

---

# Region-network hierarchical sparsity priors for high-dimensional inference in brain imaging

---

Danilo Bzdok, Michael Eickenberg, Gaël Varoquaux, Bertrand Thirion

Department of Psychiatry, Psychotherapy and Psychosomatics, RWTH Aachen, Germany

INRIA, Parietal team, Saclay, France

CEA, Neurospin, Gif-sur-Yvette, France

firstname.lastname@inria.fr

## Abstract

Structured sparsity penalization has recently improved statistical models applied to high-dimensional data in various domains. As an extension to imaging neuroscience, priors on network hierarchies of brain regions have been incorporated into logistic-regression estimators for neural activity maps. This remarries the perspectives of functional segregation and functional integration that are routinely divorced in neuroscientific research. Region-network hierarchical priors for supervised classification of 18 psychological tasks from a reference dataset are shown to outperform naïve  $\ell_1$ -norm and trace-norm penalization, as well as neurobiologically informed  $\ell_1/\ell_2$ -block-norm group sparsity. Varying the relative importance of region and network structure within the hierarchical tree prior recovered complementary aspects of the neural activity patterns with similar predictive accuracies. In sum, we demonstrate that priors of local and global neurobiological knowledge can enhance out-of-sample performance and domain interpretability by hierarchical tree sparsity.

**Keywords:** Sparsity-inducing norms, hierarchical tree sparsity, numerical optimization, systems neuroscience, functional specialization, functional integration

## 1 Introduction

Many quantitative scientific domains underwent a recent passage from the classical regime (i.e., “long data”) to the high-dimensional regime (i.e., “wide data”) [30]. Also in the brain imaging domain, many contemporary methods for acquiring brain signals yield more variables per observation than total observations per data sample. This high-dimensional scenario challenges various statistical methods from classical statistics. For instance, estimating generalized linear models without additional assumptions yields an underdetermined system of equations. Many such ill-posed estimation problems have benefited from *sparsity* assumptions [11, 24]. They act as a regularizer and can be used for model selection. Sparse supervised and unsupervised learning algorithms have proven to yield statistical relationships that can be readily estimated, reproduced, and interpreted [21]. Generally, *structured sparsity* can impose domain knowledge on the statistical estimation, thus shrinking and selecting variables guided by expected data distributions [3]. Such restrictions to complexity are an attractive plan of attack for the >100,000 variables of brain maps. Yet, what generally accepted neurobiological structure lends itself to harness the *curse of dimensionality* by structured sparsity priors?

Concepts on human brain organization have long been torn between the two extremes *functional specialization* and *functional integration*. Functional specialization emphasizes that microscopically distinguishable brain regions are responsible distinct classes of computational processes [32]. Functional integration, in turn, emphasizes that brain function is enabled by complex connections

between these distinct brain regions [45]. These notions were predominantly derived from invasive examination of anatomy (i.e., histological preparation), connectivity (i.e., axonal tracing), and functional properties (i.e., single-cell recordings) in animals. Regarding functional segregation into specialized regions, early histological investigations into the microscopic heterogeneity of the human cerebral cortex have resulted in several detailed anatomical maps [9]. Regarding axonal connections, each such cortical area has been observed to possess a unique set of incoming and outgoing connections [36, 50, 41]. Both local infrastructure and its unique global connectivity profile together are thought to realize brain function. In sum, cortical brain modules versus connections between them reflect functional specialization versus functional integration [20, 35]. Importantly, probably no existing brain analysis method acknowledges that both functional organizations are inextricably involved in the realization of mental operations [47, 40].

Functional specialization has been explored and interpreted based on many different research methods. Single-cell recordings and microscopic examination revealed, for instance, the specialization in the occipital visual cortex into V1, V2, V3, V3A/B, and V4 [25, 52]. Tissue lesion of the mid-fusiform gyrus of the visual system, in turn, was frequently reported to impair recognition of others' identity from faces [26]. The whole-brain localization of sensory, motor, and emotional functions to cortical areas was later enabled by non-invasive brain imaging with functional magnetic resonance imaging (fMRI) and positron emission tomography (PET) [19]. Further, radioactive mapping of neurotransmitter receptors rendered accessible yet another local characteristic of neuronal populations [53]. In the computational era, automatic clustering methods are increasingly employed to regionally differentiate the cerebral cortex, which can partly be more fine-grained than classical microscopical borders [7, 17]. High-throughput approaches today enable ultrahigh-resolution 3D models of brain anatomy at near-cellular scale [2]. As a crucial common point, all these methodological approaches yield neuroscientific findings that are naturally interpreted according to non-overlapping, discrete region compartments as the basic architecture of brain organization.

It is more recent that the main interpretational focus has shifted from circumscribed regions to network stratifications in systems neuroscience [51, 46]. Invasive axonal tracing studies in monkeys were complemented by diffusion MRI tractography in humans as a now frequently employed method to outline fiber bundles between brain regions [27]. Besides analyses of electrophysiological oscillations [13] and graph-theoretical properties [12], studies of functional connectivity [10] and independent component analysis (ICA) [6] became the workhorses of network discovery in neuroimaging. These revealed the important implication of canonical brain networks across psychological tasks, including the so-called "default-mode network" [38], "salience network" [42], and "dorsal attention network" [14]. Characteristic changes in the configuration of these macroscopical networks were repeatedly observed to be induced by the onset of given psychological tasks [18]. As a common point of all these methods, interpretation of findings naturally embraces cross-regional integration by overlapping network compartments as the basic architecture of brain organization, in stark contrast to methods examining regional specialization.

Building on these two major interpretational traditions in systems neuroscience, the present study proposes to incorporate established neurobiological structure underlying functional segregation and integration into supervised estimators by hierarchical structured sparsity. Learning techniques exploiting structured sparsity have recently made much progress in various application domains from processing auditory signals [16], natural images [23] and videos [31, 34] to astrophysics [48], genetics [39, 33], and conformational dynamics of protein complexes [29]. This is extended by the present work that introduced neuroscience-specific estimators capitalizing on neurobiologically plausible region and network priors. Based on the largest neuroimaging repository, we demonstrated that domain-informed supervised models gracefully tackle the curse of dimensionality, yield more human-interpretable results, and generalize better to new samples than domain-naïve estimators.

