

Interpretability of Multivariate Brain Maps in Brain Decoding: Definition and Quantification

Seyed Mostafa Kia^{1,2,3,*}

*Center for Mind/Brain Sciences (CIMEC), University of Trento, via delle Regole
101, 38123, Mattarello, TN, Italy*

Abstract

Brain decoding is a popular multivariate approach for hypothesis testing in neuroimaging. Linear classifiers are widely employed in the brain decoding paradigm to discriminate among experimental conditions. Then, the derived linear weights are visualized in the form of multivariate brain maps to further study the spatio-temporal patterns of underlying neural activities. It is well known that the brain maps derived from weights of linear classifiers are hard to interpret because of high correlations between predictors, low signal to noise ratios, and the high dimensionality of neuroimaging data. Therefore, improving the interpretability of brain decoding approaches is of primary interest in many neuroimaging studies. Despite extensive studies of this type, at present, there is no formal definition for interpretability of multivariate brain maps. As a consequence, there is no quantitative measure for evaluating the interpretability of different brain decoding methods. In this paper, first, we present a theoretical definition of interpretability in brain decoding; we show that the interpretability of multivariate brain maps can be decomposed into their reproducibility and representativeness. Second, as an application of the proposed definition, we formalize a heuristic method for approximating the interpretability of multivariate brain maps in a binary magnetoencephalography (MEG) decoding scenario. Third, we pro-

*Corresponding author:

Email address: seyedmostafa.kia@unitn.it (Seyed Mostafa Kia)

¹University of Trento, Trento, Italy

²Fondazione Bruno Kessler (FBK), Trento, Italy

³Centro Interdipartimentale Mente e Cervello (CIMEC), Trento, Italy

pose to combine the approximated interpretability and the performance of the brain decoding into a new multi-objective criterion for model selection. Our results for the MEG data show that optimizing the hyper-parameters of the regularized linear classifier based on the proposed criterion results in more informative multivariate brain maps. More importantly, the presented definition provides the theoretical background for quantitative evaluation of interpretability, and hence, facilitates the development of more effective brain decoding algorithms in the future.

Keywords: MVPA, brain decoding, brain mapping, interpretation, model selection

1. Introduction

Understanding the mechanisms of the brain has been a crucial topic throughout the history of science. Ancient Greek philosophers envisaged different functionalities for the brain ranging from cooling the body to acting as the seat of the rational soul and the center of sensation [1]. Modern cognitive science, emerging in the 20th century, provides better insight into the brain's functionality. In cognitive science, researchers usually analyze recorded brain activity and behavioral parameters to discover the answers of *where*, *when*, and *how* a brain region participates in a particular cognitive process.

To answer the key questions in cognitive science, scientists often employ mass-univariate hypothesis testing methods to test scientific hypotheses on a large set of independent variables [2, 3]. Mass-univariate hypothesis testing is based on performing multiple tests, e.g., t-tests, one for each unit of the neuroimaging data, i.e., independent variables. The high spatial and temporal granularity of the univariate tests provides fair level of interpretability. On the down side, the high dimensionality of neuroimaging data requires a large number of tests that reduces the sensitivity of these methods after multiple comparison correction. Although some techniques such as the non-parametric cluster-based permutation test [4, 5] offer more sensitivity because of the cluster assumption, they still experience low sensitivity to brain activities that are narrowly distributed in time and space [2, 6]. The multivariate counterparts of mass-univariate analysis, known generally as multivariate pattern analysis (MVPA), have the potential to overcome these deficits. Multivariate approaches are capable of identifying complex spatio-

temporal interactions between different brain areas with higher sensitivity and specificity than univariate analysis [7], especially in group analysis of neuroimaging data [8].

Brain decoding [9] is an MVPA technique that delivers a model to predict the mental state of a human subject based on the recorded brain signal. There are two potential applications for brain decoding: 1) brain-computer interfaces (BCIs) [10, 11], and 2) multivariate hypothesis testing [12]. In the first case, a brain decoder with maximum prediction power is desired. In the second case, in addition to the prediction power, extra information on the spatio-temporal nature of a cognitive process is desired. In this study, we are interested in the second application of brain decoding that can be considered a multivariate alternative for mass-univariate hypothesis testing.

In brain decoding, generally, linear classifiers are used to assess the relation between independent variables, i.e., features, and dependent variables, i.e., cognitive tasks [13, 14, 15]. This assessment is performed by solving a linear optimization problem that assigns weights to each independent variable. Currently, brain decoding is the gold standard in multivariate analysis for functional magnetic resonance imaging (fMRI) [16, 17, 18, 19] and magnetoencephalogram/electroencephalogram (MEEG) studies [20, 21, 22, 23, 24, 25, 26]. It has been shown that brain decoding can be used in combination with brain encoding [27] to infer the causal relationship between stimuli and responses [28].

Brain mapping [29] is a higher form of neuroimaging that assigns pre-computed quantities, e.g., univariate statistics or weights of a linear classifier, to the spatio-temporal representation of neuroimaging data. In MVPA, brain mapping uses the learned parameters from brain decoding to produce brain maps, in which the engagement of different brain areas in a cognitive task is visualized. Intuitively, the interpretability of a brain decoder refers to the level of information that can be reliably derived by an expert from the resulting maps. From the neuroscientific perspective, a brain map is considered *interpretable* if it enables the scientist to answer *where*, *when*, and *how* questions.

Typically, a trained classifier is a black box that predicts the label of an unseen data point with some accuracy. Valverde-Albacete and Peláez-Moreno [30] experimentally showed that in a classification task optimizing only classification error rate is insufficient to capture the transfer of crucial information from the input to the output of a classifier. It is also shown by Ramdas et al. [31] that in the case of data with small sample size using

the classification accuracy as a test statistic for two sample testing should be performed with extra cautious. Beside these limitations of classification accuracy in inference, and considering the fact that the best predictive model might not be the most informative one [32]; a classifier, taken alone, only answers the question of *what* is the most likely label of a given unseen sample [33]. This fact is generally known as knowledge extraction gap [34] in the classification context. Thus far, many efforts have been devoted to filling the knowledge extraction gap of linear and non-linear data modeling methods in different areas such as computer vision [35], signal processing [36], chemometrics [37], bioinformatics [38], and neuroinformatics [39].

Despite the theoretical advantages of MVPA, its practical application to inferences regarding neuroimaging data is limited primarily by a lack of interpretability [40, 41, 42]. Therefore, improving the interpretability of linear brain decoding and associated brain maps is a primary goal in the brain imaging literature [43]. The lack of interpretability of multivariate brain maps is a direct consequence of low signal-to-noise ratios (SNRs), high dimensionality of whole-scalp recordings, high correlations among different dimensions of data, and cross-subject variability [15, 44, 45, 14, 46, 47, 48, 49, 50, 51, 52, 41]. At present, two main approaches are proposed to enhance the interpretability of multivariate brain maps: 1) introducing new metrics into the model selection procedure and 2) introducing new penalty terms for regularization to enhance stability selection.

The first approach to improving the interpretability of brain decoding concentrates on the model selection procedure. Model selection is a procedure in which the best values for the hyper-parameters of a model are determined [14]. The selection process is generally performed by considering the generalization performance, i.e., the accuracy, of a model as the decisive criterion. Rasmussen et al. [53] showed that there is a trade-off between the spatial reproducibility and the prediction accuracy of a classifier; therefore, the reliability of maps cannot be assessed merely by focusing on their prediction accuracy. To utilize this finding, they incorporated the spatial reproducibility of brain maps in the model selection procedure. An analogous approach, using a different definition of spatial reproducibility, is proposed by Conroy et al. [54]. Beside spatial reproducibility, the stability of the classifiers [55] is another criterion that is used in combination with generalization performance to enhance the interpretability. For example, [56, 57] showed that incorporating the stability of models into cross-validation improves the interpretability of the estimated parameters (by linear models).

102 The second approach to improving the interpretability of brain decoding
 103 focuses on the underlying mechanism of regularization. The main idea be-
 104 hind this approach is two-fold: 1) customizing the regularization terms to
 105 address the ill-posed nature of brain decoding problems (where the number
 106 of samples is much less than the number of features) [58, 50] and 2) combin-
 107 ing the structural and functional prior knowledge with the decoding process
 108 so as to enhance stability selection. Group Lasso [59] and total-variation
 109 penalty [60] are two effective methods using this technique [61, 62]. Sparse
 110 penalized discriminant analysis [63], group-wise regularization [7], random-
 111 ized Lasso [47], smoothed-sparse logistic regression [64], total-variation L1
 112 penalization [65, 66], the graph-constrained elastic-net [67, 68], and random-
 113 ized structural sparsity [69] are examples of brain decoding methods in which
 114 regularization techniques are employed to improve stability selection, and
 115 thus, the interpretability of brain decoding.

116 Recently, taking a new approach to the problem, Haufe et al. questioned
 117 the interpretability of weights of linear classifiers because of the contribu-
 118 tion of noise in the decoding process [70, 39, 71]. To address this problem,
 119 they proposed a procedure to convert the linear brain decoding models into
 120 their equivalent generative models. Their experiments on the simulated and
 121 fMRI/EEG data illustrate that, whereas the direct interpretation of classifier
 122 weights may cause severe misunderstanding regarding the actual underlying
 123 effect, their proposed transformation effectively provides interpretable maps.
 124 Despite the theoretical soundness, the major challenge of estimating the em-
 125 pirical covariance matrix of the small sample size neuroimaging data [72]
 126 limits the practical application of this method.

127 In spite of the aforementioned efforts to improve the interpretability of
 128 brain decoding, there is still no formal definition for the interpretability of
 129 brain decoding in the literature. Therefore, the interpretability of different
 130 brain decoding methods are evaluated either qualitatively or indirectly (i.e.,
 131 by means of an intermediate property). In qualitative evaluation, to show
 132 the superiority of one decoding method over the other (or a univariate map),
 133 the corresponding brain maps are compared visually in terms of smooth-
 134 ness, sparseness, and coherency using already known facts (see, for exam-
 135 ple, [47, 73]). In the second approach, important factors in interpretability
 136 such as spatio-temporal reproducibility are evaluated to indirectly assess the
 137 interpretability of results (see, for example, [46, 53, 54, 74]). Despite partial
 138 effectiveness, there is no general consensus regarding the quantification of
 139 these intermediate criteria. For example, in the case of spatial reproducibil-

ity, different methods such as correlation [53, 74], dice score [46], or parameter variability [39, 54] are used for quantifying the stability of brain maps, each of which considers different aspects of local or global reproducibility.

With the aim of filling this gap, our contribution is three-fold: 1) Assuming that the true solution of brain decoding is available, we present a theoretical definition of the interpretability. Furthermore, we show that the interpretability can be decomposed into the reproducibility and the representativeness of brain maps. 2) As a proof of the concept, we propose a practical heuristic based on event-related fields for quantifying the interpretability of brain maps in MEG decoding scenarios. 3) Finally, we propose the combination of the interpretability and the performance of the brain decoding as a new Pareto optimal multi-objective criterion for model selection. We experimentally show that incorporating the interpretability into the model selection procedure provides more reproducible, more neurophysiologically plausible, and (as a result) more interpretable maps.

2. Methods

2.1. Notation and Background

Let $\mathcal{X} \in \mathbb{R}^p$ be a manifold in Euclidean space that represents the input space and $\mathcal{Y} \in \mathbb{R}$ be the output space, where $\mathcal{Y} = \Phi^*(\mathcal{X})$. Then, let $S = \{\mathbf{Z} = (\mathbf{X}, \mathbf{Y}) \mid z_1 = (x_1, y_1), \dots, z_n = (x_n, y_n)\}$ be a training set of n independently and identically distributed (iid) samples drawn from the joint distribution of $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ based on an unknown Borel probability measure ρ . In the neuroimaging context, \mathbf{X} indicates the trials of brain recording, e.g., fMRI, MEG, or EEG signals, and \mathbf{Y} represents the experimental conditions or dependent variables. The goal of brain decoding is to find the function $\Phi_S : \mathbf{X} \rightarrow \mathbf{Y}$ as an estimation of the ideal function $\Phi^* : \mathcal{X} \rightarrow \mathcal{Y}$.

As is a common assumption in the neuroimaging context, we assume the true solution of a brain decoding problem is among the family of linear functions \mathcal{H} ($\Phi^* \in \mathcal{H}$). Therefore, the aim of brain decoding reduces to finding an empirical approximation of Φ_S , indicated by $\hat{\Phi}$, among all $\Phi \in \mathcal{H}$. This approximation can be obtained by estimating the predictive conditional density $\rho(\mathbf{Y} \mid \mathbf{X})$ by training a parametric model $\rho(\mathbf{Y} \mid \mathbf{X}, \Theta)$ (i.e., a likelihood function), where Θ denotes the parameters of the model. Alternatively, Θ can be estimated by solving a risk minimization problem:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \mathcal{L}(\Phi(\mathbf{X}), \Phi_S(\mathbf{X}) + \lambda\Omega(\Theta)) \quad (1)$$

where $\mathcal{L} : \mathbf{Z} \times \mathbf{Z} \rightarrow \mathbb{R}^+$ is the loss function, $\Omega : \mathbb{R}^p \rightarrow \mathbb{R}^+$ is the regularization term, and λ is a hyper-parameter that controls the amount of regularization. There are various choices for Ω , each of which reduces the hypothesis space \mathcal{H} to $\mathcal{H}' \subset \mathcal{H}$ by enforcing different prior functional or structural constraints on the parameters of the linear decoding model (see, for example, [75, 76, 60, 77]). The amount of regularization λ is generally decided using cross-validation or other data perturbation methods in the model selection procedure.

