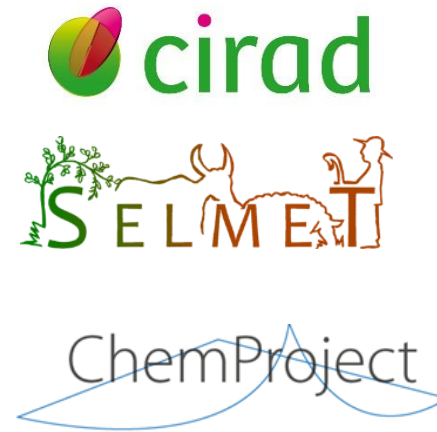# Covariance penalty criteria for selecting dimensions of PLSR models

## Illustration on NIRS data

matthieu.lesnoff@cirad.fr     ChemHouse Montpellier France, 17 November 2020
https://github.com/mlesnoff/rnirs

- **Introduction**

- **DoFs**

- **Ilustration on a NIR dataset**

- **Appendix: some theoretical details**

# Introduction

# Selection of the number of components (LV) in a PLSR model

- $a = 0, 1, ..., A$   components?

## = Model selection

# Many indicators

- **Scores, loadings,  b-coefficients**

- **Prediction errors**

    − Cross-validation, Bootstrap

    − Permutations (Not presented here)

    − Covariance penalty criteria

- **Etc.**

# Covariance penalty criteria

- $\in$ **"Information criteria"**

- Estimate of **prediction error (future)** from the training data set

- **Bias-variance** compromises

- Depend on the **model complexity** $df$

# For OLS regression

- **a well-known covariance penalty** criteria is the Mallows' Cp

$$\text{Cp} = \frac{\text{SSR}}{n} + \frac{2}{n} p \, \hat{\sigma}^2 \qquad (p = \text{number of variables})$$

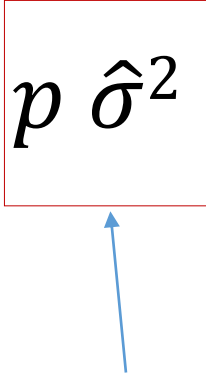**Mallows, C.L., 1973**. Some Comments on Cp. *Technometrics* 15, 661–675.
https://doi.org/10.1080/00401706.1973.10489103

- $\text{Cp} = \dfrac{\text{SSR}}{n} + \dfrac{2}{n} p \, \hat{\sigma}^2$

**Model complexity** *df*

- $\text{Cp} = \dfrac{\text{SSR}}{n} + \dfrac{2}{n} \boxed{p \, \hat{\sigma}^2}$

**Covariance penalty**

**Same interpretation** as for $\widehat{\text{MSEP}}_{\text{CV, BOOT}}$

**Zuccaro, C., 1992**. Mallows' Cp Statistic and Model Selection in Multiple Linear Regression. *International Journal of Market Research*. 34, 1–10. https://doi.org/10.1177/147078539203400204

## For OLS (and some other models)

Cp is calculated from the training set
without any simulation

$\Rightarrow$ very fast

# Cp

- Estimate of a type of **(future) prediction error**, referred to as $\mathrm{Err}_{in}$ in Hastie et al 2009

- Not exactly the same error estimated by CV or Bootstrap ($\mathrm{Err}$) but both approaches often return close models
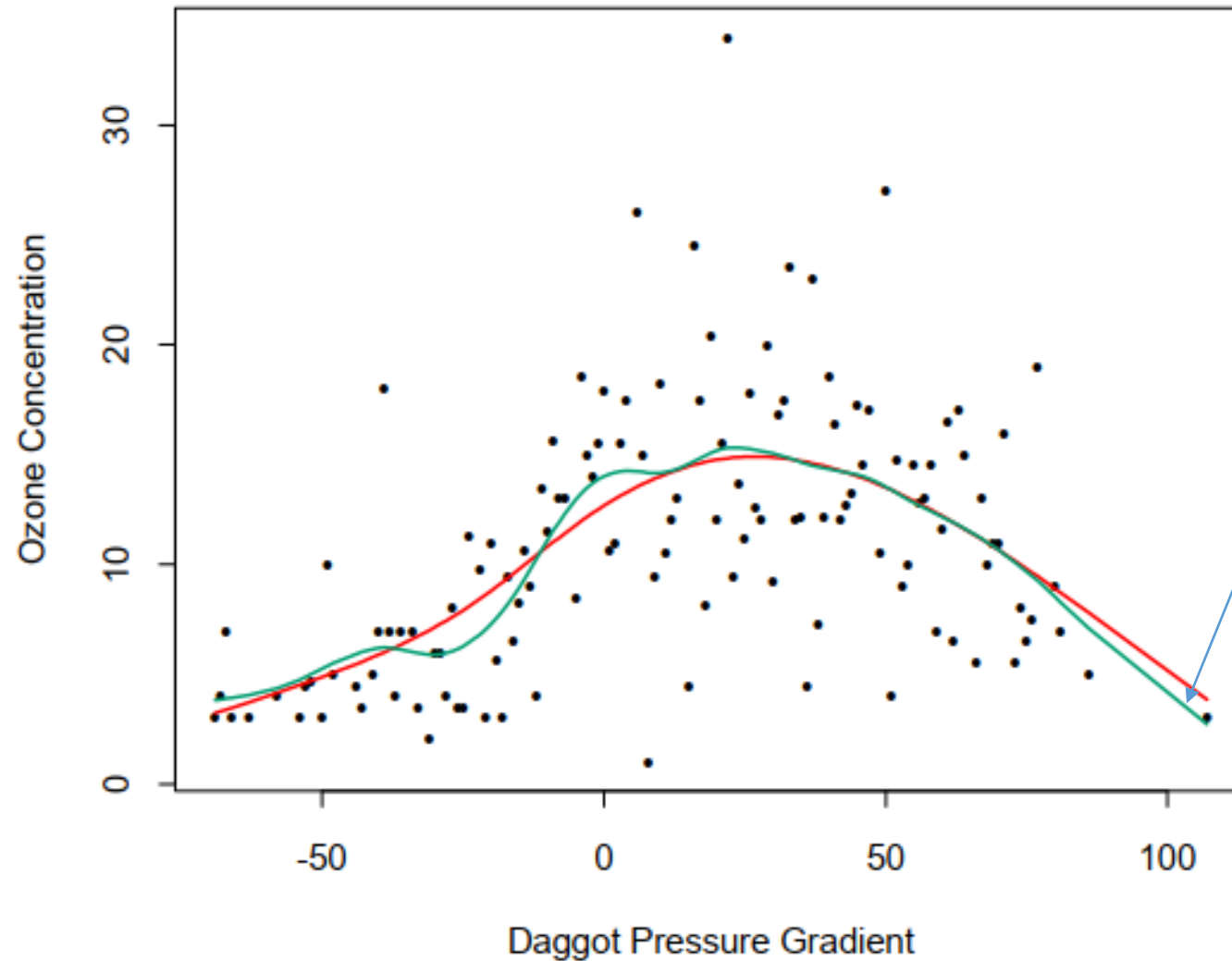
See Appendix for $\mathrm{Err}_{in}$ and $\mathrm{Err}$
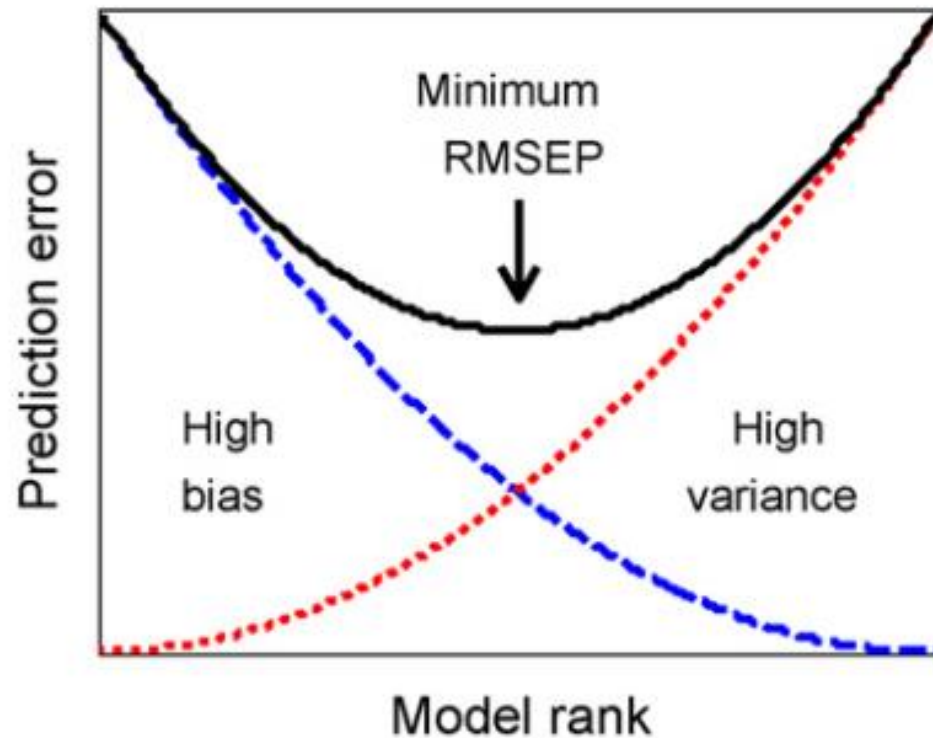
- Prediction errors   = Bias-variance compromises

$$= \text{Variance}_\tau\left(\widehat{\boldsymbol{y}}\right) + \text{Bias}_\tau\left(\widehat{\boldsymbol{y}}\right)^2 + \ldots$$

training

$\tau$: training data set
   of size $n$
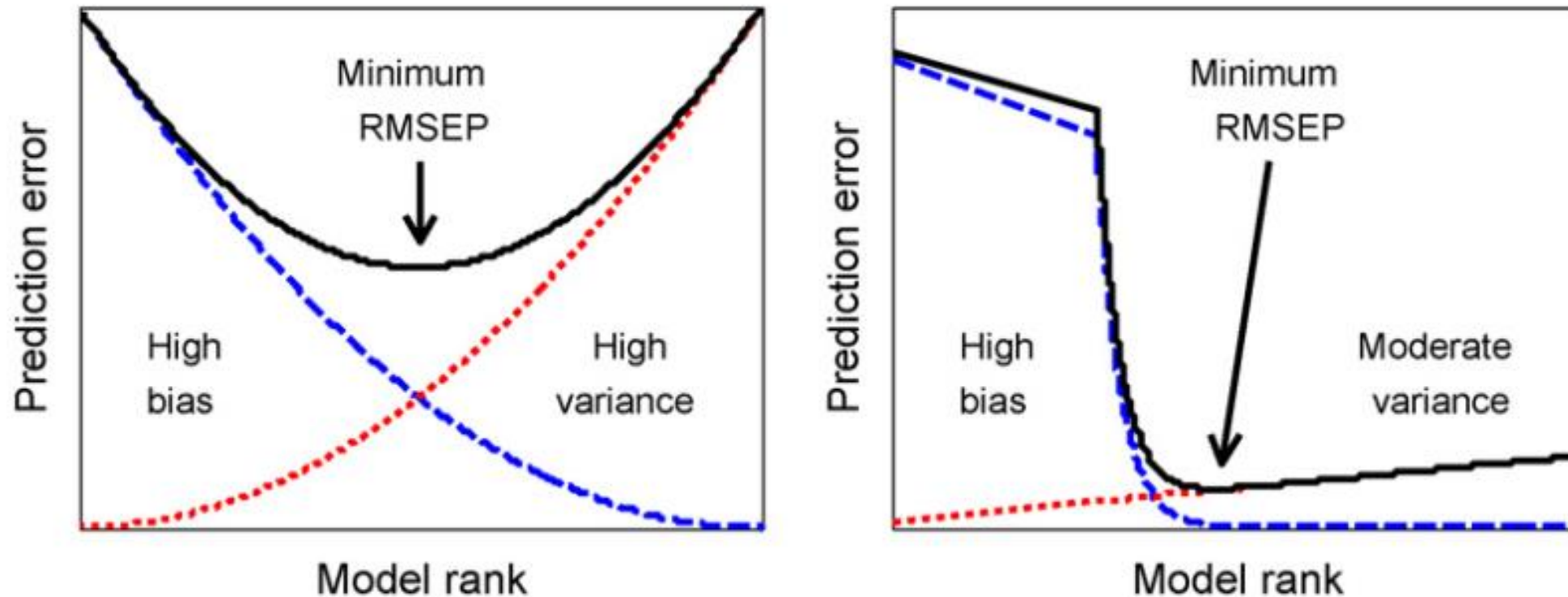
**Hastie *et al* 2009   Fig 5.7**