## 2 Methods

Our main contribution is the domain-specific adaptation of sparse structured penalties making it possible to jointly incorporate functional specialization and functional integration priors into statistical estimators. We capitalize on hierarchical group lasso as introduced by Jennatton and others in [28] to create a set of convex penalty terms. These can probe the interplay between regional specialization and inter-regional functional integration during defined psychological tasks.

**Rationale.** Three-dimensional voxel brain images as obtained by various neuroimaging techniques are very high-dimensional, but also very highly structured. While its explicit dimensionality, the number of brain voxels, which varies with the resolution of the image is generally of the order of 100,000, the number of samples available for analyses rarely exceeds hundreds or thousands. This  $n \ll p$  scenario immediately implies underdetermination of any linear model based on dot products with the voxel values.

Yet, there is abundant structure to be exploited, which can be injected into the estimation as a domain knowledge prior. Here we are interested in incorporating knowledge about both the well-established modules of functional specialization and the spatiotemporal interactions between these. Developmentally, such large-scale networks emerge during late fetal growth (Doria et al., 2010), before mental capacities mature in childhood. In adults, nodes of a such a spatiotemporally cohesive network have more similar functional profiles than nodes from different networks (Anderson et al., 2013). On the signal level, functional modules are localized structures implicating neighboring voxels in an image. A variety of known functional modules has been compiled into atlases of the brain. Different types of atlases rely on different ways of segregating modules from another. For instance, anatomical atlases rely on the variation of tissue properties across the brain, which can be at the cellular level or based on higher-order structure. Global functional networks can be discovered by statistical analysis of functional brain imaging data, for example the BOLD contrast in fMRI. A well-established method of estimating these networks is independent component analysis [?]. It robustly identifies the default mode, saliency and attention networks. Global functional networks as estimated by ICA are continuous brain maps attributing weights to each voxel. Since the regions involved in these networks are nevertheless largely disjoint spatially and roughly in correspondence with module boundaries **Danilo: Ok!**, it is possible to uniquely associate each module from a region atlas to one of the extracted global functional networks.

This correspondence makes it possible to inject established knowledge about the networks and regions into a hierarchy: A brain contains a certain number of networks, which in turn consist of the atlas modules as subregions. We can use this hierarchy to create a structural prior of expected weight maps for classifiers.

The hierarchical group lasso as introduced by [?] lends itself perfectly to represent this structure. It is based on the group lasso [?] but permits groups to contain each other in a tree structure. The first level of the tree are the network groups containing all the voxels of the modules associated to them. Then each network contains as subgroups the modules associated with it, grouping the voxels of each of these regions together. As with the group lasso, it is possible to associate an individual penalty to each group. In trading off the penalties on the network level against the penalties of the region level, we can create different regimes of estimation: Setting a low penalty on the network groups makes it probable that all of them are active in the estimated weight map. If we then select higher penalties on region groups, selection of relevant region groups is forced without the bias of the network maps. Conversely, setting low penalties on the region maps makes it possible for all voxels to be active. Selecting higher penalties on the networks then leads to a selection of networks with all regions associated to it. Taken together and varied smoothly, we can evaluate a tradeoff between the relevance of modular regions and functionally integrated networks in statistical estimation of brain activity.

**Problem formulation.** We formulate our estimation problem in the framework of regularized risk estimation applied to linear models: We would like to estimate a good predictor of cognitive task given a brain image. Let the set  $\mathcal{X} \subset \mathbb{R}^p$  represent brain images of  $p > 0$  voxels.

Then we would like to minimize the risk  $\mathcal{L}(\hat{y}, y)$ , where  $\hat{y} = X\hat{w} + \hat{b}$ , while regularizing to incorporate a useful prior. Taken together, this can be framed as an optimization problem

$$\arg \min_{w,b} \mathcal{L}(Xw + b, y) + \lambda \Omega(w),$$

where  $\lambda > 0$  and  $\Omega$  is the regularizer.

Brain regions are defined as disjoint groups of voxels. Let  $\mathcal{G}$  be a partition of  $\{1, \dots, p\}$ , i.e.

$$\bigcup_i g_i = \{1, \dots, p\} \text{ and } g_i \cap g_j = \emptyset \quad \forall i \neq j$$

Brain networks consist of regions and are thus super-regions or groups of regions. The set of brain networks  $\mathcal{H}$  is also a partition of  $\{1, \dots, p\}$  and in addition it is consistent with  $\mathcal{G}$  in the sense that

$$\text{for all } g \in \mathcal{G}, h \in \mathcal{H}, \text{ either } g \subset h \text{ or } g \cap h = \emptyset.$$

This allows a clear association of each region  $g \in \mathcal{G}$  to a network  $h \in \mathcal{H}$  and thus establishes a tree structure (up to adding a root node containing all voxels).

For a brain image  $w \in \mathbb{R}^p$  and a group  $g$ , the vector  $w_g \in \mathbb{R}^{|g|}$  is defined as the restriction of  $w$  to the coordinates in  $g$ . The structured penalty incorporating network and region information can then be written as

$$\Omega(w) = \alpha \sum_{h \in \mathcal{H}} \eta_h \|w_h\|_2 + \beta \sum_{g \in \mathcal{G}} \eta_g \|w_g\|_2.$$

According to [?] we set  $\eta_g = 1/\sqrt{|g|}$  to account for varying group size. The hierarchy-level-specific factors  $\alpha > 0$  and  $\beta > 0$  are used to trade-off region-weighted and network-weighted models against each other.

The prediction problem at hand is a multiclass classification. We choose to attack this using one-vs-rest scheme on a binary logistic regression, whose loss can be written as

$$\sum_{i=1}^n \log(1 + \exp(-y_i \langle x_i, w \rangle)) + \lambda \Omega(w),$$

if  $y \in -1, 1$  and with  $x_i \in \mathbb{R}^p$  the training sample brain images.

