# PLS-models in Practice: sparse and sparse group extensions

## Lecture 3/3

## MSc Health Data Analytics – Computational epidemiology – February 11, 2021
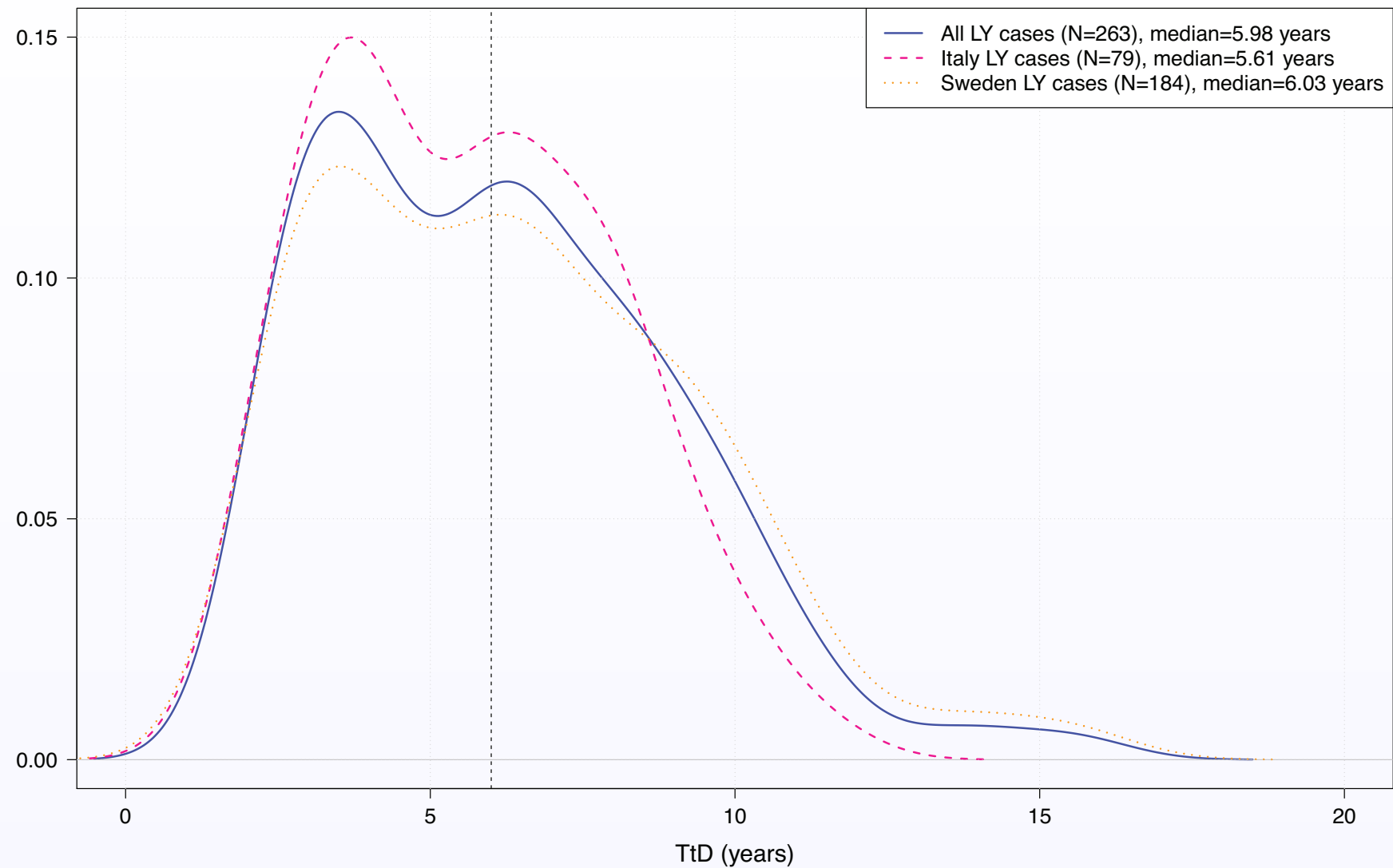
Marc Chadeau-Hyam

`m.chadeau@imperial.ac.uk`

**Imperial College**
**London**

# Lymphoma cases by subtypes and $TtD$

- EnviroGenoMarkers: a multi-OMIc study of NHL
  - Two contributing cohorts: EPIC Italy, and NHSDS
  - Transcriptomics, Proteomic (N=28) data available
- Four subtypes were identified:
  - B-cell Chronic Lymphatic Leukemia (BCLL): 14.8%
  - Diffuse Large B-cell Lymphoma (DLBCL): 15.6%
  - Follicular Lymphoma (FL): 14.4%
  - Multiple Myeloma (MM): 27.4%
- Study population:

| Subtype | $TtD<6$ | $TtD>6$ | **Total** |
|---------|---------|---------|-----------|
| BCLL | 15 | 24 | **39** |
| DLBCL | 18 | 23 | **41** |
| FL | 18 | 20 | **38** |
| MM | 42 | 30 | **72** |
| Others | 41 | 32 | **73** |
| **Total** | **93** | **97** | **263** |

