







Inférence de réseaux pour les données RNA-seq

Mélina Gallopin

Université Paris Descartes





Séminaire MIAT (INRA, Auzeville) Vendredi 26 mars 2016

1. Les données RNA-seq

2. Modélisation

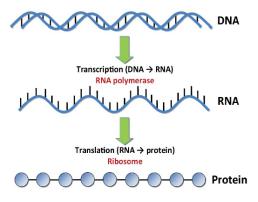
3. Réduction de dimension

1. Les données RNA-seq

2. Modélisation

3. Réduction de dimension

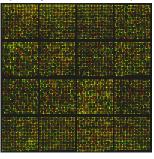
Mesurer l'expression des gènes



Mesurer l'expression des gènes

Comment?

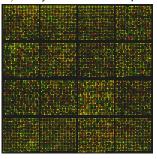
▶ les puces de ADN (1995) ⇒ hybridation des séquences d'ADNc



Mesurer l'expression des gènes

Comment?

▶ les puces de ADN (1995) ⇒ hybridation des séquences d'ADNc



▶ la technologie RNA-seq (2008) ⇒ lecture des séquences d'ADNc



Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

GATTACA, GTTTTTAGCTG, TAATTAG

Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

GATTACA, GTTTTTAGCTG, TAATTAG

4. Alignement des reads

Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

GATTACA, GTTTTTAGCTG, TAATTAG

4. Alignement des reads

Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

GATTACA, GTTTTTAGCTG, TAATTAG

4. Alignement des reads

Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

4. Alignement des reads

read aligné GATTACA génome de référence —TATTTAGCTCTGATTACAATG—

Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

GATTACA, GTTTTTAGCTG, TAATTAG

4. Alignement des reads

read aligné TTAGCTC
read aligné GATTACA
génome de référence —TATTTAGCTCTGATTACAATG—

Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

GATTACA, GTTTTTAGCTG, TAATTAG

4. Alignement des reads

read aligné read aligné read aligné GCTCTGAT TTAGCTC

GATTACA

Pour un échantillon :

- 1. Extraction de l'ARN
- 2. Retranscription ARN \Rightarrow ADNc
- 3. Lecture des brins d'ADNc, appelés reads

4. Alignement des reads

read aligné GCTCTGAT
read aligné TTAGCTC
read aligné GATTACA
génome de référence —TATTTAGCTCTGATTACAATG—

5. Comptage des reads

nombre de *reads* 45 17685 0 15 génome de référence – gène 1 — gène 2 – gène 3 – gène 4 —

		gène 1	gène 2	gène 3	gène 4	gène 5	
	porcelet 1	4	199	2987	0	65	
duodenum	porcelet 2	0	189	1806	0	29	
auoaenum	porcelet 3	6	201	1752	48	599	
	porcelet 4	4	198	2987	0	65	
	porcelet 1	0	0	1296	0	49	
1	porcelet 2	6	0	2298	0	119	
jejunum	porcelet 3	4	0	2987	0	651	
	porcelet 4	0	0	1876	0	219	
	porcelet 1	0	19931	1837	0	388	
ileum	porcelet 2	2	18319	8786	0	861	
	porcelet 3	7	23101	2237	0	76	
	porcelet 4	1	34198	9828	0	65	

		gène 1	gène 2	gène 3	gène 4	gène 5	
	porcelet 1	4	199	2987	0	65	
duodenum	porcelet 2	0	189	1806	0	29	
duodenum	porcelet 3	6	201	1752	48	599	
	porcelet 4	4	198	2987	0	65	
	porcelet 1	0	0	1296	0	49	
jejunum	porcelet 2	6	0	2298	0	119	
jejunum	porcelet 3	4	0	2987	0	651	
	porcelet 4	0	0	1876	0	219	
	porcelet 1	0	19931	1837	0	388	
ileum	porcelet 2	2	18319	8786	0	861	
	porcelet 3	7	23101	2237	0	76	
	porcelet 4	1	34198	9828	0	65	
	porceiet 4	1	34190	9020	0	- 03	

```
\mathbf{y}_i: expression pour l'échantillon i pour i=1,\ldots,n \mathbf{y}^j: expression du gène j pour j=1,\ldots,p \mathbf{y}_{ij}: expression du gène j pour l'échantillon i
```

		gène 1	gène 2	gène 3	gène 4	gène 5	
	porcelet 1	4	199	2987	0	65	
duodenum	porcelet 2	0	189	1806	0	29	
auoaenum	porcelet 3	6	201	1752	48	599	
	porcelet 4	4	198	2987	0	65	
	porcelet 1	0	0	1296	0	49	
jejunum	porcelet 2	6	0	2298	0	119	
jejunum	porcelet 3	4	0	2987	0	651	
	porcelet 4	0	0	1876	0	219	
	porcelet 1	0	19931	1837	0	388	
ileum	porcelet 2	2	18319	8786	0	861	
	porcelet 3	7	23101	2237	0	76	
	porcelet 4	1	34198	9828	0	65	

 \mathbf{y}_i : expression pour l'échantillon i pour $i=1,\ldots,n$ \mathbf{y}^j : expression du gène j pour $j=1,\ldots,p$ \mathbf{y}_{ij} : expression du gène j pour l'échantillon i \mathbf{y}_{ij} est un comptage \rightarrow Modélisations discrètes

		gène 1	gène 2	gène 3	gène 4	gène 5	
	porcelet 1	4	199	2987	0	65	
duodenum	porcelet 2	0	189	1806	0	29	
auoaenum	porcelet 3	6	201	1752	48	599	
	porcelet 4	4	198	2987	0	65	
	porcelet 1	0	0	1296	0	49	
jejunum	porcelet 2	6	0	2298	0	119	
jejunum	porcelet 3	4	0	2987	0	651	
	porcelet 4	0	0	1876	0	219	
	porcelet 1	0	19931	1837	0	388	
ileum	porcelet 2	2	18319	8786	0	861	
	porcelet 3	7	23101	2237	0	76	
	porcelet 4	1	34198	9828	0	65	