**More complex model**
**Less biased but more**
**variable** with $\tau$

Faber, N. (KLAAS) M., 1999. A closer look at the bias–variance trade-off in multivariate calibration. Journal of Chemometrics 13, 185–192. https://doi.org/10.1002/(SICI)1099-128X(199903/04)13:2<185::AID-CEM538>3.0.CO;2-N
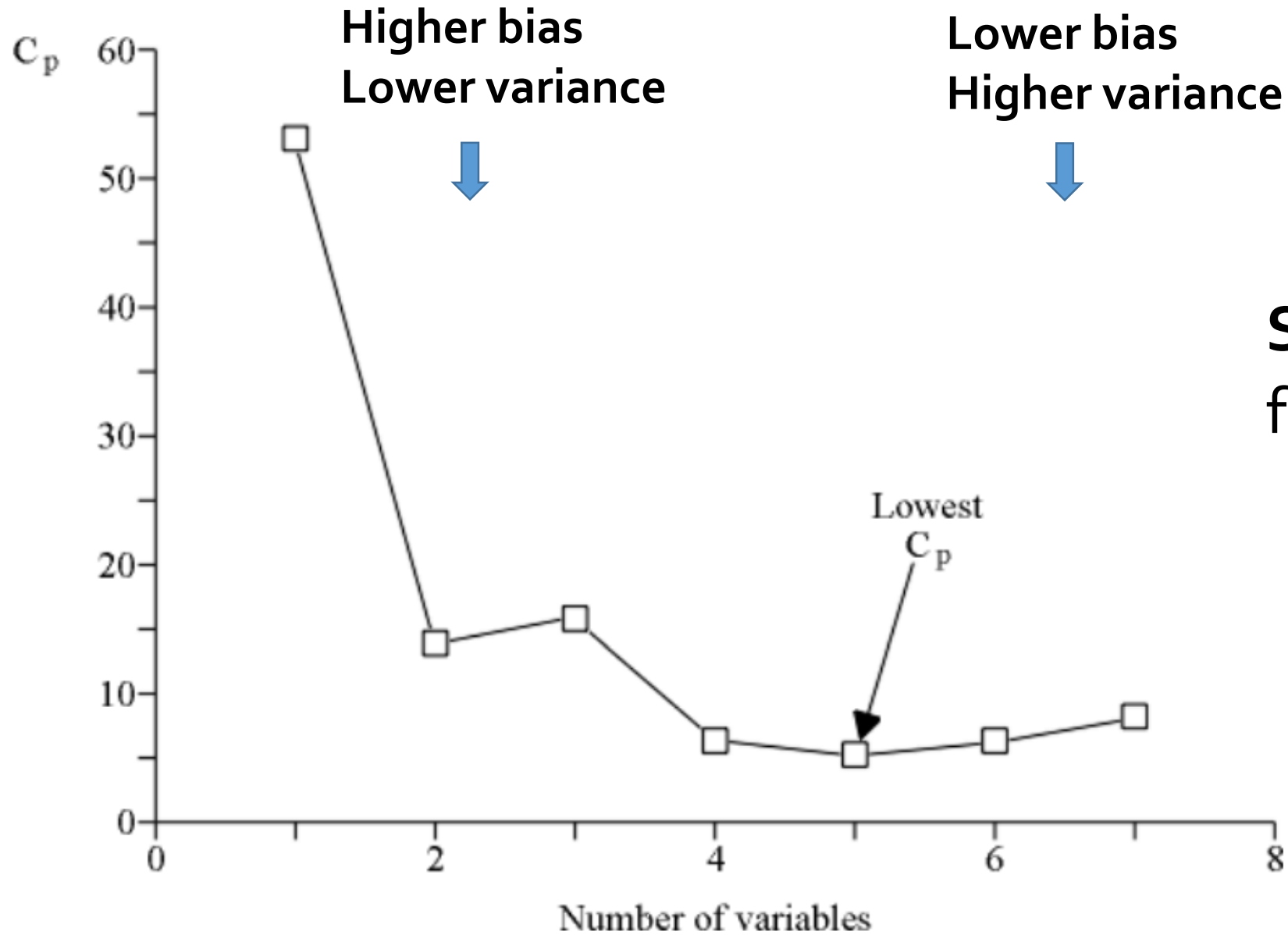
Faber, N.M., Rajkó, R., 2007. How to avoid over-fitting in multivariate calibration—The conventional validation approach and an alternative. Analytica Chimica Acta, Papers presented at the 10th International Conference on Chemometrics in Analytical Chemistry 595, 98–106. https://doi.org/10.1016/j.aca.2007.05.030

Faber, N. (KLAAS) M., 1999. A closer look at the bias–variance trade-off in multivariate calibration. Journal of Chemometrics 13, 185–192. https://doi.org/10.1002/(SICI)1099-128X(199903/04)13:2<185::AID-CEM538>3.0.CO;2-N

Faber, N.M., Rajkó, R., 2007. How to avoid over-fitting in multivariate calibration—The conventional validation approach and an alternative. Analytica Chimica Acta, Papers presented at the 10th International Conference on Chemometrics in Analytical Chemistry 595, 98–106. https://doi.org/10.1016/j.aca.2007.05.030

$$Var(\overline{x}) = Var(\hat{\mu}) = \frac{\sigma^2}{n}$$

Higher bias
Lower variance

Lower bias
Higher variance

**Same** as
for $\widehat{\text{MSEP}}_{CV,\,BOOT}$

Lowest
$C_p$

$C_p$

Number of variables

## Faber & Rajkó 2007

of a better product. A tacit assumption is that the components included in the model are ordered according to their importance for describing the $Y$-variable – the property of interest. It has been observed, however, that non-significant components can be preceded and followed by (highly) significant ones [16,23]. This phenomenon has been termed 'sandwiching' and can often be rationalized (see [24–27] for in-depth discussions of this aspect). An early component can, for example, take care of a background in the $X$-data and it consequently bears no relationship with the $Y$-variable. (Recall that PLS component 3 of the current example data set is close to being non-significant, while the preceding and following ones are highly significant.) It is clear that one should be cautious when attempting to interpret these 'sandwiched' components. We therefore recommend dis-

**Remark**
"Sandwiching" effect in PLSR

# For models with Gaussian errors

- Cp is equivalent to the **Akaike Criterion** (AIC)

$$\text{AIC} = n \log(\text{SSR}) + 2\,(p + 1) \quad \text{(after removing non-useful constants)}$$

About AIC

**Hurvich, C.M., Tsai, C.-L.,** 1989. Regression and Time Series Model Selection in Small Samples. *Biometrika* 76, 297. https://doi.org/10.2307/2336663

**Burnham, K.P., Anderson, D.R.,** 2002. Model selection and multimodel inference: a practical information-theoretic approach, 2nd ed. Springer, New York, NY, USA.

**Variants of** Cp consist in increasing the penalty coefficient "2"
for preventing overfitting → Ex: <span style="color:red">BIC penalty</span>

- $\text{Cp} = \dfrac{\text{SSR}}{n} + \dfrac{\textcolor{red}{2}}{n}\, p\, \hat{\sigma}^2$

- $\text{Cp(BIC)} = \dfrac{\text{SSR}}{n} + \dfrac{\textcolor{red}{\log(n)}}{n}\, p\, \hat{\sigma}^2$

$\log(100) = 4.6$

$\log(1000) = 6.9$

See discussion in:
 **Eubank, R.L., 1999**. Nonparametric Regression and Spline Smoothing, 2nd ed, Statistics: Textbooks and Monographs. *Marcel Dekker, Inc.*, New York, USA.

# Bias correction for $\mathrm{Cp}$ $\rightarrow$ For "small" $n$

- $\mathrm{Cp(AICc)} = \dfrac{\mathrm{SSR}}{n} + \dfrac{2}{n} p\, \hat{\sigma}^2 \dfrac{\color{red}{n}}{\color{red}{n-p-1}}$

**Hurvich, C.M., Tsai, C.-L.,** 1989. Regression and Time Series Model Selection in Small Samples. *Biometrika* 76, 297. https://doi.org/10.2307/2336663

# Generalization to other models than OLS regression

- **General form**

$$\text{Cp} = \frac{\text{SSR}}{n} + \frac{2}{n} \, df \, \hat{\sigma}^2$$

**Need to be calculated**

- *df* is easier to calculate for **linear smoothers**
  
  (OLS, Splines, Ridge, etc.)

  – $\widehat{y} = S\,y$     <span style="color:red">$y$ not involved in $S$</span>

  – Ex:    OLS regression    $\widehat{y}$    $= S\,y$    ($S$: hat matrix)
  
  $$= X(X'X)^{-1}X'\,y$$
  
  $$df = tr(2S - SS') = tr(S) = p$$

- PLSR is a non linear smoother …

- PLSR $=$ OLS regression on $a$ scores

  - But $\quad \widehat{\boldsymbol{y}} = \boldsymbol{S_{a,y}} \, \boldsymbol{y} \qquad \boldsymbol{y}$ is involved in $\boldsymbol{S}$

  - Consequence: $\quad df > a + 1$

    ("+ 1" is for the intercept)

- Same difficulty in PCA $\qquad \text{vec}(\widehat{\boldsymbol{X}}) = \boldsymbol{S}_{a,X} \text{vec}(\boldsymbol{X})$

- In PCR $\qquad df = a + 1$

# The idea $df_{\text{PLSR}} > a + 1$ is not new in chemometrics, for instance …

- **Martens, H. and Naes, T. (1989)**. Multivariate Calibration. _Wiley_, New York.

- **Frank, Ildiko E., Friedman, J.H., 1993**. A Statistical View of Some Chemometrics Regression Tools. _Technometrics_ 35, 109–135. https://doi.org/10.1080/00401706.1993.10485033

- **Voet, H. van der, 1999**. Pseudo-degrees of freedom for complex predictive models: the example of partial least squares. _Journal of Chemometrics_ 13, 195–208. https://doi.org/10.1002/(SICI)1099-128X(199905/08)13:3/4<195::AID-CEM540>3.0.CO;2-L

- **Denham, M.C., 2000**. Choosing the number of factors in partial least squares regression: estimating and minimizing the mean squared error of prediction. _Journal of Chemometrics_ 14, 351–361. https://doi.org/10.1002/1099-128X(200007/08)14:4<351::AID-CEM598>3.0.CO;2-Q

- **Krämer, N., Braun, M.L., 2007**. Kernelizing PLS, degrees of freedom, and efficient model selection, in: Proceedings of the 24th _International Conference on Machine Learning, ICML_ '07. Association for Computing Machinery, New York, NY, USA, pp. 441–448. https://doi.org/10.1145/1273496.1273552

- **Krämer, N., Sugiyama, M., 2011**. The Degrees of Freedom of Partial Least Squares Regression. _Journal of the American Statistical Association_ 106, 697–705. https://doi.org/10.1198/jasa.2011.tm10107

# … But often forgotten

Example

- **Li, B., Morris, J., Martin, E.B., 2002**. Model selection for partial least squares regression. _Chemometrics and Intelligent Laboratory Systems_ 64, 79–89. https://doi.org/10.1016/S0169-7439(02)00051-5

  **Authors** used the naïve $df = a + 1$ in the AIC criterion

Another argument often made in favor of PLS over PCR is that PCR only uses the predictor sample to choose its components, whereas PLS uses the response values as well. This argument is not unrelated to the one discussed previously. By using the response values to help determine its components, PLS uses more degrees of freedom per component and thus can fit the training data to a higher degree of accuracy than PCR with the same number of components. As a consequence, a $K$-component PLS solution will have less bias than the corresponding $K$-component PCR solution. It will, however, have greater variance, and since the mean squared prediction error is the sum of the two (bias squared plus variance) it is not clear which solution would be better in any given situation. In any case, either method is free to choose its own number of components (bias-variance trade-off) through model selection (CV).

