
Semi-supervised low-rank logistic regression for high-dimensional neuroimaging data

Danilo Bzdok, Michael Eickenberg, Olivier Grisel, Bertrand Thirion, Gaël Varoquaux

email: firstname.lastname@inria.fr

Abstract

Imaging neuroscience links human behavior to aspects of brain biology in ever-increasing datasets. Existing neuroimaging methods typically perform either discovery of neurobiological structure or evaluation of explicit hypotheses on mental tasks. Modelling mental tasks however hinges on the pertinence of the assumed neurobiological structure. We therefore propose to solve the unsupervised dimensionality reduction and supervised task classification in an identical statistical learning problem. We show that this approach yields more accurate and more interpretable neural models of psychological tasks in a reference neuroimaging dataset.

keywords: dimensionality reduction; semi-supervised learning; bioinformatics; fMRI; systems neuroscience

1 Introduction

Methods used in neuroimaging research can be grouped by discovering neurobiological structure or revealing the neural correlates associated with mental tasks. To discover coherent distributions of activation structure across time, independent component analysis (ICA; [6]) is often used to decompose the BOLD (blood-oxygen level-dependent) signals into the important modes of variation. The ensuing spatial activation patterns are believed to represent brain networks of functionally interacting brain regions. Similarly, sparse principal component analysis (SPCA; [20]) has been used to separate brain activity signals into parsimonious network components. Thus extracted brain networks have been shown to be manifestations of electrophysiological oscillation frequencies [13]. Their fundamental role in brain organization is attested by continued covariation during sleep and anesthesia. Network discovery is typically performed by applying ICA or SPCA on unlabeled "resting-state" data. These capture brain dynamics during ongoing random thought without controlled environmental input. The biggest fraction of the BOLD signals are known not to correlate with a particular behavior, stimulus, or experimental task.

On the other hand, to investigate the neural correlates underlying mental tasks, the general linear model (GLM; [10]) is the dominant approach. The contribution of individual brain voxels is estimated according to a design matrix of experimental tasks. Alternatively, psychophysiological interactions (PPI; [9]), elucidate the functional interactions between voxels as a function of experimental tasks. Dynamic causal modeling (DCM; [18]), in turn, quantifies directed, task-driven influences between regions by treating the brain as a nonlinear dynamic system with unknown neuronal states. As a last example, always more neuroimaging studies model experimental tasks by training classification algorithms on brain signals [16]. All these methods are applied to labeled task data that capture brain dynamics during stimulus-guided behavior. Two important conclusions can be drawn. First, the mentioned supervised neuroimaging analyses operate without exception in a single-voxel space. This ignores the fact that the BOLD signal exhibits coherent spatial activation patterns. Second, existing neuroimaging analyses do not acknowledge the fact that the task-induced changes of the BOLD signal amount to less than 5% of baseline activity [8]. They do thus not exploit the high

similarity of BOLD dynamics in the human brain at rest and during experimental tasks. Indeed, very similar brain networks were observed when applying ICA separately on rest and task data [17].

Both biological properties can be conjointly exploited in a semi-supervised (i.e., use rest and task data) low-rank (i.e., perform network decomposition) approach. The integration of brain-network discovery in a supervised classification goal should identify the neurobiological structure that allows for the best predictive models. Autoencoders suggest themselves because they can emulate variants of most nonsupervised learning algorithms, including PCA, SPCA, and ICA. Autoencoders are one-layered learning models that condense the input data to local and global representations by improving reconstruction from them [12]. They behave like a PCA in case of one linear hidden layer and a squared error loss [3]. This architecture yields a convex optimization objective with unique global minimum. Autoencoders behave like a SPCA if shrinkage terms are added for the matrix weights in the optimization objective. In turn, they behave like an ICA in case of a nonlinear convex function of the first-layer activation and tied weights [14]. These authors further demonstrated that ICA, sparse autoencoders, and sparse coding are mathematically equivalent under mild conditions. In this way, autoencoders can flexibly project the neuroimaging data onto the axes of main variation and thus reverse-engineer properties of the data-generating neural processes [15].

