

1

# Machine Learning for Neuroimaging with Scikit-Learn

Alexandre Abraham<sup>1,2,\*</sup>, Philippe Gervais<sup>1,2</sup>, Fabian Pedregosa<sup>1,2</sup>, Andreas Muller, Jean Kossaifi, Michael Eickenberg, Alexandre Gramfort, Bertrand Thirion<sup>1,2</sup> and Gaël Varoquaux<sup>1,2</sup>

<sup>1</sup>*Parietal Team, INRIA Saclay-Île-de-France, Saclay, France*

<sup>2</sup>*Neurospin, I<sup>2</sup>BM, DSV, CEA, 91191 Gif-Sur-Yvette, France*

Correspondence\*:

Alexandre Abraham

Parietal Team, INRIA Saclay-Île-de-France, Saclay, France,  
alexandre.abraham@inria.fr

## Research Topic

### 2 ABSTRACT

3 Statistical learning methods are increasingly used to perform neuroimaging analysis. Their  
4 main virtue for this type of application is their ability to model high-dimensional datasets, e.g.  
5 multivariate analysis of activation images, or capturing inter-subject variability. Supervised  
6 learning is typically used in decoding setting to relate brain images to behavioral or clinical  
7 observations, while unsupervised learning is typically used to uncover hidden structure in sets  
8 of images (e.g. resting state functional MRI) or to find sub-populations in large cohorts of  
9 subjects. By considering functional neuroimaging use cases, we illustrate how the Scikit-learn,  
10 a Python machine learning library, can be used to perform some key analysis steps. Scikit-learn  
11 contains a large set of statistical learning algorithms, both supervised and unsupervised, that  
12 can be applied to neuroimaging data after a proper preprocessing. Combined with other Python  
13 libraries, neuroimaging data can be loaded, processed and the results can be visualised easily.  
14 **Keywords:** Machine learning, Statistical Learning, Neuroimaging, Scikit-learn, Python

## 1 INTRODUCTION

### 1.1 SCIENTIFIC PYTHON AND NEUROIMAGING ECOSYSTEM

15 *1.1.1 Scipy and Numpy* Scipy and Numpy python packages are the basis of scientific computing in  
16 Python. They provide efficient and easy to use data representation, linear algebra, statistics, algorithms...  
17 They are the elementary bricks we use in all our algorithms.

18 *1.1.2 matplotlib* Matplotlib is a plotting library that is part of the scientific python stack. It offers a  
19 Matlab-like experience and allows to display plots, images or even 3D plots in a graphical user interface.  
20 We have used it to generate all the figures of this paper.

21 Note that a more convenient interface for matplotlib called pylab allows a procedural use of matplotlib.

22 *1.1.3 nibabel* Nibabel is a neuroimaging data loading package. Nibabel can load or save data in the  
23 most popular neuroimaging data format. This is indeed an entry point of all our scripts.

24 *1.1.4 nipy*

25 *1.1.5 scikit-learn*

## 2 SCIKIT-LEARN CONCEPTS

### 2.1 ESTIMATOR

### 2.2 DATA REPRESENTATION

26 In the scikit learn, and in the world of statistical machine learning, data are usually represented in a  
27 2-dimensional matrix of shape  $n\_samples \times n\_features$ .

### 2.3 TRANSFORMER

28 A transformer is an object that exposes a `transform` method. If the transformation can be inverted, a  
29 method called `inverse_transform` also exists.

### 2.4 CROSS VALIDATION

30 It seems more right to me to put it in this part

## 3 FROM MR VOLUMES TO A DATA MATRIX

31 As any domain specific data, MR volumes holds particular properties. Understanding them is crucial to  
32 be sure to make proper use of the data.

$$\begin{bmatrix} r_x & 0 & 0 & o_x \\ 0 & r_y & 0 & o_y \\ 0 & 0 & r_z & o_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

### 3.1 DATA PREPARATION

33 At this point, we suppose that standard preprocessings have been applied to the data. They should be  
34 registered on a common template (MNI for example). However, data is not yet ready to be processed by  
35 the scikit-learn. In fact, preprocessed data may have different shapes. Moreover, it is essential to get rid  
36 of some remaining scanner artefacts and individual trends.

37 *3.1.1 Detrending* Detrending is an essential step when dealing with fMRI data. It removes a best-fit  
38 linear trend (in the least square sense) over the time series of each voxel. It is obviously needed when you  
39 want to study the correlation between features.