## MISSING

- cross-validation scheme
- analysis of penalty tradeoff (CV scheme for that?)
- mentioning of control methods in methods. Trace norm, sparse group sparsity. Anything else? ENET?

**Implementation.** The analyses were performed in Python. We used *nilearn* to handle the large quantities of neuroimaging data [1] and *Theano* for automatic, numerically stable differentiation of symbolic computation graphs [5, 8]. All Python scripts that generated the results are accessible online for reproducibility and reuse (<http://github.com/banilo/nips2015>).

all algorithm from a same software library - $\zeta$  SPAMs

**Data.** As the currently biggest openly-accessible reference dataset, we chose resources from the Human Connectome Project (HCP) [4]. Neuroimaging task data with labels of ongoing cognitive processes were drawn from 500 healthy HCP participants (cf. Appendix for details on datasets). 18 HCP tasks were selected that are known to elicit reliable neural activity across participants (Table 1). In sum, the HCP task data incorporated 8650 first-level activity maps from 18 diverse paradigms administered to 498 participants (2 removed due to incomplete data). All maps were resampled to a common  $60 \times 72 \times 60$  space of 3mm isotropic voxels and gray-matter masked (at least 10% tissue probability). The supervised analyses were thus based on labeled HCP task maps with 79,941 voxels of interest representing z-values in gray matter.

These labeled data were complemented by unlabeled activity maps from HCP acquisitions of unconstrained resting-state activity [43]. These reflect brain activity in the absence of controlled thought. In sum, the HCP rest data concatenated 8000 unlabeled, noise-cleaned rest maps with 40 brain maps from each of 200 randomly selected participants.

We were further interested in the utility of the optimized low-rank projection in one task dataset for dimensionality reduction in another task dataset. To this end, the HCP-derived network decompositions were used as preliminary step in the classification problem of another large sample. The ARCHI dataset [37] provides activity maps from diverse experimental tasks, including auditory and visual perception, motor action, reading, language comprehension and mental calculation. Analogous to HCP data, the second task dataset thus incorporated 1404 labeled, grey-matter masked, and z-scored activity maps from 18 diverse tasks acquired in 78 participants.

sparse statistical models have only few nonzero parameters

Cognitive Task	Stimuli	Instruction for participants
1 Reward	Card game	Guess the number of a mystery card for gain/loss of money
2 Punish		
3 Shapes	Shape pictures	Decide which of two shapes matches another shape geometrically
4 Faces	Face pictures	Decide which of two faces matches another face emotionally
5 Random		
6 Theory of mind	Videos with objects	Decide whether the objects act randomly or intentionally
7 Mathematics	Spoken numbers	Complete addition and subtraction problems
8 Language	Auditory stories	Choose answer about the topic of the story
9 Tongue movement		Move tongue
10 Food movement		Squeezing of the left or right toe
11 Hand movement	Visual cues	Tapping of the left or right finger
12 Matching		
13 Relations	Shapes with textures	Decide whether two objects match in shape or texture
14 View Bodies	Pictures	Passive watching
15 View Faces	Pictures	Passive watching
16 View Places	Pictures	Passive watching
17 View Tools	Pictures	Passive watching
18 Two-Back	Various pictures	Indicate whether current stimulus is the same as two items earlier

Table 1: Description of psychological tasks to predict.

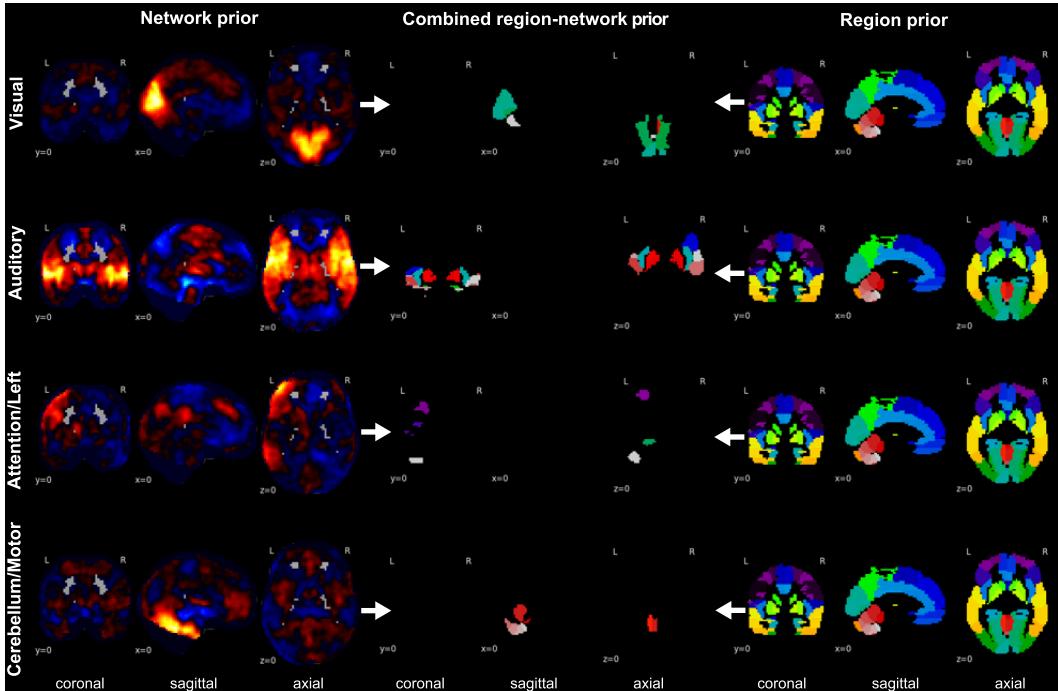
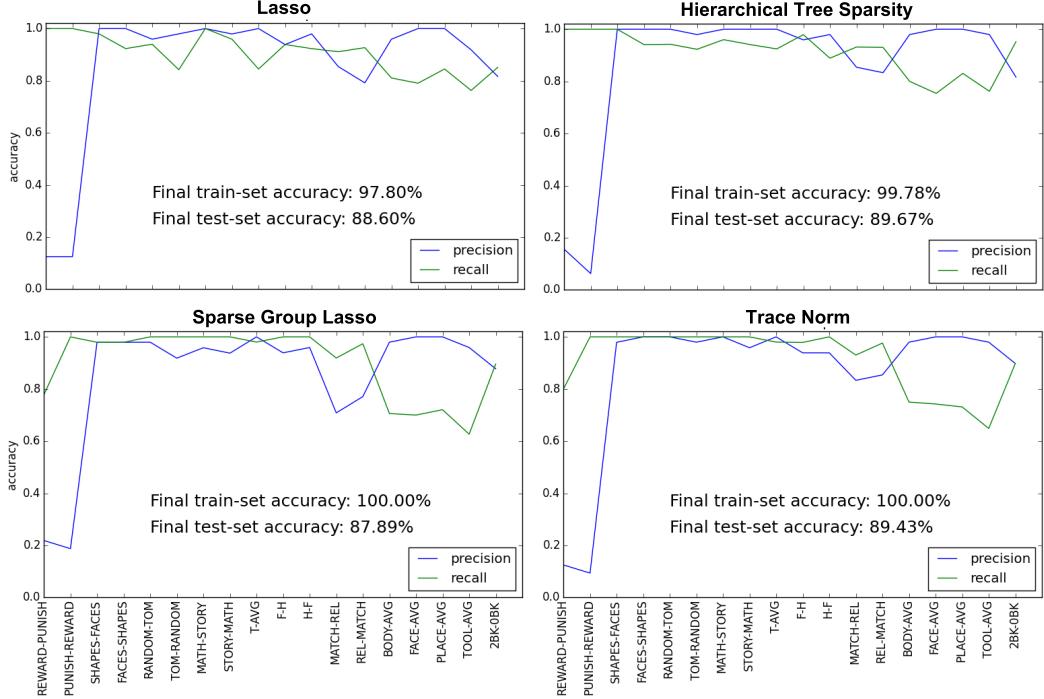


