A guided tour to multivariate models in imaging-genetics

Marco Lorenzi

Asclepios Research Group, Université Côte d'Azur, Inria





Late Onset Alzheimer's Disease

Jack et al, Lancet Neurol 2010; Frisoni et al, Nature Rev Neurol 2010



time







New opportunities (and challenges)



Data science & statistical learning



Alzheimer's research & Brain imaging



Towards novel disease models and improved diagnostic

Imaging Genetics



Genotype - phenotype



DNA is the blueprint of individuals





SNPs and individual variability



Many inter-subject differences are due to single variations in DNA

A single difference at a specific nucleotide position





SNPs and individual variability



> ~10⁷ estimated SNPs in human genome

A statistical characterization of SNPs

>1% population without same nucleotide





Haplotype

Individual's combination of SNPs





profiles



What is in there?



~500k SNPs

1st and 2nd PCA components

- Avg prediction error: 310 km
- 90% within 700 km

Novembre et al, Genes mirror geography within Europe, Nature 456, 2008



What about disease susceptibility?

Pharmacological intervention

| SNP profile | Drug Response |
|-------------|---------------|
| А | Poor |
| В | Good |
| С | None |
| D | Good |

Accu-ApoE Alzheimer Test 5 minute easy to use saliva sample collection kit

USES A SAMPLE OF YOUR SALIVA SENT TO OUR LAB FOR A DNA GENETIC ANALYSIS



This sample collection kit will allow us to test for your ApoE gene variants as to the degree of possibility of Alzheimer's Disease.

Your test results will be explained in detail.

Disease understanding

| Minor Allele | ε2 | ε3 | ε4 |
|----------------------|------|-------|-------|
| General Frequency | 8.4% | 77.9% | 13.7% |
| AD Frequency | 3.9% | 59.4% | 36.7% |

main

Genetic basis of (sporadic) Alzheimer



Inría

Is it all about APOE?



Majority of current studies based on MASS UNIVARIATE ANALYSIS (GWAS)

Very large samples required (~10-100K) Low explained variability

(nría_

GWAS in Alzheimer's disease

Genetic heritability in twins up to 80%

Gatz et al, Arch Gen Psychiatry 2006

A handful of established genes APOE, CLU, PICALM, CR1, ... overall leading to small risk factor

Table 1 Genetic loci identified by the largest GWAS in AD

| Locus | SNP | OR |
|--------|--------------------------|--------------------------|
| APOE | rs2075650 | 2.53 (2.41-2.66)* |
| CLU | rs11136000 | 0.87 (0.84-0.89)* |
| PICALM | rs3851179 | 0.87 (0.84-0.90)* |
| | rs541458 | 0.87 (0.83-0.90)* |
| CR1 | rs3818361 | 1.18 (1.13-1.23)* |
| | rs6656 <mark>4</mark> 01 | 3.5 × 10 ⁻⁹ † |
| BIN1 | rs744373 | 1.15 (1.11-1.20)* |
| | Guerreiro e | et ai, |

Biochemincal Society Trans. 2011

AD missing heritability remains extensive



Imaging-genetics: multimodal analysis of heterogeneous data

- Multivariate Modeling in Imaging Genetics
- Online learning for multicentric studies
- Genetic analysis through disease progression modeling



Imaging-genetics: multimodal analysis of heterogeneous data

- Multivariate Modeling in Imaging Genetics
- Online learning for multicentric studies
- Genetic analysis through disease progression modeling



Imaging-genetics

Identifying genetic modulators of the brain phenotype

Brain imaging



Genetics

Genetic variants (Single Nucleotide Polymorphism - SNP -)











...













TA

С

...







Iterate for > 1'000'000 variants



(nría_







Iterate for > 1'000'000 variants



Iterate for > 1'000'000 image locations



...





Iterate for > 1'000'000 variants



Iterate for > 1'000'000 image locations

- Hard interpretability
- False positive discoveries
- No interaction across brain and genetic areas



...

A recent review paper (Shen and Thompson 2019)



Li Shen and Paul Thompson, Brain Imaging Genomics: Integrated Analysis and Machine Learning, Proceedings of the IEEE, 2019.



A recent review paper (Shen and Thompson 2019)



Li Shen and Paul Thompson, Brain Imaging Genomics: Integrated Analysis and Machine Learning, Proceedings of the IEEE, 2019.



