

VariSel: An R package to perform variable selection in the multivariate linear model

Or how a Gallic village stays irreducible

Marie Perrot-Dockès, Julien Chiquet

UseR2019

A simple vision of the immune system :

Or how Astérix and Obélix can kick-off the Romans

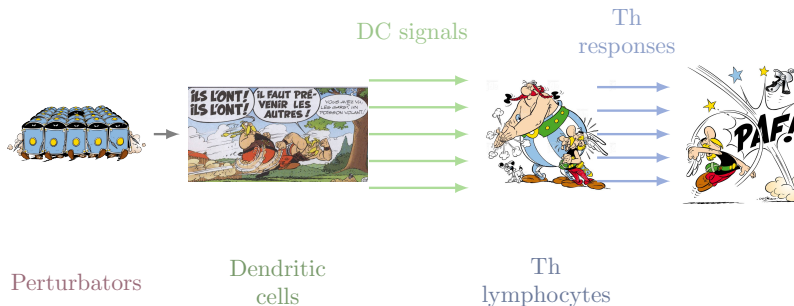


Figure 1: DC Th dialogue

Grandclaudon, M., Perrot-Dockès, M, Trichot, C et al. *A Quantitative Multivariate Model of Human Dendritic Cell-T Helper Cell Communication (March 15, 2019)*. Available at <http://dx.doi.org/10.2139/ssrn.3353217>

Experimental set up :

Ordralphabetix and Cétautomatix are two!

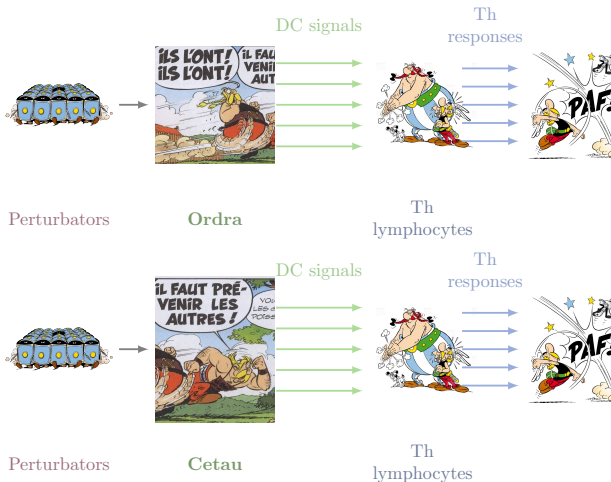


Figure 2: DC Th dialogue

► Dataset description:

- X : $n \times p$ design matrix : the DC signals



- Y : $n \times q$ response matrix : the Th responses



- **Question:** Which variables influence the responses?

- **Approach:**

- Variable selection in

$$Y = XB + E,$$

where

- B : $p \times q$ **sparse** coefficients matrix
- E : $n \times q$ error matrix with

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

- We take the dependence into account by estimating Σ_q .

► Dataset description:

- \mathbf{X} : $n \times p$ design matrix : the DC signals



- \mathbf{Y} : $n \times q$ response matrix : the Th responses



- **Question:** Which variables influence the responses?

- **Approach:**

- Variable selection in

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

where

- \mathbf{B} : $p \times q$ **sparse** coefficients matrix
- \mathbf{E} : $n \times q$ error matrix with

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

- We take the dependence into account by estimating Σ_q .

Different penalties : for different point of view

- ▶ **Lasso** : *select variables without taking into account potential links.*

$$\hat{b}_L = \text{Argmin}_b \left\{ \|y - \mathbf{X}b\|_2^2 + \lambda \|b\|_1 \right\},$$

- ▶ **Group-Lasso** : *select a group of variables.*

$$\hat{b}_G = \text{Argmin}_{b_1, \dots, b_L} \left\{ \|y - \sum_{1 \leq \ell \leq L} \mathbf{x}_{(\ell)} b_{(\ell)}\|_2^2 + \lambda \sum_{1 \leq \ell \leq L} \sqrt{p_\ell} \|b_\ell\|_2 \right\},$$

- ▶ **Fused-Lasso** : *influence a group of variables to have the same coefficient.*

$$\hat{b}_F = \text{Argmin}_b \|y - \mathbf{X}b\|_2^2 + \left\{ \lambda_1 \sum_{(i,j) \in \mathcal{G}} |b_i - b_j| + \lambda_2 \|b\|_1 \right\},$$

