

Partial least squares

Menu:	QCExpert	Predictive methods	Partial least squares
-------	----------	--------------------	-----------------------

Module PLS regression provides the user with one of the best computational tools for evaluating a pair of multidimensional variables, which is expected to have linear relationship inside one or the other multidimensional variable, and linear relationship between the two variables with each other. This computationally intensive methodology allows to explain and predict one of the variables using other group of variables. The PLS regression method found a large number of applications in the planning and management of quality in manufacturing technology, design and optimization of the characteristics of products in the development of new products, marketing studies, research in the evaluation of experiments, in clinical trials. An example might be modeling the relationship between technological parameters in the production and product quality parameters, or between the chemical composition and physical and biological characteristics. The typical questions of technological practice, which PLS can often answer include:

- It has a purity of the raw material any effect on the strength of the product?
- What happens if the temperature is increased in the process?
- Can we increase the stability of the product by reducing the speed or rotation?
- Which process parameters affect the most product strength?
- How to set the value of procedural parameters to achieve the desired product characteristics?
- What caused the decrease in the parameter?
- In what and how subsequent production batches differ?
- How to improve the stability / quality?
- How to increase the strength / value / competitiveness?
- Which input parameters are crucial for the quality?
- Which process parameters are crucial for the quality?

Mathematical basics of the PLS regression method

Let us denote $\mathbf{X}(n \times p)$ the matrix (table) of measured values of p variables (columns) with n lines and denote $\mathbf{Y}(n \times q)$ the corresponding table with the same number of lines n but with q variables. Center all columns (subtract column average from each column). To extract maximum information from the p - q - dimensional matrices to a lower dimension space, we decompose \mathbf{X} and \mathbf{Y} to the product of the orthogonal matrices $\mathbf{T}(n \times k)$ and $\mathbf{U}(n \times k)$, with coefficient matrices \mathbf{P} and \mathbf{Q}

$$\begin{aligned}\mathbf{X} &= \mathbf{TP}^T + \mathbf{E} \\ \mathbf{Y} &= \mathbf{UQ}^T + \mathbf{F}\end{aligned}$$

while maximizing the correlation between \mathbf{T} and \mathbf{U} . Required dimension k , $1 < k \leq \min(p, q)$ is chosen by user, for example, on the basis of the squares sum decrease (scree plot), see below. Noise and irrelevant information contained „litter“ in every measured data is swept into residual matrices \mathbf{E} and \mathbf{F} . Decomposition $\mathbf{U} = \mathbf{TB}$ (where \mathbf{B} is a square diagonal matrix) give us a tool for computing (estimating) \mathbf{Y} from \mathbf{X} but also \mathbf{X} from \mathbf{Y} , just by switching the \mathbf{X} and \mathbf{Y} data because the model PLS-R is symmetric.

$$\hat{\mathbf{Y}} = \mathbf{TBQ}^T,$$

\mathbf{T} is calculated from the new data \mathbf{X} , $\mathbf{T} = \mathbf{XP}^-$ (\mathbf{P}^- indicate generalized Moore-Penrose pseudoinversion of a rectangular matrix \mathbf{P}). Furthermore, there is an internal link between \mathbf{X} and \mathbf{Y} . By writing $\mathbf{W} = \mathbf{BQ}^T$, we can rewrite the original pair of relations in the form

$$\begin{aligned}\mathbf{X} &= \mathbf{TP}^T + \mathbf{E} \\ \mathbf{Y} &= \mathbf{TW} + \mathbf{F}\end{aligned}$$

so that the data \mathbf{X} and \mathbf{Y} are linked through a common scores matrix \mathbf{T} , which is actually orthogonalized original matrix \mathbf{X} in generally smaller number of dimensions, with a maximum extracted information contained in the original \mathbf{X} , with removed noise (which is moved to the matrix \mathbf{E}), while the maximum covariance with the matrix \mathbf{Y} is ensured. Using the relation

$$\mathbf{A} = \mathbf{P} \mathbf{B} \mathbf{Q}^T$$

we can reconstruct coefficients of a classical regression model with multivariate response $\mathbf{Y} = \mathbf{A}\mathbf{X}$. Columns \mathbf{a}_i of \mathbf{A} contain linear coefficients (absolute terms are zero thanks to centering data) of the models $\mathbf{y}_i = \mathbf{X}\mathbf{a}_i$, where \mathbf{y}_i is i -th column of matrix \mathbf{Y} . The coefficients are not usually fully numerically identical to coefficients obtained by classical linear regression. They are generally biased, but shrunk, which means that they have lower variances, and are generally more stable.

