

Avances dans les modèles mixtes pour la prévision et la sélection des variables dans les données de grande dimension

Sahir Rai Bhatnagar
Department of Epidemiology, Biostatistics, and Occupational Health
Department of Diagnostic Radiology
McGill University

sahirbhatnagar.com

19 novembre 2021



Analyse de données haute dimension (HD)

Classique

$$\mathbf{X}_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{12} & \cdots & x_{1p} \\ x_{31} & x_{12} & \cdots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{12} & \cdots & x_{np} \end{bmatrix}$$

Analyse de données haute dimension (HD)

Classique

$$\mathbf{X}_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{12} & \cdots & x_{1p} \\ x_{31} & x_{12} & \cdots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{12} & \cdots & x_{np} \end{bmatrix}$$

nnées HD

$$\mathbf{X}_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{12} & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & x_{np} \end{bmatrix}$$

De nouveaux défis découlent de la façon dont ces données sont *utilisées*

A		B								
y	x ₁	y	x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	x ₈
0.0	0	0	0	2	0	0	1	0	1	0
2.1	1	2.1	1	0	2	3	2	0	0	3
2.7	0	2.7	0	0	0	2	2	1	1	1
5.9	3	5.9	3	0	1	0	0	0	2	0
7.3	3	7.3	3	4	0	1	1	1	0	0
0.0	0	0.0	0	2	0	0	3	0	0	0
2.0	1	2.0	1	0	2	1	0	0	0	1

Estimated model	R_{adj}^2
$y = 0.66 + 1.92x_1$	0.83
$y = 0.22 + 1.78x_1 + 0x_2 + 0x_3 + 0x_4 + 0x_5 + 2.11x_6 + 0x_7 + 0x_8$	0.98

Objectif global: imposer une structure sur β

$$\hat{\beta} \in \arg \min_{\beta \in \mathbb{R}^p} \{ \text{DataFitting} [\mathbf{X}, \mathbf{y}, \beta] + \lambda \text{Prior} [\beta] \}$$

Objectif global: imposer une structure sur β

$$\hat{\beta} \in \arg \min_{\beta \in \mathbb{R}^p} \{ \text{DataFitting} [\mathbf{X}, \mathbf{y}, \beta] + \lambda \text{Prior} [\beta] \}$$

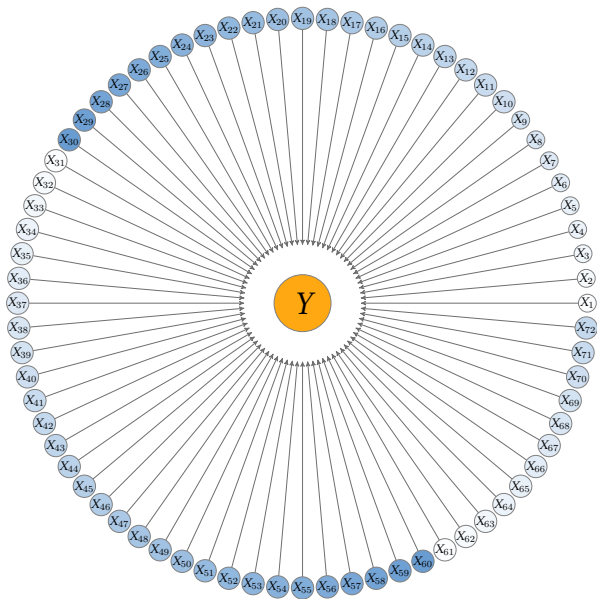
Examples:

$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \|\beta\|_0 \quad (\text{Best subset selection})$$

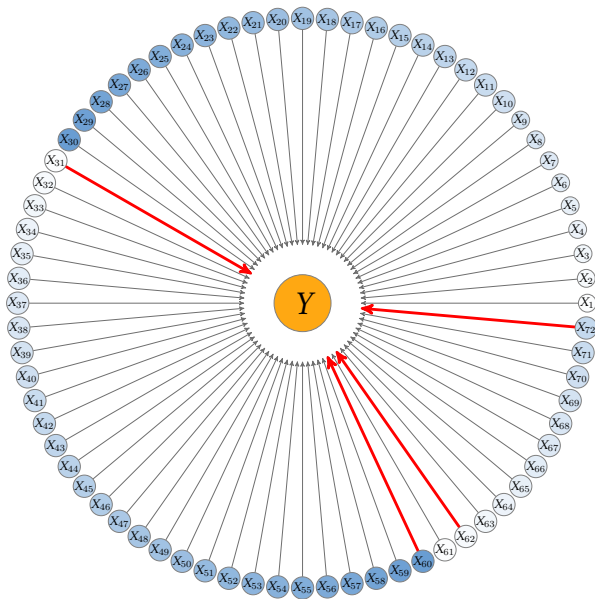
$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \quad (\text{Lasso regression})$$

$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \quad (\text{Ridge regression})$$

Miser sur la sparsité



Miser sur la sparsité



Miser sur la sparsité

Utilisez une procédure qui fonctionne bien pour les problèmes sparse, car aucune procédure ne fonctionne bien pour les problèmes denses.¹

¹The elements of statistical learning. Springer series in statistics, 2001.

Utilisez une procédure qui fonctionne bien pour les problèmes sparse, car aucune procédure ne fonctionne bien pour les problèmes denses.¹

- Un modèle statistique sparse est un modèle pour lequel seulement un petit nombre de variables explicatives jouent un rôle important.
- Hypothèse de parcimonie: peu de variables sont pertinentes pour les données de grande dimension ($N \ll p$).
- β est “creux”
- Les modèles sparse peuvent être plus rapides à calculer, plus faciles à comprendre et produire des prédictions plus stables.

¹The elements of statistical learning. Springer series in statistics, 2001.

Comment organiseriez-vous une réunion de 20 personnes?