Franck & Friedman
Technometrics 1993

Figure 1: Estimated Degrees of Freedom (stars) for the three benchmark data sets. The solid line displays the naive estimate $\mathrm{DoF}(m) = m + 1$. If the assumption of theorem 3 is fulfilled, we also display the lower bound on the Degrees of Freedom for 1 component (dashed horizontal line).

**Krämer & Sugiyama Jasa 2011**



Figure 3: Training error of PLSR and PCR. Left: Training error as a function of the number of components. Right: Training error as a function of the Degrees of Freedom.

**Can we use $C_p$ (AIC, etc.) for PLSR?**

- yes if one uses the **relevant model complexity** *df*

# Model's degree of freedom

- **Hastie, T., Tibshirani, R.J., 1990**. Generalized Additive Models, Monographs on statistics and applied probablity. *Chapman and Hall/CRC*, New York, USA.

- **Ye, J., 1998**. On Measuring and Correcting the Effects of Data Mining and Model Selection. *Journal of the American Statistical Association* 93, 120–131. https://doi.org/10.1080/01621459.1998.10474094

- **Eubank, R.L., 1999**. Nonparametric Regression and Spline Smoothing, 2nd ed, Statistics: Textbooks and Monographs. *Marcel Dekker, Inc.*, New York, USA.

- **Efron, B., 2004**. The Estimation of Prediction Error. J*ournal of the American Statistical Association* 99, 619–632. https://doi.org/10.1198/016214504000000692

- **Zou, H., Hastie, T., Tibshirani, R., 2007**. On the "degrees of freedom" of the lasso. *The Annals of Statistics* 35, 2173–2192. https://doi.org/10.1214/009053607000000127

- **Hastie, T., Tibshirani, R., Friedman, J., 2009**. The elements of statistical learning: data mining, inference, and prediction, 2nd ed. *Springer*, New York.

- **Hastie, T., Tibshirani, R., Wainwright, M., 2015**. Statistical Learning with Sparsity: The Lasso and Generalizations. *CRC Press*.

# *df*

- Model complexity

- **Effective number of parameters** in the model

For any model (linear or not) with additive homoscedastic error

the **consensus** is to consider that

$$df = \sum_{i=1}^{n} Cov_{\varepsilon}(y_i, \hat{y}_i) \big/ \sigma^2$$

Variation of the training $\boldsymbol{y}$ ($\boldsymbol{X}$ fixed)

Hastie et al 2015 p.18        *Self-influence* that each response measurement has on its prediction

Hastie *et al* 2009   Fig 5.7

**More complex model**
**A change in** $y_i$ will have more impact

**For OLS regression**

- $Cov_{\boldsymbol{\varepsilon}}(yi, \hat{y}_i) = p\,\sigma^2/n$

$$\Rightarrow \quad df = \sum_{i=1}^{n} Cov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i)/\sigma^2 = p$$

$\Rightarrow$ This covariance definition of $df$ is <span style="color:red">consistent</span>
    with the historical definition for OLS regressions $df_{\text{OLS}} = p$

For Gaussian errors $\boldsymbol{\varepsilon}$     (Stein 1981, Efron 2004)

$$Cov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i) \;=\; \sigma^2 E_{\boldsymbol{\varepsilon}}\left(\frac{\partial \hat{y}_i}{\partial y_i}\right)$$

**Stein, C.M., 1981**. Estimation of the Mean of a Multivariate Normal Distribution. *The Annals of Statistics* 9, 1135–1151.

# Leading to two usual forms for $df$

1) $df = \sum_{i=1}^{n} Cov_{\varepsilon}(y_i, \hat{y}_i)/\sigma^2$

2) $df = \sum_{i=1}^{n} E_{\varepsilon}\left(\frac{\partial \hat{y}_i}{\partial y_i}\right)$ 　　　**Gaussian error (Stein)**

Generalized df
(**Ye Jasa 1998**)

Divergence $div_i$

# Estimation of $df$

- Formal derivation (exact or approximation)

    – Either $Cov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i)$ or $\dfrac{\partial \hat{y}_i}{\partial y_i}$

- Monte Carlo simulations

# Monte Carlo

1) $\widehat{df} = \sum_{i=1}^{n} \hat{C}ov_{\varepsilon}(y_i, \hat{y}_i)/\widehat{\sigma}^2$ ← **Parametric bootstrap** (e.g. Efron 2004)

2) $\widehat{df} = \frac{1}{n}\sum_{i=1}^{n}\frac{\partial \hat{y}_i}{\partial y_i}$ (SURE) ← **Sensitivity analysis** (perturbations e.g. Ye 1998)

Stein Unbiased Risk Estimation

**Sorghum NIR data**
$n = 1006$   $p = 700$
% Crude fibers
prediction

**PCR**

$df_{\text{naive}} = a + 1$

$\widehat{df}_{\text{SURE}}$
**with** $m = 50$
**perturbation points**

df

Nb components

# PCR



$\widehat{df}_{\text{SURE}}$

**with $m$ = 1000 perturbation points**

# PLSR



$\widehat{df}_{\text{Cov}}$
with $B$ = 50 bootstrap replications

$\widehat{df}_{\text{SURE}}$
with $m$ = 50 perturbation points

$df_{\text{naive}} = a + 1$

$\widehat{df}_{\text{Cov}}$
with $B$ = **100 bootstrap**
**replications**

$\widehat{df}_{\text{SURE}}$
with $m$ = **100**
**perturbation points**

# Ratio to the naive *df*

# Krämer & Sugiyama Jasa 2011

Proposed formulas for calculating $\widehat{df}_{\text{SURE}}$

$$\hat{\boldsymbol{y}} = \boldsymbol{S}_{a,y}\,\boldsymbol{y} \qquad \text{with} \qquad \boldsymbol{S}_{a,y} = \boldsymbol{T}_{a,y}\,(\boldsymbol{T}_{a,y}'\,\boldsymbol{T}_{a,y})^{-1}\boldsymbol{T}_{a,y}'$$

$$\text{tr}\left(\frac{\partial\hat{\boldsymbol{y}}}{\partial\boldsymbol{y}}\right) = \text{tr}\left(\frac{\partial(\boldsymbol{S}_{a,y}\,\boldsymbol{y})}{\partial\boldsymbol{y}}\right)$$

R package `plsdof` (CRAN)

**Algorithm 1**

- <span style="color:red">Lanczos algorithm</span> for PLSR

- Step of derivation for each component (Kramer & Braun 2007)

- **High calculation times**

**Algorithm 2**

- $\hat{y}$ is expressed using <span style="color:red">Krylov subspaces</span>

- This simplifies the calculation of the trace of the derivative

- **Much faster calculation times**

**Sorghum NIR data**
*n* = 1006

Alg 2
(Krylov)

**Alg 1
(Lanczos)**

**Monte Carlo
Cov & SURE**

df

Nb components

**Alg 2 (Krylov) Numeric unstability?**

**(Alg 1 crashed)**

**Monte Carlo Cov & SURE**

df

Nb components

**Same problem** with another dataset **Cassava** $n = 200$ (carotenoids)

Alg 2 (Krylov)

Monte Carlo Cov

# Coming back to covariance penalty criteria

- $\text{Cp} = \dfrac{\text{SSR}}{n} + \dfrac{2}{n} \widehat{df} \; \hat{\sigma}^2$

- $\widehat{df} = \sum_{i=1}^{n} \hat{C}ov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i)/\widehat{\sigma}^2 \;=\; (\text{SURE}) \; \sum_{i=1}^{n} \dfrac{\partial \hat{y}_i}{\partial y_i}$

**Two estimates**

1) $\mathrm{Cp} = \dfrac{\mathrm{SSR}}{n} + \dfrac{2}{n}\sum_{i=1}^{n}\hat{C}ov_{\varepsilon}(y_i, \hat{y}_i)$

2) $\mathrm{Cp} = \dfrac{\mathrm{SSR}}{n} + \dfrac{2}{n}\hat{\sigma}^2\sum_{i=1}^{n}\dfrac{\partial \hat{y}_i}{\partial y_i}$      SURE estimate

# Ilustration of PLSR model selection
## on a NIR dataset

- Presentation of the data

- Separation between training vs test sets

- CV results

- Examination of the PLS scores

- Examination of the loadings and b-coefficients

- Cp results

- Sensitivity to the test set

# Data

# Sorghum (stem, leafs etc.)

Dried and grounded
Spectra FOSS 1100-2498 nm (step = 2 nm)

$n = 1206$   $p = 700$

savgol(snv(X), m = 2, n = 21, p = 2)

# Prediction of % crude fibers

PCA

# Splitting the data

# Training *vs.* Test sets

**Hypothesis**     Future =     same mechanism as the training
(probability distribution $F$)

- **Training set**         $F \rightarrow \tau = \{ (\boldsymbol{x}_1, y_1), \dots, (\boldsymbol{x}_n, y_n) \}$

- **New observation**      $F \rightarrow (\boldsymbol{x}^*, y^*)$

(Will be relaxed later)

# Examples in chemometric journals

- **Denham, M.C., 2000**. Choosing the number of factors in partial least squares regression: estimating and minimizing the mean squared error of prediction. *Journal of Chemometrics* 14, 351–361.
  https://doi.org/10.1002/1099-128X(200007/08)14:4<351::AID-CEM598>3.0.CO;2-Q

> where expectation is over the original responses **y** and the future $y(\mathbf{x}')$. To estimate this quantity, it is necessary to define the ==nature of the future observations== we shall be predicting. Here we shall assume that the explanatory variables for future observations can be regarded as being selected at random from the set of observations used in estimating the relationship. This will be appropriate when the original data have been randomly sampled from the same distribution as the samples to be predicted or when the original data have been chosen to reflect the distribution of the future samples. Under

- **Mevik, B.-H., Cederkvist, H.R., 2004**. Mean squared error of prediction (MSEP) estimates for principal component regression (PCR) and partial least squares regression (PLSR). *Journal of Chemometrics* 18, 422–429.
  https://doi.org/10.1002/cem.887

> ## 2. MSEP ESTIMATORS
>
> We assume that we have a learning data set $L = \{(\mathbf{x}_i, y_i)\}$ of $n_L$ observations and a predictor $f_L$ trained on $L$. In the present paper this will be PLSR or PCR. For the simulations we also assume that we have a test data set $T = \{(\mathbf{x}_{T,i}, y_{T,i})\}$ of size $n_T$. Both $L$ and $T$ are assumed to be random samples from a ==common distribution.==

- See also **Faber, N. (Klaas) M., 1999**. Estimating the uncertainty in estimates of root mean square error of prediction: application to determining the size of an adequate test set in multivariate calibration. *Chemometrics and Intelligent Laboratory Systems* 49, 79–89. https://doi.org/10.1016/S0169-7439(99)00027-1  **Etc. !!!**