 \mathbf{y}_i : expression pour l'échantillon i pour $i=1,\ldots,n$ \mathbf{y}^j : expression du gène j pour $j=1,\ldots,p$ \mathbf{y}_{ij} : expression du gène j pour l'échantillon i \mathbf{y}_{ii} est un comptage \rightarrow Modélisations discrètes

p gènes > 1000n échantillons < 100

		gène 1	gène 2	gène 3	gène 4	gène 5	
	porcelet 1	4	199	2987	0	65	
duodenum	porcelet 2	0	189	1806	0	29	
auoaenum	porcelet 3	6	201	1752	48	599	
	porcelet 4	4	198	2987	0	65	
	porcelet 1	0	0	1296	0	49	
jejunum	porcelet 2	6	0	2298	0	119	
jejunum	porcelet 3	4	0	2987	0	651	
	porcelet 4	0	0	1876	0	219	
	porcelet 1	0	19931	1837	0	388	
ileum	porcelet 2	2	18319	8786	0	861	
	porcelet 3	7	23101	2237	0	76	
	porcelet 4	1	34198	9828	0	65	

 \mathbf{y}_i : expression pour l'échantillon i pour $i=1,\ldots,n$ \mathbf{y}^j : expression du gène j pour $j=1,\ldots,p$ \mathbf{y}_{ij} : expression du gène j pour l'échantillon i \mathbf{y}_{ij} est un comptage \rightarrow Modélisations discrètes

 $p~g\`{e}nes > 1000$ $n~\acute{e}chantillons < 100$ Technologie RNA-seq couteuse ightarrow Manque d'échantillons

		gène 1	gène 2	gène 3	gène 4	gène 5	
	porcelet 1	4	199	2987	0	65	
duodenum	porcelet 2	0	189	1806	0	29	
auoaenum	porcelet 3	6	201	1752	48	599	
	porcelet 4	4	198	2987	0	65	
	porcelet 1	0	0	1296	0	49	
jejunum	porcelet 2	6	0	2298	0	119	
jejunum	porcelet 3	4	0	2987	0	651	
	porcelet 4	0	0	1876	0	219	
	porcelet 1	0	19931	1837	0	388	
ileum	porcelet 2	2	18319	8786	0	861	
	porcelet 3	7	23101	2237	0	76	
	porcelet 4	1	34198	9828	0	65	

```
\mathbf{y}_i: expression pour l'échantillon i pour i=1,\ldots,n \mathbf{y}^j: expression du gène j pour j=1,\ldots,p \mathbf{y}_{ij}: expression du gène j pour l'échantillon i \mathbf{y}_{ij} est un comptage \rightarrow Modélisations discrètes
```

 $p~g\`{e}nes > 1000$ $n~\acute{e}chantillons < 100$ Technologie RNA-seq couteuse \rightarrow Manque d'échantillons

► Analyse différentielle

Déterminer si un gène est plus ou moins exprimé dans deux conditions

- Analyse différentielle
 Déterminer si un gène est plus ou moins exprimé dans deux conditions
- Analyse de co-expression
 Grouper les gènes ayant des profils d'expression similaires

- Analyse différentielle
 Déterminer si un gène est plus ou moins exprimé dans deux conditions
- Analyse de co-expression
 Grouper les gènes ayant des profils d'expression similaires
- ► Inférence de réseaux Reconstituer les réseaux de régulation génique

- Analyse différentielle
 Déterminer si un gène est plus ou moins exprimé dans deux conditions
- Analyse de co-expression
 Grouper les gènes ayant des profils d'expression similaires
- ► Inférence de réseaux Reconstituer les réseaux de régulation génique

Principe

Reconstruire un graphe G = (V, E) où :

- $lackbox{ }V=\{1,\ldots,p\}$ ensemble des nœuds représentant les variables
- ▶ E ensemble des arrêtes modélisant les dépendances entre les variables

Principe

Reconstruire un graphe G = (V, E) où :

- $lackbox{ }V=\{1,\ldots,p\}$ ensemble des nœuds représentant les variables
- E ensemble des arrêtes modélisant les dépendances entre les variables

Objectif

Reconstituer des réseaux de régulation de gènes



Principe

Reconstruire un graphe G = (V, E) où :

- $lackbrack V = \{1, \dots, p\}$ ensemble des nœuds représentant les variables
- E ensemble des arrêtes modélisant les dépendances entre les variables

Objectif

Reconstituer des réseaux de régulation de gènes



Outils

- Les modèles graphiques orientés
- Les modèles graphiques non-orientés

Principe

Reconstruire un graphe G = (V, E) où :

- $lackbox{ }V=\{1,\ldots,p\}$ ensemble des nœuds représentant les variables
- E ensemble des arrêtes modélisant les dépendances entre les variables