# Time to Diagnosis ($TtD$) distribution in LY cases



$\Rightarrow TtD$=6 years is close to the median value

# Benchmark screening model

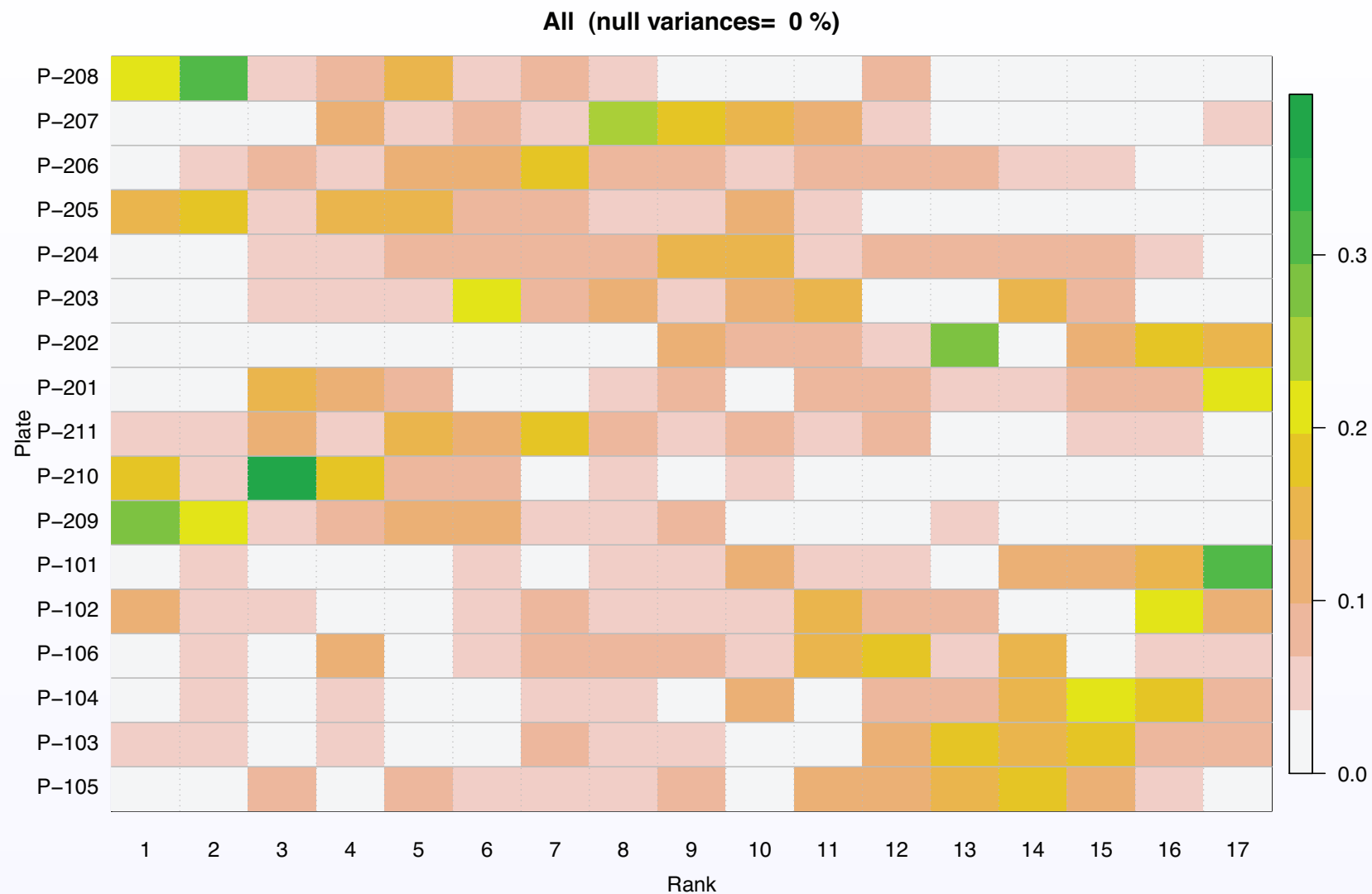- Univariate exploration of OMICs data accounting for nuisance variation

- Formulation, for individual $i$:

  - Variable of interest: $X^i$ (Ca/Co)

