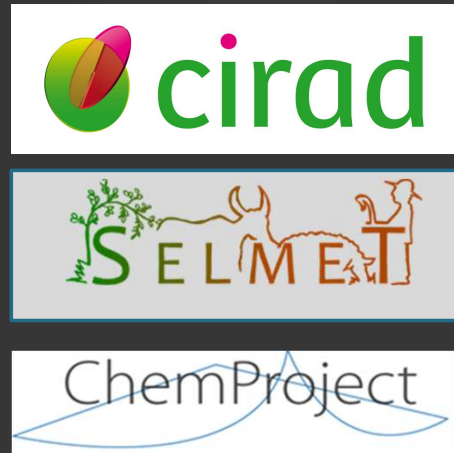# Sparse PLSR by regularized SVD

matthieu.lesnoff@cirad.fr
ChemHouse, Montpellier, 30 Sept 2025

*Statistical Applications in Genetics and Molecular Biology*

| Volume 7, Issue 1 | 2008 | Article 35 |
|---|---|---|

## A Sparse PLS for Variable Selection when Integrating Omics Data

**Kim-Anh Lê Cao,** *INRA UR 631 and Université de Toulouse*
**Debra Rossouw,** *University of Stellenbosch*
**Christèle Robert-Granié,** *INRA UR 631*
**Philippe Besse,** *Université de Toulouse*
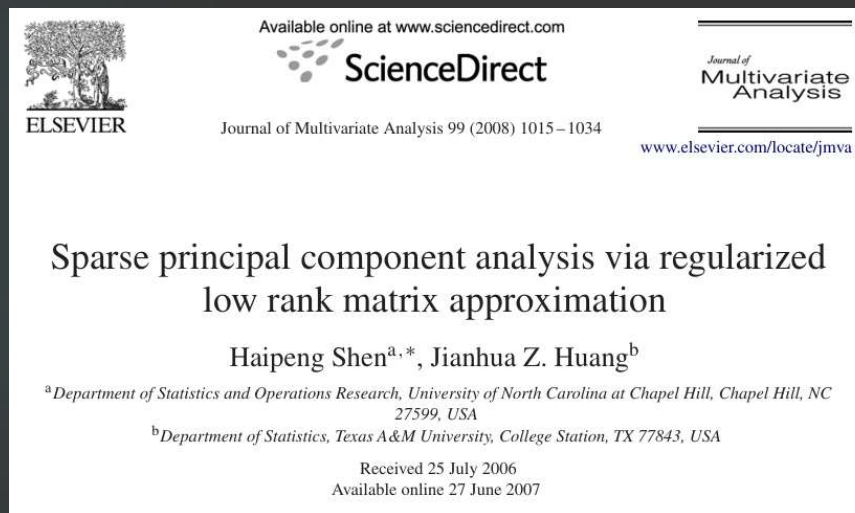
**R function `mixOmics::spls`**

Lê Cao, K.-A., Rohart, F., Gonzalez, I., Dejean, S., Abadi, A.J., Gautier, B., Bartolo, F., Monget, P., Coquery, J., Yao, F., Liquet, B., 2022. mixOmics: Omics Data Integration Project.
https://doi.org/10.18129/B9.bioc.mixOmics

# In brief about the method

- Use of a **regularized** (**instead of usual**)

SVD in the PLS algorithm

rSVD-sPCA

## Sparse principal component analysis via regularized low rank matrix approximation

Haipeng Shen[a,*], Jianhua Z. Huang[b]

[a] Department of Statistics and Operations Research, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA

[b] Department of Statistics, Texas A&M University, College Station, TX 77843, USA

**PLS   X (n, p), Y (n, q)**

- Scores   $(n, 1)$      $t_X = X\, w_X$      $t_Y = Y\, w_Y$

- Loading weights $w_X$ $(p, 1)$,  $w_Y$ $(q, 1)$

$$t_X = w_{X,1}\, x_1 + w_{X,2}\, x_2 + \ldots + w_{X,p}\, x_p$$

$$t_Y = w_{Y,1}\, y_1 + w_{Y,2}\, y_2 + \ldots + w_{Y,q}\, y_q$$

$w_X, w_Y$ such as

$$\max \mathrm{cov}^2(t_X, t_Y) \qquad \text{with} \quad \|w_X\| = \|w_Y\| = 1$$

# A usual PLS algorithm

Iterative NIPALS to get $w_X$ and $w_Y$:

1. Set $t_Y$

2. Repeat until convergence
    a) $w_X = X'\, t_Y$

    b) $w_X = w_X\, /\, norm(w_X)$

    c) $w_Y = Y'\, t_X\, /\, t_X{}'\, t_X$

    d) $w_Y = w_Y\, /\, norm(w_Y)$

    e) $t_Y = Y\, w_Y$

3. $t_X = X\, w_X$

**But $w_X$ and $w_Y$ can be computed without iteration:**

- Let $K = Y'X$  (could consider $X'Y$)

- SVD of $K \Rightarrow U\,\Delta V'$

- $w_X = V[:, 1]$

- $w_Y = U[:, 1]$

e.g., Höskuldsson 1988 Journal of Chemometrics
https://doi.org/10.1002/cem.1180020306

**(used in Jchemo: plskern, plsnipals, plscan)**

# A simple approach **to get sparse** $w_x$

Usual PLS

- $K = Y'X$

- SVD of $K \Rightarrow U \Delta V'$

- $w_X = V[:, 1]$

Sparse PLS

- $K = Y'X$

- **rSVD** (S&H 2008) of $K \Rightarrow U \Delta V'$

- $w_X = V[:, 1]$

This is the idea used by Lê Cao et al.