As mentioned above, this method is looking for a relationship between the two phenomena described by multidimensional numerical vectors. A typical example is \mathbf{X} matrix containing measured technological parameters in the production of individual units or batches and matrix \mathbf{Y} containing the relevant physical parameters of finished products, their deviations from the specifications, etc. Another example is a matrix \mathbf{X} containing climate and chemical descriptions of the various sites and \mathbf{Y} matrix with biological parameters of micro-organisms, vegetation and fauna in these locations. There are many other applications in geology, biology, toxicology, chemistry, medicine, psychiatry, behavioral sciences, pharmacology, cosmetics, food, steel industry to name just a few. With PLS prediction we can then obtain estimates of unknown quantities \mathbf{Y} on the basis of known values \mathbf{X} .

Model validation

The prediction quality of a particular PLS model can be assessed on the basis of its ability to predict the value of \mathbf{y} from the value of \mathbf{x} . This is used in various validation procedures, sometimes called cross-validation. The principle of validation of the model is the same as in the case of neural networks. We „hide“ part of the data before computation of the PLS model. This hidden data are called test or validation data. For the rest of the data, called training data, we calculate parameters of the PLS model. Then the validation are „unhidden“ and used to and check whether the model correctly predicts validation data. Validation must have the same nature and range of values \mathbf{x} , and therefore the same model as the data used for training. For the validation data we then construct diagnostic charts, which simply conclude whether the model is appropriate for all data. If the model describes well only the training data and not validation data, this usually means that we have little data (rows), or that we have chosen is too large proportion of validation data. A proportion between 10 and 40% of the validation data is usual.

Of course, even advanced PLS regression method is not miraculous and has certain restrictions, which is mainly assumption of linearity of all relationships and normality of error distribution. Along with the ability of prediction and graphical diagnostics, however, it provides a very powerful tool for analysis and prediction of multidimensional variables. In quality control, thanks to its prediction capability, PLS regression is an ideal tool for the quality planning, design of products, optimization of technologies and applied research.

Data and parameters

Module PLS regression needs two data tables, matrix \mathbf{X} with p columns and \mathbf{Y} with q columns selected as the dialog box items *Matrix X* and *Matrix Y*, see Fig. 1. The matrix columns must contain numeric data only, the number of rows must be the same for both \mathbf{X} and \mathbf{Y} . Each matrix must contain at least two columns. Columns of matrix \mathbf{X} must not appear in the matrix \mathbf{Y} . The limiting dimension k can be set by user. If the box *Dimension* is not checked, the maximum dimension is set to $k = \min(p, q)$. It is recommended to perform PLS in maximal dimension first, then optionally we can determine an appropriate value of k using the scree plot (see paragraph Graphical output below) and repeat the

computation again with new k . If the checkbox *Connect Biplot* is checked, consecutive points of the Biplot will be connected in the order of data in spreadsheet. This can help to follow a possible trajectory of the process. If the checkbox *X-Prediction* is checked, it is necessary to choose the same number of columns as in the field *Independent variable X*. These variables will be used to compute the predicted values of the dependent variable. The X for the prediction must have the same number of columns as the independent variable matrix X , but may have a different number of lines (at least 2 lines). A typical example of input matrices and data for prediction is given on Fig. 2.

