

Méthodes régularisées pour l'analyse de données multivariées en grande dimension : théorie et applications.

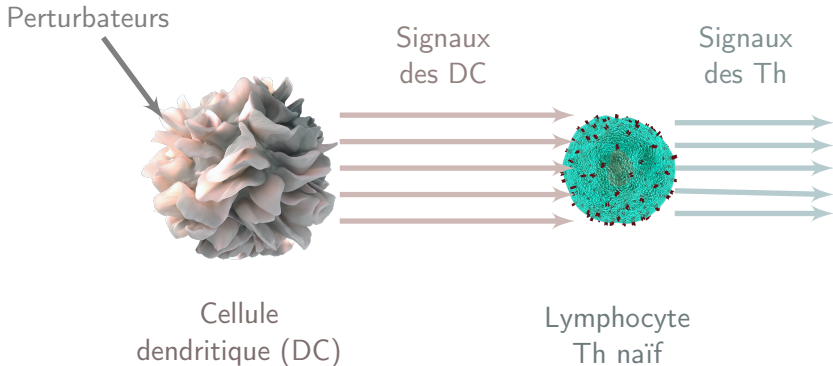
Marie Perrot-Dockès

Céline Lévy-Leduc, Julien Chiquet, Laure Sansonnet



Séminaire MIAT : 10 janvier 2020

Motivation : application en immunologie



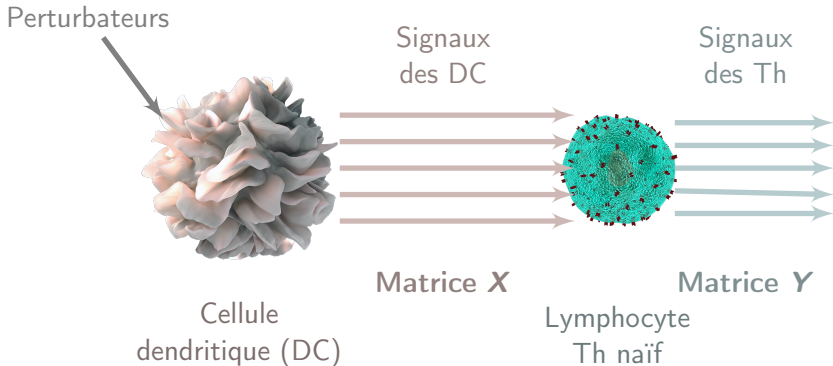
► Collaboration

Maximilien Grandclaudon, Coline Trichot, Vassili Soumelis

► Objectif

Étude du Dialogue entre DC et Th

Motivation : application en immunologie



► Collaboration

Maximilien Grandclaudon, Coline Trichot, Vassili Soumelis

► Objectif

Étude du Dialogue entre DC et Th

- ▶ **Description des données :**
 - ▶ \mathbf{X} : $n \times p$ matrice de design
contenant les signaux des cellules dendritiques
 - ▶ \mathbf{Y} : $n \times q$ matrice de réponses ($q \gg n$)
contenant les signaux des lymphocytes Th
- ▶ **Question :** Quelles variables influencent les réponses ?
- ▶ **Approche :** Sélection de variables dans le modèle linéaire général

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E},$$

où

- ▶ \mathbf{B} : $p \times q$ matrice **parcimonieuse** des coefficients
- ▶ \mathbf{E} : $n \times q$ matrice d'erreur avec

$$\forall i \in \llbracket 1, n \rrbracket, (E_{i,1}, \dots, E_{i,q}) \stackrel{iid}{\sim} \mathcal{N}(0, \Sigma)$$

en prenant en compte la dépendance en estimant Σ .

Traiter indépendamment les q modèles univariés :

$$\mathbf{Y}_{\bullet,r} = \mathbf{X}\mathbf{B}_{\bullet,r} + \mathbf{E}_{\bullet,r}, \quad \forall r \in \llbracket 1, q \rrbracket, \quad (1)$$

où $\mathbf{A}_{\bullet,r}$ désigne la r^{e} colonne de \mathbf{A} .

- ▶ **Maximum de vraisemblance** \Rightarrow pas parcimonieux !

Sélection de variables :

- ▶ AIC, BIC (Akaike, 1970, Schwarz et al., 1978)
- ▶ Tests (Mardia et al, 1980)
- ▶ **Régression pénalisée** Lasso (Tibshirani, 1996) :

$$\widehat{\mathbf{B}}_{\bullet,r}(\lambda) = \text{Argmin}_{\mathbf{B}_{\bullet,r}} \{ \|\mathbf{Y} - \mathbf{X}\mathbf{B}_{\bullet,r}\|_2^2 + \lambda \|\mathbf{B}_{\bullet,r}\|_1 \}.$$

Zhao et Yu (2006) ont montré sous certaines conditions :

$$\mathbb{P} \left(\text{sign}(\widehat{\mathbf{B}}_{\bullet,r}(\lambda)) = \text{sign}(\mathbf{B}_{\bullet,r}) \right) \rightarrow 1, \quad \text{lorsque } n \rightarrow \infty,$$

où $\text{sign}(\mathbf{x}) \in \{-1, 0, 1\}$

Cherchent à minimiser la fonction :

$$\ell(\mathbf{B}, \mathbf{\Omega}) = \text{tr} \left(\frac{1}{n} (\mathbf{Y} - \mathbf{XB}) \mathbf{\Omega} (\mathbf{Y} - \mathbf{XB})^\top \right) - \log(|\mathbf{\Omega}|), \quad (2)$$

où $\mathbf{\Omega} = \mathbf{\Sigma}^{-1}$

- ▶ **Maximum de vraisemblance** (Mardia *et al*, 1980)
- ▶ **Régressions pénalisées**
 - ▶ Rothman et al. (2010) : une méthode itérative une double pénalité , \mathbf{B} et $\mathbf{\Omega}$ parcimonieuses.
 - ▶ Lee & Liu (2012) : une étude théorique **à q fixé.**
 - ▶ Méthodes contemporaines : Zhang et al. (2017), Molstad et al. (2018).