(authors also compute sparse $w_Y$, not detailed in this presentation)

## Summary of the rSVD of Shen & Huang 2008

$$\widehat{X} = u \, \delta \, v' \qquad \|u\| = \|v\| = 1 \qquad \text{1-rank SVD of a given matrix } X$$

A usual iterative NIPALS algorithm

1. Set $u$

2. Repeat until convergence
   a) $v = X' \, u$
   b) $u = X \, v \, / \, \mathrm{norm}(X \, v)$

3. $v = v \, / \, \mathrm{norm}(v)$

Regularization of Shen & Huang

1. Set $u$

2. Repeat until convergence
   a) $v = f_{\lambda}(X' \, u)$      soft thresholding
   b) $u = X \, v \, / \, \mathrm{norm}(X \, v)$

3. $v = v \, / \, \mathrm{norm}(v)$

**Soft-thresholding**

Tuning parameter $\lambda$

$$\text{sign}(v_j)(|v_j| - \lambda)_+$$

$f(v_j)$ (y-axis)

$v_j$ (non-normed v)

$w_j$ for PLS

# Illustration:    MLNIR dataset

**All the analyses in this presentation:** done with **Jchemo.jl**
https://github.com/mlesnoff/Jchemo.jl

In particular,

function **splsr**

— Sames results as **mixOmics::spls** but

   a) Sparsity only on $w_x$ (predictive approach)

   b) Nipals PLS algorithm replaced by the **"improved kernel algorithm #1"** of Dayal & McGegor 1997 $\Rightarrow$ faster

MLNIR dataset available at: https://zenodo.org/records/16783068

N = 208 hydrocarbon samples. For each sample, a near-infrared spectrum (intensities measured at given wavelenghts in cm-1) and a density value are provided. The dataset contains:
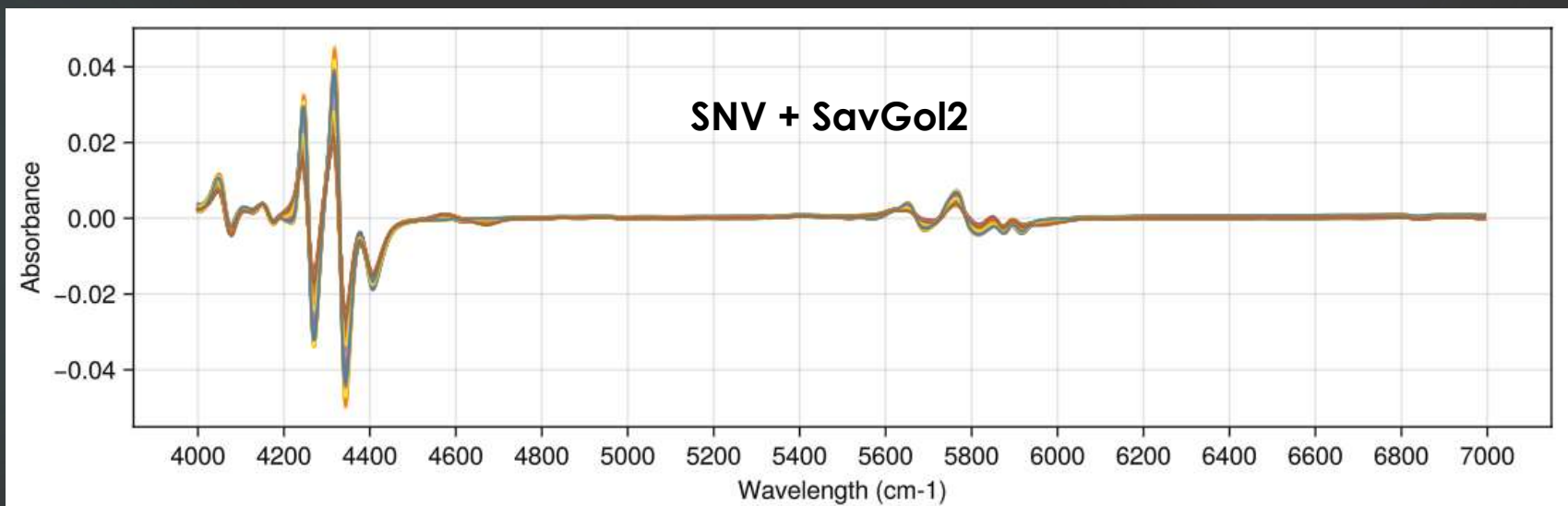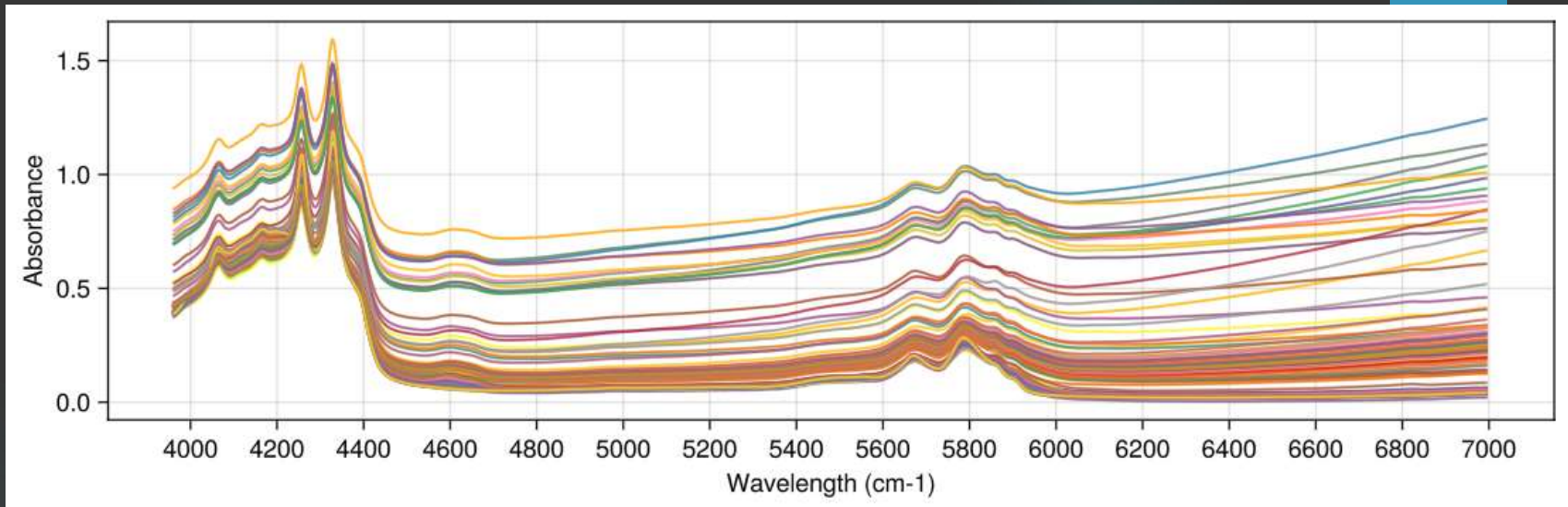
- X: Raw NIR data
- Xp: NIR data after preprocessing (SNV + 2nd derivative Savitzky-Golay)
- y: Density normalized to [0, 1] (response to predict)
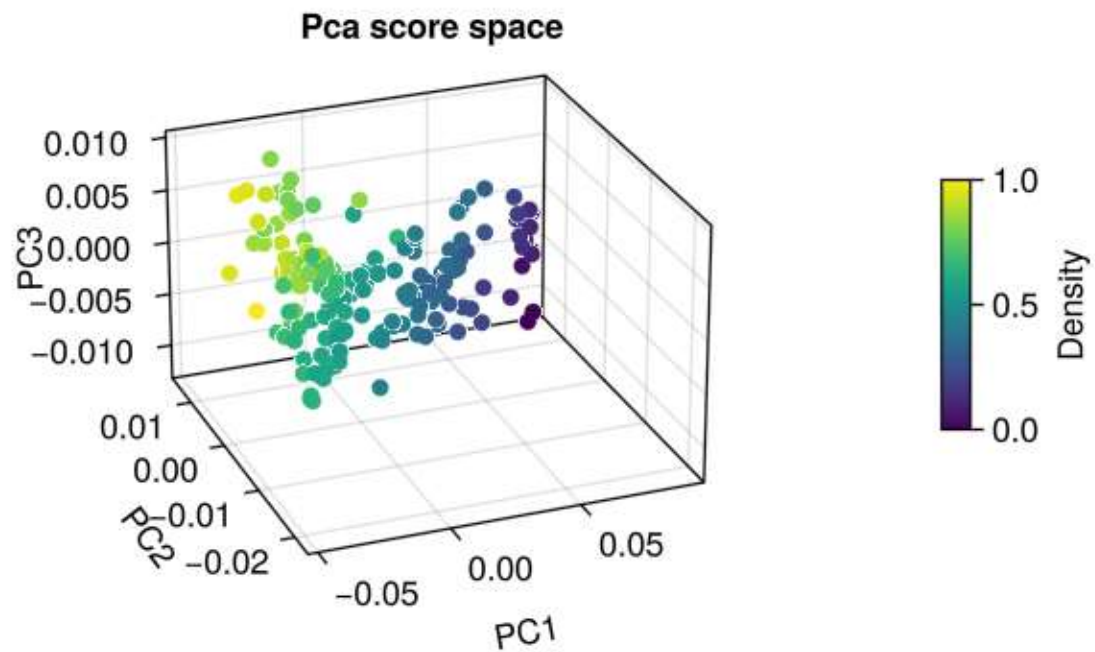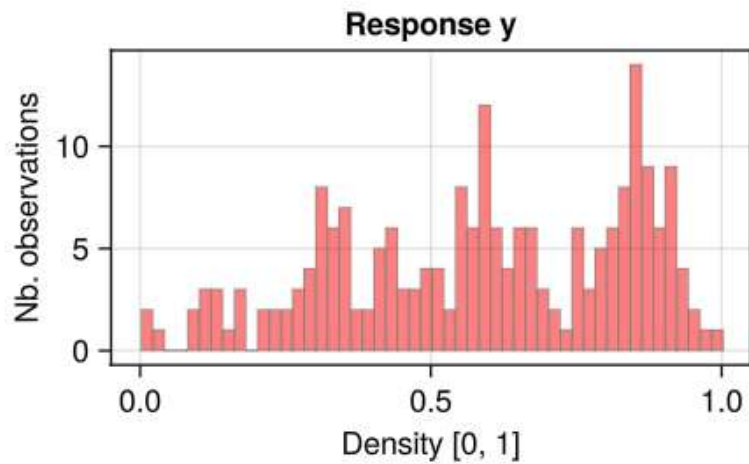
Available (JLD2) at: https://github.com/mlesnoff/**JchemoData.jl**?tab=readme-ov-file#mnist20pcts

