

Estimating the functional dimensionality of neural representations

Christiane Ahlheim^{a,*}, Bradley C. Love^{a,b},

^a*Department of Experimental Psychology, University College London, 26 Bedford Way,
London WC1H 0AP, United Kingdom*

^b*The Alan Turing Institute, United Kingdom*

Abstract

Recent advances in multivariate fMRI analysis stress the importance of information inherent to voxel patterns. Key to interpreting these patterns is estimating the underlying dimensionality of neural representations. Dimensions may correspond to psychological dimensions, such as length and orientation, or involve other coding schemes. Unfortunately, the noise structure of fMRI data inflates dimensionality estimates and thus makes it difficult to assess the true underlying dimensionality of a pattern. To address this challenge, we developed a novel approach to identify brain regions that carry reliable task-modulated signal and to derive an estimate of the signal's functional dimensionality. We combined singular value decomposition with cross-validation to find the best low-dimensional projection of a pattern of voxel-responses at a single-subject level. Goodness of the low-dimensional reconstruction is measured as Pearson correlation with a test set, which allows to test for significance of the low-dimensional reconstruction across participants. Using hierarchical Bayesian modeling, we derive the best estimate and associated uncertainty of underlying dimensionality across participants. We validated our method on simulated data of varying underlying dimensionality, showing that recovered dimensionalities match closely true dimensionalities. We then applied our method to three published fMRI data sets all involving processing of visual stimuli. The results highlight three possible applications of estimating the functional dimensionality of neural data. Firstly, it can aid

*Corresponding author

Email addresses: c.ahlheim@ucl.ac.uk (Christiane Ahlheim), b.love@ucl.ac.uk (Bradley C. Love)

evaluation of model-based analyses by revealing which areas express reliable, task-modulated signal that could be missed by specific models. Secondly, it can reveal functional differences across brain regions. Thirdly, knowing the functional dimensionality allows assessing task-related differences in the complexity of neural patterns.

Keywords: neural representations, dimensionality reduction, multivariate analysis

1 **1. Introduction**

2 A growing number of fMRI studies are investigating the representational
3 geometry of voxel response patterns. For example, using representational
4 similarity analysis (RSA; Kriegeskorte and Kievit, 2013), researchers have
5 characterized visual object representations along the ventral stream (Khaligh-
6 Razavi and Kriegeskorte, 2014) and how these representations vary across
7 tasks (Bracci et al., 2017).

8 Interpreting representational geometry in neural responses can be diffi-
9 cult. For example, RSA tests for a hypothesized representational pattern,
10 but an important and more fundamental question should be addressed first,
11 namely whether there is any dimensionality to the underlying neural pattern
12 and, if so, what that dimensionality is.

13 Knowing whether a pattern has dimensionality should be prerequisite for
14 RSA and other multivariate representational analyses because a particular
15 similarity structure can only be found when there is sufficient dimensionality
16 to represent the proposed relations. For example, searching for a flavor
17 space with dimensions sweet, sour, bitter, salty and umami would be a fool’s
18 errand in brain areas that contain little or no dimensionality. Furthermore,
19 independent of the particular geometry, the dimensionality of a neural pat-
20 tern is informative of how many features of a task are represented in a brain
21 region, which can inform our understanding of an area’s function.

22 There are many methods of dimensionality reduction and estimation,
23 most of which involve low-rank matrix approximation and aim to maximize
24 the correspondence between the original and the approximated matrix. For
25 example, two common approaches to estimate the dimensionality of an ob-
26 served neural or behavioral pattern are principal component analysis (PCA)
27 or relatedly, multidimensional scaling (MDS).

28 PCA, or the closely related factor analysis and singular value decompo-
29 sition (SVD) (Hastie et al., 2009), is widely used in the study of individual
30 differences and aids estimating how many latent components, or “factors”,
31 underlie a pattern of (item) responses across participants, as for instance
32 in the context of intelligence (Spearman, 1904) or personality tests (Cattell,
33 1947). In the context of neuroimaging, PCA has been used to identify brain
34 networks (Huth et al., 2012; Friston et al., 1993). PCA derives how much
35 variance of the observed pattern is explained by each underlying component.

36 Similarly, MDS finds the best representation of original distances in a
37 low-dimensional space (Kriegeskorte and Kievit, 2013). For example, two

38 stimuli like a chair and table that are very close to each other in the high-
39 dimensional space will be represented closely in the low-dimensional projec-
40 tion achieved by MDS, whereas two stimuli that were very distant from each
41 other, for instance a chair and a bunny, will be projected far apart. MDS has
42 been successfully applied to behavioral as well as neural data to reveal which
43 stimulus features underly observed representational geometries (Bracci and
44 Op de Beeck, 2016; Kriegeskorte and Kievit, 2013; Kriegeskorte et al., 2008),
45 though it has been questioned to which extent results from MDS are inter-
46 pretable (Goddard et al., 2017). For reasons outlined below, we will focus on
47 SVD to estimate the dimensionality of neural representations, though other
48 methods could be paired with our general approach, including nonlinear ap-
49 proaches such as Nonlinear PCA (Kramer, 1991).

50 Estimating the dimensionality of neural data brings its own unique chal-
51 lenges. In a noise-free scenario, dimensionality can be defined as the number
52 of linear orthogonal components (singular- or eigenvalues) underlying a ma-
53 trix that are larger than zero (Shlens, 2014), indicating that the component
54 fits some variance in the data. Unfortunately, actual recordings of neural ac-
55 tivity always contain noise, which inflates non-signal components above zero
56 (Fusi et al., 2016; Diedrichsen et al., 2013). This noise makes it challenging
57 to determine which areas contain signal and, if so, what the dimensionality
58 of the signal is.

59 One criterion, which we adopt in the work reported here, is to choose
60 the number of components that should maximize reconstruction accuracy
61 (measured by correlation) on new data (i.e., test data). While even for
62 data with low or moderate true dimensionality more components will always
63 increase fit for existing data (i.e., training data), performance on test data
64 (i.e., generalization, prediction) will usually be best for a moderate number of
65 components because these components largely reflect true signal as opposed
66 to noise in the observed training sample.

67 The problem of distinguishing between signal and noise in a neural pat-
68 tern is related to the bias-variance trade-off in supervised learning and model-
69 selection. Overly simple models (few components) are highly biased, fitting
70 training data poorly and not performing well on test data. These overly
71 simple models cannot pick-up on nuances in the signal. Conversely, overly
72 complex models (many components) are too sensitive to the variance in the
73 training data (i.e., overfit). Although they fit the training data very well,
74 overly complex models treat noise in the training data as signal and, there-
75 fore, generalize poorly. Thus, the sweet spot for test performance should be

76 at some moderate number of components that largely reflect true signal (see
77 Figure 1 A). Thus, identifying the true number of underlying components is
78 analogous to deciding which model best explains the data.

79 One naive way to navigate this trade-off between simple and complex
80 models is to use some arbitrary cutoff, such as including the number of com-
81 ponents that captures some amount of variance in the training data or decid-
82 ing based on visual inspection which components may carry signal (known
83 as scree plot, Cattell, 1966). In the case of fMRI, where the signal-to-noise
84 ratio depends on multiple factors like scanner settings, experimental design,
85 and physiological activity (Huettel et al., 2003), estimating the underlying
86 dimensionality based on an arbitrary cut-off criterion for explained variance
87 could be misleading. Likewise, although identifying relevant components
88 via visual inspection works for small datasets, it is not applicable to large
89 datasets as fMRI data, as it would require a manual decision for each voxel.
90 Furthermore, the size of fMRI datasets (usually thousands of voxels) calls
91 for a computationally efficient and automated approach, making estimating
92 the dimensionality for the whole brain feasible. Thus, for neuroimaging data,
93 there is a need for an efficient, systematic and objective approach that can
94 both identify areas with statistical significant dimensionality and provide a
95 useful estimate of the underlying dimensionality.

96 Previous efforts to estimate the dimensionality of neural response pat-
97 terns have applied linear classifiers to neural data to evaluate dimensionality
98 (Rigotti et al., 2013; Diedrichsen et al., 2013). Rigotti et al. (2013) were able
99 to show that dimensionality of single-cell recordings in monkey PFC is linked
100 to successful task-performance, indicating that dimensionality of neural pat-
101 terns is task-sensitive. In line with this, Diedrichsen et al. (2013) showed
102 that the dimensionality of motor cortex representations differs depending on
103 the task. Using a combination of PCA and linear Gaussian classifiers, the
104 authors showed that motor cortex representations of different force levels
105 are low dimensional, whereas usage of different fingers was associated with
106 multidimensional neural patterns (Diedrichsen et al., 2013). Notably, both
107 studies focused on estimating task-related changes in dimensionality in a pre-
108 scribed brain region, rather than estimating which areas across the brain had
109 significant dimensionality.

110 In the present work, we expand on previous contributions by evaluat-
111 ing a novel approach that, in a robust and computationally efficient manner,
112 tests for which areas display statistically significant dimensionality, estimates
113 the dimensionality, and provides an indication of the certainty of the esti-

114 mate. We combine singular value decomposition (SVD) and cross-validation
115 to identify areas across the brain with underlying dimensionality. We derive
116 which of all possible low-dimensional reconstructions of the fMRI signal is the
117 best dimensionality estimate of a held-out test run, and quantify the good-
118 ness of the low-dimensional reconstruction via Pearson correlation. Using a
119 cross-validation procedure to identify the best dimensionality estimate boosts
120 that only components that carry signal and thus generalize to new data are
121 kept. By assessing the significance of the correlation, we can distinguish be-
122 tween areas that show reliable signal with underlying dimensionality vs. areas
123 that do not show a reliable task-modulation. After establishing significant
124 functional dimensionality, we use Bayesian modeling to derive a population
125 estimate and associated uncertainty of the degree of dimensionality. We will
126 refer to this task-dependent dimensionality as functional dimensionality.