		March 2017										
		Thu 9	Fri 10	Sat 11	Sun 12	Mon 13	Tue 14	Wed 15	Thu 16	Fri 17	Sat 18	Sun 19
		5:00 PM – 9:00 PM	5:00 PM – 9:00 PM	9:00 AM – 3:00 PM	3:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM
11 participants	JayZ	✓	✓	✓		✓			✓	✓	✓	
	Evan									✓	✓	✓
	Omar	✓	✓		✓	✓			✓	✓	✓	
	Caitlin	✓	✓	✓					✓	✓	✓	
	Austin	✓	✓	✓								
	Ethan			✓	✓				✓		✓	
	Max	✓	✓	✓		✓			✓	✓	✓	
	Tycho	✓	✓	✓	✓	✓			✓	✓	✓	
	Janevi Chadha		✓	✓	✓	✓	✓			✓	✓	
	Charlotte										✓	✓
	Darshanye	✓	✓			✓			✓	✓		
	Your name	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
		5:00 PM – 9:00 PM	5:00 PM – 9:00 PM	9:00 AM – 3:00 PM	3:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM	1:00 PM – 9:00 PM
		Thu 9	Fri 10	Sat 11	Sun 12	Mon 13	Tue 14	Wed 15	Thu 16	Fri 17	Sat 18	Sun 19
		7	8	7	4	0	6	1	0	7	8	9

Les médecins misent aussi sur la sparsité



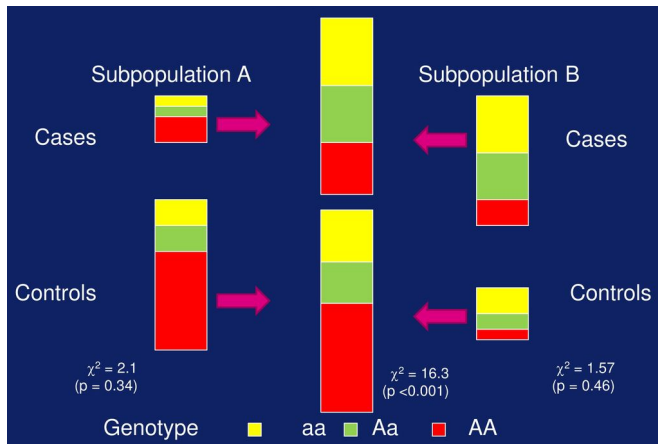


UKBiobank: Données de grande dimension ($n \ll p$)

- Données de géotypage sont issues de 500 000 individus d'origine caucasienne recrutés au Royaume-Uni
- La puce UKBioBANK comporte plus de 800 000 SNPs
- Grand nombre de variables réponses (ex. maladie, densité minérale osseuse)
- Objectif: Quelles variables explicatives sont associées à la variable réponse?

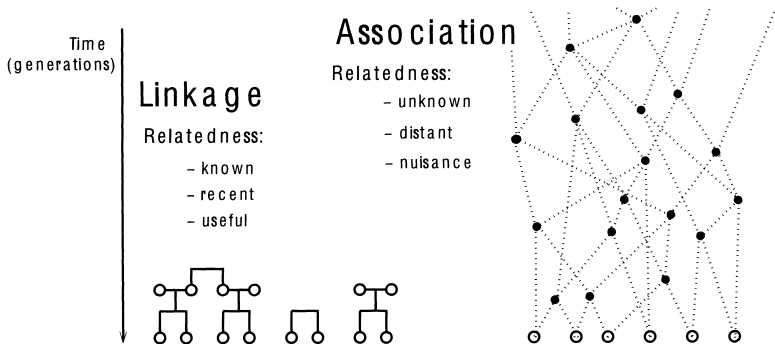


Facteur de confusion



La structure de population

- Les GWAS comparent des individus non apparentés, mais «non apparentés» en fait signifie que les relations sont **inconnues** et présumées éloignées.



¹Astle and Balding. Population structure and cryptic relatedness in genetic association studies. Statistical Science (2009)

Les observations ne sont pas indépendants

- Les observations sont **corrélées**, mais cette relation est souvent **inconnue**
- Cependant, elle peut être **estimé** à partir des données

	ID	Response	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
1	2610781	-1.255	1	2	0	0	0	1
2	4114347	-0.339	1	2	0	2	0	1
3	4399930	-0.6	1	2	1	1	0	1
4	2081319	0.809	1	2	0	1	0	2
5	1347380	0.279	2	2	0	0	0	0
6	3262449	-0.421	2	2	0	1	0	1
7	4870063	-0.454	2	2	0	0	0	2
8	1141212	1.383	2	2	1	1	1	0
9	2997954	-2.29	1	2	0	0	0	1
10	5805218	2.289	1	2	0	1	1	1

La matrice de parenté (kinship)

- Soit $kinship$ une liste de SNP utilisée pour estimer la matrice de parenté
- Soit $X_{kinship}$ une matrice de génotype normalisée $n \times q$.
- Une matrice de parenté (Φ) peut être calculée comme:

$$\Phi = \frac{1}{q-1} X_{kinship} X_{kinship}^{\top} \quad (1)$$

Test d'association avec un modèle mixte linéaire (LMM)

$$\mathbf{Y} = \sum_{j=1}^p \beta_j \cdot \text{SNP}_j + \mathbf{P} + \boldsymbol{\varepsilon} \quad (2)$$

$$\mathbf{P} \sim \mathcal{N}(0, \eta\sigma^2\boldsymbol{\Phi}) \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, (1 - \eta)\sigma^2\mathbf{I})$$

- σ^2 est la variance totale du phénotype
- $\eta \in [0, 1]$ est l'héritabilité du phénotype
- $\mathbf{Y} | (\eta, \sigma^2) \sim \mathcal{N}(\mathbf{0}, \eta\sigma^2\boldsymbol{\Phi} + (1 - \eta)\sigma^2\mathbf{I})$

Régression ridge (Hoerl & Kennard 1970, Technometrics), Lasso (Tibshirani 1996, JRSSB)

- $\widehat{\boldsymbol{\beta}}^{ridge} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_2^2$
- $\widehat{\boldsymbol{\beta}}^{lasso} = \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$

Lasso, ridge, ect. ne sont pas directement applicable au LMM

Procédure en deux étapes

- Étape 1: Ajuster un LMM sous l'hypothèse nul avec un seul effet aléatoire

$$\mathbf{Y} = \mathbf{P} + \boldsymbol{\varepsilon}$$

$$\mathbf{P} \sim \mathcal{N}(0, \eta\sigma^2\boldsymbol{\Phi}) \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, (1 - \eta)\sigma^2\boldsymbol{\mathcal{I}})$$

Procédure en deux étapes

- Étape 1: Ajuster un LMM sous l'hypothèse nul avec un seul effet aléatoire

$$\mathbf{Y} = \mathbf{P} + \boldsymbol{\varepsilon}$$
$$\mathbf{P} \sim \mathcal{N}(0, \eta\sigma^2\boldsymbol{\Phi}) \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, (1 - \eta)\sigma^2\boldsymbol{\mathcal{I}})$$

- Étape 2: Utilisez les résidus de l'étape 1 comme nouvelle réponse *indépendante*

Procédure en deux étapes

X_kinship

	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
ID1	2	2	2	2	2	2
ID2	0	2	2	2	2	2
ID3	0	2	2	2	2	2
ID4	1	2	2	2	2	2
ID5	0	2	2	2	2	2
ID6	1	2	2	2	1	2
ID7	2	2	2	2	1	2
ID8	1	2	2	2	2	2
ID9	0	2	2	2	1	2
ID10	1	2	2	1	2	2