# Remark 1)  True generating distribution $F$

− **Traditional chemometrics**

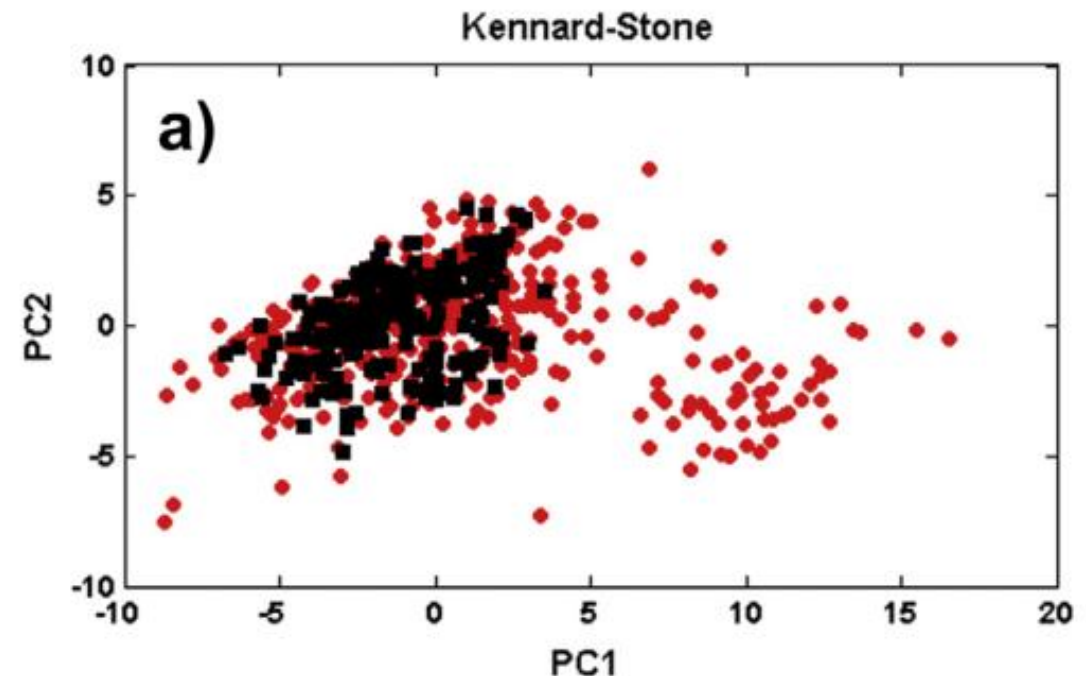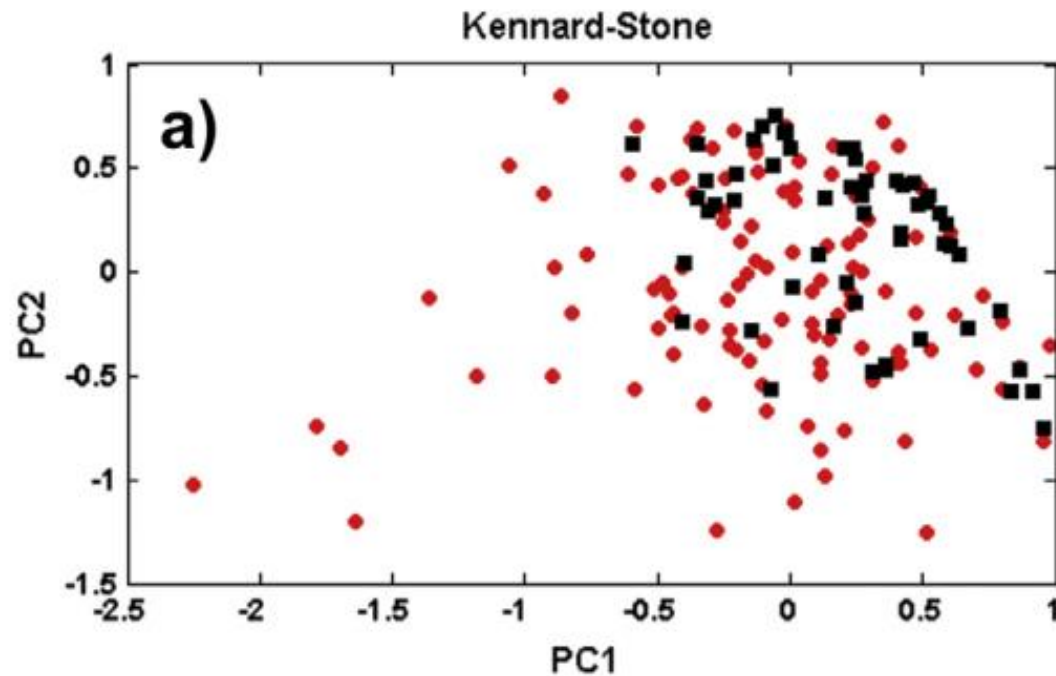  • Hidden $F$ is often considered <span style="color:red">having low dimension</span>

  • Bringing out these dimensions

− **Statistical ecology, epidemiology etc.**

  • $F$ often considered <span style="color:red">having high dimension</span> (even "infinite")

  • Model with relevant bias-variance compromise for a training set of size $n$

# Remark 2)

The Kennard-Stone algorithm does not follow the previous hypothesis, e.g.

**Westad, F., Marini, F., 2015**. Validation of chemometric models – A tutorial. *Analytica Chimica Acta* 893, 14–24. https://doi.org/10.1016/j.aca.2015.06.056

# More representative

- **Uniform random sampling**  Unbiased for $F$ but variable

- **Representative stratified sampling**  Less variable
  - Clustering and then proportional intra-cluster sampling

- **Duplex (Snee 1977)**  Alernate Kennard-Stone

- **Latin Hypercube sampling (LHS)**  Uniform over the margins

# Data separation for this example:   uniform sampling

$n$ = **1206**

- **Training set**          1006 samples

- **Test**          200 samples (15%)
                    **randomly selected**

**PCA**

y distribution

# Cross-validation

1) **K-Fold CV**    ← <span style="color:red">In the present example</span>

  – Can be repeated

2) **Test-set (or Monte Carlo) CV**

  – Alternative when K-Fold not possible (large datasets)
  – In general, returns models slightly more parcimonious
  – All samples are not seen in the VAL

# K-Fold CV



$$Performance = \frac{1}{5} \sum_{i=1}^{5} Performance_i$$

http://ethen8181.github.io/machine-learning/model_selection/model_selection.html

**Under hypothesis** $F_{\text{Train}} = F_{\text{Future}}$

$\text{Err} =$ **Expected prediction (or test) error** (See Appendix)

CV is a non parametric method for estimating Err

- $\widehat{\text{MSEP}}_{\text{CV}} = \widehat{\text{Err}}$

# Particular K-Fold CV

- Leave-One-Out    LOO-CV

  $K = n$

## Within the K-Fold strategy

- LOO-CV $(K = n)$ is the <span style="color:red">less biased</span> for Err

  but <span style="color:red">can have high variability</span> with the training $\tau$

(\*) If other training sets $\tau$ are generated from distribution $F$, LOO-CV can return quite different models

$K = n$



**MSEP_CV**

**Relative gain**

```
ncomp nbpred  msep   mad rmsep   sep      b    r2  cor2
  32    1006 5.441 1.555 2.333 2.333 -0.005 0.939 0.939
```

**Consensus in the literature     (for low to moderate biases)**

- $K$ = 5-10

- $K = n$       opt = <span style="color:red">32</span>       Wold1% = <span style="color:red">8</span>

- $K < n$    With 50 repetitions

  - $K = 10$       opt = <span style="color:red">32</span>       Wold1% = <span style="color:red">13</span>

  - $K = 5$       opt = <span style="color:red">32</span>       Wold1% = <span style="color:red">13</span>

  - $K = 2$       opt = <span style="color:red">13</span>       Wold1% = <span style="color:red">13</span>

**For the 50 repetitions, nb. occurrence for** *opt*

# Permutation tests on CV predictions

In the chemometrics literature, one implementation (within many other possible) is given in

**van der Voet, H., 1994**. Comparing the predictive accuracy of models using a simple randomization test. _Chemometrics and Intelligent Laboratory Systems_ 25, 313–323. https://doi.org/10.1016/0169-7439(94)85050-X

Test if the squared CV residuals for $a$ components and for opt components have the same distribution

$\rightarrow$ Non parametric permutation tests for matched pairs

(eventually with randomization)

**Squared CV residuals (a = 32)**

opt

= reference

**Squared CV residuals (a = 5)**

Example with the **Sign test** for $K = 5$ (with 50 repetitions)



Sign test

Other usual tests for matched pairs

• Wilcoxon signed-rank test

• Randomized permutation test (this method is used in **van der Voet 1994**)

# Alternative    =    F-tests

- Asumme Gaussian residuals

- Require calculating the good *df*s!!

# Examination of the PLS scores

# based on training set

PLSR screeplot

one row in $T$

Standart deviation

Score

Nb. components

# Randomized permutation test of Wiklund et al. 2007

- $\mathrm{Corr}(\boldsymbol{t}_a, \boldsymbol{y})$

Compare the observed distribution with a Null ($H_0$) distribution
$H_0$: $\boldsymbol{y}$ is randomly permuted    ($\boldsymbol{y}$ "scrambling")

Long calculation time for large data
since conditional permutations for each component (successive PLS1)
$a = 1, ..., A$

**Wiklund, S., Nilsson, D., Eriksson, L., Sjöström, M., Wold, S., Faber, K., 2007**. A randomization test for PLS component selection. *Journal of Chemometrics* 21, 427–439. https://doi.org/10.1002/cem.1086