In the neuroimaging context, the estimated parameters of a linear decoding model $\hat{\Theta}$ can be used in the form of a brain map so as to visualize the discriminative neurophysiological effect. Although the magnitude of $\hat{\Theta}$ is affected by the dynamic range of data and the level of regularization, it has no effect on the predictive power and the interpretability of maps. On the other hand, the direction of $\hat{\Theta}$ affects the predictive power and contains information regarding the importance of and relations among predictors. This type of relational information is very useful when interpreting brain maps in which the relation between different spatio-temporal independent variables can be used to describe how different brain regions interact over time for a certain cognitive process. Therefore, we refer to the normalized parameter vector of a linear brain decoder in the unit hyper-sphere as a multivariate brain map (MBM); we denote it by $\vec{\Theta}$ where $\vec{\Theta} = \frac{\Theta}{\|\Theta\|}$ ($\|\cdot\|$ represents the 2-norm of a vector).

As shown in Eq. 1, learning occurs using the sampled data. In other words, in the learning paradigm, we attempt to minimize the loss function with respect to Φ_S (and not Φ^*) [78]. Therefore, all of the implicit assumptions (such as linearity) regarding Φ^* might not hold on Φ_S , and vice versa (see the supplementary material for a simple illustrative example). The *irreducible error* ε is the direct consequence of sampling; it sets a lower bound on the error, where we have:

$$\Phi_S(\mathbf{X}) = \Phi^*(\mathbf{X}) + \varepsilon \quad (2)$$

The distribution of ε dictates the type of loss function \mathcal{L} in Eq. 1. For example, assuming a Gaussian distribution with mean 0 and variance σ^2 for ε implies the least squares loss function [79].

2.2. Interpretability of Multivariate Brain Maps: Theoretical Definition

In this section, we present a theoretical definition for the interpretability of linear brain decoding models and their associated MBMs. Our definition

of interpretability is based on two main assumptions: 1) the brain decoding problem is linearly separable; 2) its *unique* and neurophysiologically *plausible*¹ solution, i.e., Φ^* , is available.

Consider a linearly separable brain decoding problem in an ideal scenario where $\varepsilon = 0$ and $\text{rank}(\mathbf{X}) = p$. In this case, Φ^* is linear and its parameters Θ^* are unique and plausible. The unique parameter vector Θ^* can be computed as follows:

$$\Theta^* = \Sigma_{\mathbf{X}}^{-1} \mathbf{X}^T \mathbf{Y} \quad (3)$$

Using Θ^* as the reference, we define the *strong-interpretability* of an MBM as follows:

Definition 1. An MBM $\vec{\Theta}$ associated with a linear function Φ is “strongly-interpretable” if and only if $\vec{\Theta} \propto \Theta^*$.

It can be shown that, in practice, the estimated solution of a linear brain problem (using Eq. 1) is not strongly-interpretable because of the inherent limitations of neuroimaging data, such as uncertainty [80] in the input and output space ($\varepsilon \neq 0$), the high dimensionality of data ($n \ll p$), and the high correlation between predictors ($\text{rank}(\mathbf{X}) < p$). With these limitations in mind, even though in practice the solution of linear brain decoding is not strongly-interpretable, one can argue that some are more interpretable than others. For example, in the case in which $\Theta^* \propto [0, 1]^T$, a linear classifier where $\hat{\Theta} \propto [0.1, 1.2]^T$ can be considered more interpretable than a linear classifier where $\hat{\Theta} \propto [2, 1]^T$. This issue raises the following question:

Problem 1. Let S^1, \dots, S^m be m perturbed training sets drawn from S via a certain perturbation scheme such as jackknife, bootstrapping [81], or cross-validation [82]. Assume $\vec{\Theta}^1, \dots, \vec{\Theta}^m$ are m MBMs of a certain Φ (estimated using Eq. 1 for certain \mathcal{L} , Ω , and λ) on the corresponding perturbed training sets. How can we quantify the proximity of Φ to the strongly-intrepretable solution of brain decoding problem Φ^* ?

¹Here, neurophysiological plausibility refers to the spatio-temporal chemo-physical constraints of the underlying neural activity that is highly dependent on the acquisition device.

236 To answer this question, considering the uniqueness and the plausibility
237 of Φ^* as the two main characteristics that convey its strong-interpretability,
238 we define the interpretability as follows:

239 **Definition 2.** Let α^j ($j = 1, \dots, m$) be the angle between $\vec{\Theta}^j$ and $\vec{\Theta}^*$. The
240 “interpretability” ($0 \leq \eta_\Phi \leq 1$) of the MBM derived from a linear function
241 Φ is defined as follows:

$$\forall j \in \{1, \dots, m\}, \eta_\Phi = \mathbb{E}_S[\cos(\alpha^j)] \quad (4)$$

242 Empirically, the interpretability is the mean of cosine similarities between
243 Θ^* and MBMs derived from different samplings of the training set. In ad-
244 dition to the fact that employing cosine similarity is a common method for
245 measuring the similarity between vectors, we have another strong motivation
246 for this choice. It can be shown that, for large values of p , the distribution of
247 the dot product in the unit hyper-sphere, i.e., the cosine similarity, converges
248 to a normal distribution with 0 mean and variance of $\frac{1}{p}$, i.e., $\mathcal{N}(0, \sqrt{\frac{1}{p}})$. Due
249 to the small variance for a large enough p values, any similarity value that is
250 significantly larger than zero represents a meaningful similarity between two
251 high dimensional vectors (see the supplementary material for more details
252 about the distribution of cosine similarity).

253 In what follows, we demonstrate how the definition of interpretability is
254 geometrically related to the uniqueness and plausibility characteristics of the
255 true solution to brain decoding problem.

256 *2.3. Interpretability Decomposition into Reproducibility and Representative-* 257 *ness*

258 An alternative approach toward quantifying the interpretability is to as-
259 sess separately its uniqueness and neurophysiological plausibility. In this
260 section, we firstly define the reproducibility and representativeness as mea-
261 sures for quantifying the uniqueness and neurophysiological plausibility of
262 brain decoding model, respectively. Then we show how these definitions are
263 related to the definition of interpretability.

264 The high dimensionality and the high correlations between variables are
265 two inherent characteristics of neuroimaging data that negatively affect the
266 uniqueness of the solution of a brain decoding problem. Therefore, a certain
267 configuration of hyper-parameters may result different estimated parameters

on different portions of data. Here, we are interested in assessing this variability as a measure for uniqueness. Let θ_i^j be the i th ($i = 1, \dots, p$) element of an MBM estimated on the j th ($j = 1, \dots, m$) perturbed training set. We first define the *main multivariate brain map* as follows:

Definition 3. The “main multivariate brain map” $\vec{\Theta}^\mu \in \mathbb{R}^p$ of a linear function Φ is defined as the sum of estimated MBMs $\vec{\Theta}^j$ ($j = 1, \dots, m$) on the perturbed training sets S^j in the unit hyper-sphere:

$$\vec{\Theta}^\mu = \frac{\left[\sum_{j=1}^m \theta_1^j \quad \sum_{j=1}^m \theta_2^j \quad \dots \quad \sum_{j=1}^m \theta_p^j \right]^T}{\left\| \left[\sum_{j=1}^m \theta_1^j \quad \sum_{j=1}^m \theta_2^j \quad \dots \quad \sum_{j=1}^m \theta_p^j \right]^T \right\|} \quad (5)$$

The definition of $\vec{\Theta}^\mu$ is analogous to the main prediction of a learning algorithm [83]; it provides a reference for quantifying the reproducibility of an MBM:

Definition 4. Let $\vec{\Theta}^\mu$ be the main multivariate brain map of Φ . Then, let α^j be the angle between $\vec{\Theta}^j$ and $\vec{\Theta}^\mu$. The “reproducibility” ψ_Φ ($0 \leq \psi_\Phi \leq 1$) of an MBM derived from a linear function Φ is defined as follows:

$$\forall j \in \{1, \dots, m\}, \psi_\Phi = \mathbb{E}_S[\cos(\alpha^j)] \quad (6)$$

In fact, reproducibility provides a measure for quantifying the dispersion of MBMs, computed over different perturbed training sets, from the main multivariate brain map.

On the other hand, the coherency between the main multivariate brain map of a decoder and the true solution can be employed as a measure for the plausibility of a model. We refer to this coherency as the *representativeness* of an MBM:

Definition 5. Let $\vec{\Theta}^\mu$ be the main multivariate brain map of Φ . The “representativeness” ($0 \leq \beta_\Phi \leq 1$) is defined as the cosine similarity between $\vec{\Theta}^\mu$ and $\vec{\Theta}^*$:

$$\beta_\Phi = \frac{|\vec{\Theta}^\mu \cdot \vec{\Theta}^*|}{\|\vec{\Theta}^\mu\| \|\vec{\Theta}^*\|} \quad (7)$$

291 The following proposition shows the relationship between the presented
292 definitions for reproducibility, representativeness, and the interpretability:

293 **Proposition 1.** $\eta_{\Phi} = \beta_{\Phi} \times \psi_{\Phi}$.

294 See Appendix D for a proof. Proposition 1 indicates the interpretability
295 can be decomposed into the representativeness and the reproducibility of a
296 decoding model.

297 2.4. A Heuristic for Practical Quantification of Interpretability in Time- 298 Domain MEG decoding

299 In practice, it is impossible to evaluate the interpretability, as Φ^* is un-
300 known. In this study, to provide a practical proof of the mentioned theoret-
301 ical concepts, we propose the use of contrast event-related fields (cERFs) of
302 MEG data as neurophysiological plausible heuristics for Θ^* in a binary MEG
303 decoding scenario in the time domain.

304 The EEG/MEG data are a mixture of several simultaneous stimulus-
305 related and stimulus-unrelated brain activities. In general, unrelated-stimulus
306 brain activities are considered as Gaussian noise with zero mean and variance
307 σ^2 . One popular approach to canceling the noise component is to compute
308 the average of multiple trials. It is expected that the average will converge
309 to the true value of the signal with a variance of $\frac{\sigma^2}{n}$. The result of the av-
310 eraging process is generally known as ERF in the MEG context; separate
311 interpretation of different ERF components can be performed [84]¹.

312 Assume $\mathbf{X}^+ = \{x_i \in \mathbf{X} \mid y_i = 1\} \in \mathbb{R}^{n^+ \times p}$ and $\mathbf{X}^- = \{x_i \in \mathbf{X} \mid y_i =$
313 $-1\} \in \mathbb{R}^{n^- \times p}$. Then, the cERF brain map $\vec{\Theta}^{cERF}$ is computed as follows:

$$\vec{\Theta}^{cERF} = \frac{\frac{1}{n^+} \sum_{x_i \in X^+} x_i - \frac{1}{n^-} \sum_{x_i \in X^-} x_i}{\left\| \frac{1}{n^+} \sum_{x_i \in X^+} x_i - \frac{1}{n^-} \sum_{x_i \in X^-} x_i \right\|} \quad (8)$$

314 Using the core theory presented in [39], it can be shown that cERF is
315 the equivalent generative model for the least squares solution in a binary

¹The application of the presented heuristic to MEG data can be extended to EEG because of the inherent similarity of the measured neural correlates in these two devices. In the EEG context, the ERF can be replaced by the event-related potential (ERP).

time-domain MEG decoding scenario (see Appendix A). Using $\vec{\Theta}^{cERF}$ as a heuristic for $\vec{\Theta}^*$, the representativeness can be approximated as follows:

$$\tilde{\beta}_{\Phi} = \frac{|\vec{\Theta}^{\mu} \cdot \vec{\Theta}^{cERF}|}{\|\vec{\Theta}^{\mu}\| \|\vec{\Theta}^{cERF}\|} \quad (9)$$

Where $\tilde{\beta}_{\Phi}$ is an approximation of β_{Φ} and we have:

$$\beta_{\Phi} = \Delta_{\beta} \tilde{\beta}_{\Phi} \pm \sqrt{(1 - \tilde{\beta}_{\Phi}^2)(1 - \Delta_{\beta}^2)} \quad (10)$$

Δ_{β} represents the cosine similarity between $\vec{\Theta}^*$ and $\vec{\Theta}^{cERF}$ (see Figures B.8 and Appendix B). If $\Delta_{\beta} \rightarrow 1$ then $\tilde{\beta}_{\Phi} \rightarrow \beta_{\Phi}$.