Application du Lasso univarié

Dans le modèle $\mathcal{Y} = \mathcal{X}\mathcal{B} + \mathcal{E}$ l'estimateur Lasso est :

$$\hat{B}(\lambda) = \text{Argmin}_{\mathcal{B}} \{ \|\mathcal{Y} - \mathcal{X}\mathcal{B}\|_2^2 + \lambda \|\mathcal{B}\|_1 \}.$$

Vectorisation du modèle « blanchi »

$$\mathbf{Y}\Sigma^{-1/2} = \mathbf{X}\mathbf{B}\Sigma^{-1/2} + \mathbf{E}\Sigma^{-1/2}$$

$$\begin{aligned} \mathcal{Y} &= \text{vec}(\mathbf{Y}\Sigma^{-1/2}) = \text{vec}(\mathbf{X}\mathbf{B}\Sigma^{-1/2}) + \text{vec}(\mathbf{E}\Sigma^{-1/2}) \\ &= ((\Sigma^{-1/2})' \otimes \mathbf{X}) \text{vec}(\mathbf{B}) + \text{vec}(\mathbf{E}\Sigma^{-1/2}) \\ &= \mathcal{X}\mathcal{B} + \mathcal{E}. \end{aligned}$$



Nous ne connaissons pas Σ !

1 Estimation des erreurs : \hat{E}

Les résidus sont calculés indépendamment sur chaque colonne de \mathbf{Y}

2 Estimation de la matrice de covariance de \mathbf{E} : $\hat{\Sigma}$

- ▶ $n \gg q$: matrice de covariance empirique
- ▶ $q \gg n$: on suppose une structure particulière
 - ▶ Toeplitz symétrique,
 - ▶ par blocs.

3 « Blanchiment » : $\mathbf{Y} \hat{\Sigma}^{-1/2} = \mathbf{XB} \hat{\Sigma}^{-1/2} + \mathbf{E} \hat{\Sigma}^{-1/2}$

4 Sélection de variables en utilisant le critère Lasso et la « stability selection »

I. Estimation de matrice de covariance ($q \gg n$)

- Matrice de covariance Toeplitz
- Matrice de covariance par blocs

II. Garanties théoriques

III. Applications

- Eco-physiologie végétale
- Immunologie

IV. Conclusion et perspectives

Estimation de matrice de covariance en grande dimension ($q \gg n$)

- ▶ Matrice Toeplitz symétrique (une notion d'ordre dans les réponses)
 - ▶ la covariance ne dépend que de la distance entre deux réponses.



- ▶ Matrice par blocs (réponses groupées).



$$\forall i \in \llbracket 1, n \rrbracket, \forall t \in \mathbb{Z}, E_{i,t} - \phi_1 E_{i,t-1} = W_{i,t},$$

avec $(W_{i,t})_t \sim BB(0, 1), |\phi_1| < 1.$

$$\hat{\Sigma} = \frac{1}{1 - \hat{\phi}_1^2} \begin{pmatrix} 1 & \hat{\phi}_1 & \hat{\phi}_1^2 & \dots & \hat{\phi}_1^{q-1} \\ \hat{\phi}_1 & 1 & \hat{\phi}_1 & \dots & \hat{\phi}_1^{q-2} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \hat{\phi}_1^{q-1} & \dots & \dots & \dots & 1 \end{pmatrix},$$

Estimateur de ϕ_1 : $\hat{\phi}_1 = \frac{1}{n} \sum_{i=1}^n \hat{\phi}_1^{(i)},$

où $\hat{\phi}_1^{(i)}$ est l'estimateur de Yule-Walker de la ligne i de \mathbf{E}

En pratique on a besoin de $\hat{\Sigma}^{-1/2}$

$$\hat{\Sigma}^{-1/2} = \begin{pmatrix} \sqrt{1 - \hat{\phi}_1^2} & -\hat{\phi}_1 & 0 & \dots & 0 \\ 0 & 1 & -\hat{\phi}_1 & \dots & 0 \\ 0 & 0 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & -\hat{\phi}_1 \\ 0 & 0 & \dots & 0 & 1 \end{pmatrix},$$

Estimateur de ϕ_1 : $\hat{\phi}_1 = \frac{1}{n} \sum_{i=1}^n \hat{\phi}_1^{(i)},$

où $\hat{\phi}_1^{(i)}$ est l'estimateur de Yule-Walker de la ligne i de \mathbf{E}

Généralisation : estimateur de Σ dans le cas “Toeplitz”

$\forall i \in \llbracket 1, n \rrbracket$, on modélise $(E_{i,1}, \dots, E_{i,q})$ par un processus stationnaire.

$$\hat{\Sigma} = \begin{pmatrix} \hat{\gamma}(0) & \hat{\gamma}(1) & \cdots & \hat{\gamma}(q-1) \\ \hat{\gamma}(1) & \hat{\gamma}(0) & \cdots & \hat{\gamma}(q-2) \\ \vdots & & & \\ \hat{\gamma}(q-1) & \hat{\gamma}(q-2) & \cdots & \hat{\gamma}(0) \end{pmatrix}.$$

Estimateur de $\gamma(h)$: $\hat{\gamma}(h) = \frac{1}{n} \sum_{i=1}^n \hat{\gamma}_i(h),$

où $\hat{\gamma}_i(h)$ est un estimateur de la fonction d'autocovariance du processus $(E_{i,t})_t$ au retard h .

En pratique : on obtient $\hat{\Sigma}^{-1/2}$ à l'aide de l'inverse de Cholesky.

$$\tilde{\mathbf{E}} = \hat{\mathbf{E}} \hat{\Sigma}_q^{-1/2}$$

Sous (H0) « pour tout i dans $\{1, \dots, n\}$, $(\tilde{E}_{i,1}, \dots, \tilde{E}_{i,q})$ est un bruit blanc » on a :

$$q \sum_{i=1}^n \sum_{h=1}^H \hat{\rho}_i(h)^2 \approx \chi^2(nH), \text{ lorsque } q \rightarrow \infty,$$

où $\hat{\rho}_i(h)$ est un estimateur de la fonction d'autocovariance du processus $(\hat{\mathbf{E}}_{i,t})_t$ au retard h .

Estimation de matrice de covariance par blocs :



Estimation de Σ

Supposons qu'il existe

- ▶ \mathbf{Z} une matrice parcimonieuse de taille $q \times k$ avec $k \ll q$
- ▶ \mathbf{D} une matrice diagonale

Telles que

$$\Sigma = \mathbf{Z}\mathbf{Z}' + \mathbf{D},$$

avec les termes diagonaux de Σ égaux à 1.