## MLNIR data

with pre-selection of range 4000-7000 cm-1

1556 variables (wavelengths)



SNV + SavGol2

**Response y
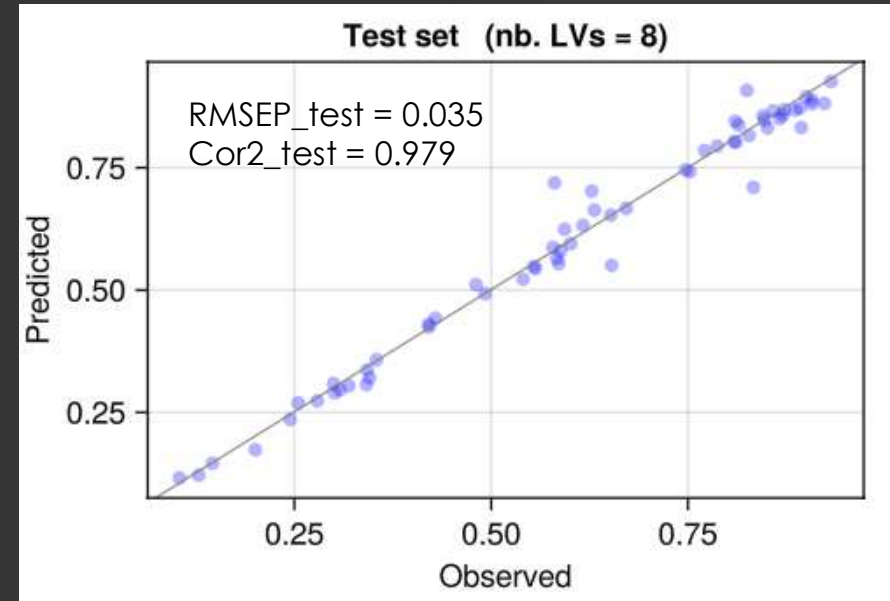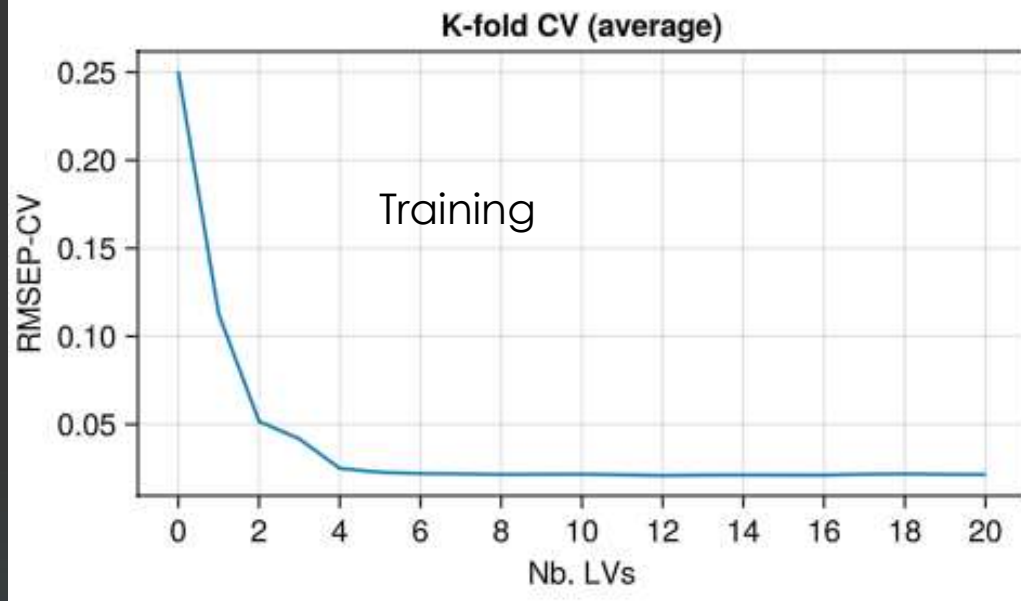Hydrocarbon density**

**Model calibration and validation**

Total   n = 208 obs.

