#### **Detecting psychiatric disorders** with statistical learning tailored to brain activity

Gaël Varoquaux PARIETAL







#### Psychiatric diseases are a major health challenge

#### $\blacksquare$ <sup>1</sup>/<sub>5</sub> American

- ■2% are on the autistic spectrum
- 1% schizophrenia
- ■7% severe depression
- Suicide is the leading causing of death in young people (15 to 29 years)

#### Psychiatric diseases are a major health challenge

# <sup>1</sup>/<sub>5</sub> American 2% are on the autistic spectrum 1% schizophrenia 7% severe depression Suicide is the leading causing of death in young people (15 to 29 years)

Behavior, not biology, guides clinical practice eg for diagnosis and prognosis DSM-IV-TR-No progress — We need measurements

Psychiatric diseases are a major health challenge



Behavior, not biology, guides clinical practice





DSM-IV-TR

#### **1** Detecting psychiatric disorders

- **2** From activity at rest to biomarkers
- 3 Factorizing huge matrices 🔊 🍯





## 1 Detecting psychiatric disorders



#### **1** An open challenge on Autism prediction

# AutismA disease "of the mind"An imperfect diagnostic

#### The challenge Incentives: winner = 3000€ Web-based: Participant submit code Hidden test set

#### Multimodal brain data: Cortical thickness & brain activity at rest

#### **1** An open challenge on Autism prediction



#### More data is the way forward



#### **1** Brain activity at rest is a marker





Autism is a "spectrum disorder" different causes under the same symptoms

Diagnostic is imperfect

Some labels are wrong

We are probably only seeing the easy cases

#### **1** Noisy labels are still useful

# Predicting brain aging ≠ chronological age Predicts age with a mean absolute error of 4.3 years

[Liem... 2016]



#### **1** Noisy labels are still useful

Predicting brain aging ≠ chronological age
Predicts age with a mean absolute error of 4.3 years

Discrepancy with chronological age correlates with cognitive impairment



[Liem... 2016]

An individual should not be reduced to a single diagnostic or behavioral quantity

[Rahim... 2017]

1 Capturing subjects psychological traits Multi-output prediction Predict jointly multiple individual phenotypes

behavioral scores
 diagnostic status
 They improve eachother's prediction

MMSE: mini mental-state examination A diagnostic exam for Alzheimer's Disease

Adding MMSE as a target improves AD prediction



[Rahim... 2017]

#### **Detecting psychiatric disorders**

Supervised learning on rest fMRI

Across thousands of subjects

 Labels are wrong but useful to define *surrogate biomarkers* Predicting multiple dimensions of individual psychology Define traits based on "biology"

### 2 From activity at rest to biomarkers



#### No salient features in rest fMRI



#### Define functional regions



# Define functional regionsLearn interactions



Define functional regionsLearn interactionsDetect differences



#### **Defining functional regions**

#### Dividing the brain in regions

#### anatomical atlases, functional atlases, region extraction methods

Some examples



#### 2 Defining regions from rest-fMRI

Clustering k-means ward

[Thirion... 2014]









...



#### 2 Defining regions from rest-fMRI

Clustering k-means ward

[Thirion... 2014]



#### **Decomposition models**



#### 2 Defining regions from rest-fMRI

Clustering k-means ward [Thirion... 2014]



#### Decomposition models ICA: seek independent

 seek independence of maps
 Sparse dictionary learning: seek sparse maps



#### 2 Region definition: resulting parcellations







**Dictionary learning** 

**Group ICA** 







K-Means clustering

#### 2 Region definition: resulting parcellations



**Dictionary learning** 



Group ICA



Ward clustering



**K-Means clustering** 

#### 2 Region definition: resulting parcellations



**Dictionary learning** 



Group ICA



Ward clustering



K-Means clustering







# Best choice of regions for prediction Defining regions functionally is important Decomposition methods work best

[Reddy in rev, ArXiv]

#### 2 Connectome: building a connectivity matrix

How to capture and represent interactions?





