



Data shared lasso: a simple L1 regularization method for subgroup analyses

Vivian Viallon
IARC

Joint work with Edouard Ollier, Cedric Garcia and Nadim Ballout

Study of association in epidemiology/ clinical research

- General setting
 - Data: $Y, \mathbf{x} \in \mathbb{R}^p$
 - Objective: estimate β^* , e.g. under a logistic model:
 $\text{logit}[\mathbb{P}(Y = 1)] = \mathbf{x}^T \beta^*$.

Study of association in epidemiology/ clinical research

- General setting
 - Data: $Y, \mathbf{x} \in \mathbb{R}^p$
 - Objective: estimate β^* , e.g. under a logistic model:
 $\text{logit}[\mathbb{P}(Y = 1)] = \mathbf{x}^T \beta^*$.

- Subgroup analyses
 - The overall population = K predefined groups (or strata)
 - groups based on “additional” covariates (e.g., gender, age categories)
 - groups based on the outcome (e.g., disease subtypes)
 - Objective: estimate $(\beta_k^*)_{k=1, \dots, K}$: association between \mathbf{x} and Y in group k
 - ⇒ to identify heterogeneities across the groups
 - ⇒ to improve prediction accuracy

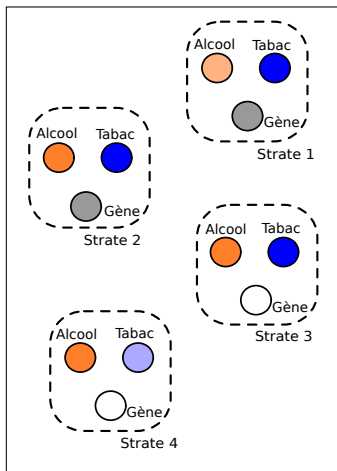
Example 1 : Linear regression on stratified data

- Association between $Y \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^p$ on K predefined strata;
 $Z = 1, \dots, K$.

- On the k -strata, n_k obs. :

$$\mathbf{Y}^{(k)} = \mathbf{X}^{(k)} \boldsymbol{\beta}_k^* + \boldsymbol{\xi}^{(k)}$$

$\Rightarrow Kp$ parameters to be estimated



Example 2 : Case-control studies with multiple subtypes of disease

(Ballout, Garcia and V., submitted)

- ~ Multinomial logistic regression
 - $Y \in \{0, 1, \dots, K\}$
 - $Y = 0$: control
 - $Y = k > 0$: case, of subtype k .
 - Example : K breast cancer histological subtypes (e.g., based on ER/PR status).

⇒ No natural order among the subtypes

$$\log \left(\frac{\mathbb{P}(Y = k | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} \right) = \mathbf{x}^T \boldsymbol{\beta}_k^*$$

Example 3 : matched case-control studies (Ballout, Garcia and V., submitted)

- Same as before, but for each case, one matched control
 - $m = n/2$ pairs of observations, $(\mathbf{x}_i^j, Y_i^j, Z_i^j)_{i=1, \dots, m}^{j=1, 2}$
 - one case, i.e. $Y_i^1 = 1$.
 - one matched control, i.e. $Y_i^2 = 0$

Example 3 : matched case-control studies (Ballout, Garcia and V., submitted)

- Same as before, but for each case, one matched control
 - $m = n/2$ pairs of observations, $(\mathbf{x}_i^j, Y_i^j, Z_i^j)_{i=1, \dots, m}^{j=1, 2}$
 - one case, i.e. $Y_i^1 = 1$.
 - one matched control, i.e. $Y_i^2 = 0$
 - $Z_i^j = k$: subtype of the case

Example 3 : matched case-control studies (Ballout, Garcia and V., submitted)

- Same as before, but for each case, one matched control
 - $m = n/2$ pairs of observations, $(\mathbf{x}_i^j, Y_i^j, Z_i^j)_{i=1, \dots, m}^{j=1, 2}$
 - one case, i.e. $Y_i^1 = 1$.
 - one matched control, i.e. $Y_i^2 = 0$
 - $Z_i^j = k$: subtype of the case

- The global study: K sub-studies
 - ① m_1 pairs: Subtype 1 BC Vs Control
 - ② m_2 pairs: Subtype 2 BC Vs Control
 - ③ ...
 - ④ m_K pairs: Subtype K BC Vs Control

Example 3 : matched case-control studies (Ballout, Garcia and V., submitted)

- Same as before, but for each case, one matched control
 - $m = n/2$ pairs of observations, $(\mathbf{x}_i^j, Y_i^j, Z_i^j)_{i=1, \dots, m}^{j=1, 2}$
 - one case, i.e. $Y_i^1 = 1$.
 - one matched control, i.e. $Y_i^2 = 0$
 - $Z_i^j = k$: subtype of the case

- The global study: K sub-studies
 - ① m_1 pairs: Subtype 1 BC Vs Control $\Rightarrow \beta_1^*$
 - ② m_2 pairs: Subtype 2 BC Vs Control
 - ③ ...
 - ④ m_K pairs: Subtype K BC Vs Control

Example 3 : matched case-control studies (Ballout, Garcia and V., submitted)

- Same as before, but for each case, one matched control
 - $m = n/2$ pairs of observations, $(\mathbf{x}_i^j, Y_i^j, Z_i^j)_{i=1, \dots, m}^{j=1, 2}$
 - one case, i.e. $Y_i^1 = 1$.
 - one matched control, i.e. $Y_i^2 = 0$
 - $Z_i^j = k$: subtype of the case

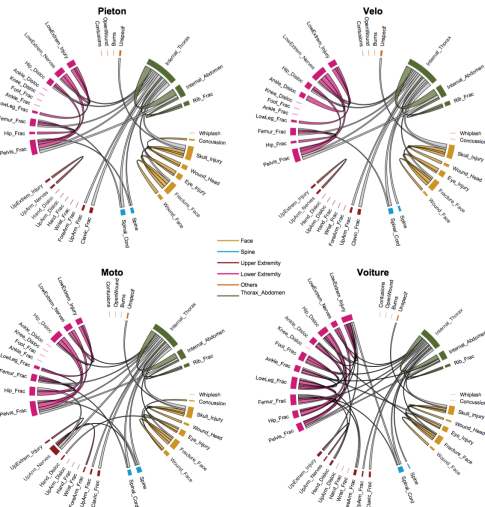
- The global study: K sub-studies
 - ① m_1 pairs: Subtype 1 BC Vs Control $\Rightarrow \beta_1^*$
 - ② m_2 pairs: Subtype 2 BC Vs Control $\Rightarrow \beta_2^*$
 - ③ ...
 - ④ m_K pairs: Subtype K BC Vs Control

Example 3 : matched case-control studies (Ballout, Garcia and V., submitted)

- Same as before, but for each case, one matched control
 - $m = n/2$ pairs of observations, $(\mathbf{x}_i^j, Y_i^j, Z_i^j)_{i=1, \dots, m}^{j=1, 2}$
 - one case, i.e. $Y_i^1 = 1$.
 - one matched control, i.e. $Y_i^2 = 0$
 - $Z_i^j = k$: subtype of the case

- The global study: K sub-studies
 - ① m_1 pairs: Subtype 1 BC Vs Control $\Rightarrow \beta_1^*$
 - ② m_2 pairs: Subtype 2 BC Vs Control $\Rightarrow \beta_2^*$
 - ③ ...
 - ④ m_K pairs: Subtype K BC Vs Control $\Rightarrow \beta_K^*$

Example 4 : K binary graphical models (Ballout and V., Statist. Med., 2019)



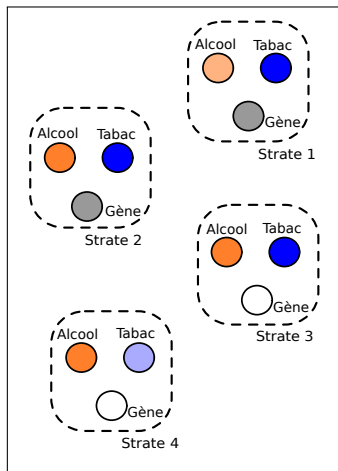
- Association among injuries suffered by victims of road accidents
- groups: \sim road user type

