Adjusted Pareto scaling for optimum prediction performance of chemometric multivariate models

Kurt Varmuza^{1,2} and Peter Filzmoser¹

TU Wien - Vienna University of Technology

- [1] Institute of Statistics and Mathematical Methods in Economics | Computational Statistics
- [2] Institute of Chemical, Environmental and Bioscience Engineering | Sustainable Technologies and Process Simulation

Multivariate models are routinely applied for the determination of a property y of samples using m measurements made on the samples.

Widely used are linear regression models $\hat{y} = \hat{b}^{T} x$, with \hat{y} for the predicted property, $\hat{\boldsymbol{b}}$ the vector with estimated regression coefficients, and x the centered vector with m measurements (variables).

Adjusted Pareto scaling $x_{\text{PARETO}} = x_{\text{C}}/s^{P}$

- $x_{\rm C}$ Centered variable; based on arithmetic mean or median.
- Spread measure of variable; classical standard deviation (sd) or s from IQR (interquartile range) or another robust measure.

Ρ **Pareto exponent** 0 ... 1, step for instance 0.1

P = 0, no scaling; P = 0.5, standard Pareto scaling; P = 1, autoscaling

The performance of such chemometric models often depends on the applied method of scaling the x-variables.

Autoscaling and Pareto scaling are widely used [1].

Here an *adjusted Pareto scaling* is presented - continuously covering the range from no scaling via standard Pareto scaling to autoscaling.

Calibration and evaluation

- PLS regression is used together with the strategy repeated double Cross Validation (rdCV) [2].
- Performance of models is characterized by the prediction errors for test set samples, obtained separately from model optimization.

• Performance measure used is SEP (standard error of prediction), the standard deviation of prediction errors, estimated from several repetitions of CV (typ. 50 values give a boxplot). The range ±2 SEP defines a ~95% confidence interval for predictions of y [3].

Experimental

Pareto exponent varied 0, 0.1, 0.2, ..., 1 SEP 50 values estimated by rdCV (50 repetitions with different random splits), presented as boxplot.

A Number of PLS components (model complexity) estimated by rdCV (well fitted, mid plot), underfitted (top plot), and over-fitted (bottom plot).

rdCV: 3 segments for test/calibration, 7 segments for optimization of A (double cross validation) Scaling: median for center; 0.7413*IQR for spread Software: R programming environment [4,5]



Conclusions based on the results in these three examples – and an outlook

Scaling of x-variables by methods based on the spread (s) may improve the performance of multivariate regression models. However, the effect has to be tested, and depends on the data. No simple general rules, related to data properties, seem to be evident.

Standard Pareto scaling $x_{PARETO} = x_C / s^P$ with P = 0.5 may be not optimal. Varying the Pareto exponent P between 0 (no scaling) and 1 (autoscaling) is recommended. A dependence of the model performance from P is clear for underfitted models, but may be low for well-fitted or over-fitted models. Other (robust) scaling methods, based on the variable spread, like range scaling or vast scaling. may give better model performances for some data sets than adjusted Pareto scaling (work in progress).

Walach J., et al.: In Jaumot J., et al.: Comprehensive analytical chemistry. Data analysis for omics sciences: Methods and applications, Elsevier, Amsterdam, p. 165-196, 2018
 Filzmoser P., et al.: J. Chemometrics, 23, 160 (2009)
 Varmuza K., Filzmoser P.: Introduction to multivariate statistical analysis in chemometrics, and the provided and provided analysis in chemometrics,

- CRC Press, Boca Raton, FL, USA, 2009 [4] R, A language and environment for statistical computing, R Development Core Team,
- Foundation for Statistical Computing, http://www.r-project.org, Vienna, Austria, 2023
- [5] Filzmoser P., Varmuza K.: R package chemometrics, Vienna, Austria, 2010, http://cran.at.r-project.org/web/packages/chemometrics/index.html
 [6] Friedl A., et al.: Anal. Chim. Acta, 544, 191 (2005)
 [7] Liebman B., et al.: Biochem. Eng. J., 52, 187 (2010)
 [8] Lee M.L., et al.: Anal. Chem., 51, 768 (1979)
 [9] Varmuza K., et al.: Comput. Struct. Biotechn. J., 5 [6], e201302007, 1 (2013)
 [10] Todeschini B., et al.: Molecular descriptors for chemoinformatics. Wiley-W

[10] Todeschini R., et al.: Molecular descriptors for chemoinformatics, Wiley-VCH, Weinheim, 2009