In a similar manner, $\vec{\Theta}^{cERF}$ can be used to heuristically approximate the interpretability as follows:

$$\tilde{\eta}_{\Phi} = \forall j \in \{1, \dots, m\}, \tilde{\eta}_{\Phi} = \mathbb{E}_S(\cos(\gamma^j)) \quad (11)$$

where $\gamma_1, \dots, \gamma_m$ are the angles between $\vec{\Theta}^1, \dots, \vec{\Theta}^m$ and $\vec{\Theta}^{cERF}$. The following equality represents the relation between η and $\tilde{\eta}$ (see Figures C.9 and Appendix C).

$$\eta_{\Phi} = \Delta_{\beta} \tilde{\eta}_{\Phi} \pm \frac{\sqrt{1 - \Delta_{\beta}^2}}{m} (\sin \gamma_1 + \dots + \sin \gamma_m) \quad (12)$$

Again, if $\Delta_{\beta} \rightarrow 1$ then $\tilde{\eta}_{\Phi} \rightarrow \eta_{\Phi}$. Notice that Δ_{β} is independent of the decoding approach used; it only depends on the quality of the heuristic. It can be shown that $\tilde{\eta}_{\Phi} = \tilde{\beta}_{\Phi} \times \psi_{\Phi}$.

Eq. 12 shows that the choice of heuristic has a direct effect on the approximation of interpretability and that an inappropriate selection of the heuristic yields a very poor estimation of interpretability because of the destructive contribution of Δ_{β} . Therefore, the choice of heuristic should be carefully justified based on accepted and well-defined facts regarding the nature of the collected data (see the supplementary material for the experimental investigation of the limitations of the proposed heuristic).

2.5. Incorporating the Interpretability into Model Selection

The procedure for evaluating the performance of a model so as to choose the best values for hyper-parameters is known as *model selection* [85]. This procedure generally involves numerical optimization of the model selection criterion. The most common model selection criterion is based on an estimator of generalization performance, i.e., the predictive power. In the context of brain decoding, especially when the interpretability of brain maps matters, employing the predictive power as the only decisive criterion in model selection is problematic in terms of interpretability [86, 53, 54]. Here, we propose a multi-objective criterion for model selection that takes into account both prediction accuracy and MBM interpretability.

Let $\tilde{\eta}_{\Phi}$ and δ_{Φ} be the approximated interpretability and the generalization performance of a linear function Φ , respectively. We propose the use of the *scalarization* technique [87] for combining $\tilde{\eta}_{\Phi}$ and δ_{Φ} into one scalar $0 \leq \zeta(\Phi) \leq 1$ as follows:

$$\zeta_{\Phi} = \begin{cases} \frac{\omega_1 \tilde{\eta}_{\Phi} + \omega_2 \delta_{\Phi}}{\omega_1 + \omega_2} & \delta_{\Phi} \geq \kappa \\ 0 & \delta_{\Phi} < \kappa \end{cases} \quad (13)$$

where ω_1 and ω_2 are weights that specify the level of importance of the interpretability and the performance, respectively. κ is a threshold on the performance that filters out solutions with poor performance. In classification scenarios, κ can be set by adding a small safe interval to the chance level of classification.

It can be shown that the hyper-parameters of a model Φ are optimized based on ζ_{Φ} are Pareto optimal [88]. In other words, there exist no other Φ' for which we obtain both $\tilde{\eta}_{\Phi'} > \tilde{\eta}_{\Phi}$ and $\delta_{\Phi'} > \delta_{\Phi}$. We expect that optimizing the hyper-parameters based on ζ_{Φ} , rather only δ_{Φ} , yields more informative MBMs.

2.6. Experimental Materials

2.6.1. Toy Dataset

To illustrate the importance of integrating the interpretability of brain decoding with the model selection procedure, we use simple 2-dimensional toy data presented in [39]. Assume that the true underlying generative function Φ^* is defined by

$$\mathcal{Y} = \Phi^*(\mathcal{X}) = \begin{cases} 1 & \text{if } x_1 = 1.5 \\ -1 & \text{if } x_1 = -1.5 \end{cases}$$

where $\mathcal{X} \in \{[1.5, 0]^T, [-1.5, 0]^T\}$; and x_1 and x_2 represent the first and the second dimension of the data, respectively. Furthermore, assume the data is contaminated by Gaussian noise with co-variance $\Sigma = \begin{bmatrix} 1.02 & -0.3 \\ -0.3 & 0.15 \end{bmatrix}$. Figure 1 shows the distribution of the noisy data.

2.6.2. MEG Data

We use the MEG dataset presented in [89]¹. The dataset was also used for the DecMeg2014 competition². In this dataset, visual stimuli consisting of famous faces, unfamiliar faces, and scrambled faces are presented to 16 subjects and fMRI, EEG, and MEG signals are recorded. Here, we are only interested in MEG recordings. The MEG data were recorded using a VectorView system (Elekta Neuromag, Helsinki, Finland) with a magnetometer and two orthogonal planar gradiometers located at 102 positions in a hemispherical array in a light Elekta-Neuromag magnetically shielded room.

Three major reasons motivated the choice of this dataset: 1) It is publicly available. 2) The spatio-temporal dynamic of the MEG signal for face vs. scramble stimuli has been well studied. The event-related potential analysis of EEG/MEG shows that *N170* occurs 130 – 200ms after stimulus presentation and reflects the neural processing of faces [90, 89]. Therefore, the *N170* component can be considered the ground truth for our analysis. 3) In the literature, non-parametric mass-univariate analysis such as cluster-based permutation tests is unable to identify narrowly distributed effects in space and time (e.g., an *N170* component) [2, 6]. These facts motivate us to employ multivariate approaches that are more sensitive to these effects.

As in [51], we created a balanced face vs. scrambled MEG dataset by randomly drawing from the trials of unscrambled (famous or unfamiliar) faces and scrambled faces in equal number. The samples in the face and scrambled face categories are labeled as 1 and -1 , respectively. The raw data is high-pass filtered at 1Hz, down-sampled to 250Hz, and trimmed from 200ms before the stimulus onset to 800ms after the stimulus. Thus, each trial has 250 time-points for each of the 306 MEG sensors (102 magnetometers and

¹The full dataset is publicly available at ftp://ftp.mrc-cbu.cam.ac.uk/personal/rik.henson/wakemandg_hensonrn/

²The competition data are available at <http://www.kaggle.com/c/decoding-the-human-brain>

204 planar gradiometers)¹. To create the feature vector of each sample, we pooled all of the temporal data of 306 MEG sensors into one vector (i.e., we have $p = 250 \times 306 = 76500$ features for each sample). Before training the classifier, all of the features are standardized to have a mean of 0 and standard-deviation of 1.

2.7. Classification and Evaluation

In all experiments, a least squares classifier with L1-penalization, i.e., Lasso [75], is used for decoding. Lasso is a very popular classification method in the context of brain decoding, mainly because of its sparsity assumption. The choice of Lasso helps us to better illustrate the importance of including the interpretability in the model selection. Lasso solves the following optimization problem:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \|\Phi(\mathbf{X}) - \Phi_S(\mathbf{X})\|_2^2 + \lambda \|\Theta\|_1 \quad (14)$$

where λ is the hyper-parameter that specifies the level of regularization. Therefore, the aim of the model selection is to find the best value for λ . Here, we try to find the best regularization parameter value among $\lambda = \{0.001, 0.01, 0.1, 1, 10, 50, 100, 250, 500, 1000, 5000, 10000, 15000, 25000, 50000\}$.

We use the out-of-bag (OOB) [91, 92] method for computing δ_Φ , ψ_Φ , $\tilde{\beta}_\Phi$, $\tilde{\eta}_\Phi$, and ζ_Φ for different values of λ . In OOB, given a training set (\mathbf{X}, \mathbf{Y}) , m replications of bootstrap [81] are used to create perturbed training sets (we set $m = 50$)². In all of our experiments, we set $\omega_1 = \omega_2 = 1$ and $\kappa = 0.6$ in the computation of ζ_Φ . Furthermore, we set $\delta_\Phi = 1 - EPE$ where EPE indicates the expected prediction error; it is computed using the procedure explained in Appendix E. Employing OOB provides the possibility of computing the bias and variance of the model as contributing factors in EPE.

To investigate the behavior of the proposed model selection criterion, we benchmark it against the commonly used performance criterion in the single-subject decoding scenario. Assuming $(\mathbf{X}_i, \mathbf{Y}_i)$ for $i = 1, \dots, 16$ are MEG trial/label pairs for subject i , we separately train a Lasso model for

¹The preprocessing scripts in python and MATLAB are available at: <https://github.com/FBK-NILab/DecMeg2014/>

²The MATLAB code used for experiments is available at <https://github.com/smkia/interpretability/>

each subject to estimate the parameter of the linear function $\hat{\Phi}_i$, where $\mathbf{Y}_i = \mathbf{X}_i \hat{\Theta}_i$. Let $\hat{\Phi}_i^\delta$ and $\hat{\Phi}_i^\zeta$ represent the optimized solution based on δ_Φ and ζ_Φ , respectively. We denote the MBM associated with $\hat{\Phi}_i^\delta$ and $\hat{\Phi}_i^\zeta$ by $\vec{\Theta}_i^\delta$ and $\vec{\Theta}_i^\zeta$, respectively. Therefore, for each subject, we compare the resulting decoders and MBMs computed based on these two model selection criteria.

3. Results

3.1. Performance-Interpretability Dilemma: A Toy Example

In the definition of Φ^* on the toy dataset discussed in Section 2.6.1, x_1 is the decisive variable and x_2 has no effect on the classification of the data into target classes. Therefore, excluding the effect of noise and based on the theory of the maximal margin classifier [93, 94], $\vec{\Theta}^* \propto [1, 0]^T$ is the true solution to the decoding problem. By accounting for the effect of noise and solving the decoding problem in (\mathbf{X}, \mathbf{Y}) space, we have $\vec{\Theta} \propto [\frac{1}{\sqrt{(5)}}, \frac{2}{\sqrt{(5)}}]^T$ as the parameter of the linear classifier. Although the estimated parameters on the noisy data yield the best generalization performance for the noisy samples, any attempt to interpret this solution fails, as it yields the wrong conclusion with respect to the ground truth (it says x_2 has twice the influence of x_1 on the results, whereas it has no effect). This simple experiment shows that the most accurate model is not always the most interpretable one, primarily because the contribution of the noise in the decoding process [39]. On the other hand, the true solution of the problem $\vec{\Theta}^*$ does not provide the best generalization performance for the noisy data.

To illustrate the effect of incorporating the interpretability in the model selection, a Lasso model with different λ values is used for classifying the toy data. In this case, because $\vec{\Theta}^*$ is known, the exact value of interpretability can be computed using Eq. 4. Table 1 compares the resultant performance and interpretability from Lasso. Lasso achieves its highest performance ($\delta_\Phi = 0.9884$) at $\lambda = 10$ with $\vec{\Theta} \propto [0.4636, 0.8660]^T$ (indicated by the magenta line in Figure 1). Despite having the highest performance, this solution suffers from a lack of interpretability ($\eta_\Phi = 0.4484$). By increasing λ , the interpretability improves so that for $\lambda = 500, 1000$ the classifier reaches its highest interpretability by compensating for 0.06 of its performance. Our observation highlights two main points:

1. In the case of noisy data, the interpretability of a decoding model is incoherent with its performance. Thus, optimizing the parameter of

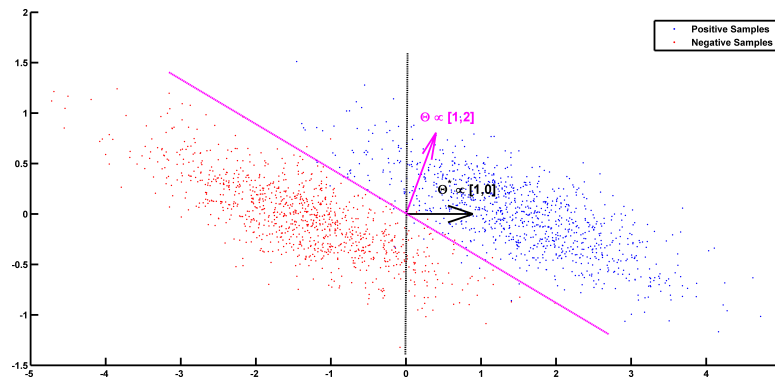


Figure 1: Noisy samples of toy data. The black line shows the true separator based on the generative model (Φ^*). The magenta line shows the most accurate classification solution. Because of the contribution of noise, any interpretation of the parameters of the most accurate classifier yields a misleading conclusion with respect to the true underlying phenomenon [39].

Table 1: Comparison between δ_Φ , η_Φ , and ζ_Φ for different λ values on the toy 2D example shows the performance-interpretability dilemma, in which the most accurate classifier is not the most interpretable one.