127 Through simulations and evaluation of three (published) fMRI datasets,
128 we find that our method successfully identifies areas with significant func-
129 tional dimensionality and provides reasonable estimates of the underlying
130 dimensionality. In the first fMRI dataset, participants performed a catego-
131 rization task which required differential attention to various stimulus features
132 (Mack et al., 2013). The second study investigated shape- and category spe-
133 cific neural responses to the presentation of natural images (Bracci and Op de
134 Beek, 2016). The third study involved categorization tasks that varied sys-
135 tematically in their attentional demands (Mack et al., 2016), which we predict
136 should affect functional dimensionality.

137 Across all three studies, we were able to identify areas carrying functional
138 dimensionality in a manner that supported and extended the original find-
139 ings. Focusing on wholebrain effects in the the first two studies, we identified
140 a consistent network of areas showing functional dimensionality during vi-
141 sual stimulus processing. This network encompassed areas that were reported
142 by the original authors as being task-relevant, identified through represen-
143 tational similarity analysis and cognitive model fitting (Bracci and Op de
144 Beek, 2016; Mack et al., 2013). Furthermore, functional dimensionality was
145 revealed in additional areas, highlighting the sensitivity of our method and
146 suggesting that reliable task-modulated signal was present that was not ex-
147 plained by the models the original authors tested. In the last study, we
148 combined a region-of-interest approach and multilevel Bayesian modeling to
149 show that dimensionality varied depending on task-requirements, which fol-
150 lows from the original authors' claims but remained untested until now (Mack
151 et al., 2016). We outline how the notion and identification of functional di-

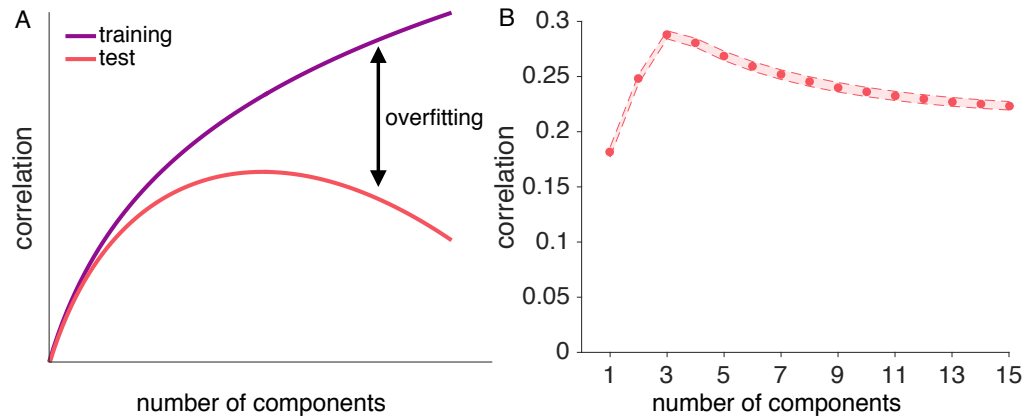


Figure 1: Illustration of the concept of overfitting and generalizability. A: As more components are added to a low-dimensional reconstruction, the correlation between the training data and the reconstruction approaches the maximum of 1 for a full-dimensional reconstruction (purple curve). Adding components is equivalent to adding model parameters to improve fit, which reduces the model’s bias and increases its variance. For the correlation between the reconstructed training and independent test data (red curve), adding components initially improves performance but at some point reduces performance due to overfit (see Parpart et al., 2017, for a related illustration). B: Reconstruction correlations achieved by all possible low-dimensional reconstructions for a simulated ground-truth dimensionality of 4. Reconstruction correlations rise as more components are added up to the point where the true dimensionality is reached, and decrease afterwards. Results are averaged across 6 runs and 1000 simulated voxel patterns with varying signal-to-noise ratios.

152 dimensionality can aid the analysis and understanding of neuroimaging data
153 in various ways.

154 2. General Methods

155 Neuroimaging data, such as fMRI, M/EEG, or single-cell recordings, can
156 be represented as a matrix of n voxels, neurons, or sensors \times m conditions.
157 For example, BOLD response patterns in the fusiform face area (FFA) to
158 3 different stimulus conditions can be expressed as a matrix Y of the size
159 n (number of voxels) \times 3 (face, house, or tool stimulus condition). The
160 maximum possible dimensionality is the minimum of n and m , which in this
161 example would be 3, assuming many voxels in FFA were included in the

162 analysis. However, functional dimensionality could be lower. For example,
163 dimensionality would be lower if the region only responded to face stimuli
164 and showed the same lower response to house and tool stimuli.

165 Various methods exist that allow to estimate a matrix’s dimensionality
166 and a review of all of them is beyond the scope of this paper. The approach we
167 present here is modular and estimates a matrix’s dimensionality by combining
168 low-rank approximation with cross-validation and significance testing. This
169 modularity allows to flexibly choose the dimensionality reduction technique
170 which best fits with ones requirements. Here, we used SVD (which is often
171 used to compute PCA solutions) because it is a well-understood, easy to
172 implement, and a computationally efficient low-rank matrix approximation.

173 The choice of SVD, as well as how the data matrix is normalized is in-
174 formed by our understanding of the underlying neural signal. Because voxels
175 differ greatly from one another in their overall activity level and activity lev-
176 els can drift over runs, we demean each row (i.e., voxel) of the data matrix by
177 run. In contrast, we do not demean each column, as would typically be done
178 with approaches that focus on the covariance of the column vectors (e.g.,
179 PCA). The reason we do not normalize by column (i.e., condition) is that we
180 are open to the possibility that different stimuli may be partially coded by
181 overall activity levels of a population of voxels. For example, imagine a brain
182 area only responds strongly to faces, but not to other stimuli. An SVD with
183 demeaned voxels (i.e., rows) would be sensitive to this dimension of represen-
184 tation, whereas a procedure that effectively worked with demeaned columns
185 would not be sensitive to this task-driven difference in neural activity (see
186 Davis et al., 2014; Hebart and Baker, 2017, for a related discussion).

187 In the following section, we describe how a combination of SVD and
188 cross-validation can be used to test whether an observed neural pattern can
189 be successfully reconstructed using a low-rank approximation, assessed as a
190 significant Pearson correlation between a low-rank approximation and a held
191 out test set, and how this technique provides an estimate of the pattern’s un-
192 derlying dimensionality (see Figure 2 for an overview of all steps). As all our
193 examples are fMRI data sets, we will describe the steps using fMRI termi-
194 nology, though the procedure could be applied to any type of neuroimaging
195 data. We provide the code and data to replicate the analyses presented here
196 and for use on other datasets at osf.io/tpq92.