  - Predictors: $Y^i$, Expression levels

  - Fixed effects: $FE^i$

  - Random Effect variables: $u^{A^i}$, where $A^i$ are nuisance variables

$$Y^i \sim \alpha + \beta_1 X^i + \beta_2 FE^i + u^{A^i} + \epsilon^i$$
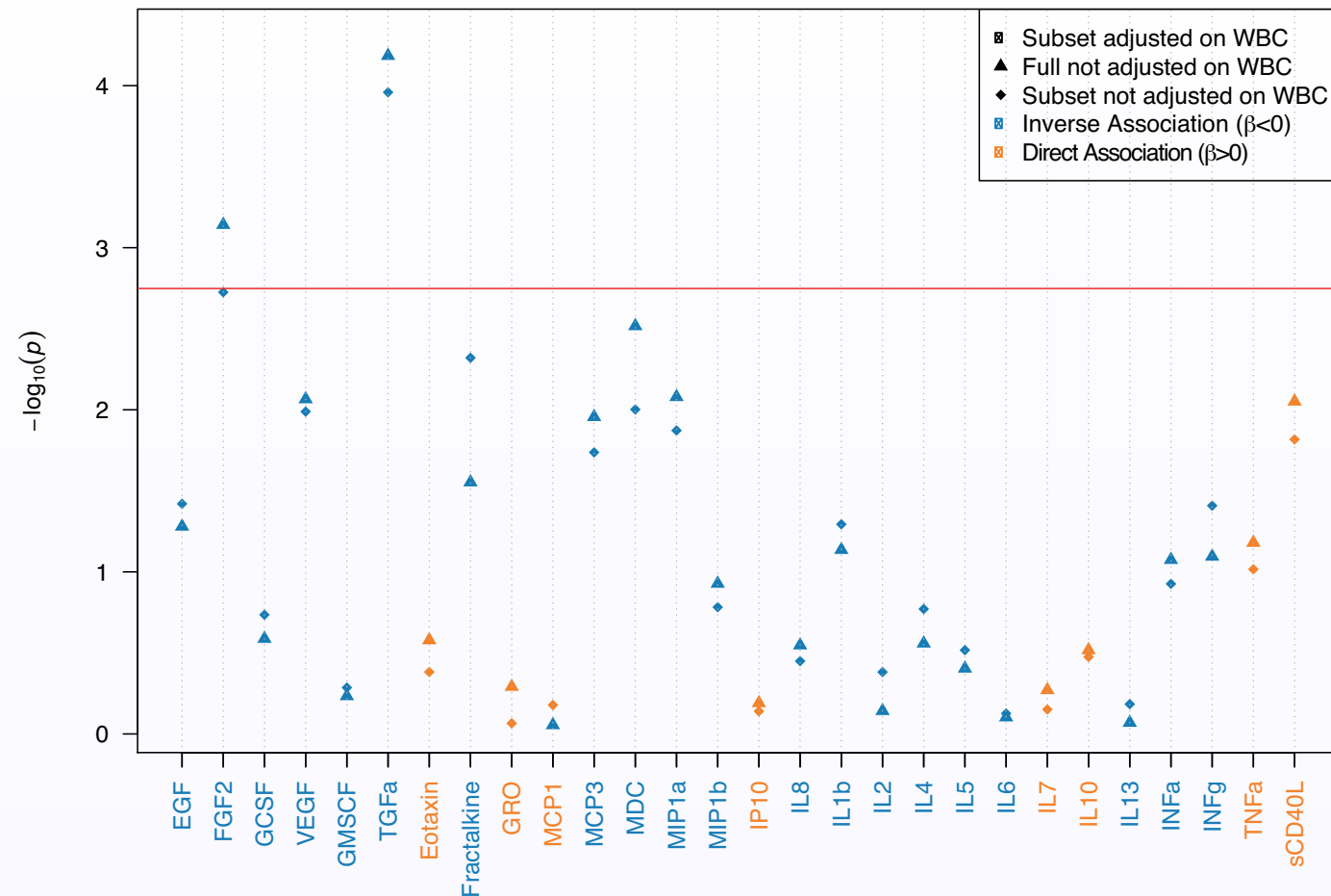
$\Rightarrow$ random intercept model

- Methodology: likelihood ratio test

  - Run the model with and without the variable of interest ($X^i$). Compare both models

    $\Rightarrow$ for each protein/probe we obtain a p-value testing the association between the probe and the disease status/or exposure

# Proteomics: Estimation of the random effects
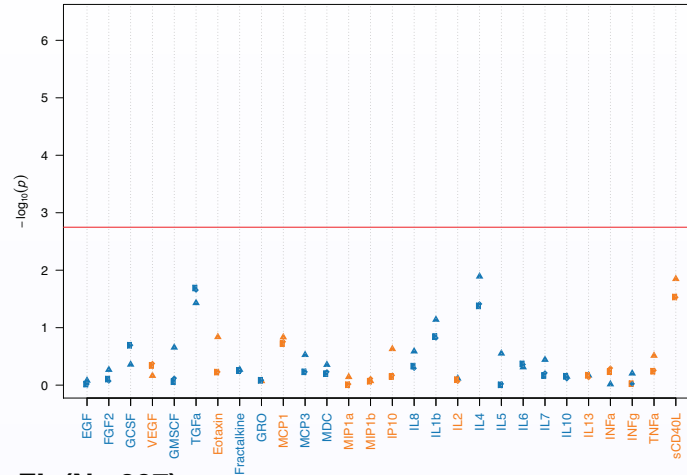


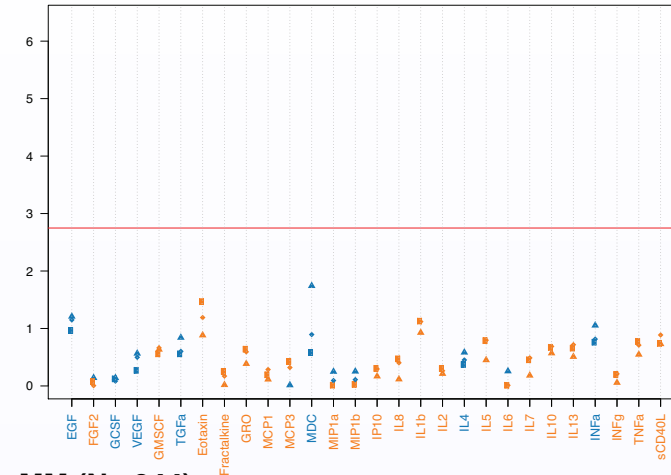⇒ P-208/209/210 lead to higher variances

# Analysis of all BCL cases



⇒ Two Bonferroni significant associations involving FGF2 & TGFα
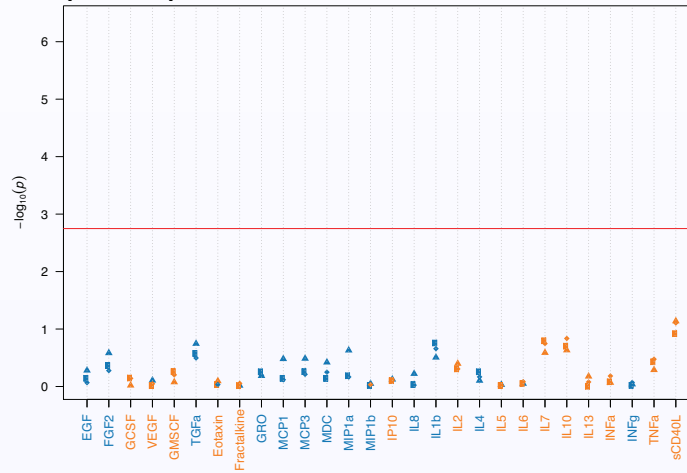⇒ weak effect of WBC adjustments