Association between SNP and brain features

statistical complexity





Multivariate modeling and dimensionality reduction



Latent variable models



The building-block: linear model

 $Y = Xw + \epsilon$



$$||Y - X\mathbf{w}||^2 = (Y - X\mathbf{w})^T (Y - X\mathbf{w}) \qquad \frac{d||Y - X\mathbf{w}||^2}{d\mathbf{w}} = -2Y^T X + \mathbf{w}^T X^T X$$
$$= Y^T Y - 2Y^T X\mathbf{w} + \mathbf{w}^T X^T X \mathbf{w} \qquad \mathbf{w} = (X^T X)^{-1} X^T Y$$

119











Geladi and Kowalski, Analytica Chimica Acta, 1985









Principal Components Analysis (case X=Y)

A **variance** maximisation problem:

$$\mathbf{w} = \operatorname{argmax}_{||\mathbf{w}||=1} (X\mathbf{w})^T (X\mathbf{w})$$
$$= \operatorname{argmax}_{||\mathbf{w}||=1} \mathbf{w}^T X^T X \mathbf{w}$$
$$= \operatorname{argmax}_{||\mathbf{w}||=1} \mathbf{w}^T S_{XX} \mathbf{w}$$

$$\nabla_{w} \mathcal{L}(\mathbf{w}, \lambda) = \nabla_{w} (\mathbf{w}^{T} S_{XX} \mathbf{w} - \lambda \mathbf{w}^{T} \mathbf{w}) = 0$$
$$\mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w} = \lambda \mathbf{w}.$$



Non-linear iterative partial least squares - NIPALS Wold 1975



Random initialization **w** 2. Solve : $\operatorname{argmin}_{\mathbf{t}}||X - \mathbf{tw}^{T}||^{2}$ $\mathbf{t} = X\mathbf{w}(\mathbf{w}^T\mathbf{w})^{-1}$ 3. Normalize : $\mathbf{t} = \frac{\mathbf{t}}{||\mathbf{t}||}$ 4. Update : $\operatorname{argmin}_{\mathbf{w}} ||X - \mathbf{tw}^T||^2$ $\mathbf{w} = X^T \mathbf{t} (\mathbf{t}^T \mathbf{t})^{-1}$ 5. Iterate 2-4 until convergence

Why it works: $4 \rightarrow const \, \mathbf{w} = X^T$ $2 \rightarrow const \, \mathbf{w} = X^T X \mathbf{w}$ Then $const \, \mathbf{w} = S_{XX} \mathbf{w}$

eigen-solution of the covariance matrix $X^T X$



Canonical Correlation Analysis

A correlation maximisation problem:

$$\mathbf{w}_{x}, \mathbf{w}_{y} = \operatorname*{argmax}_{\mathbf{w}_{x}, \mathbf{w}_{y}} \rho(\mathbf{X}\mathbf{w}_{x}, \mathbf{Y}\mathbf{w}_{y})$$
$$\mathbf{w}_{x}, \mathbf{w}_{y} \qquad \qquad \mathbf{a} \qquad \mathbf{a}^{T}\mathbf{b}$$
$$\mathbf{b} \qquad \qquad \mathbf{b}$$

$$\rho(\mathbf{X}\mathbf{w}_x, \mathbf{Y}\mathbf{w}_y) = \frac{\mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_y}{\sqrt{\mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w}_x} \sqrt{\mathbf{w}_y^T \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_y}}$$


A correlation maximisation problem:



$$\mathcal{L}(\mathbf{w}_x, \mathbf{w}_y, \lambda_x, \lambda_y) = \mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_x - \lambda_x (\mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w}_x - 1) - \lambda_y (\mathbf{w}_y^T \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_y - 1)$$



$$\mathcal{L}(\mathbf{w}_x, \mathbf{w}_y, \lambda_x, \lambda_y) = \mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_x - \lambda_x (\mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w}_x - 1) - \lambda_y (\mathbf{w}_y^T \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_y - 1)$$

$$\begin{cases} \mathbf{S}_{\mathbf{X}\mathbf{Y}\mathbf{W}_{x}} = \lambda_{\mathbf{X}}\mathbf{S}_{\mathbf{X}\mathbf{X}}\mathbf{w}_{x}, \\ \mathbf{S}_{\mathbf{Y}\mathbf{X}}\mathbf{w}_{y} = \lambda_{\mathbf{Y}}\mathbf{S}_{\mathbf{Y}\mathbf{Y}}\mathbf{w}_{y} \end{cases}$$

$$\lambda_{\mathbf{X}\mathbf{W}_x^T} \mathbf{S}_{\mathbf{X}\mathbf{X}\mathbf{W}_x} = \mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{Y}\mathbf{W}_y} = \mathbf{w}_y^T \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{w}_x = \lambda_{\mathbf{Y}} \mathbf{w}_y^T \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_y$$