Example 1 : Linear regression on stratified data

- Association between $Y \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^p$ on K predefined strata.
- On the k -strata, n_k obs. :

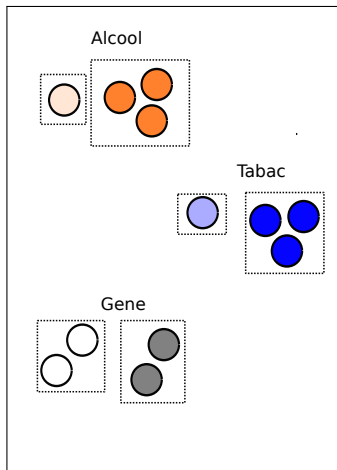
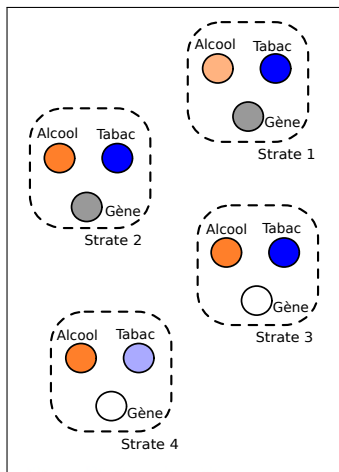
$$\mathbf{Y}^{(k)} = \mathbf{X}^{(k)} \boldsymbol{\beta}_k^* + \boldsymbol{\xi}^{(k)}$$

$\Rightarrow Kp$ parameters to be estimated



Main objective

- for each covariate: to identify the structure of its effects across the strata



Matrix formulation of the model

We have

$$\mathbf{y} = \mathbf{X}\mathbf{b}^* + \boldsymbol{\xi},$$

with

$$\mathbf{y} = \begin{pmatrix} \mathbf{Y}^{(1)} \\ \vdots \\ \mathbf{Y}^{(K)} \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} \mathbf{X}^{(1)} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{X}^{(2)} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & \mathbf{X}^{(K)} \end{pmatrix}, \quad \mathbf{b}^* = \begin{pmatrix} \beta_1^* \\ \vdots \\ \beta_K^* \end{pmatrix}$$

which are elements of \mathbb{R}^n , $\mathbb{R}^{n \times Kp}$ and \mathbb{R}^{Kp} ,
where $n = \sum_{k=1}^K n_k$.

Oracle procedure (p=1)

- $\mathbf{b}^* = (\beta_1^*, \dots, \beta_K^*) \in \mathbb{R}^K$, with
 - $\beta_1^* = 0$
 - β_2^*
 - $\beta_3^* = \dots = \beta_K^*$

$$\mathbf{y} = \begin{pmatrix} X^{(1)} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & X^{(2)} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & X^{(K)} \end{pmatrix} \begin{pmatrix} \beta_1^* \\ \vdots \\ \beta_K^* \end{pmatrix} + \boldsymbol{\xi}$$

Oracle procedure (p=1)

- $\mathbf{b}^* = (\beta_1^*, \dots, \beta_K^*) \in \mathbb{R}^K$, with
 - $\beta_1^* = 0$
 - β_2^*
 - $\beta_3^* = \dots = \beta_K^*$

$$\mathbf{y} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ X^{(2)} & \mathbf{0} \\ \mathbf{0} & X^{(3)} \\ \vdots & \vdots \\ \mathbf{0} & X^{(K)} \end{pmatrix} \begin{pmatrix} \beta_2^* \\ \mu^* \end{pmatrix} + \boldsymbol{\xi}$$

- Oracle procedure: by taking advantage of the homogeneity,
 - reduced complexity: better estimation accuracy
 - greater prediction accuracy
 - greater power to detect $\beta_2^* \neq \mu^*$

Penalized criteria

$$(\hat{\beta}_1, \dots, \hat{\beta}_K) \in \operatorname{argmin}_{(\beta_1, \dots, \beta_K)} \left\{ \sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)} \beta_k\|_2^2}{2n} + \operatorname{pen}(\beta_1, \dots, \beta_K) \right\}$$

Aim: To achieve a trade-off between

- goodness-of-fit
- match with the expected structure among vectors β_k^* :
 - Sparsity: each β_k^* is sparse
 - Homogeneity: for some (k, k', j) , $\beta_{k,j}^* = \beta_{k',j}^*$

Existing approaches... (1/2)

- Group-lasso, trace-norm:
 - ok for prediction error,
 - but not well suited for the identification of partitions

- **Generalized Fused Lasso** (Gertheiss et Tutz, 2012, V. et al. 2016):
 - oracular properties (adaptive version) when Kp is fixed ($n \rightarrow \infty$);
 - what about $Kp \rightarrow \infty$?
 - implementation is “not easy” (extensions to other models, etc.)

Existing approaches... (2/2)

- Standard approaches in epidemiology
 - either estimation on each stratum “independently”,
 - or selection of one reference stratum + interaction tests

Existing approaches... (2/2)

- Standard approaches in epidemiology
 - either estimation on each stratum “independently”,
 - or selection of one reference stratum + interaction tests

- RefLasso with $r = 3$: $\beta_k^* = \beta_3^* + \gamma_k^*$, avec $\gamma_3^* = \mathbf{0}_p$

$$\sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}(\beta_3 + \gamma_k)\|_2^2}{2n} + \lambda_1 \|\beta_3\|_1 + \sum_{k \neq 3} \lambda_{2,k} \|\gamma_k\|_1$$

Existing approaches... (2/2)

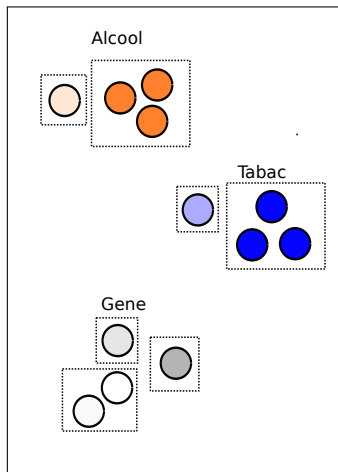
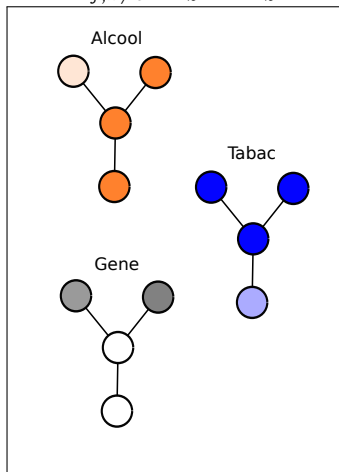
- Standard approaches in epidemiology
 - either estimation on each stratum “independently”,
 - or selection of one reference stratum + interaction tests
- RefLasso with $r = 3$: $\beta_k^* = \beta_3^* + \gamma_k^*$, avec $\gamma_3^* = \mathbf{0}_p$

$$\sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}(\beta_3 + \gamma_k)\|_2^2}{2n} + \lambda_1 \|\beta_3\|_1 + \sum_{k \neq 3} \lambda_{2,k} \|\gamma_k\|_1$$

$$\sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}\beta_k\|_2^2}{2n} + \lambda_1 \|\beta_3\|_1 + \sum_{k \neq 3} \lambda_{2,k} \|\beta_k - \beta_3\|_1$$