Objectif

Reconstituer des réseaux de régulation de gènes



Outils

- Les modèles graphiques orientés
- Les modèles graphiques non-orientés

 \mathbf{y}^j : expression du gène j pour $j=1,\ldots,p$

 \mathbf{y}^j : expression du gène j pour $j=1,\ldots,p$

Théorème (Hammersley-Clifford) La densité p des données se factorise ainsi

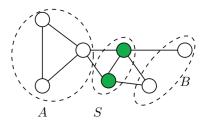
$$p(\mathbf{y}^1,\ldots,\mathbf{y}^p) = rac{1}{Z} \prod_{\mathcal{C} \in \mathfrak{C}} \psi_{\mathcal{C}}(\mathbf{y}^{\mathcal{C}}).$$

 \mathbf{y}^j : expression du gène j pour $j=1,\ldots,p$

Théorème (Hammersley-Clifford) La densité p des données se factorise ainsi

$$ho(\mathbf{y}^1,\ldots,\mathbf{y}^{
ho}) = rac{1}{Z} \prod_{\mathcal{C} \in \mathfrak{C}} \psi_{\mathcal{C}}(\mathbf{y}^{\mathcal{C}}).$$

 $ssi\left(\mathbf{y^{1}},\ldots,\mathbf{y^{p}}\right)$ vérifie la propriété de Markov par rapport au graphe G:

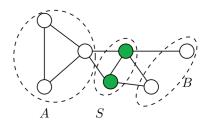


 \mathbf{y}^j : expression du gène j pour $j = 1, \dots, p$

Théorème (Hammersley-Clifford) La densité p des données se factorise ainsi

$$\rho(\mathbf{y}^1,\ldots,\mathbf{y}^\rho) = \frac{1}{Z} \prod_{\mathcal{C} \in \mathfrak{C}} \psi_{\mathcal{C}}(\mathbf{y}^{\mathcal{C}}).$$

 $ssi\left(\mathbf{y}^{1},\ldots,\mathbf{y}^{p}\right)$ vérifie la propriété de Markov par rapport au graphe G :



Exemple: Le modèle graphique gaussien (GGM) $\mathbf{y}_i \overset{\mathrm{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \Sigma)$

1. Les données RNA-seq

2. Modélisation

3. Réduction de dimension

1. Les données RNA-seq 2. Modélisation 3. Réduction de dimension

1. Les données RNA-sec

2. Modélisation

3. Reduction de dimension

Etat de l'art des méthodes développées

	Données de puces (1995) continues	Données RNA-seq (2008) discrètes
Analyse différentielle	gaussienne: limma	binomiale négative: DESeq; EdgeR
	(Smyth et al., 2005)	(Anders et al., 2010; Robinson et al., 2010),
Analyse de co-expression	gaussienne: Rmixmod	Poisson: HTSCluster
(modèle de mélange)	(Biernacki et al., 2006)	(Rau et al., 2015)
Inférence de réseaux	gaussienne: SIMoNe	
(modèle non-orienté)	(Chiquet et al., 2010)	?

Inférence de réseaux pour les données RNA-seq

Le modèle graphique gaussien (GGM)

$$\rho(\mathbf{y}^{1},...,\mathbf{y}^{p}) = \exp\left[-\frac{1}{2}\sum_{(j,j')\in E}\theta_{jj'}\mathbf{y}^{j}\mathbf{y}^{j'} - A(\Theta)\right]$$

$$où A(\Theta) = -\frac{1}{2}\log\det\left[\frac{\Theta}{2\pi}\right]$$

Inférence de réseaux pour les données RNA-seq

Le modèle graphique gaussien (GGM)

$$\rho(\mathbf{y}^{1},...,\mathbf{y}^{p}) = \exp\left[-\frac{1}{2}\sum_{(j,j')\in E}\theta_{jj'}\mathbf{y}^{j}\mathbf{y}^{j'} - A(\Theta)\right]$$

$$où A(\Theta) = -\frac{1}{2}\log\det\left[\frac{\Theta}{2\pi}\right]$$

Le modèle graphique de Poisson

$$p(\mathbf{y}^1, \dots, \mathbf{y}^p) = \exp \left[\sum_{j \in V} (\beta_j \mathbf{y}^j - \log \left(\mathbf{y}^j ! \right)) + \sum_{(j,j') \in E} \beta_{jj'} \mathbf{y}^j \mathbf{y}^{j'} - \mathbf{A}(\boldsymbol{\beta}) \right]$$

Calcul du terme de normalisation $A(\beta)$ difficile \to modélise uniquement des dépendances négatives

Inférence de réseaux pour les données RNA-seq

► Inférence par sélection de voisinage

$$p(\mathbf{y}^1, \dots, \mathbf{y}^p) = \exp \left[-\frac{1}{2} \sum_{(j,j') \in E} \theta_{jj'} \mathbf{y}^j \mathbf{y}^{j'} - A(\Theta) \right]$$

Meinshausen and Buhlmann (2006):

$$\mathbf{y}^j = \sum_{j' \neq j} \eta_{jj'} \mathbf{y}^{j'} + \epsilon_j, \,\, \epsilon_j \sim \mathcal{N}(0, \sigma_j), \,\, \eta_{jj'} = rac{ heta_{jj'}}{ heta_{jj}}.$$

- Sélection de voisinage pour les modèles discrets
 - Basée sur des lois de Poisson (Allen et al., 2013)
 - Basée sur des lois de Poisson sur-dispersées (Gallopin et al., 2013)

	Données de puces (1995) continues	Données RNA-seq (2008) discrètes
Analyse différentielle	gaussienne : limma	binomiale négative : DESeq; EdgeR
	(Smyth et al., 2005)	(Anders et al., 2010; Robinson et al., 2010),
Analyse de co-expression (modèle de mélange)	gaussienne : Rmixmod (Biernacki et al., 2006)	Poisson : HTSCluster (Rau et al., 2015)
Inférence de réseaux (modèle non-orienté)	gaussienne : SIMoNe (Chiquet et al., 2010)	Poisson; Poisson sur-dispersée (Allen et al., 2012; Gallopin et al., 2013)

