



Multivariate pattern analysis for MEG: A comparison of dissimilarity measures

Matthias Guggenmos^{a,*}, Philipp Sterzer^a, Radoslaw Martin Cichy^b

^a Visual Perception Laboratory, Charité Universitätsmedizin, Charitéplatz 1, 10117 Berlin, Germany

^b Department of Education and Psychology, Freie Universität Berlin, Habelschwerdter Allee 45, 14195 Berlin, Germany

ARTICLE INFO

Keywords:

MEG
EEG
Multi-voxel pattern analysis
Decoding
Representational similarity analysis
Cross-validation
Noise normalisation
Machine learning

ABSTRACT

Multivariate pattern analysis (MVPA) methods such as decoding and representational similarity analysis (RSA) are growing rapidly in popularity for the analysis of magnetoencephalography (MEG) data. However, little is known about the relative performance and characteristics of the specific dissimilarity measures used to describe differences between evoked activation patterns. Here we used a multisession MEG data set to qualitatively characterize a range of dissimilarity measures and to quantitatively compare them with respect to decoding accuracy (for decoding) and between-session reliability of representational dissimilarity matrices (for RSA). We tested dissimilarity measures from a range of classifiers (Linear Discriminant Analysis – LDA, Support Vector Machine – SVM, Weighted Robust Distance – WeIRD, Gaussian Naïve Bayes – GNB) and distances (Euclidean distance, Pearson correlation). In addition, we evaluated three key processing choices: 1) preprocessing (noise normalisation, removal of the pattern mean), 2) weighting decoding accuracies by decision values, and 3) computing distances in three different partitioning schemes (non-cross-validated, cross-validated, within-class-corrected). Four main conclusions emerged from our results. First, appropriate multivariate noise normalization substantially improved decoding accuracies and the reliability of dissimilarity measures. Second, LDA, SVM and WeIRD yielded high peak decoding accuracies and nearly identical time courses. Third, while using decoding accuracies for RSA was markedly less reliable than continuous distances, this disadvantage was ameliorated by decision-value-weighting of decoding accuracies. Fourth, the cross-validated Euclidean distance provided unbiased distance estimates and highly replicable representational dissimilarity matrices. Overall, we strongly advise the use of multivariate noise normalisation as a general preprocessing step, recommend LDA, SVM and WeIRD as classifiers for decoding and highlight the cross-validated Euclidean distance as a reliable and unbiased default choice for RSA.

Introduction

The investigation of the rapid neural dynamics underlying cognitive functions requires a combination of high-temporal resolution neural measurements with analytical methods that systematically and efficiently probe the information encoded in measured brain activity. A promising approach is the application of multivariate pattern analysis methods (MVPA) to magnetoencephalography (MEG), combining the sensitivity of pattern-based methods with the high temporal resolution of MEG. Two prominent MVPA methods are multivariate decoding (Cox and Savoy, 2003; Haxby et al., 2001; Haynes and Rees, 2005; Kamitani and Tong, 2005), which quantifies the discriminability of condition-specific activation patterns, and representational similarity

analysis (RSA) (Diedrichsen and Kriegeskorte, 2017; Kriegeskorte, 2009; Kriegeskorte et al., 2008a, 2008b; Kriegeskorte and Kievit, 2013). RSA characterizes the similarity of measured responses to experimental conditions in representational dissimilarity matrices (RDMs). As RDMs can in principle be computed for any measurement modality, RSA on MEG opens the way to quantitatively relate rapidly emerging brain dynamics to other data, such as fMRI (Cichy et al., 2016b, 2013) in order to localize responses; computational models (Cichy et al., 2017a, 2016a; Kietzmann et al., 2017; Pantazis et al., 2017; Seeliger et al., 2017; Su et al., 2012; Wardle et al., 2016) in order to understand the underlying algorithms and representational format; to behaviour (Cichy et al., 2017b; Furl et al., 2017; Mur et al., 2013); and across species (Cichy et al., 2014).

At the core of MVPA is the dissimilarity measure used to quantify the

* Corresponding author .

E-mail address: matthias.guggenmos@charite.de (M. Guggenmos).

discriminability (decoding) or the dissimilarity structure (RSA) of evoked activation patterns, fundamentally affecting both the accuracy and the interpretability of results. Yet little is known about the performance and characteristics of different dissimilarity measures for MEG MVPA. Inspired by previous work comparing different dissimilarity measures for fMRI (Walther et al., 2016), we conducted a comprehensive and systematic investigation of dissimilarity measures for MEG to close this gap.

To this end, we compared a set of dissimilarity metrics comprising classifiers (Linear Discriminant Analysis – LDA, Support Vector Machine – SVM, Weighted Robust Distance – WeiRD, Gaussian Naïve Bayes – GNB) and distance measures (Euclidean distance, Pearson correlation). This comparison was done qualitatively, by characterizing dissimilarity time courses, and quantitatively, by comparing decoding accuracies (decoding) and session-to-session reliabilities of RDMs (RSA). We further evaluated the effects of three main processing choices that affect dissimilarity estimation: 1) preprocessing (noise normalisation, removal of the pattern mean), 2) the use of classification decision values to preserve gradual information in classification-based MVPA, and 3) data partitioning (non-cross-validated; cross-validated; within-class-corrected, i.e. subtracting within-from between-condition distances).

Our results give rise to four straightforward recommendations for MVPA in MEG research. First, multivariate noise normalisation is strongly recommended as a general preprocessing step when considering a number of methodological intricacies. Second, for decoding we recommend LDA, SVM and WeiRD, which achieved high accuracies. Third, we show that a previously reported impairment of pattern reliability for decoding accuracy (Walther et al., 2016) can be mitigated by weighting correct and incorrect predictions with classifier decision values. Fourth and finally, concerning distance-based dissimilarity measures for RSA, we recommend the cross-validated Euclidean distance as a robust, gradual, reliable and largely unbiased default choice.

Materials and methods

Data set

The present study is based on a previously published MEG data set (Cichy et al., 2014). This data set was chosen for two reasons. First, the data set has two experimental sessions per participants, enabling us to compute inter-session reliabilities of our measures. Although it is possible to split a single experimental session into subparts to compute reliability, we reasoned that two independent sessions more realistically probe the robustness of a measure with respect to measurement quality (e.g., noise level of individual channels) or daily conditions of participants (e.g. wakefulness or motivation). Second, the employed stimulus set has been used in a number of previous studies (Cichy et al., 2016b, 2014; Cichy and Pantazis, 2016; Khaligh-Razavi and Kriegeskorte, 2014; Kiani et al., 2007; Kriegeskorte et al., 2008b; Mur et al., 2013; Walther et al., 2016), facilitating the comparison of our results with previous literature.

We briefly summarize the most relevant aspects of experimental design and acquisition underlying the present data set (for a detailed description, see Cichy et al., 2014). Participants viewed coloured images of 92 different objects on a grey background presented at the centre of a screen (2.9° visual angle, 500 ms duration), overlaid with a dark grey fixation cross. For each of two MEG sessions, participants completed 10 to 15 runs of 420 s duration each. Each image was presented twice in each MEG run in random order, with a trial onset asynchrony of 1.5 or 2 s. To control vigilance and eye blink behaviour, participants were instructed to press a button and blink their eyes in response to a paper clip that was shown randomly every 3 to 5 trials (average 4). Paper clip trials were excluded from further analysis.

During the experiment, continuous MEG signals from 306 channels (204 planar gradiometers, 102 magnetometers, Elekta Neuromag TRIUX, Elekta, Stockholm) were acquired at a sampling rate of 1000 Hz. Recorded MEG signals were filtered in a frequency range of 0.03–330 Hz

(default setting of Elekta). The lower frequency serves to remove direct current (DC) drifts and its precise value is not critical as long as it is small enough to avoid distortions of event-related responses (see Rousselle, 2012). The higher frequency serves to prevent aliasing. To protect from filter imperfections, the Elekta default value is set to 330 Hz, i.e. below the theoretical Nyquist frequency of 500 Hz. As to our knowledge there are no known informative visually evoked brain signals above the upper limit of the gamma band, i.e. 100 Hz, the precise value of the higher frequency is likewise not critical.

For spatiotemporal filtering we used the MaxFilter software (Elekta, Stockholm), which has been shown to reduce noise and remove artefacts without altering the field patterns of brain signals (Taulu et al., 2004; Taulu and Simola, 2006). We used default parameters (harmonic expansion origin in head frame = [0 0 40] mm; expansion limit for internal multipole base = 8; expansion limit for external multipole base = 3; bad channels automatically excluded from harmonic expansions = 7 s.d. above average; temporal correlation limit = 0.98; buffer length = 10 s). Intuitively, a spatial filter was applied that separated signal data from distant noise sources outside the sensor helmet. Subsequently, a temporal filter was applied that discarded time series components of the signal data that were strongly correlated with noise data.

Finally, raw MEG trials were extracted with 100 ms baseline and a 1000 ms post-stimulus period (i.e., 1101 ms length), yielding 306-dimensional pattern vectors for each time point of a trial. In addition, raw trials were down-sampled by averaging across consecutive 10 ms bins to decrease the computational costs and to increase the signal-to-noise ratio.

General analysis pipeline

We first introduce the general analysis pipelines underlying the comparison of dissimilarity measures for decoding and RSA and thereafter describe each step of the pipeline in detail. As shown in Fig. 1A, in a first step, trials were combined to *pseudo-trials* to improve the overall signal-to-noise ratio. In a second step, pseudo-trials were submitted to an optional preprocessing stage: *multivariate noise normalisation* and/or *removal of the mean pattern*. In a third step, the dissimilarity measures were applied to pseudo-trials, separately for each pairwise combination of conditions and either in a *cross-validated* procedure or a *non-cross-validated* procedure (Fig. 1B). The first three steps were performed for overall 20 randomized assignments of trials to pseudo-trials (*permutations*) and for both sessions of each participant. In a fourth and final step, dissimilarity measures were compared. For decoding, classifiers were compared based on average decoding accuracy (averaged across condition pairs, permutations and sessions). For RSA, dissimilarity measures were compared by means of the session-to-session reliability of representational dissimilarity matrices (averaged across permutations).

Pseudo-trials

To increase the signal-to-noise ratio, for each of the N_C (=92) conditions we created 5 *pseudo-trials* by dividing randomly ordered pre-processed raw trials into 5 approximately equinumerous partitions and then averaging across raw trials within partitions (Fig. 1A). To minimize effects caused by the arbitrariness of this ordering, the procedure was repeated for 20 random orderings of raw trials (henceforth referred to as *permutations*).

Optional preprocessing of pseudo-trials

Prior to MVPA, the MEG data may undergo additional preprocessing. Here, we assessed two popular preprocessing choices: 1) noise normalisation to improve the quality of the data, and 2) removal of the mean pattern to eliminate condition-nonspecific response components.