#### 2 Connectome: differences across subjects



3 controls, 1 severe stroke patient Which is which?

#### 2 Connectome: differences across subjects



Spread-out variability in correlation matrices

■ Noise in partial-correlations

Strong dependence between coefficients

[Varoquaux... 2010]

2 Information geometry: uniform-error parametrization

Estimation errors given by Fisher Information matrixCovariance matrices form a manifold

 $\Rightarrow$  project to tangent space



#### 2 Connectome: which parametrization maps differences?









#### **Connectivity matrix**

CorrelationPartial correlationsTangent space



#### **Connectivity matrix**

- CorrelationPartial correlationsTangent space
- G Varoquaux

#### 2 Machine learning for connectome prediction





#### 2 Machine learning for connectome prediction



Supervised learning Linear models

#### Predicting from brain activity at rest



1. Functional regions via linear decompositions [Abraham... 2013]

2. Tangent-space reparametrization

[Varoquaux... 2010]

(logistic regression)

**3.** Supervised linear models G Varoquaux

### **3** Factorizing huge matrices

with A. Mensch, J. Mairal, B. Thirion [Mensch... 2016, 2017]



#### **Challenge: scalability**

- 1 Intuitions
- 2 Experiments
- 3 Algorithms
- 4 Proof

#### **3** Huge matrices: recommender systems



Product ratings
 Millions of entries
 Hundreds of thousands of products and users
 Large sparse matrix



#### **3 Huge matrices: brain imaging**



voxels

Brain activity at rest  $\blacksquare$  1000 subjects with  $\sim$  100–10000 samples Images of dimensionality  $> 100\,000$ Dense matrix, large both ways voxels voxels



#### **3** Stochastic optimization

$$\min_i \sum_i I(\mathbf{x}_i \mathbf{w})$$

Many samples

Gradient descent: $\mathbf{w} \leftarrow \mathbf{w} + \alpha \nabla_{\mathbf{w}} I$ Stochastic gradient descent  $\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbb{E}[\nabla_{\mathbf{w}} I]$ 

Use a cheap estimate of  $\mathbb{E}[\nabla_w I]$  (*e.g.* subsampling)



Large matrices = terabytes of data

$$\underset{\mathbf{E},\mathbf{S}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{E} \, \mathbf{S}^{\mathsf{T}} \|_{\operatorname{Fro}}^{2} + \lambda \Omega(\mathbf{S})$$



#### Large matrices = terabytes of data

$$\underset{\mathbf{E},\mathbf{S}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{E} \mathbf{S}^{\mathsf{T}}\|_{\operatorname{Fro}}^{2} + \lambda \Omega(S)$$

Large matrices = terabytes of data

$$\underset{\mathbf{E},\mathbf{S}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{E} \, \mathbf{S}^{\mathcal{T}}\|_{\operatorname{Fro}}^{2} + \lambda \Omega(\mathbf{S})$$

Rewrite as an expectation:[Mairal... 2010] $\underset{\mathbf{E}}{\operatorname{argmin}} \sum_{i} \left( \underset{\mathbf{s}}{\min} \| \mathbf{Y}_{i} - \mathbf{E} \mathbf{s}^{T} \|_{\operatorname{Fro}}^{2} + \lambda \Omega(\mathbf{s}) \right)$  $\underset{\mathbf{E}}{\operatorname{argmin}} \mathbb{E}[f(\mathbf{E})]$ 

 $\Rightarrow$  Optimize on approximations (sub-samples)





Online matrix factorization [Mairal... 2010]







Online matrix factorization [Mairal... 2010]



#### **3** Experimental results: resting-state fMRI





#### 3 Experimental results: large images



SOMF = Subsampled Online Matrix Factorization

#### **3** Experimental results: recommender system



#### **3 Algorithm: Online matrix factorization** prior art

Stream samples  $\mathbf{x}_t$ :