Figure 1: **Building blocks of the region-network tree.** Depicts neurobiological priors introduced into the classification problem by hierarchical structured sparsity. *Left:* Continuous, partially overlapping brain network priors (*hot-colored*, taken from [44]) accommodate the functional integration perspective of brain organization. *Right:* Discrete, non-overlapping brain region priors (*single-colored*, taken from [15]) accommodate the functional segregation perspective. *Middle:* These two types of predefined voxel groups are incorporated into hierarchical priors of parent networks with their descending region nodes. *Top to bottom:* Four exemplary region-network priors are shown, including the early cortex that processes visual and sound information from the environment, a well-known attentional circuit in the left brain hemisphere, and the cerebellum that is involved in motor behavior.

### 3 Experimental Results

**Benchmarking hierarchical tree sparsity against common sparsity penalties.** Hierarchical region-network priors have been systematically evaluated against other popular choices of sparse



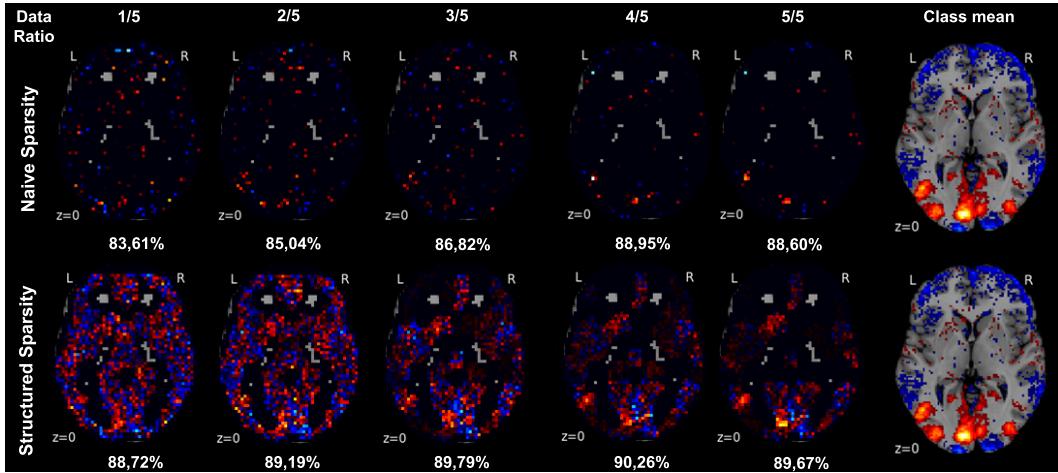
**Figure 2: Model performance across sparsity priors.** Compares the performance of logistic regression estimators with 4 different structured and unstructured sparsity terms in classifying neural activity from 18 psychological tasks. The class-wise precision and recall metrics were obtained on the same test set. Unstructured  $\ell_1$ -penalized logistic regression (*upper left*) imposed a minimum of relevant brain voxels without assuming special structure. Structured  $\ell_1/\ell_2$ -block norm with additional  $\ell_1$  term (*lower left*) imposed region compartments, but naïve to network structure. Structured trace-norm penalization (*lower right*) imposed low-rank structure with sparsity of network patterns, but naïve to region structure. Structured  $\ell_1/\ell_2$ -block-norm with a hierarchy of both region and network priors (*upper right*) exhibited the best out-of-sample performance. A priori knowledge of both region and network neighborhoods was hence most beneficial for predicting psychological tasks from brain maps.

classification algorithms in an 18-class scenario (Figure 2). Logistic regression with  $\ell_1/\ell_2$  block norm penalization incorporated a hierarchy of previously known region and network neighborhoods for a neurobiological bias of the statistical estimation. Vanilla logistic regression with  $\ell_1$ -penalization does not assume any previously known special structure. This classification estimator embraces a vision of neural activity structure that expects a minimum of topographically and functionally independent brain voxel to be relevant. Logistic regression with sparse group sparsity imposes a structured  $\ell_1/\ell_2$  block norm with additional  $\ell_1$  term with a known atlas of region voxel groups onto the statistical estimation process. This supervised estimator shrinks and selects the coefficients of topographically compact voxel groups expected to be relevant together. Logistic regression with trace-norm penalization imposed low-rank structure. This supervised classification algorithm expected a minimum of unknown network patterns to be relevant. The stratified and shuffled training data were submitted to a nested cross-validation scheme for model selection and model assessment. In the inner CV layer, the logistic regression estimators have been trained in a one-versus-rest design that distinguishes each class from the respective 17 other classes (number of maximal iterations=100, tolerance=0.001). In the outer CV layer, grid search selected among candidates for the respective  $\lambda$  parameter by searching between  $10^{-2}$  and  $10^1$  in 9 steps on a logarithmic scale. Importantly, the thus selected sparse logistic regression classifier was evaluated on an identical test set in all settings.