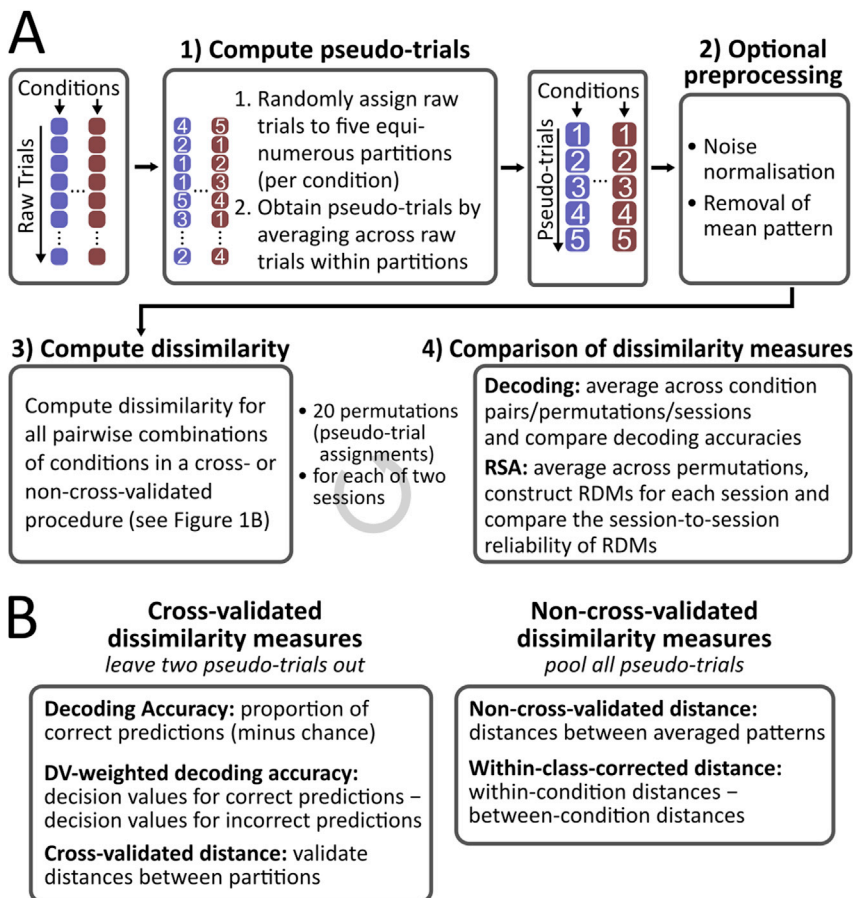


Fig. 1. Analysis pipeline and types of dissimilarity measures. (A) The analysis pipeline comprised the following steps: 1) Raw trials of each condition were assigned to five pseudo-trials through random permutation and averaged. 2) Optionally, pseudo-trials underwent noise normalisation and/or removal of the mean pattern. 3) Dissimilarities were computed on pseudo-trials. Steps 1–3 were performed for 20 permutations (i.e. assignments of trials to pseudo-trials) and for both sessions of each participant. 4) Dissimilarity measures were compared by means of averaged decoding accuracies in case of decoding or session-to-session-reliabilities of representational dissimilarity matrices (RDMs) in case of RSA. (B) Overview of dissimilarity types for the computation of dissimilarity on pseudo-trials. Cross-validated measures are decoding accuracy, decision-value (DV)-weighted decoding accuracy and cross-validated distances; non-cross-validated measures comprise non-cross-validated and within-class-corrected distances.

Multivariate noise normalisation

Commonly, MEG sensors differ in noise levels. To better exploit the information contained in multisensor MEG data, sensors with high noise levels (i.e., unreliable sensors) should be downweighted and sensors with low noise levels (i.e., reliable sensors) should be emphasized. This can be achieved by *univariate noise normalisation* (UNN), where each channel individually is normalised by an estimate of its error variance. In addition, it may be useful to emphasize or deemphasize specific spatial frequencies of MEG patterns. This can be achieved by means of *multivariate noise normalisation* (MNN), where also the error covariance between different sensors is considered. In both procedures, the MEG patterns x are normalised by means of a (co)variance matrix Σ :

$$x^* = \Sigma^{-\frac{1}{2}}x \quad (1)$$

For MNN, off-diagonal elements of Σ correspond to the respective covariances, for UNN they are set to 0.

The (co)variance matrix Σ can be obtained in several different ways characterized by data selection (i.e., baseline phase or entire epoch) and by the level of temporal specificity (i.e., whether the covariance matrix is computed separately for each condition or each time point). To determine the best method for our data set, in a first step we compared the performance of different covariance estimation methods. In brief, we computed Σ either on baseline data (*baseline method*), on the full epoch (*epoch method*) or separately for each time point (*time point method*). Moreover, since rank deficiency is often a problem for matrix inversion, we additionally tested a shrinkage transformation (Ledoit and Wolf, 2004) for the covariance matrix. The supplementary section “Comparison of noise normalisation methods” provides a detailed motivation and description of the methods.

Our main findings (summarized in Figure S1) were that 1) shrinkage improved the performance of all normalisation methods, 2) MNN was

superior to UNN for the epoch and time point method and on par for the baseline method, and 3) the two overall best normalisation methods used MNN based on the epoch or the time point method. Given the slightly higher computational costs of the time point method, for the present work we chose MNN based on the epoch method.

Removing the mean pattern (cocktail-blank removal)

One likely complication when evaluating the dissimilarity of activation patterns is that different experimental conditions often share a common response component. This has two main reasons. First, besides ample differences, experimental conditions may also have common experimental aspects, such as the fact that *any* stimulus was presented, leading to the activation of identical or overlapping neuronal populations. Second, even responses from non-overlapping, but spatially clustered, neuronal populations can produce similar activation patterns at the coarse spatial scale of neuroimaging measurements.

To account for condition-nonspecific response components, previous studies have subtracted the mean pattern from all conditions (“cocktail-blank removal”) (Op de Beeck, 2010; Pietrini et al., 2004; Williams et al., 2008, 2007). However, mean pattern removal is problematic in studies that delineate representational structure, such as RSA, as it introduces dependencies between conditions that may complicate the interpretability of RDMs (Diedrichsen et al., 2011; Garrido et al., 2013; Walther et al., 2016). Nevertheless, cocktail-blank removal is useful as a tool to determine the influence of condition-nonspecific response patterns on the reliability of dissimilarity measures in case of RSA. In the present work, we used mean pattern removal only for this aim, assessing the effect of condition-nonspecific response patterns on the reliability of the Pearson distance in RSA (note that Euclidean distances are unaffected by condition-nonspecific patterns). In detail, we computed the mean sensor pattern across conditions and subtracted it from each sample.

Data partitioning and cross-validation

Dissimilarity measurements between the pseudo trials of different condition was different depending on whether the employed dissimilarity measure was cross-validated or not. For cross-validated measures (i.e., classification, cross-validated distances), the distance of two compared conditions was computed in a stratified cross-validation procedure with one pseudo-trial of each condition in the test set. We kept a pseudo-trial of both respective conditions in the test set in order to enable the computation of test set distances between conditions (see the section on cross-validated distances below). Thus, in each permutation, all but one randomly chosen pseudo-trial of each of any two compared conditions were used for training (hereafter referred to as data partition *A*) and the two left-out pseudo-trials were used for validation (data partition *B*). Note that in this scheme permutations correspond to shuffled cross-validation folds. For non-cross-validated dissimilarity measures (i.e., non-cross-validated distances, within-class-corrected distances), all pseudo-trials of two compared conditions were pooled and the dissimilarity measure was computed once. More specifically, in the case of non-cross-validated distances, data were averaged across pseudo-trials of each condition and distances were computed on these averages. By contrast, because within-class-corrected distances required computing distances between the individual pseudo-trials within a condition, no averaging across pseudo-trials was performed.

Types of dissimilarity measures

In this section, we introduce the dissimilarity measures compared in this work, which we divide into four different groups: 1) classification, 2) non-cross-validated, 3) cross-validated and 4) within-class-corrected distances.

In the following, vectors \mathbf{x} and \mathbf{y} represent the measured MEG patterns associated with two experimental conditions. Each pattern vector has length N_S corresponding to the number of sensors of the MEG device and is computed as the average vector for a considered condition and partition of the data. The goal is to compute a distance $d(\mathbf{x}, \mathbf{y})$ between each pairwise combination of overall N_C conditions.

Decoding accuracy

In classification, an algorithm is trained to discriminate between a set of conditions on the basis of labelled training samples. Here, we used binary classification, where a classifier is applied to each possible pairwise combination of the N_C conditions. Note that decoding accuracies as reported here were always cross-validated, as in each fold classifiers were trained on one partition of the data set (partition *A*) and tested on the left-out partition (partition *B*).

We assessed three common classifiers: Linear Discriminant Analysis classifier (LDA), linear Support Vector Machine (SVM) and Gaussian Naïve Bayes (GNB). We used the implementations provided by Scikit-learn (Abraham et al., 2014) with default parameters (in particular a cost parameter of 1 for SVM and the LIBSVM backend; Chang and Lin, 2011). Exploratory analyses, in which parameters were optimized in a nested cross-validation procedure, did not yield significant benefits and were thus discarded. Note that for LDA, multivariate noise normalisation is originally an integral part of the algorithm itself. In order not to interfere with noise normalisation during preprocessing, the covariance matrix in the LDA algorithm was set to the identity matrix (equivalent to a Euclidean distance-to-centroid classifier). In addition, we included the recently proposed Weighted Robust Distance (WeiRD; Guggenmos et al., 2016; <https://github.com/m-guggenmos/weird>) as a fourth classifier. WeiRD is a distance-to-centroid classifier, where Manhattan distances are computed in a statistically weighted feature space. A more detailed description is provided in the supplementary section “Weighted robust distance”.

The dissimilarity measure based on classification was *decoding accuracy* and was defined as follows:

$$d_{\text{accuracy}} = \frac{1}{N_B} \sum_{b=1}^{N_B} \delta_{c_b, b} - 0.5 \quad (2)$$

where N_B is the number of samples in data partition *B*, δ the Kronecker delta, the condition predicted by the classifier (for sample *b*), c_b the true condition, and 0.5 the chance level. This definition of decoding accuracy has methodological advantages (expected value of zero in the absence of discriminative information; unitless) and is therefore used for all computations. However, note that for illustration purposes in figures we show decoding accuracy as percent correct classification.

Decision-value-weighted decoding accuracy

A possible drawback of decoding accuracy as a distance metric in RSA is the loss of information due to the discretization into correct and incorrect predictions (Walther et al., 2016). Yet, all four classifiers construct some form of an internal continuous decision value (DV), which can be used to ameliorate the drawbacks of discretization. For the present set of classifiers, the decision value *DV* was either the absolute distance to a decision boundary (LDA, SVM, WeiRD) or the probability for the predicted class minus the chance level 0.5 (GNB).

The dissimilarity measure based on DV-weighted decoding accuracies was defined as follows:

$$d_{\text{DV}} = \frac{1}{N_B} \sum_{b=1}^{N_B} DV_b (\delta_{c_b, b} - 0.5) \quad (3)$$

Note that the scale of decision values (and thus the scale of d_{DV}) depends on the applied classifier and the dataset.

Although there are other methods to incorporate graded information from classification, such as the area under the receiver operator curve (AUC), these require a sufficient number of samples in the cross-validation test set and are thus not suited for our analysis pipeline. For instance, given the two test samples in our cross-validation test set, the AUC could only take three values (0, 0.5, 1) and would thus not provide a finer level of differentiation compared to accuracy by itself.

Non-cross-validated distances

Non-cross-validated distances do not require partitioning of the data and can be applied to patterns \mathbf{x} and \mathbf{y} averaged across all samples in a data set. This has both advantages and disadvantages as compared to cross-validated distances. The advantages are that 1) the distance measure is often more robust because distances are computed on a larger set of data, and 2) the procedure is computationally more efficient, as the distance measure is computed only once for a comparison of two conditions. The main disadvantage is that non-cross-validated distances, in addition to estimating the true underlying distance, also capture the dissimilarity due to any source of noise. As a consequence, the measured distances are biased by noise (see also, Walther et al., 2016).