$$\mathcal{L}(\mathbf{w}_{x}, \mathbf{w}_{y}, \lambda_{x}, \lambda_{y}) = \mathbf{w}_{x}^{T} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_{x} - \lambda_{x} (\mathbf{w}_{x}^{T} \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w}_{x} - 1) - \lambda_{y} (\mathbf{w}_{y}^{T} \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_{y} - 1) \\ \begin{cases} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_{x} &= \lambda_{\mathbf{X}} \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w}_{x}, \\ \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{w}_{y} &= \lambda_{\mathbf{Y}} \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_{y} \end{cases} \\ \lambda_{\mathbf{X}} \mathbf{w}_{x}^{T} \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w}_{x} &= \mathbf{w}_{x}^{T} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_{y} = \mathbf{w}_{y}^{T} \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{w}_{x} = \lambda_{\mathbf{Y}} \mathbf{w}_{y}^{T} \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_{y} \\ \lambda_{\mathbf{X}} \mathbf{w}_{x}^{T} \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{w}_{x} &= \mathbf{w}_{x}^{T} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_{y} = \mathbf{w}_{y}^{T} \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{w}_{x} = \lambda_{\mathbf{Y}} \mathbf{w}_{y}^{T} \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \mathbf{w}_{y} \\ \begin{bmatrix} \mathbf{0} & \mathbf{S}_{\mathbf{X}\mathbf{Y}} \\ \mathbf{S}_{\mathbf{Y}\mathbf{X}} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{w}_{x} \\ \mathbf{w}_{y} \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{S}_{\mathbf{X}\mathbf{X}} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{\mathbf{Y}\mathbf{Y}} \end{bmatrix} \begin{bmatrix} \mathbf{w}_{x} \\ \mathbf{w}_{y} \end{bmatrix}} \qquad \begin{bmatrix} \mathsf{CCA} \text{ is solved as a} \\ \mathsf{generalized} \\ \mathsf{eigenvalue problem} \end{bmatrix}$$



Canonical Correlation Analysis (alternative formula)

$$\begin{bmatrix} \mathbf{0} & \mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbf{S}_{\mathbf{X}\mathbf{X}} \mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{-1/2} \\ \mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{-1/2} \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1/2} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v}_x \\ \mathbf{v}_y \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{v}_x \\ \mathbf{v}_y \end{bmatrix} \qquad \mathbf{v}_y = \mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{1/2T} \mathbf{w}_y$$
$$\mathbf{v}_y = \mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{1/2T} \mathbf{w}_y$$
$$\frac{\mathbf{u}^T}{\mathbf{v}_x \mathbf{v}_x \mathbf{v}_y} = \frac{\mathbf{v}_x^T \mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{-1/2} \mathbf{v}_y}{\sqrt{\mathbf{v}_x^T \mathbf{v}_x} \sqrt{\mathbf{v}_y^T \mathbf{v}_y}}$$

The quantity is maximized when \mathbf{v}_y parallel to \mathbf{u}



1 1000

Canonical Correlation Analysis (alternative formula)

$$\rho'(\mathbf{X}\mathbf{v}_x, \mathbf{X}\mathbf{v}_x) = \frac{\mathbf{u}^T \mathbf{u}}{\sqrt{\mathbf{v}_x^T \mathbf{v}_x}} = \frac{\mathbf{v}_x^T \mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{-1} \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1/2} \mathbf{v}_x}{\sqrt{\mathbf{v}_x^T \mathbf{v}_x}}$$

Eigen-solution for the matrix

$$\mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1/2}\mathbf{S}_{\mathbf{X}\mathbf{Y}}\mathbf{S}_{\mathbf{Y}\mathbf{Y}}^{-1}\mathbf{S}_{\mathbf{Y}\mathbf{X}}\mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1/2}$$



Partial Least Squares

A **co-variance** maximisation problem:

$$\mathbf{w}_x, \mathbf{w}_y = \operatorname*{argmax}_{\mathbf{w}_x, \mathbf{w}_y} cov(\mathbf{X}\mathbf{w}_x, \mathbf{Y}\mathbf{w}_y)$$

$$cov(\mathbf{X}\mathbf{w}_x, \mathbf{Y}\mathbf{w}_y) = \frac{\mathbf{w}_x^T \mathbf{S}_{\mathbf{X}\mathbf{Y}} \mathbf{w}_y}{\sqrt{\mathbf{w}_x^T \mathbf{w}_x} \sqrt{\mathbf{w}_y^T \mathbf{w}_y}}$$

Partial Least Squares

$$\begin{bmatrix} \mathbf{0} & \mathbf{S}_{\mathbf{X}\mathbf{Y}} \\ \mathbf{S}_{\mathbf{Y}\mathbf{X}} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{bmatrix}$$

The PLS problem is solved via singular value decomposition (SVD) of the covariance matrix $\mathbf{S}_{\mathbf{XY}}$

PLS and regularized CCA

$$\begin{bmatrix} \mathbf{0} & \mathbf{S}_{\mathbf{X}\mathbf{Y}} \\ \mathbf{S}_{\mathbf{Y}\mathbf{X}} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{S}_{\mathbf{X}\mathbf{X}} + \delta \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{\mathbf{Y}\mathbf{Y}} + \delta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{bmatrix} \xrightarrow{\delta \longrightarrow \infty} \mathsf{PLS}$$