RefLasso as a Star-Fused Lasso; example with $r = 3$

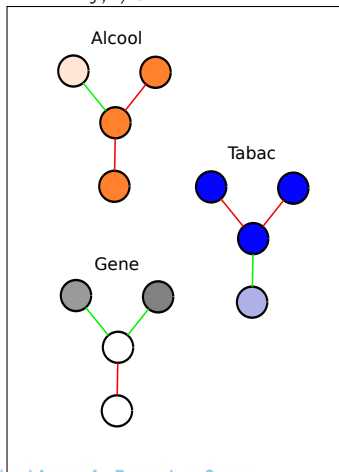
$$\sum_{j,k \neq 3} |\beta_{k,j} - \beta_{3,j}|$$



RefLasso: influence of the reference group

Reference $r = 3$

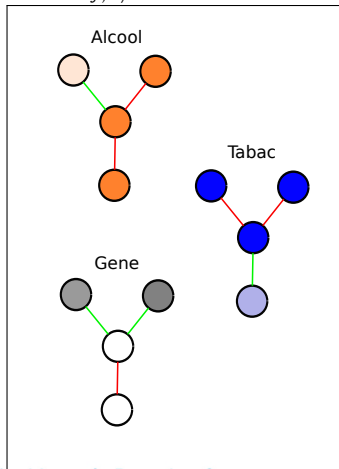
$$\sum_{j, k \neq 3} |\beta_{k,j} - \beta_{3,j}|$$



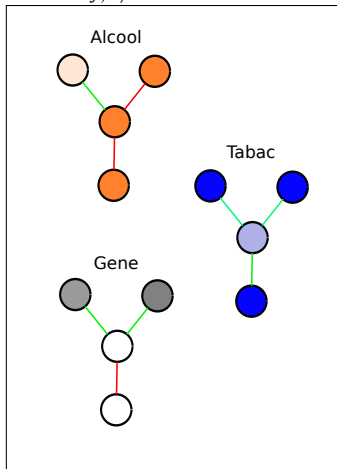
RefLasso: influence of the reference group

Reference $r = 3$

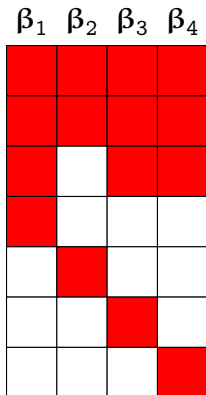
$$\sum_{j,k \neq 3} |\beta_{k,j} - \beta_{3,j}|$$

Reference $r = 4$

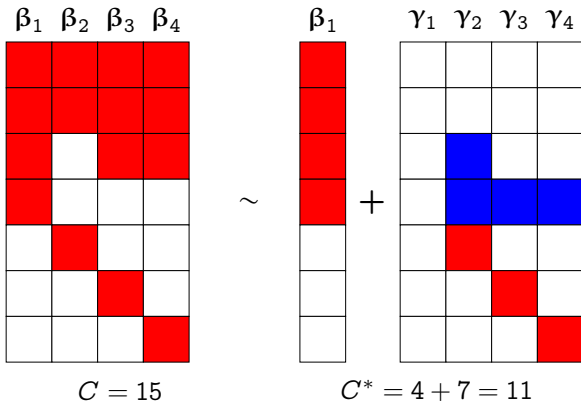
$$\sum_{j,k \neq 4} |\beta_{k,j} - \beta_{4,j}|$$



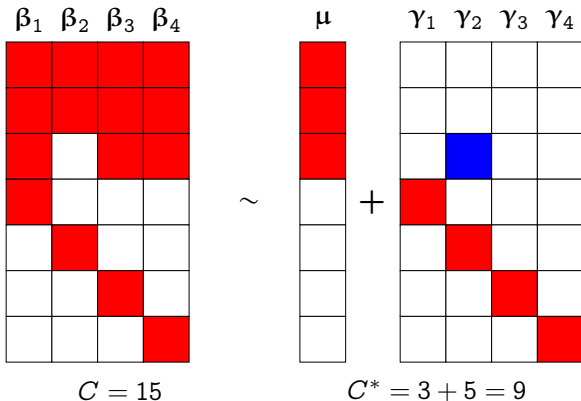
RefLasso = Rank 1 matrix + sparse matrix decomposition



RefLasso = Rank 1 matrix + sparse matrix decomposition



RefLasso = Rank 1 matrix + sparse matrix decomposition



Over-parametrized decomposition of the β_k^* 's

- Consider the **over-parametrization**

$$(\mu^*, \gamma_1^*, \dots, \gamma_K^*) \in \mathbb{R}^{(K+1)p} \quad \text{tq:} \quad \beta_k^* = \mu^* + \gamma_k^*$$

Over-parametrized decomposition of the β_k^* 's

- Consider the **over-parametrization**

$$(\boldsymbol{\mu}^*, \boldsymbol{\gamma}_1^*, \dots, \boldsymbol{\gamma}_K^*) \in \mathbb{R}^{(K+1)p} \quad \text{tq:} \quad \boldsymbol{\beta}_k^* = \boldsymbol{\mu}^* + \boldsymbol{\gamma}_k^*$$

- Particular constraints:

- $\boldsymbol{\mu}^* = \mathbf{0}_p$: independent estimations: complexity $\sum_k \|\boldsymbol{\beta}_k^*\|_0$

Over-parametrized decomposition of the β_k^* 's

- Consider the **over-parametrization**

$$(\mu^*, \gamma_1^*, \dots, \gamma_K^*) \in \mathbb{R}^{(K+1)p} \quad \text{tq:} \quad \beta_k^* = \mu^* + \gamma_k^*$$

- Particular constraints:

- $\mu^* = \mathbf{0}_p$: independent estimations: complexity $\sum_k \|\beta_k^*\|_0$
- $\gamma_1^* = \mathbf{0}_p$: reference stratum: $\mu^* = \beta_1^*$ and $\gamma_k^* = \beta_k^* - \beta_1^*$.

Over-parametrized decomposition of the β_k^* 's

- Consider the **over-parametrization**

$$(\mu^*, \gamma_1^*, \dots, \gamma_K^*) \in \mathbb{R}^{(K+1)p} \quad \text{tq:} \quad \beta_k^* = \mu^* + \gamma_k^*$$

- Particular constraints:

- $\mu^* = \mathbf{0}_p$: independent estimations: complexity $\sum_k \|\beta_k^*\|_0$
- $\gamma_1^* = \mathbf{0}_p$: reference stratum: $\mu^* = \beta_1^*$ and $\gamma_k^* = \beta_k^* - \beta_1^*$.
 - Complexity : $\|\beta_1^*\|_0 + \sum_k \|\beta_k^* - \beta_1^*\|_0$.

Over-parametrized decomposition of the β_k^* 's

- Consider the **over-parametrization**

$$(\mu^*, \gamma_1^*, \dots, \gamma_K^*) \in \mathbb{R}^{(K+1)p} \quad \text{tq:} \quad \beta_k^* = \mu^* + \gamma_k^*$$

- Particular constraints:

- $\mu^* = \mathbf{0}_p$: independent estimations: complexity $\sum_k \|\beta_k^*\|_0$
- $\gamma_1^* = \mathbf{0}_p$: reference stratum: $\mu^* = \beta_1^*$ and $\gamma_k^* = \beta_k^* - \beta_1^*$.
 - Complexity : $\|\beta_1^*\|_0 + \sum_k \|\beta_k^* - \beta_1^*\|_0$.
 - \Rightarrow generally sub-optimal

Over-parametrized decomposition of the β_k^* 's

- Consider the **over-parametrization**

$$(\mu^*, \gamma_1^*, \dots, \gamma_K^*) \in \mathbb{R}^{(K+1)p} \quad \text{tq:} \quad \beta_k^* = \mu^* + \gamma_k^*$$

- Particular constraints:

- $\mu^* = \mathbf{0}_p$: independent estimations: complexity $\sum_k \|\beta_k^*\|_0$
- $\gamma_1^* = \mathbf{0}_p$: reference stratum: $\mu^* = \beta_1^*$ and $\gamma_k^* = \beta_k^* - \beta_1^*$.
 - Complexity : $\|\beta_1^*\|_0 + \sum_k \|\beta_k^* - \beta_1^*\|_0$.
 - \Rightarrow generally sub-optimal
- Optimal constraint: $\mu^* = \mu_{r^*}^*$ where

$$\mu_{r^*,j}^* = \text{mode}(0, \beta_{1,j}^*, \dots, \beta_{K,j}^*) = \beta_{r_j^*,j}^*$$

ie. such that $\{\|\mu^*\|_0 + \sum_k \|\beta_k^* - \mu^*\|_0\}$ is minimized.