# Wiklund et al test

# Another approach
Unconditionnal randomized permutations (much faster)

# Examination of the PLS loadings and b-coefficients

# based on training set

# Loadings

# Loadings norm



Nb. components
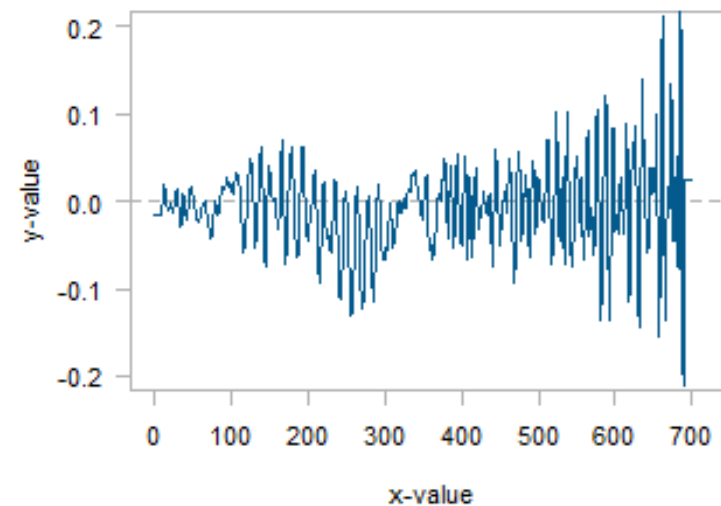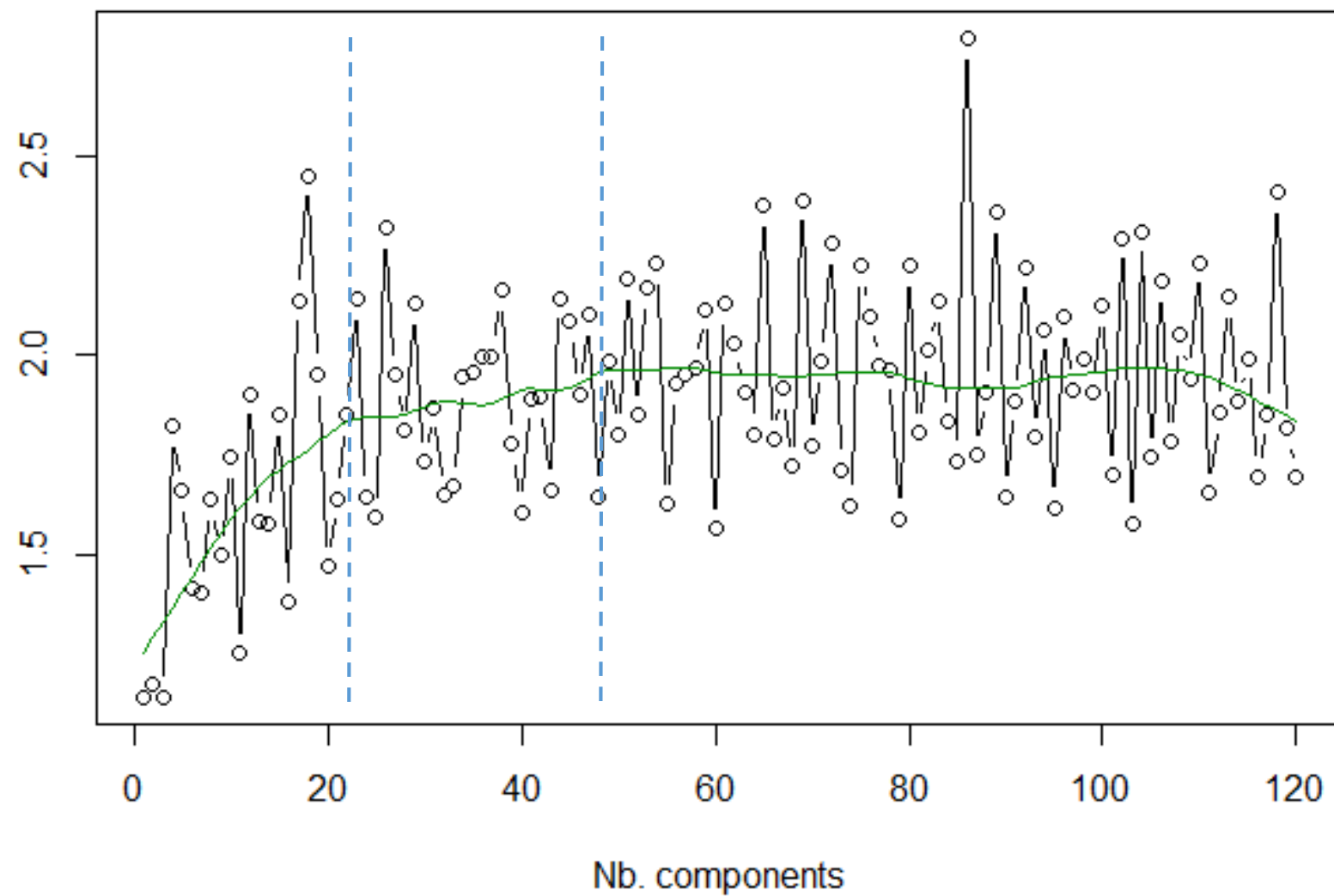
# Stability of loadings when bootstrapping

Used in PCA contexts

- **Ye, Z., Weiss, R.E., 2003**. Using the Bootstrap to Select One of a New Class of Dimension Reduction Methods. *Jasa* 98, 968–979.
https://doi.org/10.1198/016214503000000927


- **Luo, W., Li, B., 2016**. Combining eigenvalues and variation of eigenvectors for order determination. *Biometrika* 103, 875–887.
https://doi.org/10.1093/biomet/asw051

- Training $\tau \rightarrow$ Loadings $\boldsymbol{P}_{\text{obs}}$     $p \times a$   matrix

- Non parametric bootstrap $\tau^{*(b)} \rightarrow \boldsymbol{P}^{*(b)}$     $p \times a$ matrix

$\rightarrow$ Angle $(\boldsymbol{P}_{\text{obs}}, \boldsymbol{P}^{*(b)})$

When the dimension $a$ of $\boldsymbol{P}$ increases,
the last columns of $\boldsymbol{P}*$ become instable with variation of
the bootstraped training $\tau*$     (related to increasing noise)

$b = 1, ..., B \quad \rightarrow \quad$ Mean of the $B$ angles

The mean angle will tend to increase toward $pi / 2$

# How measuring angles between matrices?

- Vector correlation coefficient $q$ (Hotelling 1936) (Ye & Weiss 2003, Luo & Li 2016)

- Maxsub angle (Krzanowski 1979, Hubert et al 2005, Engelen et al. 2005)

- Multivariate coefficient correlation (El Ghaziri, E.M. Qannari 2015)

- Etc.

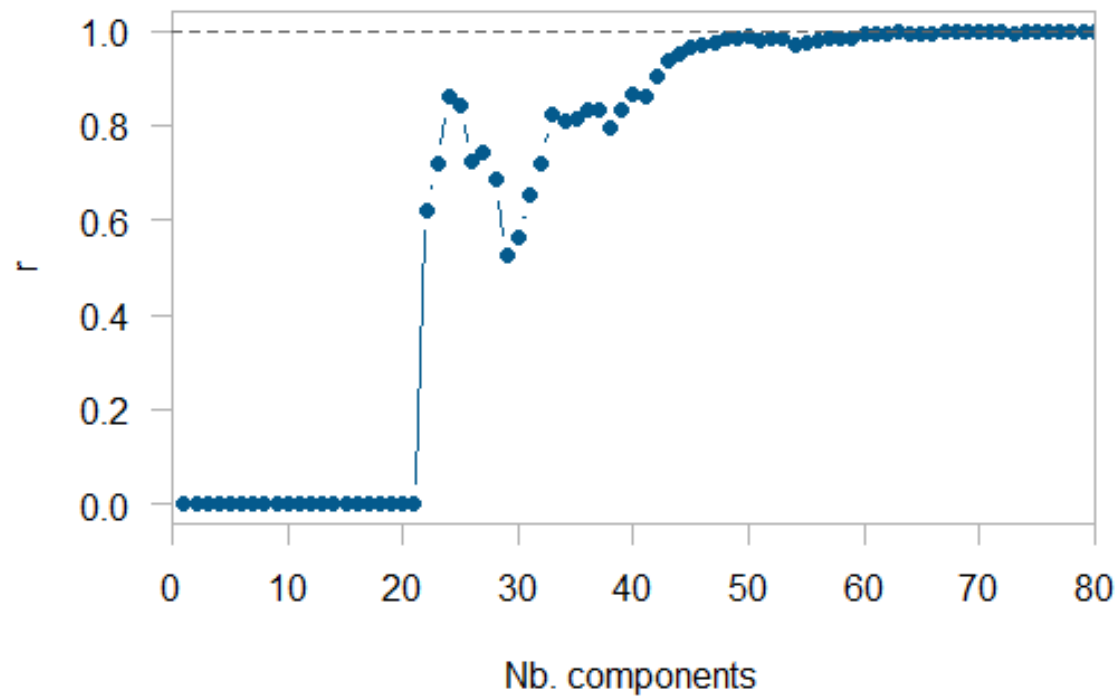**Krzanowski, W.J., 1979**. Between-Groups Comparison of Principal Components. *Journal of the American Statistical Association* 74, 703–707. https://doi.org/10.1080/01621459.1979.10481674

**Engelen, S., Hubert, M., Branden, K.V., 2005**. A Comparison of Three Procedures for Robust PCA in High Dimensions. *Austrian Journal of Statistics* 34, 117–126–117–126. https://doi.org/10.17713/ajs.v34i2.405