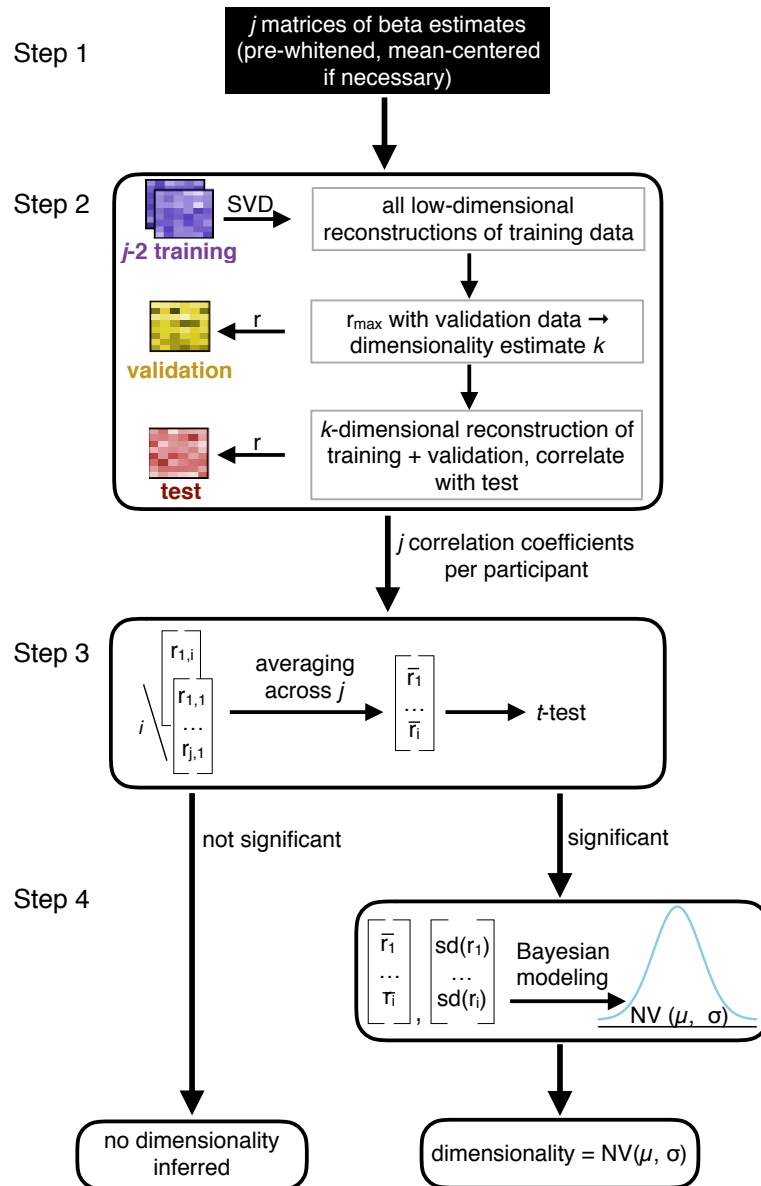


Figure 2 (*previous page*): Step 1: Prior to dimensionality estimation, raw data are pre-processed with preferred settings and software and beta estimates derived from a GLM are obtained for each condition of interest. The resulting j matrices of size n (number of voxels) $\times m$ (number of conditions) are pre-whitened and mean-centered (by row, i.e., voxel) to remove baseline differences across runs. Step 2: a combination of cross-validation and SVD is implemented to find the best dimensionality estimate k for each run j . Pearson correlations between all possible low-dimensionality reconstructions of the data and a held-out test set quantify the goodness of each reconstruction for each run j (see Figure 3 for details). Step 3: the resulting j correlations are averaged for each participant and tested for significance, for instance using t-tests, across all participants. Step 4: If the reconstruction correlations are significant across participants, a hierarchical Bayesian model can be used to derive the best estimate of the degree of functional dimensionality (see Figure 4 for details). For each participant, the average estimated dimensionality and standard deviation of this estimate is calculated and a population estimate and respective standard deviation (uncertainty in the estimate) is derived across all participants.

197 *2.1. Step 1: Data pre-processing*

198 We developed the presented method with application to fMRI data in
199 mind, though it can be easily adapted to fit requirements of single cell record-
200 ings or M/EEG data. The method takes beta estimates resulting from a
201 GLM fit to the observed BOLD response as input. In all studies presented
202 here, standard pre-processing steps were performed using SPM 12 (Wellcome
203 Department of Cognitive Neurology, London, United Kingdom), but the pre-
204 cise nature of the preprocessing and implemented GLM is not critical to our
205 method. Functional data were motion corrected, co-registered and spatially
206 normalized to the Montreal Neurological Institute (MNI) space.

207 To reduce the impact of the structured noise, which is correlated across
208 voxels, on the dimensionality estimation and to improve the reliability of
209 multivariate voxel response patterns (Walther et al., 2016), we applied mul-
210 tivariate noise-normalization, that is, spatial pre-whitening, before estimat-
211 ing the functional dimensionality. We used the residual time-series from the
212 fitted GLM to estimate the noise covariance Σ_{noise} and used regularization to
213 shrink it towards the diagonal (Ledoit and Wolf, 2004). Each $n \times m$ matrix
214 of beta estimates Y was then multiplied by $\Sigma_{noise}^{-\frac{1}{2}}$ (Walther et al., 2016).

215 In fMRI data, the baseline activation can differ across functional runs.

216 This has important implications for our approach presented here, as it can
217 bias the correlation between neural patterns across runs. To account for this,
218 we demeaned the pre-whitened beta estimates across conditions, resulting in
219 an average estimate of zero for each voxel. This demeaning reduces the
220 possible maximum dimensionality of the data to $k_{max} = m - 1$. Notably,
221 demeaning of voxels is conceptually different from demeaning conditions,
222 which would have been implemented by PCA, as it preserves differences
223 between conditions, whereas PCA would remove those.

224 *2.2. Step 2: Evaluating all possible SVD (dimensional) models*

225 The dimensionality of a matrix is defined as its number of non-zero sin-
226 gular values, identified via singular value decomposition (SVD). SVD is the
227 factorization of an observed $n \times m$ matrix M of the form $U\Sigma V^T$. U and
228 V are matrices of size $m \times m$ and $n \times n$, respectively, and Σ is an $n \times m$
229 matrix, whose diagonal entries are referred to as the singular values of M .
230 A k -dimensional reconstruction of the matrix M can be achieved by only
231 keeping the k largest singular values in Σ and replacing all others with zero,
232 resulting in $\tilde{\Sigma}$. This is known as Eckart-Young theorem (Eckart and Young,
233 1936), leading to equation 1:

$$\tilde{M} = U\tilde{\Sigma}V^T \quad (1)$$

234 To estimate the dimensionality of fMRI data, we applied SVD to j (number
235 of runs) matrices Y of n (number of voxel) \times m (number of beta estimates),
236 with the restriction of $n > m$.

237 Critically, fMRI beta estimates are noisy estimates of the true signal.
238 In the presence of noise, all singular values of a matrix will be non-zero,
239 requiring the definition of a cut-off criterion to assess the number of singular
240 values reflecting signal. Removing noise-carrying components from a matrix
241 is beneficial, as it avoids overfitting to the noise and thus, improves the
242 generalizability of the low-dimensional reconstruction to another sample (see
243 Figure 1 A for an illustration of the concept of overfitting). We aimed to avoid
244 any subjective (arbitrary) criterion as percentage of explained variance or
245 alike (Cattell, 1966). To that end, we implemented a nested cross-validation
246 procedure at the core of our method to identify singular values that carry
247 signal (see step 1 of the general overview depicted in Figure 2 and Figure 3
248 for a detailed illustration of the cross-validation approach). This allows us
249 to overcome the inflation of dimensionality of fMRI data due to noise and
250 test which areas of the brain carry signal with functional dimensionality.

251 Data are partitioned $j \times (j - 1)$ times into training (Y_{train}), validation
252 (Y_{val}), and test (Y_{test}) data. The (demeaned and pre-whitened) $j - 2$ training
253 runs are averaged, and SVD is applied to the resulting $n \times m$ matrix \bar{Y}_{train} .
254 We then build all possible low-dimensional reconstructions of the averaged
255 training data, with dimensionality ranging from 1 to $m - 1$. Low-dimensional
256 reconstructions are generated by keeping only the k highest singular values
257 and setting all others to zero. Each low-dimensional reconstruction of matrix
258 \bar{Y}_{train} is correlated with the held-out Y_{val} . This is repeated for each possible
259 partitioning in training and validation, resulting in $j - 1 \times m - 1$ correlation
260 coefficients. Correlations are Fisher's z-transformed and averaged across the
261 $j - 1$ partitionings. The dimensionality with the average highest correlation
262 is picked as best estimate k of the underlying dimensionality. As keeping
263 components that reflect noise rather than signal lowers the correlation with
264 an independent data set, the highest correlation is not necessarily achieved
265 by keeping more components. This procedure thus avoids inflated dimen-
266 sionality estimates.

267 After identifying the best dimensionality estimate k for run j , the training
268 and validation runs from 1 to $j - 1$ are averaged together and SVD is applied
269 to the averaged data. We then generate a k -dimensional reconstruction of
270 the averaged data. The quality of this final low-dimensional reconstruction
271 is measured as Pearson correlation with Y_{test} . We chose Pearson correlation
272 and not mean-square error (MSE), which is suggested by the use of SVD as
273 measure of reconstruction quality, because MSE is influenced by the variance
274 of the reconstructed data, which depends on its dimensionality k .

275 *2.3. Step 3: Determining statistical significance*

276 The approach results in j estimates of the underlying dimensionality and
277 j corresponding test correlations per participant. Under the null-hypothesis
278 of no dimensionality, and thus, only noise present in the matrix, reconstruc-
279 tion correlations averaged across runs are distributed around zero. Thus,
280 across-participants significance of the averaged reconstruction correlations
281 can be assessed using one-sample t-tests or non-parametric alternatives, as
282 for instance permutation tests (Nichols and Holmes, 2003), and established
283 correction methods for multiple comparisons, like threshold-free cluster en-
284 hancement (TFCE, see Smith and Nichols, 2009).

285 Only if a significant, k -dimensional, reconstruction correlation can be
286 established across participants, we refer to an area as showing functional di-
287 mensionality. It should be noted that a significant reconstruction correlation

288 only indicates that the underlying functional dimensionality is one or bigger.
289 However, testing for a dimensionality of two or larger can be achieved by
290 removing not only the voxel-mean before estimating the dimensionality, but
291 also the condition mean, effectively removing univariate differences between
292 conditions.

293 *2.4. Step 4: Estimating the degree of dimensionality*

294 The previously described steps allow us to identify which areas carry
295 reliable signal with functional dimensionality, but do not provide a precise
296 estimate of the degree of the underlying dimensionality. The best population
297 estimate of a region's functional dimensionality should optimally combine
298 information across participants, giving more weight to participants with more
299 reliable estimates, and should furthermore reflect how peaked the distribution
300 of underlying population estimates is, accounting for the fact that different
301 participants could express different true dimensionality.

302 Given a significant reconstruction correlation across participant, j esti-
303 mates of the degree of dimensionality are obtained (for each voxel or ROI)
304 for each participant. In a noise-free scenario, all j estimates reflect the true
305 dimensionality and thus, direct inference could be made solely based on these
306 estimates. Under noise, these estimates could over- or underestimate the true
307 dimensionality. The less reliable the j dimensionality estimates, the higher
308 the variance across them. Mere averaging of the j estimates across par-
309 ticipants would discard this information, weighting all participants equally,
310 irrespective of their reliability. Down-weighting the influence of less reliable
311 dimensionality estimates on the population estimate leads to a better popu-
312 lation estimate (Kruschke, 2014).