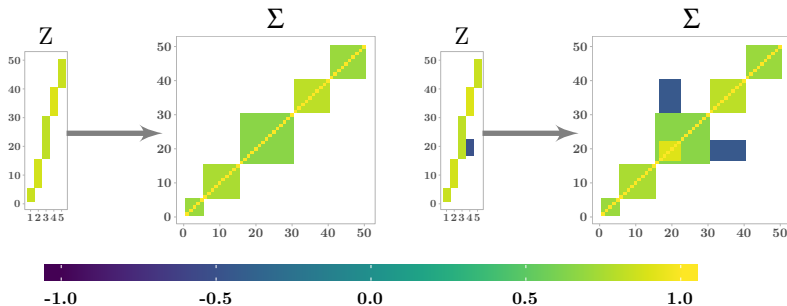


Figure – Exemples de matrices Σ générées à partir de matrices \mathbf{Z} . 16 / 44

- ▶ **Un estimateur de rang faible**

- ▶ Une matrice de rang faible contenant les termes extra-diagonaux de Σ
- ▶ Approximation de cette matrice en utilisant la décomposition en valeurs singulières.

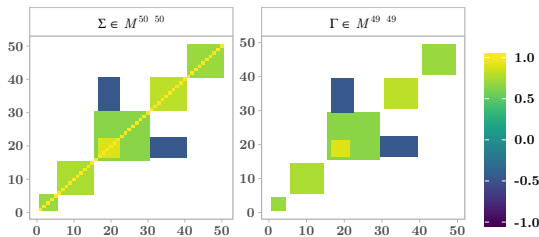
- ▶ **Un estimateur parcimonieux.**

Détecter les positions des valeurs non nulles de Σ .

- ▶ **Un estimateur défini positif.** Transformer $\tilde{\Sigma}$ en $\hat{\Sigma}$ une matrice définie positive (Higham, 2002).

Un estimateur de faible rang

► Passage de Σ à Γ



- $\text{rang}(\Sigma) = q$, $\text{rang}(\Gamma) = k \ll q$
- **En pratique** Σ est inconnu. $\tilde{\Gamma}$ est telle que

$$\tilde{\Gamma}_{i,j} = \begin{cases} R_{i,j+1} & \forall 1 \leq i \leq j \leq q-1 \\ \tilde{\Gamma}_{j,i} & \forall 1 \leq j < i \leq q-1 \end{cases},$$

où R est la matrice de corrélation empirique.

- $\tilde{\Gamma}^{(r)}$: une approximation de rang r de $\tilde{\Gamma}$ (SVD).
il faut choisir r !

Choix de r en pratique

- ▶ Critère de Cattell (Cattell, 1966)
- ▶ Méthode de permutation PA (Horn, 1965).

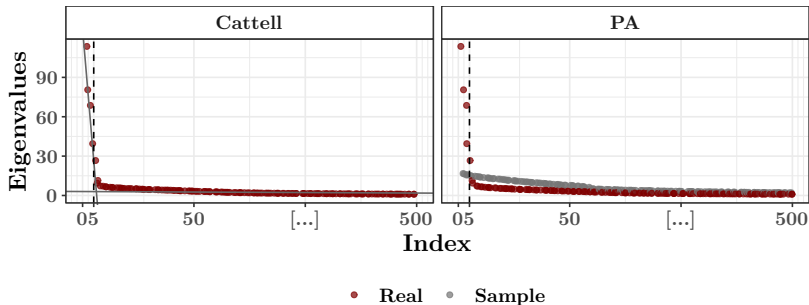


Figure – Choix de r en pratique $q = 500$ et $n = 30$, $k = 5$

Expériences numériques la méthode PA a tendance à sous-évaluer k lorsque n est faible.

1 Critère Lasso sur les valeurs de $\tilde{\Gamma}^{(r)}$

2 Ré-estimation des valeurs non nulles

⇒ Ceci revient à mettre un seuil sur les valeurs de $\tilde{\Gamma}^{(r)}$:

$$\hat{\Gamma}_{i,j}(\lambda) = \begin{cases} \tilde{\Gamma}_{j,i}^{(r)}, & \text{si } |\tilde{\Gamma}_{j,i}^{(r)}| > \frac{\lambda}{2} \\ 0, & \text{sinon} \end{cases}$$

En pratique : Comment choisir λ ?

- ▶ Le critère du coude calculé sur l'erreur $\|\hat{\Gamma}(\lambda) - \tilde{\Gamma}\|_F$
- ▶ Bickel & Levina 2008 : fondé sur la " cross-validation "

- Récupérer un estimateur de $\Sigma \Rightarrow$ on remet les 1 sur la diagonale.

$$\tilde{\Sigma}_{i,j} = \begin{cases} \hat{\Gamma}_{i,j-1}^{(r)} & \text{si } 1 \leq i < j \leq q \\ 1 & \text{si } 1 \leq i = j \leq q \\ \tilde{\Sigma}_{j,i} & \text{si } 1 \leq j < i \leq q \end{cases}$$

- Assurer sa positivité (Higham 2002) :

$$\hat{\Sigma} = \text{Argmin}_R \|\tilde{\Sigma} - R\|_F,$$

où R est une matrice de corrélation.

On veut utiliser notre méthode de sélection de variable !