Non-linear iterative partial least squares - NIPALS

scikit-learn/sklearn/cross_decomposition



Random initialization \mathbf{t}_x 2. Update \mathbf{w}_{u} : $\operatorname{argmin}_{\mathbf{w}_{y}}||Y - \mathbf{t}_{x}\mathbf{w}_{y}^{T}||^{2}$ $\mathbf{w}_y = Y^t \mathbf{t}_x (\mathbf{t}_x^T \mathbf{t}_x)^{-1}$ 3. Normalize : $\mathbf{w}_y = rac{\mathbf{w}_y}{||\mathbf{w}_y||}$ 4. $\mathbf{t}_y = Y \mathbf{w}_y$ 5. Update \mathbf{w}_x : $\begin{aligned} \operatorname{argmin}_{\mathbf{w}_x} ||X - \mathbf{t}_y \mathbf{w}_x^T||^2 \\ \mathbf{w}_x = X^T \mathbf{t}_y (\mathbf{t}_y^T \mathbf{t}_y)^{-1} \end{aligned}$ 6. Normalize : $\mathbf{w}_x = \frac{\mathbf{w}_x}{||\mathbf{w}_x||}$ 7. $\mathbf{t}_x = X \mathbf{w}_x$ Iterate 2-7 until convergence

(nría_

Geladi and Kowalski, Analytica Chimica Acta, 1985

Non-linear iterative partial least squares - NIPALS

scikit-learn/sklearn/cross_decomposition



Geladi and Kowalski, Analytica Chimica Acta, 1985

Non-linear iterative partial least squares - NIPALS Deflation

$$egin{aligned} \mathbf{X}^{(i+1)} &= \mathbf{X}^{(i)} - m{t}^{(i)} rac{m{t}^{(i)T} \mathbf{X}^{(i)}}{m{t}^{(i)T} m{t}^{(i)}}, \ \mathbf{Y}^{(i+1)} &= \mathbf{Y}^{(i)} - m{u}^{(i)} rac{m{u}^{(i)T} \mathbf{Y}^{(i)}}{m{u}^{(i)T} m{u}^{(i)}} \end{aligned}$$

Iterate until

- residual component negligible epsilon
- Difference between consecutive residual components negligible



Geladi and Kowalski, Analytica Chimica Acta, 1985

Reduced Rank Regression





 $f(\mathbf{A}, \mathbf{B}) = tr\{(\mathbf{Y} - \mathbf{X}\mathbf{A}\mathbf{B})\Gamma(\mathbf{Y} - \mathbf{X}\mathbf{A}\mathbf{B})^T\}$



Reduced Rank Regression

Solution associated to the eigen-decomposition of the matrix

$$\mathbf{R} = \Gamma^{1/2} \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \underbrace{\Gamma^{1/2}}_{\text{prior knowledge}} \overset{\text{Matrix encoding}}{\underset{\text{on Y}}{\text{prior knowledge}}}$$



Reduced Rank Regression

Solution associated to the eigen-decomposition of the matrix

$$\mathbf{R} = \Gamma^{1/2} \mathbf{S}_{\mathbf{Y}\mathbf{X}} \mathbf{S}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{S}_{\mathbf{X}\mathbf{Y}} \Gamma^{1/2} - \mathbf{Matrix encoding}$$
prior knowledge
on Y

RRR solutions:
$$\mathbf{A} = \Gamma^{-1/2} \mathbf{U}, \qquad \mathbf{B} = \mathbf{U}^T \Gamma^{1/2} \mathbf{S}_{\mathbf{YX}} \mathbf{S}_{\mathbf{XX}}^{-1}$$

Special case:

$$\Gamma = \mathbf{S}_{\mathbf{Y}\mathbf{Y}}$$
 \longrightarrow CCA

. .

Sparsity in latent variable models

$$\underset{\mathbf{w}}{\operatorname{argmin}} f(\mathbf{w}) + \lambda ||\mathbf{w}||_1$$

 $\operatorname*{argmin}_{\mathbf{w}} f(\mathbf{w}) + \lambda ||\mathbf{w}||_2^2$





Worked example: Ridge linear regression

$$\mathcal{L}(\mathbf{w}) = ||Y - X\mathbf{w}||^2 + \lambda ||\mathbf{w}||^2$$
$$\frac{d\mathcal{L}}{d\mathbf{w}} = -2Y^T X + \mathbf{w}^T X^T X + \lambda \mathbf{w}^T$$

$$\mathbf{w} = (X^T X + \lambda \boldsymbol{I} d)^{-1} X^T Y$$

Shrinkage parameter towards zero solutions

Algorithm Regularization of projections parameters \mathbf{w}_x and \mathbf{w}_y in NIPALS.

