

Joint estimation of multiple sparse regression functions with simple ℓ_1 penalties.

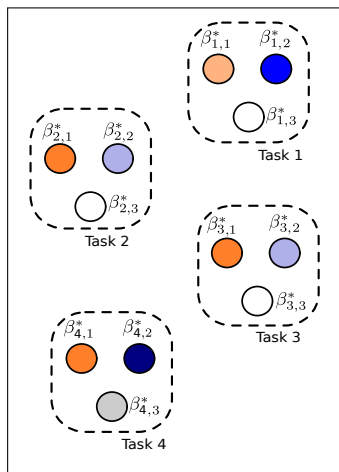
V. Viallon & E. Ollier

Univ. Lyon 1; IFSTTAR

Illustrative example

- Observations collected in several strata of a population
 - Cholesterol and gene expression levels in France, Italy, Spain and Germany for instance.
- “Naive” approaches
 - Estimate 4 independent models: does not account for the usual homogeneity: $\beta_{k,j} \simeq \beta_{k',j}$ for most genes j
 - Estimate 1 model, on all the data: masks the potential heterogeneity.

⇒ Need for dedicated approaches, that automatically adapts for the level of heterogeneity



Joint estimation of K linear regression models

- For each stratum k , data are supposed to follow the model

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

where

- $\mathbf{Y}^{(k)} = (Y_1^{(k)}, \dots, Y_{n_k}^{(k)})^T \in \mathbb{R}^{n_k}$
- $\boldsymbol{\xi}^{(k)} = (\xi_1^{(k)}, \dots, \xi_{n_k}^{(k)})^T \in \mathbb{R}^{n_k}$
- $\mathbf{X}^{(k)} = (\mathbf{x}_1^{(k)T}, \dots, \mathbf{x}_{n_k}^{(k)T})^T \in \mathbb{R}^{n_k \times p}$.

- Introduce $n = \sum_{k=1}^K n_k$

\Rightarrow K linear regression models on **fixed design**, with **gaussian and homoscedastic** errors (no intercept).

- Extensions : logistic regression, Cox model, etc.

Principle: decomposition of the β_k^* 's

- We can always write $\beta_k^* = \bar{\beta}^* + \gamma_k^*$
- $\bar{\beta}_j^*$ is the “common” effect of gene j over the strata:
 - $\bar{\beta}_j^* = \text{mode}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
 - $\bar{\beta}_j^* = \text{median}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
 - $\bar{\beta}_j^* = \text{mean}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$

Principle: decomposition of the β_k^* 's

- We can always write $\beta_k^* = \bar{\beta}^* + \gamma_k^*$
- $\bar{\beta}_j^*$ is the “common” effect of gene j over the strata:
 - $\bar{\beta}_j^* = \text{mode}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
 - $\bar{\beta}_j^* = \text{median}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
 - $\bar{\beta}_j^* = \text{mean}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
- Idea: Encourage solutions $(\hat{\beta}_1, \dots, \hat{\beta}_K)$ such that
 - $\hat{\beta}_j \approx \text{median}(\hat{\beta}_{1,j}, \dots, \hat{\beta}_{K,j})$
 - The common effect $\hat{\beta}$ is sparse;
 - Vectors $\hat{\beta}_k - \hat{\beta}$ are sparse: \Rightarrow similarity of the $\hat{\beta}_k$'s

Principle: decomposition of the β_k^* 's

- We can always write $\beta_k^* = \bar{\beta}^* + \gamma_k^*$
- $\bar{\beta}_j^*$ is the “common” effect of gene j over the strata:
 - $\bar{\beta}_j^* = \text{mode}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
 - $\bar{\beta}_j^* = \text{median}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
 - $\bar{\beta}_j^* = \text{mean}(\beta_{1,j}^*, \dots, \beta_{K,j}^*)$
- Idea: Encourage solutions $(\hat{\beta}_1, \dots, \hat{\beta}_K)$ such that
 - $\hat{\beta}_j \approx \text{median}(\hat{\beta}_{1,j}, \dots, \hat{\beta}_{K,j})$
 - The common effect $\hat{\beta}$ is sparse;
 - Vectors $\hat{\beta}_k - \hat{\beta}$ are sparse: \Rightarrow similarity of the $\hat{\beta}_k$'s
- M_1

$$\sum_{k \geq 1} \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}(\bar{\beta} + \gamma_k)\|_2^2}{2n} + \lambda_1 \|\bar{\beta}\|_1 + \sum_{k \geq 1} \lambda_{2,k} \|\gamma_k\|_1$$