40 Gal told me not to go deep into the maths, I wonder if talking about least squares is a good idea.  
Maybe I should say that a constant trend is a mean and a linear trend is simply a linear function

### 3.2 RESAMPLING

41 Resampling consists in changing the shape of the data. This is typically needed when dealing when data  
42 coming from an heterogenous dataset, as the shape depends on acquisition parameters.

43 Practically, resampling is an interpolation and thus may alterate the integrity of the data. That is why it  
44 should be used carefully. Oversampling (increasing data resolution) leads to higher memory consumption  
45 and computation resources. Downsampling is commonly used to reduce the size of the data we want to  
46 process.

47 Typical sizes are 2mm or 3mm resolutions, but the spread of high field MR scanner tends to lower these  
48 values.

- 49 • Removing confounds is necessary for some treatments

### 3.3 SIGNAL CLEANING

- 50 • Remove high frequency (scanner artefacts)

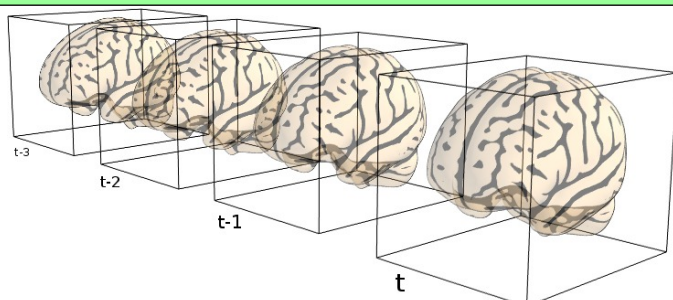
### 3.4 MASKING

3.4.1 *From 4-dimensional image to 2-dimensional array* Neuroimaging data are represented in 4 dimensions: 3 dimensions for the scans, which are positioned in a coordinate space, and one dimension for the time. Scikit-learn algorithms, on the other hand, only accept 2-dimensional data: one dimension for the features and one for the samples.

Consequently, in order to use neuroimaging data in the scikit-learn, a conversion is needed. The most simple way to achieve that would be to *flatten* the 3D scans into a 1D array. However, we know that not every voxels in a neuroimaging scan is useful. In particular, outter-brain voxels are of no use and, worse, they can bring spurious noise and scanner artefacts (such as ghosts).

To sort out voxels of interest, we will have to apply a mask on the data. Most of public datasets provide a mask, come of them even provide several, isolating different functional or anatomical brain regions.

ref to Haxby



Should tell here that some algorithms, like logistic regression, do not like colinear features.

3.4.2 *Automatically computing a mask* The simplest strategy to compute a mask is a binarization by a selected threshold. Due to the nature of the neuroimaging data, there exists some strategy to choose this threshold in order to obtain a decent segmentation.

There is a reference for the method used in Nisl. We should put it there and in the code. Add a figure with an histogram to illustrate.

Multi subject computation is simply done by intersecting subjects maps relatively to a chose threshold.

3.4.3 *Conserving geometrical structure* Applying a mask on the data obviously remove the 3-dimensional structure of the data. However, some algorithms, like the Ward, need this structural information to run.

- Speak about connectivity graphs / adjacency matrices

### 3.5 LABEL SHIFTING

Functional MRI measures brain activity by using the Blood-Oxygen-Level-Dependent contrast (BOLD). In fact, like muscles, brain regions consumes more oxygen and nutriments when stimulated. So when a part of the brain starts working, physiological mechanisms induce an oxygen-rich blood flood toward this particular region: this is called haemodynamic response.

However, this reaction takes time, usually around 6 seconds. This is the duration between the event and the reaction observed in the brain. To be able to match these two events, we will sometimes have to shift our data. This can be done by removing the two first scans of the data and the two last values of the labels (to keep an homogeneous length).

```
func = func[2:]
labels = labels[:-2]
```

## 4 DECODING

85 The process of predicting behavioral or comportamental data from fMRI scan is called decoding.

### 4.1 SVM

86 *4.1.1 Haxby dataset* For this example and the following, we will use a subset of the Haxby dataset  
 87 (Haxby et al. (2001)). Haxby dataset is from a study about face and object representation into the  
 88 brain (in particular in high level visual cortex). It is composed of 12 runs for each of the 6 subjects.  
 89 Greyscale images representing faces, houses, cats, bottles, scissors, shoes, chairs and random textured  
 90 were presented in 24 seconds blocks separated by rest periods. The repetition time (TR) between each  
 91 scan is 2.5s.