	Données de puces (1995) continues	Données RNA-seq (2008) discrètes
Analyse différentielle	gaussienne : limma	binomiale négative : DESeq; EdgeR
	(Smyth et al., 2005)	(Anders et al., 2010; Robinson et al., 2010),
		ou
		gaussienne <i>sur données transformées</i> :
		limma + "voom" (Law et al., 2014)
Analyse de co-expression	gaussienne: Rmixmod	Poisson: HTSCluster
(modèle de mélange)	(Biernacki et al., 2006)	(Rau et al., 2015)
Inférence de réseaux	gaussienne : SIMoNe	Poisson; Poisson sur-dispersée
(modèle non-orienté)	(Chiquet et al., 2010)	(Allen et al., 2012; Gallopin et al., 2013)

	Données de puces (1995) continues	Données RNA-seq (2008) discrètes
Analyse différentielle	gaussienne : limma	binomiale négative : DESeq; EdgeR
	(Smyth et al., 2005)	(Anders et al., 2010; Robinson et al., 2010),
		ou
		gaussienne <i>sur données transformées</i> :
		limma + "voom" (Law et al., 2014)
Analyse de co-expression	gaussienne: Rmixmod	Poisson: HTSCluster
(modèle de mélange)	(Biernacki et al., 2006)	(Rau et al., 2015)
		ou
		gaussienne sur données transformées
Inférence de réseaux	gaussienne : SIMoNe	Poisson; Poisson sur-dispersée
(modèle non-orienté)	(Chiquet et al., 2010)	(Allen et al., 2012; Gallopin et al., 2013)
,		

	Données de puces (1995) continues	Données RNA-seq (2008) discrètes
Analyse différentielle	gaussienne : limma	binomiale négative : DESeq; EdgeR
Analyse unferentielle	_	
	(Smyth et al., 2005)	(Anders et al., 2010; Robinson et al., 2010),
		ou
		gaussienne <i>sur données transformées</i> :
		limma + "voom" (Law et al., 2014)
Analyse de co-expression	gaussienne : Rmixmod	Poisson: HTSCluster
(modèle de mélange)	(Biernacki et al., 2006)	(Rau et al., 2015)
		ou
		gaussienne sur données transformées
Inférence de réseaux	gaussienne : SIMoNe	Poisson; Poisson sur-dispersée
(modèle non-orienté)	(Chiquet et al., 2010)	(Allen et al., 2012; Gallopin et al., 2013)
,		ou
		gaussienne sur données transformées

	Données de puces (1995) continues	Données RNA-seq (2008) discrètes
		(
Analyse différentielle	gaussienne : limma	binomiale négative : DESeq; EdgeR
	(Smyth et al., 2005)	(Anders et al., 2010; Robinson et al., 2010),
		ou
		gaussienne <i>sur données transformées</i> :
		limma + "voom" (Law et al., 2014)
Analyse de co-expression	gaussienne : Rmixmod	Poisson: HTSCluster
(modèle de mélange)	(Biernacki et al., 2006)	(Rau et al., 2015)
		ou
		gaussienne sur données transformées
Inférence de réseaux	gaussienne : SIMoNe	Poisson; Poisson sur-dispersée
(modèle non-orienté)	(Chiquet et al., 2010)	(Allen et al., 2012; Gallopin et al., 2013)
		ou
		gaussienne sur données transformées

 \rightarrow modélisations gaussiennes sont des alternatives raisonnables aux modélisations discrètes

1. Les données RNA-seq 2. Modélisation 3. Réduction de dimension

1. Les données RNA-sec

Modélisation

3. Réduction de dimension

Modèle

$$\mathbf{y}_i \overset{\mathrm{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma}) \; \mathsf{pour} \; i = 1, \dots, n$$

Chaque arête du réseau \Leftrightarrow Coefficients non nuls de $\Theta=\Sigma^{-1}$

Modèle

$$\mathbf{y}_i \overset{\mathrm{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma}) \text{ pour } i = 1, \dots, n$$

Chaque arête du réseau \Leftrightarrow Coefficients non nuls de $\Theta = \Sigma^{-1}$

Maximisation en Θ par *Graphical lasso* (Friedman et al., 2008) de

$$\mathcal{L}_{\lambda}(\Theta;S) \propto \log \det(\Theta) - \text{tr}(S\Theta) - \lambda ||\Theta||_1$$

avec S matrice de covariance empirique

Modèle

$$\mathbf{y}_i \stackrel{\mathrm{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma}) \text{ pour } i = 1, \dots, n$$

Chaque arête du réseau \Leftrightarrow Coefficients non nuls de $\Theta = \Sigma^{-1}$

Maximisation en Θ par *Graphical lasso* (Friedman et al., 2008) de

$$\mathcal{L}_{\lambda}(\Theta; S) \propto \log \det(\Theta) - \operatorname{tr}(S\Theta) - \lambda ||\Theta||_1$$

avec ${\cal S}$ matrice de covariance empirique

Algorithme Block Diagonal Screening Rule (Mazumder et Hastie, 2012)

Pour un paramètre de régularisation fixé λ

Etape 1 Seuillage de |S| au niveau $\lambda \Rightarrow$ structure en blocs

Etape 2 Graphical lasso au paramètre λ dans chaque bloc

Modèle

$$\mathbf{y}_i \stackrel{\mathrm{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma}) \text{ pour } i = 1, \dots, n$$