In the present investigation, an autoencoder will be fed by (unlabeled) rest data and integrated as bottleneck into a low-rank logistic regression fed by (labeled) task data. Using the chain rule in back-propagation, we can then solve the unsupervised data representation and a supervised classification in an identical statistical learning problem. From the perspective of dictionary learning, the first layer can be viewed as a learned set of basis functions whose linear combinations are learned in the second layer [15]. Neurobiologically, this allows delineating a low-dimensional manifold of brain network patterns and then distinguishing mental tasks by their most discriminative linear combinations. Theoretically, a big reduction of the model variance is expected by regularization using resting-state autoencoder to put probability mass on the most neurobiologically valid members of the model space. The generalization performance should consequently be improved due to much reduced Vapnik-Chervonenkis dimensions of the classification estimator. Taken together, the important modes of variation in brain dynamics and the neural correlates subserving mental operations have mostly been studied in isolation. We provide a principled computational framework to link these previously unconnected domains of systems neuroscience.

2 Methods

Data. Unsupervised projection models into a lower-dimensional space and supervised logistic-regression models are learned in concert in a same parameter gradient descent. Brain network decompositions are thus exposed that explain task-discriminative spatial patterns of neural activity.

As the currently biggest openly-accessible reference dataset, we chose the Human Connectome Project (HCP) for the present analyses [4]. Neuroimaging task data with labels of ongoing cognitive processes were drawn from 500 unrelated, healthy HCP participants. 18 HCP tasks were selected that are known to elicit reliable neural activity across participants. The task paradigms include 1) working memory/cognitive control processing, 2) incentive processing, 3) visual and somatosensory-motor processing, 4) language processing (semantic and phonological processing), 5) social cognition, 6) relational processing, and 7) emotional processing. All data were acquired on the same Siemens Skyra 3T scanner. Whole-brain EPI acquisitions were acquired with a 32 channel head coil (TR=720ms, TE=33.1ms, flip angle=52, BW=2290Hz/Px, in-plane FOV=280 × 180mm, 72 slices, 2.0mm isotropic voxels). The minimally preprocessed pipeline includes gradient unwarping, motion correction, fieldmap-based EPI distort-

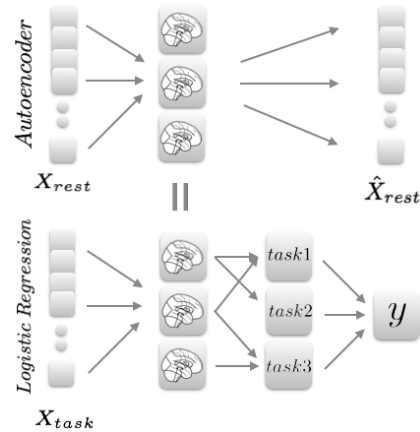


Figure 1: **Architecture**

tion correction, brain-boundary-based registration of EPI to structural T1-weighted scans, non-linear (FNIRT) registration into MNI space, and grand-mean intensity normalization. Activation maps were spatially smoothed by a Gaussian kernel of 4mm (FWHM). A GLM was implemented by FILM from the FSL suite with model regressors from convolution with a canonical hemodynamic response function and from temporal derivatives. HCP tasks were conceived to modulate activation in a maximum of different brain regions and neural systems. Indeed, at least 70% of the participants showed consistent brain activity in contrasts from the task battery, which certifies excellent coverage [4]. 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 60x72x60 space of 3mm isotropic voxels and gray-matter masked (at least 10% tissue probability). All supervised analyses were based on labeled HCP task maps of 79,941 voxels representing Z values in gray matter.

These labeled data were complemented by unlabeled activity maps from HCP acquisitions of unconstrained resting-state activity. These reflect brain activity in the absence of controlled thought. In line with the goal of the present study, acquisition of the data was specifically aimed at the study of task-rest correspondence. From each included participant, we included two time-series for left and right phase encoding with 1,200 maps of multiband, gradient-echo planar imaging acquired during a period of 15min (TR=720 ms, TE=33.1 ms, flip angle=52, FOV=280 × 180mm, and 2.0mm isotropic voxels). Besides run duration, the task acquisitions were identical to the resting-state fMRI acquisitions for maximal compatibility between task and rest data. We here drew on “minimally pre-processed” rest data from 200 randomly selected healthy, unrelated participants. PCA was applied to each set of 1,200 rest maps for denoising by keeping only the 20 main modes of variation. In sum, the HCP rest data concatenated 8000 unlabeled, noise-cleaned rest maps with 40 brain images from each of 200 randomly selected participants.