313 To account for this, we implemented a multilevel Bayesian model using the
314 software package Stan (The Stan Development Team, 2017). Given the mean
315 and standard deviation of j dimensionality estimates per participant, the
316 model derives the best estimate for the true degree of dimensionality across all
317 participants. Due to the nature of the multilevel model, individual estimates
318 are subject to shrinkage towards the estimated population mean, and the
319 degree of shrinkage is more pronounced for estimates with higher variance and
320 stronger deviation from the estimated population mean (Kruschke, 2014).

321 Additionally to the estimate of the population dimensionality, the model
322 returns estimates for the population dimensionality's variance, reflecting the
323 uncertainty of the dimensionality estimate. For each individual participant,

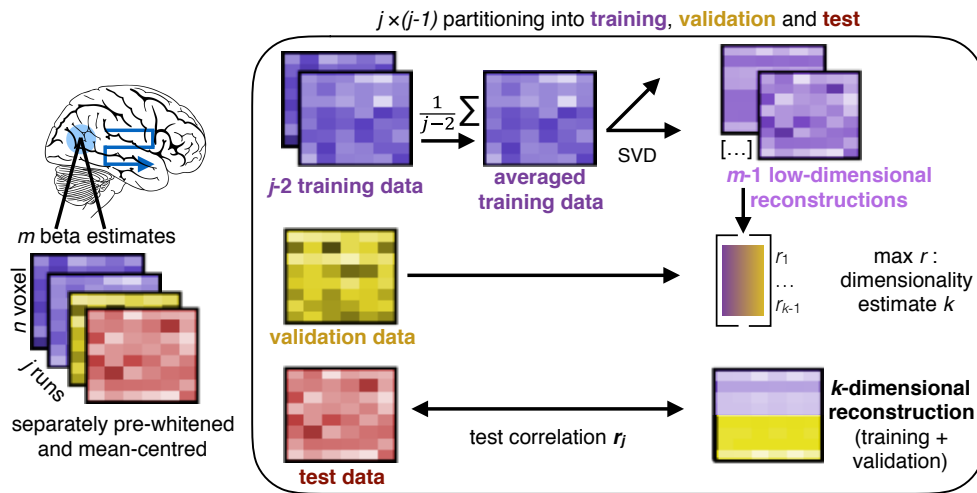


Figure 3: Illustration of the combination of SVD and cross-validation, corresponding to step 2 in Figure 2. For each searchlight or ROI, j (number of runs) n (number of voxels) \times m (number of beta estimates) matrices are used to estimate the functional dimensionality. For all possible partitions of j runs into training, validation and test data, we first average all training runs and build all possible low-dimensional reconstructions of these averaged data using SVD. All reconstructions are then correlated with the validation run, resulting in $j - 1$ correlation coefficients and respective dimensionalities. The dimensionality that results in the highest average correlation across $j - 1$ runs is picked as dimensionality estimate k for this fold and a k -dimensional reconstruction of the average of the training and validation runs is correlated with a held-out test-run, resulting in a final reconstruction correlation. In total, j reconstruction correlations are returned that can be averaged and tested for significance across participants using one-sample t-tests or alike. To derive a better estimate of the underlying dimensionality, the j dimensionality estimates per participant can be submitted to the hierarchical Bayesian model (step 4 in Figure 2)

324 the model estimates the participant's true underlying dimensionality and re-
325 turns the uncertainty of this estimate. As we did not have strong priors
326 regarding the dimensionality of the neural patterns, we implemented a uni-
327 form prior over the population dimensionality estimates, reflecting that the
328 dimensionality could be anything from 1 to $m - 1$. This can be adapted to
329 be informative for studies estimating the functional dimensionality of neural
330 patterns with stronger priors. Figure 4 shows an illustration of the model.

331 The model assumed that all individual average dimensionality estimates
332 y_i come from a truncated individual t -distribution, centered at the true indi-
333 vidual dimensionality $\hat{\mu}_i$ which comes from a common truncated normal dis-
334 tribution with mean μ and variance σ^2 , see Equation 2. We chose a truncated
335 t -distribution at the individual level to account for the fact that there is only
336 a limited number (j) of samples underlying each participant's dimensionality
337 estimate. The uniform prior distribution over the true dimensionality ranged
338 from 1 to $m - 1$.

$$\begin{aligned} y_i &\sim T(j - 1, \hat{\mu}_i, \hat{\sigma}_i), 1 \leq y_i \leq m - 1, \text{ with} \\ \hat{\mu}_i &\sim N(\mu, \sigma), 0 \leq \hat{\mu}_i \leq \sigma_{max}, \\ \hat{\sigma}_i &\sim N(\sigma_i, 1), 0 \leq \hat{\sigma}_i \leq \max(\sigma_i), \text{ and} \\ \sigma_i &\sim U(0, \max(\sigma_i)). \end{aligned} \tag{2}$$

339 The maximum population variance was defined as the expected variance
340 of this uniform distribution $\frac{1}{12}(m - 2)^2$, reflecting the prior that each partici-
341 pant could express a different, true dimensionality. On the subject-level, the
342 maximum variance was defined as

$$\max(\sigma_i^2) = \frac{j}{j - 1} * (m - 1 - \frac{m}{2})^2 \tag{3}$$

343 which corresponds to the maximum possible variance across j dimension-
344 ality estimates.

345 3. Simulations

346 Before applying our method to real fMRI data, we tested the validity of
347 our method through dimensionality-recovery studies on simulated fMRI data.
348 Estimating the dimensionality for simulated cases where the true underlying
349 dimensionality is known allowed us to assess whether our procedure results
350 in a reliable dimensionality estimate.

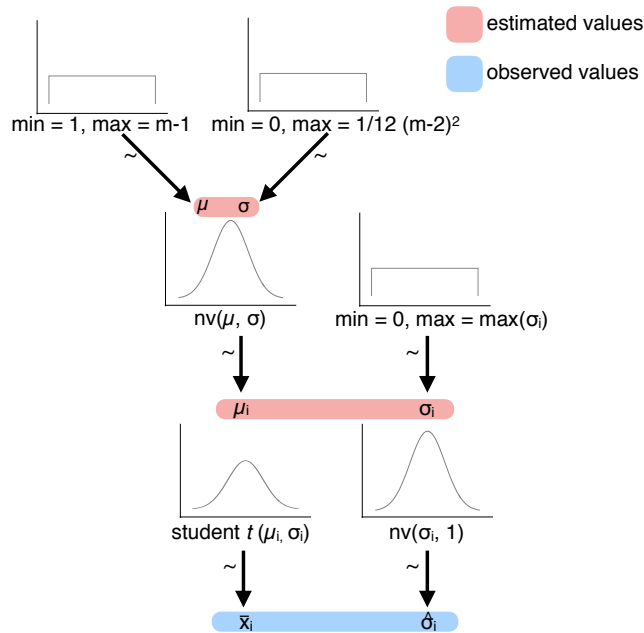


Figure 4: Illustration of the implemented multilevel model to estimate the degree of functional dimensionality, corresponding to step 4 in Figure 2. The observed averaged dimensionality estimates per participant are assumed to be sampled from an underlying subject-specific t -distribution with mean μ_i and standard deviation σ_i . The standard deviation $\hat{\sigma}_i$ of the participants' dimensionality estimates is assumed to be sampled from a normal distribution with mean σ_i and a standard deviation of 1. The subject-specific t -distributions of μ_i are assumed to come from a population distribution with a normally distributed mean μ and variance σ . Subject-specific standard deviations σ_i are assumed to come from a uniform distribution, ranging from 0 to $\max(\sigma_i)$. At the top level, a uniform prior is implemented. Mean and variance of the normal distribution of population means μ are assumed to come from a uniform distribution ranging from 1 to $m - 1$ and 0 to σ_{max} , respectively. Distributions were derived from https://github.com/rasmusab/distribution_diagrams

351 *3.1. Methods*

352 Simulated data were created using the RSA toolbox (Nili et al., 2014) and
353 custom Matlab code. Parameters of the simulation were picked in accordance
354 with the study by Mack and colleagues (2013). We simulated fMRI data of
355 presentation of 16 different stimuli, presented for 3 sec, three repetitions per
356 run, and six runs, closely matching the specifications of the original study.
357 To mimic a searchlight-approach, we defined the size of the cubic sphere
358 $4 \times 4 \times 4$ voxels, resulting in a simulated pattern of 64 voxels.

359 We simulated data with a dimensionality of 2, 4, and 6. We set the mean
360 signal to noise ratio (SNR) to match empirically observed reconstruction
361 correlation magnitudes of .25. As in the real data, reconstruction correlations
362 varied across participants, ranging from .10 to .50. Thus, participants differed
363 in their reliability of the dimensionality estimates.

364 To generate data with varying ground-truth dimensionality k , we first
365 generated true, i.e. noise-free, $n(\text{voxel}) \times m(\text{conditions})$ matrices with un-
366 derlying pre-defined dimensionality. This was achieved by applying PCA to
367 a random 16×16 matrix and building a k -dimensional reconstruction of it.
368 Rows of this matrix were added to a $n \times 16$ matrix. For each row, i.e. voxel,
369 a specific amplitude was drawn from a normal distribution and added.

