience, but the more permissive

1 Jolliffe criterion might be more appropriate when we expect to observe dimensionality under-
2 estimation. Both latent variable values and latent-to-behaviour weights are defined as random
3 uniform numbers in the range 0-1 (e.g., imagining both numbers to represent percentages of
4 maximum function / influence). We did consider other types of random distribution, but none
5 made any substantive difference to our results. And following the prior report, we employ a
6 sample size of 300. Our simulations vary the number of latent variables in the range 1-22, and
7 the number of behavioural scores in the range 22-100. We ran 1,000 simulations per parameter
8 configuration, randomly re-specifying latent variable values and latent-to-behavioural weights
9 each time, and report summary results.

11 Results

12 Analysis 1: Under-estimation with uncorrelated impairments

13 Figure 1A illustrates how the derived dimensionality of the system varies with its real
14 dimensionality, for a fixed number of (22) behavioural scores. Estimated dimensionality is
15 mostly accurate for systems with just 1 or 2 latent variables, but then grows less quickly than
16 real dimensionality – and indeed begins to fall again for higher-dimensional systems. Naturally,
17 the effect is more pronounced when using the more conservative, Kaiser criterion to count
18 derived components. Figure 1B illustrates dimensionality estimation in identical circumstances
19 to those considered in Figure 1A, with one exception: 95% of the latent-to-behavioural weights
20 are set to 10^{-6} . This change effectively ensures that most behavioural variables are mediated by
21 fewer latent variables than before: i.e., task impurity is reduced. In this case, the relationship
22 between derived and real latent system dimensionality is more intuitive, in that both grow
23 together. However, derived dimensionality still only grows about half as quickly as real
24 dimensionality, so dimensionality under-estimation still occurs.

26 Analysis 2: Under-estimation is avoided if we use many more 27 behavioural tasks

28 Figure 2 illustrates how dimensionality estimation changes as the number of behavioural tasks
29 grows. With ~6 times as many behavioural tasks as real latent variables, PCA can enumerate
30 the real latent system accurately, when using the Jolliffe criterion (~10 times for the Kaiser

1 criterion). The under-estimation problem returns as the number of latent variables increases
2 further.

3

4 **Discussion**

5 Our results suggest that dimensionality under-estimation might occur, in stroke research and
6 beyond, as a simple artefact of task impurity – or more generally, the notion that important
7 relationships between latent and observed variables might be many-to-many – even when the
8 latent variables themselves are uncorrelated. This problem appears to worsen as real latent
9 system dimensionality increases. In data derived from three or more latent variables, PCA was
10 an unreliable way to enumerate those latent variables unless: (a) task impurity was low (i.e.,
11 latent-to-behavioural weight matrices were sparse); or (b) there were many more behavioural
12 tasks than real latent variables in the system.

13 Quite how these results apply in practice, is hard to judge. First, if the latent system is
14 ‘cognition’, then we should probably allow that it might be higher-dimensional than any latent
15 system considered here. But at the same time, the *effective* dimensionality of post-stroke
16 impairments, as represented in any given sample of stroke patients, might be much lower than
17 this theoretical maximum. For example, if all of the patients in a sample have the same,
18 selective impairment, then the sample is one-dimensional. Moreover, the absolute range of
19 impairment severity might appear smaller in stroke patient samples than in the wider patient
20 population, either because standardised measures of that severity lack sensitivity, or because
21 the most severely impaired patients might struggle to travel to study sites, or tolerate the testing
22 process itself⁸. These factors make lower-dimensional estimates plausible. But since we might
23 observe lower-dimensional estimates even when they are wrong (Analysis 1), our only rational
24 choice is to treat the results of all analyses like this with caution. And further caution is called
25 for because task impurity is likely to be the norm rather than the exception in these studies⁹.

26 Our results also point to a practical solution for reducing dimensionality under-estimation:
27 using many more tasks than there are latent variables to find (Analysis 2). In practice, the
28 required ratio of tasks to variables might well vary from that observed here. From the
29 perspective of latent-space dimensionality estimation, real latent-to-behavioural weights might
30 be less efficient than those used here, so that more tasks are needed to count latent dimensions
31 accurately. On the other hand, latent-to-behaviour weights might be more informative than
32 those considered here, because batteries of tasks used in stroke research are often expressly

1 designed to vary the engagement of cognitive functions systematically and informatively. Since
2 our simulations only very rarely over-estimated real system dimensionality, one might navigate
3 this issue by adding tasks incrementally to PCA, until further additions no longer yield more
4 principal components. But of course, this approach is the opposite of what many researchers
5 might prefer, given the costs and effort required to employ extra tasks⁸.