Here we applied non-cross-validated distances after averaging the pattern vectors across all pseudo-trials of a condition. We evaluated two distance measures, the (squared) Euclidean and the Pearson distance, defined as follows:

$$d_{\text{Euclidean}}^2(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T \quad (4)$$

$$d_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = 1 - \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{\sqrt{\text{var}(\mathbf{x})\text{var}(\mathbf{y})}} \quad (5)$$

Cross-validated distances

To obtain unbiased measures of the true dissimilarity, as argued above it is necessary to compute distances in a cross-validated procedure. For this, the pattern vectors of one data partition are projected on those of an independent (validation) partition. As the noise components of the pattern vectors can be assumed to be mostly orthogonal across partitions, they become eliminated by cross-validation (mathematically, note that

the projection product of orthogonal vectors is zero).

This property of cross-validation has two key advantages. First, it can improve the reliability of distance estimates when noise levels differ between measurements. Second, unbiased distances enable testing of ratio-based hypotheses, such as whether one distance is twice as big as another distance, or whether a distance is different from zero (Walther et al., 2016). Such a test would not be sensible if distances are affected by noise as in the case of non-cross-validated distances or classification.

Equations (6) and (7) define the *cross-validated (c.v.)* variants of the (squared) Euclidean and the Pearson distance:

$$d_{\text{Euclidean, c.v.}}^2(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})_{[A]}(\mathbf{x} - \mathbf{y})_{[B]}^T \quad (6)$$

$$d_{\text{Pearson, c.v.}}(\mathbf{x}, \mathbf{y}) = 1 - \frac{\frac{1}{2}(\text{cov}(\mathbf{x}_{[A]}, \mathbf{y}_{[B]}) + \text{cov}(\mathbf{x}_{[B]}, \mathbf{y}_{[A]}))}{\sqrt{\text{cov}(\mathbf{x}_{[A]}, \mathbf{x}_{[B]})\text{cov}(\mathbf{y}_{[A]}, \mathbf{y}_{[B]})}} \quad (7)$$

where A and B denote the two partitions of the data within each cross-validation fold (A: base partition comprising all but 2 left-out pseudo-trials; B: validation partition comprising the 2 left-out pseudo-trials).

A practical issue for cross-validated Pearson distances is the fact that the variances in the denominator can be very small or even negative, leading to exceedingly negative, exceedingly positive or imaginary distances. To accommodate this issue, we regularized the cross-validated Pearson distance in three ways. First, we enforced a positive lower bound ϵ for the variance. As a sensible lower bound depends on the scaling of the data, the parameter ϵ was set to 10% of the non-cross-validated variance of data partition A (the value of 10% was determined empirically). Second, we enforced a lower bound for the denominator, for which the minimum value was set to 25% of the non-cross-validated denominator (likewise determined empirically). And third, we bounded the resulting Pearson distance between 0 and 2 (i.e., corresponding to the bounds of the non-cross-validated Pearson distance).

Within-class-corrected distances

As outlined above, condition-nonspecific response patterns are a ubiquitous phenomenon in neuroimaging measurements. Critically, such condition-nonspecific patterns can impair the accuracy of RSA when they differentially affect specific conditions in two compared modalities (e.g., MEG and behaviour). Here we propose, as a remedy for condition-nonspecific patterns, a procedure we refer to as *within-class correction*, where within-condition dissimilarities are subtracted from between-condition dissimilarities. Within-class correction has been previously used to estimate the discriminatory power of activation patterns (Golarai et al., 2007; Haxby et al., 2001; Weiner et al., 2010). In the context of RSA, it may provide a solution to remove condition-nonspecific patterns in the computation of distance-based dissimilarities.

In addition to removing condition-nonspecific responses, within-class correction can yield unbiased distance estimates under certain circumstances. However, for the two distance measures under investigation this applies only to the Euclidean distance. By contrast, the within-class-corrected Pearson distance is not unbiased in the presence of noise, as both within- and between-condition distances approach 1 with increasing noise and the within-class-corrected distance thus 0. For the Euclidean distance, within-class correction yields an unbiased measure only if the within-condition noise is at the same level as between-condition noise. The fact that within-class correction is thus not generally unbiased is a disadvantage compared to cross-validation.

Mathematically, for the *within-class-corrected distances (w.c.c.)* we subtract the average within-condition distance from the between-condition distance:

$$d_{\text{w.c.c.}}(\mathbf{x}, \mathbf{y}) = \frac{1}{N_p^2} \sum_{i=1}^{N_p} \sum_{j=1}^{N_p} d(\mathbf{x}_i, \mathbf{y}_j) - \frac{1}{N_p(N_p - 1)} \sum_{i=1}^{N_p} \sum_{j=i+1}^{N_p} (d(\mathbf{x}_i, \mathbf{x}_j) + d(\mathbf{y}_i, \mathbf{y}_j)) \quad (8)$$

where d is the given distance measure and $N_p = 5$ is the number of pseudo-trials of a condition. While the first double sum represents the average distance of all pairwise pseudo-trial combinations *between* the two compared conditions (between-condition distance), the second and subtracted double sum represents the average distance of all pairwise pseudo-trial combinations *within* each of the two conditions (within-condition distance).

Decoding

Decoding uses classification to quantify the discriminability of activation patterns pertaining to a set of experimental conditions (Cox and Savoy, 2003; Haxby et al., 2001; Haynes and Rees, 2005; Kamitani and Tong, 2005). Here, we computed decoding accuracies for all pairwise combinations of conditions and all permutations and then averaged across permutations and conditions, separately for each time point. As a result, we obtained a single average decoding accuracy time course for each participant.

Representational similarity analysis

The goal of representational similarity analysis (RSA; Kriegeskorte et al., 2008) applied to neuroimaging data is to characterize the (dis)similarity of activation patterns for a set of experimental conditions. The dissimilarities between all pairwise combinations of conditions are organized in representational dissimilarity matrices (RDMs). In the present work, to construct RDMs we computed dissimilarities for all pairwise combinations of conditions and all permutations and then averaged across permutations, separately for each time point. As a result, we obtained $N_C \times N_C$ RDMs for each time point and participant.

Reliability measures

While decoding accuracies can be directly compared to find the best classifier for decoding, such a direct comparison is not possible across the diverse full set of dissimilarity measures (decoding accuracies, decision values or distances) that can be used to construct RDMs in RSA. For this reason, we used the session-to-session reliability of RDMs as a performance measure that generalizes across types of dissimilarity measures. The rationale for using reliability was that more robust and faithful dissimilarity measures show more replicable RDMs across sessions.

Here, we computed reliabilities between sessions, exploiting the fact that the data set included two equivalent experimental sessions for each participant. In the following, we consider two types of RDM reliability.

First, using the Pearson correlation coefficient, we computed the strength of a linear relationship between two RDMs, regardless of mean and scaling. Mean and scale invariance can be desired properties, e.g. when the mean of two RDMs differs due to varying noise levels or when the scaling differs due to different transfer functions of the involved measurement devices. The *pattern reliability* was defined as follows:

$$R_{\text{Pattern}} = \frac{\text{cov}(\mathbf{d}_1, \mathbf{d}_2)}{\sqrt{\text{var}(\mathbf{d}_1)\text{var}(\mathbf{d}_2)}} \quad (9)$$

where \mathbf{d}_1 and \mathbf{d}_2 represent vectorised RDMs (i.e., the flattened lower triangular part of an RDM) pertaining to the two sessions. R_{Pattern} can take values between -1 and 1 .

On the other hand, it may often be deemed important that two RDMs are as close as possible in terms of their Euclidean distance, thus respecting both mean value and scaling. For instance, one may have reasons to believe that mean differences between RDMs reflect truthful differences that are not trivially explained by different noise levels (e.g., when noise levels were controlled or cross-validation was performed). We therefore additionally computed the normalised sum of squared differences (SSQ) between two RDMs, referred to as the *SSQ reliability*

(Walther et al., 2016):

$$R_{SSQ} = 1 - \frac{\sqrt{\sum_{i=1}^{N_c} (d_{1i} - d_{2i})^2}}{\sqrt{\sum_{i=1}^{N_c} (d_{1i}^2 + d_{2i}^2)}} \quad (10)$$

where R_{SSQ} can take values between 0 and 1.

Note that for the non-cross- and cross-validated Pearson distance, we computed the SSQ reliability on $1-d_{1i}$ and $1-d_{2i}$ in order to enforce an expected dissimilarity of zero for random activation patterns, and thus an expected SSQ reliability of zero. This ensured comparability of the Pearson distance with other dissimilarity measures in terms of SSQ reliability.

Non-cross-validated Euclidean distances likewise do not lead to a distance of zero for random activation patterns. However, since this is a fundamental property of the noise dependency, rather than a technicality as whether to use d or $1-d$, we omitted the non-cross-validated Euclidean distance when reporting SSQ reliability.

Statistical testing

Statistical tests were applied to compare dissimilarity measures either in a time-resolved manner or averaged across time points within a time window of interest. In both cases, sign permutation tests were performed against the null hypothesis of no difference. We employed a full sign permutation scheme across 16 participants. The minimal possible p -value was thus $1/2^{16} = 1.526 \cdot 10^{-5}$, which is denoted as $p < 2^{-16}$ in the text. For time-resolved statistical tests, a sign permutation test was used with Bonferroni-correction for the number of time points.

Availability of data and code

The MEG data analysed here are available for download online from the project page of the original publication (http://userpage.fu-berlin.de/rmcichy/nn_project_page/main.html) or upon request from the authors. In addition, this article is accompanied by an online tutorial with code (Python, MATLAB) and instructions to reproduce key results, available at <https://github.com/m-guggenmos/megmvp>.

Results

Decoding accuracy: comparing classifier performance

Overall, four different classifiers (LDA, SVM, WeiRD and GNB) with and without multivariate noise normalisation (MNN) were evaluated, where each classifier predicted stimulus category labels from MEG data in a time-resolved fashion. Time courses of decoding accuracy were computed by comparison with true category labels (Fig. 2).

An analysis of decoding accuracy time courses showed that the accuracy was above chance for all tested classifiers at each time point of the time window of interest ($p < 0.05$, sign permutation test, Bonferroni-corrected for 100 time points). The accuracy curves peaked around 100 ms (95% confidence intervals illustrated as black bars in Fig. 3 over curves) and decreased afterwards.

For a summary statistical comparison of classifiers, we averaged accuracies across time (from 50 to 550 ms, i.e. accounting for an ~ 50 ms offset between stimulus presentation and cerebral processing). In a first step, we investigated the effect of MNN and found that applying MNN prior to classification boosted classification performance between 5 and 20% percent correct classification. This makes MNN an indispensable preprocessing step. Next, we compared decoding accuracy across noise-normalised classifiers. Noise-normalised LDA, SVM and WeiRD performed comparably well with peak accuracies of over 90% correct classification, while GNB performed markedly worse (Figure S2). Thus, LDA, WeiRD and SVM with MNN were suitable choices for MEG pattern classification.

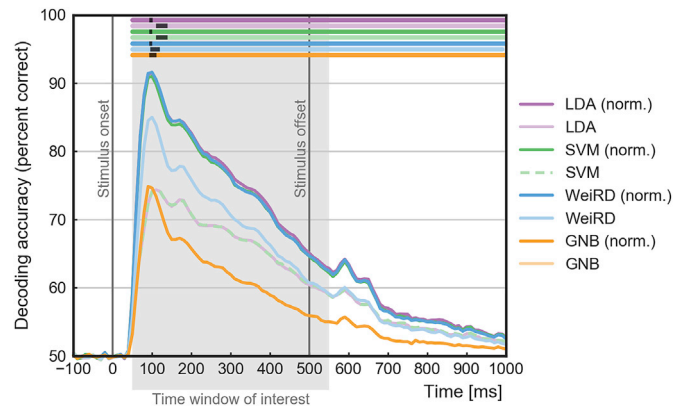


Fig. 2. Decoding: time-resolved MEG decoding accuracies. Color-coded curves (as in legend) report decoding accuracy (here percent correct) for the tested set of classifiers (LDA, SVM, WeiRD, GNB) with and without multivariate noise normalisation (norm.). Horizontal lines at the top of the figure mark statistically significant time points (sign permutation test, Bonferroni-corrected for the number of time points (100), corrected significance level $p < 0.05$). Semi-transparent black bars indicate the bootstrapped 95% confidence interval for peaks (10^6 samples with replacement).