370 In the next step, we calculated the dot-product of the generated beta
371 matrices and generated design matrices, which were HRF convolved. This
372 resulted in noise-free fMRI time series.

373 A noise matrix was generated by randomly sampling from a Gaussian dis-
374 tribution. The $n(\text{voxel}) \times t(\text{timesteps})$ matrix was then spatially smoothed
375 and temporally smoothed with a Gaussian kernel of 4 FWHM. Finally, this
376 temporally and spatially smoothed noise matrix was added to the noise-free
377 time-series and the design matrix was fit the the resulting data using a GLM.
378 This resulted in a (noisy) voxel \times conditions beta matrix for each simulated
379 run. The generated beta matrices were then passed on to the dimensionality
380 estimation.

381 We capitalized on our prior knowledge of possible dimensionalities that
382 could underly the pattern and thus tested only for reconstruction correlations
383 that could be achieved by keeping either 2, 4, or 6 components. This resulted
384 in three reconstruction correlations per run. Reconstruction correlations were
385 averaged across runs and we assessed how often each of the possible models
386 of dimensionality achieved the highest correlation across subjects, for each
387 respective ground-truth. Ideally, for each participant, the highest reconstruc-
388 tion correlation would be achieved by the k -dimensional reconstruction that

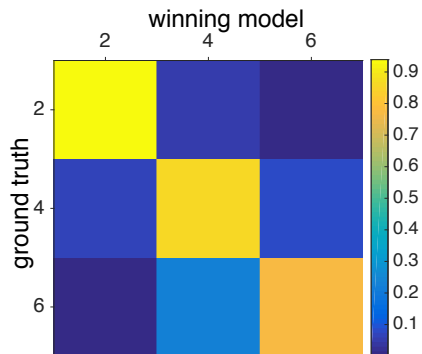


Figure 5: Results from the simulation. Confusion matrix indicating with which percentage a 2, 4, or 6 dimensional model was picked as best model given a ground-truth dimensionality of 2, 4, or 6. Values in the diagonal, that is, where the correct model was picked for a given ground-truth, were consistently higher than off-diagonal values.

389 fits with the underlying ground truth. However, due to noise, deviations from
390 this are possible and we aimed to assess how likely those deviations could be
391 expected to occur.

392 To gather a reliable estimate of the performance of our procedure, we ran
393 a total of 1000 of these simulations for each ground-truth dimensionality.

394 3.2. Results

395 Across 1000 simulations of data with a ground-truth dimensionality of 2,
396 4, or 6, we found that the highest reconstruction correlations were generally
397 achieved by the low-dimensional reconstruction of the data that matched
398 the ground-truth (see Figure 5), with 93.9%, 85.7%, and 76.7% correctly
399 classified, respectively.

400 As described earlier, key to our method is the fact that keeping more
401 components than actually underly an observed pattern results in a reduced
402 reconstruction correlation, thereby allowing us to identify the best dimen-
403 sionality estimate based on the achieved reconstruction correlation. Figure
404 1 B illustrates how the reconstruction correlations drop as components are
405 added that do not carry signal for the case of a true underlying dimensionality
406 of 4.

407 3.3. Discussion

408 By applying our procedure to simulated fMRI data with different under-
409 lying ground truth dimensionality, we tested how well estimated dimension-
410 alities match with true dimensionalities. The results confirm the validity of
411 our approach, showing that for data with reasonable signal-to-noise ratio,
412 estimated dimensionalities match closely the underlying ground truth. One
413 observation is that estimates become more confusable at higher dimension-
414 alities.

415 4. Data sets

416 Following the successful tests of our procedure with simulated data, we
417 applied our method to three different, previously published fMRI datasets, all
418 employing visual stimuli and testing healthy populations. We tested three
419 core aims of our method: 1) Identifying areas carrying functional dimen-
420 sionality, 2) Using functional dimensionality to assess sensitivity to stimulus
421 features, and 3) Measuring task-dependent differences in dimensionality.

422 4.1. Identifying areas carrying functional dimensionality

423 Using data from a category learning study by Mack et al. (2013), we aimed
424 to identify areas carrying functional dimensionality and compare them with
425 the areas found by the original authors' model-based analysis. Model-based
426 analyses make specific assumptions about representational geometry that
427 our approach does not. Furthermore, these analyses require some underlying
428 dimensionality to identify an area. Therefore, we expected our method to
429 reveal significant functional dimensionality in all areas that were reported in
430 the original study, as well as additional areas that were reliably modulated by
431 the task in a way that was not captured by the model tested in the original
432 publication.

433 4.1.1. Methods

434 Participants were trained on categorizing nine objects that differed on four
435 binary dimensions: shape (circle/triangle), color (red/green), size (large/small),
436 and position (left/right). During the fMRI session, participants were pre-
437 sented with the set of all 16 possible stimuli and had to perform the same
438 categorization task. Out of 23 participants, 20 were included in the final
439 analysis presented here, with 19 participants completing 6 runs composed of
440 48 trials and one participant completing 5 runs.

441 Standard pre-processing steps were carried out using SPM12 (Penny et al.,
442 2006) and beta estimates were derived from a GLM containing one regres-
443 sor per stimulus (16 in total, see Supplemental Materials for details). The
444 dataset was retrieved from `osf.io/62rgs`.

445 We ran a whole-brain searchlight with a 7mm radius sphere to estimate
446 which brain areas carry signal with functional dimensionality, that is, signal
447 that could be reliably predicted across runs based on a low-dimensional recon-
448 struction. For each searchlight, data were pre-whitened and mean-centered
449 as described above. Dimensionality estimation was performed as previously
450 described and the resulting j correlations and dimensionality estimates were
451 ascribed to the center of the searchlight. The code for the searchlight was
452 based on the RSA toolbox (Nili et al., 2014).

453 For each voxel, the j correlation coefficients were averaged and their sig-
454 nificance was assessed via non-parametric one-sample t-tests across subjects
455 using FSL’s `randomise` function (Winkler et al., 2014). Results were family-
456 wise error (FWE) corrected using a TFCE threshold of $p < .05$.

457 In their original analysis, the authors fit a cognitive model to participants
458 classification behavior to estimate attention-weights to the single stimulus
459 features. Based on these attention weights, they derived model-based simi-
460 larities between stimuli and used RSA to examine which brain regions show
461 a representational geometry that matches with these predictions. We repli-
462 cated this analysis using the same beta estimates that were passed on to the
463 dimensionality estimation in order to maximize comparability of the two ap-
464 proaches. As for estimating the dimensionality, we ran a whole-brain search-
465 light with a 7mm radius sphere (based on the RSA toolbox, Nili et al., 2014).
466 We averaged voxel response patterns across runs and calculated the repre-
467 sentational distance matrices (RDM) as all pairwise 1–Pearson correlation
468 distance. We assessed correspondence of these RDMs with the model-based
469 distance matrices via Spearman correlation. The resulting Spearman corre-
470 lation for each participant was assigned to the center of the searchlight and
471 their significance was assessed via non-parametric one-sample t-tests across
472 subjects using FSL’s `randomise` function (Winkler et al., 2014). Results were
473 family-wise error (FWE) corrected using a TFCE threshold of $p < .05$.

474 *4.1.2. Results*

475 We aimed to identify areas that show functional dimensionality and ex-
476 amine how those overlap with the authors’ original findings implementing a
477 model-based analysis. We found significant dimensionality (i.e., reconstruc-

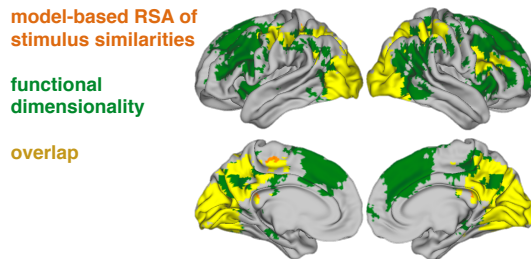


Figure 6: Areas that showed significant functional dimensionality (green), significant fit with the RSA comparing neural representational similarity with model-based predictions of stimulus similarity (orange), or both (yellow). FWE-corrected using a TFCE threshold of $p < .05$. Notably, our method identifies large clusters of functional dimensionality in prefrontal cortex, indicating that areas here were consistently engaged by the task, though their patterns did not fit with the implemented cognitive model.

478 tion correlations) in an extended network of occipital, parietal and prefrontal
479 areas (see Figure 6). In these areas, signal was reliable across runs and showed
480 functional dimensionality.

481 As can be seen in Figure 6, our method successfully identified all areas
482 that were found in the original model-based analysis, which bolsters the
483 authors original interpretation of their results. Notably, we were able to
484 identify further areas that did not show a fit with the implemented attention-
485 based model, suggesting that signal changes in those areas reflect a different
486 aspect of the task space than captured by the cognitive model.

