BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	0000000

Sparse group latent variables and penalized regression methods: a comparative study

Soufiane Ajana

Niazi Acar, Lionel Bretillon, Boris Hejblum, Hélène Jacqmin-Gadda, Cécile Delcourt

31 January 2018



université **BORDEAUX**



LEHA / Lifelong exposures, health and aging

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	000	00000	0	0000000



High dimensional issues

Predictive models

Penalized regression methods Latent variable methods

Model evaluation and comparison

Repeated double cross-validation Prediction indices

Results

Discussion

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
•0	0	000	000	00000	0	0000000

$\mathsf{BLISAR}\ \mathsf{study}$

High dimensional issues

Predictive models

Model evaluation and comparison

Results

Discussion



Objective : predict retinal n-3 PUFA status from the circulating biomarkers

- Samples of retina, plasma and red blood cells collected from human donors free of retinal diseases
- Lipid profile \Rightarrow gas chromatography (GC)
- Structural analyses of red blood cells \Rightarrow liquid chromatography mass spectrometry (LCMS)
- Circulating biomarkers of retinal n-3 PUFA status \Rightarrow 5 groups :
 - 1. Cholesteryl Esters (CE)
 - 2. Phosphatidylcholines (PC)
 - 3. Total plasma (PI)
 - 4. Red blood cells (GR) assessed by GC
 - 5. Red blood cells (GR) assessed by LMCS

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
0•	0	000	0 00	00000	0	0000000



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	•	000	0 00	00000	0	0000000

High dimensional issues

Predictive models

Model evaluation and comparison

Results

Discussion

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	•	000	0 00	00000	0	0000000

High dimensional issues (p>n)

• Multicollinearity

 $\hat{\theta}^{OLS} = (X^{\top}X)^{-1}X^{\top}Y$

 $\Rightarrow X^{\top}X \text{ is not invertible} \\\Rightarrow \text{Variance of } \hat{\theta}^{OLS} \text{ inflated}$

• Overfitting

 \Rightarrow Model complexity vs data ressources

- Data structure
 - \Rightarrow Natural grouping effect
 - \Rightarrow Select groups of variables



J : number of groups

Solution : Add constraints

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	• 00 00	0 00	00000	0	0000000

High dimensional issues

Predictive models

Penalized regression methods Latent variable methods

Model evaluation and comparison

Results

Discussion

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	• 00 00	0 00	00000	0	0000000

Penalized regression methods

With group structure :

$$\operatorname{argmin}_{\theta \in \mathbb{R}^{p}} \left\{ \left\| y - \sum_{l=1}^{L} X^{(l)} \theta^{(l)} \right\|_{2}^{2} + \lambda \sum_{l=1}^{L} [(1-\alpha)\sqrt{p_{l}} \| \theta^{(l)} \|_{2} + \alpha \| \theta \|_{1}] \right\}$$

Without group structure :

$$\operatorname{argmin}_{\theta \in \mathbb{R}^{p}} \left\{ \parallel y - X\theta \parallel_{2}^{2} + \lambda[(1 - \alpha) \parallel \theta \parallel_{2}^{2} + \alpha \parallel \theta \parallel_{1}] \right\}$$

- y : continuous outcome
- $X^{(l)}$: submatrix of X with columns corresponding to the predictors in group I
- $\theta^{(\prime)}$: coefficient vector for the group I
- $\sqrt{p_l}$: accounts for size of the group l
- λ : penalty parameter (tuning parameter)
- α : controls the combination between L1 and L2 norms (tuning parameter)

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	0000000

Penalized regression methods

With group structure :

$$argmin_{\theta \in \mathbb{R}^{p}} \left\{ \left\| y - \sum_{l=1}^{L} X^{(l)} \theta^{(l)} \right\|_{2}^{2} + \lambda \sum_{l=1}^{L} \left[(1-\alpha) \sqrt{p_{l}} \left\| \theta^{(l)} \right\|_{2} + \alpha \left\| \theta \right\|_{1} \right] \right\}$$

- $\alpha = 0 \Rightarrow$ Group Lasso (gLasso)
- $\alpha = 1 \Rightarrow Lasso$
- $0 < \alpha < 1 \Rightarrow$ Sparse Group Lasso (sgLasso)

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0	00000	0	0000000

Penalized regression methods

Without group structure :

 $argmin_{\theta \in \mathbb{R}^{p}} \left\{ \left\| y - X\theta \right\|_{2}^{2} + \lambda \left[(1 - \alpha) \right\| \theta \right\|_{2}^{2} + \alpha \left\| \theta \right\|_{1} \right] \right\}$

- $\alpha = \mathbf{0} \Rightarrow \mathsf{Ridge}$
- $\alpha = 1 \Rightarrow Lasso$
- $0 < \alpha < 1 \Rightarrow$ Elastic net