Given current estimates of \mathbf{w}_x and \mathbf{w}_y . While not converged do:

- 1. compute $\mathbf{t} = \mathbf{X}\mathbf{w}_x$,
- 2. compute $\mathbf{u} = \mathbf{Y}\mathbf{w}_y$,
- 3. compute $\overline{\mathbf{w}_x}$ by solving the Elastic-Net regression:

$$\overline{\mathbf{w}_x} = \operatorname*{arg\,min}_{\mathbf{v}} \left(\mathbf{t} - \mathbf{X}\mathbf{v}\right)^2 + \lambda_{x2} \|\mathbf{v}\|_2^2 + \lambda_{x1} \|\mathbf{v}\|_1,$$

4. compute $\overline{\mathbf{w}_y}$ by solving the Elastic-Net regression:

$$\overline{\mathbf{w}_{y}} = \operatorname*{arg\,min}_{\mathbf{v}} \left(\mathbf{u} - \mathbf{Y}\mathbf{v}\right)^{2} + \lambda_{y2} \|\mathbf{v}\|_{2}^{2} + \lambda_{y1} \|\mathbf{v}\|_{1},$$

3. Normalize
$$\overline{\mathbf{w}_x}$$
 and $\overline{\mathbf{w}_x}$,

4. Set
$$\mathbf{w}_x = \overline{\mathbf{w}_x}, \ \mathbf{w}_y = \overline{\mathbf{w}_y}.$$

S. Waaijenborg, A. H. Zwinderman, Penalized canonical correlation analysis to quantify the association between gene expression and dna markers, in: BMC proceedings, Vol. 1, BioMed Central, 2007



Group-wise penalization



$$\mathbf{a}_{1} = \{snp_{1}^{1}, snp_{2}^{1}, \dots, snp_{k}^{1}\}$$
$$\mathbf{a}_{2} = \{snp_{1}^{2}, snp_{2}^{2}, \dots, snp_{l}^{2}\}$$

$$f(\mathbf{a}) = ||\mathbf{y} - \mathbf{X}\mathbf{a}||_2^2 + \lambda \sum_{l=1}^L w_l ||\mathbf{a}_l||_2$$

Vonou et al, NeuroImage 2010; Silver et al, NeuroImage, 2012; Zhu et al, 2017; ...



Group-wise penalization

Reduced-rank regression proposed by Silver et al, 2012:

Imaging features genetic data

$$f(\mathbf{a}, \mathbf{b}) = tr\{(\mathbf{Y} - \mathbf{Xab})(\mathbf{Y} - \mathbf{Xab})^T\} + \lambda \sum_{l=1}^{L} w_l ||\mathbf{a}_l||_2$$
Mapping from genetics to latent space

| Rank | KEGG pathway name | π^{path} | Size (# SNPs) | Lasso selected genes in pathway ¹ | Known AD genes ² in pathway |
|------|--------------------------------------|--------------|------------------|--|---|
| 1. | Chemokine signaling pathway | 0.261 | 2769 | PRKCB PIK3R3 PIK3CG ADCY8 ADCY2 ITK GNAI1 XCL1 GNG2 GRK5 | CCR2 IL8 |
| 2. | Jak stat signaling pathway | 0.234 | 1311 | PIK3R3 PIK3CG IL2RA | |
| 3. | Tight junction | 0.227 | 3332 | PRKCB PRKCA YES1 ACTN1 GNAI1 CTNNA2 | |
| 4. | Insulin signaling pathway | 0.218 | 1517 | PIK3R3 PIK3CG HK2 G6PC ACACA | |
| 5. | Leukocyte transendothelial migration | 0.213 | 2289 | PRKCB PIK3R3 PRKCA PIK3CG ACTN1 ITK GNAI1 CTNNA2 | |
| 6. | Leishmania infection | 0.204 | 620 | CR1 PRKCB | CR1 |
| 7. | Calcium signaling pathway | 0.202 | 5111 | PRKCB PRKCA ADCY8 ADCY2 MYLK ATP2B2 RYR2 SLC8A1 | |













Antelmi, Ayache, Robert and Lorenzi, ICML 2019



Decoding: data reconstruction from the latent representation



Antelmi, Ayache, Robert and Lorenzi, ICML 2019



Decoding: data reconstruction from the latent representation Encoding: latent representation from the data



Antelmi, Ayache, Robert and Lorenzi, ICML 2019



Decoding: data reconstruction from the latent representation Encoding: latent representation from the data



Antelmi, Ayache, Robert and Lorenzi, ICML 2019



Decoding: data reconstruction from the latent representation Encoding: latent representation from the data



Antelmi, Ayache, Robert and Lorenzi, ICML 2019



Decoding: data reconstruction from the latent representation Encoding: latent representation from the data



Antelmi, Ayache, Robert and Lorenzi, ICML 2019



Decoding: data reconstruction from the latent representation Encoding: latent representation from the data



Antelmi, Ayache, Robert and Lorenzi, ICML 2019

Latent variable models via Variational Autoencoders

Kingma & Welling, 2014; Rezende et al. 2014

$$\mathbf{z} \longrightarrow \mathbf{X}$$
Posterior $p(\mathbf{z}|\mathbf{x}) \quad p(\mathbf{x}|\mathbf{z})$ Likelihood
$$p(\mathbf{z}|\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$
Difficult to compute
$$\underbrace{-p(\mathbf{z}|\mathbf{x})}_{-q(\mathbf{z}|\mathbf{x})} \underbrace{-p(\mathbf{z}|\mathbf{x})}_{-q(\mathbf{z}|\mathbf{x})} \underbrace{-p(\mathbf{z}|\mathbf{x})}_{-q(\mathbf{z}|$$

nnia

- 63

Latent variable models via Variational Autoencoders

Kingma & Welling, 2014; Rezende et al. 2014



$D_{KL}[q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})] = \mathbf{E}_{\mathbf{z}\sim q} \log[q(\mathbf{z}|x)] - \mathbf{E}_{\mathbf{z}\sim q} \log[p(\mathbf{z}|x)]$