⇒ Il nous faut un estimateur de $\Sigma^{-1/2}$

- $\hat{\Sigma}$ est symétrique donc il existe \mathbf{U} orthogonale et \mathbf{D} diagonale telles que

$$\hat{\Sigma} = \mathbf{U}\mathbf{D}\mathbf{U}'.$$

- En pratique on propose l'estimateur

$$\hat{\Sigma}^{-1/2} = \mathbf{U}\mathbf{D}_t^{-1/2}\mathbf{U}',$$

où

$$\mathbf{D}_t^{-1/2}{}_{i,i} = \begin{cases} \frac{1}{\sqrt{\mathbf{D}_{i,i}}} & \text{si } \mathbf{D}_{i,i} \geq t \\ 0 & \text{sinon.} \end{cases}$$

Comparaison avec des méthodes existantes : $\Sigma^{-1/2}$

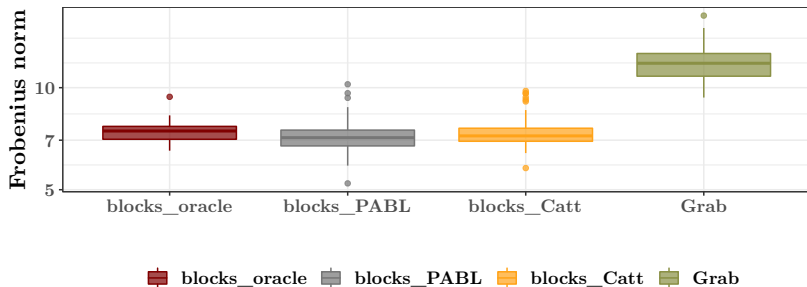


Figure – Comparaison de la norme de Frobenius $\|\hat{\Sigma}^{-1/2}\Sigma\hat{\Sigma}^{-1/2} - \text{Id}\|_F$ dans le cas **Extra-Diagonal-Equal** pour $n = 30$ et $q = 100$.

I. Estimation de matrice de covariance ($q \gg n$)

- Matrice de covariance Toeplitz
- Matrice de covariance par blocs

II. Garanties théoriques

III. Applications

- Eco-physiologie végétale
- Immunologie

IV. Conclusion et perspectives

Rappel

- ▶ $\mathbf{Y} = \mathbf{XB} + \mathbf{E}, \forall i \in \llbracket 1, n \rrbracket, (E_{i,1}, \dots, E_{i,q}) \stackrel{iid}{\sim} \mathcal{N}(0, \Sigma)$
- ▶ objectif : un estimateur parcimonieux de \mathbf{B}

Vectorisation du modèle « blanchi »

$$\mathbf{Y}\hat{\Sigma}^{-1/2} = \mathbf{XB}\hat{\Sigma}^{-1/2} + \mathbf{E}\hat{\Sigma}^{-1/2}$$

$$\begin{aligned}\mathcal{Y} &= \text{vec}(\mathbf{Y}\hat{\Sigma}^{-1/2}) = \text{vec}(\mathbf{XB}\hat{\Sigma}^{-1/2}) + \text{vec}(\mathbf{E}\hat{\Sigma}^{-1/2}) \\ &= ((\hat{\Sigma}^{-1/2})' \otimes \mathbf{X}) \text{vec}(\mathbf{B}) + \text{vec}(\mathbf{E}\hat{\Sigma}^{-1/2}) \\ &= \mathcal{X}\mathcal{B} + \mathcal{E}.\end{aligned}$$

Application du Lasso univarié

$$\hat{\mathcal{B}}(\lambda) = \text{Argmin}_{\mathcal{B}} \{ \|\mathcal{Y} - \mathcal{X}\mathcal{B}\|_2^2 + \lambda \|\mathcal{B}\|_1 \}.$$

Théorème (Perrot-Dockès et al, 2018)

Supposons qu'il existe M_1, M_2, M_3 et M_4, M_5 telles que

- ▶ $\|(\mathbf{X}^\top \mathbf{X})/n\|_\infty \leq M_1$, $\lambda_{\min}((\mathbf{X}^\top \mathbf{X})/n) \geq M_2$
- ▶ $\lambda_{\max}(\Sigma^{-1}) \leq M_3$, $\lambda_{\min}(\Sigma^{-1}) \geq M_4$
- ▶ Conditions d'irreprésentabilité sur \mathcal{X} construit avec Σ .
- ▶ Il existe c_1, c_2 telles que $0 < c_1 + c_2 < \frac{1}{2}$ qui satisfont
 - ▶ $s = O_{\mathbb{P}}(q^{c_1})$ où s est le cardinal du support J de \mathcal{B} ,
 - ▶ $q^{c_2} \min_{j \in J} |\mathcal{B}_j| \geq M_3$.
- ▶ $\|\Sigma^{-1} - \hat{\Sigma}^{-1}\|_\infty = \mathcal{O}_{\mathbb{P}}((nq)^{-1/2})$, $\rho(\Sigma - \hat{\Sigma}) = \mathcal{O}_{\mathbb{P}}((nq)^{-1/2})$

Alors, pour tout λ tel que $\frac{\lambda}{\sqrt{n}} \rightarrow \infty$ et $\frac{\lambda}{n} = o(q^{-(c_1+c_2)})$, lorsque $n \rightarrow \infty$

où $q = q_n = o\left(n^{\frac{1}{2(c_1+c_2)}}\right) = o(n^k)$ si $c_1 + c_2 = \frac{1}{2k}$, on a

$$\mathbb{P}\left(\text{sign}(\hat{\mathcal{B}}(\lambda)) = \text{sign}(\mathcal{B})\right) \rightarrow 1, \text{ lorsque } n \rightarrow \infty.$$

Un cas simple où les conditions sont vérifiées

- ▶ \mathbf{X} telle que $\mathbf{X}^\top \mathbf{X} = \nu \mathbf{I}$
(ex : matrice d'ANOVA à 1 facteur équilibré)
- ▶ $\forall i \in \llbracket 1, n \rrbracket$ E_i processus $AR(1)$

$$\forall i \in \{1, \dots, n\}, \forall t \in \mathbb{Z}, \mathbf{E}_{i,t} - \phi_1 \mathbf{E}_{i,t-1} = W_{i,t},$$

avec $(W_{i,t})_t \sim BB(0, 1), |\phi_1| < 1$.

Dans ce cas :

- ▶ si on estime ϕ_1 comme

$$\hat{\phi}_1 = \frac{\sum_{i=1}^n \sum_{\ell=2}^q \hat{\mathbf{E}}_{i,\ell} \hat{\mathbf{E}}_{i,\ell-1}}{\sum_{i=1}^n \sum_{\ell=1}^{q-1} \hat{\mathbf{E}}_{i,\ell}^2},$$

- ▶ si $j \in J, j + p$ ou $j - p$ n'est pas dans J (pour (IC))
alors, les conditions du théorème sont vérifiées !