Chaque arête du réseau \Leftrightarrow Coefficients non nuls de $\Theta = \Sigma^{-1}$

Maximisation en Θ par Graphical lasso (Friedman et al., 2008) de

$$\mathcal{L}_{\lambda}(\Theta; S) \propto \log \det(\Theta) - \operatorname{tr}(S\Theta) - \lambda ||\Theta||_1$$

avec ${\it S}$ matrice de covariance empirique

Algorithme Block Diagonal Screening Rule (Mazumder et Hastie, 2012)

Pour un paramètre de régularisation fixé λ

Etape 1 Seuillage de $\mid S \mid$ au niveau $\lambda \Rightarrow$ structure en blocs

Etape 2 Graphical lasso au paramètre λ dans chaque bloc



Phénomène d'ultra-grande dimension (Verzelen, 2012)

$$\frac{d\log(\frac{p}{d})}{n} \geq \frac{1}{2}$$

où d est le degré maximal du réseau

Exemple :
$$n = 50, p = 200, d \ge 8$$

Phénomène d'ultra-grande dimension (Verzelen, 2012)

$$\frac{d\log(\frac{p}{d})}{n} \geq \frac{1}{2}$$

où d est le degré maximal du réseau

Exemple :
$$n = 50, p = 200, d \ge 8$$

Solutions

1. Restreindre le nombre de gènes à l'aide d'informations externes

Phénomène d'ultra-grande dimension (Verzelen, 2012)

$$\frac{d\log(\frac{p}{d})}{n} \geq \frac{1}{2}$$

où d est le degré maximal du réseau

Exemple :
$$n = 50, p = 200, d \ge 8$$

Solutions

- 1. Restreindre le nombre de gènes à l'aide d'informations externes
- 2. Se focaliser uniquement sur les gènes les plus variables

Phénomène d'ultra-grande dimension (Verzelen, 2012)

$$\frac{d\log(\frac{p}{d})}{n} \geq \frac{1}{2}$$

où d est le degré maximal du réseau

Exemple :
$$n = 50, p = 200, d \ge 8$$

Solutions

- 1. Restreindre le nombre de gènes à l'aide d'informations externes
- 2. Se focaliser uniquement sur les gènes les plus variables
- 3. Sélectionner automatiquement les gènes clefs

Phénomène d'ultra-grande dimension (Verzelen, 2012)

$$\frac{d\log(\frac{p}{d})}{n} \geq \frac{1}{2}$$

où d est le degré maximal du réseau

Exemple :
$$n = 50, p = 200, d \ge 8$$

Solutions

- 1. Restreindre le nombre de gènes à l'aide d'informations externes
- 2. Se focaliser uniquement sur les gènes les plus variables
- 3. Sélectionner automatiquement les gènes clefs

Sélection automatique des gènes clefs

Cluster graphical lasso (Tan et al., 2015)

- 1. Détection de K groupes par classification hiérarchique des variables
- 2. Glasso dans chaque groupe avec différents paramètres de régularisation

Sélection automatique des gènes clefs

Cluster graphical lasso (Tan et al., 2015)

- Détection de K groupes par classification hiérarchique des variables
 ⇒ Choix de K par validation croisée
- 2. Glasso dans chaque groupe avec différents paramètres de régularisation

Sélection automatique des gènes clefs

Cluster graphical lasso (Tan et al., 2015)

- Détection de K groupes par classification hiérarchique des variables
 ⇒ Choix de K par validation croisée
- 2. Glasso dans chaque groupe avec différents paramètres de régularisation

Notre procédure (Devijver & Gallopin, soumis)

- 1. Choix automatique des groupes par un critère de sélection de modèles
- 2. Glasso dans chaque groupe avec différents paramètres de régularisation

Détection des groupes de gènes

Hypothèse:
$$\mathbf{y}_i \stackrel{\mathrm{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \Sigma_B)$$
 avec $\Sigma_B = \begin{pmatrix} \Sigma^1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Sigma^K \end{pmatrix}$

La structure en K blocs de Σ_B définit les groupes de gènes $B=(B_1,\ldots,B_K)$

Détection des groupes de gènes

Hypothèse:
$$\mathbf{y}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \Sigma_B)$$
 avec $\Sigma_B = \begin{pmatrix} \Sigma^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Sigma^K \end{pmatrix}$

La structure en K blocs de Σ_B définit les groupes de gènes $B=(B_1,\ldots,B_K)$

$$F_B = \{f_B = \mathcal{N}_p(0, \Sigma_B) \text{ with } \Sigma_B \in S_B\}$$

$$S_B = \left\{ \Sigma_B \in \mathbb{S}_p^{++}(\mathbb{R}) \middle| \Sigma_B = P_\sigma \begin{pmatrix} \Sigma^1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Sigma^K \end{pmatrix} P_\sigma^{-1}, \Sigma^k \in \mathbb{S}_{p_k}^{++}(\mathbb{R}), \forall k \in 1, \dots, p \right\}$$

Détection des groupes de gènes

Hypothèse:
$$\mathbf{y}_i \stackrel{\mathrm{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \Sigma_B)$$
 avec $\Sigma_B = \begin{pmatrix} \Sigma^1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Sigma^K \end{pmatrix}$