Data Shared Lasso

(Ollier and V., *Biometrika*, 2017), Gross and Tibshirani (*CSDA*, 2016)

- Objective: to mimic this "oracle".

$$\sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}(\boldsymbol{\mu} + \boldsymbol{\gamma}_k)\|_2^2}{2n} + \lambda_1 \|\boldsymbol{\mu}\|_1 + \sum_k \lambda_{2,k} \|\boldsymbol{\gamma}_k\|_1$$

Data Shared Lasso

(Ollier and V., *Biometrika*, 2017), Gross and Tibshirani (CSDA, 2016)

- Objective: to mimic this "oracle".

or, equivalently:

$$\sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)} \boldsymbol{\beta}_k\|_2^2}{2n} + \lambda_1 \|\boldsymbol{\mu}\|_1 + \sum_k \lambda_{2,k} \|\boldsymbol{\beta}_k - \boldsymbol{\mu}\|_1$$

$$\begin{aligned} \hat{\mu}_j &= \min_m \left\{ \lambda_1 |m| + \sum_k \lambda_{2,k} |\hat{\beta}_{k,j} - m| \right\} \\ &= \text{WSmedian}(\hat{\beta}_{1,j}, \dots, \hat{\beta}_{K,j}) \end{aligned}$$

Implementation

- Set $\gamma_k = \beta_k - \mu$, $\tau_k = \lambda_{2,k}/\lambda_1$, and introduce

$$\mathbf{x}_0 = \begin{pmatrix} \mathbf{X}^{(1)} & \frac{\mathbf{x}^{(1)}}{\tau_1} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & \mathbf{0} & \cdots & \frac{\mathbf{x}^{(K)}}{\tau_K} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\theta}_0 = \begin{pmatrix} \mu \\ \tau_1 \gamma_1 \\ \vdots \\ \tau_K \gamma_K \end{pmatrix}$$

which belong to $\mathbb{R}^{n \times (K+1)p}$ and $\mathbb{R}^{(K+1)p}$, respectively.

Implementation

- Set $\gamma_k = \beta_k - \mu$, $\tau_k = \lambda_{2,k}/\lambda_1$, and introduce

$$\mathbf{X}_0 = \begin{pmatrix} \mathbf{X}^{(1)} & \frac{\mathbf{X}^{(1)}}{\tau_1} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & \mathbf{0} & \dots & \frac{\mathbf{X}^{(K)}}{\tau_K} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\theta}_0 = \begin{pmatrix} \boldsymbol{\mu} \\ \tau_1 \boldsymbol{\gamma}_1 \\ \vdots \\ \tau_K \boldsymbol{\gamma}_K \end{pmatrix}$$

which belong to $\mathbb{R}^{n \times (K+1)p}$ and $\mathbb{R}^{(K+1)p}$, respectively.

- Data Shared Lasso minimizes

$$\sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}(\boldsymbol{\mu} + \boldsymbol{\gamma}_k)\|_2^2}{2n} + \lambda_1 \|\boldsymbol{\mu}\|_1 + \sum_k \lambda_{2,k} \|\boldsymbol{\gamma}_k\|_1.$$

Implementation

- Set $\gamma_k = \beta_k - \mu$, $\tau_k = \lambda_{2,k}/\lambda_1$, and introduce

$$\mathbf{X}_0 = \begin{pmatrix} \mathbf{X}^{(1)} & \frac{\mathbf{x}^{(1)}}{\tau_1} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & \mathbf{0} & \dots & \frac{\mathbf{x}^{(K)}}{\tau_K} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\theta}_0 = \begin{pmatrix} \mu \\ \tau_1 \gamma_1 \\ \vdots \\ \tau_K \gamma_K \end{pmatrix}$$

which belong to $\mathbb{R}^{n \times (K+1)p}$ and $\mathbb{R}^{(K+1)p}$, respectively.

- Data Shared Lasso minimizes

$$\frac{\|\mathbf{y} - \mathbf{X}_0 \boldsymbol{\theta}_0\|_2^2}{2n} + \lambda_1 \|\boldsymbol{\theta}_0\|_1.$$

⇒ Implementation: straightforward, under a variety of models
(GLMs, survival models, ...)

RefLasso: given one reference stratum (for each j)

- Let $\mathbf{r} = (r_1, \dots, r_p)$ be one vector, indicating the reference stratum for each j .
 - example 1 : $\mathbf{r} = (1, \dots, 1)$
 - example 2 : $\mathbf{r} = \mathbf{r}^* = (r_1^*, \dots, r_p^*)$

RefLasso: given one reference stratum (for each j)

- Let $\mathbf{r} = (r_1, \dots, r_p)$ be one vector, indicating the reference stratum for each j .
 1. example 1 : $\mathbf{r} = (1, \dots, 1)$
 2. example 2 : $\mathbf{r} = \mathbf{r}^* = (r_1^*, \dots, r_p^*)$
- Consider the decomposition $\beta_k^* = \mu_{\mathbf{r}}^* + \gamma_{\mathbf{r},k}^*$ with
 - $\mu_{\mathbf{r}j}^* = \beta_{r_j,j}^*$
 - $\gamma_{\mathbf{r},k}^* = \beta_k^* - \mu_{\mathbf{r}}^*$

RefLasso: given one reference stratum (for each j)

- Let $\mathbf{r} = (r_1, \dots, r_p)$ be one vector, indicating the reference stratum for each j .
 - example 1 : $\mathbf{r} = (1, \dots, 1)$
 - example 2 : $\mathbf{r} = \mathbf{r}^* = (r_1^*, \dots, r_p^*)$
- Consider the decomposition $\beta_k^* = \mu_{\mathbf{r}}^* + \gamma_{\mathbf{r},k}^*$ with
 - $\mu_{r_j}^* = \beta_{r_j,j}^*$
 - $\gamma_{\mathbf{r},k}^* = \beta_k^* - \mu_{\mathbf{r}}^*$
- Estimation of these parameters can be obtained by minimizing

$$\sum_{k=1}^K \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}(\boldsymbol{\mu} + \boldsymbol{\gamma}_k)\|_2^2}{2n} + \lambda_1 \|\boldsymbol{\mu}\|_1 + \sum_{k=1}^K \lambda_{2,k} \|\boldsymbol{\gamma}_k\|_1,$$

under the constraints $\gamma_{r_j,j} = 0$ for all $j \in [p]$

- Recall $\gamma_k = \beta_k - \mu$, $\tau_k = \lambda_{2,k}/\lambda_1$

$$\mathbf{x}_0 = \begin{pmatrix} \mathbf{X}^{(1)} & \frac{\mathbf{X}^{(1)}}{\tau_1} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & \mathbf{0} & \cdots & \frac{\mathbf{X}^{(K)}}{\tau_K} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\theta}_0 = \begin{pmatrix} \mu \\ \tau_1 \gamma_1 \\ \vdots \\ \tau_K \gamma_K \end{pmatrix}$$

qui appartiennent à $\mathbb{R}^{n \times (K+1)p}$ et $\mathbb{R}^{(K+1)p}$.

- Data Shared Lasso (AutoRefLasso), AprioriRefLasso (r chosen a priori) et OracleRefLasso ($r = r^*$ optimal) all minimize

$$\frac{\|\mathbf{y} - \mathbf{x}\boldsymbol{\theta}\|_2^2}{2n} + \lambda_1 \|\boldsymbol{\theta}\|_1$$

for one particular choice " $\mathbf{X} \subseteq \mathbf{X}_0$ ".