Remarks

- M_1

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

- At optimum, $\widehat{\boldsymbol{\beta}}_j$: shrunk version of median($\widehat{\boldsymbol{\beta}}_{1,j}, \dots, \widehat{\boldsymbol{\beta}}_{K,j}$): penalization yields “identifiability”.
- If λ_1 is large enough, $\widehat{\boldsymbol{\beta}} = \mathbf{0}_p$ and $\widehat{\boldsymbol{\beta}}_k = \widehat{\boldsymbol{\gamma}}_k$.
- If the $\lambda_{2,k}$'s are large enough, $\widehat{\boldsymbol{\gamma}}_k = \mathbf{0}_p$ and $\widehat{\boldsymbol{\beta}}_k = \widehat{\boldsymbol{\beta}}$.

Adaptive version of M1 when $n_k \gg p$ and K is odd

- Denote by $\tilde{\beta}_k$ initial OLS estimates of the β_k 's.
- For all j , further define ℓ_j such that $\tilde{\beta}_{\ell_j, j} = \text{median}(\tilde{\beta}_{1, j}, \dots, \tilde{\beta}_{K, j})$.
- And set $\bar{\beta}_j^{\text{init}} = \tilde{\beta}_{\ell_j, j}$.

Adaptive version of M1 when $n_k \gg p$ and K is odd

- Denote by $\tilde{\beta}_k$ initial OLS estimates of the β_k 's.
- For all j , further define ℓ_j such that $\tilde{\beta}_{\ell_j, j} = \text{median}(\tilde{\beta}_{1, j}, \dots, \tilde{\beta}_{K, j})$.
- And set $\bar{\beta}_j^{\text{init}} = \tilde{\beta}_{\ell_j, j}$.
- Then setting $w_{k, j} = 1/|\tilde{\beta}_{k, j}|$ and $v_{k, j} = 1/|\tilde{\beta}_{k, j} - \bar{\beta}_j^{\text{init}}|$ consider the problem

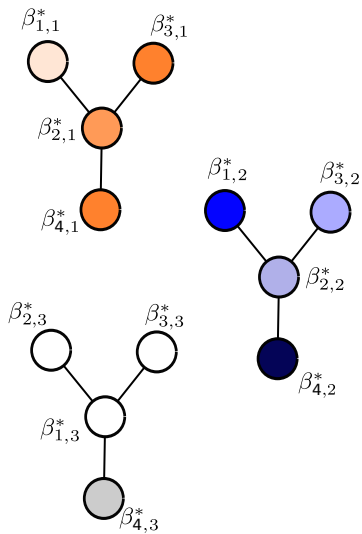
$$\sum_{k \geq 1} \left\{ \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)} \beta_k\|_2^2}{2n} + \lambda_1 \sum_j w_{\ell_j, j} |\bar{\beta}_j| + \sum_{k \geq 1} \lambda_{2, k} \sum_j v_{k, j} |\gamma_{k, j}| \right\}.$$

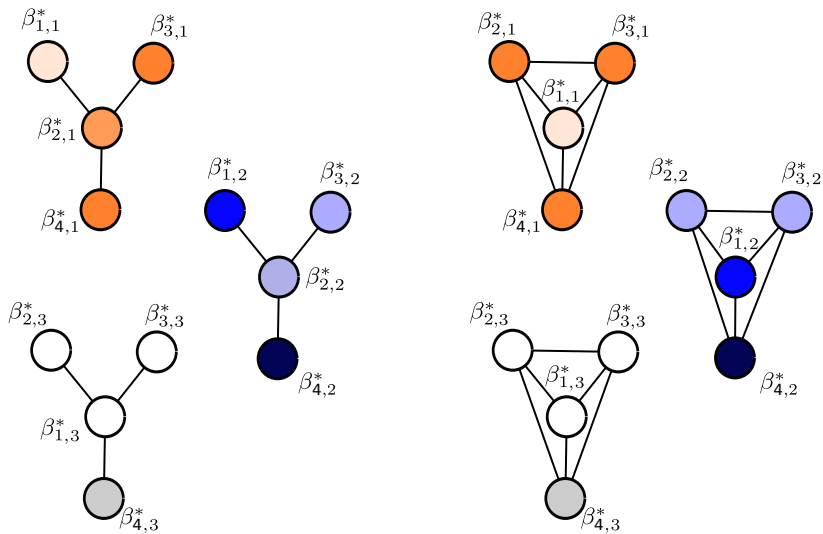
- This problem can be rewritten as

$$\sum_{k \geq 1} \left\{ \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)} \beta_k\|_2^2}{2n} + \lambda_1 \sum_j w_{\ell_j, j} |\beta_{\ell_j, j}| + \sum_{k \geq 1} \lambda_{2, k} \sum_j v_{k, j} |\beta_{k, j} - \beta_{\ell_j, j}| \right\}.$$