La structure en K blocs de Σ_B définit les groupes de gènes $B=(B_1,\ldots,B_K)$

$$F_{\mathcal{B}} = \left\{ f_{\mathcal{B}} = \mathcal{N}_{\mathcal{P}}(0, \Sigma_{\mathcal{B}}) \text{ with } \Sigma_{\mathcal{B}} \in \mathcal{S}_{\mathcal{B}} \right\}$$

$$S_{\mathcal{B}} = \left\{ \Sigma_{\mathcal{B}} \in \mathbb{S}_{\mathcal{P}}^{++}(\mathbb{R}) \middle| \Sigma_{\mathcal{B}} = P_{\sigma} \begin{pmatrix} \Sigma^{1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Sigma^{K} \end{pmatrix} P_{\sigma}^{-1}, \Sigma^{k} \in \mathbb{S}_{\mathcal{P}_{k}}^{++}(\mathbb{R}), \forall k \in 1, \dots, \mathcal{P} \right\}$$

Sélection par l'heuristique de pente (Birgé and Massart, 2007)

$$\hat{B} = \underset{B}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{i=1}^{n} \log(\hat{f}_{B}(\mathbf{y}_{i})) + \operatorname{pen}(B) \right\},$$

$$\operatorname{pen}(B) = \kappa D_{B}.$$

 \mathcal{B} : ensemble de toutes les partitions des variables possibles \Rightarrow Exploration exhaustive de \mathcal{B} impossible

 \mathcal{B} : ensemble de toutes les partitions des variables possibles \Rightarrow Exploration exhaustive de \mathcal{B} impossible

 \mathcal{B}^{\wedge} : ensemble des partitions obtenues par seuillage de $\mid S \mid$

 \mathcal{B} : ensemble de toutes les partitions des variables possibles \Rightarrow Exploration exhaustive de \mathcal{B} impossible

 \mathcal{B}^{\wedge} : ensemble des partitions obtenues par seuillage de $\mid S \mid$

$$\hat{B} = \underset{B \in \mathcal{B}^{\Lambda}}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{i=1}^{n} \log(\hat{f}_{B}(\mathbf{y}_{i})) + \operatorname{pen}(B) \right\},$$

$$\operatorname{pen}(B) = \kappa D_{B}.$$

 \mathcal{B} : ensemble de toutes les partitions des variables possibles \Rightarrow Exploration exhaustive de \mathcal{B} impossible

 \mathcal{B}^{\wedge} : ensemble des partitions obtenues par seuillage de $\mid S \mid$

$$\hat{B} = \underset{B \in \mathcal{B}^{\Lambda}}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{i=1}^{n} \log(\hat{f}_{B}(\mathbf{y}_{i})) + \operatorname{pen}(B) \right\},$$

$$\operatorname{pen}(B) = \kappa D_{B}.$$

ightarrow Procédure garantie par des résultats non-asymptotiques

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\operatorname{pen}(B) \geq \kappa' \frac{D_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{D_B(\frac{D_B}{n}c^2 \wedge 1)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C\mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{B}^{\textbf{borné}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\operatorname{pen}(B) \geq \kappa' \frac{\mathsf{D}_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{\mathsf{D}_B\left(\frac{\mathsf{D}_B}{n}c^2 \wedge 1\right)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C \mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{B}^{\textbf{borné}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\mathrm{pen}(B) \geq \kappa' \frac{\mathsf{D}_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{\mathsf{D}_B\left(\frac{\mathsf{D}_B}{n}c^2 \wedge 1\right)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C \mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{\hat{B}}^{\textbf{born\'e}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Borne minimax: Pour tout $B \in \mathcal{B}$, il existe une constante $C_1 > 0$ telle que:

$$\inf_{\hat{f}_B} \sup_{f \in F_B^{\text{borné}}} \mathbb{E} \big(d_H^2 \big(\hat{f}_B, f \big) \big) \geq C_1 \frac{D_B}{n} \big(1 + \log \left(\frac{C_2}{D_B^2} \right) \big).$$

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\mathrm{pen}(B) \geq \kappa' \frac{\mathsf{D}_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{\mathsf{D}_B\left(\frac{\mathsf{D}_B}{n}c^2 \wedge 1\right)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C \mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{\hat{B}}^{\textbf{born\'e}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Borne minimax: Pour tout $B \in \mathcal{B}$, il existe une constante $C_1 > 0$ telle que:

$$\inf_{\hat{f}_B} \sup_{f \in F_B^{\text{borné}}} \mathbb{E}(d_H^2(\hat{f}_B, f)) \ge C_1 \frac{D_B}{n} (1 + \log\left(\frac{C_2}{D_B^2}\right)).$$

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\mathrm{pen}(B) \geq \kappa' \frac{\mathsf{D}_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{\mathsf{D}_B\left(\frac{\mathsf{D}_B}{n}c^2 \wedge 1\right)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C \mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{B}^{\textbf{born\acute{e}}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Borne minimax: Pour tout $B \in \mathcal{B}$, il existe une constante $C_1 > 0$ telle que:

$$\inf_{\hat{f}_B} \sup_{f \in F_B^{\text{borné}}} \mathbb{E} \big(d_H^2 \big(\hat{f}_B, f \big) \big) \geq C_1 \frac{\mathsf{D}_B}{n} \big(1 + \log \left(\frac{C_2}{\mathsf{D}_B^2} \right) \big).$$

→ Procédure de sélection de modèles optimale

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\mathrm{pen}(B) \geq \kappa' \frac{\mathsf{D}_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{\mathsf{D}_B\left(\frac{\mathsf{D}_B}{n}c^2 \wedge 1\right)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C\mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{B}^{\textbf{born\'e}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Borne minimax: Pour tout $B \in \mathcal{B}$, il existe une constante $C_1 > 0$ telle que:

$$\inf_{\hat{f}_B} \sup_{f \in \mathcal{F}_B^{\text{borné}}} \mathbb{E}(d_H^2(\hat{f}_B, f)) \geq C_1 \frac{D_B}{n} (1 + \log\left(\frac{C_2}{D_B^2}\right)).$$