Different penalties : for different point of view

- ▶ **Lasso** : *select variables without taking into account potential links.*

$$\hat{b}_L = \text{Argmin}_b \left\{ \|y - \mathbf{X}b\|_2^2 + \lambda \|b\|_1 \right\},$$

- ▶ **Group-Lasso** : *select a group of variables.*

$$\hat{b}_G = \text{Argmin}_{b_1, \dots, b_L} \left\{ \|y - \sum_{1 \leq \ell \leq L} \mathbf{x}_{(\ell)} b_{(\ell)}\|_2^2 + \lambda \sum_{1 \leq \ell \leq L} \sqrt{p_\ell} \|b_\ell\|_2 \right\},$$

- ▶ **Fused-Lasso** : *influence a group of variables to have the same coefficient.*

$$\hat{b}_F = \text{Argmin}_b \|y - \mathbf{X}b\|_2^2 + \left\{ \lambda_1 \sum_{(i,j) \in \mathcal{G}} |b_i - b_j| + \lambda_2 \|b\|_1 \right\},$$

VariSel for one model type

```
mod <- train_VariSel( Y = T_resp,  
                     regressors = DC_sign,  
                     group = dc,  
                     type = "group_multi_regr")
```

```
X <- model.matrix(~DC_sign:dc -1)  
mod <- train_VariSel( Y = T_resp,  
                     X = X,  
                     sepx = ":",  
                     type = "group_multi_regr")
```

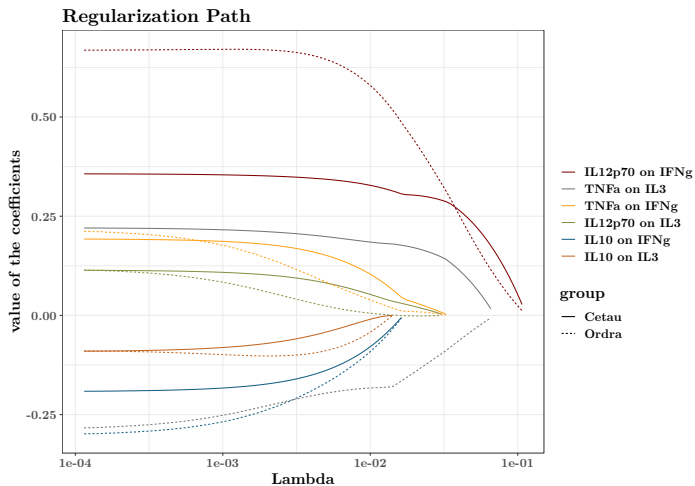

VariSel for one model type

```
mod <- train_VariSel( Y = T_resp,  
                     regressors = DC_sign,  
                     group = dc,  
                     type = "group_multi_regr")
```

```
X <- model.matrix(~DC_sign:dc -1)  
mod <- train_VariSel( Y = T_resp,  
                     X = X,  
                     sepx = ":",  
                     type = "group_multi_regr")
```

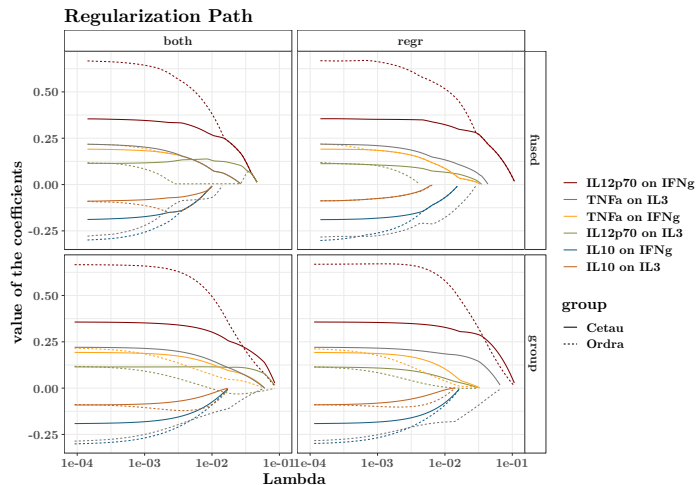
VariSel for one model type : Outcome

```
plot(mod)
```