487 *4.1.3. Discussion*

488 Within the first dataset, we showed that by identifying areas with signifi-
489 cant functional dimensionality, it is possible to reveal areas that can plausibly
490 be tested for correspondence with a hypothesized representational similarity
491 structure, as for instance derived from a cognitive model. More specifically,
492 we were able to identify all areas that have been reported in the original
493 analysis by Mack et al. (2013) to show a representational similarity as pre-
494 dicted by a cognitive model. Additionally, we found further areas that had
495 not been revealed in the original analysis to show functional dimensionality.
496 This indicates that those areas have a reliable functional dimensionality but
497 reflect cognitive processes or task-aspects that are not captured by the cog-
498 nitive model. For instance, activation in the medial BA 8 has been found

499 to correlate with uncertainty and task-difficulty (Volz et al., 2005; Huettel,
500 2005; Crittenden and Duncan, 2014), suggesting that the neural patterns in
501 this region in the current task might reflect processes related to the difficulty
502 or category uncertainty of the categorization decision for each stimulus. To-
503 gether, the findings highlight the potential of our procedure to aid evaluation
504 of model performance and identify areas ahead of model-fitting.

505 *4.2. Using functional dimensionality to assess sensitivity to stimulus features*

506 Using data from a study with real-world categories and photographic
507 stimuli by Bracci and Op de Beeck (2016), we tested whether different
508 brain regions show functional dimensionality in response to different stimulus
509 groupings (i.e., depending on how the stimulus-space is summarized). For ex-
510 ample, the columns in the data matrix may be organized along either visual
511 categories or shape. In this fashion, our technique could be useful in eval-
512 uating general hypotheses regarding the nature and basis of the functional
513 dimensionality in brain regions.

514 *4.2.1. Methods*

515 During the experiment, participants were presented repeatedly with 54
516 different natural images that were of nine different shapes and belonged to six
517 different categories (minerals, animals, fruit/vegetables, music instruments,
518 sport instruments, tools), allowing the authors to dissociate between neural
519 responses reflecting shape or category information.

520 Standard pre-processing of the data was carried out using SPM12 (see
521 Supplemental Material for details). In line with the authors original analysis,
522 we tested for differences depending on whether the stimuli were averaged to
523 emphasize their category or shape information. To that end, we constructed
524 two separate GLMs. The first GLM (catGLM) was composed of one regressor
525 per category (six in total), thus averaging across objects shapes. The second
526 GLM (shapeGLM) consisted of nine different regressors, one for each shape,
527 averaging neural responses across object categories. In both GLMs, regres-
528 sors were convolved with the HRF and six motion-regressors as covariates of
529 no interest were included.

530 Dimensionality was estimated separately for both GLMs. We ran a whole-
531 brain searchlight with a 7mm sphere on the beta estimates of the respective
532 GLM, again pre-whitening and mean-centering voxel patterns within each
533 searchlight before estimating the dimensionality. Reconstruction correlations
534 were averaged across runs for each participant and tested for significance

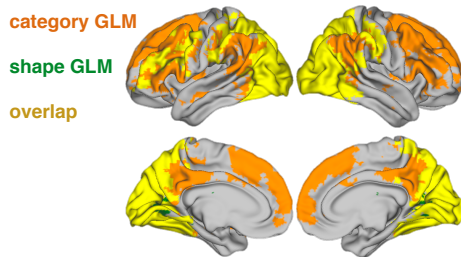


Figure 7: Areas showing significant functional dimensionality for the shape GLM (green), the category GLM (orange), or both (yellow). Results are FWE-corrected using an TFCE threshold of $p < .05$. Across both GLMs, posterior and parietal regions show functional dimensionality. Prefrontal regions show more pronounced functional dimensionality for the category GLM, in line with the original findings.

535 across participants using FSL’s randomise function (Winkler et al., 2014).
536 Results were FWE corrected using a TFCE threshold of $p < .05$.

537 4.2.2. Results

538 When testing for functional dimensionality for the shape-sensitive GLM,
539 we found significant reconstruction correlations in bilateral posterior occipito-
540 temporal and parietal regions, indicating functional dimensionality in these
541 areas. Additionally, a significant cluster was revealed in the left lateral pre-
542 frontal cortex (see Figure 7). Testing for functional dimensionality for the
543 category-sensitive GLM also revealed strong significant correlations in occip-
544 ital and posterior-temporal regions, but notably showed more pronounced
545 correlations in bilateral lateral and medial prefrontal areas as well. This is
546 in line with the authors original findings that showed that neural patterns in
547 parietal and prefrontal ROIs correlated more strongly with a model reflect-
548 ing category similarities, whereas shape similarities were largely restricted to
549 occipital and posterior temporal ROIs.

550 4.2.3. Discussion

551 With the second dataset, we tested whether different areas are identi-
552 fied to express significant functional dimensionality depending on how the
553 underlying task-space is summarized. In line with the original authors’
554 findings (Bracci and Op de Beeck, 2016), we found more pronounced func-
555 tional dimensionality in prefrontal regions for the GLM emphasizing the
556 category-information across stimuli, compared to the one focusing on shape-

557 information. Likewise, functional dimensionality in occipital regions was
558 more pronounced for the shape-based GLM.

559 However, compared to the authors' original findings, we did not find a
560 sharp dissociation between shape and category. For example, we find both
561 shape and category dimensionality present in early visual regions and shape
562 dimensionality extending into frontal areas. As discussed in the previous
563 section, our method provides a general test of dimensionality whereas the
564 original authors evaluate specific representational accounts that make ad-
565 ditional assumptions about shape and category similarity structure. Com-
566 paring results suggest that to some degree the dissociation found in Bracci
567 and Op de Beek (2016) rests on these specific assumptions. A more gen-
568 eral test of functional dimensionality, for stimuli organized along shape or
569 category, provides additional information to assist in interpreting the cog-
570 nitive function of these brain regions, which complements testing more specific
571 representational accounts.

572 *4.3. Measuring task-dependent differences in dimensionality*

573 In this third dataset, we consider whether the underlying dimensionality
574 of neural representations changes as a function of task. In Mack et al.
575 (2016), participants learned a categorization rule over a common stimulus
576 set that either depended on one or two stimulus dimensions. We predicted
577 that the estimated functional dimensionality, as measured by our hierarchi-
578 cal Bayesian method, should be higher for the more complex categorization
579 problem, extending the original authors' findings.

580 *4.3.1. Methods*

581 Participants learned to classify bug stimuli that varied on three binary
582 dimensions (mouth, antenna, legs) into two contrasting categories based on
583 trial-and-error learning. Over the course of the experiment, participants
584 completed two learning problems (in counterbalanced order). Correct classi-
585 fication in type I problem required attending to only one of the bugs features,
586 whereas classification in type II problem required combining information of
587 two features in an exclusive-or manner.

588 Previous research has shown that neural dimensionality appropriate for
589 the problem at hand is linked to successful task performance (Rigotti et al.,
590 2013). Thus, we hypothesized that dimensionality of the neural response
591 would be higher for type II compared to type I in areas known to process
592 visual features, as for instance lateral occipito-temporal cortex (LOC; see e.g.

593 Eger et al., 2008). We included data from 22 participants in our analysis (one
594 participant was excluded due to artifacts in the fMRI data, please refer to
595 the Supplemental Material for further details on the experiment and data
596 preprocessing). The dataset was retrieved from osf.io/5byhb.

597 In order to infer the degree of functional dimensionality, we estimated it
598 across ROIs encompassing LOC in the left and right hemisphere separately
599 for the two categorization tasks. Because the relevant stimulus dimensions
600 were learned through trial-and-error learning, we excluded the first functional
601 run (early learning) of each problem and analyzed the remaining three runs
602 for each problem.

603 Prior to estimating the dimensionality, data were pre-whitened and mean-
604 centered. Dimensionality was estimated across all voxels for each ROI and
605 problem, resulting in 3 (runs) \times 2 (ROIs) \times 2 (problems) correlation co-
606 efficients and dimensionality estimates. Correlation coefficients were aver-
607 aged per participant, ROI and problem and tested for significance using
608 one-sample *t*-tests. To derive the best population estimate for the under-
609 lying dimensionality for each ROI and problem, we implemented the above
610 described hierarchical Bayesian model. To that end, we calculated mean and
611 standard deviation of each participant's dimensionality estimate per ROI
612 and problem and used those summary statistics to estimate the degree of
613 underlying dimensionality for each ROI and problem.

614 *4.3.2. Results*

615 Estimating dimensionality across two different ROIs in LOC and two
616 different tasks allowed us to test whether the estimated dimensionality differs
617 across problems with different task-demands. As participants had to pay
618 attention to one stimulus feature in the type I problem and two stimulus
619 features in the the type II problem, we hypothesized that dimensionality of
620 the neural response would be higher for type II compared to type I in an
621 LOC ROI.

622 Both ROIs showed significant reconstruction correlations across both tasks
623 (lLOC, type I: $t_{21} = 3.08, p = .006$; rLOC, type I: $t_{21} = 2.21, p = .038$; lLOC,
624 type II: $t_{21} = 3.03, p = .006$; rLOC, type II: $t_{21} = 3.37, p = .003$). This shows
625 that signal in the LOC showed reliable functional dimensionality across runs
626 for both problem types, which is a prerequisite for estimating the degree of
627 functional dimensionality.

