# Model Interpretation

Janaina Mourao-Miranda,

Machine Learning and Neuroimaging Lab,

University College London, UK

- In machine learning:

    Which features are driving the predictions?

- In neuroscience:

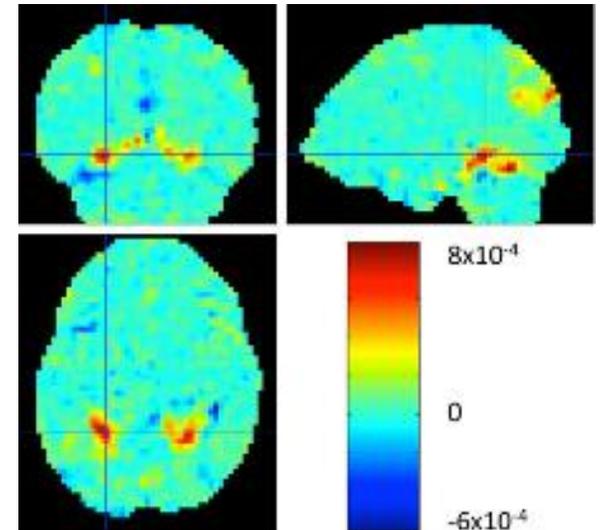    Which brain regions (or time windows) are driving the predictions?

# Linear Predictive Models

- Linear predictive models (classifier or regression) are parameterized by a weight vector **w** and a bias term $b$.

- The general equation for making predictions for a new test example $\mathbf{x}_*$ is:

$$f(\mathbf{x}_*) = \mathbf{w} \times \mathbf{x}_* + b$$

- **w** has the same dimensionality of the input data and can be plotted as an image.
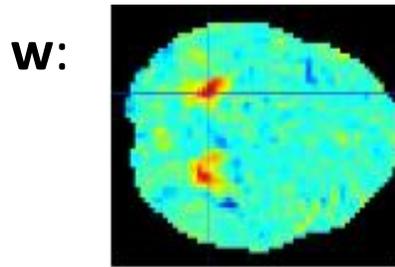
Weight map or predictive pattern
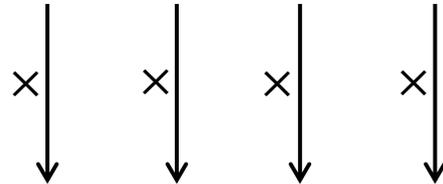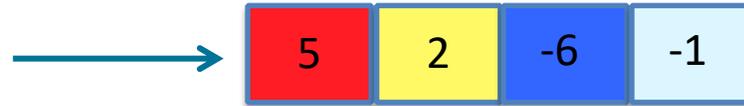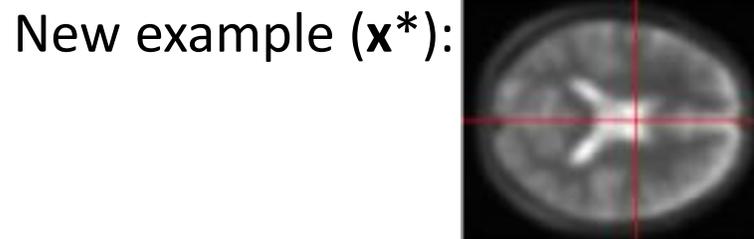


8x10⁻⁴

0

-6x10⁻⁴

(Haxby dataset, S1, Faces vs Houses)

# Predictive function

Linear predictive function:  $f(\mathbf{x}_*) = \mathbf{w} \times \mathbf{x}_* + b$

Estimated parameters:

**w**:



| 5 | 2 | -6 | -1 |

$\times$  $\times$  $\times$  $\times$

*b:*  | 0 |

New example (**x***):



| 2 | 1 | 2 | -1 |

$f(\mathbf{x}_*) = (5 \times 2) + (2 \times 1) + (-6 \times 2) + (-1 \times -1) + 0$

$f(\mathbf{x}_*) = 10 + 2 - 12 + 1 = 1$

$f(\mathbf{x}_*)$ is the predicted score for regression or the distance to the decision boundary for classification models.