　　　Train   n = 146      ⟹  Replicated K-fold CV   K = 3, nrep = 30

　　　test      n = 62       30% by random sampling

**PLSR**



K-fold CV (average)

Training



Test set   (nb. LVs = 8)

RMSEP_test = 0.035
Cor2_test = 0.979

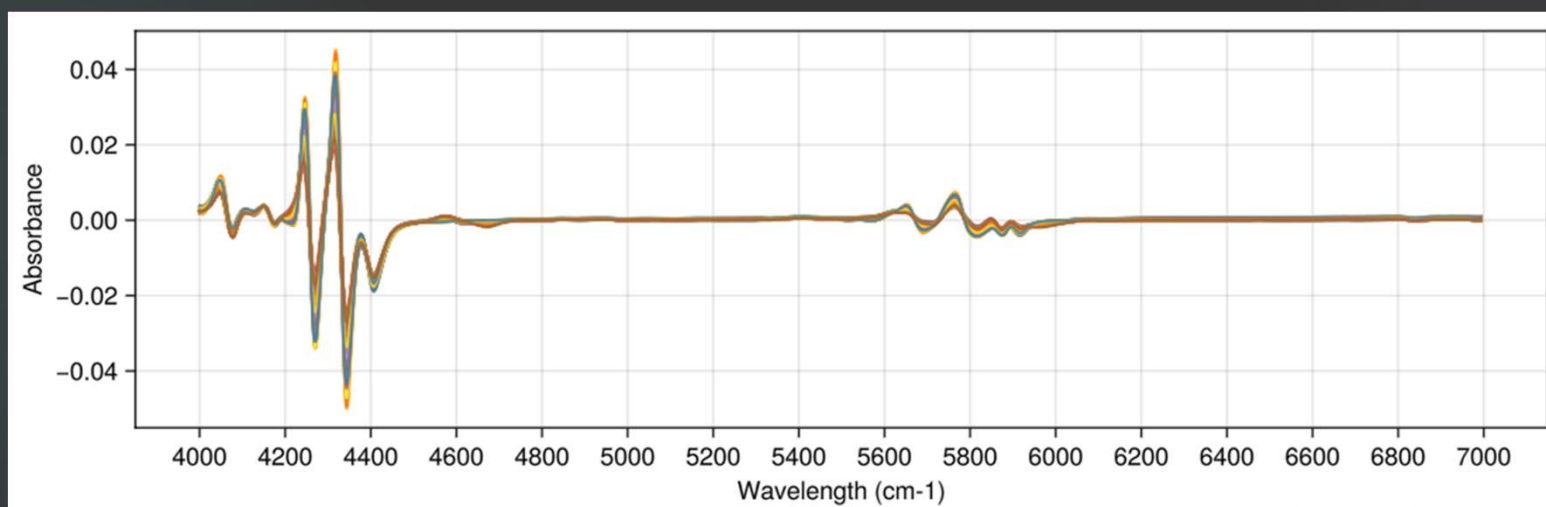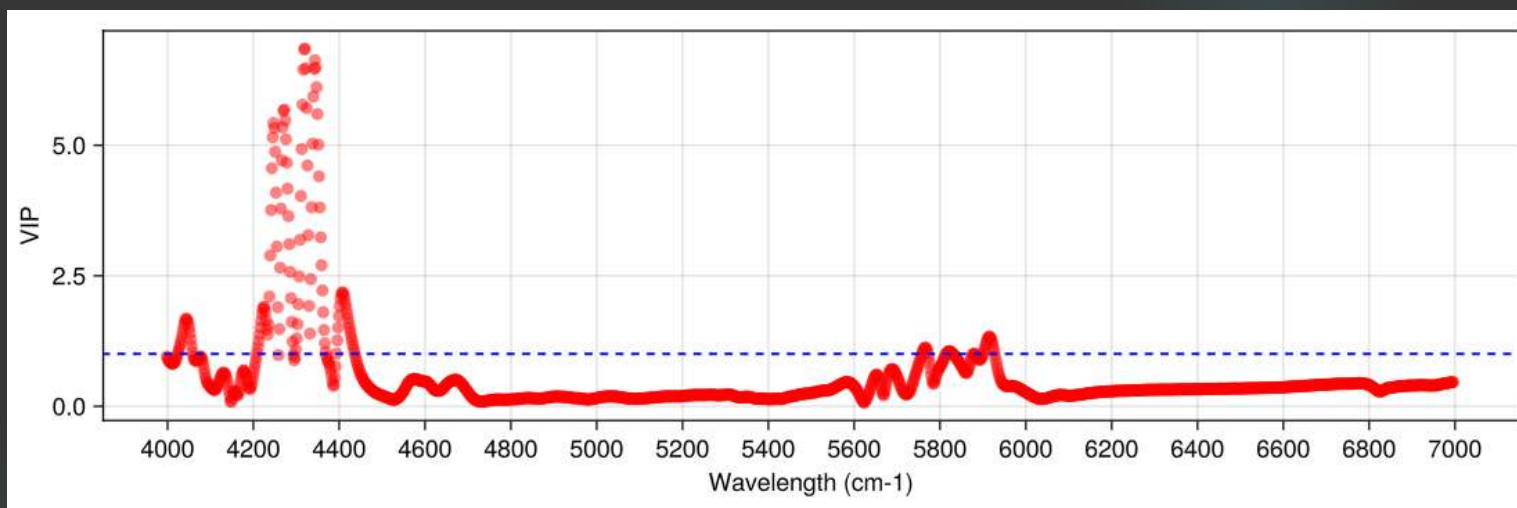**VIPs**