6 Notably, our results are also at least potentially consistent with Sperber and colleagues' account
7 of latent space dimensionality under-estimation, at least in stroke outcomes research. They
8 highlighted spatial correlations in natural stroke-induced lesion distributions. This factor might
9 operate in tandem with the task impurity on which our analyses are focused. We hope that our
10 results will encourage caution in the interpretation of results, derived from post-stroke
11 impairment severity score batteries, via PCA.

12

13 **Data availability**

14 No empirical data were used in the reported analyses. MATLAB scripts used to run the
15 analyses can be downloaded from: https://github.com/tmhopegit/pca_dim_under-estimation

16

17 **Competing interests**

18 The authors report no competing interests.

19

20 **Funding**

21 This work was supported by the Wellcome Trust (205103/Z/16/Z and 203147/Z/16/Z to
22 C.J.P.); and the Medical Research Council (MR/V031481/1 to A.D.H.).

23

24 **References**

25 1. Halai AD, Woollams AM, Lambon Ralph MA. Using principal component analysis to
26 capture individual differences within a unified neuropsychological model of chronic post-
27 stroke aphasia: Revealing the unique neural correlates of speech fluency, phonology and

- 1 semantics. *Single Brain Model Sufficient*. 2017;86:275-289. doi:10.1016/j.cortex.2016.04.016
- 2 2. Ramsey LE, Siegel JS, Lang CE, Strube M, Shulman GL, Corbetta M. Behavioural clusters
3 and predictors of performance during recovery from stroke. *Nat Hum Behav*. 2017;1(3):0038.
4 doi:10.1038/s41562-016-0038
- 5 3. Schumacher R, Halai AD, Lambon Ralph MA. Assessing and mapping language, attention
6 and executive multidimensional deficits in stroke aphasia. *Brain J Neurol*. 2019;142(10):3202-
7 3216. doi:10.1093/brain/awz258
- 8 4. Akkad H, Hope TMH, Howland C, et al. Mapping spoken language and cognitive deficits
9 in post-stroke aphasia. *NeuroImage Clin*. 2023;39:103452. doi:10.1016/j.nicl.2023.103452
- 10 5. Sperber C, Gallucci L, Umarova R. The low dimensionality of post-stroke cognitive deficits:
11 it's the lesion anatomy! *Brain*. 2023;146(6):2443-2452. doi:10.1093/brain/awac443
- 12 6. Kaiser HF. The application of electronic computers to factor analysis. *Educ Psychol Meas*.
13 1960;20(1):141-151.
- 14 7. Jolliffe IT. Discarding variables in a principal component analysis. I: Artificial data. *J R Stat*
15 *Soc Ser C Appl Stat*. 1972;21(2):160-173.
- 16 8. Halai AD, De Dios Perez B, Stefaniak JD, Lambon Ralph MA. Efficient and effective
17 assessment of deficits and their neural bases in stroke aphasia. *Cortex*. 2022;155:333-346.
18 doi:10.1016/j.cortex.2022.07.014
- 19 9. Burgess PW. Theory and methodology in executive function research. In: *Methodology of*
20 *Frontal and Executive Function*. Routledge; 2004:87-121.

21
22

1 **Figure legends**

2 **Figure 1 Dimensionality under-estimation.** Both panels illustrate the mean and standard
3 deviation of derived dimensionality when analysing 22 behavioural scores, derived from latent
4 systems including 1-22 latent variables. Accurate dimensionality estimation occurs when the
5 lines intersect the diagonal of either panel, where estimated dimensionality equals real system
6 dimensionality. In panel A, latent-to-behaviour weights are random uniform numbers in the
7 range 0-1. In this case, dimensionality estimation is mostly accurate for systems with 1-3 latent
8 variables, but then becomes less accurate as real dimensionality increases. In panel B, 95% of
9 the weights are set to 10^{-6} , reducing task impurity, so that estimated dimensionality grows as
10 real dimensionality grows (albeit that the former only grows about half as quickly as the latter).

11

12 **Figure 2 Latent space dimensionality estimation as the numbers of tasks increases from**
13 **40 to 100.** As in Figure 1, derived dimensionality equals the real dimensionality of the system
14 along the dotted diagonal line. The number of systems that can be accurately estimated
15 increases with the number of task scores. Once the real dimensionality increases beyond $\sim 1/6$
16 of the number of task scores ($\sim 1/10$ for Kaiser criterion), derived dimensionality again under-
17 estimates the true the dimensionality of the system: i.e., the estimation problem returns. Panels
18 A and B illustrate the effect when using the Joliffe and Kaiser criteria, respectively: for a given
19 real system dimensionality, more tasks are required for accurate dimensionality estimation
20 when using the more conservative (Kaiser) criterion to count those estimated dimensions.