Procédure en deux étapes

X_kinship

	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
ID1	2	2	2	2	2	2
ID2	0	2	2	2	2	2
ID3	0	2	2	2	2	2
ID4	1	2	2	2	2	2
ID5	0	2	2	2	2	2
ID6	1	2	2	2	1	2
ID7	2	2	2	2	1	2
ID8	1	2	2	2	2	2
ID9	0	2	2	2	1	2
ID10	1	2	2	1	2	2



X_kinship **X**_kinship^T

	ID1	ID2	ID3	ID4	ID5	ID6	ID7	ID8	ID9	ID10
ID1	0.97	0	0	0	-0.02	0.03	0.02	-0.01	-0.02	0.03
ID2	0	1	0	-0.01	0	-0.01	-0.01	0	0	0
ID3	0	0	0.98	0.01	0.01	0.01	0	0.03	-0.01	-0.01
ID4	0	-0.01	0.01	1.03	0.04	0.01	-0.01	0.01	0.01	-0.01
ID5	-0.02	0	0.01	0.04	0.97	-0.01	-0.01	0.01	0.03	0.03
ID6	0.03	-0.01	0.01	0.01	-0.01	1.02	0	0	0	0.01
ID7	0.02	-0.01	0	-0.01	-0.01	0	1	0.02	0.02	0
ID8	-0.01	0	0.03	0.01	0.01	0	0.02	1.01	0.01	0
ID9	-0.02	0	-0.01	0.01	0.03	0	0.02	0.01	1.04	0.01
ID10	0.03	0	-0.01	-0.01	0.03	0.01	0	0	0.01	0.95

Procédure en deux étapes

Step 1:

Y		P										
Response		ID1	ID2	ID3	ID4	ID5	ID6	ID7	ID8	ID9	ID10	
-1.255	~	ID1	0.97	0	0	-0.02	0.03	0.02	-0.01	-0.02	0.03	
-0.339		ID2	0	1	0	-0.01	0	-0.01	-0.01	0	0	
-0.6		ID3	0	0	0.98	0.01	0.01	0.01	0	0.03	-0.01	-0.01
0.809		ID4	0	-0.01	0.01	1.03	0.04	0.01	-0.01	0.01	0.01	-0.01
0.279		ID5	-0.02	0	0.01	0.04	0.97	-0.01	-0.01	0.01	0.03	0.03
-0.421		ID6	0.03	-0.01	0.01	0.01	-0.01	1.02	0	0	0	0.01
-0.454		ID7	0.02	-0.01	0	-0.01	-0.01	0	1	0.02	0.02	0
1.383		ID8	-0.01	0	0.03	0.01	0.01	0	0.02	1.01	0.01	0
-2.29		ID9	-0.02	0	-0.01	0.01	0.03	0	0.02	0.01	1.04	0.01
2.289		ID10	0.03	0	-0.01	-0.01	0.03	0.01	0	0	0.01	0.95

+ E_1

Step 2:

Residuals
from Step 1

~

	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
ID1	2	2	2	2	2	2
ID2	0	2	2	2	2	2
ID3	0	2	2	2	2	2
ID4	1	2	2	2	2	2
ID5	0	2	2	2	2	2
ID6	1	2	2	2	1	2
ID7	2	2	2	2	1	2
ID8	1	2	2	2	2	2
ID9	0	2	2	2	1	2
ID10	1	2	2	1	2	2

+ E_2

Procédure en deux étapes

Step 1:

Y		P										
Response		ID1	ID2	ID3	ID4	ID5	ID6	ID7	ID8	ID9	ID10	
-1.255	~	ID1	0.97	0	0	-0.02	0.03	0.02	-0.01	-0.02	0.03	
-0.339		ID2	0	1	0	-0.01	0	-0.01	-0.01	0	0	
-0.6		ID3	0	0	0.98	0.01	0.01	0.01	0	0.03	-0.01	-0.01
0.809		ID4	0	-0.01	0.01	1.03	0.04	0.01	-0.01	0.01	0.01	-0.01
0.279		ID5	-0.02	0	0.01	0.04	0.97	-0.01	-0.01	0.01	0.03	0.03
-0.421		ID6	0.03	-0.01	0.01	0.01	-0.01	1.02	0	0	0	0.01
-0.454		ID7	0.02	-0.01	0	-0.01	-0.01	0	1	0.02	0.02	0
1.383		ID8	-0.01	0	0.03	0.01	0.01	0	0.02	1.01	0.01	0
-2.29		ID9	-0.02	0	-0.01	0.01	0.03	0	0.02	0.01	1.04	0.01
2.289		ID10	0.03	0	-0.01	-0.01	0.03	0.01	0	0	0.01	0.95

+ E_1

Step 2:

Residuals
from Step 1

~

	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
ID1	2	2	2	2	2	2
ID2	0	2	2	2	2	2
ID3	0	2	2	2	2	2
ID4	1	2	2	2	2	2
ID5	0	2	2	2	2	2
ID6	1	2	2	2	1	2
ID7	2	2	2	2	1	2
ID8	1	2	2	2	2	2
ID9	0	2	2	2	1	2
ID10	1	2	2	1	2	2

+ E_2

- Dans les tests d'association, on sait qu'il souffre d'énormes pertes de puissance (Oualkacha et al. Gene. Epi. (2013))

Notre proposition

- Nous proposons, `ggmix`, une procédure en **une seule étape** qui contrôle simultanément les populations structurées et effectue une sélection de variables dans les modèles mixtes linéaires

PLOS GENETICS

RESEARCH ARTICLE

Simultaneous SNP selection and adjustment for population structure in high dimensional prediction models

Sahir R. Bhatnagar^{1,2*}, Yi Yang³, Tianyuan Lu^{4,5}, Erwin Schurr⁶, JC Loredo-Osti⁷, Marie Forest⁸, Karim Oualkacha⁹, Celia M. T. Greenwood^{1,4,5,10,11}

1 Department of Epidemiology, Biostatistics and Occupational Health, McGill University, Montréal, Québec, Canada, **2** Department of Diagnostic Radiology, McGill University, Montréal, Québec, Canada, **3** Department of Mathematics and Statistics, McGill University, Montréal, Québec, Canada, **4** Quantitative Life Sciences, McGill University, Montréal, Québec, Canada, **5** Lady Davis Institute, Jewish General Hospital, Montréal, Québec, Canada, **6** Department of Medicine, McGill University, Montréal, Québec, Canada, **7** Department of Mathematics and Statistics, Memorial University, St. John's, Newfoundland and Labrador, Canada, **8** École de Technologie Supérieure, Montréal, Québec, Canada, **9** Département de Mathématiques, Université du Québec à Montréal, Montréal, Québec, Canada, **10** Gerald Bronfman Department of Oncology, McGill University, Montréal, Québec, Canada, **11** Department of Human Genetics, McGill University, Montréal, Québec, Canada