**Hubert, M., Rousseeuw, P.J., Vanden Branden, K., 2005**. ROBPCA: A New Approach to Robust Principal Component Analysis. *Technometrics* 47, 64–79. https://doi.org/10.1198/004017004000000563
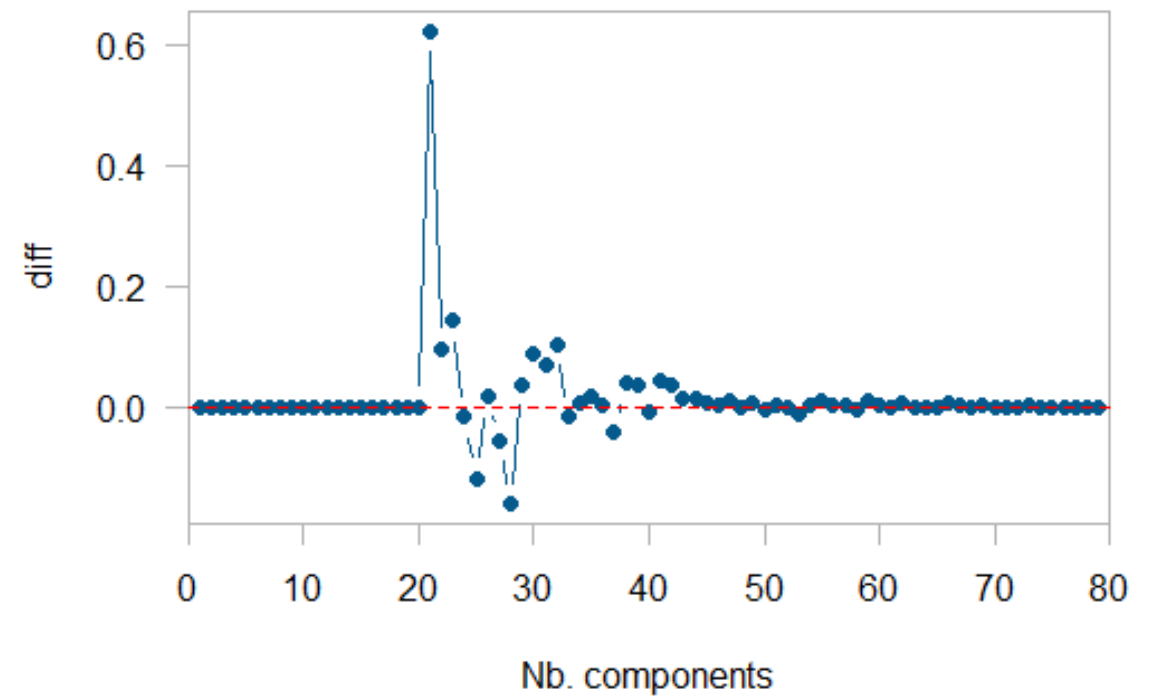
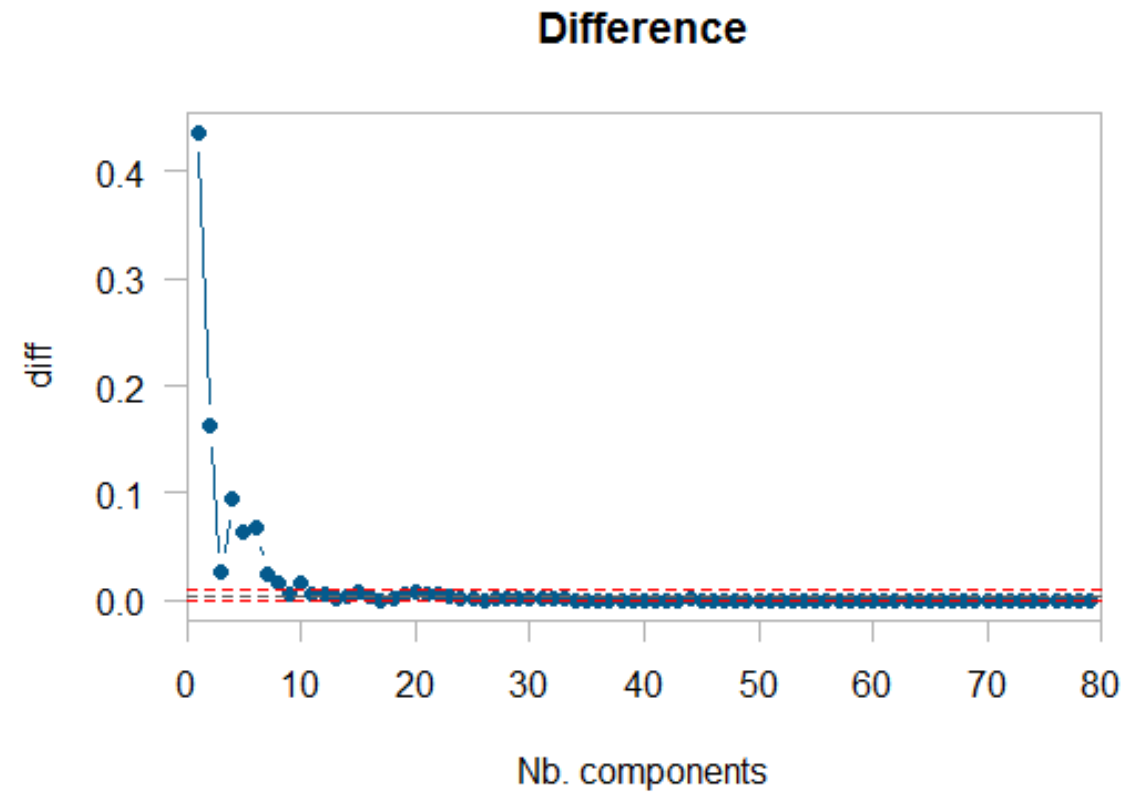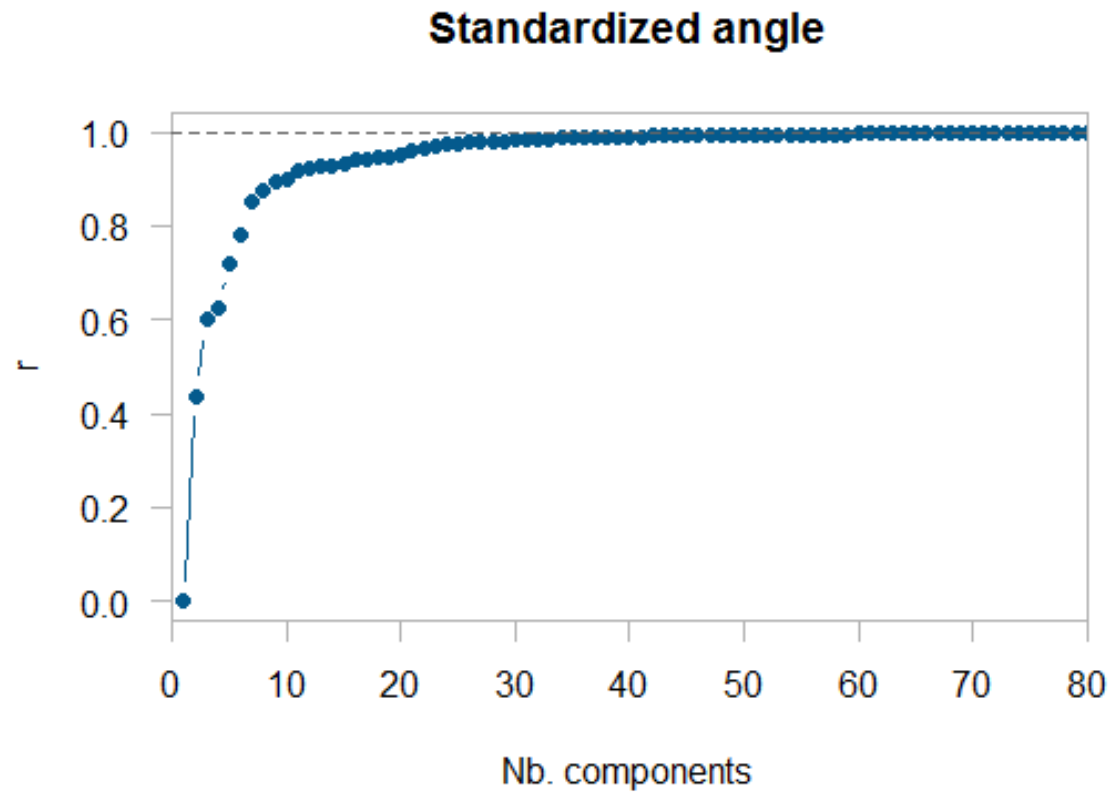**El Ghaziri, A., Qannari, E.M., 2015**. Measures of association between two datasets; Application to sensory data. *Food Quality and Preference* 40, 116–124. https://doi.org/10.1016/j.foodqual.2014.09.010
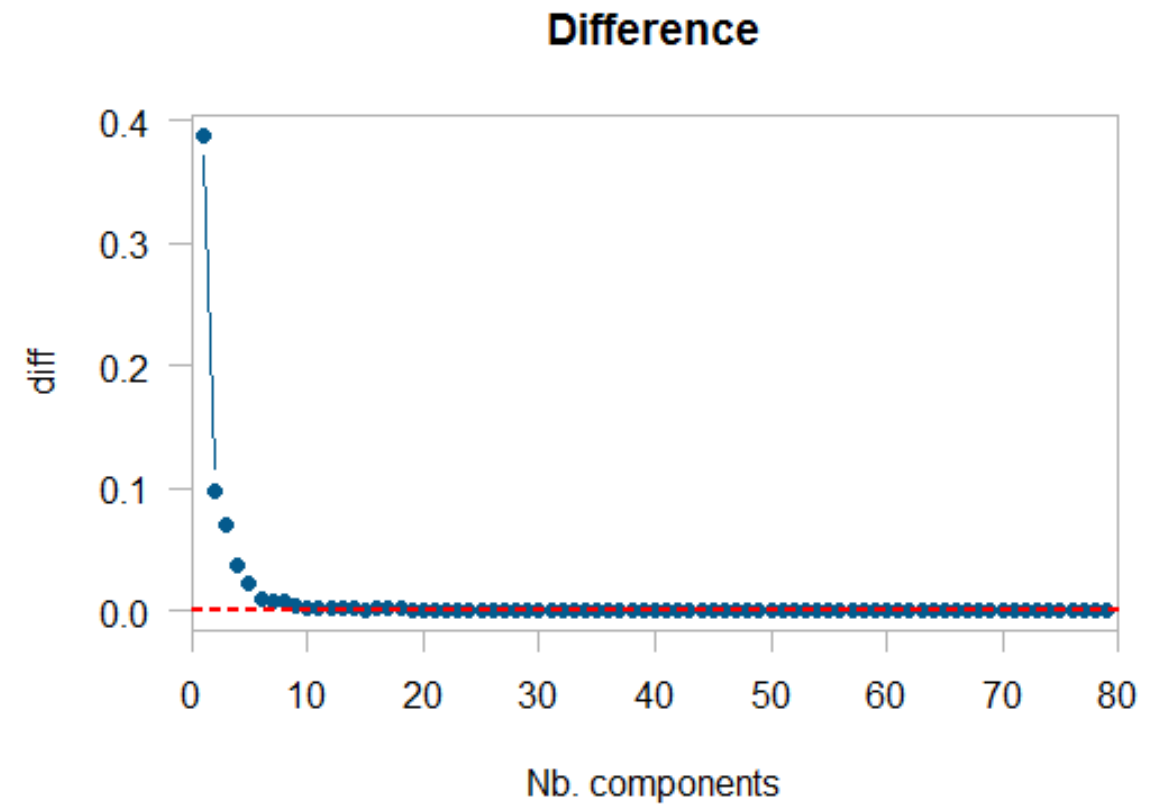
# Hotelling $q$

# Maxsub

# Mult corr

## Another possible method!! $\quad$ SVD( $[\boldsymbol{p}_a^{*(1)}, ..., \boldsymbol{p}_a^{*(B)}]$ ) $\quad \lambda_1 / \text{sum}(\boldsymbol{\lambda})$

## 1 – Colinearity index



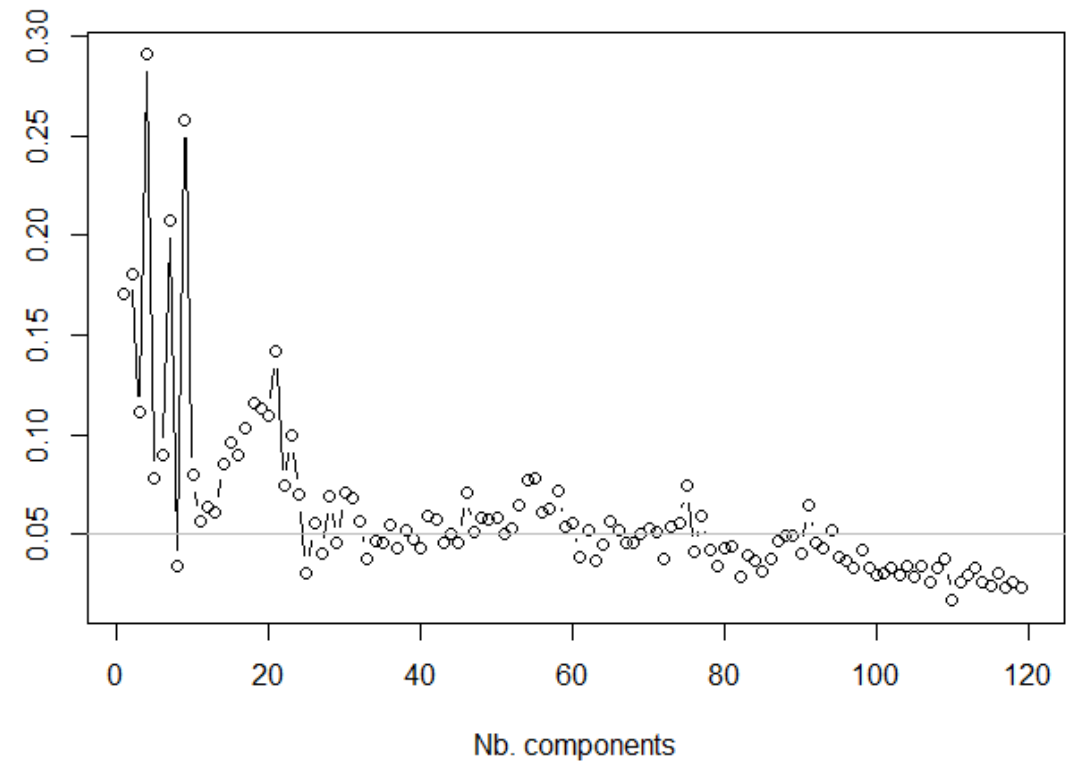**Standardized angle**

**Difference**

# *b*-Coefficients



Norm of PLSR b-coeff
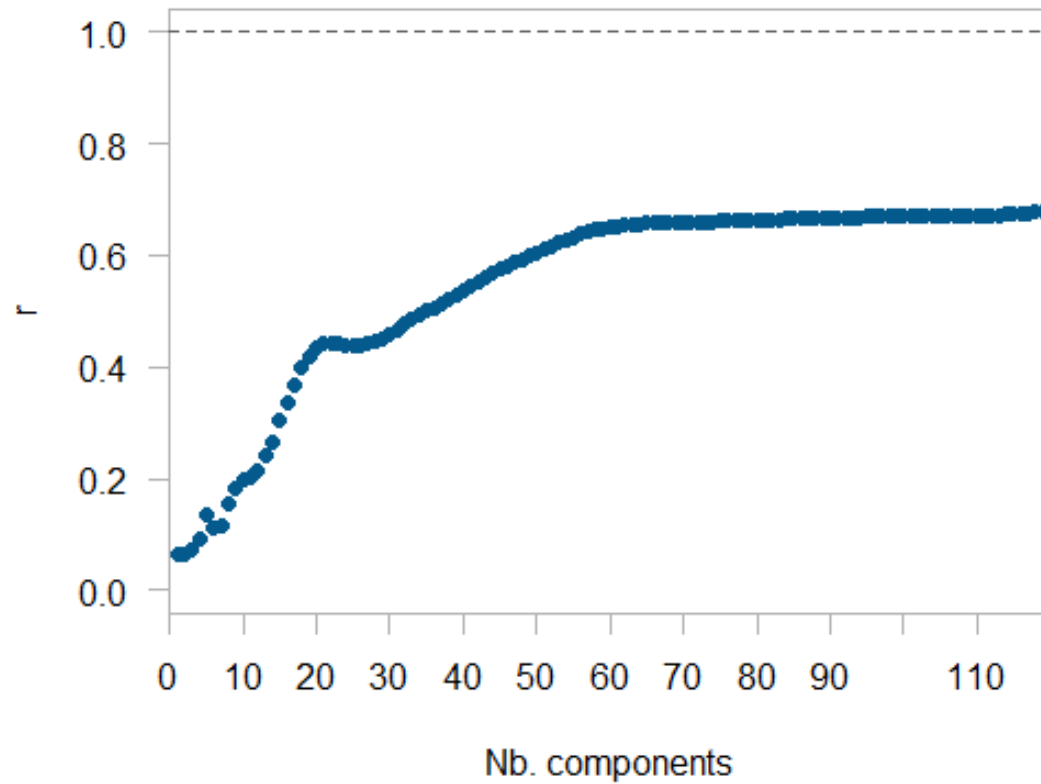
Difference(log(Norm))