C. M. Bishop, Pattern Recognition and Machine Learning, Ch.10, Ed. 2006

Latent variable models via Variational Autoencoders

Kingma & Welling, 2014; Rezende et al. 2014



 $D_{KL}[q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})] = \mathbf{E}_{\mathbf{z}\sim q} \log[q(\mathbf{z}|x)] - \mathbf{E}_{\mathbf{z}\sim q} \log[p(\mathbf{z}|x)]$

Evidence lower bound (ELBO)

$$\mathcal{L} = \mathbf{E}_{\mathbf{z} \sim q} \log[p(\mathbf{x}|\mathbf{z})] - D_{KL}[q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$$

reconstruction

regularization



C. M. Bishop, Pattern Recognition and Machine Learning, Ch.10, Ed. 2006



minimize

$$dist(q(z | X_c), p(z | X_1, X_2, ..., X_c))$$





Evidence Lower bound (ELBO)

$$\frac{1}{C} \sum_{c=1}^{C} E_{q(z|X_c)} [\sum_{i} \ln p(X_i | z)] - DKL(q(z | X_c) || p(z))$$

minimize

$$dist(q(z | X_c), p(z | X_1, X_2, ..., X_c))$$





Evidence Lower bound (ELBO)

$$\frac{1}{C} \sum_{c=1}^{C} E_{q(z|X_c)} \left[\sum_{i} \ln p(X_i | z) \right] - DKL \left(q(z | X_c) \| p(z) \right)$$

Encoding for given channel

minimize

$$dist(q(z | X_c), p(z | X_1, X_2, ..., X_c))$$





minimize I $dist(q(z | X_c), p(z | X_1, X_2, ..., X_c))$

Evidence Lower bound (ELBO)

 $\frac{1}{C} \sum_{c=1}^{C} E_{q(z|X_c)} \left[\sum_{i} \ln p(X_i | z) \right] - DKL \left(q(z | X_c) \| p(z) \right)$

Encoding for given channel Reconstruction of all channels



Antelmi, Ayache, Robert and Lorenzi, ICML 2019



minimize I $dist(q(z | X_c), p(z | X_1, X_2, ..., X_c))$

Evidence Lower bound (ELBO)

 $\frac{1}{C} \sum_{c=1}^{C} E_{q(z|X_c)} \left[\sum_{i} \ln p(X_i | z) \right] - DKL(q(z | X_c) || p(z)) :$

Encoding for given channel Reconstruction of all channels Regularization: sparsity inducing prior

[Kingma et al, NIPS, 2015; Molchanov et al, ICML 2017]



Prediction from latent space



Antelmi, Ayache, Robert and Lorenzi, ICML 2019

Generation from latent space


An example of PLS in imaging-genetics



Multivariate Association studies

Maximizing the joint relationship between genetic variants and brain features



Partial least squares (PLS) $\max_{p,q} Cov(X \cdot p, Y \cdot q)$



Liu et al, Front in Neuroinformatics, 2014; Silver et al, NeuroImage 2012; Szymczak et al, Genetic Epidemiology 2009; ...

Multivariate Association studies

Maximizing the joint relationship between genetic variants and brain features



Liu et al, Front in Neuroinformatics, 2014; Silver et al, NeuroImage 2012; Szymczak et al, Genetic Epidemiology 2009; ...



Multivariate Association studies

Maximizing the joint relationship between genetic variants and brain features



Liu et al, Front in Neuroinformatics, 2014; Silver et al, NeuroImage 2012; Szymczak et al, Genetic Epidemiology 2009; ...



Inría



















Study cohort



Healthy AD Ν 401 238 Age (years) 74.45 74.72 Sex (% females) 49 45 MMSE 29.1 23.2 Apoe4 (% 0/1/2) 72/26/2 31/48/21





Phenotype features

- Freesurfer brain cortical thickness maps (327,684 mesh points)
- Radial distance of hippocampi and amygdalae (27,120 mesh points) [Gutman et al, NeuroImage 2013]

Genotype features

• Individuals' minor allele counts for 1,167,126 SNPs in chromosomes 1 to 22

Standard quality control: MAF < 0.01, Genotype Call Rate <95%, Hardy-Weinberg Equilibrium < 1x10^{-6.} Imputation to HapMap III reference panel, quality controlled (MAF > 0.01 and R-squared > 0.3)





(nría_

Lorenzi, Altmann, Gutman, Wray, Arber, et al, PNAS, 115 (12), 2018

Investigating biological mechanisms through Meta-analysis

PLS statistical result





Investigating biological mechanisms through Meta-analysis

PLS statistical result



Querying gene annotation databases





McLaren et al. The Ensembl Variant Effect Predictor. Genome Biology, 2016



Investigating biological mechanisms through Meta-analysis



S. Wray



148 SNP-gene combinations

6 tested tissues

hippocampus, whole blood, Adipose subcutaneous, artery tibia, nerve tibial, treated fibroblast