* sahir.bhatnagar@mcgill.ca



¹R package: sahirbhatnagar.com/ggmix, <https://cran.r-project.org/package=ggmix>

ggmix: une procédure en une seule étape

Y

Response
-1.255
-0.339
-0.6
0.809
0.279
-0.421
-0.454
1.383
-2.29
2.289

~

X

	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
ID1	2	2	2	2	2	2
ID2	0	2	2	2	2	2
ID3	0	2	2	2	2	2
ID4	1	2	2	2	2	2
ID5	0	2	2	2	2	2
ID6	1	2	2	2	1	2
ID7	2	2	2	2	1	2
ID8	1	2	2	2	2	2
ID9	0	2	2	2	1	2
ID10	1	2	2	1	2	2

+

P

	ID1	ID2	ID3	ID4	ID5	ID6	ID7	ID8	ID9	ID10
ID1	0.97	0	0	0	-0.02	0.03	0.02	-0.01	-0.02	0.03
ID2	0	1	0	-0.01	0	-0.01	-0.01	0	0	0
ID3	0	0	0.98	0.01	0.01	0.01	0	0.03	-0.01	-0.01
ID4	0	-0.01	0.01	1.03	0.04	0.01	-0.01	0.01	0.01	-0.01
ID5	-0.02	0	0.01	0.04	0.97	-0.01	-0.01	0.01	0.03	0.03
ID6	0.03	-0.01	0.01	0.01	-0.01	1.02	0	0	0	0.01
ID7	0.02	-0.01	0	-0.01	-0.01	0	1	0.02	0.02	0
ID8	-0.01	0	0.03	0.01	0.01	0	0.02	1.01	0.01	0
ID9	-0.02	0	-0.01	0.01	0.03	0	0.02	0.01	1.04	0.01
ID10	0.03	0	-0.01	-0.01	0.03	0.01	0	0	0.01	0.95

+ **E**

¹R package: sahirbhatnagar.com/ggmix, <https://cran.r-project.org/package=ggmix>

Data and Model

- Phenotype: $\mathbf{Y} = (y_1, \dots, y_n) \in \mathbb{R}^n$
- SNPs: $\mathbf{X} = (\mathbf{X}_1; \dots, \mathbf{X}_n)^T \in \mathbb{R}^{n \times p}$, where $p \gg n$
- Twice the Kinship matrix or Realized Relationship matrix: $\Phi \in \mathbb{R}^{n \times n}$
- Regression Coefficients: $\beta = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$
- Polygenic random effect: $\mathbf{P} = (P_1, \dots, P_n) \in \mathbb{R}^n$
- Error: $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$
- We consider the following LMM with a single random effect:

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{P} + \varepsilon$$
$$\mathbf{P} \sim \mathcal{N}(0, \eta\sigma^2\Phi) \quad \varepsilon \sim \mathcal{N}(0, (1 - \eta)\sigma^2\mathcal{I})$$

- σ^2 is the phenotype total variance
- $\eta \in [0, 1]$ is the phenotype heritability (narrow sens)
- $\mathbf{Y} | (\beta, \eta, \sigma^2) \sim \mathcal{N}(\mathbf{X}\beta, \eta\sigma^2\Phi + (1 - \eta)\sigma^2\mathcal{I})$

Likelihood

- The negative log-likelihood is given by

$$-\ell(\Theta) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \log(\det(\mathbf{V})) + \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\beta)$$

$$\mathbf{V} = \eta \Phi + (1 - \eta) \mathcal{I}$$

- Assume the spectral decomposition of Φ

$$\Phi = \mathbf{U} \mathbf{D} \mathbf{U}^T$$

- \mathbf{U} is an $n \times n$ orthogonal matrix and \mathbf{D} is an $n \times n$ diagonal matrix
- One can write

$$\mathbf{V} = \mathbf{U}(\eta \mathbf{D} + (1 - \eta) \mathcal{I}) \mathbf{U}^T = \mathbf{U} \mathbf{W} \mathbf{U}^T$$

with $\mathbf{W} = \text{diag}(w_i)_{i=1}^n$, $w_i = \eta \mathbf{D}_{ii} + (1 - \eta)$

Likelihood

- Projection of \mathbf{Y} (and columns of \mathbf{X}) into $\text{Span}(\mathbf{U})$ leads to a simplified correlation structure for the transformed data: $\tilde{\mathbf{Y}} = \mathbf{U}^\top \mathbf{Y}$
- $\tilde{\mathbf{Y}} | (\boldsymbol{\beta}, \eta, \sigma^2) \sim \mathcal{N}(\tilde{\mathbf{X}}\boldsymbol{\beta}, \sigma^2 \mathbf{W})$, with $\tilde{\mathbf{X}} = \mathbf{U}^\top \mathbf{X}$
- The negative log-likelihood can then be expressed as

$$-\ell(\boldsymbol{\Theta}) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^n \log(w_i) + \frac{1}{2\sigma^2} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})^\top \mathbf{W}^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})$$

- For fixed σ^2 and η , solving for $\boldsymbol{\beta}$ is a weighted least squares problem

Penalized Maximum Likelihood Estimator

- Define the objective function:

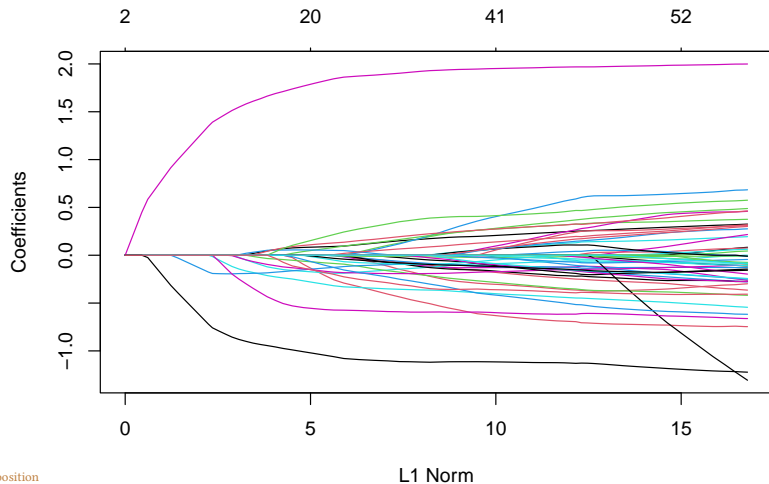
$$Q_\lambda(\Theta) = -\ell(\Theta) + \lambda \sum_j p_j(\beta_j)$$