92 To make this example easier, we will work on a subset of this dataset. We will consider only one subject  
 93 and will try to classify faces versus houses.

94 *4.1.2 Feature selection: ANOVA F-Test* Even if the resolution of brain-imaging data seems low (3mm  
 95 cubes, around 100000 neurones), from a computational point of view, this is a huge. For example, Haxby  
 96 dataset has a resolution of  $64 \times 64 \times 40 = 163840$  voxels. After applying the mask, only 39912 voxels  
 97 are left, which is still high.

98 In order to reduce the number of features, we can aggregate them (in regions of interest for example) or  
 99 we can select only the most relevant ones (those who correlates most with the task). As we expect a lot of  
 100 feature to be irrelevant for our task, we opt for a feature selection method.

101 In supervised learning, the most popular feature selection method is the ANalysis Of VAriance  
 102 (ANOVA) F-Test. This is a generalization of the t-test to more than 2 features. Basically, ANOVA  
 103 compares several groups to determine if they are similar (ie randomly drawn from the same population,  
 104 this is the null hypothesis). We use it to compare the distributions of the features values across the classes.

105 `sklearn.feature_selection` contains a panel of feature selection strategies. One can choose to  
 106 take a percentile of the features (`SelectPercentile`), or a fixed number of features (`SelectKBest`)  
 107 for example. All these objects are implemented as transformers. Here we use a fixed number of features  
 108 and we use the `f_classif` function (ANOVA F-Test) for scoring.

```
109 from sklearn.feature_selection import SelectKBest, f_classif
110
111 ### Define the dimension reduction to be used.
112 # Here we use a classical univariate feature selection based on F-test,
113 # namely Anova. We set the number of features to be selected to 500
114 feature_selection = SelectKBest(f_classif, k=500)
```

115 *4.1.3 Support Vector Classifier* A Support Vector Classifier (SVC) is a simple classifier that finds a  
 116 linear hyperplane that separates the samples. Classifying a new example boils down to seeing on which  
 117 side of the hyperplane the example is. SVC has the advantage to give reliable results even when the  
 118 number of dimensions is greater than the number of samples.

119 The decision is taken based upon a subset of training data called support vectors. We can say that these  
 120 support vectors holds the information allowing to discriminate the two classes, this is why we will display  
 121 them and try to see if they match some neuroscientific knowledge.

```
122 from sklearn.svm import SVC
123 clf = SVC(kernel='linear', C=1.)
124
125 ### Look at the discriminating weights
126 svc = clf.support_vectors_
127 # reverse feature selection
128 svc = feature_selection.inverse_transform(svc)
```

129 **4.1.4 Pipeline** The workflow described above (feature selection + estimator) is a standard one. In fact,  
 130 in most cases, the workflow will consist in atomic steps *linked* together (the output of a step is the input  
 131 of the next one). For this purpose, scikit-learn offers a pipeline object that allows such linking. A pipeline  
 132 is simply a list of scikit-learn objects through which the input data will be conveyed. The function to  
 133 call for each object (transform, fit...) depends on its type. This allow the developpers to write a complete  
 134 processing as a one-liner.

```
135 from sklearn.pipeline import Pipeline
136 anova_svc = Pipeline([('anova', feature_selection), ('svc', clf)])
```

137 **4.1.5 Displaying the results**

138  
 139 Should we define a visualization function once and for all?

## 4.2 SEARCHLIGHT

- 140 • Present the Searchlight problem
- 141 • Say it is less a pain to implement thanks to scikit-learn bricks (estimator and cross\_val). Plus it is
- 142 easily customizable.

## 4.3 CLASSIFICATION OF M/EEG SENSOR SPACE DATA

### 4.4 ORTHOGONAL MATCHING PURSUIT

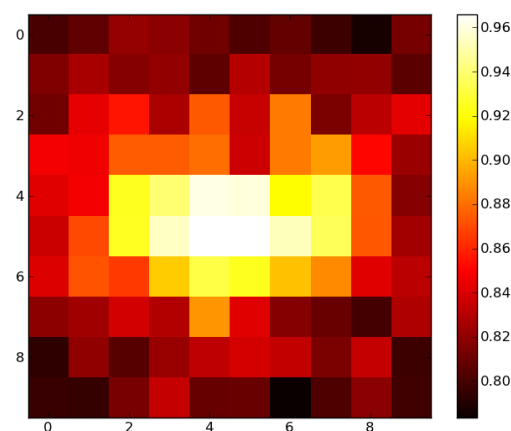
143 **4.4.1 Kamitani dataset** Kamitani dataset is based on a visual task like Haxby. In this experiment,  
 144 several series of  $10 \times 10$  binary images are presented to two subjects. Our goal will be to use the scikit-learn  
 145 to learn a correlation between brain activation and voxel color.