As an aside, for SVM it is common practise to standardize data by z-scoring before fitting. For MEG and EEG data, this is often done in the form of normalizing epoched data channel-wise by subtracting the mean and dividing by the standard deviation of the baseline period. In the terminology of the current manuscript, this amounts to a form of univariate noise normalization. For SVM specifically, we compared this standardization to the MNN procedure employed here (Figure S3), and found that it leads to lower peak decoding accuracy than MNN, and lowers accuracy when combined with MNN compared to MNN alone. Thus, consistent with all other analyses we did not standardize data before submitting it to SVM classification.

In conclusion, MNN is a highly recommended preprocessing step prior to classification. In terms of classifiers, our results indicate that LDA, SVM and WeiRD are all suitable and powerful classifiers and are thus recommended for multivariate decoding in MEG research.

Distance measures: characteristics and the effect of cross-validation and within-class correction

Before investigating reliabilities, it is helpful to identify special properties of different distance measures, as those impact the proper interpretation of session-to-session reliabilities. We thus begin our analysis with a characterization of raw dissimilarity time courses. Readers familiar with these properties or only interested in reliability may skip this section.

Euclidean distance: noise dependency and a remedy through cross-validation

Without cross-validation, the non-cross-validated Euclidean distance (Fig. 3A, dashed yellow curve) remained at a relatively high level from around 100 ms onwards, even though our classification analysis indicated a decrease of informative category-related signals during these later time points. How does cross-validation affect the Euclidean distance estimates across time? We found that cross-validation led to a gradual decrease of Euclidean distance estimates after 100 ms (Fig. 3A, solid yellow curve), consistent with the decrease in decoding accuracy. This result confirms the notion that cross-validation yields unbiased distance measures in the presence of noise (Walther et al., 2016) and thus addresses the noise inflation observed without cross-validation.

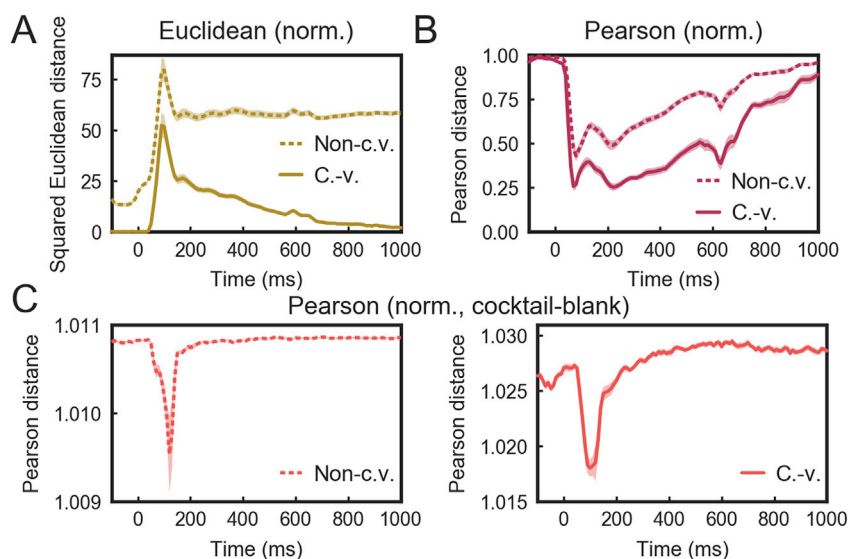


Fig. 3. Raw distance measures across time with (C.-v.) and without (Non-c.v.) cross-validation. (A) Euclidean distance. (B) Pearson distance. (C) Pearson distance with cocktail-blank removal. Note that the Pearson distance exceeded a value of 1, which is indicative of a negative bias caused by cocktail-blank removal. For brevity and clarity, all distance measures are only shown with prior multivariate noise normalisation of the data. Shaded areas indicate bootstrapped standard errors across participants (10^5 samples with replacement).

Pearson distance: the effect of cross-validation and the critical issue of condition-nonspecific signal influences

Different from the Euclidean distance, the Pearson distance between entirely random patterns has an expected value of 1, regardless of the general noise level. Thus, already without cross-validation, the Pearson distance has a meaningful zero point (i.e., when subtracting a constant value of 1). Yet, the Pearson distance is by no means unaffected by noise: as correlation coefficients measuring non-zero correlations decrease with increasing noise, the measured Pearson distance will always show a smaller deviation from 1 compared to the true Pearson distance. This raises the question whether cross-validation, just like in the case of the Euclidean distance, could enable more faithful estimates of the Pearson metric.

The results for our data conformed to these expectations. During the baseline phase (0–100 ms) for which we can assume random MEG patterns, the Pearson distance was 1 both without and with cross-validation (Fig. 3B). With the onset of the stimulus, however, the cross-validated Pearson distance (solid curve) became markedly lower than the non-cross-validated Pearson distance (dashed curve). This indicates that the non-cross-validated Pearson distance underestimated the true (positive) correlation between MEG patterns, a noise bias that was moderated by cross-validation.

Both without and with cross-validation, the time course of the Pearson distance showed a shape that may seem surprising at first: Pearson distances became *smaller* after stimulus onset, i.e. the patterns of different conditions became *more similar*. Indeed, the Pearson distance time courses showed an almost inverted shape compared to the cross-validated Euclidean distance or decoding accuracies. The reason is that, although different conditions differed in stimulus content, they nevertheless shared a number of commonalities, not least the fact that *any stimulus* was presented or that there was always a *stimulus offset* (peak between 600 and 700 ms). However, when activation patterns are rendered highly *similar* through a dominance of condition-nonspecific signals, the goal of RSA is at risk, i.e. mapping the *specific dissimilarity* structure between conditions.

To correct for condition-nonspecific signals, we subtracted the mean pattern from all conditions (cocktail-blank removal) for the Pearson distance (Fig. 3C). We found that cocktail-blank removal was indeed successful in removing the bulk of condition-nonspecific signal contributions. However, confirming earlier work (Diedrichsen et al., 2011; Garrido et al., 2013; Walther et al., 2016), cocktail-blank removal led to negative correlations between activation patterns of conditions (i.e., Pearson distances greater than 1), even in the baseline phase. Because of

these artefactual dependencies, in the present work we limited the application of cocktail-blank removal to assessing a potential inflationary effect of condition-nonspecific response patterns on RDM reliability.

Within-class correction reveals the condition-specific component of non-cross-validated distances

Like cocktail-blank removal, the goal of *within-class correction* is to eliminate condition-nonspecific components from activation patterns (e.g., Haxby et al., 2001). This is achieved by subtracting within-condition distances from between-condition distances. The underlying premise is that condition-nonspecific signals not only affect distances between *different* conditions, but also between repeated measurements of the *same* condition. Within-class correction thus removes signal and noise components unrelated to the difference between conditions.

We found that, when using within-class correction, the raw dissimilarity time courses based on the Euclidean and Pearson distance metrics were more similar to the time courses of decoding accuracy and thus had closer bearing to the discriminatory power of condition-specific activation patterns (Fig. 4). For the Euclidean distance, within-class correction largely eliminated the substantial *condition-nonspecific noise component* that caused high Euclidean distances long after stimulus onset (Fig. 4A). In contrast, for the Pearson distance, within-class correction removed the previously dominating *condition-nonspecific signal component* (Fig. 4B). In particular, within-class correction revealed that the Pearson distance of condition-specific signal components indeed increased with stimulus presentation, just as in the case of the Euclidean distance. We thus conclude that within-class correction accounted for the effect of *condition-nonspecific response* components on dissimilarity, which otherwise affected the non-cross-validated Euclidean and Pearson distance as well as the cross-validated Pearson distance.

Reliability: finding the best dissimilarity measure for representational similarity analysis

We used the reliability of RDMs across measurements to compare the performance of all dissimilarity measures under investigation. Two different reliability measures were applied: *pattern reliability*, a correlation-based measure assessing pattern similarity irrespective of mean and scaling, and *SSQ reliability*, a Euclidean-distance-based measure assessing the similarity taking mean and scaling into account.

We report reliability time courses for all investigated dissimilarity measures: decoding accuracies (Fig. 5A), DV-weighted decoding

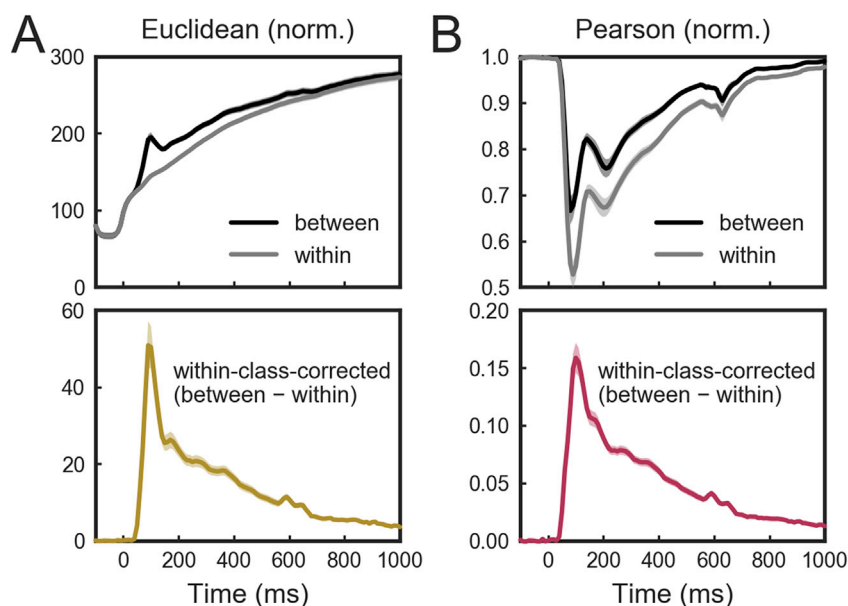


Fig. 4. Within-class-correction of non-cross-validated distances. Top panels show within- and between-condition distances, bottom panels show difference curves (between minus within, referred to as within-class-corrected). Multivariately noise normalised (A) Euclidean distance and (B) Pearson distance. Note that the between-condition Euclidean and Pearson distances shown here are not identical to the non-cross-validated distances of Fig. 3, as they are based on the average distance between all pairwise combinations of pseudo-trials of the two conditions (non-cross-validated distances are computed on the averages across pseudo-trials of each condition). Shaded areas indicate bootstrapped standard errors across participants (10^5 samples with replacement).