# How is w computed?

- **w** is estimated by solving an optimization problem consisting of a data fit term *E* and a penalty/regularization term *J*.
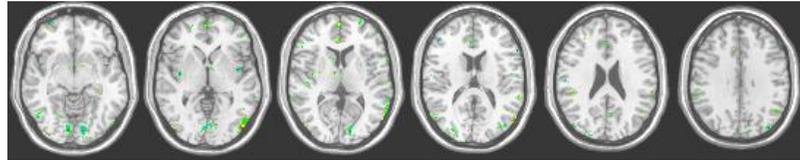
$$\min_{\mathbf{w} \hat{} \ R^p} \left\{ E(\mathbf{w}) + /\, J(\mathbf{w}) \right\}$$
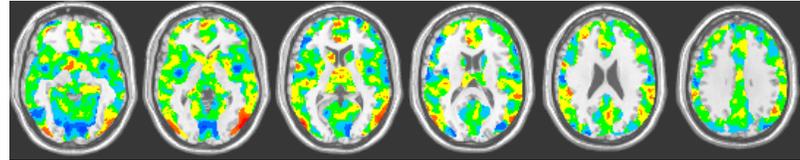
$$\downarrow$$

Regularization
parameter

- Data fit term or loss function (*E*): denotes the price we pay when we make mistakes in the predictions (e.g. squared loss, Hinge loss).

- Regularization term (*J*): favours certain properties (e.g. sparsity) and improves the generalisation over unseen examples (e.g. L2-norm, L1-norm).

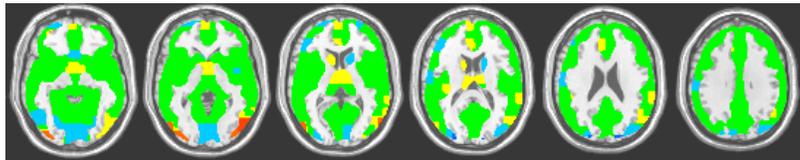- Many machine learning algorithms are particular choices of *E* and *J*.
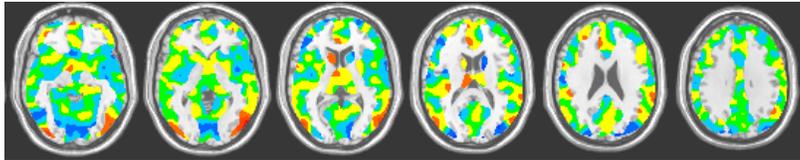
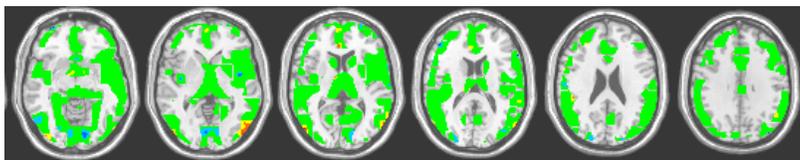# Impact of the regularization on w
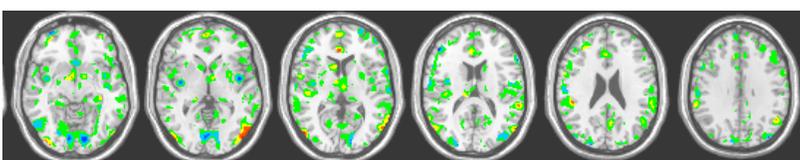


LASSO
86.31%

Elastic Net
88.02%

Total Variation (TV)
85.79%

Laplacian (LAP)
83.71%

Sparse TV
85.86%

Sparse LAP
87.05%

- Weight maps for classifying fMRI images during visualization of pleasant vs. unpleasant pictures.