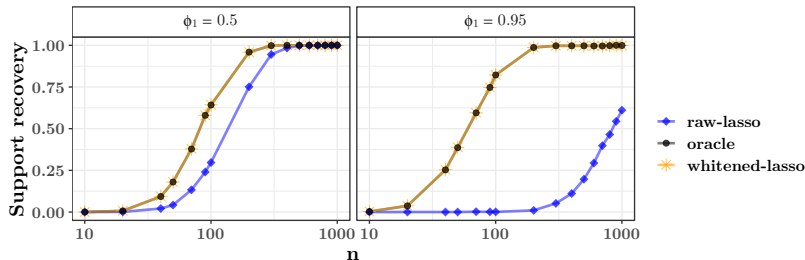
Expériences numériques : retrouver le support

- **Étude** : Fréquence des cas où

$$\exists \lambda, \text{sign}(\hat{B}(\lambda)) = \text{sign}(B)$$

- **Données** :

- $q = 1000$
- \mathbf{X} matrice d'ANOVA à 1 facteur à 2 modalités équilibré
- $\forall i \in \llbracket 1, n \rrbracket E_i$ processus $AR(1)$
- Dans le théorème on veut $q = q_n = o(n^k)$, ici $k = 2$



Expériences numériques : au-delà des hypothèses

- **Étude** : Fréquence des cas où

$$\exists \lambda, \text{sign}(\hat{B}(\lambda)) = \text{sign}(B)$$

- **Données** :

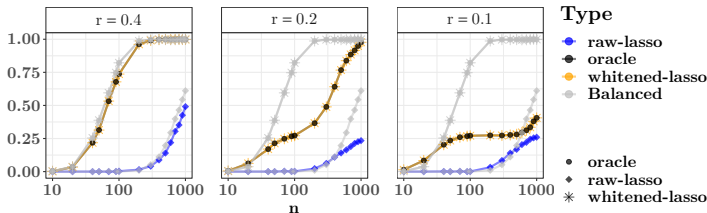
- $q = 1000$

- \mathbf{X} matrice d'ANOVA à 1 facteur à 2 modalités **déséquilibré**

$$r = \frac{\text{taille groupe 1}}{\text{taille totale}}$$

- $\forall i \in \llbracket 1, n \rrbracket$ E_i processus $AR(1)$

- Rappel $q = q_n = o\left(n^{\frac{1}{2(c_1+c_2)}}\right) = o(n^k)$, ici $k = 2$

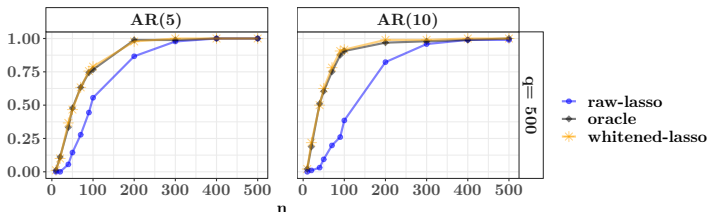


- **Étude** : Fréquence des cas où

$$\exists \lambda, \text{sign}(\hat{B}(\lambda)) = \text{sign}(B)$$

- **Données** :

- $q = 1000$
- \mathbf{X} matrice d'ANOVA à 1 facteur à 2 modalités équilibré
- $\forall i \in \llbracket 1, n \rrbracket$ E_i processus $AR(p)$
- Rappel $q = q_n = o\left(n^{\frac{1}{2(c_1+c_2)}}\right) = o(n^k)$, ici $k = 2$



- 1 Estimation des erreurs : \hat{E}
- 2 Estimation de la matrice de covariance de E : $\hat{\Sigma}$
- 3 « Blanchiment » : $Y \hat{\Sigma}^{-1/2} = XB \hat{\Sigma}^{-1/2} + E \hat{\Sigma}^{-1/2}$
- 4 Sélection de variables en utilisant le critère Lasso et la « stability selection »
 - ▶ Validation croisée pour sélectionner λ_{CV}
 - ▶ N tirage de taille $n/2$: soit F_i la fréquence où chaque variable i est sélectionnée
 - ▶ on garde les variable i telles que $F_i > \text{seuil}$

I. Estimation de matrice de covariance ($q \gg n$)

- Matrice de covariance Toeplitz
- Matrice de covariance par blocs

II. Garanties théoriques

III. Applications

- Eco-physiologie végétale
- Immunologie

IV. Conclusion et perspectives

Application en écophysiologie végétale

Étude de l'impact de la température de production sur la qualité des graines

Froid : 14-16 °C Standard : 18-22 °C Chaud : 25-28 °C



► Collaboration

Gwendal Cueff, Loic Rajjou

► Objectif

Recherche de biomarqueurs

► Données

► $X : 9 \times 3$

gammes de températures

► $Y : 9 \times 199$

accumulations des métabolites

► En pratique

► Covariance Toeplitz symétrique

► Seuil de « stability selection »
0.93

Estimation de B

Effet de la température sur la qualité des graines

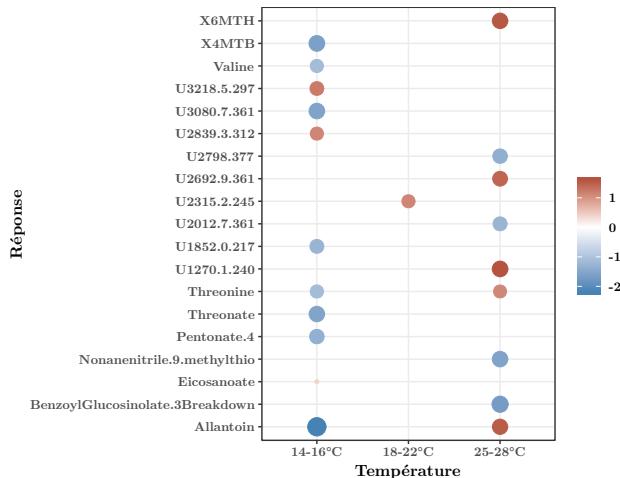
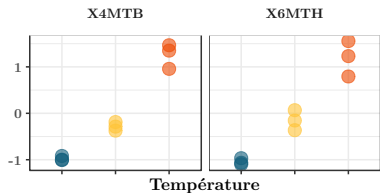


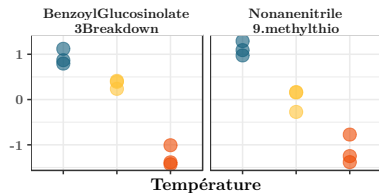
Figure – Estimation des coefficients $B_{i,j}$ pour les métabolites sélectionnés avec un seuil égal à 0.93.