# Histological subtype analyses



CLL (N= 310)

DLBCL (N= 312)

FL (N= 307)

MM (N= 344)

$\Rightarrow$ 8 (strong) and inverse associations for MM

$\Rightarrow$ no association for the other subtypes

# All BCL excluding MM



$\Rightarrow$ Both BCL-related associations lose significance upon exclusion of MM cases

$\Rightarrow$ MM may have driven the BCL associations

# PLS analyses: Rationale and plan

- The 28 proteins can be classified in three functional groups
  - Growth Factors (N=6)
  - Chemokines (N=10)
  - Cytokines (N=12)
- Research questions
  - Do proteins jointly concur to BCL (and subtypes) onset?
  - Is the functional grouping relevant to the disease?
  - Are there groups (and proteins within each group) more associated to disease?

One Million $ question: $\Rightarrow$ How can we use PLS?

# PLS analyses: Rationale and plan

- The 28 proteins can be classified in three functional groups
  - Growth Factors (N=6)
  - Chemokines (N=10)
  - Cytokines (N=12)
- Additional versions of PLS:
  - Sparsity achieved through penalisation
  - Grouping signals a priori (e.g. pathways, genes)

► sparse PLS components (sPLS)

$$C^k = u_1 \times X_1 + \underbrace{u_2}_{=0} \times X_2 + \underbrace{u_3}_{=0} \times X_3 + \ldots + u_p \times X_p$$

► group PLS components (gPLS)

$$C^k = \overbrace{\underbrace{u_1}_{=0} X_1 + \underbrace{u_2}_{=0} X_2}^{module_1} + \overbrace{\underbrace{u_3}_{\neq 0} X_3 + \underbrace{u_4}_{\neq 0} X_1 + \underbrace{u_5}_{\neq 0} X_5}^{module_2} \ldots \overbrace{\underbrace{u_{p-1}}_{=0} X_{p-1} + \underbrace{u_p}_{=0} X_p}^{module_K}$$

► sparse group PLS components (sgPLS)

$$C^k = \overbrace{\underbrace{u_1}_{=0} X_1 + \underbrace{u_2}_{=0} X_2}^{module_1} + \overbrace{\underbrace{u_3}_{\neq 0} X_3 + \underbrace{u_4}_{=0} X_4 + \underbrace{u_5}_{=0} X_5}^{module_2} \ldots \overbrace{\underbrace{u_{p-1}}_{=0} X_{p-1} + \underbrace{u_p}_{=0} X_p}^{module_K}$$

# PLS analyses: Rationale and plan

- The 28 proteins can be classified in three functional groups
  - Growth Factors (N=6)
  - Chemokines (N=10)
  - Cytokines (N=12)
- Research questions
  - Do proteins jointly concur to BCL (and subtypes) onset?
  - Is the functional grouping relevant to the disease?
  - Are there groups (and proteins within each group) more associated to disease?

<span style="color:red">Two Million $ question: which models????</span>

# PLS analyses: Rationale and plan

- The 28 proteins can be classified in three functional groups
  - Growth Factors (N=6)
  - Chemokines (N=10)
  - Cytokines (N=12)
- Research questions
  - Do proteins jointly concur to BCL (and subtypes) onset? – (s)PLS
  - Is the functional grouping relevant to the disease? – gPLS
  - Are there groups (and proteins within each group) more associated to disease? – sgPLS

Two Million $ response

# PLS analyses: Rationale and plan

- The 28 proteins can be classified in three functional groups

  - Growth Factors (N=6)
  - Chemokines (N=10)
  - Cytokines (N=12)

- Research questions

  - Do proteins jointly concur to BCL (and subtypes) onset? – (s)PLS
  - Is the functional grouping relevant to the disease? – gPLS
  - Are there groups (and proteins within each group) more associated to disease? – sgPLS

- Analytical Plan: all PLS variants to analyse

  - All BCL
  - Each subtype separately
  - In cases only: the time to diagnosis

# gPLS: Penalty function and calibration

- For each component:

$$\min_{||u||=1, \; ||v||=1} \sum_{k=1}^{K} \sum_{l=1}^{L} \underbrace{||X^{(k)^T} Y^{(l)} - u^{(k)} v^{(l)^T}||_F^2}_{\text{covariances between } k^{th} \text{ and } l^{th} \text{ block}} + P_{\lambda_1}(u) + P_{\lambda_2}(v)$$

where

$$P_{\lambda_1}(u) = \lambda_1 \sum_{k=1}^{K} \sqrt{p_k} \; \underbrace{||u^{(k)}||_2}_{\substack{\text{loadings of} \\ k^{th} \text{ block in X}}} \qquad P_{\lambda_2}(v) = \lambda_2 \sum_{l=1}^{L} \sqrt{q_l} \; \underbrace{||v^{(l)}||_2}_{\substack{\text{loadings of} \\ l^{th} \text{ block in Y}}}$$

- Penalisation adapts to the number of variables in each group $(p_k, q_l)$
- Calibration: Number of selected groups in $X$ and $Y$ via cross-validation using MSEP