Connection with the “reference stratum” strategy:

- This corresponds to the decomposition: $\beta_{k,j} = \beta_{\ell_j,j} + \gamma_{k,j}$, for $k \neq \ell_j$.
- Can be seen as a refinement of the following “standard” approach:
 - Select a reference stratum ℓ
 - Write $\beta_k = \beta_\ell + \delta_k$
 - Select the non-zero components in β_ℓ and δ_k (e.g., with the lasso).
- Advantages of our approach:
 - the reference ℓ_j stratum is automatically selected from the initial OLS estimates;
 - it does not have to be the same for all gene j .

M_1 vs GenFused

M_1 vs GenFused

Implementation of M_1 : weighted lasso on “augmented” data.

- Set

$$\mathbf{y} = \begin{pmatrix} \mathbf{Y}^{(1)} \\ \vdots \\ \mathbf{Y}^{(K)} \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} \mathbf{X}^{(1)} & \mathbf{X}^{(1)} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & \mathbf{0} & \dots & \mathbf{X}^{(K)} \end{pmatrix} \quad \boldsymbol{\theta} = \begin{pmatrix} \bar{\boldsymbol{\beta}} \\ \gamma_1 \\ \vdots \\ \gamma_K \end{pmatrix}$$

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

Implementation of M_1 : weighted lasso on “augmented” data.

- Set

$$\mathbf{y} = \begin{pmatrix} \mathbf{Y}^{(1)} \\ \vdots \\ \mathbf{Y}^{(K)} \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} \mathbf{X}^{(1)} & \mathbf{X}^{(1)} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}^{(K)} & \mathbf{0} & \dots & \mathbf{X}^{(K)} \end{pmatrix} \quad \boldsymbol{\theta} = \begin{pmatrix} \bar{\boldsymbol{\beta}} \\ \gamma_1 \\ \vdots \\ \gamma_K \end{pmatrix}$$

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

- Setting $\mathbf{v} = (1, \dots, \lambda_{2,1}/\lambda_1, \dots, \lambda_{2,K}/\lambda_1)$, the objective function minimized under M_1 can be written as

$$\frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2}{2n} + \lambda_1 \|\boldsymbol{\theta}\|_{\mathbf{v}(1)},$$

where $\|\boldsymbol{\theta}\|_{\mathbf{v}(1)} = \sum_{j=1}^{(K+1)p} v_j |\theta_j| : \Rightarrow$ **Weighted Lasso**

- Asymptotic oracle properties of the adaptive versions
 - p and K are held fixed, while $n \rightarrow \infty$
 - General assumptions:
 - n_k/n_κ , with $0 < \kappa < 1$
 - $(\mathbf{X}^{(k)T} \mathbf{X}^{(k)})/n_k \rightarrow \mathbf{C}^{(k)}$ with $\lambda_{\min}(\mathbf{C}^{(k)}) > 0$
 - ...
 - Easy to derive from results obtained for GenFused (*e.g.*, V. et al., 2014)

- Non-asymptotic properties of M_1
 - Generally, \mathcal{X} does not enjoy RIP, RE, or the Irrepresentability condition.
 - In some cases though, it does !
 - If the β_k^* 's are all equal
 - If $K > 49\alpha^2 s^2$, where $\alpha > 1$ and $s \geq |\text{supp}(\theta^*)|$.