- **All models used a square loss + regularization.**

Baldassarre L, Pontil M, Mourao-Miranda J (2017)

# Potential strategies to improve model interpretability

1. Feature selection (e.g. DeMartino et al. 2008; Chu et al. 2012; Rondina et al. 2013)

2. Searchlight mapping (e.g. Krigeskorte et al, 2006; Allefed and Haynes 2014)

3. Sparse algorithms (e.g. Gramfort et al. 2013; Grosenick et al. 2013; Baldassarre et al. 2017)

4. Atlas based weight summarization (e.g. Schrouff et al. 2013)

5. Multiple Kernel Learning (e.g. Schrouff et al. 2018)

6. Permutation test (e.g. Mourao-Miranda et al, 2005; Gaonkar and Davatzikos, 2013)

7. Transforming weights into activation patterns (e.g. Heufe at al., 2014)

# Potential strategies to improve model interpretability in PRoNTo

- Feature selection (e.g. DeMartino et al. 2008; Chu et al. 2012; Rondina et al. 2013)

- Searchlight mapping (e.g. Krigeskorte et al, 2006; Allefed and Haynes 2014)

1. Sparse algorithms (e.g. Gramfort et al. 2013; Grosenick et al. 2013; Baldassarre et al. 2017)

2. Atlas based weight summarization (e.g. Schrouff et al. 2013)

3. Multiple Kernel Learning (e.g. Schrouff et al. 2018)

- Permutation test (e.g. Mourao-Miranda et al, 2005; Gaonkar and Davatzikos, 2013) with extra code

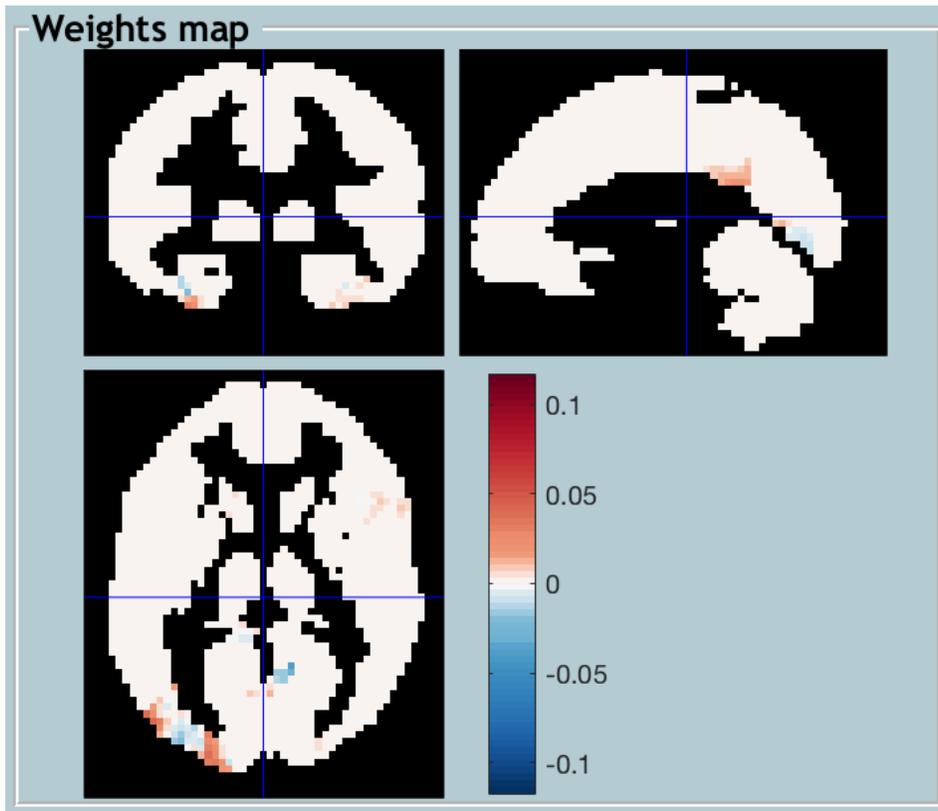- Transforming weights into activation patterns (e.g. Heufe at al., 2014)
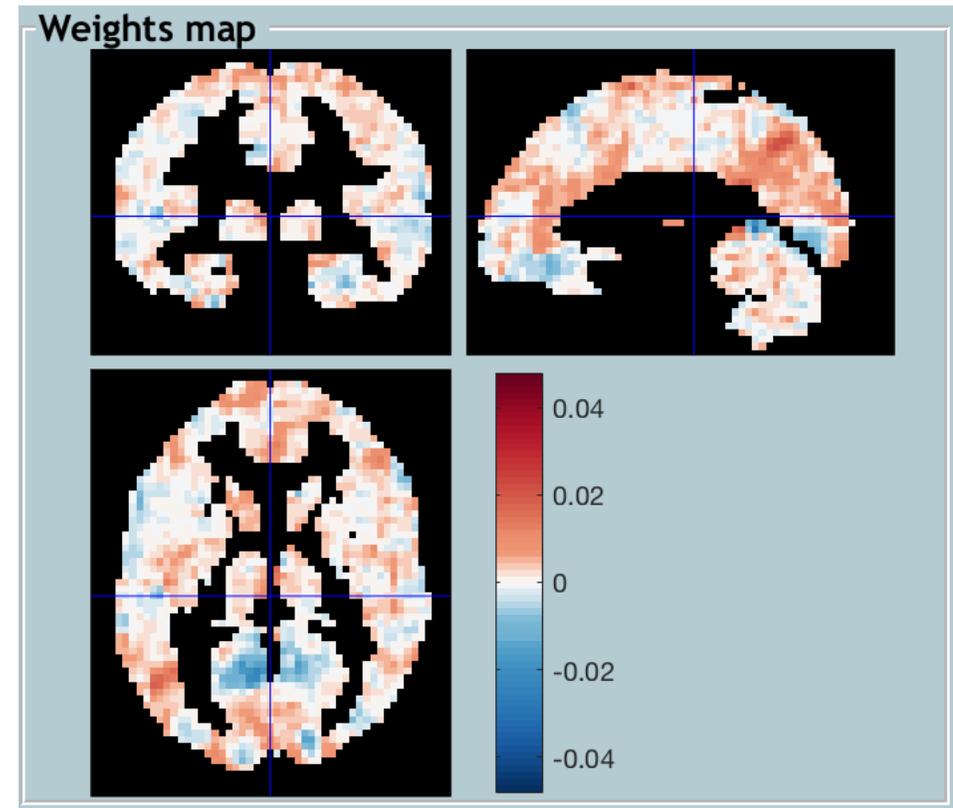
# 1- Sparse algorithms