[Mairal... 2010]

#### 1. Compute code

$$\alpha_t = \operatorname*{argmin}_{\alpha \in \mathbb{R}^k} \|\mathbf{x}_t - \mathbf{D}_{t-1}\alpha\|_2^2 + \lambda \Omega(\alpha_t)$$

2. Update the surrogate function

$$g_t(\mathbf{D}) = \frac{1}{t} \sum_{i=1}^t \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2 = \operatorname{trace}\left(\frac{1}{2}\mathbf{D}^\top \mathbf{D}\mathbf{A}_t - \mathbf{D}^\top \mathbf{B}_t\right)$$
$$\mathbf{A}_t = (1 - \frac{1}{t})\mathbf{A}_{t-1} + \frac{1}{t}\alpha_t\alpha_t^\top \qquad \mathbf{B}_t = (1 - \frac{1}{t})\mathbf{B}_{t-1} + \frac{1}{t}\mathbf{x}_t\alpha_t^\top$$

# 3. Minimize surrogate $D_t = \underset{D \in C}{\operatorname{argmin}} g_t(D)$

$$\nabla g_t = \mathbf{D}\mathbf{A}_t - \mathbf{B}_t$$

34

#### **3 Algorithm: Online matrix factorization** prior art

Stream samples  $\mathbf{x}_t$ :

[Mairal... 2010]

#### **1.** Compute code

$$\alpha_t = \operatorname*{argmin}_{\alpha \in \mathbb{R}^k} \|\mathbf{x}_t - \mathbf{D}_{t-1}\alpha\|_2^2 + \lambda \Omega(\alpha_t)$$

2. Update the surrogate function

$$g_t(\mathbf{D}) = \frac{1}{t} \sum_{i=1}^t \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2 = \operatorname{trace}\left(\frac{1}{2}\mathbf{D}^\top \mathbf{D}\mathbf{A}_t - \mathbf{D}^\top \mathbf{B}_t\right)$$
$$\mathbf{A}_t = (1 - \frac{1}{t})\mathbf{A}_{t-1} + \frac{1}{t}\alpha_t\alpha_t^\top \qquad \mathbf{B}_t = (1 - \frac{1}{t})\mathbf{B}_{t-1} + \frac{1}{t}\mathbf{x}_t\alpha_t^\top$$

$$g_t(\mathbf{D}) \stackrel{\text{surrogate}}{=} \sum_{\mathbf{x}} l(\mathbf{x}, \mathbf{D}) \quad \alpha_i \text{ is used, and not } \alpha$$

 $\Rightarrow$  Stochastic Majorization-Minimization

No nasty hyper-parameters

**3 Algorithm: Online matrix factorization** prior art Stream samples  $\mathbf{x}_t$ : [Mairal... 2010] **1.** Compute code complexity depends on p  $\alpha_t = \operatorname{argmin} \|\mathbf{x}_t - \mathbf{D}_{t-1}\alpha\|_2^2 + \lambda \Omega(\alpha_t)$  $\alpha \in \mathbb{R}^k$ 2. Update the surrogate function  $\mathcal{O}(p)$  $g_t(\mathbf{D}) = \frac{1}{t} \sum_{i=1}^t \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2 = \operatorname{trace}\left(\frac{1}{2}\mathbf{D}^\top \mathbf{D}\mathbf{A}_t - \mathbf{D}^\top \mathbf{B}_t\right)$  $\mathbf{A}_t = (1 - \frac{1}{t})\mathbf{A}_{t-1} + \frac{1}{t}\alpha_t \alpha_t^\top \qquad \mathbf{B}_t = (1 - \frac{1}{t})\mathbf{B}_{t-1} + \frac{1}{t}\mathbf{x}_t \alpha_t^\top$ 3. Minimize surrogate  $\mathcal{O}(p)$  $\mathbf{D}_t = \operatorname*{argmin}_{\mathbf{D} \in \mathcal{C}} g_t(\mathbf{D})$  $\nabla g_t = \mathbf{D}\mathbf{A}_t - \mathbf{B}_t$ 

