

Size Dependence of Validation in Linear Modeling

Péter Király
2nd year PhD report
Supervisor: Gergely Tóth

May 21, 2020

- Sample size dependence of validation parameters (VPs)
 - studied on multiple linear regression (MLR) and multi-target partial least squares regression (PLS2) models
 - Scramble or randomize?
 - Leave-one-out (LOO) or leave-many-out (LMO) CV?
 - R^2 is superior to CCC (concordance correlation coefficient)
 - Assessment of predictivity on 2D Roy-Ojha-Kovács diagrams
 - Rank correlation of VPs

Internal validation

Y_{sample}

X_{sample}

goodness of fit, robustness

External validation

$$Y = X\mathbf{B} + E$$

Y_{test}

 Y_{train}



$\mathbf{B}_{\text{train}}$

X_{test}

 X_{train}

predictivity

Validation parameters (VPs)

internal VPs

$$R^2 = \frac{1}{d} \sum_{i=1}^d 1 - \frac{\text{RSS}_i}{\text{TSS}_i} \quad \bigg| \quad R_{\text{inc}}^2 = 1 - \frac{\sum_{i=1}^d \text{RSS}_i}{\sum_{i=1}^d \text{TSS}_i}$$

$$Q_{\text{LOO}}^2 = \frac{1}{d} \sum_{i=1}^d 1 - \frac{\text{PRESS}_i}{\text{TSS}_i}$$

$$\text{CCC} = \frac{1}{d} \sum_{i=1}^d \frac{2\sigma(y_i, \hat{y}_i)}{\text{RSS}_i + \text{MSS}_i + n_i (\bar{y}_i - \hat{\bar{y}})^2}$$

$$\text{RMSE} = \frac{1}{d} \sum_{i=1}^d \sqrt{\frac{\text{RSS}_i}{n_i}} \quad \dots$$

external VPs

$$Q_{F1}^2 = \frac{1}{d} \sum_{i=1}^d 1 - \frac{\text{PRESS}_{\text{test},i}}{\text{TSS}_{\text{test},i} (\bar{y}_{\text{train},i})}$$

$$Q_{F2}^2 = \frac{1}{d} \sum_{i=1}^d 1 - \frac{\text{PRESS}_{\text{test},i}}{\text{TSS}_{\text{test},i}}$$

$$Q_{F3}^2 = \frac{1}{d} \sum_{i=1}^d 1 - \frac{\text{PRESS}_{\text{test},i}/n_{\text{test},i}}{\text{TSS}_{\text{train},i}/n_{\text{train},i}}$$

⋮

Sample size dependence of VPs

Datasets

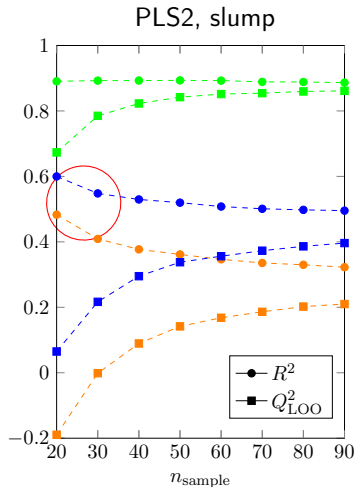
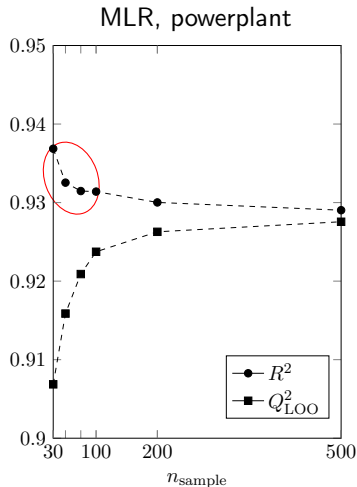
dataset	targets	features	instances
powerplant	1	4	9568
concrete	1	8	1030
grammatical1-3	1-1-1	4	369
temperature	2	26(1)	400
slump	3	7(7)	103
air pollution	6	38(19)	363

Specialities for multi-target sets

- sets are modeled with *multi-target partial least squares regression* (PLS2)
- predictor and response matrices standardized
- number of latent variables is set by *repeated double cross validation* (rdCV), R^2

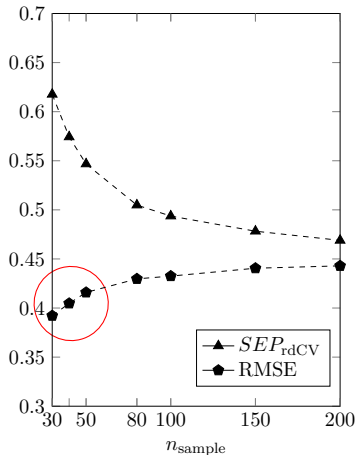
Tasks for size dependence study

- random sampling from population
- *training/test* split of sample for external validation (80/20 for training)
- linear modeling and computation of *validation parameters* (VPs)
- tasks above repeated $1000\times$ for given sample size, VPs averaged

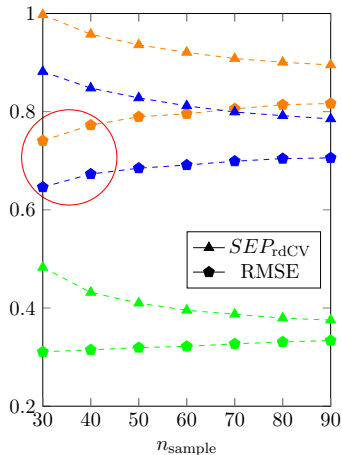


Is the smaller model the better one?