La température de production sur les glucosinolates

Glucosinolates



Produits du catabolisme des glucosinolates

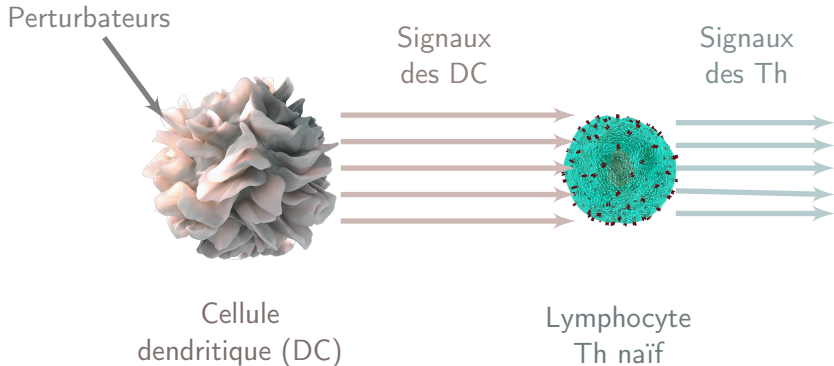


Température ● 14-16°C ● 18-22°C ● 25-28°C

La température modifie le métabolisme des glucosinolates qui

- ▶ luttent contre les ravageurs,
- ▶ sont antifongiques et antioxydants.

⇒ modification de la qualité biochimique et physiologique des graines.



► **Collaboration**

Maximilien Grandclaudon, Coline Trichot, Vassili Soumelis

► **Objectif**

Étude du Dialogue entre DC et Th

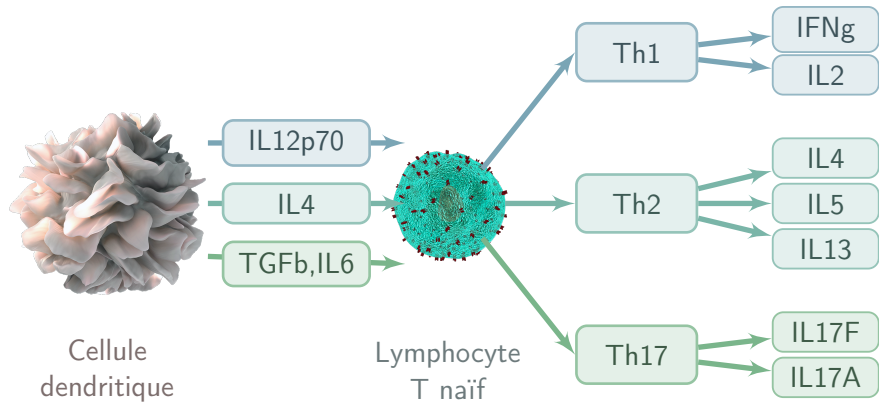


Figure – Les différents profils Th

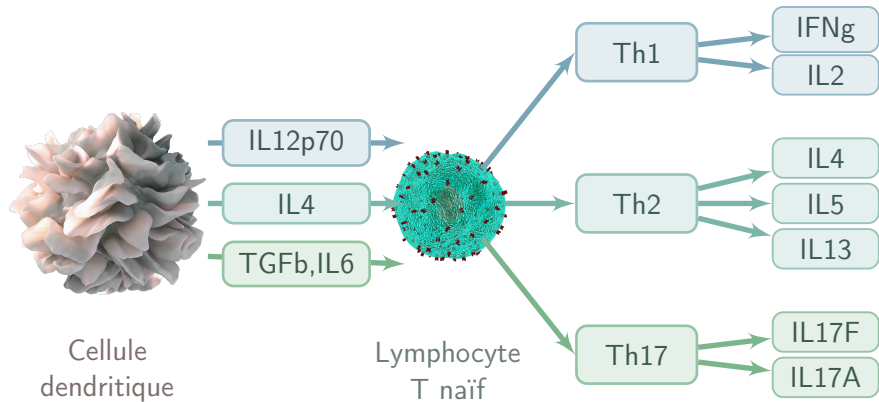
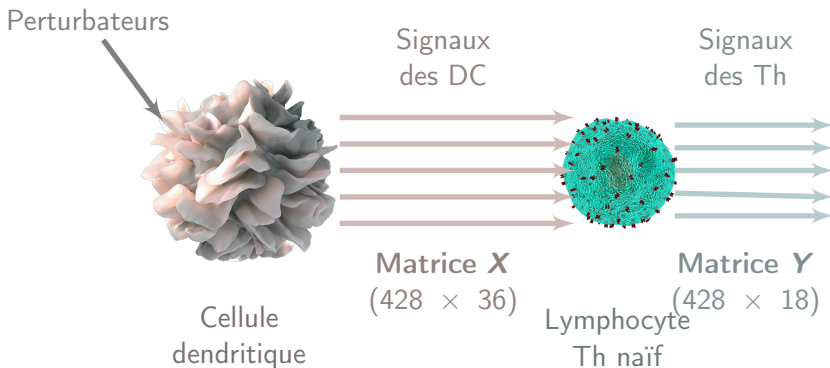


Figure – Les différents profils Th

Notre approche



► Données

► X : 428 × 36 signaux des DC

► Y : 428 × 18 signaux des Th

► En pratique

► Covariance empirique

► Seuil de « stability selection » 0.65

On retrouve les profils Th !