#### **3** Sub-sample features

- **Data stream**:  $(\mathbf{x}_t)_t \rightarrow \text{masked}$  $(\mathbf{M}_t \mathbf{x}_t)_t$
- **Dimension**:  $p \rightarrow s$
- Use only  $\mathbf{M}_t \mathbf{x}_t$  in computation  $\rightarrow$  complexity in  $\mathcal{O}(s)$





Modify all steps to work on <i>s</i> features		
Code	Surrogate	Surrogate
computation	update	minimization

#### **3** Sub-sample features

#### **Original online MF 1.** Code computation

$$\alpha_{t} = \underset{\alpha \in \mathbb{R}^{k}}{\operatorname{argmin}} \|\mathbf{x}_{t} - \mathbf{D}_{t-1}\alpha\|_{2}^{2}$$
$$+ \lambda \mathbf{O}(\alpha_{t})$$

2. Surrogate aggregation

$$\begin{split} \mathbf{A}_t &= \frac{1}{t} \sum_{i=1}^t \alpha_i \alpha_i^\top \\ \mathbf{B}_t &= \mathbf{B}_{t-1} + \frac{1}{t} (\mathbf{x}_t \alpha_t^\top - \mathbf{B}_{t-1}) \end{split}$$

3. Surrogate minimization

$$\mathbf{D}^{j} \leftarrow 
ho_{\mathcal{C}_{j}^{\prime}}^{\perp}(\mathbf{D}^{j} \!-\! rac{1}{(\mathbf{A}_{t})_{j,j}}(\mathbf{D}\mathbf{A}_{t}^{j} \!-\! \mathbf{B}_{t}^{j}))$$

#### Our algorithm

1. Approximate code computation: masked

$$\begin{split} \boldsymbol{\beta}_{t}^{(i)} &\leftarrow (1-\gamma) \mathbf{G}_{t-1}^{(i)} + \gamma \mathbf{D}_{t-1}^{\top} \mathbf{M}_{t} \mathbf{x}^{(i)} \\ \mathbf{G}_{t}^{(i)} &\leftarrow (1-\gamma) \mathbf{G}_{t-1}^{(i)} + \gamma \mathbf{D}_{t-1}^{\top} \mathbf{M}_{t} \mathbf{D}_{t-1} \\ \boldsymbol{\alpha}_{t} &\leftarrow \operatorname*{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^{k}} \frac{1}{2} \boldsymbol{\alpha}^{\top} \mathbf{G}_{t} \boldsymbol{\alpha} - \boldsymbol{\alpha}^{\top} \boldsymbol{\beta}_{t} + \lambda \, \Omega(\boldsymbol{\alpha}). \end{split}$$

2. Surrogate aggregation, averaging

$$\mathbf{A}_{t} = \frac{1}{w_{t}} \alpha_{t} \alpha_{t}^{\top} + (1 - \frac{1}{w_{t}}) \mathbf{A}_{t-1}$$
$$\mathbf{P}_{t} \mathbf{\bar{B}}_{t} \leftarrow (1 - w_{t}) \mathbf{P}_{t} \mathbf{\bar{B}}_{t-1} + w_{t} \mathbf{P}_{t} \mathbf{x}_{t} \alpha_{t}^{\top}$$

3. Surrogate minimization

 $\mathbf{P}_{t}\mathbf{D}_{t} \leftarrow \operatorname*{argmin}_{\mathbf{D}^{r}\in\mathcal{C}^{r}} \frac{1}{2} \operatorname{tr}(\mathbf{D}^{r^{\top}}\mathbf{D}^{r}\bar{\mathbf{A}}_{t}) - \operatorname{tr}(\mathbf{D}^{r^{\top}}\mathbf{P}_{t}\dot{\mathbf{B}}_{t})$  $\mathbf{P}_{t}^{\perp}\bar{\mathbf{B}}_{t} \leftarrow (1 - w_{t})\mathbf{P}_{t}^{\perp}\bar{\mathbf{B}}_{t-1} + w_{t}\mathbf{P}_{t}^{\perp}\mathbf{x}_{t}\alpha_{t}^{\top}.$ 