MLR, gramatica3



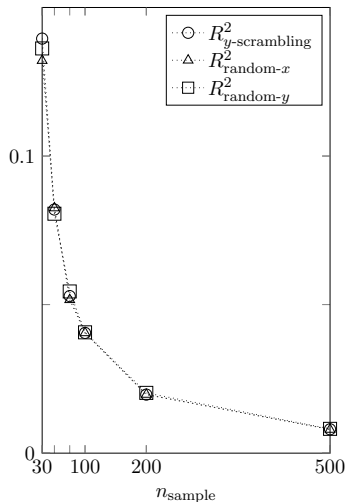
PLS2, slump



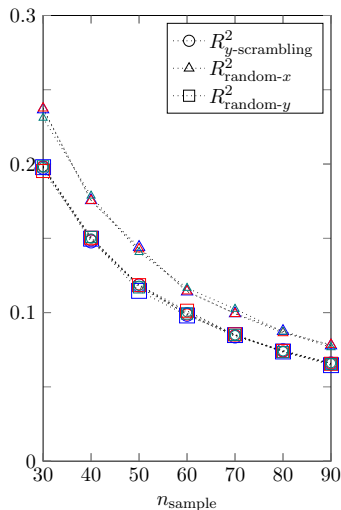
Is the smaller model the better one?

Internal VPs - scrambling, randomization

MLR, powerplant

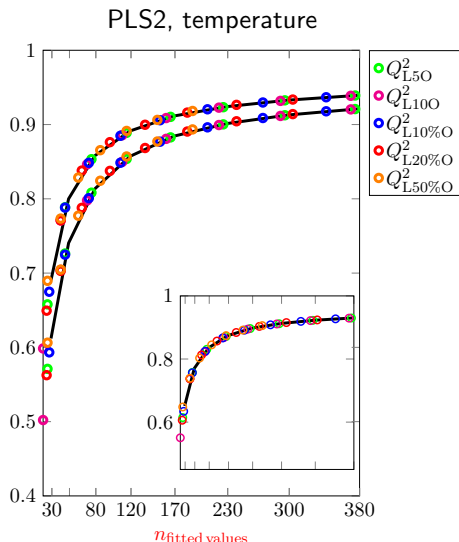
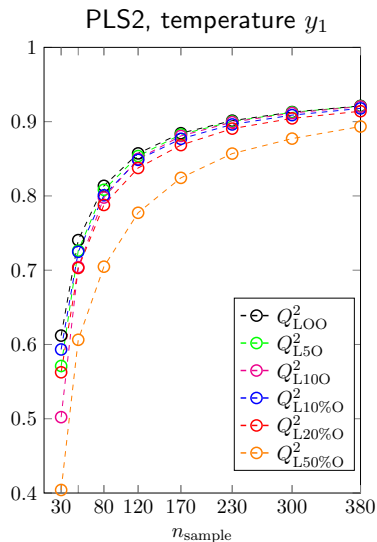


PLS2, slump



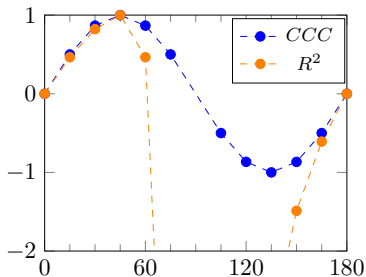
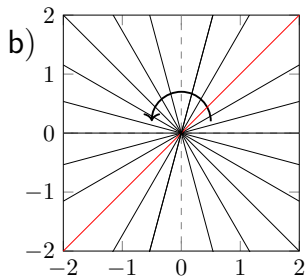
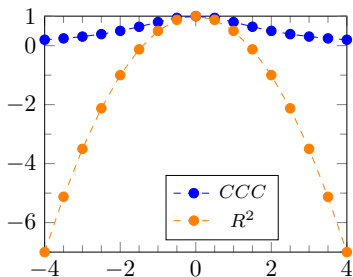
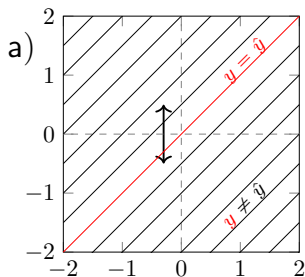
y-scrambling easiest to calculate