- $p_j(\cdot)$ is a penalty term on β_1, \dots, β_p
- An estimate of the model parameters $\hat{\Theta}_\lambda$ is obtained by

$$\hat{\Theta}_\lambda = \arg \min_{\Theta} Q_\lambda(\Theta)$$

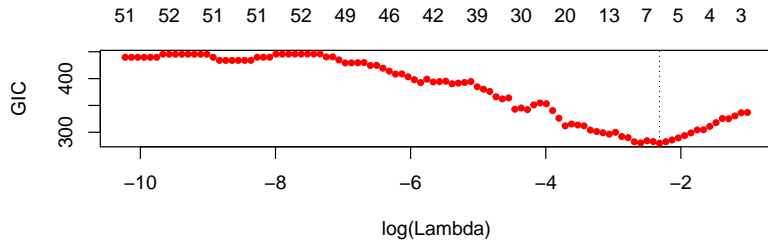
ggmix R package

```
library(ggmix)
data("admixed")
fit <- ggmix(x = admixed$xtrain,
            y = admixed$ytrain,
            kinship = admixed$kin_train)
plot(fit)
```



ggmix R package

```
hdbic <- gic(fit)
plot(hdbic)
```



```
coef(hdbic, type = "nonzero")
```

```
##                1
## (Intercept) -0.03598164
## X302        -0.17617815
## X524         1.34917874
## X538        -0.72073279
## eta          0.99000000
## sigma2      1.60477653
```


Real data applications

1. UK Biobank

- ▶ 10,000 LD-pruned SNPs (Essentially un-correlated variables) to predict standing height in 18k related individuals
- ▶ Standing height is highly polygenic (many variables associated with response)

Real data applications

1. UK Biobank

- ▶ 10,000 LD-pruned SNPs (Essentially un-correlated variables) to predict standing height in 18k related individuals
- ▶ Standing height is highly polygenic (many variables associated with response)

2. GAW20 Simulated dataset

- ▶ 50,000 SNPs (all on chromosome 1) to predict high-density lipoproteins in 679 related individuals
- ▶ Not much correlation between causal SNP and others
- ▶ Very sparse signals (only 1 causal variant)

Real data applications

1. UK Biobank

- ▶ 10,000 LD-pruned SNPs (Essentially un-correlated variables) to predict standing height in 18k related individuals
- ▶ Standing height is highly polygenic (many variables associated with response)

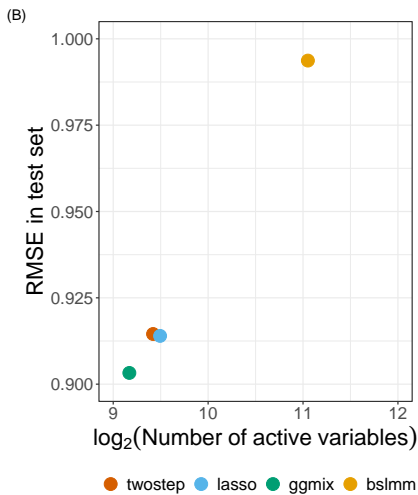
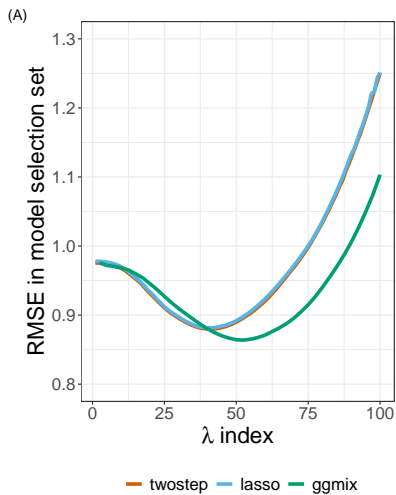
2. GAW20 Simulated dataset

- ▶ 50,000 SNPs (all on chromosome 1) to predict high-density lipoproteins in 679 related individuals
- ▶ Not much correlation between causal SNP and others
- ▶ Very sparse signals (only 1 causal variant)

3. Mouse Crosses

- ▶ Find loci associated with mouse sensitivity to mycobacterial infection
- ▶ 189 samples, and 625 microsatellite markers
- ▶ Highly correlated variables

Results: UK Biobank



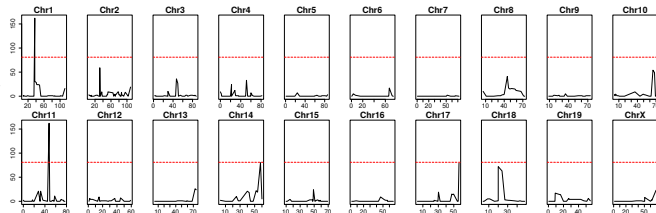
Results: GAW20

Method	Median number of active variables (Inter-quartile range)	RMSE (SD)
twostep	1 (1 - 11)	0.3604 (0.0242)
lasso	1 (1 - 15)	0.3105 (0.0199)
ggmix	1 (1 - 12)	0.3146 (0.0210)
BSLMM	40,737 (39,901 - 41,539)	0.2503 (0.0099)

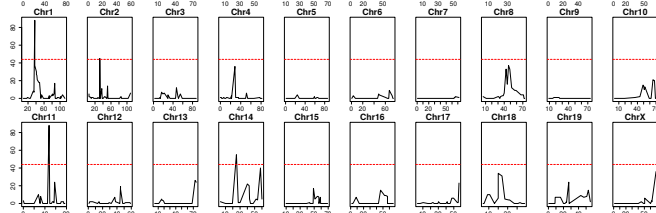
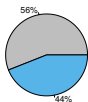
Table: Summary of model performance based on 200 GAW20 simulations. Five-fold cross-validation root-mean-square error was reported for each simulation replicate.

