

# Une PLS multi-blocs parcimonieuse pour données hétérogènes incomplètes

Hadrien Lorenzo<sup>1</sup>, Jérôme Saracco<sup>2</sup>, Rodolphe Thiébaud<sup>1</sup>

<sup>1</sup>SISTM (Inserm, U1219, Bordeaux Population Health and Inria, Talence, France) and Vaccine Research Institute, Creteil, France.

<sup>2</sup>CQFD (INRIA Bordeaux Sud-Ouest, France), CNRS (UMR5251)

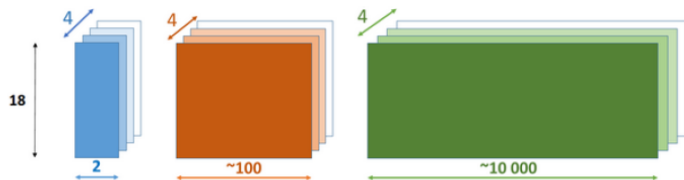
JDS 2018, 19 novembre 2018



# Motivation

## Essai rVSV-ZEBOV Ebola de phase 1 avec doses échelonnées

- Premier vaccin à présenter une efficacité depuis la survenue de la maladie [HENAQ-RESTREPO et al., *The Lancet*, 2017 ]



Réponse  
anticorps

Jours 28, 56, 84, 180

Fonctionnalité  
cellulaire

Jours 0, 1, 3, 7

Expression  
génétique

Jours 0, 1, 3, 7

## Echantillons manquants : données génétiques

$t_1$															
$t_2$															
$t_3$															
$t_4$															

TABLE – Missing path du dataset Ebola rVSV-ZEBOV RNA-Seq où  $t_1 = jour_0$ ,  $t_2 = jour_1$ ,  $t_3 = jour_3$  et  $t_4 = jour_7$ . Colonnes pour les participants.

- ▶ 30% de données/échantillons manquants,
- ▶ Lien "Missing structure"/"time structure"

## Objectif

Prédire la réponse anticorps de façon parcimonieuse en gérant efficacement les données manquantes

## Modèle général

Combinent des alternances d'estimation :

0. Initialiser les valeurs pour les données manquantes,
1. Estimer une factorisation des données complétées,
2. Estimer les données manquantes,
3. Recommencer en 1. jusqu'à convergence.

... en attente de stabilisation.

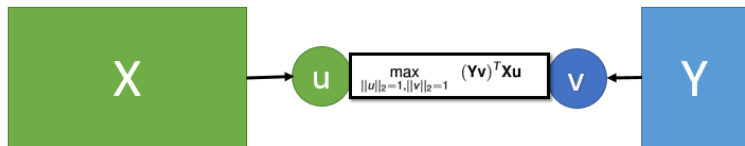
—> D'autant plus vrai dans le cas de modèles parcimonieux.

Côté utilisateur : difficile à optimiser

## Contrainte majeure

Très peu d'individus : la stabilisation est plus difficile à trouver.

# Approches PLS [Wold père et fils, 1983]



Équivalent à une recherche de sous-espaces propres (SVD). On appelle :

- ▶ **Poids** ou **weights** ou **loadings**  $u$  et  $v$  : importance donnée d'une variable de  $X$ , via  $u$ , et de  $Y$ , via  $v$ .
- ▶ **Scores** ou **variates**  $Xu$  et  $Yv$  : projections de  $X$  et de  $Y$  dans les sous-espaces définis par  $u$  et  $v$ .

⇒ Rechercher dans  $X$  l'information qui est très liée à  $Y$ .