⇒ Unified framework to study and compare them (and implement them)

Irrepresentability conditions (Ollier and V., Biometrika, 2017)

- A priori : $\theta_r^* = (\mu_r^{*T}, \tau_1 \gamma_{r,1}^{*T}, \dots, \tau_K \gamma_{r,K}^{*T})^T \in \mathbb{R}^{Kp}$.
- Oracle : $\theta_{r^*}^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{r^*,1}^{*T}, \dots, \tau_K \gamma_{r^*,K}^{*T})^T \in \mathbb{R}^{Kp}$
- Auto : $\theta_0^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{0,1}^{*T}, \dots, \tau_K \gamma_{0,K}^{*T})^T \in \mathbb{R}^{(K+1)p}$

Irrepresentability conditions (Ollier and V., Biometrika, 2017)

- A priori : $\theta_r^* = (\mu_r^{*T}, \tau_1 \gamma_{r,1}^{*T}, \dots, \tau_K \gamma_{r,K}^{*T})^T \in \mathbb{R}^{Kp}$.
 - Oracle : $\theta_{r^*}^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{r^*,1}^{*T}, \dots, \tau_K \gamma_{r^*,K}^{*T})^T \in \mathbb{R}^{Kp}$
 - Auto : $\theta_0^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{0,1}^{*T}, \dots, \tau_K \gamma_{0,K}^{*T})^T \in \mathbb{R}^{(K+1)p}$
 - $\text{supp}(\theta_r^*) = J_r$ and $\text{supp}(\theta_0^*) = J_0$
- $\Leftrightarrow S_r = \{j : \beta_{r_j,j}^* \neq 0\}, T_r = \{(k,j) \in [K] \times [p] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}$
- $\theta_{r^* J_{r^*}}^* = \theta_{0 J_0}^*$ and so $\mathcal{X}_{r^* J_{r^*}} = \mathcal{X}_{0 J_0}$

Irrepresentability conditions (Ollier and V., Biometrika, 2017)

- A priori : $\theta_r^* = (\mu_r^{*T}, \tau_1 \gamma_{r,1}^{*T}, \dots, \tau_K \gamma_{r,K}^{*T})^T \in \mathbb{R}^{Kp}$.
 - Oracle : $\theta_{r^*}^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{r^*,1}^{*T}, \dots, \tau_K \gamma_{r^*,K}^{*T})^T \in \mathbb{R}^{Kp}$
 - Auto : $\theta_0^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{0,1}^{*T}, \dots, \tau_K \gamma_{0,K}^{*T})^T \in \mathbb{R}^{(K+1)p}$
 - $\text{supp}(\theta_r^*) = J_r$ and $\text{supp}(\theta_0^*) = J_0$
- $\Leftrightarrow S_r = \{j : \beta_{r_j,j}^* \neq 0\}, T_r = \{(k,j) \in [K] \times [p] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}$
- $\theta_{r^* J_{r^*}}^* = \theta_{0 J_0}^*$ and so $\mathcal{X}_{r^* J_{r^*}} = \mathcal{X}_{0 J_0}$

$$(IC)_r : \Lambda_{\min}(\mathcal{X}_{r J_r}^T \mathcal{X}_{r J_r}) \geq C_r > 0 \text{ et}$$

$$c_r = \max_{j \notin J_r} \|(\mathcal{X}_{r J_r}^T \mathcal{X}_{r J_r})^{-1} \mathcal{X}_{r J_r}^T \mathcal{X}_{r j}\|_1 < 1,$$

$$(IC)_0 : \Lambda_{\min}(\mathcal{X}_{0 J_0}^T \mathcal{X}_{0 J_0}) \geq C_0 > 0 \text{ et}$$

$$c_0 = \max_{j \notin J_0} \|(\mathcal{X}_{0 J_0}^T \mathcal{X}_{0 J_0})^{-1} \mathcal{X}_{0 J_0}^T \mathcal{X}_{0 j}\|_1 < 1.$$

Irrepresentability conditions (Ollier and V., Biometrika, 2017)

- A priori : $\theta_r^* = (\mu_r^{*T}, \tau_1 \gamma_{r,1}^{*T}, \dots, \tau_K \gamma_{r,K}^{*T})^T \in \mathbb{R}^{Kp}$.
 - Oracle : $\theta_{r^*}^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{r^*,1}^{*T}, \dots, \tau_K \gamma_{r^*,K}^{*T})^T \in \mathbb{R}^{Kp}$
 - Auto : $\theta_0^* = (\mu_{r^*}^{*T}, \tau_1 \gamma_{0,1}^{*T}, \dots, \tau_K \gamma_{0,K}^{*T})^T \in \mathbb{R}^{(K+1)p}$
 - $\text{supp}(\theta_r^*) = J_r$ and $\text{supp}(\theta_0^*) = J_0$
- $\Leftrightarrow S_r = \{j : \beta_{r_j,j}^* \neq 0\}, T_r = \{(k,j) \in [K] \times [p] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}$
- $\theta_{r^* J_{r^*}}^* = \theta_{0 J_0}^*$ and so $\mathcal{X}_{r^* J_{r^*}} = \mathcal{X}_{0 J_0}$

$$(IC)_r : \Lambda_{\min}(\mathcal{X}_{r J_r}^T \mathcal{X}_{r J_r}) \geq C_r > 0 \text{ et}$$

$$c_r = \max_{j \notin J_r} \|(\mathcal{X}_{r J_r}^T \mathcal{X}_{r J_r})^{-1} \mathcal{X}_{r J_r}^T \mathcal{X}_{r j}\|_1 < 1,$$

$$(IC)_0 : \Lambda_{\min}(\mathcal{X}_{0 J_0}^T \mathcal{X}_{0 J_0}) \geq C_0 > 0 \text{ et}$$

$$c_0 = \max_{j \notin J_0} \|(\mathcal{X}_{0 J_0}^T \mathcal{X}_{0 J_0})^{-1} \mathcal{X}_{0 J_0}^T \mathcal{X}_{0 j}\|_1 < 1.$$

The particular case where $n_k = n/K$ and $(\mathbf{X}^{(k)T} \mathbf{X}^{(k)})/n_k = \mathbf{I}_{n_k}$

- $S_r = \{j : \beta_{r_j,j}^* \neq 0\}$, $T_r = \{(k, j) \in [K] \times [p] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}$
- $\mathcal{D}_{r,0} = \max_{j \notin S_r} |\{k \in [K] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}|$ (or 0)
- $\mathcal{D}_{r,1} = \max_{j \in S_r} |\{k \in [K] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}|$ (or $-\infty$)

The particular case where $n_k = n/K$ and $(\mathbf{X}^{(k)T} \mathbf{X}^{(k)})/n_k = \mathbf{I}_{n_k}$

- $S_r = \{j : \beta_{r_j,j}^* \neq 0\}$, $T_r = \{(k, j) \in [K] \times [p] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}$
- $\mathcal{D}_{r,0} = \max_{j \notin S_r} |\{k \in [K] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}|$ (or 0)
- $\mathcal{D}_{r,1} = \max_{j \in S_r} |\{k \in [K] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}|$ (or $-\infty$)
- Let $\tau_k = \tau_0 K^{-1/2}$ for some $\tau_0 > 0$,

$$(sIC)_r : \quad 0 \leq \frac{K^{1/2}}{K - 2\mathcal{D}_{r,1}} < \tau_0 < \frac{K^{1/2}}{\mathcal{D}_{r,0}}.$$

$$(sIC)_0 : \quad 0 \leq \frac{K^{1/2}}{K - 2\mathcal{D}_{r^*,1}} < \tau_0 < \frac{K^{1/2}}{\mathcal{D}_{r^*,0}}.$$

The particular case where $n_k = n/K$ and $(\mathbf{X}^{(k)T} \mathbf{X}^{(k)})/n_k = \mathbf{I}_{n_k}$