# Colinearity method $(\lambda_1 / \mathrm{sum}(\lambda))$ as previous but now on the b-coefficient vector

# Gaussian $b$?
## Shapiro-Wilk statistic



Shapiro-Wilk statistic

# QQplot method



QQnorm plot for $b_{a=5 \text{ LV}}$

→ Linearity index in QQnorm plots

# Autocorrelated $b$? (order 1)



Durbin-Watson statistic



Difference

On this example DW seems very not informative

# Cp

- $\quad \text{Cp} = \dfrac{\text{SSR}}{n} + \dfrac{2}{n}\sum_{i=1}^{n}\hat{C}ov_{\varepsilon}(y_i, \hat{y}_i)$

30 replications bootsrap

**AIC**

Opt = 32

**Relative gain**

# AICC



Opt = 24

# Relative gain

**BIC**

Opt = 13

Value — Nb. components

**Relative gain**

R — Nb. components

**MSEP_CV K = 5**

AIC

AICc

BIC

- *CV*
  - *LOO*        opt = 32        Wold1% = 8
  - $K = 10$      opt = 32        Wold1% = 13
  - $K = 5$      opt = 32        Wold1% = 13
  - $K = 2$      opt = 13        Wold1% = 13

- Covariance penalty Cp
  - AIC      opt = 32        Wold1% = 13
  - AICc      opt = 24        Wold1% = 14
  - BIC      opt = 32        Wold1% = 12

# Remark　　　Weights for model averaging (stacking)

AIC weights　　(e.g. Burnham & Anderson 2002, 2004; Hastie et al 2009)

- $\Delta\mathrm{AIC}_a = \mathrm{AIC}_a - \min(\mathrm{AIC})$

- $w_a = \exp(-.5\ \Delta\mathrm{AIC}_a) / \mathrm{sum}(\ \exp(-.5\ \Delta\mathrm{AIC}_a))$

Same for AICc, BIC, $\mathrm{MSEP}_{CV}$, etc.

**Burnham, K.P., Anderson, D.R., 2002.** Model selection and multimodel inference: a practical information-theoretic approach, 2nd ed. *Springer*, New York, NY, USA.
**Burnham, K.P., Anderson, D.R., 2004**. Multimodel Inference: Understanding AIC and BIC in Model Selection. *Sociological Methods & Research* 33, 261–304. https://doi.org/10.1177/0049124104268644

**delta AICC**        **Model weights**

# Results on the test set

# MSEP_TEST

Opt    CV    = 32    [13]
       AIC   = 32    [15]
       AICc  = 24    [14]
       BIC   = 13    [12]

opt = 35

Same as for $a$ = 12

msep

ncomp

**PLSR**

| ncomp | nbpred | msep | mad | rmsep | sep | b | r2 | cor2 |
|---|---|---|---|---|---|---|---|---|
| **35** | 200 | 3.93 | 1.704 | 1.982 | 1.978 | 0.127 | 0.947 | 0.947 |

**PCR**

| ncomp | nbpred | msep | mad | rmsep | sep | b | r2 | cor2 |
|---|---|---|---|---|---|---|---|---|
| **79** | 200 | 3.891 | 1.8 | 1.973 | 1.971 | 0.069 | 0.947 | 0.947 |

# Sensitivity to the test set

**Hypothesis**    Future $\neq$    Mechanism (probability distribution) generating the training

- **Training set**    $F \rightarrow \boldsymbol{\tau} = \{ (\boldsymbol{x}_1 , y_1) , \ldots , (\boldsymbol{x}_n , y_n) \}$

- **New observation**    $F^* \textcolor{red}{(\neq F)} \rightarrow (\boldsymbol{x}^*, y^*)$

**Example 1**

**Data split by the <span style="color:red">Kennard-Stone algorithm</span> (inverted)**

- **Training set**        Lower dispersion    $n = 1006$

- **New observation**       Higher dispersion   $m = 200$

**PCA**

Score distance

Orthogonal distance

Group
- Train
- Test

# MSEP_TEST



Opt
| | | | |
|---|---|---|---|
| CV (5) | = | 30 | [12] |
| AIC | = | 25 | [14] |
| AICc | = | 25 | [10] |
| BIC | = | 12 | [10] |

```
ncomp nbpred    msep    mad rmsep    sep      b    r2  cor2
   25       200 10.314 2.182 3.212 3.207 -0.18 0.889 0.891
```

ncomp

**Example 2**

**Data split based on the** <span style="color:red">**Score-Orthogonal distance**</span>

- **Training set**       Lower distances   $n = 1156$

- **New observation**       Higher distances  $m = 50$

**PCA**

**MSEP_TEST**

| Opt | CV | = 29 | [13] |
| | AIC | = 29 | [13] |
| | AICc | = 26 | [13] |
| | BIC | = 13 | [11] |

```
ncomp nbpred   msep    mad rmsep    sep      b    r2  cor2
  10        50 7.408 2.102 2.722 2.711 -0.241 0.942 0.942
```

# Some conclusion points

1) Cp (AIC, **etc.) was consistent** with CV and other preliminary analyses

2) Also with results on the **test sets**

3) For the sorghum data, **if $n$ increases**

- the variance-bias compromise will probably allow to increase the nb. components

- Guess: opt would probably increase up to around 70-80 (if one believes the graphical methods on loadings and b-coefficents)

4) **Choice between** AIC vs BIC penalty

# 5) Cp **requires Monte Carlo for estimating** $df_{\mathrm{PLSR}}$

- Could we use a rule of thumb $df_{\mathrm{approx}}$ for escaping simulations?

- This would be a very fast procedure for approximating what returns $\mathrm{MSEP_{CV}}$ selection

**AICC**

Opt = 25

**Relative gain**

**BIC**

Opt = 16

**Relative gain**

with $df_{\mathrm{approx}} = 2.8\,a + 1$

**But the ratio pattern can vary with the data …**
Ex: Cassava data ($n = 200$)

6) **Generating distributions**

– If hypothesis $F = F^*$ is not accepted,

one should specify $F^*$

and optimize the model for this specific future

(personal opinion: there is no generic rule for model selection if $F^*$ is not specified)

# APPENDIX – Some details

# Statistical model

**"True" generating distribution $F$**

- **Training set**        $F \rightarrow \boldsymbol{\tau} = \{ (\boldsymbol{x}_1, y_1), \dots, (\boldsymbol{x}_n, y_n) \}$

- **New observation**     $F \rightarrow (\boldsymbol{x}^*, y^*)$

**Two sources of variations**    $\boldsymbol{\tau}$ and $(\boldsymbol{x}^*, y^*)$

**Same distribution** for training set and new observations

# A convenient way is to consider a "true" generating model $g$

- **Training set** $\qquad\qquad y \mid \boldsymbol{x} \qquad = g(\boldsymbol{x}, \gamma) + \varepsilon$

- **New observation** $\qquad y^* \mid \boldsymbol{x}^* \quad = g(\boldsymbol{x}^*, \gamma) + \varepsilon^*$

Irreducible error

$\varepsilon \qquad$ iid $\qquad E(\varepsilon){=}0 \qquad\qquad Var(\varepsilon) = \sigma^2$

$\varepsilon^* \qquad$ iid $\qquad E(\varepsilon^*){=}0 \qquad\quad Var(\varepsilon^*) = \sigma^2$

$g \quad$ unknown (and will stay unkown)

**Prediction model** $f(\boldsymbol{x}, \boldsymbol{\theta})$     Ex:     $f$ = PLSR model
$$f \neq g \text{ (unknown)}$$

- Training set     $\tau \;\; \rightarrow \;\; \widehat{\boldsymbol{\theta}}$

- Predictions     $\hat{y} \mid \boldsymbol{x} = f(\boldsymbol{x}, \widehat{\boldsymbol{\theta}})$

  $\hat{y}^* \mid \boldsymbol{x}^* = f(\boldsymbol{x}^*, \widehat{\boldsymbol{\theta}})$

- Residual $\qquad e\,|\,\boldsymbol{x}_i = y\,|\,\boldsymbol{x}_i - f(\boldsymbol{x}_i, \hat{\theta})$ $\qquad$ calibration error

- Prediction error $\qquad e^*\,|\,\boldsymbol{x}^* = y^*\,|\,\boldsymbol{x}^* - f(\boldsymbol{x}^*, \hat{\boldsymbol{\theta}})$ $\qquad$ non-observable

$\qquad = (g(\boldsymbol{x}^*, \gamma) + \varepsilon^*) - f(\boldsymbol{x}^*, \hat{\boldsymbol{\theta}})$ $\qquad$ two sources of variation

# Prediction errors
# for model selection

**Notations used in**

**Hastie, T., Tibshirani, R., Friedman, J., 2009**.
The elements of statistical learning: data
mining, inference, and prediction, 2nd ed.
*Springer*, New York.

# Conditional extra-sample error $\mathrm{Err}_\tau$

*Other names:* **Conditional generalization (or test) error**

For a quadratic loss:

$$\mathrm{Err}_\tau = E_{\boldsymbol{x}^*, y^*} \left( \{ y^* | \boldsymbol{x}^* - f(\boldsymbol{x}^*, \widehat{\boldsymbol{\theta}}) \}^2 \mid \tau \right)$$

Future
(distribution $F$)

New
observation

Fit
based
on $\tau$

Training

**Ideally for model selection**, one would expect comparing the models based on their $\text{Err}_\tau$

**Problem:** $\text{Err}_\tau$ is very difficult to estimate from the training data $\tau$

In general, $\mathrm{Err}_\tau$ is estimated *a posteriori* **from a test set**
(data not used in the training!)