Across analyses, hierachial tree sparsity was most successful in distinguishing unseen neural activity maps from 18 psychological tasks (89.67%, mean recall XX.YY, mean precision XX.YY). It was closely followed by logistic regression structured by trace-norm regularization (89.43%, mean recall XX.YY, mean precision XX.YY). Lasso featured an average performance comparing to the other sparse estimators (88.60%, mean recall XX.YY, mean precision XX.YY). Introducing a priori knowledge of brain region compartments by sparse group sparsity performed worst (87.89%, mean recall XX.YY, mean precision XX.YY). In an important subanalysis, the gain of the combined region-network prior was also confirmed by selectively zeroing the  $\eta_g$  coefficients of all region groups or all network groups in the hierarchical prior. Removing region structure from the prior achieved 88.84% accuracy (mean recall XX.YY, mean precision XX.YY), while removing network structure from the prior achieved 87.05% accuracy (mean recall XX.YY, mean precision XX.YY). These results from partial priors are indeed outperformed the full region-network tree prior at 89.67% accuracy (mean recall XX.YY, mean precision XX.YY). In sum, biasing sparse model selection by domain knowledge of region-network hierarchies outcompeted other types of frequently used sparse penalization techniques.

**Sample complexity of naïve versus informed sparse model selection.** Subsequently, the sample complexity of  $\ell_1$ -penalized and hierarchical-tree-penalized logistic regression were quantitatively compared (Figure 3). Region-network priors should bias model selection towards more neurobiologically plausible classification estimators. This should yield better out-of-sample generalization and support recovery than neurobiology-naïve  $\ell_1$ -constrained logistic regression in the data-scarce and data-rich scenarios. The HCP task data with examples from 18 psychological tasks were first divided into 90% of training set (i.e., 7584 neural activity maps) and 10% of test set (i.e., 842 neural activity maps). Both learning algorithms were fitted based on the training set at different subsampling fractions: 20% (1516 maps), 40% (3033 maps), 60% (4550 maps), 80% (6067 maps), and 100% (7584 maps). The stratified and shuffled training data were submitted to a nested cross-validation scheme for model selection and model assessment. In the inner CV layer, the logistic regression estimators have been trained in a one-versus-rest design that distinguishes each class from the respective 17 other classes (number of maximal iterations=100, tolerance=0.001). In the outer CV layer, grid search selected among candidates for the respective  $\lambda$  parameter by searching between  $10^{-2}$  and  $10^1$  in 9 steps on a logarithmic scale. Importantly, the thus selected sparse logistic regression classifier was evaluated on an identical test set in all settings.

Three observations have been made. In the data-scarce scenario (i.e., 1/5 of available training data), hierarchical tree sparsity achieved the biggest advantage in out-of-sample performance by 5.11% as well as better support recovery with weight maps already much closer to the class averages. In the case of scarce training data, which is typical for the brain imaging domain, regularization by region-network priors thus allowed for more effective extraction of classification-relevant structure from the neural activity maps. Across training data fractions, the weight maps from ordinary logistic regression exhibited higher variance and more zero coefficients than hierarchical tree logistic regression. Given the usually high multicollinearity in neuroimaging data, this observation is likely