# sgPLS: Penalty function and calibration

- For each component:

$$\min_{||u||=1,\ ||v||=1} \sum_{k=1}^{K} \sum_{l=1}^{L} \underbrace{||X^{(k)^T} Y^{(l)} - u^{(k)} v^{(l)^T}||_F^2}_{\text{covariances between } k^{th} \text{ and } l^{th} \text{ block}} + P_{\lambda_1}(u) + P_{\lambda_2}(v)$$
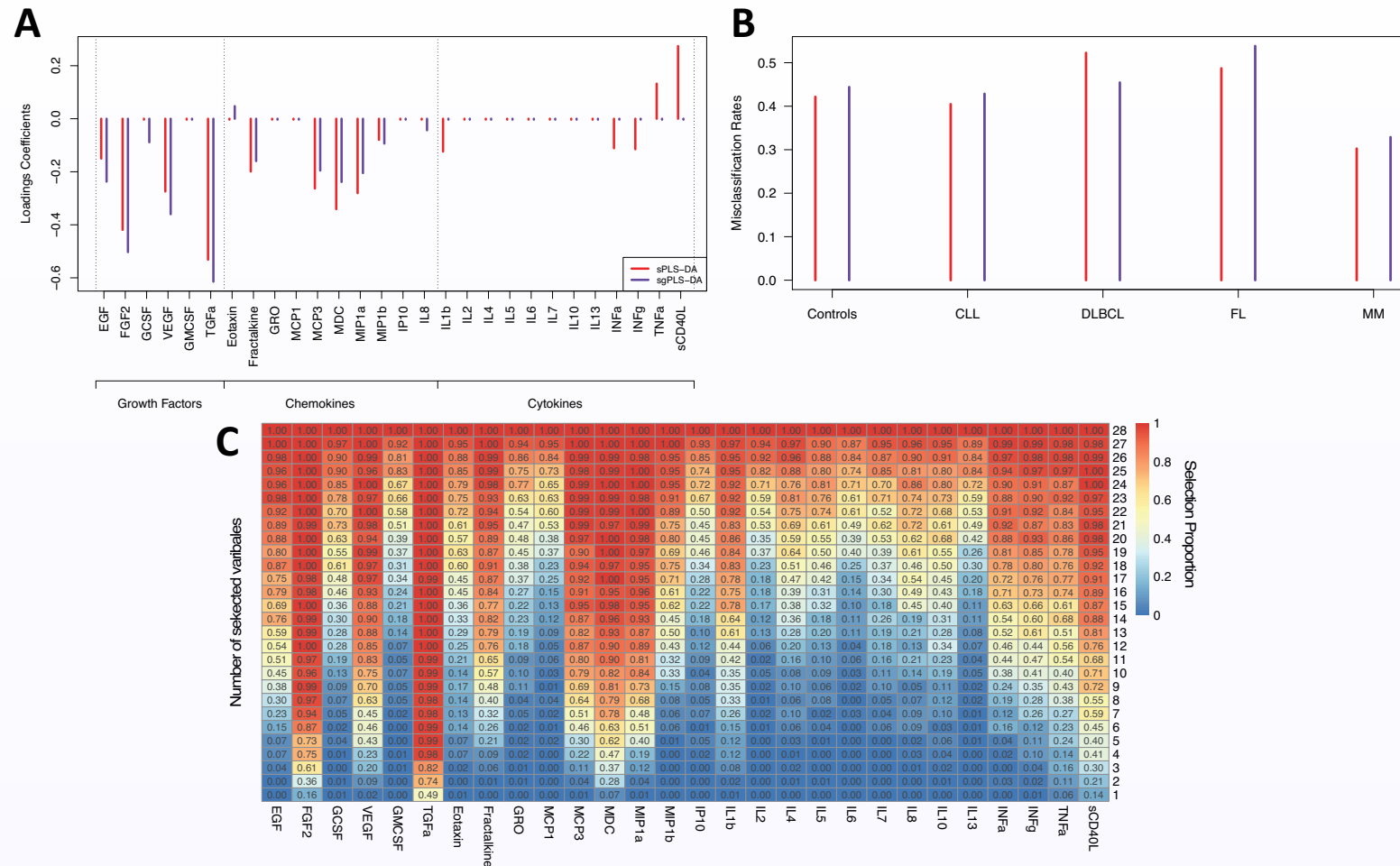
- Adding a LASSO penalty within each group:

$$P_{\lambda_1}(u) = \lambda_1 \sum_{k=1}^{K} \sqrt{p_k} ||u^{(k)}||_2 + \underbrace{\alpha_1 \lambda_1 ||u||_1}_{\substack{\text{sparsity} \\ \text{in X}}}$$

$$P_{\lambda_2}(v) = \lambda_2 \sum_{l=1}^{L} \sqrt{q_l} ||v^{(l)}||_2 + \underbrace{\alpha_2 \lambda_2 ||v||_1}_{\substack{\text{sparsity} \\ \text{in Y}}}$$