# Résolution du problème de PLS

Utilisation du formalisme lagrangien :

$$\max_{u,v,\alpha_x,\alpha_y} (\mathbf{Y}v)^T \mathbf{X}u - \alpha_x/2(\|u\|_2^2 - 1) - \alpha_y/2(\|v\|_2^2 - 1),$$

$\mathbf{X}$  et  $\mathbf{Y}$  les matrices échantillons, centrées, des covariables et des variables à prédire.  $\alpha_x$  et  $\alpha_y$  les coefficients de Lagrange. Alors :

**Système :**

$$\begin{cases} \partial_{u.} : & \alpha_x u = \mathbf{X}^T \mathbf{Y} v \\ \partial_{v.} : & \alpha_y v = \mathbf{Y}^T \mathbf{X} u \\ \partial_{\alpha_x.} : & \|u\|_2^2 = 1 \\ \partial_{\alpha_y.} : & \|v\|_2^2 = 1 \end{cases}$$

**Optimisation :**

1.  $u \leftarrow \mathbf{X}^T \mathbf{Y} v$
2.  $u \leftarrow u / \|u\|_2$
3.  $v \leftarrow \mathbf{Y}^T \mathbf{X} u$
4.  $v \leftarrow v / \|v\|_2$

**Régression :**

$$\mathbf{Y} \approx \mathbf{X} \mathbf{B}$$
$$\mathbf{B} = \frac{v^T \mathbf{Y}^T \mathbf{X} u}{\| \mathbf{X} u \|_2^2} u v^T$$

**Classification :**

LDA sur  $(\mathbf{X}u, \mathbf{Y})$

## Matrice de variance-covariance

Elle est au centre des approches PLS, via  $\mathbf{Y}^T \mathbf{X}$ !

# La sélection de variables en PLS → sparse PLS

## Principe, intérêt et pistes explorées

- ▶ Peu de mesures biologiques nécessaires en prédiction.
- ▶ Pénalisations  $\mathcal{L}_1$  des poids  
⇒ Sélection des variables et régularisation des données.

## Des PLS parcimonieuses

- ▶ [LÊ CAO et al., 2008 ], **2 paramètres/axe** :

$$\min_{u,v} \|\mathbf{Y}^T \mathbf{X} - v u^T\|_F^2 + \lambda_x \|u\|_1 + \lambda_y \|v\|_1$$

- ▶ [CHUN et KELEŞ, 2010 ],  $M = \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}$ , **3 paramètres/axe** :

$$\min_{w,c} -\kappa w^T M w + (1 - \kappa)(c - w)^T M (c - w) + \lambda_1 \|c\|_1 + \lambda_2 \|c\|_2$$

subj. to  $w^T w = 1$ ,

# Data Driven sPLS (dd-sPLS)

## Idée

Travailler directement sur la matrice de variance-covariance en effaçant les liens entre les variables de X et de Y.

avec

**Solution utilisée :**  
Seuillage doux,  $S_\lambda$ , de  $\frac{\mathbf{Y}^T \mathbf{X}}{n-1}$ .

$$S_\lambda(t) = \begin{cases} 0 & \text{si } |t| < \lambda \\ t - \lambda & \text{si } t \geq \lambda \\ t + \lambda & \text{si } t \leq -\lambda \end{cases}$$

## Intérêts

- ▶ Sélectionner sur X et sur Y avec un seul paramètre,  $\lambda$ ,
- ▶ Considérer un problème très bien connu : SVD

## dd-sPLS

$$\max_{\substack{\mathbf{u} \in \mathbb{R}^{p \times R} \\ \mathbf{u}^T \mathbf{u} = \mathbb{I}_R}} \left\| S_\lambda \left( \frac{\mathbf{Y}^T \mathbf{X}}{n-1} \right) \mathbf{u} \right\|_F^2,$$



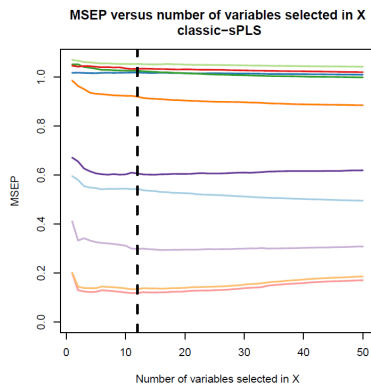


# Application : Liver Toxicity Dataset via classic sPLS

Voir [heinloth2004gene]. 64 Souris droguées. Expression de RNA. 10 variables réponses restituant l'état du foie,  $\mathbf{X}_{64 \times 3116}$  and  $\mathbf{Y}_{64 \times 10}$ .

Comparaison **Classique sPLS** / **dd-sPLS**

- ▶ Paramètre de parcimonie en  $\mathbf{Y}$  fixé arbitrairement à 2,
- ▶ Minimum d'erreur pour 12 covariables sélectionnées.