Internal VPs - leave many out



Calculation of Q^2_{LOO} is straightforward

R^2 superior to CCC

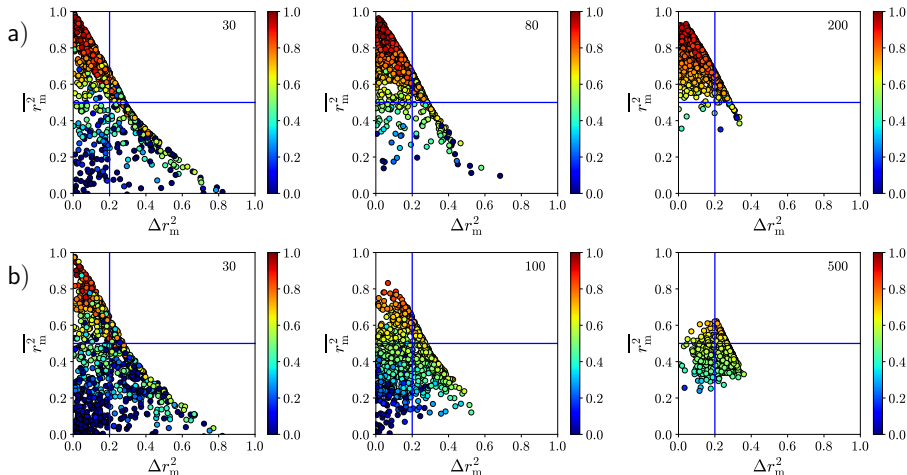


Roy-Ojha-Kovács diagram (test set)

a) MLR, gramatica3

b) MLR, concrete

} models colored according to $Q_{F2}^2(R_{\text{test}}^2)$



Rank correlations

$$\underbrace{R^2, CCC}_{\text{internal (goodness-of-fit)}}$$

internal (goodness-of-fit)

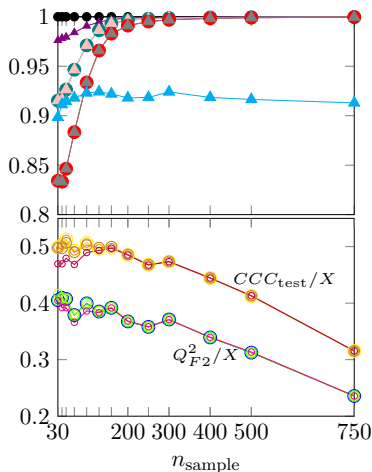
$$\underbrace{Q^2_{\text{LOO}}, CCC_{\text{LOO}}}_{\text{internal (robustness)}}$$

internal (robustness)

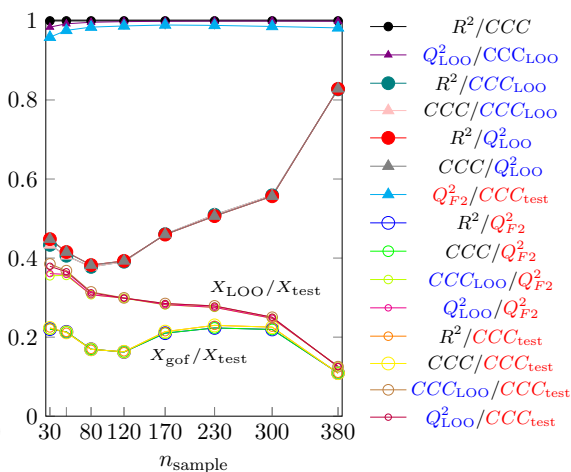
$$\underbrace{Q^2_{F2}, CCC_{\text{test}}}_{\text{external (predictivity)}}$$

external (predictivity)

MLR, concrete



PLS2, temperature mean



- R^2/CCC
- $Q^2_{\text{LOO}}/CCC_{\text{LOO}}$
- R^2/CCC_{LOO}
- CCC/CCC_{LOO}
- R^2/Q^2_{LOO}
- CCC/Q^2_{LOO}
- $Q^2_{F2}/CCC_{\text{test}}$
- R^2/Q^2_{F2}
- CCC/Q^2_{F2}
- $CCC_{\text{LOO}}/Q^2_{F2}$
- $Q^2_{\text{LOO}}/Q^2_{F2}$
- R^2/CCC_{test}
- CCC/CCC_{test}
- $CCC_{\text{LOO}}/CCC_{\text{test}}$
- $Q^2_{\text{LOO}}/CCC_{\text{test}}$

- model validation: goodness of fit (R^2), robustness (Q_{LOO}^2), predictivity (Q_{F2}^2)
- sample size dependence of R^2 , RMSE show anomaly
- scrambling can be used instead of randomization
- Q_{LOO}^2 is sufficient to use instead of LMO variants
- model predictivity is categorized on Roy-Ojha-Kovács diagrams
- internal and external schemes have different information content for all sample sizes
- Q_{F1-3}^2 for external validation show different sensitivity for different sample size
- rdCV redesigned for PLS2

We thank Prof. Imre Salma for providing us the air pollution dataset.

Co-workers: Ramóna Kiss, Dániel Kovács, Gergely Tóth