- Include a regularization term that enforces sparsity (e.g. LASSO, Elastic-net).

- Examples of sparse algorithms in PRoNTo:
  - ✓ L1- Multi-kernel Learning
  - ✓ L1- Support Vector Machine (non-kernel)
  - ✓ L1- Logistic regression (non-kernel)

- Other sparse models can be added to PRoNTo using custom machine option (see Chapter 22 in the PRoNTo manual).

# 1- Sparse algorithms



Sparse model – L1MKL
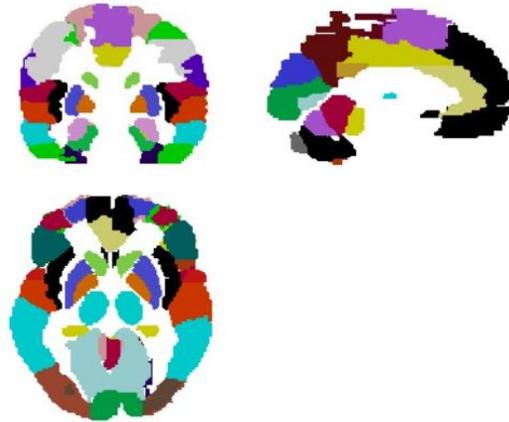Haxby data, S1, faces vs. houses

Non sparse model – SVM
Haxby data, S1, faces vs. houses

- Summarize whole brain weights using a pre-defined atlas (Schrouff et al, 2013).

- Average the absolute value of the weights within each regions to create a rank according to their contribution to the decision function.

$$NW_{ROI} = \frac{\overset{\circ}{a}_{v\hat{I} \ ROI} |w_v|}{\#\{v \ \hat{I} \ ROI\}}$$

AAL atlas (Tzourio-Mazoyer, 2002)

# 2- Atlas based summarization



Non sparse model – SVM
Haxby data, S1, faces vs. houses

Non sparse model – SVM with atlas based summarization
Haxby data, S1, faces vs. houses

• Learn simultaneously the decision function and the relative contribution of the different kernels (e.g. brain regions or time windows, Schrouff et al, 2018).

• L1-MKL (SimpleMKL, Rakotomamonjy, et al. 2008): sparsity on the kernel combination -> selects a subset of kernels (e.g. brain regions or time windows) that are relevant to the model.