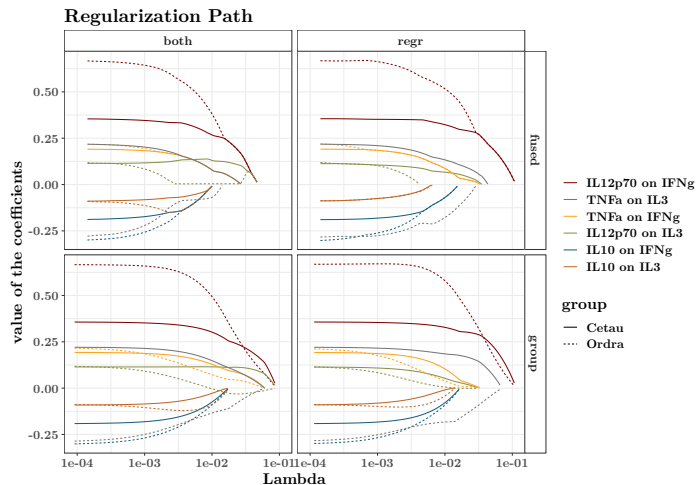
Different modelling strategy

```
compar_path(mods = list(mod,m2,m3,m4))
```



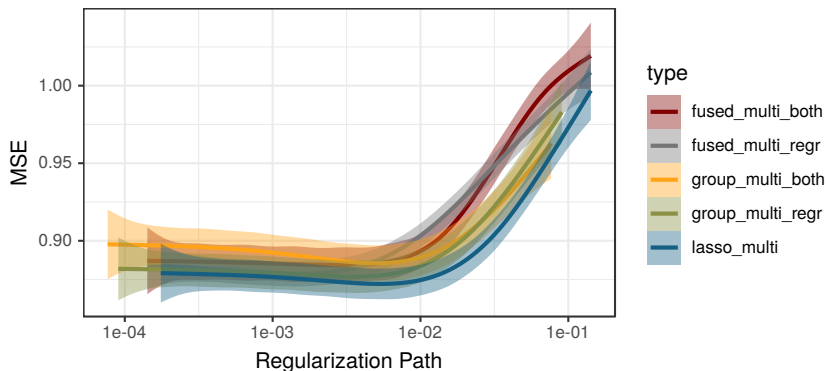
Different modelling strategy

```
compar_path(mods = list(mod,m2,m3,m4))
```



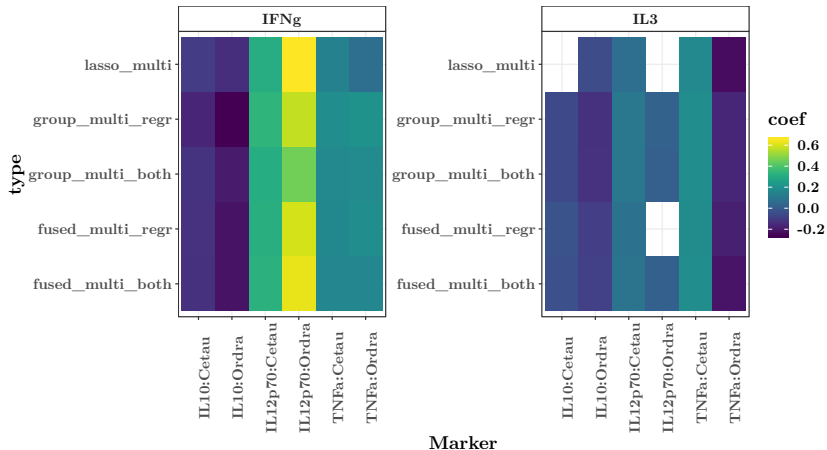
Models selection

```
ct <- compar_type( Y = T_resp, regressors = DC_sign,
                  group = dc,
                  types = c("group_multi_regr" , "group_multi_both" ,
                           "fused_multi_regr" , "fused_multi_both" ,
                           "lasso_multi" ), times = 10)
```



Best models representation

```
bm <- get_best_models(ct, criterion = "MSE_boot")  
plot_md(bm)
```



Conclusion

This is an R package to perform variable selection in multivariate linear models. It can

- ▶ Associate explicative variables
- ▶ Associate responses
- ▶ Associate both explicative variables and responses
- ▶ Let all variables 'free', without associating any of them

Come and see the vignette!

<https://github.com/Marie-PerrotDockes/VariSel>