628 To estimate whether the dimensionality differed across problems, we ana-
629 lyzed the data by implementing a multilevel Bayesian model using Stan (The

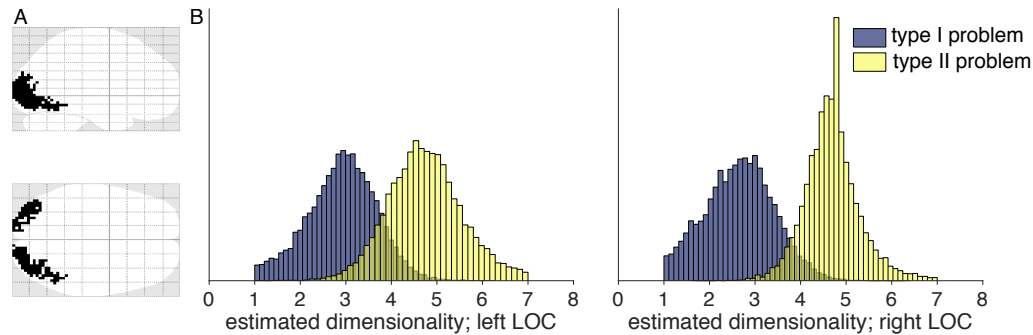


Figure 8: Results of estimating functional dimensionality for two different categorization problems. A: Outline of the two ROIs in left and right LOC. B: Histograms of posterior distributions of estimated dimensionalities in left and right LOC for the type I and II problems. Dimensionalities were estimated by implementing separate multilevel models for each ROI and model using Stan. Across both ROIs, the peak of the posterior distributions of the estimated dimensionality for type II was higher than for type I, mirroring the structure of the two problems.

630 Stan Development Team, 2017), see Figure 2 for an illustration of the model.
631 As hypothesized, the estimated underlying dimensionality was higher for the
632 type II problem compared to type I (type I: $\mu_{left} = 2.92$ (CI 95% : 1.33, 4.33),
633 $\mu_{right} = 2.66$ (CI 95% : 1.23, 4.14); type II: $\mu_{left} = 4.74$ (CI 95% : 3.20, 6.46),
634 $\mu_{right} = 4.69$ (CI 95% : 3.56, 6.06), see Figure 8).

635 4.3.3. Discussion

636 Besides knowing which areas show neural patterns with functional di-
637 mensionality, an important question concerns the degree of the underlying
638 dimensionality. Using data from a categorization task where participants
639 had to attend to either one or two features of a stimulus, we demonstrate
640 how our method can be used to test whether the degree of underlying dimen-
641 sionality of neural patterns varies with task demands. A notable strength
642 of the dataset for our research question is that the authors used the same
643 stimuli in a within-subject paradigm, counterbalancing the order of the two
644 categorization tasks across subjects. This allowed us to investigate how the
645 dimensionality of a neural pattern changes with task, while controlling for
646 possible effects due to differences in signal-to-noise ratios across participants
647 or brain regions.

648 Our results show that, as expected, the degree of underlying functional
649 dimensionality is higher when the task required attending to two stimulus

650 features instead of only one. Notably, this assumption was implicit to the
651 conclusions drawn by the authors in the original publication (Mack et al.,
652 2016). The authors analyzed neural patterns in hippocampus and imple-
653 mented a cognitive model to show that stimulus-specific neural patterns were
654 stretched across relevant compared to irrelevant dimensions. Thus, irrelevant
655 dimensions were compressed and the dimensionality of the neural pattern
656 was reduced the less dimensions were relevant to the categorization problem.
657 Our approach allows to directly assess this effect without the need of fitting
658 a cognitive model.

659 **5. General Discussion**

660 Multivariate and model-based analyses of fMRI data have deepened our
661 understanding of the human brain and its representational spaces (Norman
662 et al., 2006; Kriegeskorte and Kievit, 2013; Haxby et al., 2014; Turner et al.,
663 2017). However, before evaluating specific representational accounts, it is
664 sensible to first ask the more basic question of whether brain areas displays
665 functional dimensionality more generally. Here, we presented a novel ap-
666 proach to estimate an area’s functional dimensionality by a combined SVD
667 and cross-validation procedure. Our procedure identifies areas with signif-
668 icant functional dimensionality and provides an estimate, reflecting uncer-
669 tainty, of the degree of underlying dimensionality. Across three different data
670 sets, we confirmed and extended the findings from the original contributions.

671 After verifying the operation of the method with a synthetic (simulated)
672 dataset in which the ground truth dimensionality was known, we applied
673 our method to three published fMRI datasets. In each case, the procedure
674 confirmed and extended the authors’ original findings, advancing our un-
675 derstanding of the function of the brain regions considered. Each of three
676 datasets highlighted a potential use of estimating functional dimensionality.

677 In the first study, working with data from Mack et al. (2013), we demon-
678 strated that testing for functional dimensionality can complement model-
679 based fMRI analyses that evaluate more specific representational hypothe-
680 ses. First, one cannot find a rich relationship between model representations
681 and brain measures when there is no functional dimensionality in regions
682 of interest. Second, there might be additional areas that display significant
683 functional dimensionality that do not show correspondence with the model.

684 These additional areas invite further analysis as they might implement
685 processes and representations outside the scope of the tested model. Func-

686 tional dimensionality can indicate interesting unexplained signal. For exam-
687 ple, in the first dataset examined, functional dimensionality was found in
688 all the areas identified by Mack et al. (2013), plus medial BA 8, which is a
689 candidate region for task difficulty and response conflict (see Alexander and
690 Brown, 2011, for a model of medial prefrontal cortex function), which was
691 not the authors' original focus but may merit further study.

692 In the second study, working with data from Bracci and Op de Beeck
693 (2016), we demonstrated how stimuli could be grouped or organized in differ-
694 ent fashions to explore how dimensional organization varies across the brain.
695 In this case, the data matrix was either organized along shape or category.
696 We found neural patterns of shape and category selectivity consistent with
697 the authors' original results. However, we found the selectivity to be more
698 mixed in our analyses and identified additional responsive regions, mirroring
699 our results when we considered data from Mack et al. (2013).

700 Our method may have been more sensitive to signal because it makes
701 fewer assumptions about the underlying representational structure and al-
702 lows for individual differences in the underlying dimensions. In this sense,
703 assessing functional complexity complements existing analysis procedures.
704 Indeed, our approach could be used to evaluate multiple stimulus groupings
705 to inform feature selection in encoding models (Diedrichsen and Kriegeskorte,
706 2017; Naselaris et al., 2011).

707 In a third study, working with data from Mack et al. (2016), we evaluated
708 whether our method could identify changes in task-driven dimensionality. By
709 combining estimates of functional dimensionality with a hierarchical Bayesian
710 model, we found that the functional dimensionality in LOC was higher when
711 a category decision required using two features rather than one. These results
712 are consistent with the original authors' theory but were hitherto untestable.

713 In summary, assessing functional dimensionality across these three studies
714 complemented the original analyses and revealed additional nuances in the
715 data. In each case, our understanding of the neural function was further
716 constrained. Moreover, comparing the results to those from model-based
717 and other multivariate approaches was informative in terms of understanding
718 underlying assumptions and their importance.

719 Of course, as touched upon in the Introduction, there are many possible
720 ways to assess dimensional structure in brain measures and progress has been
721 made on this challenge Rigotti et al. (2013); Machens et al. (2010); Rigotti
722 and Fusi (2016); Diedrichsen et al. (2013); Bhandari et al. (2017). Here, our
723 aim was to specify a general, computational efficient, robust, and relatively

724 simple and interpretable procedure that can easily be applied to whole brain
725 data to first test for statistical significant functional dimensionality and, if
726 found, to provide an estimate of its magnitude using Bayesian hierarchical
727 modeling to make clear the uncertainty in that estimate.

728 We hope our contribution is useful to researches interested in further
729 exploring their data, whether it be fMRI, MEG, EEG, or single-cell record-
730 ings. Researchers may consider variants of our method. For example, as
731 mentioned in the Introduction, the SVD could be substituted with another
732 procedure depending on the needs and assumptions of the researchers. There
733 is no magic bullet to the difficult problems of estimating the underlying di-
734 mensionality of noisy neural data, but we have made progress on this issue
735 both theoretically and practically. In doing so, we have also provided addi-
736 tional insights into the brain basis of visual categorization. We hope that
737 by demonstrating the merits of estimating the functional dimensionality of
738 neural data that we motivate others to take advantage of this additional and
739 complementary viewpoint on neural function.

740 **6. Data availability**

741 A Matlab toolbox for estimating functional dimensionality of fMRI data
742 as well as data needed to replicate the analyses presented here will be made
743 available after publication. Nifti files and code for the analyses presented
744 here are available from the authors upon request.

745 **7. Acknowledgments**

746 This work was funded by the National Institutes of Health [grant number
747 1P01HD080679]; the Leverhulme Trust [grant number RPG-2014-075]; and
748 Wellcome Trust Senior Investigator Award [WT106931MA] to Bradley C.
749 Love. Correspondences regarding this work can be sent to c.ahlheim@ucl.ac.uk
750 or b.love@ucl.ac.uk. The authors are grateful to all study authors for sharing
751 their data and wish to thank all members of the LoveLab and Jan Balaguer
752 for valuable input. Declarations of interest: none.