λ	0	0.001	0.01	0.1	1	10	50	100	250	500	1000
$\delta(\Phi)$	0.9883	0.9883	0.9883	0.9883	0.9883	0.9884	0.9880	0.9840	0.9310	0.9292	0.9292
$\eta(\Phi)$	0.4391	0.4391	0.4391	0.4392	0.4400	0.4484	0.4921	0.5845	0.9968	1	1
$\zeta(\Phi)$	0.7137	0.7137	0.7137	0.7137	0.7142	0.7184	0.7400	0.7842	0.9639	0.9646	0.9646
$\tilde{\Theta} \propto$	$\begin{bmatrix} 0.4520 \\ 0.8920 \end{bmatrix}$	$\begin{bmatrix} 0.4520 \\ 0.8920 \end{bmatrix}$	$\begin{bmatrix} 0.4520 \\ 0.8920 \end{bmatrix}$	$\begin{bmatrix} 0.4521 \\ 0.8919 \end{bmatrix}$	$\begin{bmatrix} 0.4532 \\ 0.8914 \end{bmatrix}$	$\begin{bmatrix} 0.4636 \\ 0.8660 \end{bmatrix}$	$\begin{bmatrix} 0.4883 \\ 0.8727 \end{bmatrix}$	$\begin{bmatrix} 0.5800 \\ 0.8146 \end{bmatrix}$	$\begin{bmatrix} 0.99 \\ 0.02 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

the model based on its performance does not necessarily improve its interpretability. This observation confirms the previous finding by Rasmussen et al. [53] regarding the trade-off between the spatial reproducibility (as a measure for the interpretability) and the prediction accuracy in brain decoding.

2. If the right criterion is used in the model selection, employing proper regularization technique (sparsity prior, in this case) leads to more interpretability for the decoding models.

3.2. Mass-Univariate Hypothesis Testing on MEG Data

Results show that non-parametric mass-univariate analysis is unable to detect narrowly distributed effects in space and time (e.g., an $N170$ component) [2, 6]. To illustrate the advantage of the proposed decoding framework

for spotting these effects, we performed a non-parametric cluster-based permutation test [5] on our MEG dataset using Fieldtrip toolbox [95]. In a single subject analysis scenario, we considered the trials of MEG recordings as the unit of observation in a between-trials experiment. Independent-samples t-statistics are used as the statistics for evaluating the effect at the sample level and to construct spatio-temporal clusters. The maximum of the cluster-level summed t-value is used for the cluster level statistics; the significance probability is computed using a Monte Carlo method. The minimum number of neighboring channels for computing the clusters is set to 2. Considering 0.025 as the two-sided threshold for testing the significance level and repeating the procedure separately for magnetometers and combined-gradiometers, no significant result is found for any of the 16 subjects. This result motivates the search for more sensitive (and, at the same time, more interpretable) alternatives for hypothesis testing.

3.3. Single-Subject Decoding on MEG Data

In this experiment, we aim to compare the multivariate brain maps of brain decoding models when δ_Φ and ζ_Φ are used as the criteria for model selection. Figure 2(a) represents the mean and standard-deviation of the performance and interpretability of Lasso across 16 subjects for different λ values. The performance and interpretability curves further illustrate the performance-interpretability dilemma in the single-subject decoding scenario in which increasing the performance delivers less interpretability. The average performance across subjects is improved when λ approaches 1, but on the other side, the reproducibility and the representativeness of models declines significantly [see Figure 2(b)].

One possible reason behind the performance-interpretability dilemma is illustrated in Figure 3. The figure shows the mean and standard deviation of bias, variance, and EPE of Lasso across 16 subjects. The plot proposes that the effect of variance is overwhelmed by bias in the computation of EPE, where the best performance (minimum EPE) at $\lambda = 1$ has the lowest bias, its variance is higher than for $\lambda = 0.001, 0.01, 0.1$. While this tiny increase in the variance is not reflected in EPE but Figure 2(b) shows a significant effect on the reproducibility.

Table 2 summarizes the performance, reproducibility, representativeness, and interpretability of $\hat{\Phi}_i^\delta$ and $\hat{\Phi}_i^\zeta$ for 16 subjects. The average result over 16 subjects shows that employing ζ_Φ instead of δ_Φ in model selection provides significantly higher reproducibility, representativeness, and (as a result)

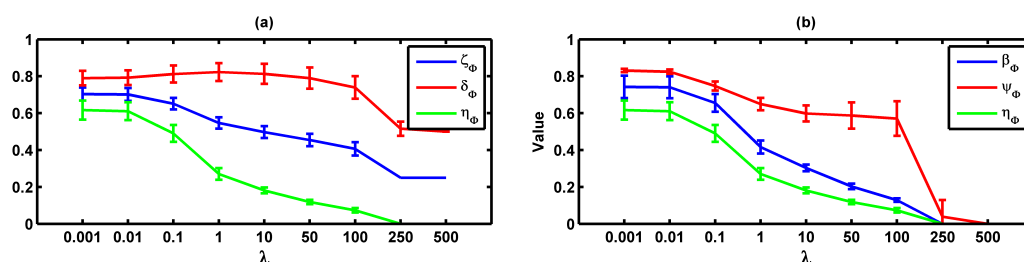


Figure 2: (a) Mean and standard-deviation of the performance, interpretability, and plausibility of Lasso over 16 subjects. The performance and interpretability become incoherent as λ increases. (b) Mean and standard-deviation of the reproducibility, representativeness, and interpretability of Lasso over 16 subjects. The interpretability declines because of the decrease in both reproducibility and representativeness.

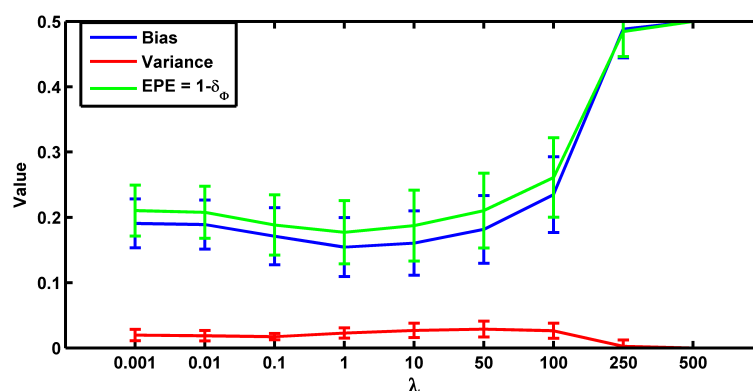


Figure 3: Mean and standard-deviation of the bias, variance, and EPE of Lasso over 16 subjects. The effect of variance on the EPE is overwhelmed by bias.

Table 2: The performance, reproducibility, representativeness, and interpretability of $\hat{\Phi}_i^\delta$ and $\hat{\Phi}_i^\zeta$ over 16 subjects.

Subj	Criterion: δ_Φ					Criterion: ζ_Φ				
	δ_Φ	ζ_Φ	$\tilde{\eta}_\Phi$	β_Φ	ψ_Φ	δ_Φ	ζ_Φ	$\tilde{\eta}_\Phi$	β_Φ	ψ_Φ
1	0.81	0.53	0.26	0.42	0.62	0.78	0.70	0.63	0.76	0.83
2	0.80	0.70	0.60	0.72	0.83	0.80	0.70	0.60	0.72	0.83
3	0.81	0.63	0.45	0.64	0.71	0.78	0.71	0.64	0.78	0.83
4	0.84	0.52	0.20	0.31	0.66	0.76	0.70	0.64	0.77	0.83
5	0.80	0.54	0.29	0.44	0.65	0.78	0.69	0.61	0.73	0.83
6	0.79	0.52	0.24	0.39	0.63	0.74	0.67	0.61	0.74	0.82
7	0.84	0.55	0.27	0.40	0.66	0.81	0.70	0.59	0.71	0.84
8	0.87	0.55	0.24	0.35	0.68	0.85	0.68	0.52	0.61	0.84
9	0.80	0.55	0.31	0.46	0.67	0.77	0.67	0.57	0.69	0.82
10	0.79	0.53	0.26	0.41	0.64	0.77	0.68	0.58	0.70	0.83
11	0.74	0.65	0.56	0.68	0.82	0.74	0.65	0.56	0.68	0.82
12	0.80	0.55	0.29	0.46	0.64	0.79	0.70	0.61	0.74	0.83
13	0.83	0.50	0.18	0.29	0.61	0.77	0.70	0.63	0.76	0.82
14	0.90	0.58	0.27	0.39	0.68	0.81	0.78	0.74	0.89	0.84
15	0.92	0.63	0.34	0.48	0.71	0.89	0.78	0.66	0.77	0.86
16	0.87	0.55	0.23	0.37	0.62	0.81	0.74	0.67	0.81	0.83
Mean	0.83±0.05	0.57 ± 0.05	0.31 ± 0.12	0.45 ± 0.13	0.68 ± 0.07	0.79 ± 0.04	0.70±0.04	0.62±0.05	0.74±0.06	0.83±0.01

interpretability compensating for 0.04 of performance.

These results are further analyzed in Figure 4 where $\hat{\Phi}_i^\delta$ and $\hat{\Phi}_i^\zeta$ are compared subject-wise in terms of their performance and interpretability. The comparison shows that adopting ζ_Φ instead of δ_Φ as the criterion for model selection yields significantly better interpretable models by compensating a negligible degree of performance in 14 out of 16 subjects. Figure 4(a) shows that employing δ_Φ provides on average slightly higher accurate models (Wilcoxon rank sum test p-value= 0.012) across subjects (0.83 ± 0.05) than using ζ_Φ (0.79 ± 0.04). On the other side, Figure 4(b) shows that employing ζ_Φ and compensating by 0.04 in the performance provides (on average) substantially higher (Wilcoxon rank sum test p-value= 5.6×10^{-6}) interpretability across subjects (0.62 ± 0.05) compared to δ_Φ (0.31 ± 0.12). For example, in the case of subject 1 (see table 2), using δ_Φ in model selection to select the best λ value for the Lasso yields a model with $\delta_\Phi = 0.81$ and $\tilde{\eta}_\Phi = 0.26$. In contrast, using ζ_Φ delivers a model with $\delta_\Phi = 0.78$ and $\tilde{\eta}_\Phi = 0.63$.

The advantage of the exchange between the performance and the interpretability can be seen in the quality of MBMs. Figure 5a and 5b show $\tilde{\Theta}_1^\delta$ and $\tilde{\Theta}_1^\zeta$ of subject 1, i.e., the spatio-temporal multivariate maps of the Lasso models with maximum values of δ_Φ and ζ_Φ , respectively. The maps are plotted for 102 magnetometer sensors. In each case, the time course of weights of classifiers associated with the MEG2041 and MEG1931 sensors are plotted. Furthermore, the topographic maps represent the spatial pat-

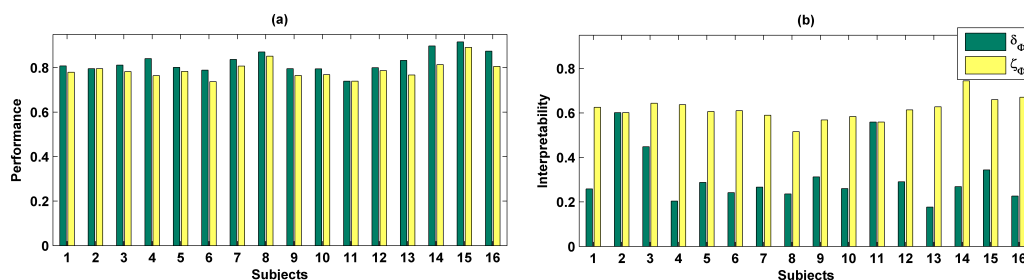
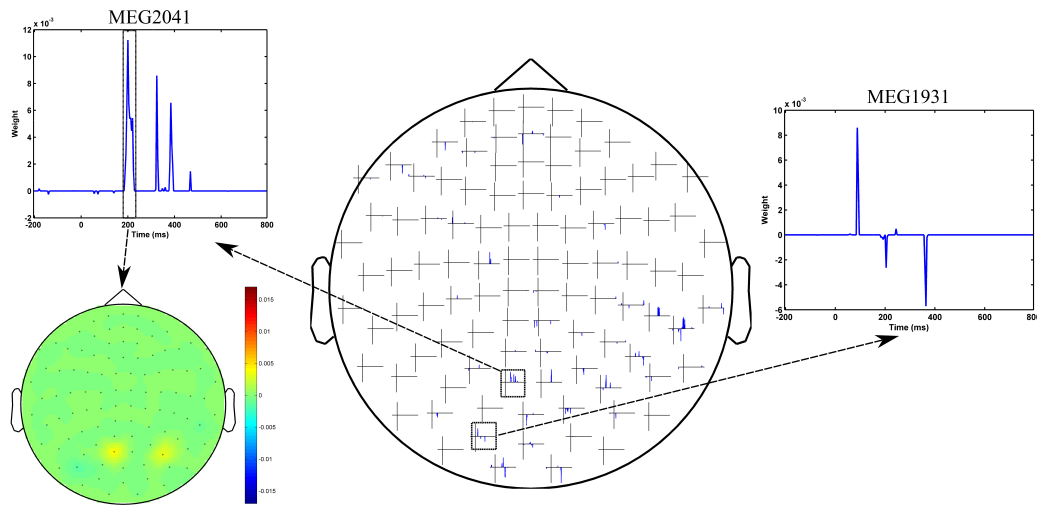


Figure 4: a) Comparison between performance of $\hat{\Phi}_i^\delta$ and $\hat{\Phi}_i^\zeta$. Adopting ζ_Φ instead of δ_Φ in model selection yields (on average) 0.04 less accurate classifiers over 16 subjects. b) Comparison between interpretability of $\hat{\Phi}_i^\delta$ and $\hat{\Phi}_i^\zeta$. Adopting ζ_Φ instead of δ_Φ in model selection yields on average 0.31 more interpretable classifiers over 16 subjects.