- Calibration: Number of selected groups in $X$ and $Y$ via cross-validation and the components sparsity parameter (not the number of variables)
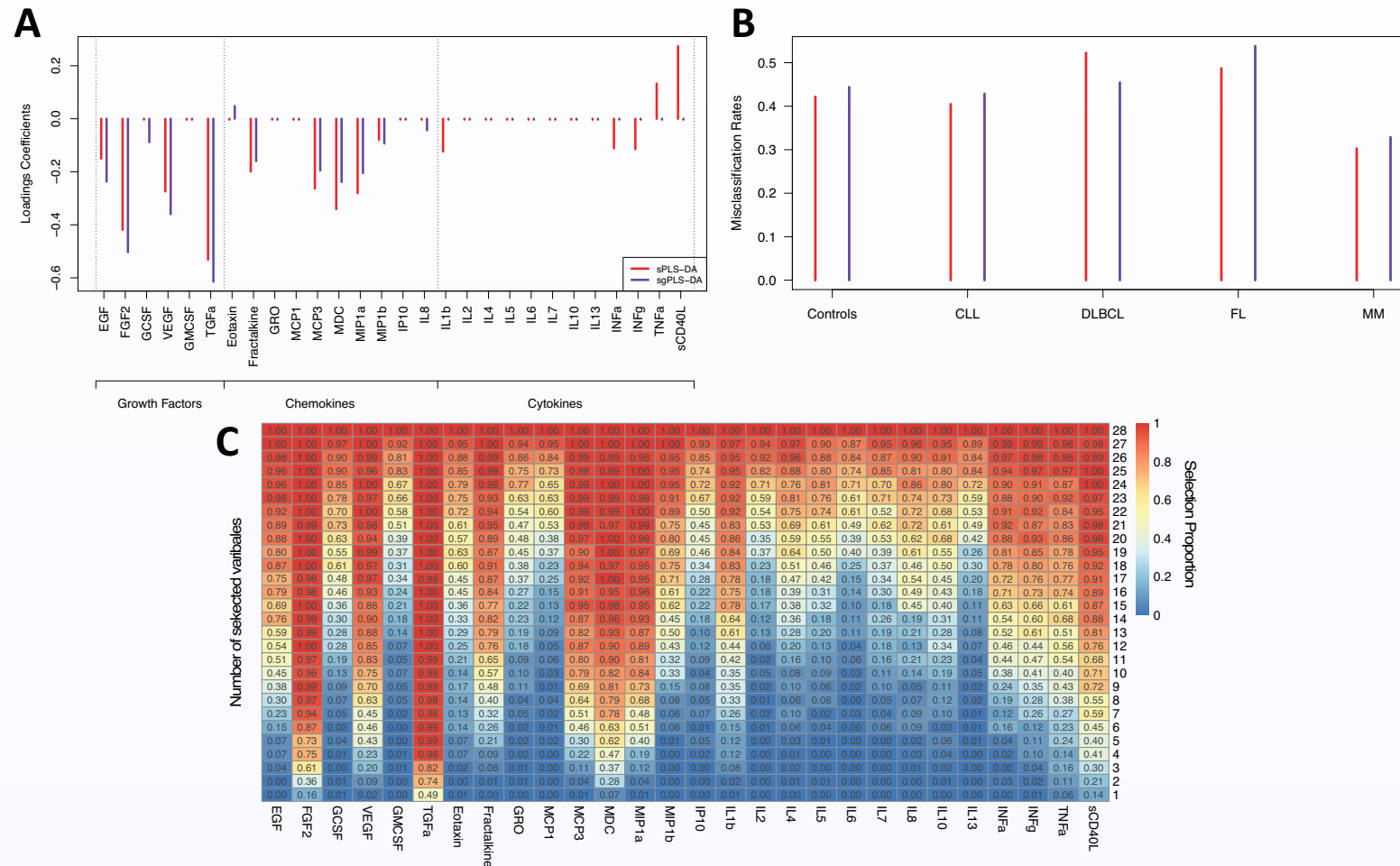
# PLS analyses: All BCL