Results: Mouse crosses

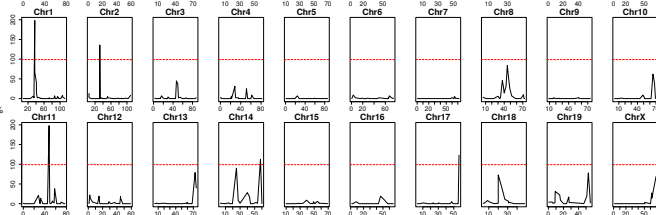
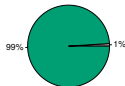
(a) twostep



(b) lasso



(c) ggmix



Discussion

- La procédure en deux étapes conduit à un grand nombre de faux positifs et de faux négatifs
- L'ajustement de la composante principale dans lasso peut ne pas être suffisant pour contrôler la confusion, en particulier lorsqu'il y a beaucoup de corrélation entre les observations
- `ggmix` fonctionne bien même lorsque les variables causales sont utilisées dans le calcul de la matrice de parenté
- `ggmix` a montré la plus grande amélioration par rapport à `twostep` et `lasso` quand il y avait des variables hautement corrélées avec beaucoup de structure (exemple de croix de souris)

Example 1: Prediction study with UK Biobank

- Osteoporosis screening identifies only a small proportion of the screened population to be eligible for intervention to prevent osteoporosis-related fractures
- Much of the screening expenditure is spent on individuals who will not qualify for intervention

PLOS MEDICINE

RESEARCH ARTICLE

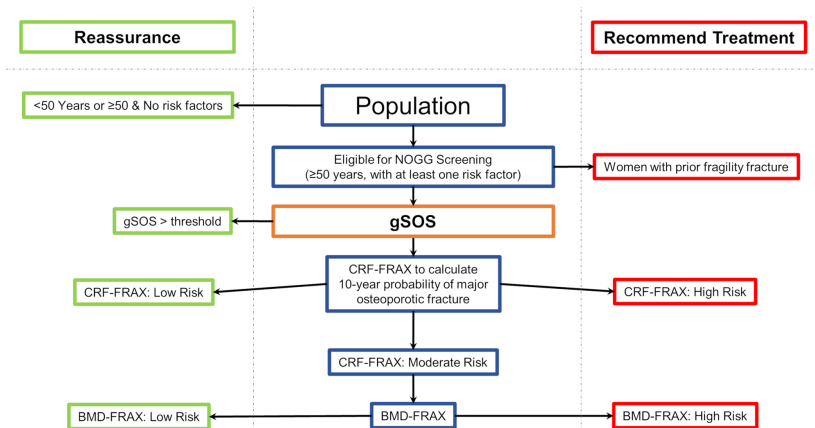
Development of a polygenic risk score to improve screening for fracture risk: A genetic risk prediction study

Vincenzo Forgetta^{1,4}, Julyan Keller-Baruch^{2,4}, Marie Forest¹, Audrey Durand³, Sahir Bhatnagar¹, John P. Kemp^{4,5}, Maria Nethander^{6,7}, Daniel Evans⁹, John A. Morris¹, Douglas P. Kiel⁹, Fernando Rivadeneira¹⁰, Helena Johansson^{11,12}, Nicholas C. Harvey^{13,14}, Dan Mellström⁷, Magnus Karlsson¹⁵, Cyrus Cooper^{13,14,16}, David M. Evans^{4,5}, Robert Clarke¹⁷, John A. Kanis^{11,12}, Eric Orwoll^{18,19}, Eugene V. McCloskey²⁰, Claes Ohlsson⁷, Joelle Pineau³, William D. Leslie²¹, Celia M. T. Greenwood^{1,2,22,23}, J. Brent Richards^{1,2,24*}



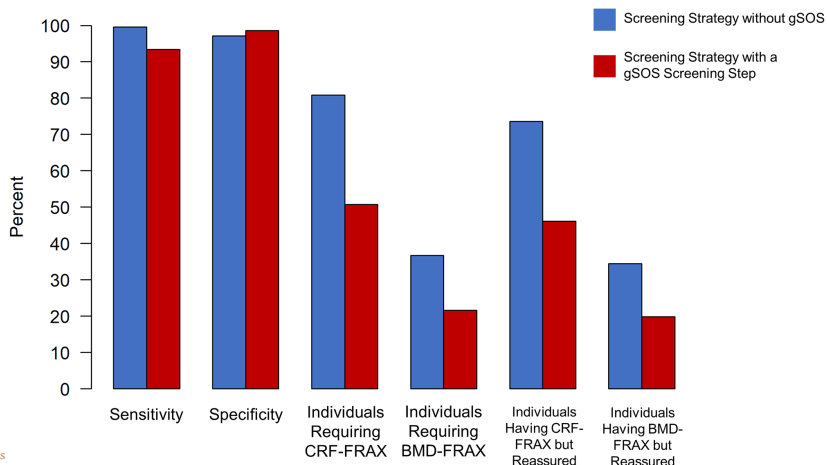
gSOS: a PRS for heel quantitative ultrasound speed of sound (SOS)

- SOS is a heritable risk factor for osteoporotic fracture
- The prediction of SOS using PRS could decrease the number of screened individuals by reassuring those with low genetic risk (*negative screening*)



gSOS: a cheap *negative screening* tool

- 81% of the population required expensive testing to achieve 99.6% sensitivity and 97.1% specificity
- Our polygenic risk score (gSOS) consisting of 21,717 genetic variants, only required 51% of the population while maintaining similar sensitivity and specificity.



Example 2: Identify individuals with rare variants

- An LDL-C PRS could be used to identify individuals with a higher probability of harboring FH variants
- We find that those with a low LDL-C PRS had a 21-fold higher probability of carrying an FH variant compared with those with a high LDL-C PRS

Circulation: Genomic and Precision Medicine

ORIGINAL ARTICLE

Polygenic Risk Score for Low-Density Lipoprotein Cholesterol Is Associated With Risk of Ischemic Heart Disease and Enriches for Individuals With Familial Hypercholesterolemia

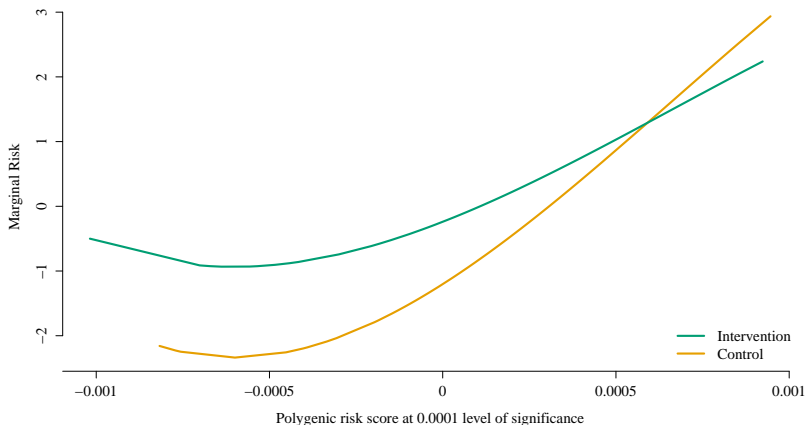
Haoyu Wu¹, MSc; Vincenzo Forgetta², PhD; Sirui Zhou, PhD; Sahir R. Bhatnagar³, PhD; Guillaume Paré⁴, MD; J. Brent Richards⁵, MD