We were further interested in the utility of the optimal 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-cohort dataset. The ARCHI dataset [?] provides activity maps from diverse experimental tasks, including auditory and visual perception, motor action, reading, language comprehension and mental calculation. 81 right-handed healthy participants (3 not included in present analyses due to incomplete data) without psychiatric or neurological history participated in four fMRI sessions acquired under different experimental paradigms. The functional maps were warped into the MNI space and resampled to isotropic 3mm resolution. Whole-brain EPI data were acquired with the same Siemens Trio with a 32 channel head coil (TR=2400ms, TE=30ms, flip angle=60, in-plane FOV=19.2 × 19.2cm, 40 slices, 3.0mm isotropic voxels). Standard preprocessing was performed with Nipype [11], including slice timing, motion correction, alignment, and spatial normalization. Activation maps were spatially smoothed by a Gaussian kernel of 5mm (FWHM). Analogous to HCP data, the second task dataset incorporated 1404 labeled, grey-matter masked, and z-scored activity maps from 18 diverse tasks acquired in 78 participants.

Unsupervised layer. The labeled and unlabeled data were fed into a linear statistical model. The model is composed of an autoencoder and a low-rank logistic regression. On the one hand, the affine autoencoder takes the input x and projects it into coordinates of a latent representation z by

$$\begin{aligned} z &= W_0 x + b_{W_0} \\ x' &= W_1 z + b_{W_1} \end{aligned} \tag{1}$$

where $x \in \mathbb{R}^d$ denotes the vector of $d = 79,941$ voxel from each rest map, $z \in \mathbb{R}^n$ is the n hidden dimensions (i.e., distributed macroscopical activity patterns), and $x' \in \mathbb{R}^d$ is the reconstruction vector of the original activity map from the hidden variables. Further, W_0 denotes the weight matrix that transforms from input space into the hidden space (i.e., encoder), W_1 is the weight matrix for back-projection from the hidden space to the output space (i.e., decoder), b_{W_0} and b_{W_1} are bias vectors. Note that W_0 and W_1 are tied $W_0 = W_1^T$. The optimal model parameters W_0, b_{W_0}, b_{W_1} are found by minimizing the reconstruction difference by squared error

$$\mathcal{L}_{\mathcal{AE}}(x, x') = || -' ||^2 \tag{2}$$

This reconstruction error criterion equates with maximizing a lower bound on the mutual information between input and the learned representation. Non-linearities were not applied on the transformation results.

Supervised layer. On the other hand, such lossy compression by a low-dimensional bottleneck is also imposed by the first layer of the low-rank logistic regression.

$$P(Y = i|x, V_0, V_1, b_{V_0}, b_{V_1}) = \text{softmax}_i(V_1(V_0x + b_{V_0}) + b_{V_1}) \\ = \frac{e^{V_{1i}(V_{0i}x + b_{V_{0i}}) + b_{V_{1i}}}}{\sum_j e^{V_{1j}(V_{0j}x + b_{V_{0j}}) + b_{V_{1j}}}} \quad (3)$$

$$\mathcal{L}_{\mathcal{R}}(x, y) = \sum_{i=0}^{N_{X_{task}}} \log(P(Y = y^{(i)}|x^{(i)}; V_1(V_0x + b_{V_0}) + b_{V_1})) \quad (4)$$

Layer combination. Importantly, the optimization problem of the linear autoencoder and the low-rank logistic regression are linked one two levels. First, their transformation matrices mapping from input to the latent space are identical. $V_0 = W_0$ We thus search for a compression of the 79,941 voxel values into n latent components that represent an optimal encoding for both rest and task activity data. Second, the objectives of the autoencoder and the low-rank logistic regression are combined into a same loss function.

$$\arg \min_{\theta} \mathcal{L}(\theta, \lambda) = \lambda \frac{1}{N_{X_{task}}} \mathcal{L}_{\mathcal{R}} + (1 - \lambda) \frac{1}{N_{X_{rest}}} \mathcal{L}_{\mathcal{A}}$$

In so doing, we search for the combined set of model parameters $\theta = \{V_0, V_1, b_{V_0}, b_{V_1}, b_{V_0}, b_{V_1}\}$ with respect to the (unsupervised) reconstruction error and the (supervised) task classification. For parameter shrinkage, we impose a combination of ℓ_1 and ℓ_2 penalties (i.e., *elasticnet*) on the parameters θ .