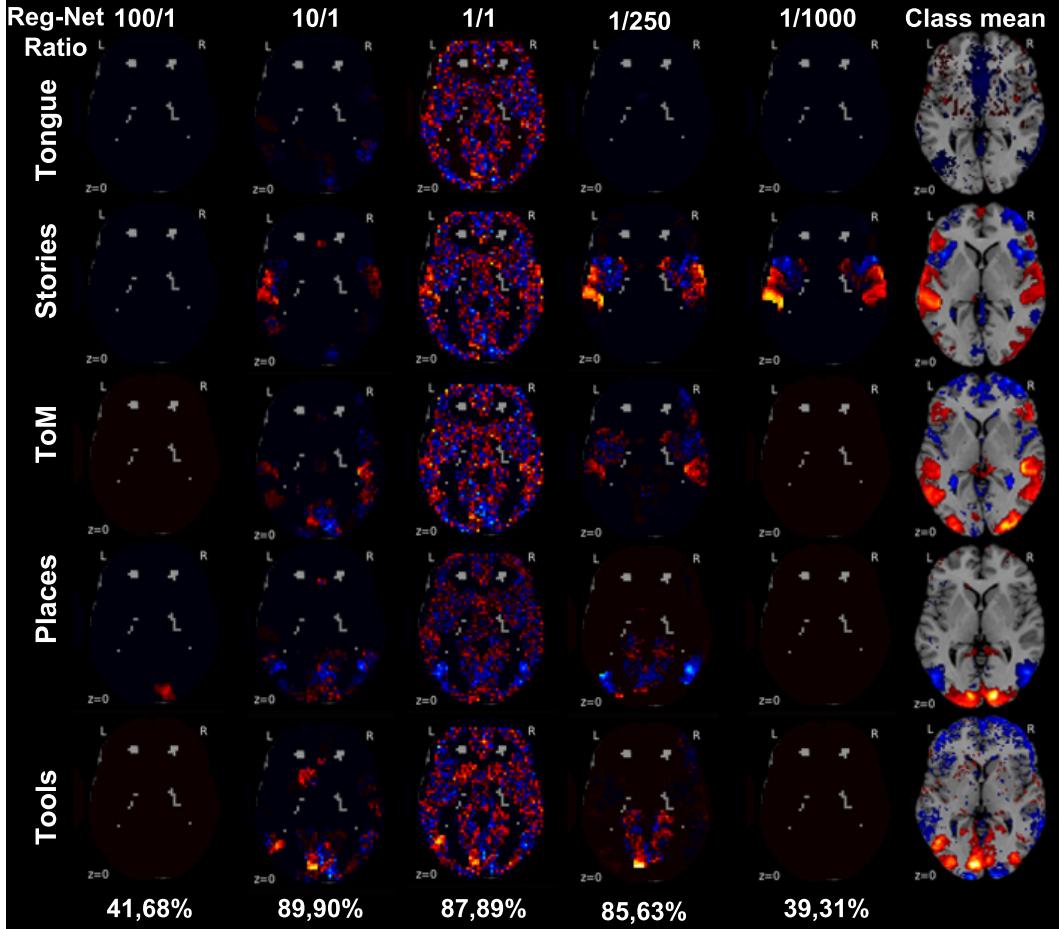


**Figure 3: Naïve versus informed sparse model selection across training set sizes.** Ordinary  $\ell_1$ -penalized logistic regression (*upper row*) is compared to hierarchical-tree-penalized logistic regression (*lower row*) with increasing fraction of the available training data (*left to right columns*). For one example (i.e., ‘‘View tools’’) from 18 psychological tasks, unthresholded axial maps of model weights are shown for comparison against the sample average of that class (*rightmost column*, thresholded at the 75<sup>th</sup> percentile). The out-of-sample accuracies for predicting all 18 psychological tasks is given in percent. In the data-scarce scenario, typical for brain imaging, hierarchical tree sparsity achieves much better support recovery with the biggest difference in model performance. In the data-rich scenario, neurobiologically informed logistic regression profits more from the available information quantities than neurobiologically naïve logistic regression.

to reflect instable selection of representatives among class-responsive predictor groups due to the  $\ell_1$ -norm penalization. In the data-rich scenario (i.e., entire training data used for model fitting), neurobiologically informed logistic regression profited more from the increased information quantities than neurobiologically naïve logistic regression. That is, the region-network priors actually further enhance the similarity to the weight maps even in abundant input data. This was the case although the maximal classification performance of  $\approx 90\%$  has already been reached with small training data fractions by the structured estimator. In contrast, the unstructured estimator reached this generalization performance only with bigger input data quantities.

**Support recovery as a function of region and network emphasis.** Finally, the relative importance of the region and network priors within the hierarchical tree prior was quantified (Figure 4). The group weight  $\eta_g$  of region priors was multiplied with a region-network ratio, while the group weight  $\eta_g$  of network priors was divided by that region-network ratio. For instance, a region-network ratio of 3 increased the relative importance of known region structure by multiplying  $\frac{3}{1}$  to  $\eta_g$  of all region group penalties and multiplying  $\frac{1}{3}$  to  $\eta_g$  of all network group penalties. The data splitting cross-validation scheme was identical to the above modelling experiments.

As the most important observation, a range between region-dominant and network-dominant structured penalties yielded quantitatively almost identical generalization to new data but qualitatively different decision functions manifested in the weight maps (Figure 4, second and forth column). Classification models with many zero coefficients but high absolute coefficients in either region compartments or network compartments can similarly extrapolate to unseen neural activity maps. Second, these achieve classification performance comparable to equilibrated region-network priors that set less voxel coefficients to zero and spread the probability mass with lower absolute coefficients across the whole brain (Figure 4, third column in the middle). Third, overly strong emphasis on either level of the hierarchical prior can yield the neurobiologically informative results with maps of the most necessary region or network structure for statistically significant out-of-sample performance (Figure 4, leftmost and rightmost columns). In sum, stratifying the hierarchical tree penalty between region and network emphasis suggests that *class-specific region-network weights* might offer more performant and more interpretable classification models in the future.



**Figure 4: Support recovery as a function of region and network emphasis.** The relative impact of the region and network priors on model selection is systematically varied against each other. This region-network ratio (*upper fractions*) weighted voxel groups to privilege sparse models in function space that acknowledge known brain region neighborhoods (*left columns*) or known brain networks neighborhoods (*right columns*). Among the 18 classes, the model weights are shown for the psychological tasks (*from top to bottom*): tongue movement, listening stories, taking somebody else’s perspective (ToM, “theory of mind”), as well as viewing locations and tools. The 18-class out-of-sample accuracy is indicated on the *bottom* and the class-wise mean neural activity (*right-most column*, thresholded at the 75<sup>th</sup> percentile). Different emphasis on regions versus networks in hierarchical structured sparsity can yield comparable model performance. Favoring region versus network structure during model selection recovers complementary aspects of the neural activity pattern. Equal region and network emphasis yields more dispersed, less interpretable predictive model choices.

## 4 Discussion

Relevant structure in neuroimaging data has long been investigated according to two separate organizational principles: functional segregation into discrete brain regions [36] and functional integration by interregional brain networks [45]. This proof-of-concept study demonstrates the simultaneous exploitation of both these neurobiological compartments for sparse variable selection and high-dimensional prediction in a reference dataset. Introducing existing domain knowledge into model selection allowed privileging members of the function space that are most neurobiologically plausible. Domain-informed hierarchical structured sparsity is shown to enhance both model interpretability and generalization performance, although these statistical-learning goals are typically in conflict.

The present approach has important advantages over previous analysis strategies that rely on dimensionality reduction of the neuroimaging data to harness the curse of dimensionality. They often use preliminary pooling functions within regions or regression against network templates for subsequent supervised learning on the aggregated feature space. Such lossy approaches divided into feature engineering and inference steps *i*) can only satisfy the specialization or integration account of brain organization, *ii*) depend on the ground truth being a region or network effect, and *iii*) cannot issue individual coefficients for every brain voxels. Hierarchical region-network sparsity addresses these shortcomings by estimating individual voxel contributions while benefitting from their functional segregation and integration to restrict statistical complexity. Viewed from the bias-variance trade-off, our modification to logistic regression estimators entailed a large decrease in model variance but only a modest increase in model bias. Viewed from the Vapnik-Chervonenkis dimensions, this entailed a healthy decrease in the complexity capacity of the prediction model with a higher chance of generalizing to unobserved data.