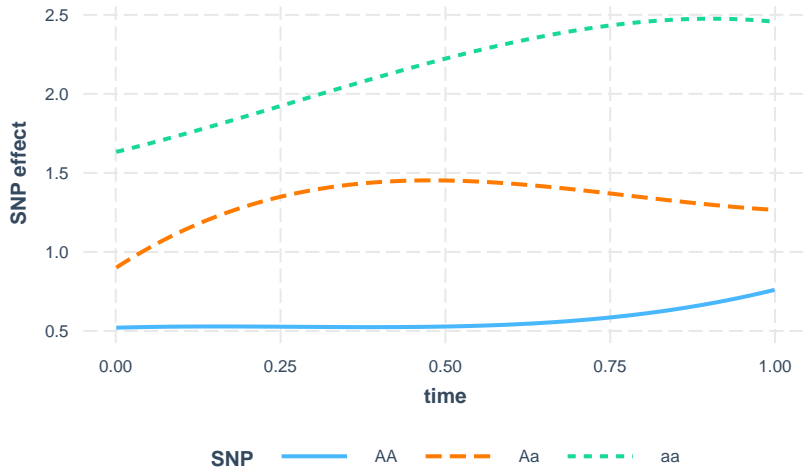
My former MSc student
Haoyu Wu

Nurse Family Partnership (CSDA, 2021+)

- Early intervention in young children has been shown to positively impact intellectual abilities.
- Genome-wide association studies (GWAS) suggest that 20% of the variance in educational attainment (years of education) may be accounted for by common genetic variation.



Canadian Longitudinal Study on Aging (CLSA)

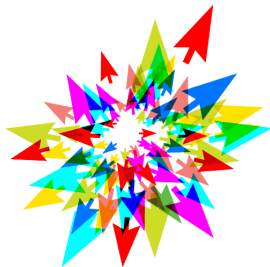


Remerciements

- Haoyu Wu (McGill)
- Tianyuan Lu (McGill)
- Yi Yang (McGill)
- Karim Oualkacha (UQÀM)
- Celia Greenwood (Lady Davis Institute)
- Brent Richards (Lady Davis Institute)
- UK Biobank Resource under project number 27449. We appreciate the generosity of UK Biobank volunteers



compute | **calcul**
canada | canada

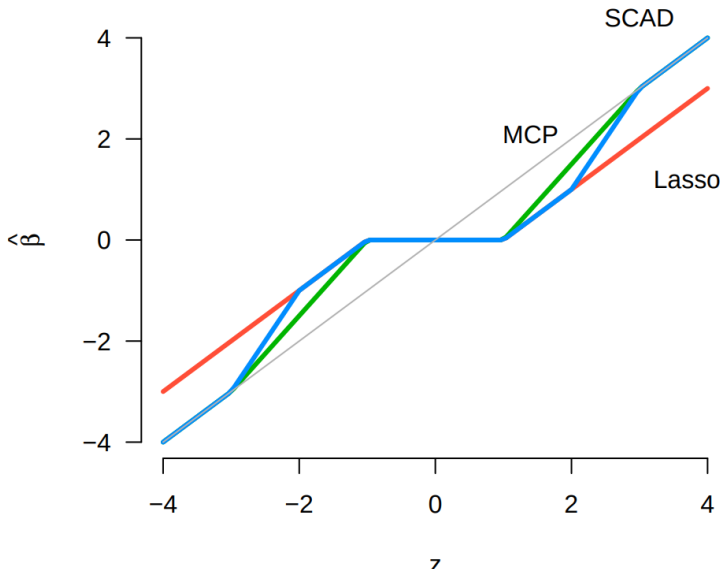


References

1. Bhatnagar SR, Yang Y, Lu T, Schurr E, Loredó-Osti JC, Forest M, Oualkacha K, Greenwood CMT (2020). Simultaneous SNP selection and adjustment for population structure in high dimensional prediction models. *PLoS Genetics* 16(5): e1008766. DOI [10.1371/journal.pgen.1008766](https://doi.org/10.1371/journal.pgen.1008766).
2. Wu H, Forgetta V, Zhou S, Bhatnagar SR, Paré G and Richards JB (2020). A Polygenic Risk Score for Low-density Lipoprotein Cholesterol is Associated with Risk of Ischemic Heart Disease and Enriches for Individuals with Familial Hypercholesterolemia. *Circulation: Genomic and Precision Medicine*
3. Forgetta V, Keller-Baruch J, Forest M, Durand A, Bhatnagar SR, ..., Richards JB (2020). Development of a polygenic risk score to improve screening for fracture risk: A genetic risk prediction study. *Plos Medicine*. DOI [10.1371/journal.pmed.1003152](https://doi.org/10.1371/journal.pmed.1003152)

sahirbhatnagar.com

SCAD (Fan et Li, JASA, 2001), MCP (Zhang, Ann. Stat., 2010)



Computational challenges

- Past approaches for optimization for SCAD/MCP relies upon descent method, first- or second- order
- e.g., sparsenet (Mazumder et al. 2011) uses coordinate descent with full step size, whose coordinate update cycles through $\tilde{\beta}_j = S_{\gamma_k} \left(\sum_{i=1}^n (y_i - \tilde{y}_i^j) x_{ij}, \lambda_\ell \right)$, where $\tilde{y}_i^j = \sum_{k \neq j} x_{ik} \tilde{\beta}_k$
- However, coordinate descent is difficult to vectorize, and rate of convergence is difficult of establish – though past literature suggests $O(1/k)$ rate of convergence for ISTA

Our proposal: Accelerated gradient (AG) method

Improving Convergence for Nonconvex Composite Programming

Kai Yang · Masoud Asgharian · Sahir Bhatnagar

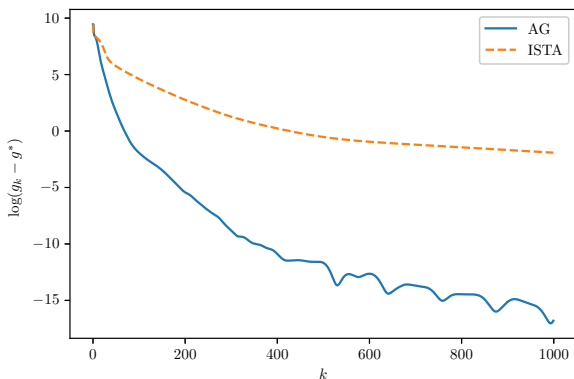
Received: date / Accepted: date

Abstract High-dimensional nonconvex composite problems are popular in today's machine learning and statistical genetics research. Recently, Ghadimi and Lan [1] proposed an algorithm to optimize nonconvex high-dimensional problems. There are several parameters in their algorithm that are to be set before running the algorithm. It is not trivial how to choose these parameters nor there is, to the best of our knowledge, an explicit rule how to select the parameters to make the algorithm converges faster. We analyze Ghadimi and Lan's algorithm to gain an interpretation based on the inequality constraints for convergence and the upper bound for the norm of the gradient analogue. Our interpretation of their algorithm suggests this to be a damped accelerated gradient scheme. Based on this, we propose an approach how to select the parameters to improve convergence of the algorithm. Our numerical studies using high-dimensional nonconvex sparse learning problems, motivated by image denoising and statistical genetics applications, show that convergence can be made, on average, considerably faster than that of the conventional ISTA algorithm for such optimization problems with over 10000 variables should the parameters be chosen using our proposed approach.