21

22

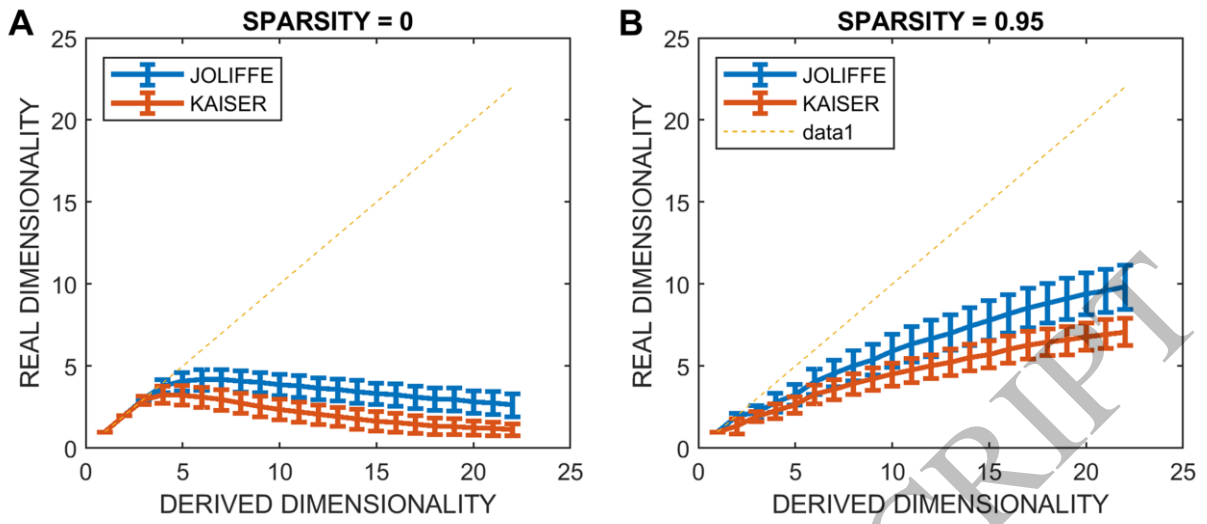
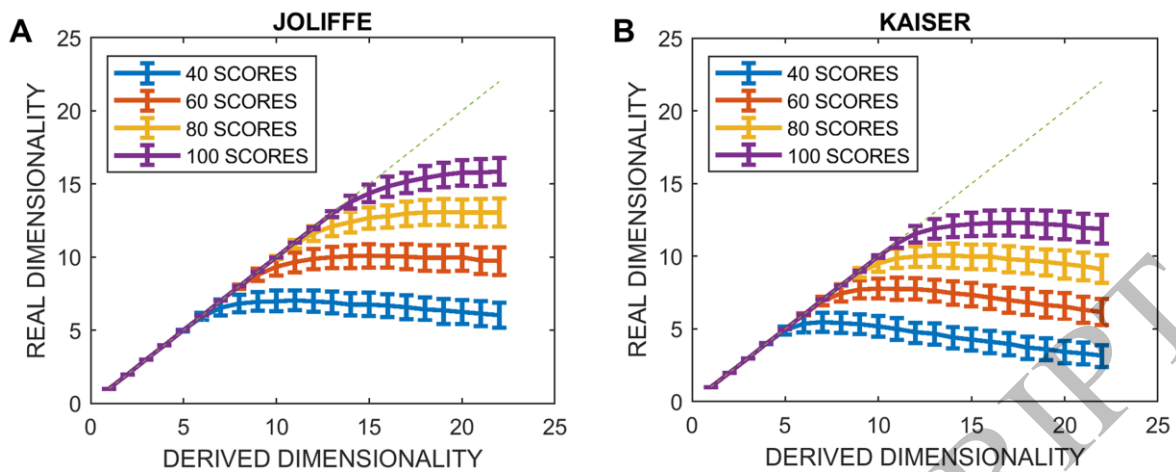


Figure 1
159x69 mm (x DPI)

1
2
3
4

ACCEPTED MANUSCRIPT



1
2
3

Figure 2
159x62 mm (x DPI)

ACCEPTED MANUSCRIPT