The Dirty Model of Jalali et al., 2013

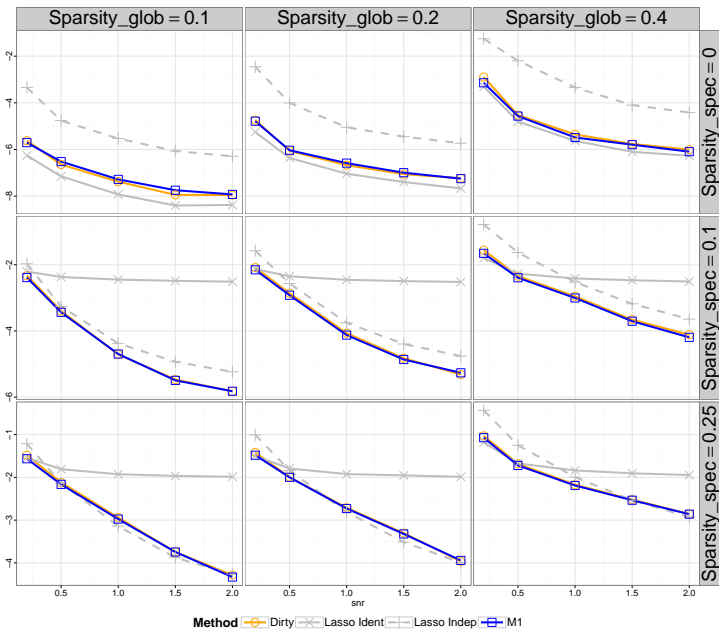
- Jalali et al. consider the following objective function

$$\sum_k \left\{ \frac{\|\mathbf{Y}^{(k)} - \mathbf{X}^{(k)}(\mathbf{R}^{(k)} + \mathbf{S}^{(k)})\|_2^2}{2n_k} + \lambda_S \|\mathbf{S}^{(k)}\|_1 \right\} + \lambda_R \sum_{j=1}^p \|R_j\|_\infty$$

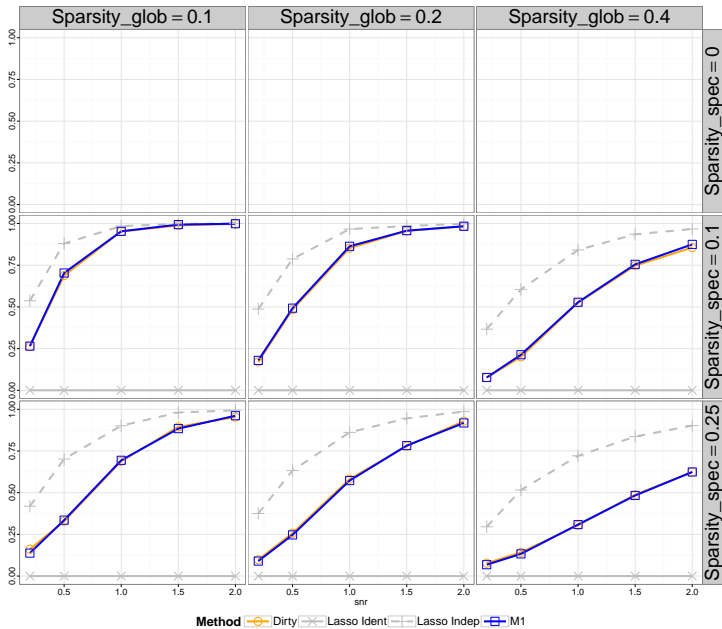
avec $R_j = (R_j^{(1)}, \dots, R_j^{(K)})^T$.

- The optimal solution $\widehat{\mathbf{B}} = \widehat{\mathbf{R}} + \widehat{\mathbf{S}}$ is then the sum of a row-sparse and a sparse matrix.
- Our approaches can be seen as simplified versions of Dirty: in particular, M_1
 - is much faster to run for linear regression models (many fast functions for the lasso)
 - is ready to use under many other regression models:
 - GLMs (logistic, Poisson, multinomial, ...), Cox models
 - conditional logistic regression, Cox models with competing risks, ...
- When $K = 2$, and if $\beta_{1j}^* \beta_{2j}^* \geq 0$, $M_1 \geq$ Dirty for support recovery!

Log-Prediction Error ($p = 15, K = 5, n_k = 225$)



Heterogeneity detection ($p = 15, K = 5, n_k = 225$)



Conclusion/ Discussion

- Our approach has deep connections with Ref^{ce}. Stratum, GenFused, and Dirty:
 - Refined version of the Ref^{ce}. Stratum strategy: the selection of the reference stratum is automatic, and “gene”-dependent.
 - Simplified version of the Dirty; similar empirical performance
 - Ready to use under a variety of models: linear, logistic, Cox, conditional logistic, ...
- Perspectives:
 - Extend our non-asymptotic results (lasso on correlated designs, Dalalyan et al. 2014).
 - Apply the approach on real data sets: Breast Cancer (IARC, WHO), Drivers Responsibility (IFSTTAR), ...
 - Extend the method to graphical models, ...