Parcellation using anatomical atlas

Region 1    ...    Region 116

Multiple Kernel Learning

$$f(\mathbf{x}_*) = \sum_{m=1}^{116} d_m f_m(\mathbf{x}_*)$$

# 3- Multiple Kernel Learning



Sparse model – L1MKL
Haxby data, S1, faces vs. houses

Sparse model – L1MKL
Multimodal face data, S1, famous vs. scrambled

# Summary: weight maps interpretation



- ✓ Spatial representation of the predictive function.

- ✓ Show the contribution of each feature/voxel for the prediction.

- ✓ Multivariate pattern -> All voxels with weights different from zero contribute to the final prediction (no arbitrary threshold should be applied).

- ✓ The choice of regularization affects the sparseness and smoothness of **w.**

- ✓ Different strategies have been proposed to improve their interpretation.

# Recommended reading

- Schrouff J, Mourao-Miranda J. *Interpreting weight maps in terms of cognitive or clinical neuroscience: nonsense?* International Workshop on Pattern Recognition in Neuroimaging (PRNI) (2018b).
- Schrouff J, Monteiro JM, Portugal L, Rosa MJ, Phillips C, Mourao-Miranda J. *Embedding anatomical or functional knowledge in whole-brain multiple kernel learning models.* Neuroinformatics (2018a).
- Baldassarre, L., Pontil, M. & Mourão-Miranda, J. **Sparsity is better with stability: combining accuracy and stability for model selection in brain decoding.** Frontiers in Neuroscience: Brain Imaging Method (2017).
- Kia, S.M., Vega-Pons, S., Weisz, N. & Passerini, A. **Interpretability of Multivariate Brain Maps in Linear Brain Decoding: Definition, and Heuristic Quantification in Multivariate Analysis of MEG Time-Locked Effects.** Frontiers in Neuroscience (2017).
- Rondina J., Hahn T., de Oliveira L., Marquand A., Dresler T., Leitner T., Fallgatter A., Shawe-Taylor J. & Mourao-Miranda J. **SCoRS - a method based on stability for feature selection and mapping in neuroimaging.** IEEE Trans Med Imaging (2014).
- Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J. D., Blankertz, B. & Bießmann, F. **On the interpretation of weight vectors of linear models in multivariate neuroimaging.** NeuroImage (2014).
- Allefeld C, Haynes J-D. **Searchlight-based multi-voxel pattern analysis of fMRI by cross-validated MANOVA.** Neuroimage (2014).
- Grosenick, L., Klingenberg, B., Katovich, K., Knutson, B. & Taylor, J.E. **Interpretable whole-brain prediction analysis with GraphNet.** NeuroImage (2013)

# Recommended reading

- Schrouff, J., Cremers, J., Garraux, G., Baldassarre, L., Mourão-Miranda, J. & Phillips, C. **Localizing and comparing weight maps generated from linear kernel machine learning models.** International Workshop on Pattern Recognition in Neuroimaging (PRNI) (2013).
- Gaonkar, B. & Davatzikos, C. **Analytic estimation of statistical significance maps for support vector machine based multivariate image analysis and classification.** NeuroImage (2013).
- Gramfort, A., Thirion, B., and Varoquaux, G. **Identifying predictive regions from fMRI with TV-l1 prior.** International Workshop on Pattern Recognition in Neuroimaging (PRNI) (2013).
- Rakotomamonjy, A., Bach, F., Canu, S. & Grandvalet, Y. **SimpleMKL.** Journal of Machine Learning (2008).
- De Martino, F., Valente, G., Staeren, N., Ashburner, J., Goebel, R., & Formisano, E. **Combining multivariate voxel selection and support vector machines for mapping and classification of fMRI spatial patterns.** NeuroImage (2008).
- Kriegeskorte, N., Rainer, G. & Bandettini, P. **Information-based functional brain mapping.** *PNAS* 103 (2006).
- Mourao-Miranda J, Bokde AL, Born C, Hampel H, Stetter M. **Classifying brain states and determining the discriminating activation patterns: support vector machine on functional MRI data**. NeuroImage (2005).

# Thank you!

Questions?