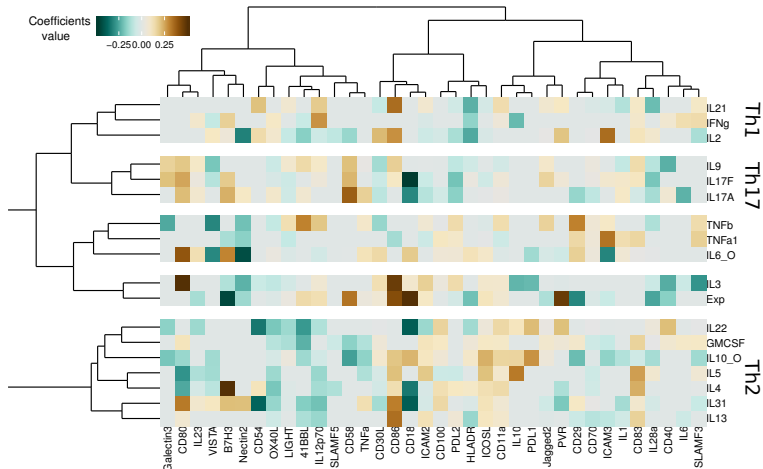
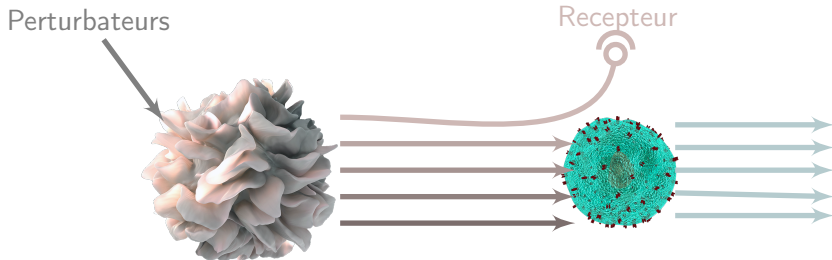


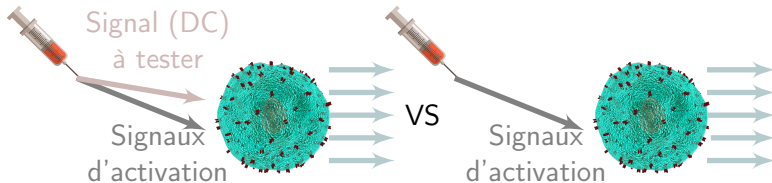
Figure – Coefficients de la modélisation des signaux des lymphocytes Th par les signaux des cellules dendritiques avec un seuil de 0.65.

Exemples d'expériences

► Un récepteur particulier



► Directement au lymphocytes Th



Les apports de cette présentation

- ▶ Un estimateur parcimonieux des coefficients :
 - ▶ Résultats théoriques : consistance en signe,
 - ▶ Simulations numériques.
- ▶ Des estimateurs de matrice de covariance
 - ▶ Dépendance de processus stationnaire
 - ▶ Résultats théoriques : vérification des hypothèses,
 - ▶ Études par simulations numériques,
 - ▶ Par blocs
 - ▶ Études par simulations numériques
- ▶ Applications :
 - ▶ à un problème d'écophysiologie végétale (métabolomique ciblée)
 - ▶ à un problème immunologique
validation de nombreuses associations importantes

Pour aller plus loin

- ▶ Vers d'autres cas vérifiant les conditions de consistance en signe de notre estimateur :
 - ▶ le cas des $ARMA(p, q)$ (Haddad, 2004),
 - ▶ les matrices de covariance par blocs diagonaux.
- ▶ Adaptation d'autres matrice de design
 - ▶ **En pratique** : R package VariSel qui permet de
 - ▶ regrouper des coefficients (group-lasso),
 - ▶ fusionner des coefficients (fused-lasso).
 - ▶ **En théorie** : adaptation au cas multivarié du
 - ▶ group-lasso (Bach, 2008),
 - ▶ fused-lasso (Rinaldo et al., 2009).
 - ▶ **Application**
 - ▶ Prendre en compte le type de cellule dendritique dans le dialogue avec les lymphocytes Th
- ▶ Développer des tests pour trouver la meilleure modélisation de la dépendance

Merci !

► Article publiés

- **Journal of Multivariate Analysis** M. Perrot-Dockès, C. Lévy-Leduc, L. Sansonnet, J. Chiquet, *"Variable selection in multivariate linear models with high-dimensional covariance matrix estimation"* , 166 :78 – 97, 2018.
- **Statistical Applications in Genetics and Molecular Biology** M. Perrot-Dockès, C. Lévy-Leduc, J. Chiquet, L. Sansonnet, M. Brégère, M.-P. Étienne, S. Robin, G. Genta-Jouve *"A variable selection approach in the multivariate linear model : An application to LC-MS metabolomics data"* 17(5), 2018.
- **Cell** M. Grandclaudon*, M. Perrot-Dockès*, C. Trichot,* O. Mostafa-Abouzid, W. Abou-Jaoudé, F. Berger, P. Hupé, D. Thieffry, L. Sansonnet, J. Chiquet, C. Lévy-Leduc, V. Soumelis *A quantitative multivariate model of human dendritic cell-T helper cell communication.* , 2019

* : ces auteurs ont contribué de manière égale à cette publication

► Article soumis

M. Perrot-Dockès, C. Lévy-Leduc *"Estimation of large block structured covariance matrices : Application to "multi-omic" approaches to study seed quality"*

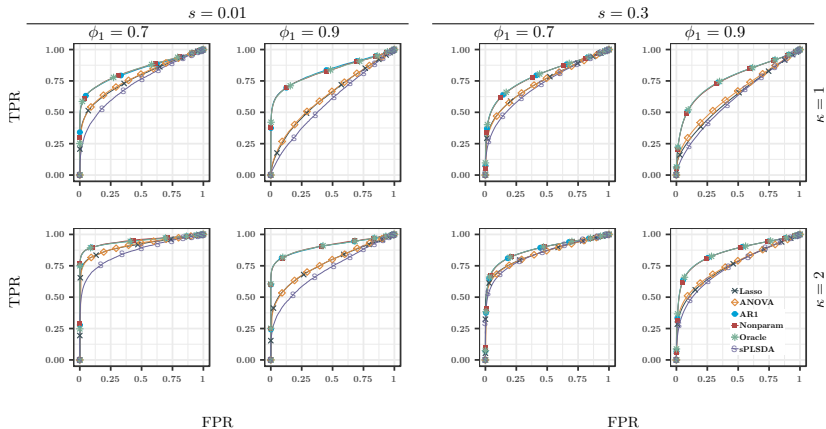
► Chapitre de livre (à paraître prochainement)

Wiley M. Perrot-Dockès, C. Lévy-Leduc *"Estimation of large block structured covariance matrices : Application to "multi-omic" approaches to study seed quality"*