- sPLS mainly selects variables in GF an chemokines groups
- Two cytokines proteins selected with larger loadings (TNF-$\alpha$, sCD40)
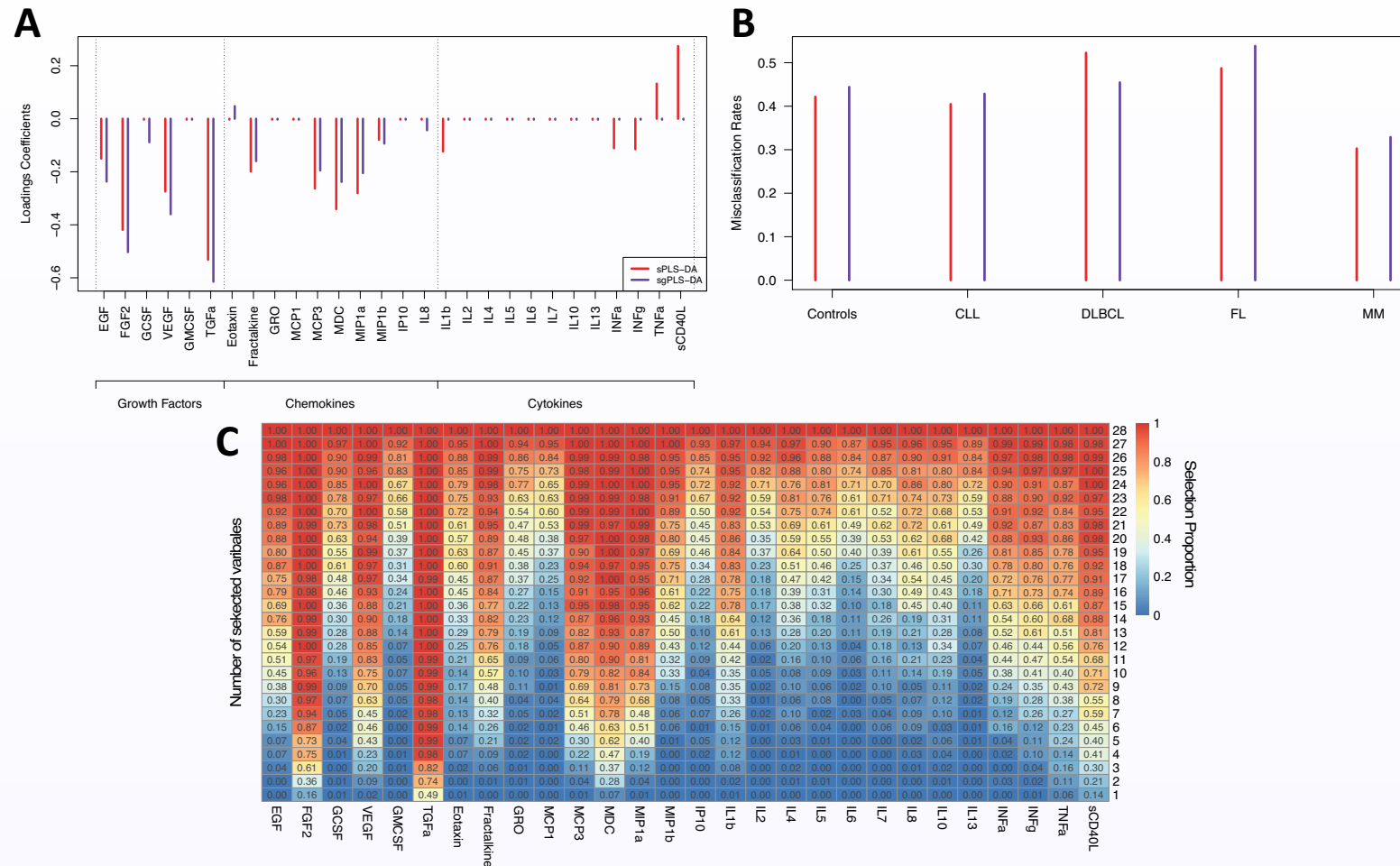- sgPLS selects the two group with more non zero loadings
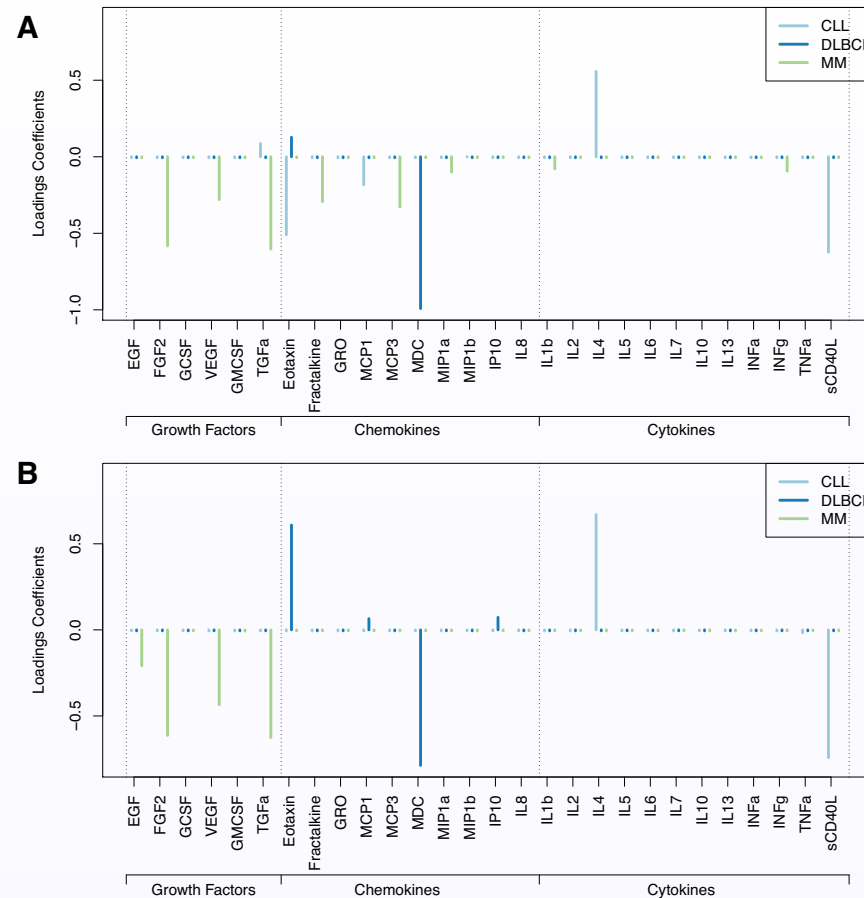
# PLS analyses: All BCL



- sPLS and sgPLS yield comparable misclassification rates (unimportant exclusion of cytokines)