PLS 8 LVs

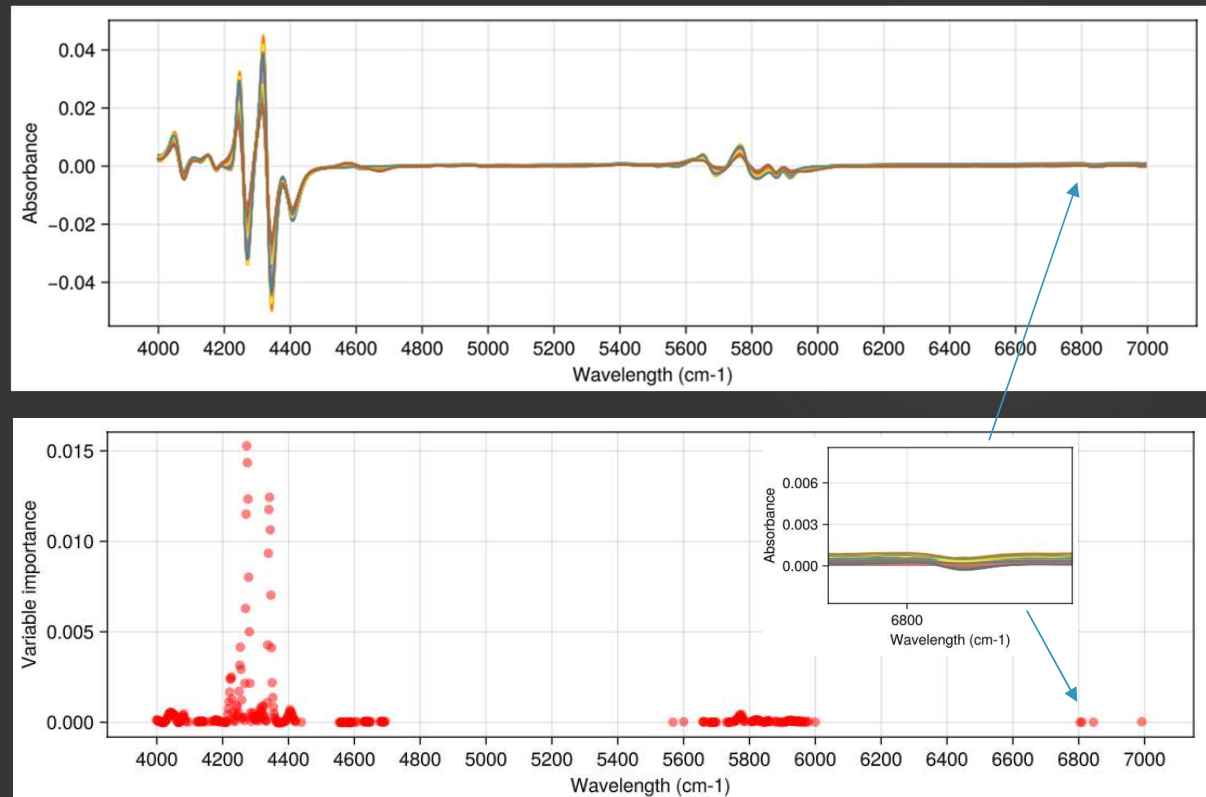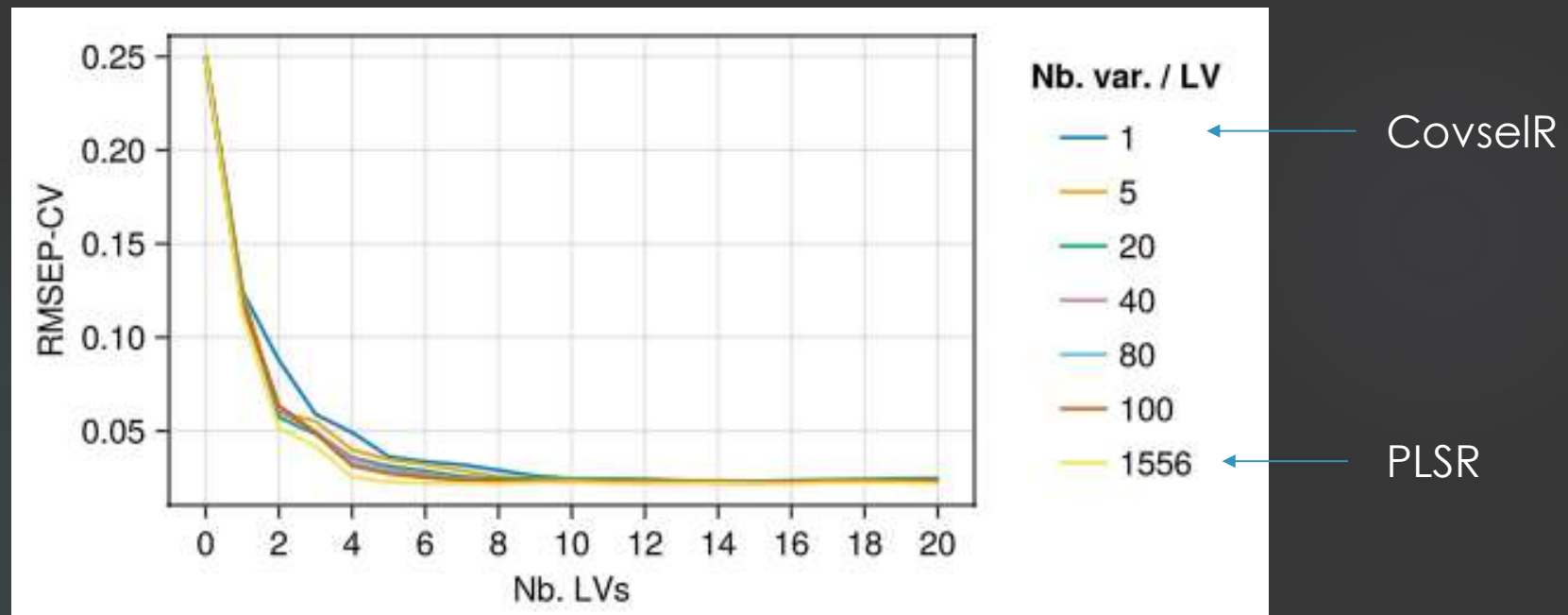**Alternative to VIPs (very simple and efficient):** **Variable importance by permutation**