In the future, region-network sparsity priors could be incorporated into various pattern-learning methods in systems neuroscience. This includes supervised methods for whole-brain classification and regression with one or several target variables. The principled regularization scheme could even inform unsupervised structure-discovery methods, such as principal component analysis [29] and k-means clustering [49]. Additionally, model regularization by hierarchical structured sparsity could be extended from the spatial domain of neural activity to priors of coherent spatiotemporal activity structure [22]. Ultimately, successful high-dimensional inference is an important prerequisite for predicting diagnosis, disease trajectories, and treatment response in personalized psychiatry and neurology.

**Acknowledgment.** The research leading to these results has received funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement no. 604102 (Human Brain Project). Data were provided by the Human Connectome Project. Further support was received from the German National Academic Foundation (D.B.), the German Research Foundation (BZ2/2-1 and BZ2/3-1 to D.B.), and the MetaMRI associated team (B.T., G.V.).

## References

- [1] Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort, A., Thirion, B., Varoquaux, G.: Machine learning for neuroimaging with scikit-learn. *Front Neuroinform* 8, 14 (2014)
- [2] Amunts, K., Lepage, C., Borgeat, L., Mohlberg, H., Dickscheid, T., Rousseau, M.E., Bludau, S., Bazin, P.L., Lewis, L.B., Oros-Peusquens, A.M., et al.: Bigbrain: an ultrahigh-resolution 3d human brain model. *Science* 340(6139), 1472–1475 (2013)
- [3] Bach, F., Jenatton, R., Mairal, J., Obozinski, G.: Optimization with sparsity-inducing penalties. *Foundations and Trends® in Machine Learning* 4(1), 1–106 (2012)
- [4] Barch, D.M., Burgess, G.C., Harms, M.P., Petersen, S.E., Schlaggar, B.L., Corbetta, M., Glasser, M.F., Curtiss, S., Dixit, S., Feldt, C.: Function in the human connectome: task-fmri and individual differences in behavior. *Neuroimage* 80, 169–189 (2013)
- [5] Bastien, F., Lamblin, P., Pascanu, R., Bergstra, J., Goodfellow, I., Bergeron, A., Bouchard, N., Warde-Farley, D., Bengio, Y.: Theano: new features and speed improvements. *arXiv preprint arXiv:1211.5590* (2012)
- [6] Beckmann, C.F., DeLuca, M., Devlin, J.T., Smith, S.M.: Investigations into resting-state connectivity using independent component analysis. *Philos Trans R Soc Lond B Biol Sci* 360(1457), 1001–13 (2005)

- [7] Behrens, T.E., Johansen-Berg, H., Woolrich, M.W., Smith, S.M., Wheeler-Kingshott, C.A., Boulby, P.A., Barker, G.J., Sillery, E.L., Sheehan, K., Ciccarelli, O., Thompson, A.J., Brady, J.M., Matthews, P.M.: Non-invasive mapping of connections between human thalamus and cortex using diffusion imaging. *Nat Neurosci* 6(7), 750–7 (2003)
- [8] Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., Turian, J., Warde-Farley, D., Bengio, Y.: Theano: a cpu and gpu math expression compiler. Proceedings of the Python for scientific computing conference (SciPy) 4, 3 (2010)
- [9] Brodmann, K.: Vergleichende Lokalisationslehre der Groshirnrinde (1909)
- [10] Buckner, R.L., Krienen, F.M., Yeo, B.T.: Opportunities and limitations of intrinsic functional connectivity mri. *Nature neuroscience* 16(7), 832–837 (2013)
- [11] Bühlmann, P., Van De Geer, S.: Statistics for high-dimensional data: methods, theory and applications. Springer Science & Business Media (2011)
- [12] Bullmore, E., Sporns, O.: Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature Reviews Neuroscience* 10(3), 186–198 (2009)
- [13] Buzsáki, G., Draguhn, A.: Neuronal oscillations in cortical networks. *science* 304(5679), 1926–1929 (2004)
- [14] Corbetta, M., Patel, G., Shulman, G.L.: The reorienting system of the human brain: from environment to theory of mind. *Neuron* 58(3), 306–24 (2008)
- [15] Craddock, R.C., James, G.A., Holtzheimer, P. E., r., Hu, X.P., Mayberg, H.S.: A whole brain fmri atlas generated via spatially constrained spectral clustering. *Hum Brain Mapp* 33(8), 1914–28 (2012)
- [16] Daudet, L.: Sparse and structured decompositions of audio signals in overcomplete spaces. In: In International Conference on Digital Audio Eects (2004)
- [17] Eickhoff, S.B., Thirion, B., Varoquaux, G., Bzdok, D.: Connectivity-based parcellation: Critique and implications. *Hum Brain Mapp* (2015)
- [18] Fransson, P.: How default is the default mode of brain function? further evidence from intrinsic bold signal fluctuations. *Neuropsychologia* 44, 28362845 (2006)
- [19] Friston, K.J.: Imaging cognitive anatomy. *Trends in cognitive sciences* 1(1), 21–27 (1997)
- [20] Friston, K.: Beyond phrenology: what can neuroimaging tell us about distributed circuitry? *Annual review of neuroscience* 25(1), 221–250 (2002)
- [21] Giraud, C.: Introduction to High-Dimensional Statistics. CRC Press (2014)
- [22] Gramfort, A., Papadopoulo, T., Baillet, S., Clerc, M.: Tracking cortical activity from m/eeg using graph cuts with spatiotemporal constraints. *NeuroImage* 54(3), 1930–1941 (2011)
- [23] Harzallah, H., Jurie, F., Schmid, C.: Combining efficient object localization and image classification. In: Computer Vision, 2009 IEEE 12th International Conference on. pp. 237–244. IEEE (2009)
- [24] Hastie, T., Tibshirani, R., Wainwright, M.: Statistical Learning with Sparsity: The Lasso and Generalizations. CRC Press (2015)
- [25] Hubel, D.H., Wiesel, T.N.: Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology* 160(1), 106 (1962)
- [26] Iaria, G., Fox, C.J., Waite, C.T., Aharon, I., Barton, J.J.: The contribution of the fusiform gyrus and superior temporal sulcus in processing facial attractiveness: neuropsychological and neuroimaging evidence. *Neuroscience* 155(2), 409–22 (2008)
- [27] Jbabdi, S., Behrens, T.E.: Long-range connectomics. *Annals of the New York Academy of Sciences* 1305(1), 83–93 (2013)
- [28] Jenatton, R., Audibert, J.Y., Bach, F.: Structured variable selection with sparsity-inducing norms. *The Journal of Machine Learning Research* 12, 2777–2824 (2011)
- [29] Jenatton, R., Obozinski, G., Bach, F.: Structured sparse principal component analysis. arXiv preprint arXiv:0909.1440 (2009)
- [30] Jordan, M.I.: Frontiers in massive data analysis. National Academies Report (2015)
- [31] Kang, J.W.: Structured sparse representation of residue in screen content video coding. *Electronics Letters* 51(23), 1871–1873 (2015)
- [32] Kanwisher, N.: Functional specificity in the human brain: a window into the functional architecture of the mind. *Proceedings of the National Academy of Sciences* 107(25), 11163–11170 (2010)
- [33] Kim, S., Xing, E.P., et al.: Tree-guided group lasso for multi-response regression with structured sparsity, with an application to eqtl mapping. *The Annals of Applied Statistics* 6(3), 1095–1117 (2012)