- $\widehat{\mathrm{MSEP}}_{\mathrm{Test}} = \widehat{\mathrm{Err}}_\tau$

- Can not be used for model selection

**Note:** Double repeated CV expects to estimate both the errors $\mathrm{Err}_\tau$ and $\mathrm{Err}$
(see thereafter) in the same time

**Instead of $\mathrm{Err}_\tau$ , statistical methods for model selection target two other error measures, e.g.**

    1) CV, Bootstrap                                 $\mathrm{Err}$

    2) Covariance penalty criteria (Cp, AIC, BIC, etc.)       $\mathrm{Err}_{\mathrm{in}}$

Both measures are derived from $\mathrm{Err}_\tau$

1) **Expected prediction (or test) Error**    $\mathrm{Err} = E_\tau(\mathrm{Err}_\tau)$

- CV, Bootstrap    $\widehat{\mathrm{MSEP}}_{\mathrm{CV}}$ , $\widehat{\mathrm{MSEP}}_{\mathrm{Boot}} = \widehat{\mathrm{Err}}$

2) **Conditional in-sample Error**    $\mathrm{Err}_{\mathrm{in}} = \sum_{i=1}^{n} \mathrm{Err}_\tau(\boldsymbol{x}_i)/n$

- Simplif. 1: Plug-in of $\mathrm{Err}_\tau$ on the training set $\boldsymbol{\tau}$
- Simplif. 2: $\boldsymbol{\tau}$ variations comes only from $\boldsymbol{\varepsilon}$ (design $\boldsymbol{X}$ assumed fixed)

- Covariance penalty criteria    Cp, etc.  =  $\widehat{\mathrm{Err}}_{\mathrm{in}}$

**Note:** $\mathrm{Err}_{\mathrm{in}}$ is noted "Err" in Efron Jasa 2004 (!!)

Both Err and $\text{Err}_{\text{in}}$ have been shown

- effective for model selection

- estimable from the training set $\tau$

- compromises between variance and bias

# Variance-bias compromise for $\mathrm{Err}$

$$\mathrm{Err} = \sigma^2 + \mathrm{Var}_\tau(\widehat{\boldsymbol{y}}) + \mathrm{Bias}_\tau(\widehat{\boldsymbol{y}})^2$$

$\uparrow$ when model complexity $\uparrow$

$\downarrow$ when model complexity $\uparrow$

Model selection expects to find a compromise well afforded by the training of size $n$

# Variance-bias compromise for $\text{Err}_{\text{in}}$

$E_{\varepsilon}\,(\text{Err}_{\text{in}})\ =\text{MSEP}\qquad =\frac{1}{n}\,\sum_{i=1}^{n}\text{MSEP}(\pmb{x}_i)$ $\qquad\qquad$ Only $\pmb{\varepsilon}$ varies (X assumed fixed)

$\qquad\qquad = \sigma^2 + \frac{1}{n}\,\sum_{i=1}^{n}\text{MSE}(\pmb{x}_i) = \sigma^2 + \text{MSE}$

$\qquad\qquad = \sigma^2 + \frac{1}{n}\sum_{i=1}^{n} Var_{\varepsilon}(f(\pmb{x}_i,\widehat{\pmb{\theta}})) + \frac{1}{n}\sum_{i=1}^{n}\alpha(\pmb{x}_i)^2$

$\qquad\qquad\quad$ with $\alpha(x_i)= g(\pmb{x}_i,\,\gamma) - E_{\varepsilon}\,(f(\pmb{x}_i,\widehat{\pmb{\theta}})) = -\,\text{Bias}_{\varepsilon}(f(\pmb{x}_i,\widehat{\pmb{\theta}}))$

$\qquad = \overline{\sigma}_*{}^2 + \frac{1}{n}\,\pmb{\alpha}'\pmb{\alpha}$

$\qquad\qquad\quad$ with $\ \overline{\sigma}_*{}^2\ =\frac{1}{n}\sum_{i=1}^{n}\sigma_*{}^2(\pmb{x}_i) =\sum_{i=1}^{n}\big(\sigma^2 + Var_{\varepsilon}(f(\pmb{x}_i,\widehat{\pmb{\theta}})\big)$

# Hastie *et al* 2009   Fig 7.1



Err

$$\overline{err} \quad = \text{SSR}/n$$
$$= \text{MSEC}$$

# Important relation for $\text{Err}_{\text{in}}$

Expected **model optimism** $\omega$

Taking the expectation over $\boldsymbol{\tau}$ for $\boldsymbol{x}_1, \dots, \boldsymbol{x}_n$ fixed

$$E_{\boldsymbol{\varepsilon}}(\text{Err}_{\text{in}}) = E_{\boldsymbol{\varepsilon}}\left(\frac{\text{SSR}}{n}\right) + \frac{2}{n}\sum_{i=1}^{n} Cov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i)$$

$\boldsymbol{\tau}$

= MSEP

= MSEC

$f(\boldsymbol{x}_i, \widehat{\boldsymbol{\theta}})$

**Efron, B., 2004**. The Estimation of Prediction Error. *Journal of the American Statistical Association* 99, 619–632. https://doi.org/10.1198/016214504000000692

$$\omega = \frac{2}{n} \sum_{i=1}^{n} Cov_{\varepsilon}(y_i, \widehat{y}_i)$$

Higher is the effect of a variation of $y_i$ on its prediction $\widehat{y}_i$

higher is $\omega$

For Gaussian errors $\boldsymbol{\varepsilon}$     (Stein 1981, Efron 2004)

$$Cov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i) \;=\; \sigma^2 E_{\boldsymbol{\varepsilon}}\left(\frac{\partial \hat{y}_i}{\partial y_i}\right)$$

$$\Rightarrow \;\; \omega = \frac{2}{n}\sum_{i=1}^{n} Cov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i) = \frac{2}{n}\sigma^2 \sum_{i=1}^{n} E_{\boldsymbol{\varepsilon}}\left(\frac{\partial \hat{y}_i}{\partial y_i}\right)$$

**Stein, C.M., 1981**. Estimation of the Mean of a Multivariate Normal Distribution. _The Annals of Statistics_ 9, 1135–1151.

# Covariance penalty criteria

$$E_{\boldsymbol{\varepsilon}}(\mathrm{Err}_{\mathrm{in}}) = E_{\boldsymbol{\varepsilon}}\left(\frac{\mathrm{SSR}}{n} + \frac{2}{n}\sum_{i=1}^{n} Cov_{\boldsymbol{\varepsilon}}(y_i, \hat{y}_i)\right)$$

$\Rightarrow$ **Natural unbiased estimator** for $\mathrm{Err_{in}}$

(in the sense $\mathrm{E}_{\boldsymbol{\varepsilon}}(\widehat{\mathrm{Err}}_{\mathrm{in}}) = \mathrm{E}_{\boldsymbol{\varepsilon}}(\mathrm{Errin})$ )

$$\widehat{\mathrm{Err}}_{\mathrm{in}} = \frac{\mathrm{SSR}}{n} + \frac{2}{n}\sum_{i=1}^{n}\hat{C}ov_{\varepsilon}(y_i, \hat{y}_i)$$

MSEC

Covariance penalty

$\rightarrow$ = Covariance penalty criteria

$\rightarrow$ = Mallows' Cp (or "AIC", Hastie et al 2009 p.231) family criteria

# Using Stein's equality when assuming Gaussian error

$$\widehat{\mathrm{Err}}_{\mathrm{in}} = \frac{\mathrm{SSR}}{n} \quad + \quad \frac{2}{n}\hat{\sigma}^2 \sum_{i=1}^{n} \frac{\partial \hat{y}_i}{\partial y_i}$$

Divergence *i*

$\rightarrow$ = Stein Unbiased Risk Estimate (SURE) for $\mathrm{Err}_{\mathrm{in}}$

**Can be estimated using Monte Carlo simulations**

- $\widehat{\text{Err}}_{\text{in}} = \dfrac{\text{SSR}}{n} \quad + \quad \dfrac{2}{n}\sum_{i=1}^{n}\hat{C}ov_{\varepsilon}(y_i, \hat{y}_i)$

  Parametric bootstrap
  (See e.g. Efron 2004)

- $\widehat{\text{Err}}_{\text{in}} = \dfrac{\text{SSR}}{n} \quad + \quad \dfrac{2}{n}\hat{\sigma}^2\sum_{i=1}^{n}\dfrac{\partial \hat{y}_i}{\partial y_i} \quad \text{(SURE)}$

  Sensitivity analysis
  (perturbations)

# Model's degrees of freedom

# Simple case:     Linear smoothers

- $\widehat{y} = \boldsymbol{S}\,\boldsymbol{y}$       $\color{red}{\boldsymbol{y}\ \text{not involved in}\ \boldsymbol{S}}$

- Ex:   OLS regression     $\widehat{y} = \boldsymbol{S}\,\boldsymbol{y} = \boldsymbol{H}\,\boldsymbol{y} = \boldsymbol{X}(\boldsymbol{X'X})^{-1}\boldsymbol{X'}\,\boldsymbol{y}$

Then

$$\sum_{i=1}^{n} Cov_{\varepsilon}(y_i, \hat{y}_i) = tr(2\boldsymbol{S} - \boldsymbol{SS}')\sigma^2$$

If $\boldsymbol{S}$ is idempotent $p \times p$

$$\sum_{i=1}^{n} Cov_{\varepsilon}(y_i, \hat{y}_i) = tr(\boldsymbol{S})\sigma^2 = p\sigma^2$$

- $\widehat{\mathrm{Err}}_{\mathrm{in}} = \dfrac{\mathrm{SSR}}{n} \quad + \quad \dfrac{2}{n} p \hat{\sigma}^2$        Mallows' Cp

- $\widehat{\mathrm{Err}}_{\mathrm{in}} = \dfrac{\mathrm{SSR}}{n} \quad + \quad \dfrac{\log(n)}{n} p \hat{\sigma}^2$        BIC ("Cp" form)

**For orthogonal linear smoother such as OLS regression, it is accepted that** $df = p$

If $df$ is defined by

$$df = \sum_{i=1}^{n} Cov_{\varepsilon}(y_i, \hat{y}_i)/\sigma^2$$

then

$$df = (p\sigma^2)/\sigma^2 = p$$

which is consistent with the usual accepted $df$

**This has led to the general definition**

For any model (linear or not), the consensus is to consider that

$$df = \sum_{i=1}^{n} Cov_{\varepsilon}(y_i, \hat{y}_i)/\sigma^2$$

First we need to define precisely what we mean by the degrees of freedom of an adaptively fitted model. Suppose we have an additive-error model, with

$$y_i = f(x_i) + \epsilon_i, \ i = 1, \ldots, N, \tag{2.16}$$

for some unknown $f$ and with the errors $\epsilon_i$ iid $(0, \sigma^2)$. If the $N$ sample predictions are denoted by $\widehat{\mathbf{y}}$, then we define

$$\mathrm{df}(\widehat{\mathbf{y}}) := \frac{1}{\sigma^2} \sum_{i=1}^{N} \mathrm{Cov}\left(\widehat{y}_i, y_i\right). \tag{2.17}$$

The covariance here is taken over the randomness in the response variables $\{y_i\}_{i=1}^N$ with the predictors held fixed. Thus, the degrees of freedom corresponds to the total amount of *self-influence* that each response measurement has on its prediction. The more the model fits—that is, adapts—to the data, the larger the degrees of freedom. In the case of a fixed linear model, using $k$ predictors chosen independently of the response variable, it is easy to show that $\mathrm{df}(\widehat{\mathbf{y}}) = k$ (Exercise 2.7). However, under adaptive fitting, it is typically the case that the degrees of freedom is larger than $k$.

As for covariance penalty, *df* can be estimated from parametric bootstrap or sensitivity analysis

- $\widehat{df} = \sum_{i=1}^{n} \hat{C}ov_{\varepsilon}(y_i, \hat{y}_i)/\hat{\sigma}^2$    **Efron Jasa 2004**

- $\widehat{df} = \frac{1}{n}\sum_{i=1}^{n}\frac{\partial \hat{y}_i}{\partial y_i}$    (SURE)    **Ye Jasa 1998**    *Generalized df*

**Ye, J., 1998**. On Measuring and Correcting the Effects of Data Mining and Model Selection. *Journal of the American Statistical Association* 93, 120–131. https://doi.org/10.1080/01621459.1998.10474094

The GDF is an extension of (7) to general modeling procedures. It is defined to be the sum of the *average* sensitivities of the fitted value $\hat{\mu}_i(\mathbf{Y})$ to a small change in $y_i$. Thus it measures the flexibility of the modeling procedure $\mathcal{M}$. If $\mathcal{M}$ is highly flexible, then the fitted values tend to be close to the observed values. Thus the sensitivity of the fitted values to the observed values would be high, and the GDF would be large.