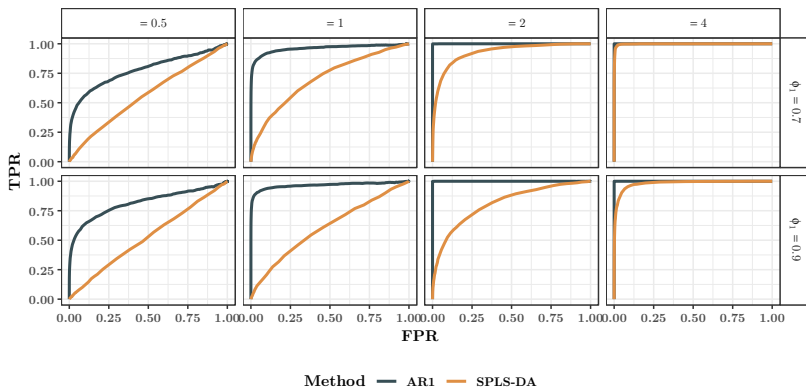
► Package R

- MultiVarSel disponible sur le CRAN.
- BlockCov disponible sur le CRAN.
- VariSel disponible sur github.

Comparaison courbe ROC



Comparaison dans le cas "classification"



$$|(\mathcal{X}^\top \mathcal{X})_{J^c, J} \{(\mathcal{X}^\top \mathcal{X})_{J, J}\}^{-1} \text{sign}(\mathcal{B}_J)| \leq \mathbf{1} - \eta,$$

où l'inégalité est vrai pour tous les éléments. Notons que :

$$\begin{aligned}(\mathcal{X}^\top \mathcal{X})_{J, J} &= \{((\Sigma^{-1/2})^\top \otimes \mathbf{X})^\top ((\Sigma^{-1/2})^\top \otimes \mathbf{X})\}_{J, J} \\&= (\Sigma^{-1/2} (\Sigma^{-1/2})^\top \otimes \mathbf{X}^\top \mathbf{X})_{J, J} \\&= (\Sigma^{-1} \otimes \mathbf{X}^\top \mathbf{X})_{J, J}.\end{aligned}$$

Donc : $S = \mathcal{X}^\top \mathcal{X} = \Sigma^{-1} \otimes \mathbf{X}^\top \mathbf{X}$.

$$\|S_{J^c, J} (S_{J, J})^{-1} \text{sign}(\mathcal{B}_J)\|_\infty \leq 1 - \eta,$$

Simulations numériques : choix de r

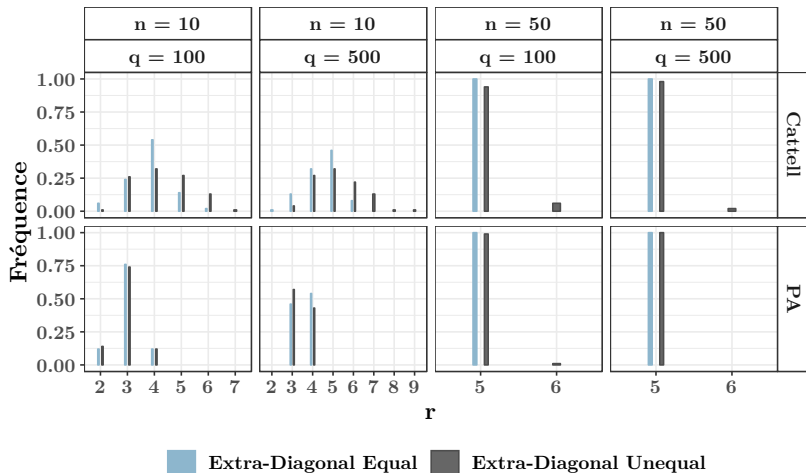


Figure – Choix de r en pratique ($k = 5$)

Sélection de variable : Comment choisir λ

- 1 Découper le jeux de données en 10 sous-groupes
Soit \mathcal{G}^v les données privées du v^e sous-groupe
- 2 Pour tout \mathcal{G}^v ,
 - ▶ **validation croisée** pour sélectionner λ_{CV} .
 - ▶ **Stability selection** au niveau λ_{CV} avec N réplifications
→ N_i^v le nombre de fois ou la variable i est sélectionnée dans \mathcal{G}^v
- 3 Garder les variables i telle que $F_i = \sum_{v=1}^{10} N_i^v / (10 \times N) > \text{seuil}$

Sélection de variable : choix du seuil en pratique

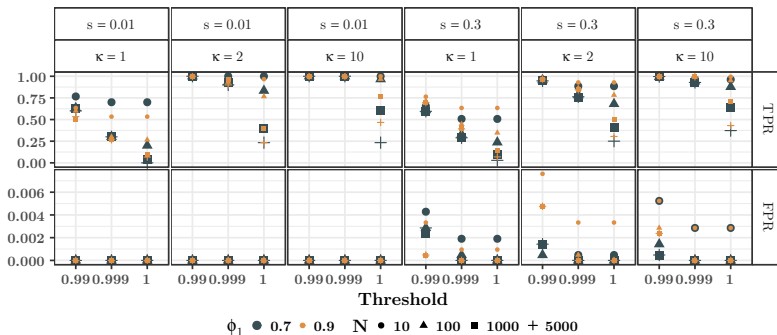


Figure – Influence du nombre de réplifications N et du seuil.

► Qu'est ce qu'un modèle à facteurs ?

$$\mathbf{E}_i = f_i \mathbf{Z}_f^\top + \mathbf{U}_i, \forall i \in \llbracket 1, n \rrbracket$$

où

- \mathbf{E}_i est la ligne i de \mathbf{E} ,
- \mathbf{Z}_f^\top est une matrice $k \times q$,
- f_i est un vecteur iid de taille k ,
- \mathbf{U}_i est un vecteur d'erreur de taille q (indep de f_i).

Sous cette hypothèse d'indépendance on a :

$$\Sigma = \mathbf{Z}_f^\top \text{Cov}(f) \mathbf{Z}_f^\top + \Sigma_u,$$

► Gérer les modèles à facteurs

- Blum et al. (2016) : un estimateur parcimonieux de \mathbf{B}_f donc de Σ lorsque $\forall i \in \llbracket 1, n \rrbracket, (f_{i,1}, \dots, f_{i,k}) \stackrel{iid}{\sim} \mathcal{N}(0, \mathbf{I})$.
- Hosseini & Lee (2016) : un estimateur parcimonieux de Ω à l'aide de modèle à facteur.