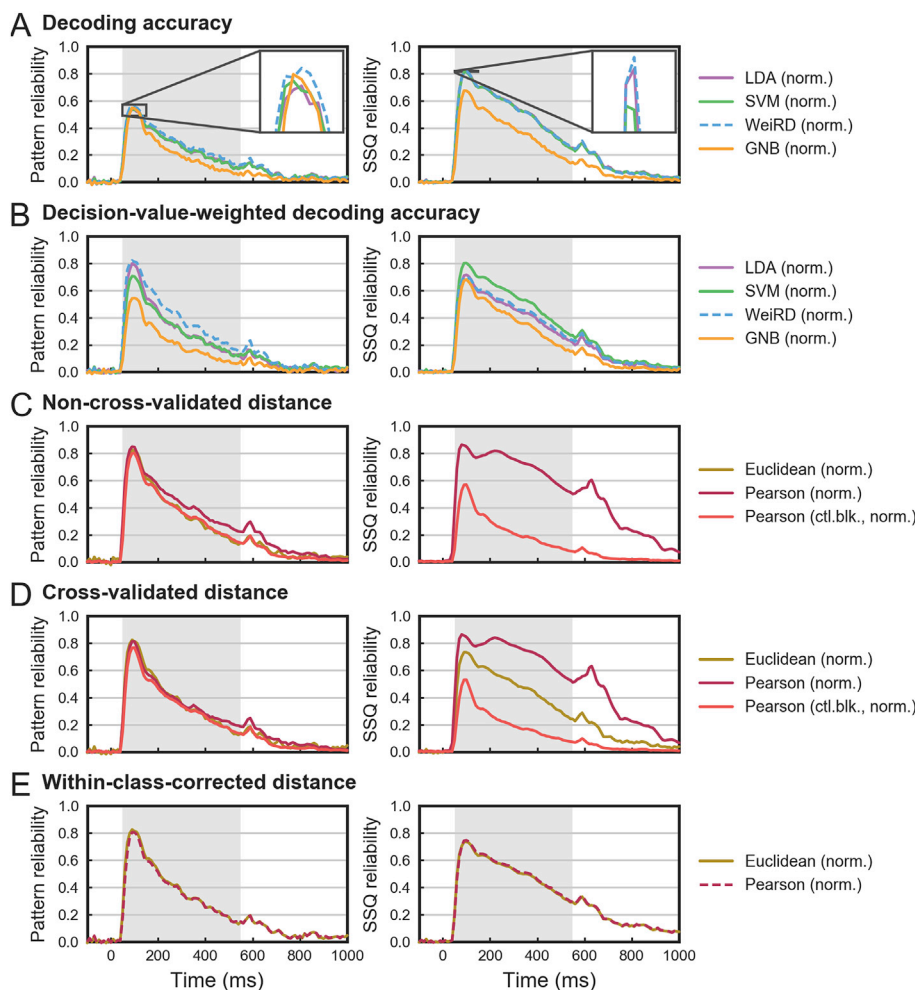


Fig. 5. Time-resolved pattern and SSQ reliability of RDMs across sessions. For clarity, all measures are depicted only with multivariate noise normalisation. The shaded grey area depicts the time window of interest [50 ms; 550 ms]. (A) Decoding accuracies. (B) Decision-value-weighted decoding accuracies. (C) Non-cross-validated distances. Note that for the Euclidean distance, the SSQ reliability cannot be meaningfully computed and is therefore omitted. (D) Cross-validated distances. (E) Within-class-corrected distances. Abbreviations: norm. = multivariate noise normalisation; ctl.blk. = cocktail-blank removal.

accuracies (Fig. 5B), non-cross-validated distances (Fig. 5C), cross-validated distances (Fig. 5D) and within-class-corrected distances (Fig. 5E). For summary purposes, we additionally computed the average

and the maximum reliability over the time window of interest (Fig. 6; Figure S4 for direct comparisons).

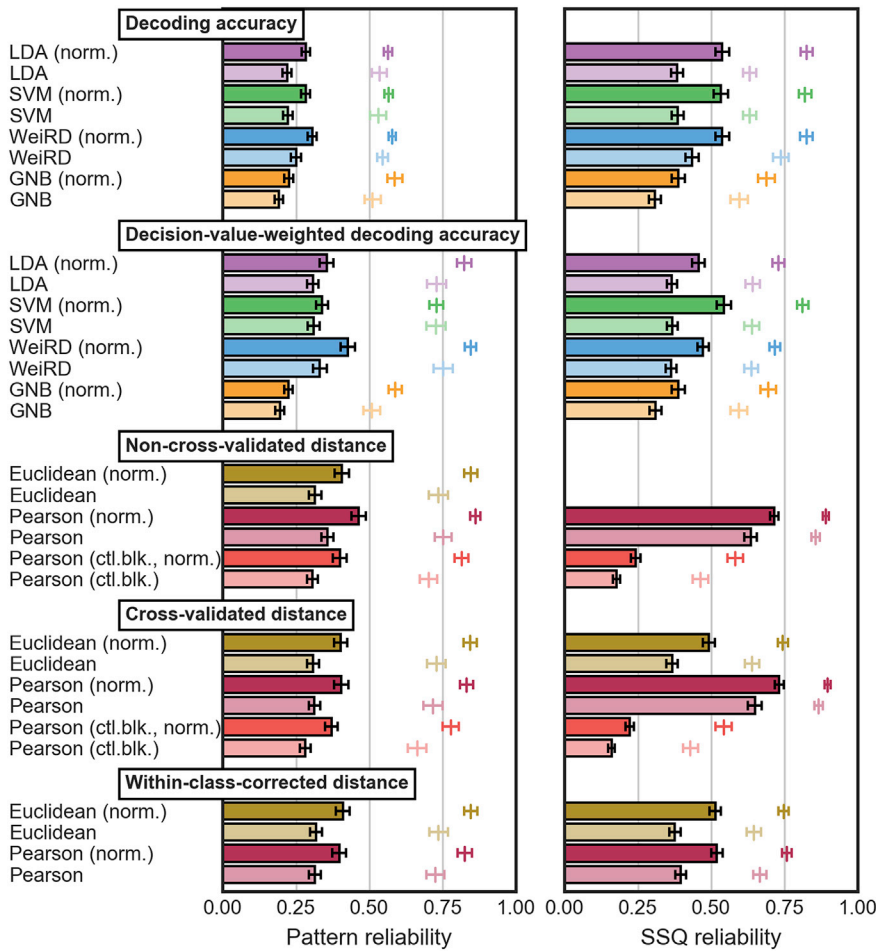


Fig. 6. Average and maximum reliability of RDMs across sessions. Bars represent the average reliability, free-floating lines the maximum reliability within the time window [50 ms; 550 ms]. Error bars indicate bootstrapped standard errors across participants (10^6 bootstrap samples). Abbreviations: norm. = multivariate noise normalisation; ctl.blk. = cocktail-blank removal.

Multivariate noise normalisation improves reliability

As evident from Fig. 6, MNN substantially improved the reliability of all dissimilarity measures. Indeed, in most cases the choice of a dissimilarity measure itself was less critical than whether MNN was performed or not. Across dissimilarity measures, MNN led to an average gain of $\Delta R_{\text{Pattern}} = 0.07$ and $\Delta R_{\text{SSQ}} = 0.11$. In view of these generally beneficial effect of MNN on reliability, below we restricted all further analyses to cases where MNN was applied.

The most accurate classifiers are also the most reliable classifiers

We showed earlier that LDA, SVM and WeiRD were comparable in accuracy, while GNB performed markedly worse. We found that this pattern of results was matched by our analysis of reliability. LDA, SVM and WeiRD showed only small reliability differences, while GNB fell off by a large margin. Between LDA, SVM and WeiRD, WeiRD slightly outperformed LDA ($\Delta R_{\text{Pattern}} = 0.023$, $p < 2^{-16}$, sign permutation test) and SVM ($\Delta R_{\text{Pattern}} = 0.023$, $p < 2^{-16}$) in terms of pattern reliability and SVM in terms of SSQ reliability ($\Delta R_{\text{SSQ}} = 0.004$, $p = 0.036$); all other comparisons were not significant.

We conclude that LDA, SVM and WeiRD were all suitable choices with respect to RDM reliability. In our data set, WeiRD was slightly more reliable than LDA and SVM. More generally, the fact that the reliability analysis yielded much the same conclusion as a comparison of decoding accuracies adds validation to the rationale of using reliability as a performance criterion for dissimilarity measures.

The non-cross-validated Euclidean and Pearson distance are equally reliable when accounting for condition-nonspecific response components

Judging between the Euclidean and the Pearson distance, which is the

more reliable distance measure when no cross-validation or within-class correction is performed? We found that there are two answers to this question.

When considering pattern reliability,¹ the Pearson distance outperformed the Euclidean distance by a large margin ($\Delta R_{\text{Pattern}} = 0.057$, $p < 2^{-16}$). Indeed, in this view, the non-cross-validated Pearson distance was the most reliable measure across all tested measures.

However, given the strong contribution of condition-nonspecific response components to the Pearson distance, the question arises to what degree these shared components might inflate the reliability. With respect to the goals of RSA this is a critical question, because reliability gains based on response components that are shared between conditions are of no use. Evaluating the non-cross-validated Pearson distance after cocktail-blank removal revealed that the reliability of the non-cross-validated Pearson distance became indistinguishable from the non-cross-validated Euclidean distance ($\Delta R_{\text{Pattern}} = -0.007$, $p = 0.17$). This suggests that the reliability advantage of the non-cross-validated Pearson distance was based on condition-nonspecific response components.

In sum, the non-cross-validated Euclidean and Pearson distance are equally reliable when condition-nonspecific response components are subtracted through cocktail-blank removal.

¹ Note that the SSQ reliability cannot be meaningfully computed for the non-cross-validated Euclidean distance and thus is not considered in this comparison.

Cross-validation is robust for the Euclidean distance, but not for the Pearson distance

Two aspects of cross-validation may impact reliability in opposite ways. On the one hand, cross-validation provides faithful distance measures that are largely unbiased in the presence of noise. This could benefit reliability in cases where noise levels differ between measurements. On the other hand, the split of the data into cross-validation folds might negatively impact the robust estimation of distances and thus reliability.

In our data, we found a marginal reduction in pattern reliability when cross-validating the Euclidean distance as compared to no cross-validation ($\Delta R_{\text{Pattern}} = -0.005$, $p = 0.068$) and a significant loss in case of the Pearson distance (without cocktail-blank: $\Delta R_{\text{Pattern}} = -0.059$, $p < 2^{-16}$; with cocktail-blank: $\Delta R_{\text{Pattern}} = -0.029$, $p < 2^{-16}$). In terms of SSQ reliability, cross-validation had a negative effect on the Pearson distance without ($\Delta R_{\text{SSQ}} = -0.017$, $p = 0.0001$), but not with cocktail-blank removal ($\Delta R_{\text{SSQ}} = 0.021$, $p < 2^{-16}$). Detrimental effects of cross-validation on pattern reliability were confirmed by simulation (Figure S5), suggesting that they were generic and not specific to our data set.

Taken together, while cross-validation has a unique advantage in terms of providing unbiased distance estimates, it is robust only for the Euclidean distance and reduces reliability for the Pearson distance.

Within-class correction provides condition-specific estimates of the Pearson distance at high reliability

Within-class correction addresses the issue of condition-nonspecific signal and noise components that are shared between activation patterns of different conditions. Yet, how reliable are the RDMs generated through within-class-corrected distances compared to non-cross- and cross-validated distances?

For the Euclidean distance, we found equal reliability of within-corrected distances and non-cross-validated distances ($\Delta R_{\text{Pattern}} = -0.001$, $p = 0.33$) and a slight advantage of within-class correction over cross-validation ($\Delta R_{\text{Pattern}} = 0.007$, $p = 0.0002$; $\Delta R_{\text{SSQ}} = 0.022$, $p < 2^{-16}$). Thus, in terms of reliability only, within-class correction of the Euclidean distance is not associated with large changes in reliability.