- Successively for each variable:

    a) Obs. are **randomly permuted**

    b) The model (here PLSR 8 LVs) is fitted on Xtrain and used to predict Xtest

    c) Rmsep_test (or other indicator) is computed and compared to the original without permutation

- Grouped-approaches are very easy to implement (simultaneous permutation of sets of variables)
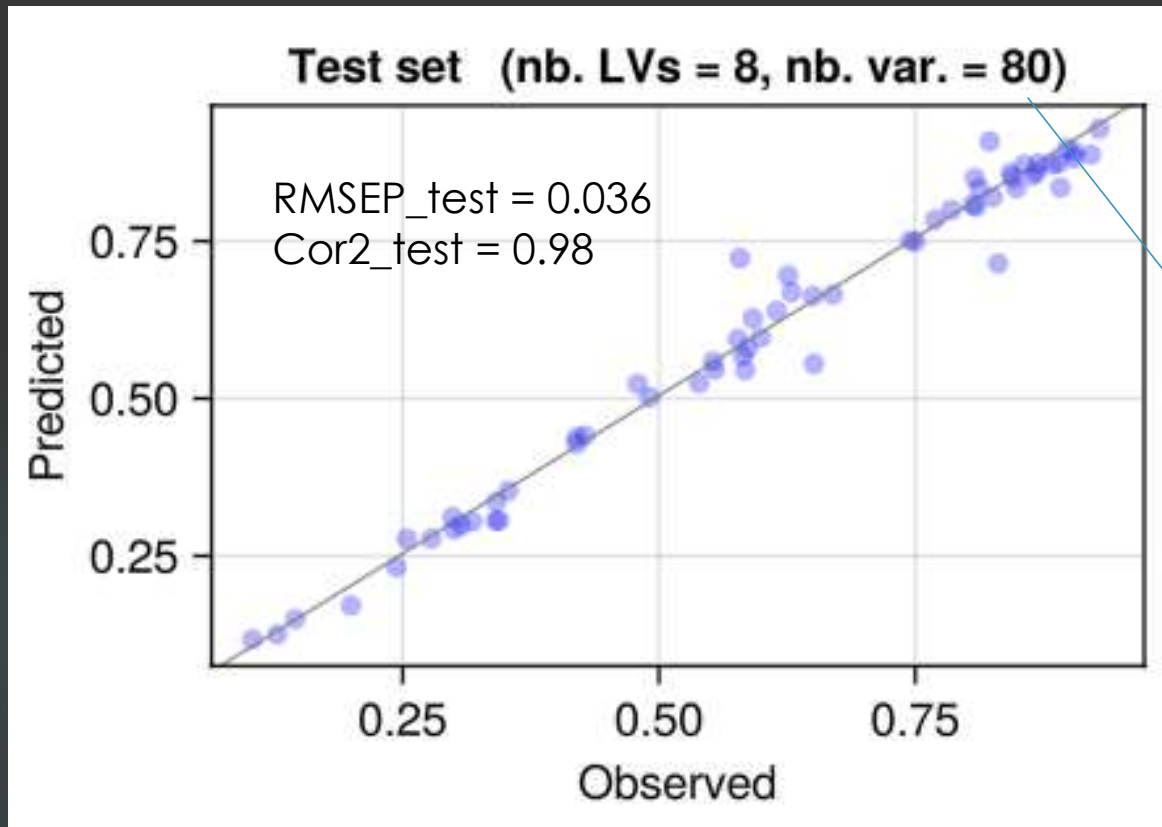
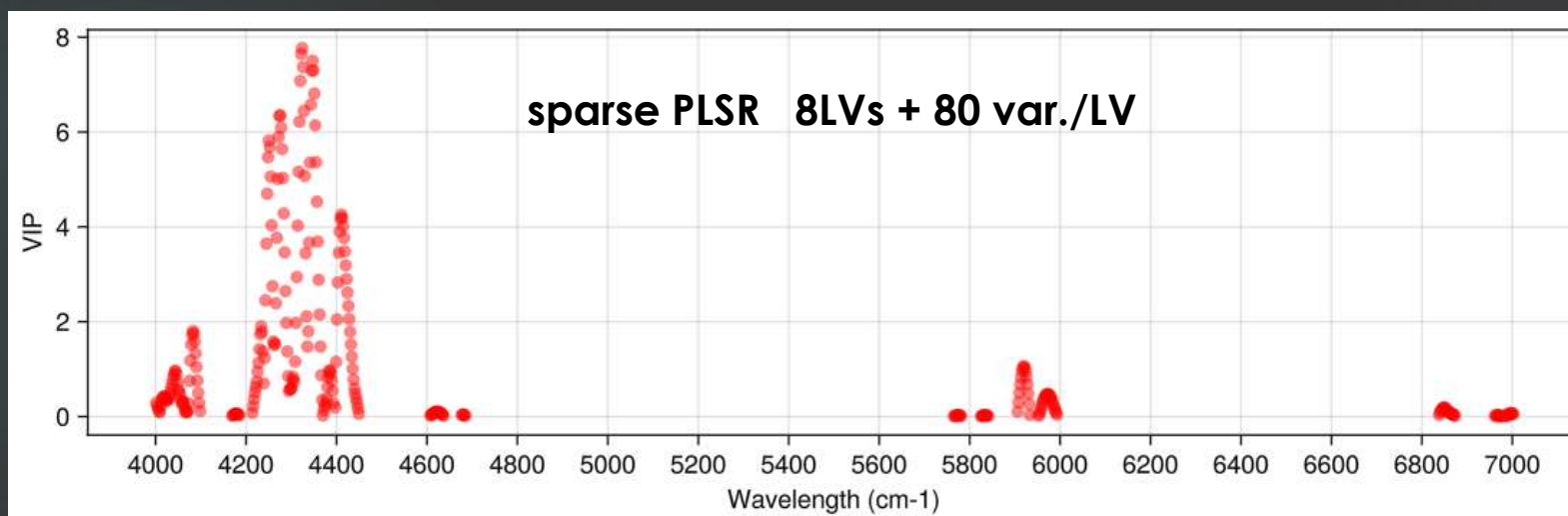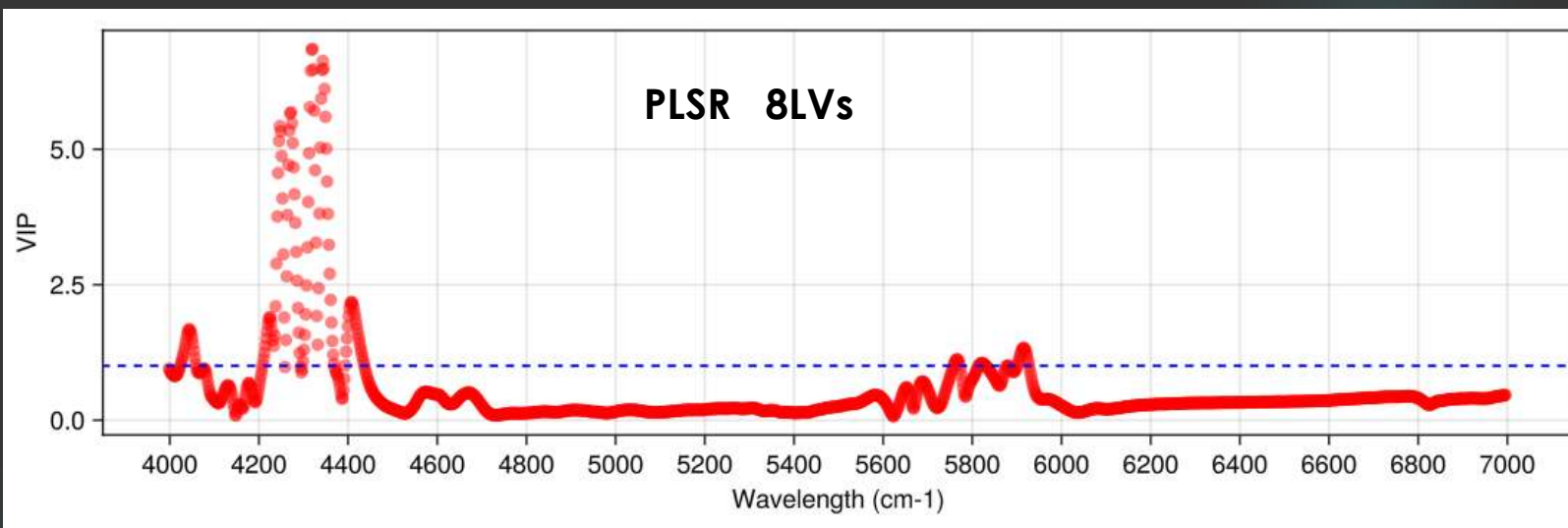- Better than Interval-PLS

**Sparse PLSR**

$\Rightarrow$ CV to tune the nb. LVs and the nb. variables
that are selected for each LV