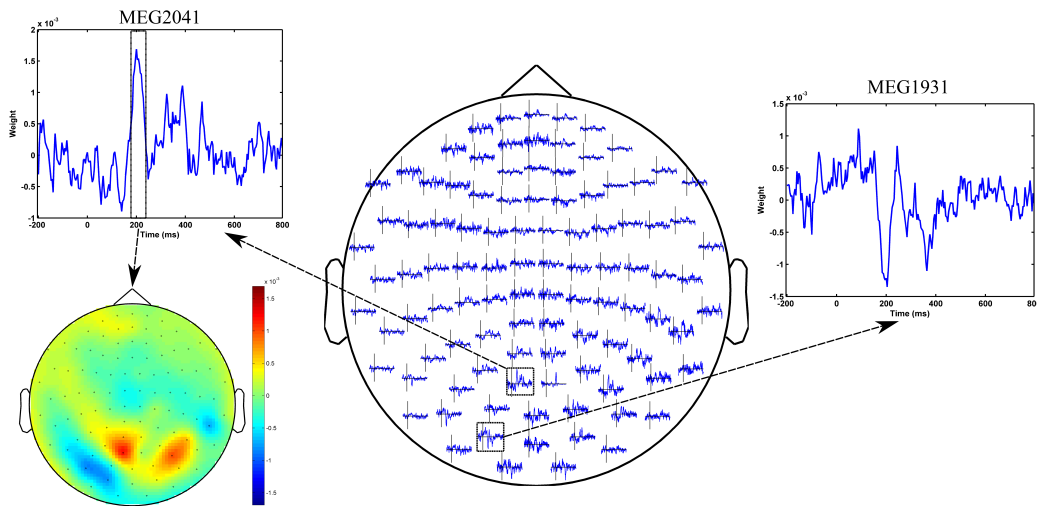
terns of weights averaged between 184ms and 236ms after stimulus onset¹. While $\vec{\Theta}_1^\delta$ is sparse in time and space, it fails to accurately represent the spatio-temporal dynamic of the N170 component. Furthermore, the multicollinearity problem arising from the correlation between the time course of the MEG2041 and MEG1931 sensors causes extra attenuation of the N170 effect in the MEG1931 sensor. Therefore, the model is unable to capture the spatial pattern of the dipole in the posterior area. In contrast, $\vec{\Theta}_1^\zeta$ represents the dynamic of the N170 component in time (see Figure 6). In addition, it also shows the spatial pattern of two dipoles in the posterior and temporal areas. In summary, $\vec{\Theta}_1^\zeta$ suggests a more representative pattern of the underlying neurophysiological effect than $\vec{\Theta}_1^\delta$.

In addition, optimizing the hyper-parameters of brain decoding based on ζ_Φ offers more reproducible brain decoders. According to table 2, using ζ_Φ instead of δ_Φ provides (on average) 0.15 more reproducibility over 16 subjects. To illustrate the advantage of higher reproducibility on the interpretability of maps, Figure 7 visualizes $\vec{\Theta}_1^\delta$ and $\vec{\Theta}_1^\zeta$ over 4 perturbed training sets. The spatial maps [Figure 7(a) and Figure 7(c)] are plotted for the magnetometer sensors averaged in the time interval between 184ms and 236ms after stimulus onset. The temporal maps [Figure 7(b) and Figure 7(d)] are showing

¹The bounds of colorbars are symmetrized based on the maximum absolute value of parameters



(a) Spatio-temporal pattern of $\vec{\Theta}_1^\delta$.



(b) Spatio-temporal pattern of $\vec{\Theta}_1^\zeta$.

Figure 5: Comparison between spatio-temporal multivariate maps of the most accurate (5a) and the most interpretable (5b) classifiers for Subject 1. $\vec{\Theta}_1^\zeta$ provides more spatio-temporal representativeness of the N170 effect than $\vec{\Theta}_1^\delta$.

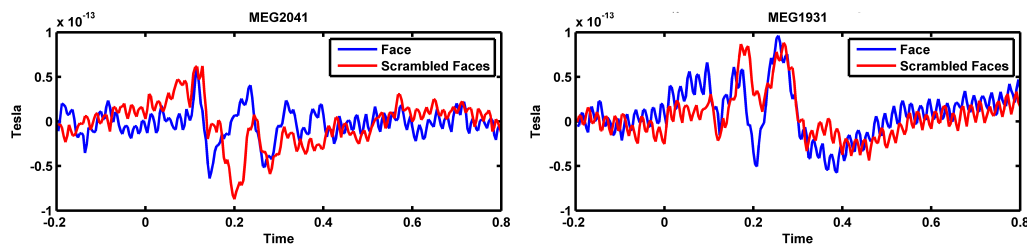


Figure 6: Event related fields (ERFs) of face and scrambled face samples for MEG2041 and MEG1931 sensors.

the multivariate temporal maps of MEG1931 and MEG2041 sensors. While $\hat{\Theta}_1^\delta$ is unstable in time and space across the 4 perturbed training sets, $\hat{\Theta}_1^\zeta$ provides more reproducible maps.

4. Discussions

4.1. Defining Interpretability: Theoretical Advantages

An overview of the brain decoding literature shows frequent co-occurrence of the terms interpretation, interpretable, and interpretability with the terms model, classification, parameter, decoding, method, feature, and pattern (see the quick meta-analysis on the literature in the supplementary material); however, a formal formulation of the interpretability is never presented. In this study, our primary interest is to present a theoretical definition of the interpretability of linear brain decoding models and their corresponding MBMs. Furthermore, we show the way in which interpretability is related to the reproducibility and neurophysiological representativeness of MBMs. Our definition and quantification of interpretability remains theoretical, as we assume that the true solution of the brain decoding problem is available. Despite this limitation, we argue that the presented definition provides a concrete framework of a previously abstract concept and that it establishes a theoretical background to explain an ambiguous phenomenon in the brain decoding context. We support our argument using an example in time-domain MEG decoding in which we show how the presented definition can be exploited to heuristically approximate the interpretability. This example shows how

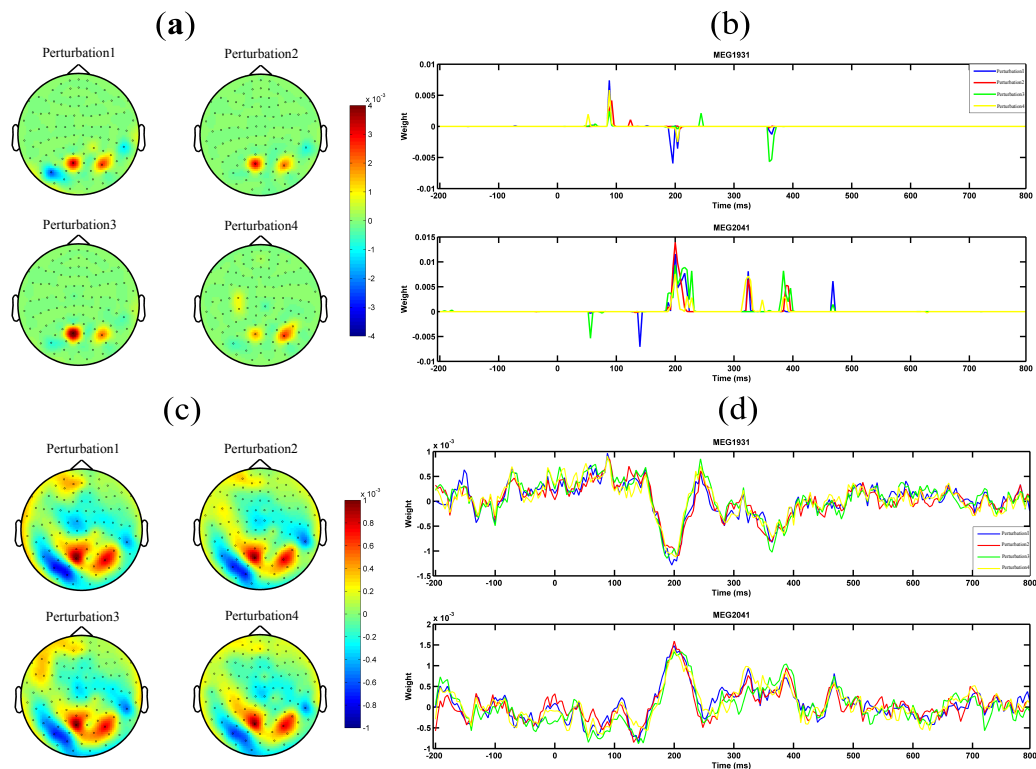


Figure 7: Comparison of the reproducibility of Lasso when δ_Φ and ζ_Φ are used in the model selection procedure. (a) and (b) show the spatio-temporal patterns represented by $\vec{\Theta}_1^\delta$ across the 4 perturbed training sets. (c) and (d) show the spatio-temporal patterns represented by $\vec{\Theta}_1^\zeta$ across the 4 perturbed training sets. Employing ζ_Φ instead of δ_Φ in the model selection yields more reproducible MBMs.

573 partial prior knowledge¹ regarding underlying brain activity can be used to
 574 find more plausible multivariate patterns in data. Furthermore, the proposed
 575 decomposition of the interpretability of MBMs into their reproducibility and
 576 representativeness explains the relationship between the influential coopera-
 577 tive factors in the interpretability of brain decoding models and highlights the
 578 possibility of indirect and partial evaluation of interpretability by measuring
 579 these effective factors.

580 *4.2. Application in Model Evaluation*

581 Discriminative models in the framework of brain decoding provide higher
 582 sensitivity and specificity than univariate analysis in hypothesis testing of
 583 neuroimaging data. Although multivariate hypothesis testing is performed
 584 based solely on the generalization performance of classifiers, the emergent
 585 need for extracting reliable complementary information regarding the un-
 586 derlying neuronal activity motivated a considerable amount of research on
 587 improving and assessing the interpretability of classifiers and their associated
 588 MBMs. Despite ubiquitous use, the generalization performance of classifiers
 589 is not a reliable criterion for assessing the interpretability of brain decoding
 590 models [53]. Therefore, considering extra criteria might be required. How-
 591 ever, because of the lack of a formal definition for interpretability, different
 592 characteristics of brain decoding models are considered as the main objec-
 593 tive in improving their interpretability. Reproducibility [53, 54], stability
 594 selection [7, 47, 69], sparsity [96], and neurophysiological plausibility [97] are
 595 examples of related criteria.

596 Our definition of interpretability helped us to fill this gap by introducing
 597 a new multi-objective model selection criterion as a weighted compromise be-
 598 tween interpretability and generalization performance of linear models. Our
 599 experimental results on single-subject decoding showed that adopting the
 600 new criterion for optimizing the hyper-parameters of brain decoding models
 601 is an important step toward reliable visualization of learned models from
 602 neuroimaging data. It is not the first time in the neuroimaging context that
 603 a new metric is proposed in combination with generalization performance for
 604 the model selection. Several recent studies proposed the combination of the
 605 reproducibility of the maps [53, 54, 43] or the stability of the classifiers [56, 57]

¹The partial knowledge can be based on already known facts regarding the timing and location of neural activity.

with the performance of discriminative models to enhance the interpretability of decoding models. Our definition of interpretability supports the claim that the reproducibility is not the only effective factor in interpretability. Therefore, our contribution can be considered a complementary effort with respect to the state of the art of improving the interpretability of brain decoding at the model selection level.

Furthermore, this work presents an effective approach for evaluating the quality of different regularization strategies for improving the interpretability of MBMs. As briefly reviewed in Section 1, there is a trend in research within the brain decoding context in which prior knowledge is injected into the penalization term as a technique to improve the interpretability of decoding models. Thus far, in the literature, there is no ad-hoc method to compare these different methods. Our findings provide a further step toward direct evaluation of interpretability of the currently proposed penalization strategies. Such an evaluation can highlight the advantages and disadvantages of applying different strategies on different data types and facilitates the choice of appropriate methods for a certain application.

4.3. Regularization and Interpretability

Haufe et al. [39] demonstrated that the weight in linear discriminative models are unable to accurately assess the relationship between independent variables, primarily because of the contribution of noise in the decoding process. The problem is primarily caused by the decoding process that minimizes the classification error only considering the uncertainty in the output space [80, 98, 99] and not the uncertainty in the input space (or noise). The authors concluded that the interpretability of brain decoding cannot be improved using regularization. Our experimental results on the toy data (see Section 3.1) shows that if the right criterion is used for selecting the best values for hyper-parameters, appropriate choice of the regularization strategy can still play significant role in improving the interpretability of results. For example, in this case, the true generative function behind the sampled data is sparse (see Section 2.6.1), but because of the noise in the data, the sparse model is not the most accurate one. Using a more comprehensive criterion (in this case, ζ_Φ) shows the advantage of selecting correct prior assumptions about the distribution of the data via regularization. This observation encourages the modification of the conclusion in [39] as follows: if the performance of the model is the only criterion in the model selection, then the interpretability cannot necessarily be improved by means of regularization.