#### **3** Sub-sample features – variance reduction



<u>G</u> Varoquaux

36

t I

#### 3 Why does it work?

Objective:

 $\mathbf{D} = \operatorname*{argmin}_{\mathbf{D} \in \mathcal{C}} \sum_{\mathbf{x}} I(\mathbf{x}, \mathbf{D}) \quad \text{where } I(\mathbf{x}, \mathbf{D}) = \min_{\alpha} f(\mathbf{x}, \mathbf{D}, \alpha)$ 

Algorithm (online matrix factorization)

 $g_t(\mathbf{D}) \stackrel{\text{majorant}}{=} \sum_{\mathbf{x}} l(\mathbf{x}, \mathbf{D}) \quad \alpha_i \text{ is used, and not } \alpha^*$  $\Rightarrow$  Stochastic Majorization-Minimization [Mairal 2013]

#### 3 Why does it work?



#### **3** Stochastic Approximate Majorization-Minimization



#### Massive matrix factorization via subsampling

•Subsampling features  $\Rightarrow$  doubly stochastic

■10x speed ups on a fast algorithm

 Analysis via stochastic approximate majorization-minization

Conclusive on various high-dimensional problems



Detecting psychiatric disorders with statistical learning tailored to brain activity Improving psychiatry

Not a formal problem Learning to map brain to behavior More data is better



**Detecting psychiatric disorders** with statistical learning tailored to brain activity epistomology & sociology Improving psychiatry Huge data stochastic computation Factorization: costly in large-p, large-n Sub-sampling p gives huge speed ups Stochastic Approximate Majorization-Minimization https://github.com/arthurmensch/modl



Detecting psychiatric disorders<br/>with statistical learning tailored to brain activityImproving psychiatryepistomology & sociology

Huge data

Software

Scikit-learn

stochastic computation



nilearn

Putting new methods in the hands of everybody



#### **References** I

- A. Abraham, E. Dohmatob, B. Thirion, D. Samaras, and G. Varoquaux. Extracting brain regions from rest fMRI with total-variation constrained dictionary learning. In *MICCAI*, page 607. 2013.
- F. Liem, G. Varoquaux, J. Kynast, F. Beyer, S. K. Masouleh, J. M. Huntenburg, L. Lampe, M. Rahim, A. Abraham, R. C. Craddock, ... Predicting brain-age from multimodal imaging data captures cognitive impairment. *NeuroImage*, 2016.
- J. Mairal. Stochastic majorization-minimization algorithms for large-scale optimization. In *Advances in Neural Information Processing Systems*, 2013.
- J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11:19, 2010.
- A. Mensch, J. Mairal, B. Thirion, and G. Varoquaux. Dictionary learning for massive matrix factorization. In *ICML*, 2016.

#### **References II**

- A. Mensch, J. Mairal, B. Thirion, and G. Varoquaux. Stochastic subsampling for factorizing huge matrices. *IEEE Transactions on Signal Processing*, 66(1):113–128, 2017.
- M. Rahim, B. Thirion, D. Bzdok, I. Buvat, and G. Varoquaux. Joint prediction of multiple scores captures better individual traits from brain images. *Neuroimage*, in rev, 2017.
- B. Thirion, G. Varoquaux, E. Dohmatob, and J. Poline. Which fMRI clustering gives good brain parcellations? *Name: Frontiers in Neuroscience*, 8:167, 2014.
- G. Varoquaux, F. Baronnet, A. Kleinschmidt, P. Fillard, and B. Thirion. Detection of brain functional-connectivity difference in post-stroke patients using group-level covariance modeling. In *MICCAI*. 2010.