Optimization. The statistical learning problem was approximated in the neuroimaging data by computing an updating the parameter derivatives of the semi-supervised low-rank logistic regression. The required gradients are easily obtained by using the chain rule to backpropagate error derivatives first through the decoder network and then through the encoder network. As solver, we chose *rmsprop* [19], a mini-batch version of *rprop*. This procedure dictates an adaptive learning rate for each model parameter by scaled gradients from a running average. Gradient normalization by RMSprop is known to effectively exploit curvature information. We opted for a small batch size of 100, given the high degree of redundancy in X_{rest} and X_{task} . The matrix parameters were initialized by Gaussian random values multiplied by a gain of 0.004. The bias parameters were initialized to be zero. With a slight abuse of notation, let θ denote a component of θ . The normalization factor and the update rule for θ are given by

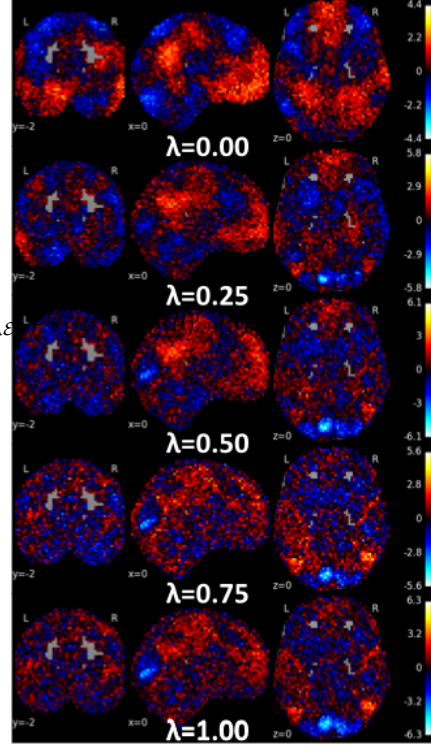


Figure 2: **Equilibrium between autoencoder and low-rank logistic regression** One learned decomposition component (out of 20) between the only-autoencoder (*uppermost panel*) and only-logistic-regression (*lowermost panel*) scenarios. The congruent structure (*red, middle column*) in the posterior cingulate cortex, posterior mid-cingulate cortex, and medial prefrontal cortex that is hence important in decompositions of the rest and task neuroimaging data. With increasing weight on the supervised learning, non-congruent structure emerges in the early (*blue*) and lateral (*red*) visual cortex (*right column*). Matrix weights were z-scored.

$$\begin{aligned}
v^{(t+1)} &= \rho v^{(t+1)} + (1 - \rho) \left(\frac{\partial f}{\partial \theta} \right)^2 \\
\theta^{(t+1)} &= \theta^{(t)} + \alpha \frac{\nabla f(\theta^{(t)})}{\sqrt{v^{(t+1)} + \epsilon}},
\end{aligned} \tag{6}$$

where $0 < \rho < 1$ constitutes the decay rate. ρ was set to 0.9 to deemphasize the magnitude of the gradient. Further α is the learning rate and ϵ a global damping factor. The hyper-parameter α was set to 0.00001 by prior manual cross-validation and ϵ was set to 1^{-6} . Note that we have also experimented with other solvers (stochastic gradient descent, adadelta, and adagrad) but found that rmsprop converged faster and with higher generalization performance.

Hints. In fact, the constraint by a rest-data autoencoder qualifies as a *hint* rather than regularization [2]. Its purpose is not to prevent overfitting but to introduce prior knowledge on known properties of the unknown target function f . Rather than relying only on input-output pairs in the learning process, we thus narrow our hypothesis set to the biologically most plausible solutions. That is, we reduce the search space in a way that is compatible with the expected representation of BOLD activity signals.

Implementation. The analyses were performed in Python. We used *nilearn* to handle the high-dimensional neuroimaging data [1] and *theano* for automatic differentiation of symbolic computation graphs [5, 7]. All Python scripts that generated the results are accessible online for reproducibility and reuse <http://github.com/banilo/nips2015>.

3 Experimental Results

This was evaluated by the prediction accuracy on a validation set (20% of the training data) at each iteration.

all vectors are column vectors

affine encoder and decoder

Low-rank regression outperformed serial ICA/SparsePCA and logistic regression.

reduction of n gray-matter voxels to n components

20 components: high bias/low variance 100 components: low bias/high variance

Generating samples from the learned statistical model

Modifications of the model that do not improve the generalization performance: drop-out, input corruption, adding non-linearities (sigmoid, tanh),

introducing a nonlinearity (sigmoid, tanh) into the system did not improve predictive accuracy but elastic did - the most useful decomposition has characteristics of a SPCA and not PCA or ICA

outperforms plain vanilla LR inject prior domain knowledge into the learning process

class weights: the right model should have high probability in only some places

4 Discussion and Conclusion

There is an increasing occurrence of high-dimensional problems in the neuroimaging domain. This calls for new statistical learning algorithms that behave well in large-cohort settings. Ideally, they should acknowledge and exploit existing widely-accepted neuroscientific knowledge. In the present work, we propose such an estimator that learns dimensionality reduction in a neurobiologically valid and interpretable fashion.