Keywords Accelerated Gradient · Composite Optimization · Nonconvex Optimization

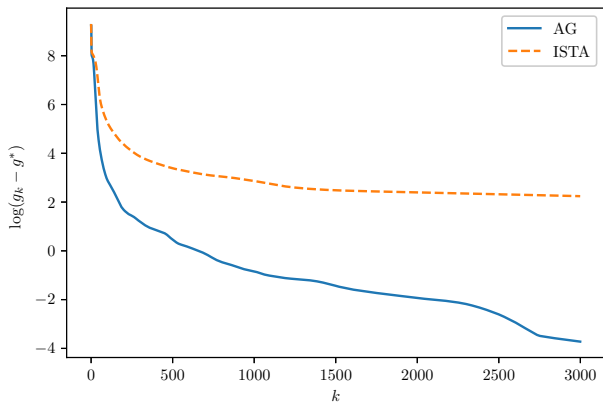
¹<https://arxiv.org/abs/2009.10629>

Numerical Study for SCAD



$\mathbf{x}_i \stackrel{i.i.d.}{\sim} N(\mathbf{0}, \mathbf{I})$, $\varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, $\mathbf{y} = \mathbf{X}\boldsymbol{\tau}_{\text{generate}} + \boldsymbol{\varepsilon}$, $\sigma^2 = \frac{\|\boldsymbol{\tau}_{\text{generate}}\|^2}{3}$,
 $\boldsymbol{\tau}_{\text{generate}} \in \mathbb{R}^{10006}$ is a sparse constant vector with 6 values of 1.23(intercept), 3, 4, 5, 6, 59 as true effect coefficients and 10000 values of 0. Start point: $\boldsymbol{\tau}_0 = \mathbf{1}_{10006}$, $a = 3.7$, $\lambda = 0.6$.

Numerical Study for MCP



Simulation settings here is same as before in SCAD, $\gamma = 2.5$, $\lambda = 0.6$.



Journal of Statistical Software

MMMMMM YYYY, Volume VV, Issue II.

doi: 10.18637/jss.v000.i00

casebase: An Alternative Framework For Survival Analysis and Comparison of Event Rates

Sahir Rai Bhatnagar*
McGill University

Maxime Turgeon*
University of Manitoba

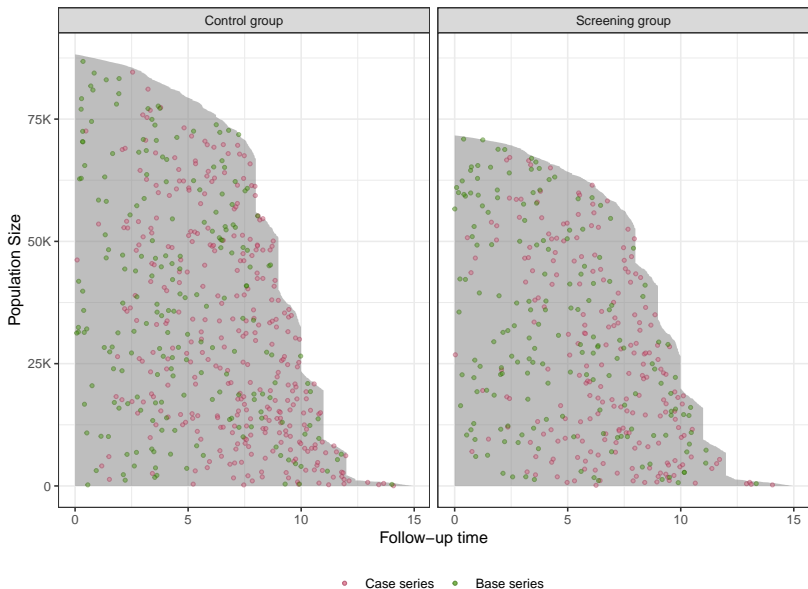
Jesse Islam
McGill University

James A. Hanley
McGill University

Olli Saarela
University of Toronto

¹<https://arxiv.org/abs/2009.10264>,
<https://cran.r-project.org/package=casebase>

Case-base sampling



Case-base sampling

- The unit of analysis is a person-moment.
- Case-base sampling reduces the model fitting to a familiar logistic regression.
- The sampling process is taken into account using an offset term.
- By sampling a large base series, the information loss eventually becomes negligible.
- This framework can easily be used with time-varying covariates (e.g. time-varying exposure). We can fit any hazard λ of the following form:

$$\log \lambda(t; \alpha, \beta) = g(t; \alpha) + \beta X$$

- Different choices of the function g leads to familiar parametric families:
 - ▶ Exponential: g is constant.
 - ▶ Gompertz: $g(t; \alpha) = \alpha t$.
 - ▶ Weibull: $g(t; \alpha) = \alpha \log t$

Orientations futures

- `ggmix` est limité par le nombre d'individus (ne s'applique pas à l'ensemble de la cohorte UK Biobank de 500k) → approximations de rang inférieur de la matrice de parenté
- Problèmes de mémoire lorsque le nombre de covariables dans le modèle dépasse 50k → stratégies de mappage de mémoire (par exemple `biglasso` de Zeng et Breheny (2017))
- Extension aux données multivariées, longitudinales, combinaisons de plusieurs cohortes → Plusieurs effets aléatoires.

Session Info

```
R version 4.1.1 (2021-08-10)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Pop!_OS 21.04

Matrix products: default
BLAS:   /usr/lib/x86_64-linux-gnu/openblas-pthread/libblas.so.3
LAPACK: /usr/lib/x86_64-linux-gnu/openblas-pthread/libopenblas-p-r0.3.13.so

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods    base

other attached packages:
[1] ggmix_0.0.1 knitr_1.36

loaded via a namespace (and not attached):
 [1] lattice_0.20-41  codetools_0.2-18  glmnet_4.1-2     foreach_1.5.1
 [5] grid_4.1.1      magrittr_2.0.1   evaluate_0.14   highr_0.9
 [9] stringi_1.7.5   Matrix_1.3-2     splines_4.1.1   iterators_1.0.13
[13] tools_4.1.1     stringr_1.4.0    survival_3.2-13  xfun_0.26
[17] compiler_4.1.1  shape_1.4.6
```