For the Pearson distance, the effect of within-class correction strongly depended on whether non-cross- and cross-validated distances were computed with prior cocktail-blank removal. Without cocktail-blank removal, within-class-corrected distances were generally inferior to non-cross-validated ($\Delta R_{\text{Pattern}} = -0.066$, $p < 2^{-16}$; $\Delta R_{\text{SSQ}} = -0.196$, $p < 2^{-16}$) and cross-validated distances ($\Delta R_{\text{Pattern}} = -0.007$, $p = 0.100$; $\Delta R_{\text{SSQ}} = -0.213$, $p < 2^{-16}$). By contrast, when the within-class-corrected Pearson distance was compared to the cocktail-blank-corrected non-cross- and cross-validated distances this pattern was largely reversed. Here, within-class-corrected distances showed strong reliability advantages compared to cocktail-blank-corrected non-cross-validated ($\Delta R_{\text{SSQ}} = 0.277$, $p < 2^{-16}$; but $\Delta R_{\text{Pattern}} = -0.002$, $p = 0.33$) and cross-validated ($\Delta R_{\text{Pattern}} = 0.027$, $p = 0.0004$; $\Delta R_{\text{SSQ}} = 0.299$, $p < 2^{-16}$) distances. Thus, not only did within-class correction rectify the effect of condition-nonspecific response components on the Pearson distance in a more valid manner than cocktail-blank removal, it also did so at higher reliability. Note that the much reduced SSQ reliability of non-cross- and cross-validated Pearson distances caused by cocktail-blank removal is explained by the fact that the removal of condition-nonspecific components leads to a substantial decrease of correlation coefficients (denominator in Eq. (10)) that is not matched by an analogue decrease of session-to-session variability (nominator in Eq. (10)).

Overall, our data showed that within-class correction had a small but positive influence on the reliability of Euclidean distances, and a strong positive effect on the reliability of the Pearson distances when those were corrected for condition-nonspecific response components. Of note, this pattern of results was confirmed by simulation (Figure S6). As a main conclusion, we recommend to use within-class correction for the Pearson

distances when condition-nonspecific responses are deemed problematic.

Decision-value weighting narrows the gap in pattern reliability between classification-based and distance-based RSA

Disadvantages of decoding accuracies for RSA are that 1) gradual information contained in the decision values of classifier predictions is lost by using a binary measure for the correctness of predictions, and 2) decoding accuracies naturally decrease with noise and are thus not unbiased. Walther et al. (2016) previously showed that these disadvantages can substantially impair pattern reliability of fMRI data. We observed results consistent with this prediction for MEG, i.e. decoding accuracies exhibited generally lower pattern reliability than distances (Fig. 6). Notably this pattern was reversed for SSQ reliability, where decoding accuracies generally outperformed distance measures. In this case, the discretization of decoding accuracies has a positive effect, likely because discretization constrains individual dissimilarity estimates to -0.5 and 0.5 which prevents unreasonably large differences between activation patterns (e.g., caused by measurement noise). Nevertheless, as the impairment of pattern reliability for decoding accuracies is substantial, raw decoding accuracies should be avoided for RSA.

A potential remedy is to weigh the correctness of the predictions by decision values (DVs) (see Fig. S5 for raw time courses of DV-weighted decoding accuracies). We found that for most classifiers, DV-weighting substantially increased the pattern reliability compared to raw decoding accuracies (Fig. 6B). Yet, DV-weighting impaired the SSQ reliability for most classifiers (exception: SVM), likely because decision values were much more volatile and prone to extreme values than bounded decoding accuracies. Thus, while DV-weighting can lead to considerable improvements in pattern reliability, this may have to be traded off against an impairment in SSQ reliability.

Across classifiers, we found that the pattern reliability of DV-weighted decoding accuracies was higher for WeiRD than LDA ($\Delta R_{\text{Pattern}} = 0.072$, $p < 2^{-16}$) and SVM ($\Delta R_{\text{Pattern}} = 0.088$, $p < 2^{-16}$), and higher for LDA than SVM ($\Delta R_{\text{Pattern}} = 0.016$, $p = 0.028$). In terms of SSQ reliability, SVM outperformed LDA ($\Delta R_{\text{SSQ}} = 0.087$, $p < 2^{-16}$) and WeiRD ($\Delta R_{\text{SSQ}} = 0.071$, $p < 2^{-16}$), and WeiRD outperformed LDA ($\Delta R_{\text{SSQ}} = 0.022$, $p = 0.012$). Thus, when classification-based dissimilarity measures are desired for RSA, DV-weighted decoding accuracies of WeiRD and SVM may be particularly attractive choices depending on whether one prioritizes pattern and SSQ reliability, respectively.

How does the reliability of DV-weighted decoding accuracies compare with distance measures? We here report the comparison with the distance schemes for the Euclidean and Pearson metric that we recommended based on the results in the previous two sections: cross-validated Euclidean distances and within-class-corrected Pearson distances. Our results show that LDA and SVM, despite the improvements in pattern reliability through DV-weighting, could not fully reach the pattern reliability of these two distance measures. However, DV-weighting of WeiRD decoding accuracies yielded higher pattern reliability than both the cross-validated Euclidean distance ($\Delta R_{\text{Pattern}} = 0.025$, $p < 2^{-16}$) and the within-class-corrected Pearson distance ($\Delta R_{\text{Pattern}} = 0.029$, $p < 2^{-16}$). In terms of SSQ reliability, while DV-weighting caused a lower SSQ reliability for WeiRD and LDA compared to the reference distances, the SSQ reliability of SVM became superior to the cross-validated Euclidean distance ($\Delta R_{\text{Pattern}} = 0.052$, $p < 2^{-16}$) and the within-class-corrected Pearson distance ($\Delta R_{\text{Pattern}} = 0.024$, $p = 0.0001$). Thus, DV-weighting narrows the gap in pattern reliability between classification- and distance-based dissimilarity measures, but largely diminishes the advantage in SSQ reliability for classifiers relative to distances.

In sum, DV-weighting ameliorated the impairment in pattern reliability due to the lossy binarization in correct and incorrect responses and is thus a recommendable procedure for classification-based RSA. Regarding the choice of the classifier for DV-weighted accuracies, we recommend WeiRD for high pattern reliability and SVM for high SSQ.

Discussion

Summary

We assessed and compared the reliability of dissimilarity measures for representational similarity analysis of MEG data. In brief, we found that 1) multivariate noise normalisation of the data strongly improved the accuracy of classifiers and the reliability of all dissimilarity measures, 2) distances were in general superior to classifiers in terms of pattern reliability, a difference that 3) could be largely ameliorated through decision-value weighting of decoding accuracies, 4) in terms of reliability the Euclidean metric was en par with or better than the Pearson metric when correcting for condition-nonspecific response components, 5) cross-validation provided robust unbiased distance estimates for the Euclidean distance, but came at the cost of slight reliability reductions and was unstable for the Pearson distance, and 6) within-class correction addressed the problematic influence of condition-nonspecific response components on Pearson distances.

Multivariate noise normalisation improves decoding accuracies and reliability

Multivariate noise normalisation substantially improved decoding accuracies and the reliability across classifiers and distance measures. Importantly, the success of noise normalisation depended on a number of important methodological details.

First, the efficacy of noise normalisation critically depended on the specific method used to compute the cross-sensor covariance matrix Σ . For instance, computing Σ just on the baseline data (*baseline method*) degraded RDM reliabilities, which may either suggest that using the baseline phase provided insufficient data for a robust estimate of Σ , or that the baseline phase was not representative of the noise characteristics during stimulus presentation. We believe that our data speak to the latter interpretation, as another method that estimated Σ separately for each time point (*time point method*) led to equal improvements in reliability. As a bottom line, we thus recommend to always include data from the stimulus period in the estimation of Σ , although care is required to ensure that Σ is not compromised by signal information. Here, this was achieved by computing Σ for each condition separately and subsequent averaging across conditions.

Second, it should be noted that for our method of choice (*epoch method*), univariate noise normalisation was already responsible for a large gain in reliability. Thus, the increase in reliability was largely based on a univariate down-weighting of noisy sensors and up-weighting of sensors with little noise variance.

Yet third, there was nevertheless a significant advantage of MNN relative to UNN. Thus, emphasizing spatial frequencies of the MEG patterns with lower variance and de-emphasizing frequencies with higher variance (i.e., the *multivariate* aspect of MNN) yielded a further gain in reliability.

Overall, multivariate noise normalisation is a highly recommended preprocessing step irrespective of other analytic choices. However, the specific implementation of noise normalisation has to be chosen carefully and should consider potentially different noise structures between baseline and stimulus-driven activity.

Choosing a classifier for decoding

When using multivariate decoding to characterize brain representations it is in most cases desirable to maximize decoding accuracy. The main reason is that most neuroimaging measurements including MEG suffer from decreased sensitivity, i.e. only a fraction of the information that is encoded by the underlying neuronal ensembles can be decoded at the level of MEG or at the level of other non-invasive techniques. It is thus desirable to fully exploit every bit of information that is contained in these measurements, thereby also increasing the chance to meaningfully

test more subtle comparisons between experimental conditions.

Here, we compared four classifiers (LDA, SVM, WeiRD, GNB) with and without multivariate noise normalisation and found two main results. First, whether or not to perform MNN was more critical than the choice of a classifier itself, yielding improvements in (peak) decoding accuracy between 5 and 20%. Second, comparing classifiers with prior MNN, we found that LDA, SVM and WeiRD achieved all high and comparable levels of accuracy (>90% peak accuracy in our data set), while GNB performed much worse. A likely reason for the relatively low performance of GNB is the non-independence of sensors in our data set, which violates the assumption of feature independence implicit to GNB. This is consistent with the observation that GNB classification performance improves when features are decorrelated by means of principal component analysis (see Grootswagers et al., 2017). Overall, for analysis pipelines similar to ours, we strongly advise the use of multivariate noise normalisation and recommend LDA, SVM and WeiRD as suitable classifiers for MEG decoding.

Decision-value weighting boosts the pattern reliability of decoding accuracies

Confirming previous work, we found that the pattern RDM reliability of decoding accuracies was markedly impaired compared to distance measures (Walther et al., 2016). One likely reason for this clear handicap of decoding accuracies is the loss of precision through the binarization of predictions (Walther et al., 2016). Another reason could be the fact that decoding accuracies naturally decrease with noise and are thus not unbiased. Decoding accuracies are thus not recommended for RSA.

To address these problems of decoding accuracy, in the present work we introduced DV-weighting of decoding accuracies as a potential solution. By weighting the correctness of single predictions with classifier decision values, gradual information is reintroduced into classification-based dissimilarity measures. Our results showed that DV-weighting indeed rectified the loss in pattern reliability for classification, advancing it close to the level of distance measures.

As noted above, there are other methods to incorporate graded information from classifier decision values (e.g., AUC) which could not be used due to the limited number of samples (pseudo-trials) in our cross-validation test sets. However, in settings with more samples or different cross-validation procedures these methods become feasible. Future research could compare the performance of DV-weighted accuracies to these alternative methods and evaluate their relative advantages and disadvantages.

In sum, our results discourage the use of raw decoding accuracies for RSA and instead advocate DV-weighting of accuracies if classification-based RSA is desired.

Deciding between the Pearson and the Euclidean distance metric

If reliability were the only criterion, our results would strongly favour the non-cross-validated Pearson distance over the non-cross-validated Euclidean distance. However, this conclusion should be critically vetted in view of an important caveat: while the Euclidean distance is invariant to condition-nonspecific response components, the Pearson distance is strongly affected by such signals. Indeed, in most cases the Pearson distance will be a mixture of the dissimilarity due to both condition-specific signals (i.e., the signal of interest) and condition-nonspecific signals. This hampers interpretation not only of Pearson distance per se, but also of its reliability which can be inflated by condition-nonspecific response components.