753 **8. References**

- 754 Alexander, W. H. and Brown, J. W. (2011). Medial prefrontal cortex as an
755 action-outcome predictor. *Nature Neuroscience*, 14(10):1338–1344.
- 756 Bhandari, A., Rigotti, M., Gagne, C., Fusi, S., and Badre, D. (2017). Char-
757 acterizing human prefrontal cortex representations with fMRI. In *Society*
758 *for Neuroscience*, Washington, DC:.
- 759 Bracci, S., Daniels, N., and Op de Beeck, H. (2017). Task Context Overrides
760 Object- and Category-Related Representational Content in the Human
761 Parietal Cortex. *Cerebral Cortex*, pages 1–12.
- 762 Bracci, S. and Op de Beeck, H. (2016). Dissociations and associations be-
763 tween shape and category representations in the two visual pathways. *Jour-*
764 *nal of Neuroscience*, 36(2):432–444.
- 765 Cattell, R. B. (1947). Confirmation and clarification of primary personality
766 factors. *Psychometrika*, 12(3):197–220.
- 767 Cattell, R. B. (1966). The Scree Test For The Number Of Factors. *Multi-*
768 *variate Behavioral Research*, 1(2):245–276.
- 769 Crittenden, B. M. and Duncan, J. (2014). Task difficulty manipulation re-
770 veals multiple demand activity but no frontal lobe hierarchy. *Cerebral*
771 *Cortex*, 24(2):532–540.
- 772 Davis, T., LaRocque, K. F., Mumford, J. A., Norman, K. A., Wagner, A. D.,
773 and Poldrack, R. A. (2014). What do differences between multi-voxel and
774 univariate analysis mean? How subject-, voxel-, and trial-level variance
775 impact fMRI analysis. *NeuroImage*, 97:271–283.
- 776 Diedrichsen, J. and Kriegeskorte, N. (2017). Representational models: A
777 common framework for understanding encoding, pattern-component, and
778 representational-similarity analysis. *PLoS Computational Biology*, 13(4):1–
779 33.
- 780 Diedrichsen, J., Wiestler, T., and Ejaz, N. (2013). A multivariate method
781 to determine the dimensionality of neural representation from population
782 activity. *NeuroImage*, 76:225–235.

- 783 Eckart, C. and Young, G. (1936). The approximation of one matrix by
784 another of lower rank. *Psychometrika*, 1(3):211–218.
- 785 Eger, E., Ashburner, J., Haynes, J.-D., Dolan, R. J., and Rees, G. (2008).
786 fMRI activity patterns in human LOC carry information about object
787 exemplars within category. *Journal of cognitive neuroscience*, 20(2):356–
788 370.
- 789 Friston, K. J., Frith, C. D., Liddle, P. F., and Frackowiak, R. S. J. (1993).
790 Functional Connectivity: The Principal-Component Analysis of Large
791 (PET) Data Sets. *Journal of Cerebral Blood Flow & Metabolism*, 13(1):5–
792 14.
- 793 Fusi, S., Miller, E. K., and Rigotti, M. (2016). Why neurons mix: High di-
794 mensionality for higher cognition. *Current Opinion in Neurobiology*, 37:66–
795 74.
- 796 Goddard, E., Klein, C., Solomon, S. G., Hogendoorn, H., Thomas, A., and
797 Klein, C. (2017). Interpreting the dimensions of neural feature represen-
798 tations revealed by dimensionality reduction. *NeuroImage*.
- 799 Hastie, T., Tibshirani, R., and Friedman, J. (2009). Unsupervised Learn-
800 ing. In *The Elements of Statistical Learning*, chapter 14, pages 485–585.
801 Springer, New York, NY, second edition.
- 802 Haxby, J. V., Connolly, A. C., and Guntupalli, J. S. (2014). Decoding Neu-
803 ral Representational Spaces Using Multivariate Pattern Analysis. *Annual*
804 *review of neuroscience*, pages 435–456.
- 805 Hebart, M. N. and Baker, C. I. (2017). Deconstructing multivariate decoding
806 for the study of brain function. *NeuroImage*.
- 807 Huettel, S. A. (2005). Decisions under Uncertainty: Probabilistic Context
808 Influences Activation of Prefrontal and Parietal Cortices. *Journal of Neu-*
809 *roscience*, 25(13):3304–3311.
- 810 Huettel, S. A., Song, A. W., and McCarthy, G. (2003). *Functional magnetic*
811 *resonance imaging*. Sinauer Associates, Inc., Sunderland, Massachusetts,
812 USA, second edition.

- 813 Huth, A. G., Nishimoto, S., Vu, A. T., and Gallant, J. L. (2012). A Continu-
814 ous Semantic Space Describes the Representation of Thousands of Object
815 and Action Categories across the Human Brain. *Neuron*, 76(6):1210–1224.
- 816 Khaligh-Razavi, S. M. and Kriegeskorte, N. (2014). Deep Supervised, but
817 Not Unsupervised, Models May Explain IT Cortical Representation. *PLoS*
818 *Computational Biology*, 10(11).
- 819 Kramer, M. A. (1991). Nonlinear principal component analysis using autoas-
820 sociative neural networks. *AIChE Journal*, 37(2):233–243.
- 821 Kriegeskorte, N. and Kievit, R. A. (2013). Representational geometry: In-
822 tegrating cognition, computation, and the brain. *Trends in Cognitive Sci-*
823 *ences*, 17(8):401–412.
- 824 Kriegeskorte, N., Mur, M., Ruff, D. A., Kiani, R., Bodurka, J., Esteky, H.,
825 Tanaka, K., and Bandettini, P. A. (2008). Matching Categorical Object
826 Representations in Inferior Temporal Cortex of Man and Monkey. *Neuron*,
827 60(6):1126–1141.
- 828 Kruschke, J. K. (2014). *Doing Bayesian data analysis : a tutorial with R,*
829 *JAGS, and Stan.*
- 830 Ledoit, O. and Wolf, M. (2004). A well-conditioned estimator for large-
831 dimensional covariance matrices. *Journal of Multivariate Analysis*,
832 88(2):365–411.
- 833 Machens, C. K., Romo, R., and Brody, C. D. (2010). Functional, but not
834 anatomical, separation of what and when in prefrontal cortex. *Journal of*
835 *Neuroscience*, 30(1):350–360.
- 836 Mack, M. L., Love, B. C., and Preston, A. R. (2016). Dynamic updating
837 of hippocampal object representations reflects new conceptual knowledge.
838 *Proceedings of the National Academy of Sciences of the United States of*
839 *America*, 113(46):13203–13208.
- 840 Mack, M. L., Preston, A. R., and Love, B. C. (2013). Decoding the brain’s
841 algorithm for categorization from its neural implementation. *Current Bi-*
842 *ology*, 23(20):2023–2027.

- 843 Naselaris, T., Kay, K. N., Nishimoto, S., and Gallant, J. L. (2011). Encoding
844 and decoding in fMRI. *NeuroImage*, 56(2):400–410.
- 845 Nichols, T. and Holmes, A. (2003). Nonparametric Permutation Tests
846 for Functional Neuroimaging. *Human Brain Function: Second Edition*,
847 15(1):887–910.
- 848 Nili, H., Wingfield, C., Walther, A., Su, L., Marslen-Wilson, W., and
849 Kriegeskorte, N. (2014). A Toolbox for Representational Similarity Anal-
850 ysis. *PLoS Computational Biology*, 10(4):e1003553.
- 851 Norman, K. A., Polyn, S. M., Detre, G. J., and Haxby, J. V. (2006). Be-
852 yond mind-reading: multi-voxel pattern analysis of fMRI data. *Trends in*
853 *Cognitive Sciences*, 10(9):424–430.
- 854 Parpart, P., Jones, M., and Love, B. (2017). Heuristics as Bayesian inference
855 under extreme priors. *Cognitive Psychology*, in press.
- 856 Penny, W., Friston, K., Ashburner, J., Kiebel, S., and Nichols, T. (2006).
857 *Statistical Parametric Mapping: The Analysis of Functional Brain Images:*
858 *The Analysis of Functional Brain Images*, volume 8. Academic press.
- 859 Rigotti, M., Barak, O., Warden, M. R., Wang, X.-J., Daw, N. D., Miller,
860 E. K., and Fusi, S. (2013). The importance of mixed selectivity in complex
861 cognitive tasks. *Nature*, 497(7451):1–6.
- 862 Rigotti, M. and Fusi, S. (2016). Estimating the dimensionality of neural
863 responses with fMRI Repetition Suppression. *Arxiv*.
- 864 Shlens, J. (2014). A tutorial on principal component analysis. *arXiv*.
- 865 Smith, S. M. and Nichols, T. E. (2009). Threshold-free cluster enhancement:
866 Addressing problems of smoothing, threshold dependence and localisation
867 in cluster inference. *NeuroImage*, 44(1):83–98.
- 868 Spearman, C. (1904). General Intelligence, Objectively Determined and Mea-
869 sured. *American Journal of Psychology*, 15(2):201–292.
- 870 The Stan Development Team (2017). MatlabStan: the MATLAB interface
871 to Stan.

- 872 Turner, B. M., Forstmann, B. U., Love, B. C., Palmeri, T. J., and Van Maa-
873 nen, L. (2017). Approaches to analysis in model-based cognitive neuro-
874 science. *Journal of Mathematical Psychology*, 76:65–79.
- 875 Volz, K. G., Schubotz, R. I., and Von Cramon, D. Y. (2005). Variants of
876 uncertainty in decision-making and their neural correlates. *Brain Research*
877 *Bulletin*, 67(5):403–412.
- 878 Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., and Diedrich-
879 sen, J. (2016). Reliability of dissimilarity measures for multi-voxel pattern
880 analysis. *NeuroImage*, 137:188–200.
- 881 Winkler, A. M., Ridgway, G. R., Webster, M. A., Smith, S. M., and Nichols,
882 T. E. (2014). Permutation inference for the general linear model. *Neu-*
883 *roImage*, 92:381–397.