- $S_r = \{j : \beta_{r_j,j}^* \neq 0\}$, $T_r = \{(k, j) \in [K] \times [p] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}$
- $\mathcal{D}_{r,0} = \max_{j \notin S_r} |\{k \in [K] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}|$ (or 0)
- $\mathcal{D}_{r,1} = \max_{j \in S_r} |\{k \in [K] : \beta_{k,j}^* \neq \beta_{r_j,j}^*\}|$ (or $-\infty$)
- Let $\tau_k = \tau_0 K^{-1/2}$ for some $\tau_0 > 0$,

$$(sIC)_r : \quad 0 \leq \frac{K^{1/2}}{K - 2\mathcal{D}_{r,1}} < \tau_0 < \frac{K^{1/2}}{\mathcal{D}_{r,0}}.$$

$$(sIC)_0 : \quad 0 \leq \frac{K^{1/2}}{K - 2\mathcal{D}_{r^*,1}} < \tau_0 < \frac{K^{1/2}}{\mathcal{D}_{r^*,0}}.$$

- $(sIC)_0 = (sIC)_{r^*}$
- If $\mathcal{D}_{r,0} = \mathcal{D}_{r,1} = \mathcal{D}_r$ then $(sIC)_r \Rightarrow \mathcal{D}_r < K/3$
- $(sIC)_r$ generally stronger than $(sIC)_{r^*}$

Theorem 1 (Ollier and V., Biometrika, 2017)

- Assume that $\varepsilon_i^{(k)}$: i.i.d. centered sub-Gaussian variables with parameter $\sigma > 0$.
- Under $(sIC)_0$, introduce:

$$\gamma = \min \left\{ 1 - \mathcal{D}_0 \tau_0 K^{-1/2}, 1 - \frac{K^{1/2} + \mathcal{D}_1 \tau_0}{(K - \mathcal{D}_1) \tau_0} \right\}$$

$$C_{\min} = \min \left(1, \tau_0^{-2}, \frac{1}{2} \left[(\tau_0^{-2} + 1) - \left\{ (\tau_0^{-2} - 1)^2 + \frac{4\mathcal{D}_1}{\tau_0^2 K} \right\} \right] \right)$$

$$\lambda_1 > \frac{2}{\gamma \min(1, \tau_0)} \left[\frac{2\sigma^2 \log\{(K+1)p\}}{n} \right]^{1/2}, \quad \lambda_{2,k} = \tau_k \lambda_1$$

$$\beta_{\min} = \lambda_1 [(|S_{r^*}| + |T_{r^*}|)^{1/2} C_{\min}^{-1} + 4\sigma C_{\min}^{-1/2}],$$

- Assume that $|\beta_{r_j^*, j}^*| > \beta_{\min}$ and $|\beta_{k,j}^* - \beta_{r^*, j}^*| > \frac{K^{1/2} \beta_{\min}}{\tau_0}$.

Then, with high probability,

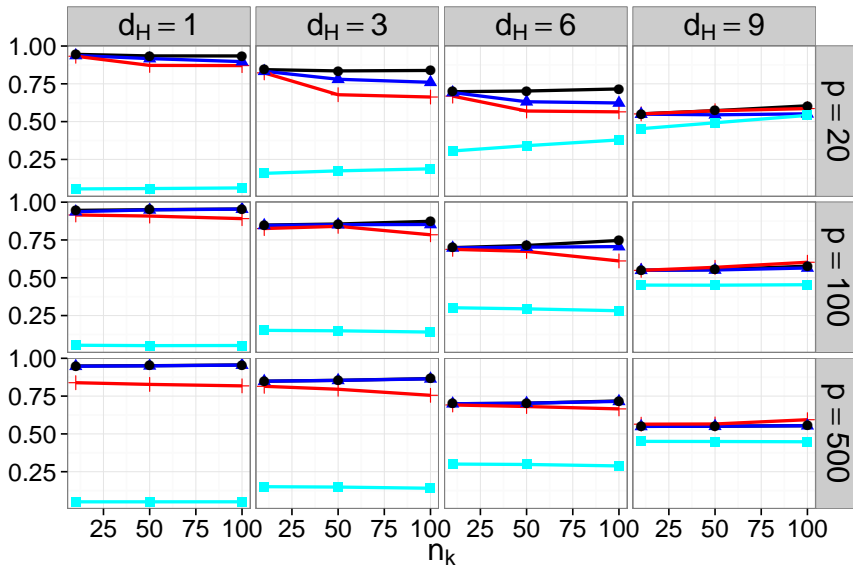
- S_{r^*} and T_{r^*} are both recovered;
- $\|\hat{\theta}_{0J_0} - \theta_{0J_0}^*\|_{\infty} \leq \beta_{\min}$.

Simulations in the Linear regression case

- $K = 20$; $\mathbf{x}_i^{(k)} \sim \mathcal{N}(\mathbf{0}_p, \Sigma)$.
- $P_0 \subseteq [p]$, with $|P_0| = 15$

$$\beta_{k,j}^* = \begin{cases} 0 & \forall j \in [p] \setminus P_0 \\ 1 & j \in P_0 \text{ \& } k > d_H \\ 1 + K^{1/2} & j \in P_0 \text{ \& } k \leq d_H \end{cases}$$

- Accuracy regarding:
 - $T_{P_0}^* = \{(k, j) \in [K] \times P_0 : \beta_{k,j}^* \neq \beta_{r_j^*, j}^*\}$
 - $T_{1, P_0}^* = \{(k, j) \in [K] \times P_0 : \beta_{k,j}^* \neq \beta_{1,j}^*\}$



ORefLasso
 AutoRefLasso
 RefLasso
 CliqueFused

Sparse multinomial Logistic regression (Ballout, Garcia and V., submitted)

- In unmatched designs, if
 - $Y = 0$: control [“natural” reference category]
 - $Y = k$: case of subtype k , $k = 1, \dots, K$
 - no natural order among the subtypes.

⇒ Sparse multinomial logistic regression is a natural extension of sparse (binary) logistic regression.

Sparse multinomial Logistic regression (Ballout, Garcia and V., submitted)

- In unmatched designs, if
 - $Y = 0$: control [“natural” reference category]
 - $Y = k$: case of subtype k , $k = 1, \dots, K$
 - no natural order among the subtypes.

⇒ Sparse multinomial logistic regression is a natural extension of sparse (binary) logistic regression.

- Two formulations exist. No clear guidance in the literature on which one to chose in practice.

Two formulations

- ① Select a reference category (say 0) and then assume the existence of K parameter vectors $\boldsymbol{\gamma}_1^*, \dots, \boldsymbol{\gamma}_K^* \in \mathbb{R}^p$ s.t.

$$\mathbb{P}(Y = \ell | \mathbf{x}) = \frac{\{\exp(\mathbf{x}^T \boldsymbol{\delta}_\ell^*)\}^{\mathbf{I}(\ell \neq 0)}}{1 + \sum_{k=1}^K \exp(\mathbf{x}^T \boldsymbol{\delta}_k^*)}.$$

Two formulations

- 1 Select a reference category (say 0) and then assume the existence of K parameter vectors $\boldsymbol{\gamma}_1^*, \dots, \boldsymbol{\gamma}_K^* \in \mathbb{R}^p$ s.t.

$$\mathbb{P}(Y = \ell | \mathbf{x}) = \frac{\{\exp(\mathbf{x}^T \boldsymbol{\delta}_\ell^*)\}^{\mathbf{I}(\ell \neq 0)}}{1 + \sum_{k=1}^K \exp(\mathbf{x}^T \boldsymbol{\delta}_k^*)}.$$

SparseMultinomialRef0 maximizes

$$L(\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_K) - \lambda \sum_{k=1}^K \|\boldsymbol{\gamma}_k\|_1.$$

Two formulations

- ① Select a reference category (say 0) and then assume the existence of K parameter vectors $\boldsymbol{\gamma}_1^*, \dots, \boldsymbol{\gamma}_K^* \in \mathbb{R}^p$ s.t.