14 Significantly expressed genes

TM2D1 (amyloid-beta binding protein), IL10RA (increase in hippo in mouse model), TRIB3

(neuronal cell death, modulates PSEN1 stability, interacts with APP)

| | Significance (p-value) | |
|--------------|------------------------|---------|
| | training | testing |
| TM2D1 | 0.005 | 0.053 |
| IL10RA | 0.107 | 0.620 |
| TRIB3 | 0.003 | 0.003 |
| ZBTB7A | 0.036 | 0.913 |
| LYSMD4 | 0.000 | 0.206 |
| CRYL1 | 0.621 | 0.118 |
| FAM135B | 0.000 | 0.559 |
| ІР6КЗ | 0.000 | 0.465 |
| ITGA1 | 0.099 | 0.731 |
| KIN | 0.001 | 0.206 |
| LAMC1 | 0.002 | 0.062 |
| LINC00941 | 0.000 | 0.690 |
| RBPMS2 | 0.000 | 0.215 |
| RP11-181K3.4 | 0.002 | 0.053 |



Imaging-genetics: multimodal analysis of heterogeneous data

- Multivariate Modeling in Imaging Genetics
- Online learning for multicentric studies
- Genetic analysis through disease progression modeling



Large multicentric clinical studies





Data for ~100'000 individuals



Big Data in medicine

Single hospital: 100s – 1'000s patients

Data from many hospitals needed



Access to multiple centers data falls into General Data Protection Regulation (GDPR): Privacy, confidentiality, security, ...

Data cannot be gathered in a single centre!

Standard learning algorithms cannot be used in multicentric data



Big Data in medicine

Circumventing the problem of data access **Federated-analysis (or meta-analysis)**





Is the association significant?



- No data sharing
- Ok for standard statistical testing (p-values, effect size)
- No complex modeling possible

































Federated analysis toolkit

A methodology for distributed



Allows a federated framework for several key statistical operations:

Data standardization, accounting for covariates, dimensionality reduction, ...



Standard statistical pipeline in multivariate analysis

- Data standardization
- Confounding effect correction
- Multivariate analysis



Federated moment estimation: Mean





Federated moment estimation: SD













$$L(Y \mid X, W) = \left\| Y - XW^T \right\|^2$$

$$W = (X^T X)^{-1} X^T Y$$





$$L(Y \mid X, W) = \left\| Y - XW^T \right\|^2$$

$$L(Y_{c} | X_{c}, W_{c}) = ||Y_{c} - X_{c}W_{c}^{T}||^{2}$$







Alternating direction method of multipliers

$$L_{\rho}(W_{c},\tilde{W},\alpha) = \sum_{c} L(Y_{c} \mid X_{c},W_{c}) + \left\langle \alpha_{c},W_{c} - \tilde{W} \right\rangle + \frac{\rho}{2} \left\| W_{c} - \tilde{W} \right\|_{2}^{2}$$


Alternating direction method of multipliers

$$L_{\rho}(W_{c},\tilde{W},\alpha) = \sum_{c} L(Y_{c} \mid X_{c},W_{c}) + \left\langle \alpha_{c},W_{c} - \tilde{W} \right\rangle + \frac{\rho}{2} \left\| W_{c} - \tilde{W} \right\|_{2}^{2}$$

Iteratively:

$$W_{c}^{(k+1)} = \operatorname{argmin}_{W_{c}} L_{\rho}(W_{c}, \tilde{W}^{(k)}, \alpha_{c}^{(k)}) = (X_{c}^{T}X_{c} + \frac{\rho}{2}Id)^{-1}(X_{c}^{T}Y_{c} - \frac{1}{2}\alpha_{c}^{(k)} + \frac{\rho}{2}\tilde{W}^{(k)})$$





Alternating direction method of multipliers

$$L_{\rho}(W_{c},\tilde{W},\alpha) = \sum_{c} L(Y_{c} \mid X_{c},W_{c}) + \left\langle \alpha_{c},W_{c} - \tilde{W} \right\rangle + \frac{\rho}{2} \left\| W_{c} - \tilde{W} \right\|_{2}^{2}$$

Iteratively:

$$\tilde{W}^{(k+1)} = \operatorname{argmin}_{\tilde{W}} L_{\rho}(W_{c}^{(k+1)}, \tilde{W}, \alpha_{c}^{(k)}) = \frac{1}{C} \sum \frac{\alpha_{c}^{(k)}}{\rho} + W_{c}^{(k+1)}$$





Alternating direction method of multipliers

$$L_{\rho}(W_{c},\tilde{W},\alpha) = \sum_{c} L(Y_{c} \mid X_{c},W_{c}) + \left\langle \alpha_{c},W_{c} - \tilde{W} \right\rangle + \frac{\rho}{2} \left\| W_{c} - \tilde{W} \right\|_{2}^{2}$$