643 4.4. *Advantage over Mass-Univariate Analysis*

644 Mass-univariate hypothesis testing methods are among the most popular
 645 tools in neuroscience research because they provide significance checks and
 646 a fair level of interpretability via univariate brain maps. Mass-univariate
 647 analyses consist of univariate statistical tests on single independent variables
 648 followed by multiple comparison correction. Generally, multiple compari-
 649 son correction reduces the sensitivity of mass-univariate approaches because
 650 of the large number of univariate tests involved. Cluster-based permuta-
 651 tion testing [5] provides a more sensitive univariate analysis framework by
 652 making the cluster assumption in the multiple comparison correction. Un-
 653 fortunately, this method is not able to detect narrow spatio-temporal effects
 654 in the data [2]. As a remedy, brain decoding provides a very sensitive tool
 655 for hypothesis testing; it has the ability to detect multivariate patterns, but
 656 suffers from a low level of interpretability. Our study proposes a possible
 657 solution for the interpretability problem of classifiers, and therefore, it facili-
 658 tates the application of brain decoding in the analysis of neuroimaging data.
 659 Our experimental results for the MEG data demonstrate that, although the
 660 non-parametric cluster-based permutation test is unable to detect the N170
 661 effect in MEG data, employing ζ_{Φ} instead of δ_{Φ} in model selection not only
 662 detects the stimuli-relevant information in the data, but also assures both
 663 reproducible and representative spatio-temporal mapping of the timing and
 664 the location of underlying neurophysiological effect.

665 4.5. *Limitations and Future Directions*

666 Despite theoretical and practical advantages, the proposed definition and
 667 quantification of interpretability suffer from some limitations. All of the
 668 presented concepts are defined for linear models, with the main assumption
 669 that $\Phi^* \in \mathcal{H}$ (where \mathcal{H} is a class of linear functions). This fact highlights
 670 the importance of linearizing the experimental protocol in the data collection
 671 phase [27]. Extending the definition of interpretability to non-linear models
 672 demands future research into the visualization of non-linear models in the
 673 form of brain maps. Currently, our findings cannot be directly applied to
 674 non-linear models. Furthermore, the proposed heuristic for the time-domain
 675 MEG data applies only to binary classification. One possible solution in mul-
 676 ticlass classification is to separate the decoding problem into several binary
 677 sub-problems. In addition the quality of the proposed heuristic is limited for
 678 the small sample size datasets (see supplementary material). Finding phys-

679 iologically relevant heuristics for other acquisition modalities such as fMRI
680 can be also considered in future work.

681 5. Conclusions

682 We presented a novel theoretical definition for the interpretability of linear
683 brain decoding and associated multivariate brain maps. We demonstrated
684 how the interpretability relates to the representativeness and reproducibility
685 of brain decoding. Although it is theoretical, the presented definition pro-
686 vides a first step toward practical solution for filling the knowledge extraction
687 gap in linear brain decoding. As an example of this major breakthrough,
688 and to provide a proof of concept, a heuristic approach based on the contrast
689 event-related field is proposed for practical evaluation of the interpretability
690 in time-domain MEG decoding. We experimentally showed that adding the
691 interpretability of brain decoding models as a criterion in the model selec-
692 tion procedure yields significantly higher interpretable models by sacrificing
693 a negligible amount of performance. Our methodological and experimental
694 achievements can be considered a complementary theoretical and practical
695 effort that contributes to researches on enhancing the interpretability of mul-
696 tivariate pattern analysis.

697 Acknowledgments

698 The author wishes to thank Sandro Vega-Pons and Nathan Weisz for
699 valuable discussions and comments.

700 Appendix A. cERF and its Generative Nature

701 According to [39], for a linear discriminative model with parameters Θ ,
702 the unique equivalent generative model can be computed as follows:

$$A \propto \Sigma_{\mathbf{X}} \Theta \quad (\text{A.1})$$

703 In a binary ($\mathbf{Y} = \{1, -1\}$) least squares classification scenario, we have:

$$A \propto \Sigma_{\mathbf{X}} \Sigma_{\mathbf{X}}^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{X}^T \mathbf{Y} = \mu^+ - \mu^- \quad (\text{A.2})$$

704 where $\Sigma_{\mathbf{X}}$ represents the covariance of the input matrix \mathbf{X} , and μ^+ and μ^-
705 are the means of positive and negative samples, respectively. Therefore,

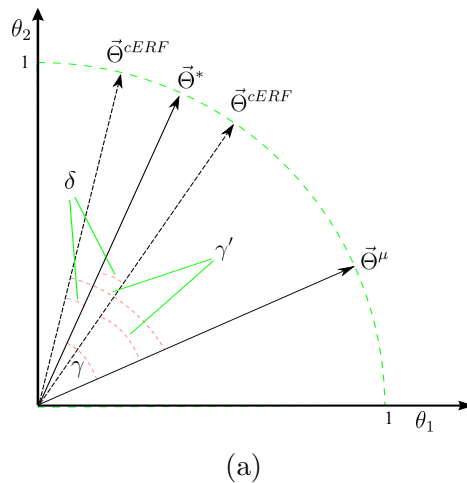


Figure B.8: Misrepresentation of $\vec{\Theta}^{cERF}$ with respect to $\vec{\Theta}^*$.

the equivalent generative model for the above classification problem can be derived by computing the difference between the mean of samples in two classes that is equivalent to the definition of cERF in time-domain MEG data.

Appendix B. Relation between β_Φ and $\tilde{\beta}_\Phi$ (Eq. 10)

Let γ be the angle between $\vec{\Theta}^\mu$ and $\vec{\Theta}^*$. Let γ' be the angle between $\vec{\Theta}^\mu$ and $\vec{\Theta}^{cERF}$. Furthermore, assume that δ is the angle between $\vec{\Theta}^*$ and $\vec{\Theta}^{cERF}$ and that $\Delta_\beta = \cos(\delta)$. We consider both cases in which β_Φ is underestimated/overestimated by $\tilde{\beta}_\Phi$ (see Figure B.8 as an example in 2-dimensional space). Then, we have:

$$\begin{aligned} \gamma = \gamma' \pm \delta &\Rightarrow \cos(\gamma) = \cos(\gamma' \pm \delta) \\ &= \cos(\gamma') \cos(\delta) \pm \sin(\gamma') \sin(\delta) = \tilde{\beta}_\Phi \Delta_\beta \pm \sqrt{(1 - \tilde{\beta}_\Phi^2)(1 - \Delta_\beta^2)} \end{aligned} \quad (\text{B.1})$$

Appendix C. Relation between η_Φ and $\tilde{\eta}_\Phi$ (Eq. 12)

Let $\alpha_1, \dots, \alpha_m$ be the angles between $\vec{\Theta}^1, \dots, \vec{\Theta}^m$ and $\vec{\Theta}^*$, and $\gamma_1, \dots, \gamma_m$ be the angles between $\vec{\Theta}^1, \dots, \vec{\Theta}^m$ and $\vec{\Theta}^{cERF}$. Furthermore, assume that

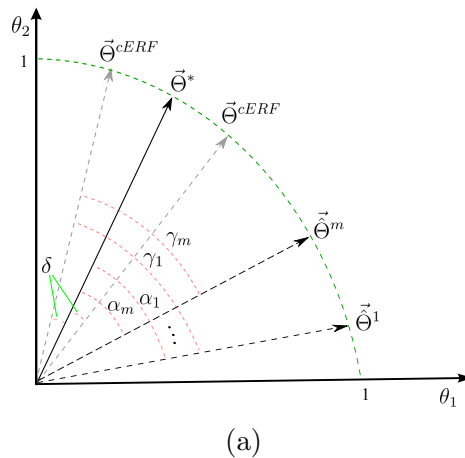


Figure C.9: Relation between η_{Φ} and $\tilde{\eta}_{\Phi}$.

δ is the angle between $\vec{\Theta}^*$ and $\vec{\Theta}^{cERF}$. We consider both cases in which η_{Φ} is underestimated/overestimated by $\tilde{\eta}_{\Phi}$ (see Figure C.9 as an example in 2-dimensional space).

$$\begin{aligned}
 \eta_{\Phi} &= \frac{\cos(\alpha_1) + \dots + \cos(\alpha_m)}{m} = \frac{\cos(\gamma_1 \pm \delta) + \dots + \cos(\gamma_m \pm \delta)}{m} \\
 &= \frac{\cos(\gamma_1) \cos(\delta) \pm \sin(\gamma_1) \sin(\delta) + \dots + \cos(\gamma_m) \cos(\delta) \pm \sin(\gamma_m) \sin(\delta)}{m} \\
 &\xrightarrow{\Delta_{\beta} = \cos(\delta)} = \frac{\Delta_{\beta} [\cos(\gamma_1) + \dots + \cos(\gamma_m)] \pm \sin(\delta) [\sin(\gamma_1) + \dots + \sin(\gamma_m)]}{m} \quad (C.1) \\
 &\xrightarrow{\tilde{\eta}_{\Phi} = \frac{\cos(\gamma_1) + \dots + \cos(\gamma_m)}{m}} \eta_{\Phi} = \Delta_{\beta} \tilde{\eta}_{\Phi} \pm \frac{\sqrt{1 - \Delta_{\beta}^2}}{m} (\sin(\gamma_1) + \dots + \sin(\gamma_m))
 \end{aligned}$$

Appendix D. Proof of Proposition 1

Throughout this proof, we assume that all of the parameter vectors are normalized in the unit hypersphere (see Figure D.10 as an illustrative example in 2 dimensions). Let $T = \{\vec{\Theta}^1, \dots, \vec{\Theta}^m\}$ be a set m MBMs, for m perturbed training sets where $\vec{\Theta}^i \in \mathbb{R}^p$. Now, consider any arbitrary $p - 1$ -dimensional hyperplane \mathcal{A} that contains $\vec{\Theta}^{\mu}$. Clearly, \mathcal{A} divides the p -dimensional parameter space into 2 subspaces. Let ∇ and \blacktriangledown be binary

operators where $\vec{\Theta}^i \nabla \vec{\Theta}^k$ indicates that $\vec{\Theta}^i$ and $\vec{\Theta}^k$ are in the same subspace, and $\vec{\Theta}^i \blacktriangledown \vec{\Theta}^k$ indicates that they are in different subspaces. Now, we define $T_U = \{\vec{\Theta}^i \mid \vec{\Theta}^i \nabla \vec{\Theta}^*\}$ and $T_L = \{\vec{\Theta}^i \mid \vec{\Theta}^i \blacktriangledown \vec{\Theta}^*\}$. Let the cardinality of T_L denoted by $n(T_L)$ be j ($n(T_L) = j$). Thus, $n(T_U) = m - j$. Now, assume that $\angle(\vec{\Theta}^i, \mathcal{A}) = \alpha_1, \dots, \alpha_j$ are the angles between $\vec{\Theta}^i \in T_L$ and \mathcal{A} , and (similarly) $\alpha_{j+1}, \dots, \alpha_m$ for $\vec{\Theta}^i \in T_U$ and \mathcal{A} . Based on Eq. 5, let $\vec{\Theta}_L^\mu$ and $\vec{\Theta}_U^\mu$ be the main maps of T_L and T_U , respectively. Therefore, we obtain $\vec{\Theta}^\mu = \frac{\vec{\Theta}_L^\mu + \vec{\Theta}_U^\mu}{\|\vec{\Theta}_L^\mu + \vec{\Theta}_U^\mu\|}$ and $\angle(\vec{\Theta}_L^\mu, \mathcal{A}) = \angle(\vec{\Theta}_U^\mu, \mathcal{A}) = \alpha$. Furthermore, assume $\angle(\vec{\Theta}^*, \mathcal{A}) = \gamma$. As a result, $\psi_\Phi = \cos(\alpha)$ and $\beta_\Phi = \cos(\gamma)$. According to Eq. 4 and using a cosine similarity definition, we have:

$$\begin{aligned} \eta_\Phi &= \frac{1}{m} \sum_{j=1}^m |\vec{\Theta}^* \cdot \vec{\Theta}^j| \\ &= \frac{\cos(\gamma + \alpha_1) + \dots + \cos(\gamma + \alpha_j) + \cos(\gamma - \alpha_{j+1}) + \dots + \cos(\gamma - \alpha_m)}{m} \\ &= \frac{\cos(\gamma + \alpha) + \cos(\gamma - \alpha)}{2} \\ &= \frac{\cos(\gamma) \cos(\alpha) - \sin(\gamma) \sin(\alpha) + \cos(\gamma) \cos(\alpha) + \sin(\gamma) \sin(\alpha)}{2} \\ &= \cos(\gamma) \cos(\alpha) = \beta_\Phi \times \psi_\Phi. \end{aligned} \tag{D.1}$$

A similar procedure can be used to prove $\tilde{\eta}_\Phi = \tilde{\beta}_\Phi \times \psi_\Phi$ by replacing $\vec{\Theta}^*$ with $\vec{\Theta}^{cERF}$.