Empirically, an inflation of reliability through condition-nonspecific response components was supported by two facts. First, the reliability of the Pearson metric was dramatically reduced when the mean pattern was removed prior to distance computation (cocktail-blank removal). Acknowledging that this procedure introduced new dependencies – a negative correlation – between conditions, the result conclusively

indicates that the mean pattern was a driving factor for the high reliability of the Pearson distance. Second, the reliability of the Pearson distance was likewise reduced when distances were subjected to within-class correction, while this was not the case for the Euclidean distance. In both cases, the reliability of the Pearson distance became indistinguishable from the Euclidean distance to the point of being reversed in favour of the Euclidean distance. These two analyses thus strongly suggest that the high reliability of the non-cross-validated Pearson distance was partly based on condition-nonspecific response components.

Overall, in the absence of a priori reasons to prefer a mean/scale-invariant dissimilarity measure, we thus recommend the Euclidean distance over the Pearson distance for better interpretability and higher condition-specific reliability. If mean/scale invariance is desired, it is advisable to carefully check the impact of condition-nonspecific response components on the goal of RSA. In a best-case scenario, these components cancel out in a comparison of RDMs across modalities and have little effect on accuracy. However, if the contribution of condition-nonspecific components substantially differs between modalities in a condition-pair-specific manner, the accuracy of RSA might be severely impaired.

The case for cross-validation

The unique advantage of cross-validation is that it provides unbiased estimates for distances between conditions. Indeed, unbiased distance estimates are critical for a number of research questions. Consider that we construct an RDM using the Euclidean distances measure and are interested in whether certain parts of the RDM are different from zero, i.e. we want to test whether the conditions in these parts are meaningfully different in terms of their neural representations. This test would be impossible for the non-cross-validated Euclidean distance, which is inflated by noise and produces non-zero distance estimates even for identical neural representations. Cross-validation enables testing this hypothesis by cancelling out noise between partitions of the data, thereby introducing a meaningful zero point. In a similar vein, cross-validation allows testing for ratios between distances e.g. whether distance A is twice as big as distance B (Walther et al., 2016).

Despite this key advantage, our data and simulations suggest that cross-validation may come at a cost in terms of reliability. While the reliability of the Euclidean distance was unchanged by cross-validation in our data, the reliability of the Pearson distance showed a significant decrease. These results imply that potential positive effects of cross-validation on the reliability of the Pearson distance (e.g., increased robustness with respect to varying noise levels between measurements), if at all present, could not outweigh the disadvantageous effects of data splitting in cross-validation. Although a negative effect of data splitting should be mitigated through averaging across cross-validation folds, this may not always be fully compensatory.

Despite this potential hit on reliability, for two reasons our results nevertheless lend support to cross-validation specifically for the Euclidean distance. First, the non-cross-validated Euclidean distance, much more than the Pearson distance, was severely distorted in the presence of noise, which in some cases rendered distance estimates almost uninterpretable. Thus, the need of cross-validation was especially pressing for the Euclidean metric. Second, cross-validation was robust for the Euclidean distance as it did not suffer from instability issues affecting the cross-validated Pearson distance (see below). Overall, cross-validation is thus a recommended procedure for the Euclidean distance.

In contrast, we recommend to forgo cross-validation for the Pearson distance for two reasons. First, the non-cross-validated Pearson distance is invariant to the level of noise in phases where the true signal is zero (e.g., baseline). This is because the expected correlation coefficient is zero in the absence of any signal, irrespective of the noise level. And second, cross-validation of the Pearson distance was associated with a consistent and often marked loss in reliability both in our data in simulation. A likely reason for this relatively strong negative impact of cross-

validation on reliability are the unstable cross-validated variances which form the denominator of the cross-validated Pearson distance. These variances can easily reach close-to-zero or negative values in realistic data sets and thus dramatically distort the resulting Pearson distances. To address this instability issue, careful regularization is required. Yet, even if regularization is successful, the techniques and parameters of regularization may differ between modalities, introducing arbitrary choices and complicating cross-modal comparisons. Together, these results and considerations advise against cross-validating the Pearson distance.

The case for within-class correction

In our data, within-class correction was slightly more reliable than cross-validation for the Euclidean distance and markedly more reliable for the Pearson distance (with cocktail-blank removal). Simulation corroborated this empirical result. Nevertheless, there are additional considerations that should guide the decision of whether to use a within-class-corrected processing scheme.

First, within-class correction provides unbiased distance estimates only for the Euclidean distance metric and only if noise affects within- and between-condition distances to the same degree. Given highly similar shapes of cross-validated (Fig. 3A) and within-class-corrected (Fig. 4A) Euclidean distances, this condition seemed to be met in our data set, but might not generally be the case. Given that the reliability gains of within-class correlation relative to cross-validation are very small in the case of the Euclidean distance ($\Delta R_{\text{Pattern/SSQ}} \sim 0.01$), cross-validation therefore remains the recommended procedure for the Euclidean distance.

Second, in most cases within-class correction will be computationally costlier than cross-validation. The within-class-corrected algorithm applied here assessed the full permutation scheme by computing distances for all combinations of pseudo-trials within a condition and between conditions. This was feasible due to the relatively small number of pseudo-trials per condition. However, as the number of permutations grows quadratically with the number of (pseudo-)trials, random subsampling schemes will become necessary at some point, which adds algorithmic complexity and potentially a certain degree of randomness.

And third, while these first two considerations might swing the pendulum to cross-validation, there nevertheless remains the fact that only within-class correction removes condition-nonspecific signal components. This aspect is important for the Pearson distance, which, as noted above, can be strongly affected by such signals. Thus, if one wishes to use the Pearson distance without the influence of condition-nonspecific signal components and without the bias of cocktail-blank removal, within-class correction is the recommended procedure.

In sum, within-class-corrected distances come both with advantages (eliminated condition-nonspecific signal contributions, meaningful zero point, often higher reliability than cross-validation) and disadvantages (distances are not generally unbiased, increased computational complexity). We recommend to use within-class correction in cases where condition-nonspecific signal contributions are a severe issue, as in the case of the Pearson metric.

Limitations and roadmap for future research

Several limitations apply to our investigation that point towards future research efforts. First, the present work was focused on MVPA for MEG and cannot answer to what degree these results generalize to electroencephalography (EEG). Similarly, it remains an open question how MEG and EEG can be combined best in the framework of MVPA. A recent study observed that MVPA methods applied to MEG and EEG yielded largely convergent results, and that combining MEG with EEG data revealed more information than either imaging modality alone (Cichy and Pantazis, 2016). This suggests that the results of the present work will generalize well to EEG data, and that MEG and EEG can be fruitfully combined in the MVPA framework. Future investigation is

necessary to confirm these predictions.

Second, it is unknown whether our results fully transfer to analysis of brain responses from experimental settings that differ in task, sensory modality and cognitive function. MEG MVPA has thus far been mostly applied in studies of visual perception, but is rapidly expanding and increasingly applied to other modalities (e.g. auditory stimuli; Akram et al., 2016; Kocagoncu et al., 2017) and in different tasks contexts (Hebart et al., 2017). Future studies that compare analysis options on data sets from a diverse set of experimental settings are needed to determine the generality of the findings. To facilitate this process, we provide the code for the major processing steps in both MATLAB and Python format.

Third, we investigated MEG MVPA in the sensor rather than in source space. The reason was that here we focussed on MEG's high resolution in revealing the temporal rather than the spatial dynamics of neural activity. Our results here are unlikely to trivially carry over to source space analysis, as source localization algorithms make additional assumptions that likely impact the outcome of subsequent analysis choices. In particular, many source localization methods depend on the estimate of a covariance matrix across sensors, a step which will likely interact with subsequent multivariate noise normalization. Similarly, it is likely that different source localization methods will differently affect subsequent multivariate analysis. Research efforts that systematically investigate the interaction between source localization methods and MEG MVPA are required to shed light on these issues.

Fourth, here we focussed on the analysis of evoked, rather than induced responses. The analysis of induced responses requires resolution of MEG data in frequency space that might impact subsequent processing choices. We hope that the current results can serve as a pointer for future studies dedicated to the investigation of MEG MVPA in frequency space.

Fifth, reliability of a measure across independent data sets is a widely acknowledged quality criterion in that it indicates how reproducible results it. In this sense, session-to-session reliability as used here is scientifically useful. However, note that we cannot prove from first principles and with certainty that reliability selects the measure for MEG RSA that in fact captures the space of representations best. Such verification requires knowledge of some ground truth to which results could be compared, and which is lacking in our case. Future empirical investigations that relate MEG data to some ground truth, and intricate computational analysis of RSA measures on artificial data might shed further light on this open question.

Sixth, here we did not systematically evaluate the effect of feature selection techniques, such as principal component analysis (PCA) or nested feature selection (e.g., recursive feature elimination; Guyon et al., 2002). PCA in particular has shown promise as a feature selection step in a recent study (Grootswagers et al., 2017) and it would be interesting to examine the interaction of PCA with the preprocessing steps and dissimilarity measures considered here. Here, we examined feature selection in a basic analysis by investigating the role of gradiometers and magnetometers (Supplementary section 7 and Figures S7 and S8). Our results indicate that planar gradiometers capture most of the information that is relevant for decoding and RDM reliability. Future studies that systematically evaluate the feature selection techniques in the framework of MEG MVPA are needed.

Conclusion and final recommendations

For both decoding and RSA we strongly recommend to apply multivariate noise normalisation on the data, which provided a considerable boost in terms of decoding accuracy and reliability. Importantly, a key finding of our work is that multivariate noise normalisation is most effective if covariance matrices are computed based on data that include stimulus periods. By contrast, we do not recommend to use cocktail-blank removal as a preprocessing step, which failed at its primary task – removing dependencies between conditions – by introducing other dependencies, confirming earlier reports (Diedrichsen et al., 2011; Garrido

et al., 2013; Walther et al., 2016).

If the goal of MVPA is decoding of MEG signals in an information-based framework (i.e., decoding accuracy), we recommend to use LDA, SVM or WeIRD as classifiers, which yielded comparable accuracy. For MEG RSA, weighing in empirical results and considerations of practicality, and assuming there are no a priori constraints on the choice of a dissimilarity metric, we recommend the cross-validated Euclidean distance as a default choice, which provides a gradual, unbiased and yet reliable dissimilarity measure. Moreover, in comparison to the Pearson distance, the Euclidean distance is much less affected by condition-nonspecific signals common to different conditions and does not suffer from instability in cross-validated processing schemes. Nevertheless, when condition-nonspecific signals are not problematic or when the use of within-class correction is feasible, the Pearson distance is likewise a reliable and recommended choice for MEG RSA. If classification-based RSA is desired, we recommend DV-weighting of correct and incorrect predictions, which led to considerable boosts for the pattern reliability of RDMs.

Acknowledgements

We thank Martin N. Hebart and Dimitrios Pantazis for providing insightful comments to a previous version of the manuscript. This research was funded by an Emmy Noether Award (CI 241/1-1) of the German Research Foundation (Deutsche Forschungsgemeinschaft; DFG) to RMC, and by the DFG research group 1617, grants STE 1430/6-1 and STE 1430/6-2. Computing resources were provided by the high performance computing facilities at ZEDAT, FU Berlin.

Appendix A. Supplementary data

Supplementary data related to this article can be found at <https://doi.org/10.1016/j.neuroimage.2018.02.044>.