-hypothesis space includes sparse PCA and PCA but not ICA since no linearity
 respect structure in the fMRI data

- if linearity, then would be closer to the notion of 1-hidden layer neural network rather than low-rank logistic regression

We hope that these results stimulate the development of even more powerful semi-supervised classification methods

improve computational tractability, prediction accuracy, and interpretability neuroimaging datasets
 does it produce testable predictions? -¿ we can test predictive value of network-network architectures across mental domains.

open window to study the correspondence between brain dynamics during every-day mind-wandering and task-focused brain states.

classifier that operates in a (sparse) network space, rather than in a voxel space. domain-specific classification algorithms

repertoire of mental operations that the human brain can perform

large quantities of neuroimaging data logistic regression in an autoencoder paradigm

automatically learn a mapping to and from a brain-network space

statistical structure

PCA: captures structure in the data that is well described by Gaussian clouds, linear pair-wise correlations are most important form of statistical dependence/orthogonal components

-¿ BOLD images can readily be reduced to linear combinations of sparse spatial structures

scale different classification architectures to large neuroimaging data

task data might concentrate near a low-dimensional manifold of brain networks

neurobiological fidelity

entities of the neuroimaging domain

There is much uncertainty about the most pertinent representation of neural activation information
 put probability mass where we expected neurobiological structure

canonical set of brain networks

FUTURE link the resting-state component to the pattern of cognitive processes
 reduce information from resting-state data in neurological and psychiatric populations for discovery of neurobiological sub-groups as well as prediction of disease trajectories and drug responses

Acknowledgment Data were provided the Human Connectome Project. The study was supported by the German National Academic Foundation (D.B.).

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] Abu-Mostafa, Y.S.: Learning from hints. *Journal of Complexity* 10(1), 165–178 (1994)
- [3] Baldi, P., Hornik, K.: Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks* 2(1), 53–58 (1989)
- [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] 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)
- [8] Fox, D.F., Raichle, M.E.: Spontaneous fluctuations in brain activity observed with functional magnetic resonance imaging. *Nat Rev Neurosci* 8, 700–711 (2007)
- [9] Friston, K.J., Buechel, C., Fink, G.R., Morris, J., Rolls, E., Dolan, R.J.: Psychophysiological and modulatory interactions in neuroimaging. *Neuroimage* 6(3), 218–29 (1997)
- [10] Friston, K.J., Holmes, A.P., Worsley, K.J., Poline, J.P., Frith, C.D., Frackowiak, R.S.: Statistical parametric maps in functional imaging: a general linear approach. *Hum Brain Mapp* 2(4), 189–210 (1994)
- [11] Gorgolewski, K., Burns, C.D., Madison, C., Clark, D., Halchenko, Y.O., Waskom, M.L., Ghosh, S.S.: Nipype: a flexible, lightweight and extensible neuroimaging data processing framework in python. *Front Neuroinform* 5, 13 (2011)
- [12] Hinton, G.E., Salakhutdinov, R.R.: Reducing the dimensionality of data with neural networks. *Science* 313(5786), 504–507 (2006)
- [13] Hipp, J.F., Siegel, M.: Bold fmri correlation reflects frequency-specific neuronal correlation. *Curr Biol* (2015)
- [14] Le, Q.V., Karpenko, A., Ngiam, J., Ng, A.Y.: Ica with reconstruction cost for efficient overcomplete feature learning. In: *Advances in Neural Information Processing Systems*. pp. 1017–1025 (2011)
- [15] Olshausen, B.A., et al.: Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature* 381(6583), 607–609 (1996)
- [16] Poldrack, R.A., Halchenko, Y.O., Hanson, S.J.: Decoding the large-scale structure of brain function by classifying mental states across individuals. *Psychol Sci* 20(11), 1364–72 (2009)
- [17] 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)
- [18] Stephan, K.E.: On the role of general system theory for functional neuroimaging. *J Anat* 205(6), 443–70 (2004)
- [19] Tieleman, T., Hinton, G.: Lecture 6.5rmsprop: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning (2012)
- [20] Varoquaux, G., Gramfort, A., Pedregosa, F., Michel, V., Thirion, B.: Multi-subject dictionary learning to segment an atlas of brain spontaneous activity. *Information Processing in Medical Imaging* pp. 562–573 (2011)