In this example, the same nb. variables are selected
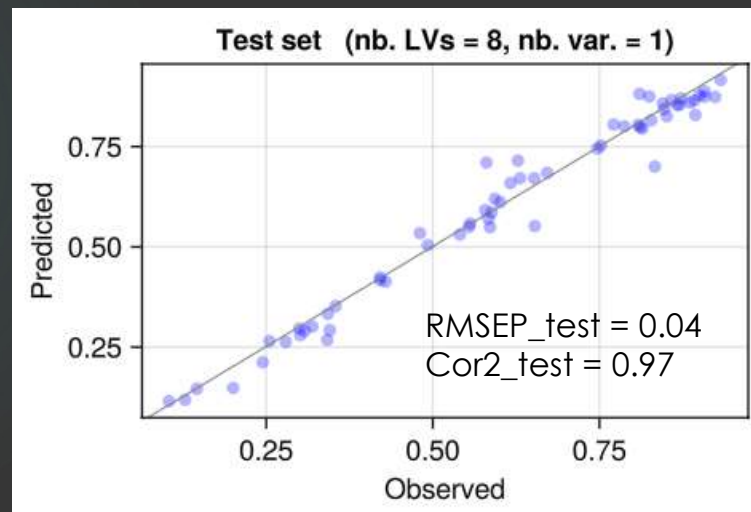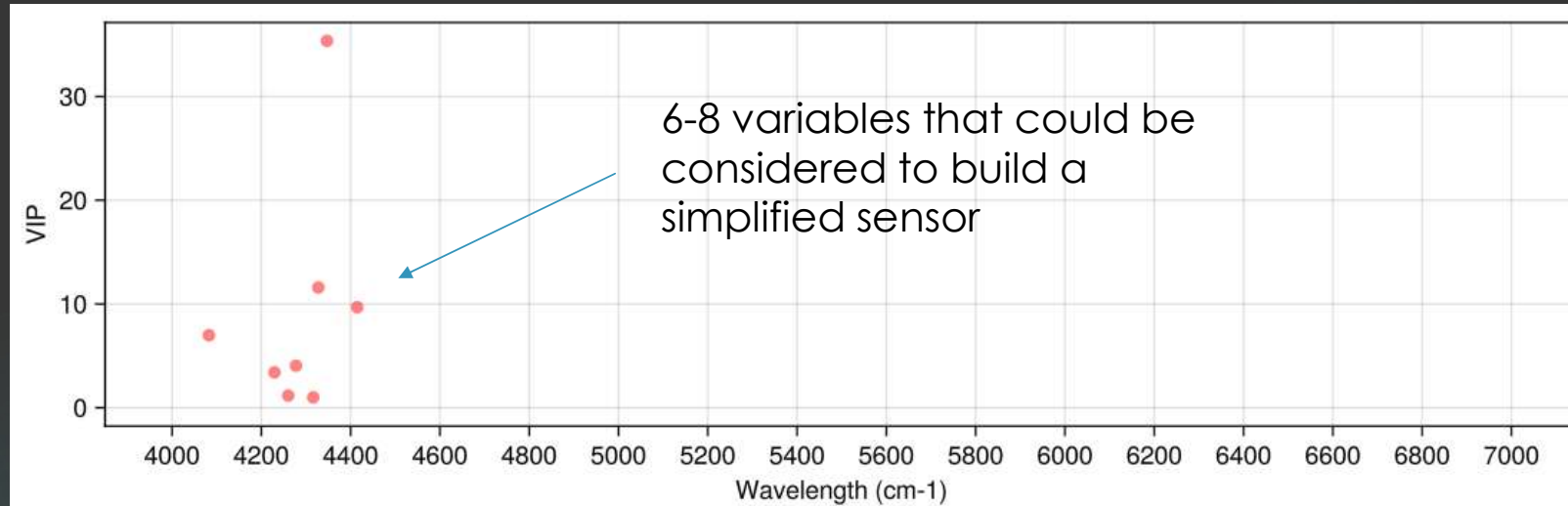for all the LVs, but this can be relaxed

Test set   (nb. LVs = 8, nb. var. = 80)

RMSEP_test = 0.036
Cor2_test = 0.98

⇒ Total nb. selected variables = 256 (over 1556)

PLSR   8LVs

sparse PLSR   8LVs + 80 var./LV

## With only 1 variable selected / LV   (CovSelR)

6-8 variables that could be considered to build a simplified sensor
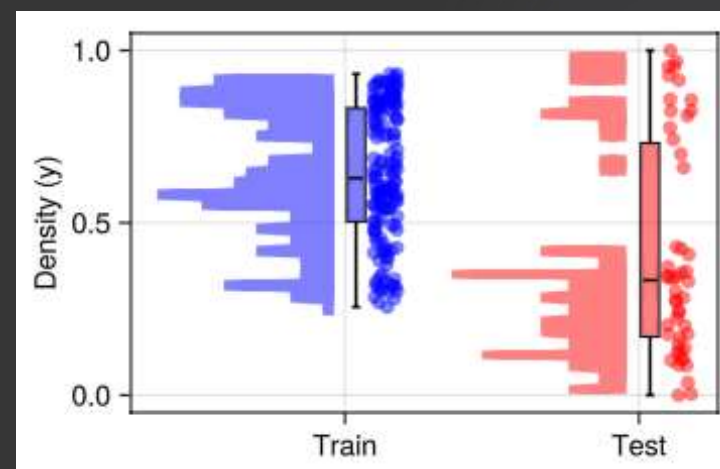


RMSEP_test = 0.04
Cor2_test = 0.97

# And what about predictive robustness?

### New split Train/Test

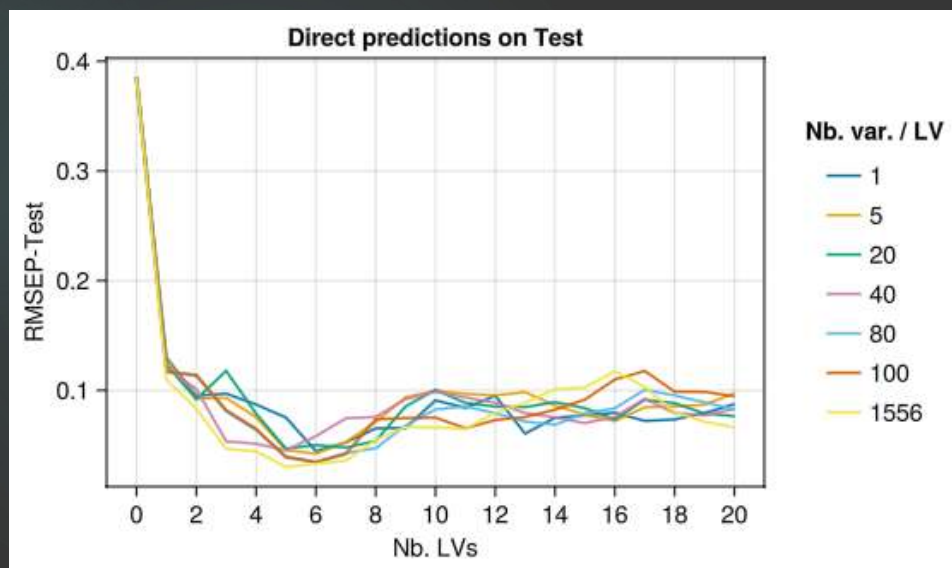- Test (n = 50):  Mainly in extrapolation

Cross-validation on Train

ComSelR

PLSR

sparse PLSR   8LVs   80 var. / LV



RMSEP_test = 0.06,
Cor2_test = 0.97