- Better misclassification rates for MM

# PLS analyses: All BCL



- Assessing the sensitivity to calibration via stability analyses
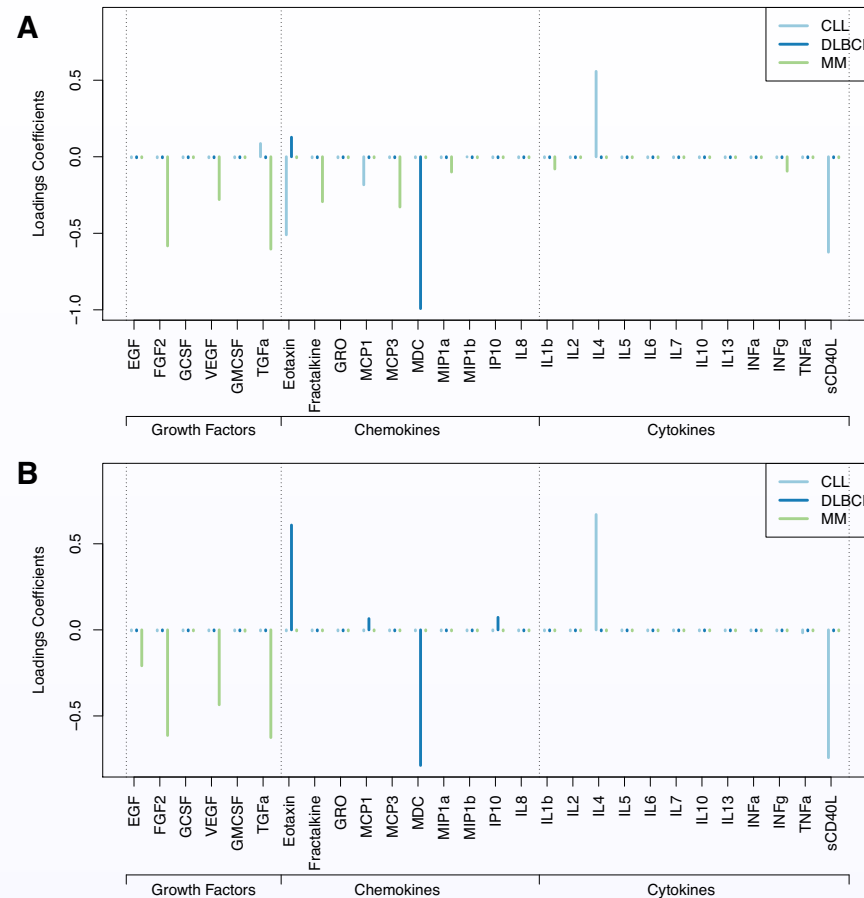- The largest loadings are the first and most frequently selected (sPLS)

# PLS analyses: subtype analyses



sPLS analyses select:

- **MM**: proteins mainly in chemokines and growth factors
- **CLL**: chemokines and cytokines (though only 2/12 proteins)
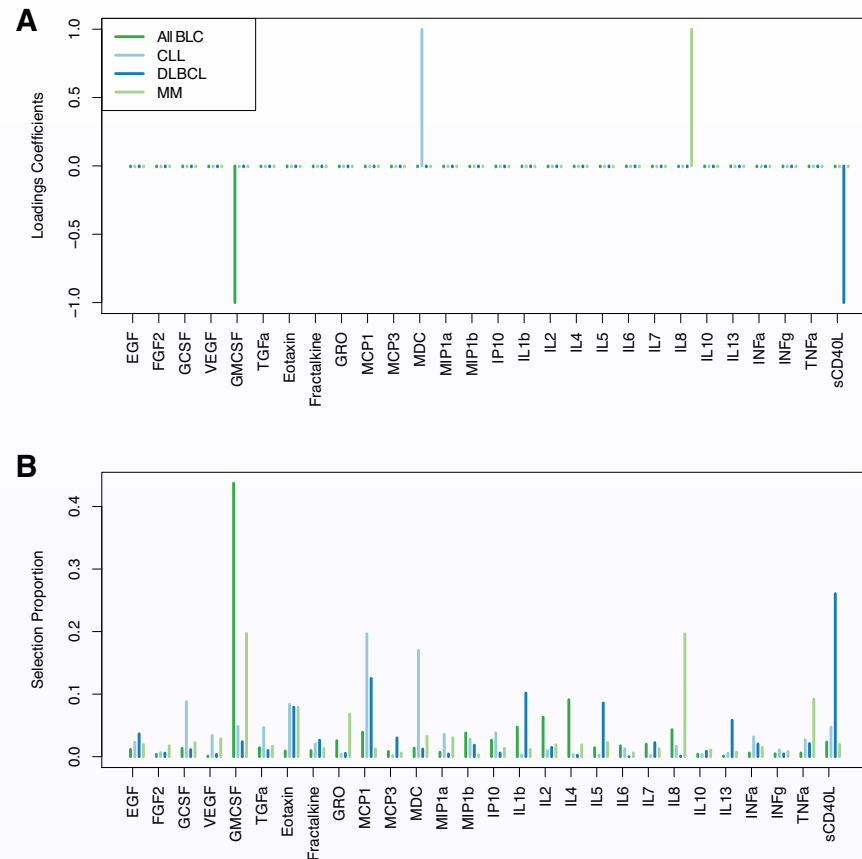- **DLBCL**: 2 chemokines are selected

# PLS analyses: subtype analyses



sgPLS analyses select:

- **MM**: growth factors and within the group the same variables as sPLS
- **CLL**: cytokines and both the sPLS proteins
- **DLBCL**: Chemokines are selected (including the the 2 sPLS proteins)

# Cases-only analyses



sPLS analyses of TtD (continuous outcome) selects:

- A single and specific protein for each subtype
- For all BCL, GMCSF is mostly selected, for other subtypes, several candidates compete
- For DLBCL, sCD40 seem to be more frequently selected.

# Wrap-up summary

- PLS analyses were able to identify associations the were not detected by univariate models

  ⇒ these potential markers were supported by external biological evidence

- Inclusion of groups allows to account for correlations across proteins and select the most informative sets of predictors

  ⇒ contribution to the sparsity and interpretability of the results

- Limitation: sensitivity to the grouping strategy

  ⇒ grouping is defining a prior hypothesis

- Extensions:
  - s-g-sg-PLS can accommodate large block of data (e.g. gene expression)
  - OMICs integration via sgsPLS
  - Computational Optimisation: bigPLS

# Acknowledgments

**Imperial College London**

Universiteit Utrecht

- M Chadeau-Hyam
- B Bodinier
- M Karimi

- R Vermeulen
- L Portengen
- F Saberi

- Collaborators: B Liquet (University of Pau)
- Financial Support:
  - EU FP7 EXPOsOMICS
  - CR-UK Population Research Committee project 'Mechanomics'