146 Kamitani training set is composed of random images (where black and white pixels are balanced). The  
 147 testing set is composed of structured images containing geometric shape (square, cross...) and letters  
 148 (spelling the word *neuro*).

149 There are two ways to establish a link between brain voxels and image pixels: we can either try to  
 150 reconstruct image from brain voxel activation, this is called decoding, or we can try to predict brain  
 151 activation from an image, this is called encoding.

152 In the present example, we will do both encoding and decoding and see if the results match.

153 **4.4.2 Preprocessing** Classical preprocessing (detrend...) have been applied to the data.



154 **4.4.3 Decoding: reconstructing image from brain activity**

## 5 ENCODING

After talking with Michael, he told me that he could make a fairly simple example for encoding, which I think is a plus for the paper. The example will be integrated in Nisl.

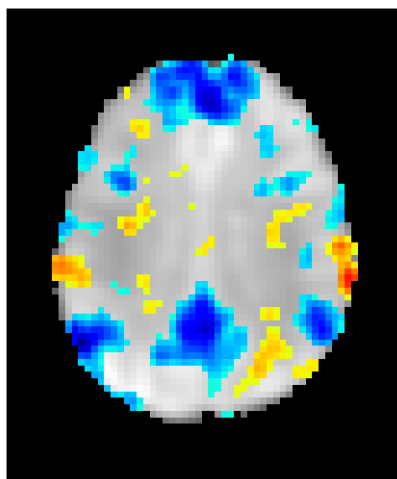
## 6 FUNCTIONAL CONNECTIVITY

Should we speak of correlation matrices to represent interaction between regions?

### 6.1 INDEPENDENT COMPONENT ANALYSIS (ICA)

**6.1.1 Intuition** ICA is a blind source separation method. Its principle is to separate a multivariate signal into several components by maximizing their non-gaussianity. A typical example is the *cocktail party problem* where ICA separates the voices of people using signal from several mikes.

It is historically the reference method to extract networks from resting state fMRI Biswal and Ulmer (1999).



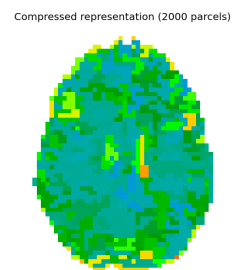
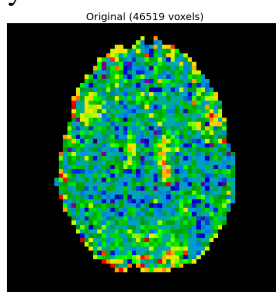
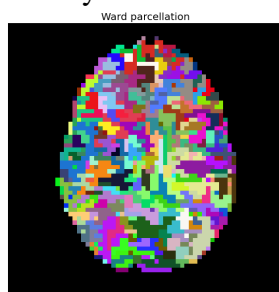
#### 6.1.2 Application

### 6.2 CLUSTERING

Make an example with Ward Clustering. Indicate then that other algorithms can be used such as KMeans and Spectral clustering and only give results.

We use a PCA here to reduce dimensionality.

Bonus: may be used as dimensionality reduction



## DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

168 The authors declare that the research was conducted in the absence of any commercial or financial  
169 relationships that could be construed as a potential conflict of interest.

## ACKNOWLEDGEMENT

170 Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text  
171 Text Text Text Text Text.  
172 *Funding:* Text Text Text Text Text Text Text Text.

## SUPPLEMENTAL DATA

173 Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text  
174 Text Text Text Text Text Text.

## REFERENCES

175 Haxby, J. V., Gobbini, I. M., Furey, M. L., Ishai, A., Schouten, J. L., and Pietrini, P. (2001) Distributed  
176 and overlapping representations of faces and objects in ventral temporal cortex. *Science* 293 2425.  
177 Biswal, B. and Ulmer, J. (1999) Blind source separation of multiple signal sources of fMRI data sets using  
178 independent component analysis. *Journal of computer assisted tomography* 23 265.