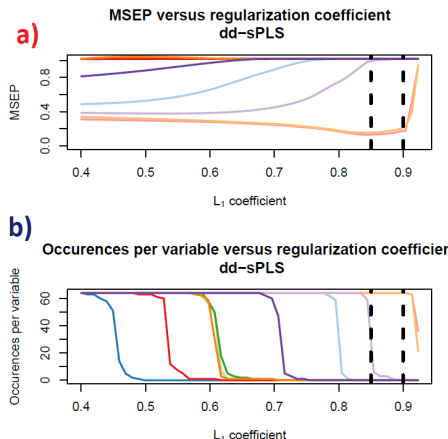
# Application : Liver Toxicity Dataset via dd-sPLS

Deux visualisations disponibles :

- ▶ **a)** : L'erreur en prédiction,
- ▶ **b)** : Le nombre d'occurrences de chaque Y dans les modèles de validation

Observations :

- ▶ Via **a)** ,  $\lambda \approx 0.85$  : 2 variables Y sélectionnées (?),
- ▶ Via **b)** :  $\lambda \approx 0.9$  exactement 2 variables Y sélectionnées.



# Liver Toxicity : Variables sélectionnées dans X

Classic-sPLS $keep_X = 12$		dd-sPLS			
		$\lambda = 0.85$		$\lambda = 0.9$	
Var	Coeff	Var	Coeff	Var	Coeff
A_43_P11724	0.172	A_43_P11724	0.172	<b>A_43_P14131</b>	<b>-0.862</b>
A_42_P802628	-0.117	A_42_P705413	-0.026	<b>A_42_P620915</b>	<b>-0.507</b>
A_43_P10606	-0.14	A_42_P802628	-0.117		
<b>A_43_P14131</b>	<b>-0.6</b>	A_43_P10606	-0.14		
A_42_P675890	-0.175	A_43_P22616	-0.012		
A_43_P23376	-0.213	<b>A_43_P14131</b>	<b>-0.6</b>		
<b>A_42_P620915</b>	<b>-0.515</b>	A_42_P675890	-0.175		
A_42_P758454	-0.175	A_43_P23376	-0.213		
A_42_P578246	-0.143	<b>A_42_P620915</b>	<b>-0.515</b>		
A_43_P17415	-0.331	A_42_P758454	-0.175		
A_42_P610788	-0.072	A_42_P578246	-0.143		
A_42_P840776	-0.264	A_43_P17415	-0.331		
		A_42_P610788	-0.072		
		A_42_P840776	-0.264		

TABLE – Comparaison des résultats de l'analyse du jeu de données de Liver Toxicity

# Cas multi-blocs : mdd-sPLS avec données manquantes

Alternance de 2 étapes :

**Construction du modèle** et estimation des données manquantes.

Construction du modèle : Solution en deux étapes

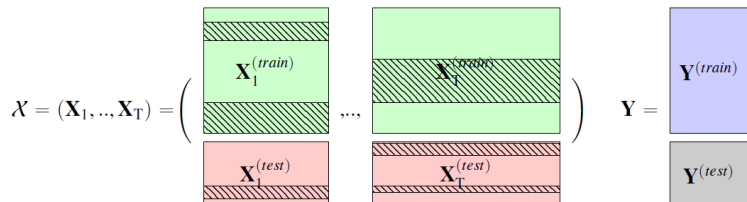
$$\max_{u_t^{(r)}, \beta_t^{(r)}} \sum_{t=1}^T \sum_{r=1}^R \beta_t^{(r)2} \left\| S_\lambda \left( \frac{\mathbf{Y}^T \mathbf{X}_t}{n-1} \right) u_t^{(r)} \right\|_2^2 \quad (1)$$

1.  $\forall t = 1..T, (u_t^{(1)}, \dots, u_t^{(R)}) = \arg \max_{\mathbf{u}^T \mathbf{u} = \mathbb{I}_R} \left\| S_\lambda \left( \frac{\mathbf{Y}^T \mathbf{X}_t}{n-1} \right) \mathbf{u} \right\|_F^2,$
2. Résolution de (1) via une SVD de R composantes

Estimation des données manquantes

Ceci grâce au modèle précédemment construit et à ce qui suit...