- [34] Kim, T., Shakhnarovich, G., Urtasun, R.: Sparse coding for learning interpretable spatio-temporal primitives. In: Advances in neural information processing systems. pp. 1117–1125 (2010)
- [35] Mesulam, M.M.: From sensation to cognition. *Brain* 121, 1013–52 (1998)
- [36] Passingham, R.E., Stephan, K.E., Kotter, R.: The anatomical basis of functional localization in the cortex. *Nat Rev Neurosci* 3(8), 606–16 (2002)
- [37] Pinel, P., Thirion, B., Meriaux, S., Jobert, A., Serres, J., Le Bihan, D., Poline, J.B., Dehaene, S.: Fast reproducible identification and large-scale databasing of individual functional cognitive networks. *BMC Neurosci* 8, 91 (2007)
- [38] Raichle, M.E., MacLeod, A.M., Snyder, A.Z., Powers, W.J., Gusnard, D.A., Shulman, G.L.: A default mode of brain function. *Proc Natl Acad Sci U S A* 98(2), 676–82 (2001)
- [39] Rapaport, E., Barillot, E., Vert, J.P.: Classification of arraycgh data using fused svm. *Bioinformatics* 24(13), i375–i382 (2008)
- [40] Saygin, Z.M., Osher, D.E., Koldewyn, K., Reynolds, G., Gabrieli, J.D., Saxe, R.R.: Anatomical connectivity patterns predict face selectivity in the fusiform gyrus. *Nat Neurosci* 15(2), 321–7 (2012)
- [41] Scannell, J.W., Blakemore, C., Young, M.P.: Analysis of connectivity in the cat cerebral cortex. *J Neurosci* 15(2), 1463–83 (1995)
- [42] Seeley, W.W., Menon, V., Schatzberg, A.F., Keller, J., Glover, G.H., Kenna, H., Reiss, A.L., Greicius, M.D.: Dissociable intrinsic connectivity networks for salience processing and executive control. *J Neurosci* 27(9), 2349–2356 (2007)
- [43] Smith, S.M., Beckmann, C.F., Andersson, J., Auerbach, E.J., Bijsterbosch, J., Douaud, G., Duff, E., Feinberg, D.A., Griffanti, L., Harms, M.P., et al.: Resting-state fmri in the human connectome project. *Neuroimage* 80, 144–168 (2013)
- [44] Smith, S.M., Fox, P.T., Miller, K.L., Glahn, D.C., Fox, P.M., Mackay, C.E., Filippini, N., Watkins, K.E., Toro, R., Laird, A.R., Beckmann, C.F.: Correspondence of the brain’s functional architecture during activation and rest. *Proc Natl Acad Sci U S A* 106(31), 13040–5 (2009)
- [45] Sporns, O.: Contributions and challenges for network models in cognitive neuroscience. *Nat Neurosci* 17(5), 652–60 (2014)
- [46] Stephan, K.E., Friston, K.J., Frith, C.D.: Dysconnection in schizophrenia: from abnormal synaptic plasticity to failures of self-monitoring. *Schizophr Bull* 35(3), 509–27 (2009)
- [47] Tononi, G., Edelman, G.M., Sporns, O.: Complexity and coherency: integrating information in the brain. *Trends in cognitive sciences* 2(12), 474–484 (1998)
- [48] Vinci, G., Freeman, P., Newman, J., Wasserman, L., Genovese, C.: Estimating the distribution of galaxy morphologies on a continuous space. arXiv preprint arXiv:1406.7536 (2014)
- [49] Witten, D.M., Tibshirani, R.: A framework for feature selection in clustering. *Journal of the American Statistical Association* 105(490) (2010)
- [50] Young, M.P.: The organization of neural systems in the primate cerebral cortex. *Proc Biol Sci* 252(1333), 13–8 (1993)
- [51] Yuste, R.: From the neuron doctrine to neural networks. *Nat Rev Neurosci* 16(8), 487–497 (2015)
- [52] Zeki, S.M.: Functional specialisation in the visual cortex of the rhesus monkey. *Nature* 274(5670), 423–428 (1978)
- [53] Zilles, K., Amunts, K.: Receptor mapping: architecture of the human cerebral cortex. *Current opinion in neurology* 22(4), 331–339 (2009)