Iteratively:

$$\alpha_{c}^{(k+1)} = \alpha_{c}^{(k)} + \rho(W_{c}^{(k+1)} - \tilde{W}^{(k+1)})$$



Results on synthetic tests





Covariance estimation and eigen-decomposition







Silva S., Gutman B., Romero E., Thompson P., Altmann A. and Lorenzi M. ISBI 2019, arXiv:1810.08553

Covariance estimation and eigen-decomposition





Silva S., Gutman B., Romero E., Thompson P., Altmann A. and Lorenzi M. ISBI 2019, arXiv:1810.08553

Testing



Mean and sd of dot product

Absolute feature-wise error





Federated analysis of subcortical brain regions in dementia

| ADNI | PPMI | UK Biobank | Miriad |
|-------------|-------------|------------|-------------|
| Alzheimer's | Parkinson's | Healthy | Alzheimer's |
| 802 | 232 | 208 | 68 |

Projection on latent components

Brain subcortical components







Meta-ImaGen

future steps

| Project ID: 10414 | k. | |
|----------------------------------|--|---|
| No license. All rights res | erved 🗢 4 Commits 🕴 1 Branch 🥏 0 Tags 🙆 225 KB File | s |
| re, you can download th | e MetalmaGen Pipeline for different platforms. | |
| member Docker is nece | ssary to properly execute the software ([See install depend | encies]) After downloading you just need to |
| scute: command to exec | tute | |
| aster – m | etaimagen-cli | History Q Find file Q - |
| README added to | all the platforms | 222a5bbr G |
| Santiago Smith aut | lored 4 months ago | inerer d |
| | | |
| README | | |
| | | |
| Name | Last commit | Last update |
| Name In Linux | Last commit README added to all the platforms | Last update 4 months ago |
| Name In Linux In Mac | Last commit README added to all the platforms README added to all the platforms | Last update 4 months ago 4 months ago |
| Name Linux Mac Windows | Last commit README added to all the platforms README added to all the platforms README added to all the platforms | Last update 4 months ago 4 months ago 4 months ago |
| Name Linux Mac Windows README.md | Last commit README added to all the platforms README added to all the platforms README added to all the platforms Initial commit | Last update 4 months ago 4 months ago 4 months ago 4 months ago 4 months ago |

metaimagen-cli

Here, you can download the MetalmaGen Pipeline for different platforms.

- · Remember Docker is necessary to properly execute the software ([See install dependencies])
- After downloading you just need to execute: command to execute

- Software freely released
- Dedicated server purchase (expected May 2019)
- Application to large scale imaging-genetics analysis



UCA IPMC, FR Poston Lab – Stanford, USA IRCCS Santa Lucia, IT UCL, UK

• Industry application: starting collaboration with Accenture Labs

Innia

Imaging-genetics: multimodal analysis of heterogeneous data

- Multivariate Modeling in Imaging Genetics
- Online learning for multicentric studies
- Genetic analysis through disease progression modeling



Modeling the natural history of neurodegeneration



Inría







Lorenzi, Filippone, Frisoni, Alexander & Ourselin, NeuroImage, 2017



- 123



Statistical disease progression model via monotonic Gaussian Processes (GP)



- Multivariate non-parametric random effects modeling
- Monotonic GP [Riihimäki & Vehtari, PMLR, 2010; Lorenzi & Filippone, ICML, 2018]
- Time reparameterization [Jedynak et al, NeuroImage 2012; Durrleman et al, IJCV, 2013; Schiratti et al, NIPS 2015]



Highlighting dynamics and relationship between biomarkers







Lorenzi, Filippone, Frisoni, Alexander & Ourselin, NeuroImage, 2017

gpprogressionmodel.inria.fr

| Try it now table_APOEposRID.csv Instructions: Data should be in csv format (comma separated) | Browse | Upload |
|---|-------------------------|--------|
| table_APOEposRID.csv Instructions: Data should be in .csv format (comma separated) | Browse | Upload |
| Instructions: | , | |
| When GP Progression Model completes the estimation the user will receive a notification with a link for downloading the Acknowledgments | e results. | |
| If you found GP Progression Model useful for your work, please cite the following papers: | | |
| Marco Lorenzi, Maurizio Filippone, Giovanni B. Frisoni, Daniel C. Alexander, Sebastien Ourselin. Probabilistic disease pr | progression modeling to | |



Disease staging from cortical amyloid and hippocampal volume







Disease staging as composite biomarker



GWAS results





Chromosome

Hippocampal volume

Amyloid burden



Disease staging



GWAS results













N. Ayache

P. Robert V. Manera



ILLINOIS INSTITU

L. Antelmi





S.S. Silva





M. Milanesio

UCA-Ville de Nice Young Researcher award



S. Garbarino







J. Banus

M. Sermesant







Hôpitaux Universitaires Genève





G.B Frisoni





EURECOM

M. Filippone

Thank you