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000 •0	000	00000	0	0000000

Latent variable methods

$$\operatorname{argmax}_{u \in \mathbb{R}^{p}, \|u\|_{2}=1} \left\{ \operatorname{cov}(X^{(l)}u^{(l)}, y) - \lambda[(1-\alpha)\sum_{l=1}^{L} \sqrt{p^{(l)}} \| u^{(l)} \|_{2} + \alpha \| u \|_{1}] \right\}$$

- y : continuous outcome
- $X^{(l)}$: submatrix of X with columns corresponding to the predictors in group I
- $u^{(l)}$: loading vector for the group I
- $\sqrt{p^{(l)}}$: accounts for size of the group I
- λ : penalty parameter (tuning parameter)
- α : controls the combination between L1 and L2 norms (tuning parameter)

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000 0•	0 00	00000	0	0000000

Summary



Model	High correlation	Group structure
Lasso		
Group Lasso (gLasso)		x
Sparse Group Lasso (sgLasso)		x
Elastic net	х	
Sparse Partial Least Squares (sP	LS) X	
Groupe Partial Least Squares (gP	LS) X	x
Sparse Groupe Partial Least Squa	ares (sgPLS) X	x

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	00	00000	0	0000000

High dimensional issues

Predictive models

Model evaluation and comparison

Repeated double cross-validation Prediction indices

Results

Discussion

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	•	00000	0	0000000

Repeated double cross-validation



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	○ ●O	00000	0	0000000

Prediction indices

Mean Squared Error of Prediction (MSEP) :

$$MSEP = rac{\sum_{i=1}^{n_{test}} (y_i - \hat{y}_i)^2}{n_{test}}$$

• MSEP =
$$0 \Rightarrow$$
 Perfect prediction

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	○ ○●	00000	0	0000000

```
Prediction indices
```

Goodness of fit (R^2) :

$$R^2 = cor(y_{test}, \hat{y}_{test})^2$$

• $R^2 = 1 \Rightarrow$ Perfect prediction

•
$$R^2 \searrow \Rightarrow$$
 Prediction performance \searrow

Important note : R^2 calculated via cross-validation on the test data, assesses the quality of predictions on independent (unseen) sets

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	000	•0000	0	0000000

High dimensional issues

Predictive models

Model evaluation and comparison

Results

Discussion

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	•0000	0	0000000

Comparison of the multivariable regression methods for 10 random divisions with 100 runs (N=46, p=332)

Method	Test data R ² (SD)	Test data RMSEP (SD)	Number of selected variables*	Selected groups*
Lasso	0.14 (0.05)	2.73 (0.14)	4	CE, PC
sgLasso	0.20 (0.05)	2.72 (0.16)	143	CE, PC, LCMS
gLasso	0.21 (0.05)	2.69 (0.15)	285	CE, PC, PI, LCMS
Elastic net	0.18 (0.05)	2.65 (0.12)	23	CE, PC, PI, LCMS
SPLS	0.36 (0.03)	2.32 (0.05)	8	CE, PI
gPLS	0.30 (0.03)	2.43 (0.05)	32	CE
sgPLS	0.38 (0.02)	2.27 (0.04)	7	CE

* In at least 60% of the samples



Venn diagrams of the variables selected by latent variable and penalized regression methods





Boxplots of the difference in RMSEP between ${\bf sgPLS}$ and the other six methods over the 100 runs





Lipids frequency selection with sgPLS over 100 runs. The vertical dashed line correspond to a frequency selection of 60 %





Lipids frequency selection with sgLasso over 100 runs. The vertical dashed line correspond to a frequency selection of 60 %



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	000	00000	•	0000000

High dimensional issues

Predictive models

Model evaluation and comparison

Results

Discussion

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	•	0000000

Discussion

- Prediction performance \Rightarrow sPLS and sgPLS
- Variable selection \Rightarrow sgPLS (only one group selected)
 - \Rightarrow sparse data with some relevant predictors
 - \Rightarrow biologically relevant predictors
 - \Rightarrow lower costs
- Generally, Latent variables methods outperformed penalized regression methods in this application
- Small sample size + High correlations + Group structure \Rightarrow sgPLS +++

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	000	00000	0	•000000

High dimensional issues

Predictive models

Model evaluation and comparison

Results

Discussion

BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	•000000

Lasso



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	000000

Elastic net



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	0000000

gLasso



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	0000000

sPLS



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	0000000

 gPLS



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	000	00000	0	0000000

AdaLasso



BLISAR study	High dimensional issues	Predictive models	Model evaluation and comparison	Results	Discussion	Annexes
00	0	000	0 00	00000	0	000000

AdagLasso