→ Procédure de sélection de modèles optimale

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\mathrm{pen}(B) \geq \kappa' \frac{\mathsf{D}_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{\mathsf{D}_B\left(\frac{\mathsf{D}_B}{n}c^2 \wedge 1\right)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C\mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{\hat{B}}^{\textbf{borné}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Borne minimax: Pour tout $B \in \mathcal{B}$, il existe une constante $C_1 > 0$ telle que:

$$\inf_{\hat{f}_B} \sup_{f \in \mathcal{F}_B^{\text{borné}}} \mathbb{E}(d_H^2(\hat{f}_B, f)) \ge C_1 \frac{D_B}{n} (1 + \log\left(\frac{C_2}{D_B^2}\right)).$$

- → Procédure de sélection de modèles optimale
- → Calibration de la pénalité à partir des données

Inégalité Oracle : On suppose qu'il existe une constante $\kappa'>0$ telle que pour chaque partition $B\in\mathcal{B}$,

$$\mathrm{pen}(B) \geq \kappa' \frac{\mathsf{D}_B}{n} \left[2c^2 + \rho \log \left(\frac{1}{\mathsf{D}_B\left(\frac{\mathsf{D}_B}{n}c^2 \wedge 1\right)} \right) + (1 \vee \tau) \log \left(\frac{0.792p}{\log(p+1)} \right) \right],$$

où c est une constante. Alors $\hat{f}_{\hat{B}}$ satisfait, pour une constante C :

$$\mathbb{E}(d_H^2(f^*,\hat{f}_{\hat{B}})) \leq C \mathbb{E}\left(\inf_{B \in \mathcal{B}^{\Lambda}} \left(\inf_{t \in \mathcal{F}_{\hat{B}}^{\textbf{born\'e}}} \mathsf{KL}(f^*,t) + \mathrm{pen}(B)\right) + (1 \vee \tau)\frac{1}{n}\right).$$

Borne minimax: Pour tout $B \in \mathcal{B}$, il existe une constante $C_1 > 0$ telle que:

$$\inf_{\hat{f}_B} \sup_{f \in \mathcal{F}_B^{\text{borné}}} \mathbb{E}(d_H^2(\hat{f}_B, f)) \geq C_1 \frac{D_B}{n} (1 + \log\left(\frac{C_2}{D_B^2}\right)).$$

- → Procédure de sélection de modèles optimale
- → Calibration de la pénalité à partir des données
- \rightarrow Pénalité de la forme pen $(B) = \kappa D_B$
 - à l'aide du package R capushe (Baudry et al.),

Calibration du coefficient κ dans $pen(B) = \kappa D_B$

ightharpoonup Illustrations sur données simulées : p=100, n=70 et $K^\star=15$

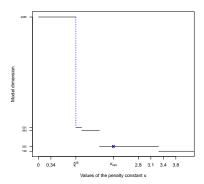
Calibration du coefficient κ dans $pen(B) = \kappa D_B$

- ▶ Illustrations sur données simulées : p=100, n=70 et $K^{\star}=15$
- Méthodes de calibration utilisées en pratique implémentées dans le package R capushe (Baudry et al., 2012)

Calibration du coefficient κ dans $pen(B) = \kappa D_B$

- ▶ Illustrations sur données simulées : p = 100, n = 70 et $K^* = 15$
- Méthodes de calibration utilisées en pratique implémentées dans le package R capushe (Baudry et al., 2012)

Méthode 1 : SHDJ Slope Heuristics Dimension Jump



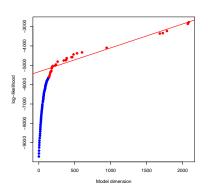
Calibration du coefficient κ dans $pen(B) = \kappa D_B$

- ▶ Illustrations sur données simulées : p = 100, n = 70 et $K^* = 15$
- Méthodes de calibration utilisées en pratique implémentées dans le package R capushe (Baudry et al., 2012)

Méthode 1 : SHDJ Slope Heuristics Dimension Jump

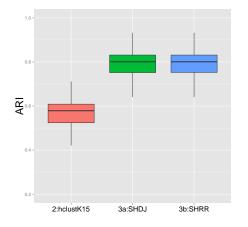
Values of the penalty constant is

Méthode 2 : SHRR Slope Heuristics Robust Regression



Données simulées

p=100, n=70 et Σ diagonale par blocs et $K^*=15$.



Adjusted Rand Index entre la vraie partition et la partition détectée

Données simulées : p = 100, n = 70 et Σ diagonale par blocs avec $K^* = 15$.

▶ (1) Glasso graphical lasso sélectionné par BIC sur l'ensemble des variables

Données simulées : p = 100, n = 70 et Σ diagonale par blocs avec $K^* = 15$.

- (1) Glasso graphical lasso sélectionné par BIC sur l'ensemble des variables
- ▶ (2) Cluster Graphical Lasso (Tan et al., 2015)

```
Etape 1 : classification hiérarchique des variables, pour K = K^* fixé Etape 2 : \rho_1, \ldots, \rho_{K^*} d'après le corollaire de Tan 2015.
```

Données simulées : p = 100, n = 70 et Σ diagonale par blocs avec $K^* = 15$.