References

- Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort, A., Thirion, B., Varoquaux, G., 2014. Machine learning for neuroimaging with scikit-learn. *Front. Neuroinf.* 8, 14. <https://doi.org/10.3389/fninf.2014.00014>.
- Akram, S., Presacco, A., Simon, J.Z., Shamma, S.A., Babadi, B., 2016. Robust decoding of selective auditory attention from MEG in a competing-speaker environment via state-space modeling. *Neuroimage* 124, 906–917. <https://doi.org/10.1016/j.neuroimage.2015.09.048>.
- Chang, C., Lin, C., 2011. LIBSVM: a library for support vector machines. *ACM Trans. Intell. Syst. Technol.* 2 (27), 1–27, 27.
- Cichy, R.M., Khosla, A., Pantazis, D., Oliva, A., 2017a. Dynamics of scene representations in the human brain revealed by magnetoencephalography and deep neural networks. *Neuroimage* 153, 346–358. <https://doi.org/10.1016/j.neuroimage.2016.03.063>.
- Cichy, R.M., Khosla, A., Pantazis, D., Torralba, A., Oliva, A., 2016a. Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence. *Sci. Rep.* 6, 27755. <https://doi.org/10.1038/srep27755>.
- Cichy, R.M., Kriegeskorte, N., Jozwik, K.M., Van den Bosch, J.J.F., Charest, I., 2017b. Neural dynamics of real-world object vision that guide behaviour. *BioRxiv*. <https://doi.org/10.1101/147298>.
- Cichy, R.M., Pantazis, D., 2016. Multivariate pattern analysis of MEG and EEG: a comparison of representational structure in time and space. *BioRxiv* 1–30. <https://doi.org/10.1101/095620>.
- Cichy, R.M., Pantazis, D., Oliva, A., 2016b. Similarity-based fusion of MEG and fMRI reveals spatio-temporal dynamics in human cortex during visual object recognition. *Cerebr. Cortex* 26, 3563–3579. <https://doi.org/10.1093/cercor/bhw135>.
- Cichy, R.M., Pantazis, D., Oliva, A., 2014. Resolving human object recognition in space and time. *Nat. Neurosci.* 17, 455–462. <https://doi.org/10.1038/nn.3635>.
- Cichy, R.M., Sterzer, P., Heinze, J., Elliott, L.T., Ramirez, F., Haynes, J.-D., 2013. Probing principles of large-scale object representation: category preference and location encoding. *Hum. Brain Mapp.* 34, 1636–1651. <https://doi.org/10.1002/hbm.22020>.
- Cox, D.D., Savoy, R.L., 2003. Functional magnetic resonance imaging (fMRI) “brain reading”: detecting and classifying distributed patterns of fMRI activity in human visual cortex. *Neuroimage* 19, 261–270. [https://doi.org/10.1016/S1053-8119\(03\)00049-1](https://doi.org/10.1016/S1053-8119(03)00049-1).
- Diedrichsen, J., Kriegeskorte, N., 2017. Representational models: a common framework for understanding encoding, pattern-component, and representational-similarity analysis. *PLoS Comput. Biol.* <https://doi.org/10.1371/journal.pcbi.1005508>.

- Diedrichsen, J., Ridgway, G.R., Friston, K.J., Wiestler, T., 2011. Comparing the similarity and spatial structure of neural representations: a pattern-component model. *Neuroimage* 55, 1665–1678. <https://doi.org/10.1016/j.neuroimage.2011.01.044>.
- Furl, N., Lohse, M., Pizzorni-Ferrarese, F., 2017. Low-frequency oscillations employ a general coding of the spatio-temporal similarity of dynamic faces. *Neuroimage*. <https://doi.org/10.1016/j.neuroimage.2017.06.023>.
- Garrido, L., Vaziri-Pashkam, M., Nakayama, K., Wilmer, J., 2013. The consequences of subtracting the mean pattern in fMRI multivariate correlation analyses. *Front. Neurosci.* 7, 1–4. <https://doi.org/10.3389/fnins.2013.00174>.
- Golarai, G., Ghahremani, D.G., Whitfield-Gabrieli, S., Reiss, A., Eberhardt, J.L., Gabrieli, J.D.E., Grill-Spector, K., 2007. Differential development of high-level visual cortex correlates with category-specific recognition memory. *Nat. Neurosci.* 10, 512. <https://doi.org/10.1038/nn1865>.
- Grootswagers, T., Wardle, S.G., Carlson, T.A., 2017. Decoding dynamic brain patterns from evoked responses: a tutorial on multivariate pattern analysis applied to time series neuroimaging data. *J. Cognit. Neurosci.* 29, 677–697.
- Guggenmos, M., Schmack, K., Sterzer, P., 2016. WeiRD - a fast and performant multivoxel pattern classifier. In: 2016 Int. Work. Pattern Recognit. Neuroimaging, pp. 1–4. <https://doi.org/10.1109/PRNI.2016.7552349>.
- Guyon, I., Weston, J., Barnhill, S., Vapnik, V., 2002. Gene selection for cancer classification using support vector machines. *Mach. Learn.* 46, 389–422. <https://doi.org/10.1023/A:1012487302797>.
- Haxby, J.V., Gobbini, M.I., Furey, M.L., Ishai, A., Schouten, J.L., Pietrini, P., 2001. Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science* (80-.) 293, 2425–2430. <https://doi.org/10.1126/science.1063736>.
- Haynes, J.-D., Rees, G., 2005. Predicting the orientation of invisible stimuli from activity in human primary visual cortex. *Nat. Neurosci.* 8, 686–691. <https://doi.org/10.1038/nn1445>.
- Hebart, M.N., Bankson, B.B., Harel, A., Baker, C.I., Cichy, R.M., 2017. The representational dynamics of task and object category processing in humans. *eLife* 7, e32816. <https://doi.org/10.1101/153684>.
- Kamitani, Y., Tong, F., 2005. Decoding the visual and subjective contents of the human brain. *Nat. Neurosci.* 8, 679–685. <https://doi.org/10.1038/nn1444>.
- Khaligh-Razavi, S.M., Kriegeskorte, N., 2014. Deep supervised, but not unsupervised, models may explain cortical representation. *PLoS Comput. Biol.* 10 <https://doi.org/10.1371/journal.pcbi.1003915>.
- Kiani, R., Esteky, H., Mirpour, K., Tanaka, K., 2007. Object category structure in response patterns of neuronal population in monkey inferior temporal cortex. *J. Neurophysiol.* 97, 4296–4309. <https://doi.org/10.1152/jn.00024.2007>.
- Kietzmann, T.C., Gert, A.L., Tong, F., König, P., 2017. Representational dynamics of facial viewpoint encoding. *J. Cognit. Neurosci.* 29, 637–651. https://doi.org/10.1162/jocn_a_01070.
- Kocagoncu, E., Clarke, A., Devoreux, B.J., Tyler, L.K., 2017. Decoding the cortical dynamics of sound-meaning mapping. *J. Neurosci.* 37, 1312–1319. <https://doi.org/10.1523/JNEUROSCI.2858-16.2016>.
- Kriegeskorte, N., 2009. Relating population-code representations between man, monkey, and computational models. *Front. Neurosci.* 3, 363–373. <https://doi.org/10.3389/neuro.01.035.2009>.
- Kriegeskorte, N., Kievit, R.A., 2013. Representational geometry: integrating cognition, computation, and the brain. *Trends Cognit. Sci.* 17, 401–412. <https://doi.org/10.1016/j.tics.2013.06.007>.
- Kriegeskorte, N., Mur, M., Bandettini, P., 2008a. Representational similarity analysis - connecting the branches of systems neuroscience. *Front. Syst. Neurosci.* 2, 1–28. <https://doi.org/10.3389/neuro.06.004.2008>.
- Kriegeskorte, N., Mur, M., Ruff, D. a, Kiani, R., Bodurka, J., Esteky, H., Tanaka, K., Bandettini, P. a, 2008b. Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron* 60, 1126–1141. <https://doi.org/10.1016/j.neuron.2008.10.043>.
- Ledoit, O., Wolf, M., 2004. A well-conditioned estimator for large-dimensional covariance matrices. *J. Multivariate Anal.* 88, 365–411. [https://doi.org/10.1016/S0047-259X\(03\)00096-4](https://doi.org/10.1016/S0047-259X(03)00096-4).
- Mur, M., Meys, M., Bodurka, J., Goebel, R., Bandettini, P.A., Kriegeskorte, N., 2013. Human object-similarity judgments reflect and transcend the primate-IT object representation. *Psychol.* 4, 1–22. <https://doi.org/10.3389/fpsyg.2013.00128>.
- Op de Beeck, H.P., 2010. Against hyperacuity in brain reading: spatial smoothing does not hurt multivariate fMRI analyses? *Neuroimage*. <https://doi.org/10.1016/j.neuroimage.2009.02.047>.
- Pantazis, D., Fang, M., Qin, S., Mohsenzadeh, Y., Li, Q., Cichy, R.M., 2017. Decoding the orientation of contrast edges from MEG evoked and induced responses. *BioRxiv*. <https://doi.org/10.1101/148056>.
- Pietrini, P., Furey, M.L., Ricciardi, E., Gobbini, M.I., Wu, W.-H.C., Cohen, L., Guazzelli, M., Haxby, J.V., 2004. Beyond sensory images: object-based representation in the human ventral pathway. *Proc. Natl. Acad. Sci. U. S. A* 101, 5658–5663. <https://doi.org/10.1073/pnas.0400707101>.
- Rousselet, G.A., 2012. Does filtering preclude us from studying ERP time-courses? *Front. Psychol.* 3, 1–9. <https://doi.org/10.3389/fpsyg.2012.00131>.
- Seeliger, K., Fritsche, M., Güçlü, U., Schoenmakers, S., Schoffelen, J.-M., Bosch, S.E., Van Gerven, M.A.J., 2017. CNN-based encoding and decoding of visual object recognition in space and time. *BioRxiv*. <https://doi.org/10.1101/118091>.
- Su, L., Fonteneau, E., Marslen-Wilson, W., Kriegeskorte, N., 2012. Spatiotemporal searchlight representational similarity analysis in EMEG source space. In: Proc. - 2012 2nd Int. Work. Pattern Recognit. NeuroImaging, PRNI 2012, pp. 97–100. <https://doi.org/10.1109/PRNI.2012.26>.
- Taulu, S., Kajola, M., Simola, J., 2004. Suppression of interference and artifacts by the signal space separation method. *Brain Topogr.* 16, 269–275. <https://doi.org/10.1023/B:BRAT.0000032864.93890.f9>.
- Taulu, S., Simola, J., 2006. Spatiotemporal signal space separation method for rejecting nearby interference in MEG measurements. *Phys. Med. Biol.* 51, 1759–1768. <https://doi.org/10.1088/0031-9155/51/7/008>.
- Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., Diedrichsen, J., 2016. Reliability of dissimilarity measures for multi-voxel pattern analysis. *Neuroimage* 137, 188–200. <https://doi.org/10.1016/j.neuroimage.2015.12.012>.
- Wardle, S.G., Kriegeskorte, N., Grootswagers, T., Khaligh-Razavi, S.M., Carlson, T.A., 2016. Perceptual similarity of visual patterns predicts dynamic neural activation patterns measured with MEG. *Neuroimage* 132, 59–70. <https://doi.org/10.1016/j.neuroimage.2016.02.019>.
- Weiner, K.S., Sayres, R., Vinberg, J., Grill-Spector, K., 2010. fMRI-adaptation and category selectivity in human ventral temporal cortex: regional differences across time scales. *J. Neurophysiol.* 103, 3349–3365. <https://doi.org/10.1152/jn.01108.2009>.
- Williams, M.A., Baker, C.I., Op de Beeck, H.P., Mok Shim, W., Dang, S., Triantafyllou, C., Kanwisher, N., 2008. Feedback of visual object information to foveal retinotopic cortex. *Nat. Neurosci.* 11, 1439–1445. <https://doi.org/10.1038/nn.2218>.
- Williams, M.A., Dang, S., Kanwisher, N.G., 2007. Only some spatial patterns of fMRI response are read out in task performance. *Nat. Neurosci.* 10, 685–686. <https://doi.org/10.1038/nn1900>.