Appendix E. Computing the Bias and Variance in Binary Classification

Here, using the out-of-bag (OOB) technique, and based on procedures proposed by [83] and [100], we compute the expected prediction error (EPE) for a linear binary classifier Φ under bootstrap perturbation of the training set. Let m be the number of perturbed training sets resulting from partitioning (X, Y) into (X_{tr}, Y_{tr}) and (X_{ts}, Y_{ts}) , i.e., training and test sets. If $\hat{\Phi}^j$ is the linear classifier estimated from the j th perturbed training set, then the main prediction $\Phi^\mu(\mathbf{x}_i)$ for each sample in the dataset can be computed as follows:

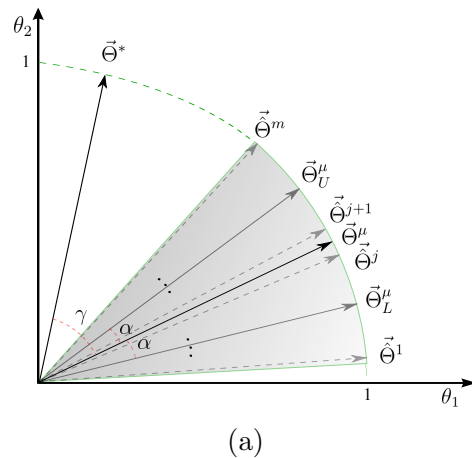


Figure D.10: Relation between representativeness, reproducibility, and interpretability in 2 dimensions.

$$\Phi^\mu(\mathbf{x}_i) = \begin{cases} 1 & \text{if } \frac{1}{k_i} \sum_{j=1}^{k_i} \hat{\Phi}^j(\mathbf{x}_i) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases} \quad (\text{E.1})$$

751 where k_i is the number of times that x_i is present in the test set^{1.1}

752 The computation of bias is challenging because the optimal model Φ^*
753 is unknown. According to [101], misclassification error is one of the loss
754 measures that satisfies a Pythagorean-type equality, and:

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\Phi^\mu(\mathbf{x}_i), \Phi^*(\mathbf{x}_i)) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \Phi^\mu(\mathbf{x}_i)) - \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \Phi^*(\mathbf{x}_i)) \quad (\text{E.2})$$

755 Because all terms of the above equation are positive, the mean loss be-
756 tween the main prediction and the actual labels can be considered as an
757 upper-bound for the bias:

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\Phi^\mu(\mathbf{x}_i), \Phi^*(\mathbf{x}_i)) \leq \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \Phi^\mu(\mathbf{x}_i)) \quad (\text{E.3})$$

¹It is expected that each sample $\mathbf{x}_i \in X$ appears (on average) $k_i \approx \frac{m}{3}$ times in the test sets.

Therefore, a pessimistic approximation of bias $B(\mathbf{x}_i)$ can be calculated as follows:

$$B(\mathbf{x}_i) = \begin{cases} 0 & \text{if } \Phi^\mu(\mathbf{x}_i) = y_i \\ 1 & \text{otherwise} \end{cases} \quad (\text{E.4})$$

Then, the unbiased and biased variances (see [83] for definitions) in each training set can be calculated by:

$$V_u^j(\mathbf{x}_i) = \begin{cases} 1 & \text{if } B(\mathbf{x}_i) = 0 \text{ and } \Phi^\mu(\mathbf{x}_i) \neq \hat{\Phi}^j(\mathbf{x}_i) \\ 0 & \text{otherwise} \end{cases} \quad (\text{E.5})$$

$$V_b^j(\mathbf{x}_i) = \begin{cases} 1 & \text{if } B(\mathbf{x}_i) = 1 \text{ and } \Phi^\mu(\mathbf{x}_i) \neq \hat{\Phi}^j(\mathbf{x}_i) \\ 0 & \text{otherwise} \end{cases} \quad (\text{E.6})$$

Then, the expected prediction error of Φ can be computed as follows (ignoring the irreducible error):

$$\begin{aligned} EPE_\Phi(X) &= \underbrace{\frac{1}{n} \sum_{i=1}^n B(\mathbf{x}_i)}_{\text{Bias}} + \\ &\underbrace{\frac{1}{nm} \sum_{j=1}^m \sum_{i=1}^n [V_u^j(\mathbf{x}_i) - V_b^j(\mathbf{x}_i)]}_{\text{Variance}} \end{aligned} \quad (\text{E.7})$$

References

- [1] E. Crivellato, D. Ribatti, Soul, mind, brain: Greek philosophy and the birth of neuroscience, Brain research bulletin 71 (2007) 327–336.
- [2] D. M. Groppe, T. P. Urbach, M. Kutas, Mass univariate analysis of event-related brain potentials/fields i: A critical tutorial review, Psychophysiology 48 (2011) 1711–1725.
- [3] E. Maris, Statistical testing in electrophysiological studies, Psychophysiology 49 (2012) 549–565.

- 772 [4] E. Bullmore, M. Brammer, S. C. Williams, S. Rabe-Hesketh, N. Janot,
773 A. David, J. Mellers, R. Howard, P. Sham, Statistical methods of esti-
774 mation and inference for functional mr image analysis, *Magnetic Resonance*
775 *in Medicine* 35 (1996) 261–277.
- 776 [5] E. Maris, R. Oostenveld, Nonparametric statistical testing of eeg-and meg-
777 data, *Journal of neuroscience methods* 164 (2007) 177–190.
- 778 [6] D. M. Groppe, T. P. Urbach, M. Kutas, Mass univariate analysis of event-
779 related brain potentials/fields ii: Simulation studies, *Psychophysiology* 48
780 (2011) 1726–1737.
- 781 [7] M. van Gerven, C. Hesse, O. Jensen, T. Heskes, Interpreting single trial data
782 using groupwise regularisation, *NeuroImage* 46 (2009) 665–676.
- 783 [8] T. Davis, K. F. LaRocque, J. A. Mumford, K. A. Norman, A. D. Wagner,
784 R. A. Poldrack, What do differences between multi-voxel and univariate
785 analysis mean? how subject-, voxel-, and trial-level variance impact fmri
786 analysis, *NeuroImage* 97 (2014) 271–283.
- 787 [9] J.-D. Haynes, G. Rees, Decoding mental states from brain activity in hu-
788 mans, *Nature Reviews Neuroscience* 7 (2006) 523–534.
- 789 [10] J. R. Wolpaw, N. Birbaumer, D. J. McFarland, G. Pfurtscheller, T. M.
790 Vaughan, Brain–computer interfaces for communication and control, *Clinical*
791 *neurophysiology* 113 (2002) 767–791.
- 792 [11] L. F. Nicolas-Alonso, J. Gomez-Gil, Brain computer interfaces, a review,
793 *Sensors* 12 (2012) 1211–1279.
- 794 [12] D. Bzdok, Classical statistics and statistical learning in imaging neuro-
795 science, *arXiv preprint arXiv:1603.01857* (2016).
- 796 [13] F. Pereira, T. Mitchell, M. Botvinick, Machine learning classifiers and fMRI:
797 a tutorial overview., *NeuroImage* 45 (2009) 199–209.
- 798 [14] S. Lemm, B. Blankertz, T. Dickhaus, K.-R. Müller, Introduction to machine
799 learning for brain imaging, *Neuroimage* 56 (2011) 387–399.
- 800 [15] M. Besserve, K. Jerbi, F. Laurent, S. Baillet, J. Martinerie, L. Garnero,
801 Classification methods for ongoing eeg and meg signals, *Biological research*
802 40 (2007) 415–437.

- 803 [16] J. V. Haxby, M. I. Gobbini, M. L. Furey, A. Ishai, J. L. Schouten, P. Pietrini,
804 Distributed and Overlapping Representations of Faces and Objects in Ven-
805 tral Temporal Cortex, *Science* 293 (2001) 2425–2430.
- 806 [17] D. D. Cox, R. L. Savoy, Functional magnetic resonance imaging (fmri) brain
807 reading: detecting and classifying distributed patterns of fmri activity in
808 human visual cortex, *Neuroimage* 19 (2003) 261–270.
- 809 [18] T. M. Mitchell, R. Hutchinson, R. S. Niculescu, F. Pereira, X. Wang, M. Just,
810 S. Newman, Learning to decode cognitive states from brain images, *Machine*
811 *Learning* 57 (2004) 145–175.
- 812 [19] K. A. Norman, S. M. Polyn, G. J. Detre, J. V. Haxby, Beyond mind-reading:
813 multi-voxel pattern analysis of fmri data, *Trends in cognitive sciences* 10
814 (2006) 424–430.
- 815 [20] L. Parra, C. Alvino, A. Tang, B. Pearlmutter, N. Yeung, A. Osman, P. Sajda,
816 Single-trial detection in EEG and MEG: Keeping it linear, *Neurocomputing*
817 52-54 (2003) 177–183.
- 818 [21] J. W. Rieger, C. Reichert, K. R. Gegenfurtner, T. Noesselt, C. Braun, H.-J.
819 Heinze, R. Kruse, H. Hinrichs, Predicting the recognition of natural scenes
820 from single trial meg recordings of brain activity, *Neuroimage* 42 (2008)
821 1056–1068.
- 822 [22] M. K. Carroll, G. A. Cecchi, I. Rish, R. Garg, A. R. Rao, Prediction and
823 interpretation of distributed neural activity with sparse models, *NeuroImage*
824 44 (2009) 112–122.
- 825 [23] A. M. Chan, E. Halgren, K. Marinkovic, S. S. Cash, Decoding word and
826 category-specific spatiotemporal representations from meg and eeg, *Neu-*
827 *roimage* 54 (2011) 3028–3039.
- 828 [24] H. Huttunen, T. Manninen, J.-P. Kauppi, J. Tohka, Mind reading with
829 regularized multinomial logistic regression, *Machine vision and applications*
830 24 (2013) 1311–1325.
- 831 [25] D. Vidaurre, C. Bielza, P. Larrañaga, A survey of l1 regression, *International*
832 *Statistical Review* 81 (2013) 361–387.
- 833 [26] M. Abadi, R. Subramanian, S. Kia, P. Avesani, I. Patras, N. Sebe, De-
834 caf: Meg-based multimodal database for decoding affective physiological re-
835 sponses, *IEEE Transactions on Affective Computing* 6 (2015) 209–222.

- 836 [27] T. Naselaris, K. N. Kay, S. Nishimoto, J. L. Gallant, Encoding and decoding
837 in fmri, *Neuroimage* 56 (2011) 400–410.
- 838 [28] S. Weichwald, T. Meyer, O. Özdenizci, B. Schölkopf, T. Ball, M. Grosse-
839 Wentrup, Causal interpretation rules for encoding and decoding models in
840 neuroimaging, *NeuroImage* 110 (2015) 48–59.
- 841 [29] N. Kriegeskorte, R. Goebel, P. Bandettini, Information-based functional
842 brain mapping, *Proceedings of the National Academy of Sciences of the*
843 *United States of America* 103 (2006) 3863–3868.
- 844 [30] F. J. Valverde-Albacete, C. Peláez-Moreno, 100% classification accuracy
845 considered harmful: The normalized information transfer factor explains
846 the accuracy paradox, *PLOS ONE* 9 (2014) e84217.
- 847 [31] A. Ramdas, A. Singh, L. Wasserman, Classification accuracy as a proxy for
848 two sample testing, *arXiv preprint arXiv:1602.02210* (2016).
- 849 [32] R. Turner, A model explanation system, 2015.
- 850 [33] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, K.-R.
851 Müller, How to explain individual classification decisions, *The Journal of*
852 *Machine Learning Research* 11 (2010) 1803–1831.
- 853 [34] A. Vellido, J. Martin-Guerrero, P. Lisboa, Making machine learning models
854 interpretable, in: *Proceedings of the 20th European Symposium on Arti-*
855 *ficial Neural Networks, Computational Intelligence and Machine Learning*
856 *(ESANN)*. Bruges, Belgium, 2012, pp. 163–172.
- 857 [35] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek,
858 On pixel-wise explanations for non-linear classifier decisions by layer-wise
859 relevance propagation, *PloS one* 10 (2015).
- 860 [36] G. Montavon, M. Braun, T. Krueger, K.-R. Müller, Analyzing local struc-
861 ture in kernel-based learning: Explanation, complexity, and reliability as-
862 sessment, *Signal Processing Magazine, IEEE* 30 (2013) 62–74.
- 863 [37] D. Yu, S. J. Lee, W. J. Lee, S. C. Kim, J. Lim, S. W. Kwon, Classification
864 of spectral data using fused lasso logistic regression, *Chemometrics and*
865 *Intelligent Laboratory Systems* 142 (2015) 70–77.
- 866 [38] K. Hansen, D. Baehrens, T. Schroeter, M. Rupp, K.-R. Müller, Visual
867 interpretation of kernel-based prediction models, *Molecular Informatics* 30
868 (2011) 817–826.