- (1) Glasso graphical lasso sélectionné par BIC sur l'ensemble des variables
- ▶ (2) Cluster Graphical Lasso (Tan et al., 2015)

Etape 1 : classification hiérarchique des variables, pour $K = K^*$ fixé Etape 2 : $\rho_1, \ldots, \rho_{K^*}$ d'après le corollaire de Tan 2015.

(3) Notre procédure

Etape 1: (3a) SHRR partition (3b) SHDJ partition

Etape 2 : graphical lassos sélectionnés par BIC

Données simulées : p = 100, n = 70 et Σ diagonale par blocs avec $K^* = 15$.

- (1) Glasso graphical lasso sélectionné par BIC sur l'ensemble des variables
- ▶ (2) Cluster Graphical Lasso (Tan et al., 2015)

```
Etape 1 : classification hiérarchique des variables, pour K = K^* fixé Etape 2 : \rho_1, \ldots, \rho_{K^*} d'après le corollaire de Tan 2015.
```

► (3) Notre procédure

```
Etape 1: (3a) SHRR partition (3b) SHDJ partition
```

Etape 2 : graphical lassos sélectionnés par BIC

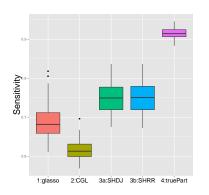
 (4) Partition des variables connues graphical lassos sélectionnés par BIC dans chaque bloc

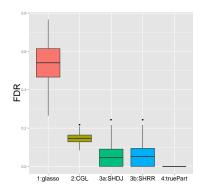
Performance des stratégies sur données simulées

p=100, n=70 et Σ diagonale par blocs avec $K^*=15$.

Sensibilité =
$$\frac{TP}{(TP + FN)}$$

$$FDR = \frac{FP}{(TP + FP)}$$

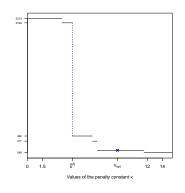


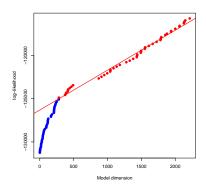


Sur 100 jeux de données simulés

Données réelles

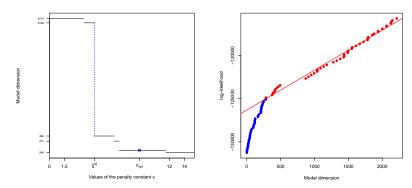
- ▶ Pickrell et al. (2010) : étude des lymphoblastes chez n = 69 individus
- Présélection des p = 200 gènes les plus variables





Données réelles

- ▶ Pickrell et al. (2010) : étude des lymphoblastes chez n = 69 individus
- ightharpoonup Présélection des p=200 gènes les plus variables



→ Les partitions SHRR et SHDJ coïncident

Graphical lasso

ightharpoonup D = 19900 paramètres à estimer

Graphical lasso

▶ D = 19900 paramètres à estimer

Partition détectée par l'heuristique des pentes \hat{B}

$$D_{\hat{B}_{SH}} = 283$$

Graphical lasso

▶ D = 19900 paramètres à estimer

Partition détectée par l'heuristique des pentes \hat{B}

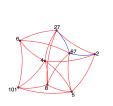
- $D_{\hat{B}_{SH}} = 283$
- $ightharpoonup \hat{K}_{SH} = 150 \text{ blocs}$

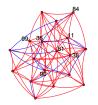
Graphical lasso

ightharpoonup D = 19900 paramètres à estimer

Partition détectée par l'heuristique des pentes \hat{B}

- ▶ $D_{\hat{B}_{SH}} = 283$
- $\hat{K}_{SH} = 150 \text{ blocs}$
- ▶ 140 blocs de taille 1, 2 blocs de taille 2, 4 blocs de taille 3 et 4 blocs de taille 18, 13, 8 et 5









Conclusion et perspectives

Conclusion

- Une procédure non-asymptotique pour réduire la dimension du problème d'inférence (Devijver et Gallopin, soumis)
- implémentée dans le package shock disponible sur le CRAN, développé sur Github https://github.com/Gallopin/shock

Conclusion et perspectives

Conclusion

- Une procédure non-asymptotique pour réduire la dimension du problème d'inférence (Devijver et Gallopin, soumis)
- implémentée dans le package shock disponible sur le CRAN, développé sur Github https://github.com/Gallopin/shock

Perspectives

- Applications sur des données réelles du département de génétique animale de l'INRA
- Inclusion d'informations externes pour améliorer la sélection des gènes : exemple, les données Hi-C

1. Les données RNA-seq 2. Modélisation 3. Réduction de dimension

Merci pour votre attention

Bibliographie

- J. Whittaker, Graphical Models in Applied Multivariate Statistics, Wiley Publishing, 2009
- ▶ J. Friedman, T. Hastie, R. Tibshirani, Sparse inverse covariance estimation with the Lasso, Biostatistics, 2008
- R. Mazumder, T. Hastie, Exact Covariance Thresholding into Connected Components for Large-Scale Graphical Lasso, Journal of Machine Learning Reasearch, 2012
- D. Witten, J. Friedman, N. Simon, New Insights and Faster Computations for the Graphical Lasso, Journal of Computational and Graphical Statistics, 2011
- K. Tan, D. Witten, A. Shojaie, The Cluster Graphical Lasso for improved estimation of Gaussian graphical models, Computation Statistics & Data Analysis, 2015
- ▶ L. Birgé, P. Massart, Minimal penalties for Gaussian model selection, Probab. Theory Related Fields, 2007
- P. Massart, Concentration inequalities and model selection, Lecture Notes in Mathematics. Springer, 33, 2003, Saint-Flour, Cantal.