$$\mathbb{P}(Y = \ell | \mathbf{x}) = \frac{\{\exp(\mathbf{x}^T \boldsymbol{\delta}_\ell^*)\}^{\mathbf{I}(\ell \neq 0)}}{1 + \sum_{k=1}^K \exp(\mathbf{x}^T \boldsymbol{\delta}_k^*)}.$$

- ② Use a symmetric over-parametrization instead: assume the existence of $K + 1$ vectors $\{\boldsymbol{\beta}_0^*, \dots, \boldsymbol{\beta}_K^*\}$ s.t.

$$\mathbb{P}(Y = \ell | \mathbf{x}) = \frac{\exp(\mathbf{x}^T \boldsymbol{\beta}_\ell^*)}{\sum_{k=0}^K \exp(\mathbf{x}^T \boldsymbol{\beta}_k^*)}.$$

$$(\boldsymbol{\delta}_k^* = \boldsymbol{\beta}_k^* - \boldsymbol{\beta}_0^*)$$

Two formulations

- ① Select a reference category (say 0) and then assume the existence of K parameter vectors $\boldsymbol{\gamma}_1^*, \dots, \boldsymbol{\gamma}_K^* \in \mathbb{R}^p$ s.t.

$$\mathbb{P}(Y = \ell | \mathbf{x}) = \frac{\{\exp(\mathbf{x}^T \boldsymbol{\delta}_\ell^*)\}^{\mathbf{I}(\ell \neq 0)}}{1 + \sum_{k=1}^K \exp(\mathbf{x}^T \boldsymbol{\delta}_k^*)}.$$

- ② Use a symmetric over-parametrization instead: assume the existence of $K + 1$ vectors $\{\boldsymbol{\beta}_0^*, \dots, \boldsymbol{\beta}_K^*\}$ s.t.

$$\mathbb{P}(Y = \ell | \mathbf{x}) = \frac{\exp(\mathbf{x}^T \boldsymbol{\beta}_\ell^*)}{\sum_{k=0}^K \exp(\mathbf{x}^T \boldsymbol{\beta}_k^*)}.$$

$$(\boldsymbol{\delta}_k^* = \boldsymbol{\beta}_k^* - \boldsymbol{\beta}_0^*)$$

`glmnet` maximizes

$$\mathcal{L}(\boldsymbol{\beta}_0, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K) - \lambda \sum_{k=0}^K \|\boldsymbol{\beta}_k\|_1.$$

Link between the two formulations

- Data Shared Version of MultinomRef0:
 - write $\delta_k = \mu + \gamma_k$
 - and maximize

$$L(\mu + \gamma_1, \dots, \mu + \gamma_K) - \lambda \|\mu\|_1 - \lambda \sum_{k=1}^K \|\gamma_k\|_1$$

Link between the two formulations

- Data Shared Version of MultinomRef0:
 - write $\delta_k = \mu + \gamma_k$
 - and maximize

$$\begin{aligned}
 & L(\mu + \gamma_1, \dots, \mu + \gamma_K) - \lambda \|\mu\|_1 - \lambda \sum_{k=1}^K \|\gamma_k\|_1 \\
 & = \mathcal{L}(\beta_0, \beta_1, \dots, \beta_K) - \lambda \sum_{k \geq 0} \|\beta_k\|_1
 \end{aligned}$$

Link between the two formulations

- Data Shared Version of MultinomRef0:
 - write $\delta_k = \mu + \gamma_k$
 - and maximize

$$\begin{aligned}
 & L(\mu + \gamma_1, \dots, \mu + \gamma_K) - \lambda \|\mu\|_1 - \lambda \sum_{k=1}^K \|\gamma_k\|_1 \\
 & = \mathcal{L}(\beta_0, \beta_1, \dots, \beta_K) - \lambda \sum_{k \geq 0} \|\beta_k\|_1
 \end{aligned}$$

- glmnet returns $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_K)$, which is exactly the estimates $(-\hat{\mu}, \hat{\gamma}_1, \dots, \hat{\gamma}_K)$ we would obtain by performing MultinomRef0 with Data Shared lasso

Link between the two formulations

- Data Shared Version of MultinomRef0:
 - write $\delta_k = \mu + \gamma_k$
 - and maximize

$$\begin{aligned} L(\mu + \gamma_1, \dots, \mu + \gamma_K) - \lambda \|\mu\|_1 - \lambda \sum_{k=1}^K \|\gamma_k\|_1 \\ = \mathcal{L}(\beta_0, \beta_1, \dots, \beta_K) - \lambda \sum_{k \geq 0} \|\beta_k\|_1 \end{aligned}$$

- glmnet returns $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_K)$, which is exactly the estimates $(-\hat{\mu}, \hat{\gamma}_1, \dots, \hat{\gamma}_K)$ we would obtain by performing MultinomRef0 with Data Shared lasso

⇒ The symmetric formulation (glmnet) encourages homogeneity

Symmetric formulation

	β_0	β_1	β_K
1					
⋮					
⋮					
p					

Ref = 0

Ref = 1

Standard formulation

with $\delta_k = \beta_k - \beta_0$

	δ_1	δ_K
1				
⋮				
⋮				
p				

$$C(0) = p(K-1)$$

with $\delta_k = \beta_k - \beta_1$

	δ_0	δ_2	...	δ_K
1				
⋮				
⋮				
p				

$$C(1) = p$$

Decomposition targeted by

MultinomDataSharedRef

μ	γ_1	γ_K

~

$$C(0) = p + 0 = p$$

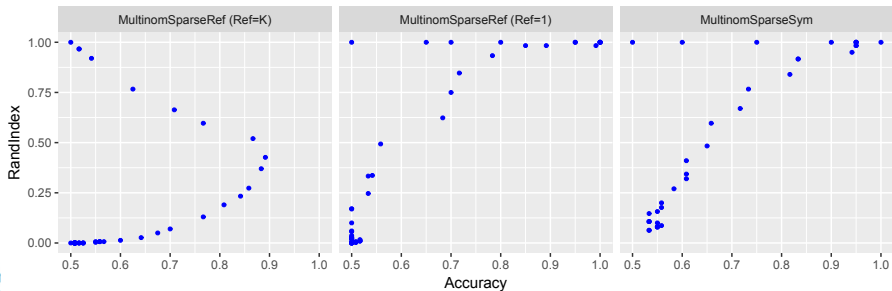
μ	γ_0	γ_2	...	γ_K

~

$$C(1) = 0 + p = p$$

A simple example in the homogenous case: $\beta_0^* = 0, \beta_1^* = \dots = \beta_K^* \neq 0$

- We make λ vary and compute, for each value
 - Accuracy : accuracy in the identification of null/non-null entries in $B^* = (\beta_1^*, \dots, \beta_K^*)$
 - Rand Index : \sim accuracy in the identification of groups of disease subtypes for which any one given metabolite shows the same level of association.