- 869 [39] S. Haufe, F. Meinecke, K. Görgen, S. Dähne, J.-D. Haynes, B. Blankertz,
870 F. Bießmann, On the interpretation of weight vectors of linear models in
871 multivariate neuroimaging, *NeuroImage* (2013).
- 872 [40] M. R. Sabuncu, A universal and efficient method to compute maps from
873 image-based prediction models, *Medical Image Computing and Computer-*
874 *Assisted Intervention–MICCAI 2014* (2014) 353–360.
- 875 [41] J.-D. Haynes, A primer on pattern-based approaches to fmri: Principles,
876 pitfalls, and perspectives, *Neuron* 87 (2015) 257–270.
- 877 [42] T. Naselaris, K. N. Kay, Resolving ambiguities of mvpa using explicit models
878 of representation, *Trends in cognitive sciences* 19 (2015) 551–554.
- 879 [43] S. C. Strother, P. M. Rasmussen, N. W. Churchill, K. Hansen, *Stability and*
880 *Reproducibility in fMRI Analysis*, New York: Springer-Verlag, 2014.
- 881 [44] A. Anderson, J. S. Labus, E. P. Vianna, E. A. Mayer, M. S. Cohen, Com-
882 mon component classification: What can we learn from machine learning?,
883 *Neuroimage* 56 (2011) 517–524.
- 884 [45] K. H. Brodersen, F. Haiss, C. S. Ong, F. Jung, M. Tittgemeyer, J. M.
885 Buhmann, B. Weber, K. E. Stephan, Model-based feature construction for
886 multivariate decoding, *NeuroImage* 56 (2011) 601–615.
- 887 [46] G. Langs, B. H. Menze, D. Lashkari, P. Golland, Detecting stable distributed
888 patterns of brain activation using gini contrast, *NeuroImage* 56 (2011) 497–
889 507.
- 890 [47] G. Varoquaux, A. Gramfort, B. Thirion, Small-sample brain mapping:
891 sparse recovery on spatially correlated designs with randomization and clus-
892 tering, in: *Proceedings of the 29th International Conference on Machine*
893 *Learning (ICML-12)*, 2012, pp. 1375–1382.
- 894 [48] J.-P. Kauppi, L. Parkkonen, R. Hari, A. Hyvärinen, Decoding magnetoen-
895 cephalographic rhythmic activity using spectrospatial information, *NeuroIm-*
896 *age* 83 (2013) 921–936.
- 897 [49] S. Taulu, J. Simola, J. Nenonen, L. Parkkonen, Novel noise reduction meth-
898 ods, *Magnetoencephalography* (2014) 35–71.
- 899 [50] G. Varoquaux, B. Thirion, How machine learning is shaping cognitive neu-
900 roimaging, *GigaScience* 3 (2014) 28.

- 901 [51] E. Olivetti, S. M. Kia, P. Avesani, Meg decoding across subjects, in: Pattern
902 Recognition in Neuroimaging, 2014 International Workshop on, IEEE, 2014.
- 903 [52] S. Haufe, S. Dähne, V. V. Nikulin, Dimensionality reduction for the analysis
904 of brain oscillations, *NeuroImage* (2014).
- 905 [53] P. M. Rasmussen, L. K. Hansen, K. H. Madsen, N. W. Churchill, S. C.
906 Strother, Model sparsity and brain pattern interpretation of classification
907 models in neuroimaging, *Pattern Recognition* 45 (2012) 2085–2100.
- 908 [54] B. R. Conroy, J. M. Walz, P. Sajda, Fast bootstrapping and permutation
909 testing for assessing reproducibility and interpretability of multivariate fmri
910 decoding models, *PloS one* 8 (2013) e79271.
- 911 [55] O. Bousquet, A. Elisseeff, Stability and generalization, *The Journal of*
912 *Machine Learning Research* 2 (2002) 499–526.
- 913 [56] B. Yu, Stability, *Bernoulli* 19 (2013) 1484–1500.
- 914 [57] C. Lim, B. Yu, Estimation stability with cross validation (escv), *Journal of*
915 *Computational and Graphical Statistics* (2015).
- 916 [58] N. Mørch, L. K. Hansen, S. C. Strother, C. Svarer, D. A. Rottenberg,
917 B. Lautrup, R. Savoy, O. B. Paulson, Nonlinear versus linear models in
918 functional neuroimaging: Learning curves and generalization crossover, in:
919 *Information processing in medical imaging*, Springer Berlin Heidelberg, 1997,
920 pp. 259–270.
- 921 [59] M. Yuan, Y. Lin, Model selection and estimation in regression with grouped
922 variables, *Journal of the Royal Statistical Society: Series B (Statistical*
923 *Methodology)* 68 (2006) 49–67.
- 924 [60] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, K. Knight, Sparsity and
925 smoothness via the fused lasso, *Journal of the Royal Statistical Society:*
926 *Series B (Statistical Methodology)* 67 (2005) 91–108.
- 927 [61] E. P. Xing, M. Kolar, S. Kim, X. Chen, High-dimensional sparse structured
928 input-output models, with applications to gwas, *Practical Applications of*
929 *Sparse Modeling* (2014) 37.
- 930 [62] I. Rish, G. A. Cecchi, A. Lozano, A. Niculescu-Mizil, *Practical Applications*
931 *of Sparse Modeling*, MIT Press, 2014.

- 932 [63] L. Grosenick, S. Greer, B. Knutson, Interpretable classifiers for fmri improve
933 prediction of purchases, *Neural Systems and Rehabilitation Engineering*,
934 *IEEE Transactions on* 16 (2008) 539–548.
- 935 [64] M. de Brecht, N. Yamagishi, Combining sparseness and smoothness improves
936 classification accuracy and interpretability, *NeuroImage* 60 (2012) 1550–
937 1561.
- 938 [65] V. Michel, A. Gramfort, G. Varoquaux, E. Eger, B. Thirion, Total variation
939 regularization for fmri-based prediction of behavior, *Medical Imaging*, *IEEE*
940 *Transactions on* 30 (2011) 1328–1340.
- 941 [66] A. Gramfort, B. Thirion, G. Varoquaux, Identifying predictive regions from
942 fmri with tv-l1 prior, in: *Pattern Recognition in Neuroimaging (PRNI)*, 2013
943 *International Workshop on*, *IEEE*, 2013, pp. 17–20.
- 944 [67] L. Grosenick, B. Klingenberg, S. Greer, J. Taylor, B. Knutson, Whole-brain
945 sparse penalized discriminant analysis for predicting choice, *NeuroImage* 47
946 (2009) S58.
- 947 [68] L. Grosenick, B. Klingenberg, K. Katovich, B. Knutson, J. E. Taylor, In-
948 terpretable whole-brain prediction analysis with graphnet, *NeuroImage* 72
949 (2013) 304–321.
- 950 [69] Y. Wang, J. Zheng, S. Zhang, X. Duan, H. Chen, Randomized structural
951 sparsity via constrained block subsampling for improved sensitivity of dis-
952 criminative voxel identification, *NeuroImage* (2015).
- 953 [70] F. Bießmann, S. Dähne, F. C. Meinecke, B. Blankertz, K. Görden, K.-R.
954 Müller, S. Haufe, On the interpretability of linear multivariate neuroimaging
955 analyses: filters, patterns and their relationship, in: *Proceedings of the 2nd*
956 *NIPS Workshop on Machine Learning and Interpretation in Neuroimaging*,
957 2012.
- 958 [71] S. Haufe, F. Meinecke, K. Görden, S. Dähne, J.-D. Haynes, B. Blankertz,
959 F. Bießmann, Parameter interpretation, regularization and source localiza-
960 tion in multivariate linear models, in: *Pattern Recognition in Neuroimaging*,
961 2014 *International Workshop on*, *IEEE*, 2014, pp. 1–4.
- 962 [72] D. A. Engemann, A. Gramfort, Automated model selection in covariance
963 estimation and spatial whitening of meg and eeg signals, *NeuroImage* 108
964 (2015) 328–342.

- 965 [73] Z. Li, Y. Wang, Y. Wang, X. Wang, J. Zheng, H. Chen, A novel feature
966 selection approach for analyzing high dimensional functional mri data, arXiv
967 preprint arXiv:1506.08301 (2015).
- 968 [74] S. M. Kia, S. Vega-Pons, E. Olivetti, P. Avesani, Multi-task learning for
969 interpretation of brain decoding models, in: NIPS Workshop on Machine
970 Learning and Interpretation in Neuroimaging (MLINI), 2014, Springer Lec-
971 ture Notes on Artificial Intelligence Series, In press.
- 972 [75] R. Tibshirani, Regression shrinkage and selection via the lasso, Journal of
973 the Royal Statistical Society. Series B (Methodological) (1996) 267–288.
- 974 [76] H. Zou, T. Hastie, Regularization and variable selection via the elastic net,
975 Journal of the Royal Statistical Society: Series B 67 (2005) 301–320.
- 976 [77] R. Jenatton, J.-Y. Audibert, F. Bach, Structured variable selection with
977 sparsity-inducing norms, arXiv preprint arXiv:0904.3523 (2009).
- 978 [78] T. Poggio, C. Shelton, On the mathematical foundations of learning, Amer-
979 ican Mathematical Society 39 (2002) 1–49.
- 980 [79] M. C.-K. Wu, S. V. David, J. L. Gallant, Complete functional characteri-
981 zation of sensory neurons by system identification, Annu. Rev. Neurosci. 29
982 (2006) 477–505.
- 983 [80] C. C. Aggarwal, P. S. Yu, A survey of uncertain data algorithms and appli-
984 cations, Knowledge and Data Engineering, IEEE Transactions on 21 (2009)
985 609–623.
- 986 [81] B. Efron, Bootstrap methods: another look at the jackknife, The annals of
987 Statistics (1979) 1–26.
- 988 [82] R. Kohavi, et al., A study of cross-validation and bootstrap for accuracy
989 estimation and model selection, in: Ijcai, volume 14, 1995, pp. 1137–1145.
- 990 [83] P. Domingos, A unified bias-variance decomposition for zero-one and squared
991 loss, AAAI/IAAI 2000 (2000) 564–569.
- 992 [84] M. D. Rugg, M. G. Coles, Electrophysiology of mind: Event-related brain
993 potentials and cognition., Oxford University Press, 1995.
- 994 [85] T. Hastie, R. Tibshirani, J. Friedman, The elements of statistical learning,
995 volume 2, Springer, 2009.

- 996 [86] A. Gramfort, G. Varoquaux, B. Thirion, Beyond brain reading: randomized
997 sparsity and clustering to simultaneously predict and identify, in: Machine
998 Learning and Interpretation in Neuroimaging, Springer, 2012, pp. 9–16.
- 999 [87] M. Caramia, P. Dell’ Olmo, Multi-objective optimization, Multi-objective
1000 Management in Freight Logistics: Increasing Capacity, Service Level and
1001 Safety with Optimization Algorithms (2008) 11–36.
- 1002 [88] R. T. Marler, J. S. Arora, Survey of multi-objective optimization methods
1003 for engineering, Structural and multidisciplinary optimization 26 (2004)
1004 369–395.
- 1005 [89] R. N. Henson, D. G. Wakeman, V. Litvak, K. J. Friston, A Parametric Em-
1006 pirical Bayesian framework for the EEG/MEG inverse problem: generative
1007 models for multisubject and multimodal integration, Frontiers in Human
1008 Neuroscience 5 (2011).
- 1009 [90] S. Bentin, T. Allison, A. Puce, E. Perez, G. McCarthy, Electrophysiological
1010 studies of face perception in humans, Journal of cognitive neuroscience 8
1011 (1996) 551–565.
- 1012 [91] D. H. Wolpert, W. G. Macready, An efficient method to estimate bagging’s
1013 generalization error, Machine Learning 35 (1999) 41–55.
- 1014 [92] L. Breiman, Random forests, Machine learning 45 (2001) 5–32.
- 1015 [93] V. N. Vapnik, S. Kotz, Estimation of dependences based on empirical data,
1016 volume 40, Springer-verlag New York, 1982.
- 1017 [94] V. Vapnik, The nature of statistical learning theory, Springer Science &
1018 Business Media, 2013.
- 1019 [95] R. Oostenveld, P. Fries, E. Maris, J.-M. Schoffelen, Fieldtrip: open source
1020 software for advanced analysis of meg, eeg, and invasive electrophysiological
1021 data, Computational intelligence and neuroscience 2011 (2010).
- 1022 [96] S. Dash, D. M. Malioutov, K. R. Varshney, Learning interpretable classifi-
1023 cation rules using sequential rowsampling, in: Acoustics, Speech and Signal
1024 Processing (ICASSP), 2015 IEEE International Conference on, IEEE, 2015,
1025 pp. 3337–3341.
- 1026 [97] B. Afshin-Pour, H. Soltanian-Zadeh, G.-A. Hossein-Zadeh, C. L. Grady, S. C.
1027 Strother, A mutual information-based metric for evaluation of fmri data-
1028 processing approaches, Human brain mapping 32 (2011) 699–715.

- 1029 [98] J. B. T. Zhang, Support vector classification with input data uncertainty,
1030 Advances in neural information processing systems 17 (2005) 161.
- 1031 [99] C. Tzelepis, V. Mezaris, I. Patras, Linear maximum margin classifier for
1032 learning from uncertain data, arXiv preprint arXiv:1504.03892 (2015).
- 1033 [100] G. Valentini, T. G. Dietterich, Bias-variance analysis of support vector
1034 machines for the development of svm-based ensemble methods, The Journal
1035 of Machine Learning Research 5 (2004) 725–775.
- 1036 [101] R. Tibshirani, Bias, variance and prediction error for classification rules,
1037 University of Toronto, Department of Statistics, 1996.