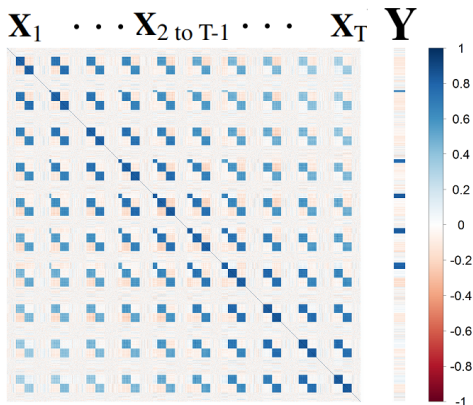
# Gestion des données manquantes



- ▶ Le dataset d'**apprentissage** est imputé par bouclage sur lui-même (sans résultat théorique pour l'instant sur la convergence).
- ▶ Le dataset de **test** est imputé sans bouclage car l'apprentissage se fait sur la partie d'**apprentissage**.

# Simulations

Construire un data-set de  $T$  blocs corrélé entre eux avec certaines variables de certains blocs corrélés à une variable réponse, univariée pour la comparaison à d'autres méthodes.



# Résultats de simulations

Comparaisons à des approches en 2 temps :

- ▶ Imputation :
  - ▶ missMDA (PCA + a-priori de groupe [HUSSON et JOSSE, 2013 ])
  - ▶ softImpute (modèle de prédiction [HASTIE et al., 2015 ])
  - ▶ variable mean value
  - ▶ `nipals` fonction de `mixOmics` (+ sPLS classique).
- ▶ Prédiction sur données complétées : **mdd-sPLS** et **Lasso**

Protocole similaire à [CHE et al., 2018 ] mais pour de la régression.  
Challenge de réseaux récurrents.

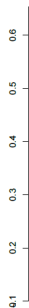
# Résultats de simulation

20 échantillons de 100 individus pour 10 blocs de 160 variables avec trois dimensions dont une seule est corrélée avec la réponse univariée.

2% of missing values



30% of missing values



- mdd-sPLS
- mixOmics sPLS
- imputeMFA + mdd-sPLS
- ▲ imputeMFA + Lasso
- softImpute + mdd-sPLS
- softImpute + Lasso
- ▲ Mean + mdd-sPLS
- Mean + Lasso



# Conclusion et futurs travaux

- ▶ Bons résultats de la méthode pour  $p_{NA} \geq 20\%$ ,
- ▶ Etude du comportement, prédiction, sélection, convergence, pour  $n$  plus faible,
- ▶ Application aux données Ebola rVSV.

**Merci !**

# References

-  Zhengping CHE et al. “Recurrent neural networks for multivariate time series with missing values”. In : *Scientific reports* 8.1 (2018), p. 6085.
-  Hyonho CHUN et Sündüz KELEŞ. “Sparse partial least squares regression for simultaneous dimension reduction and variable selection”. In : *Journal of the Royal Statistical Society : Series B (Statistical Methodology)* 72.1 (2010), p. 3–25.
-  Yash DESHPANDE et Andrea MONTANARI. “Sparse PCA via covariance thresholding”. In : *Advances in Neural Information Processing Systems*. 2014, p. 334–342.
-  Trevor HASTIE et al. “Matrix completion and low-rank svd via fast alternating least squares”. In : *J. Mach. Learn. Res* 16.1 (2015), p. 3367–3402.
-  Ana Maria HENAO-RESTREPO et al. “Efficacy and effectiveness of an rVSV-vectored vaccine in preventing Ebola virus disease : final results from the Guinea ring vaccination, open-label, cluster-randomised trial (Ebola Ça Suffit !)” In : *The Lancet* 389.10068 (2017), p. 505–518.
-  François HUSSON et Julie JOSSE. “Handling missing values in multiple factor analysis”. In : *Food quality and preference* 30.2 (2013), p. 77–85.
-  Kim-Anh LÊ CAO et al. “A sparse PLS for variable selection when integrating omics data”. In : *Statistical applications in genetics and molecular biology* 7.1 (2008).