Conclusion/Discussion

- Data Shared Lasso has some nice properties
 - easy to implement, under a variety of models (~ Lasso with interaction terms)
 - \approx Oracle RefLasso regarding theoretical properties
 - \approx A priori RefLasso regarding implementation
 - can partially describe the structure of the effects in situations where CliqueFused may fail to fully describe it.
- ⇒ Nice alternative/complement to “A priori RefLasso”
- Provides insights into the sparse multinomial regression models

CliqueFused

▶ Back

- Penalized criterion:

$$\frac{\|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2}{2n} + \lambda_1 \|\mathbf{b}\|_1 + \lambda_2 \sum_{j_1 \sim j_2} |b_{j_1} - b_{j_2}|$$

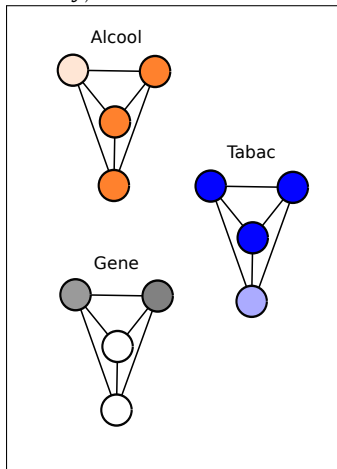
CliqueFused

▶ Back

- Penalized criterion:

$$\frac{\|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2}{2n} + \lambda_1 \|\mathbf{b}\|_1 + \lambda_2 \sum_{j_1 \sim j_2} |b_{j_1} - b_{j_2}|$$

$$\sum_{j, k > k'} |\beta_{k,j} - \beta_{k',j}|$$

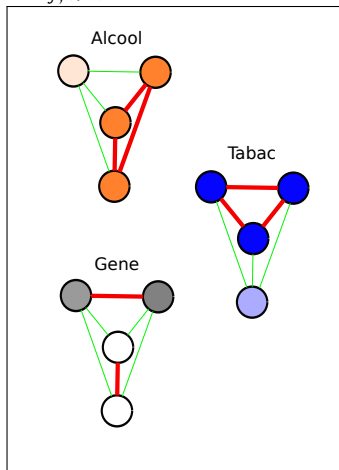


Adaptive Clique Fused (V. et al., Stat. Comp., 2016)

▶ Back

$$\sum_{j, k > k'} v_{k, k', j} |\beta_{k, j} - \beta_{k', j}|$$

$$v_{k, k', j} = |\tilde{\beta}_{k, j} - \tilde{\beta}_{k', j}|^{-1}$$

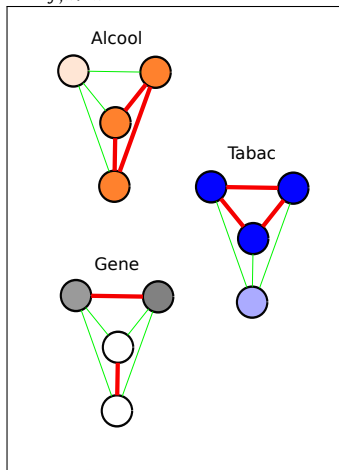


International Agency for Research on Cancer

Adaptive Clique Fused (V. et al., Stat. Comp., 2016)

▶ Back

$$\sum_{j,k>k'} v_{k,k',j} |\beta_{k,j} - \beta_{k',j}|$$



$$v_{k,k',j} = |\tilde{\beta}_{k,j} - \tilde{\beta}_{k',j}|^{-1}$$

Corollaire 1

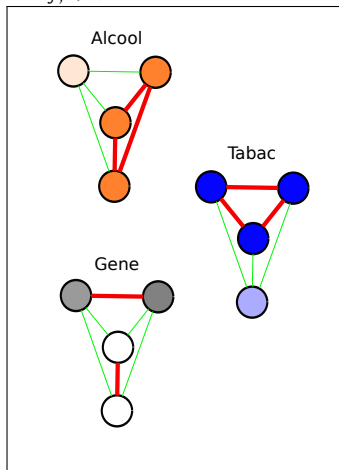
If Kp is fixed, and under some regularity conditions,

- ① $\mathbb{P}(\text{structure identified, for all } j) \rightarrow 1;$
- ② $\sqrt{n}(\hat{\mathbf{b}}_{\mathcal{A}}^{ad} - \mathbf{b}_{\mathcal{A}}^*) \rightarrow_d \mathcal{N}(\mathbf{0}_{s_0}, \sigma^2 \tilde{\mathbf{C}}_{\mathcal{A}}^{-1}).$

Adaptive Clique Fused (V. et al., Stat. Comp., 2016)

▶ Back

$$\sum_{j, k > k'} v_{k, k', j} |\beta_{k, j} - \beta_{k', j}|$$



$$v_{k, k', j} = |\tilde{\beta}_{k, j} - \tilde{\beta}_{k', j}|^{-1}$$

Corollaire 1

If Kp is fixed, and under some regularity conditions,

- ① $\mathbb{P}(\text{structure identified, for all } j) \rightarrow 1;$
- ② $\sqrt{n}(\hat{\mathbf{b}}_{\mathcal{A}}^{ad} - \mathbf{b}_{\mathcal{A}}^*) \rightarrow_d \mathcal{N}(\mathbf{0}_{s_0}, \sigma^2 \tilde{\mathbf{C}}_{\mathcal{A}}^{-1}).$

- If $n_k \gg p$, Adaptive Clique Fused ok

Pénalités classiques en apprentissage multi-tâche

- (sparse) GroupLasso : soit $\beta^j = (\beta_{1,j}, \dots, \beta_{K,j})$, et

$$\text{pen}(\beta_1, \dots, \beta_K) = \sum_j \|\beta^{(j)}\|_2$$

- Nuclear(Trace)-Norm (\sim reduced-rank regression). Soit \mathbf{B} la matrice de dimension (p, K) avec $(\mathbf{B})_{k,j} = \beta_{k,j}$, et

$$\text{pen}(\beta_1, \dots, \beta_K) = \text{tr}(\sqrt{\mathbf{B}^T \mathbf{B}}).$$

Pénalités classiques en apprentissage multi-tâche

- (sparse) GroupLasso : soit $\beta^j = (\beta_{1,j}, \dots, \beta_{K,j})$, et

$$\text{pen}(\beta_1, \dots, \beta_K) = \sum_j \|\beta^{(j)}\|_2$$

- Nuclear(Trace)-Norm (\sim reduced-rank regression). Soit \mathbf{B} la matrice de dimension (p, K) avec $(\mathbf{B})_{k,j} = \beta_{k,j}$, et

$$\text{pen}(\beta_1, \dots, \beta_K) = \text{tr}(\sqrt{\mathbf{B}^T \mathbf{B}}).$$

- ⇒ Ok pour prédiction, mais pas adaptées à notre objectif d'identification des partitions (structure des effets) pour chaque covariable

- Soit $\tilde{\mathbf{X}}_r^{(k)} = \mathbf{X}_{P_r^{(k)}}^{(k)}$ où $P_r^{(k)} = \{j \in [p] : k \neq r_j\}$
- Soit alors la matrice de taille $n \times Kp$:

$$\mathbf{x}_r = \begin{pmatrix} \mathbf{X}^{(1)} & \tilde{\mathbf{X}}_\ell^{(1)}/\tau_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & 0 & \dots & \tilde{\mathbf{X}}_\ell^{(K)}/\tau_K \end{pmatrix}$$

- Soit $\tilde{\mathbf{X}}_r^{(k)} = \mathbf{X}_{P_r^{(k)}}^{(k)}$ où $P_r^{(k)} = \{j \in [p] : k \neq r_j\}$
- Soit alors la matrice de taille $n \times Kp$:

$$\mathbf{x}_r = \begin{pmatrix} \mathbf{X}^{(1)} & \tilde{\mathbf{X}}_\ell^{(1)}/\tau_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & 0 & \dots & \tilde{\mathbf{X}}_\ell^{(K)}/\tau_K \end{pmatrix}$$

de manière équivalente, on peut minimiser en $\boldsymbol{\theta}_r \in \mathbb{R}^{Kp}$,

$$\frac{1}{2n} \|\mathbf{y} - \mathbf{x}_r \boldsymbol{\theta}_r\|_2^2 + \lambda_1 \|\boldsymbol{\theta}_r\|_1$$

qui retourne un estimateur de

$$\boldsymbol{\theta}_r^* = (\boldsymbol{\mu}_r^{*T}, \tau_1 \boldsymbol{\gamma}_{r,1}^{*T}, \dots, \tau_K \boldsymbol{\gamma}_{r,K}^{*T})^T \in \mathbb{R}^{Kp}.$$

avec $\boldsymbol{\gamma}_{r,k}^{*T} = (\boldsymbol{\beta}_k^* \boldsymbol{\mu}_r^*)_{P_r^{(k)}}$