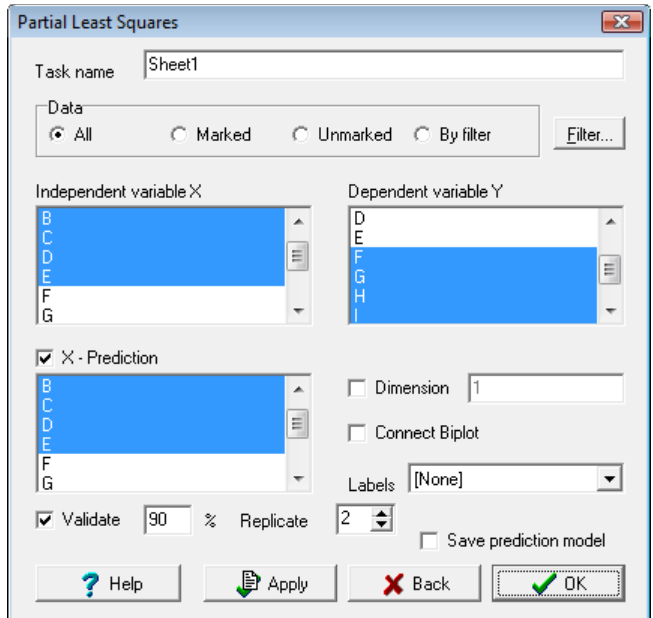


Fig. 1 Dialog box for PLS regression

X ($n \times p$)	Y ($n \times q$)	X for prediction ($n_1 \times p$)
<pre> 0.41957 0.17436 0.66634 0.08041 0.63552 0.78197 0.384139 0.518022 0.046317 0.00993 0.895669 0.73931 0.148208 0.813652 0.708706 0.089473 0.523150 0.089165 0.41656 0.24622 0.3246 0.398454 0.289135 0.81532 0.46589 0.43877 0.38323 0.625482 0.253899 0.265526 0.041343 0.244389 0.343002 0.034870 0.531379 0.795859 0.879075 0.470874 0.688881 0.29432 0.52163 0.143877 0.827791 0.23048 0.659749 0.508735 0.317254 0.965795 0.287525 0.948185 0.991974 0.382204 0.768407 0.168896 0.778956 0.177135 0.338466 0.952114 0.624442 0.054468 0.186456 0.952228 0.361623 0.004229 0.501774 0.787195 0.461342 0.052546 0.786247 0.947771 0.742377 0.270243 0.386353 0.011079 0.303803 0.723919 0.832865 0.151775 0.774244 0.470488 0.579035 0.912297 0.881181 0.481162 0.796337 0.704171 0.87371 0.027992 0.167786 0.038761 0.683003 0.271599 0.862357 0.001443 0.57207 0.390827 0.912076 0.62066 0.508872 0.711773 0.534762 0.108233 0.961997 0.881641 0.9115 0.333276 0.877707 0.953954 0.461549 0.881932 0.862015 0.678119 0.919633 0.09148 0.637207 0.22063 0.062526 0.882136 0.437341 0.108867 0.332447 0.898979 0.636402 0.311754 0.848994 0.819589 0.731862 0.572933 0.800603 0.372942 0.674465 0.606622 0.313375 0.732101 0.883259 0.683365 0.745699 0.641006 0.119153 0.411284 0.178413 0.428442 0.588339 0.702752 0.940482 0.887129 0.538465 0.229779 0.843725 0.638897 0.894809 0.134846 0.634876 0.008897 0.096223 0.482486 0.691526 0.521136 0.288826 0.21891 0.128328 0.377603 0.518915 0.516567 0.938186 0.138659 0.101258 0.68274 0.646466 0.903459 0.078803 0.119693 0.603076 0.934846 0.316316 0.142384 0.932677 0.174841 0.327421 0.898946 0.071027 0.643903 0.898957 0.252485 0.048652 0.086438 0.386959 0.203895 0.652022 0.811387 0.380691 0.209899 0.868808 0.383776 0.71162 0.007885 0.088165 0.403377 0.863116 0.183863 0.053725 0.254237 0.658447 0.738345 0.843084 0.100515 0.37384 0.07614 0.788695 0.261827 0.398673 0.161171 0.817737 0.206701 0.270314 0.708894 0.289775 0.264204 0.32061 0.14389 0.338965 0.03007 0.430713 0.192943 0.043976 0.107245 0.465384 0.848733 0.147961 0.650191 0.638742 0.087267 0.110787 0.206888 0.710395 0.882889 0.48886 0.231446 0.472584 0.232383 0.591558 0.044796 0.887559 0.410267 0.440395 0.486071 0.3047 0.665359 0.689781 0.038413 0.308811 0.041401 0.555135 0.888465 0.919953 0.097316 0.62282 0.28913 0.81987 0.194749 0.511953 0.089379 0.217046 0.62328 0.825541 0.846223 0.841471 0.164885 0.487345 0.43092 0.276475 0.11732 0.634624 0.300355 0.768264 0.419584 0.418424 0.91146 0.658605 0.302522 0.048149 0.378785 0.541781 0.038563 0.683707 0.369605 0.121518 0.084378 0.644105 0.977254 0.316469 0.480253 0.711488 0.488693 0.548594 0.118045 0.533381 0.548461 0.024603 0.27952 0.688805 0.954876 0.153879 0.441086 0.219264 0.684501 0.019389 0.619342 0.048188 0.772905 0.884643 0.438963 0.183421 0.017334 0.289669 0.988392 0.183486 0.765208 0.321553 0.180121 0.18373 0.622575 </pre>	<pre> 0.9011 0.7586 0.7487 0.8821 0.2387 0.1121 0.1276 0.9637 0.0569 0.5877 0.8882 0.6635 0.6144 0.3198 0.1558 0.6881 0.5292 0.6715 0.7795 0.8879 0.0245 0.1189 0.8221 0.6520 0.3788 0.0105 0.8057 0.8684 0.6037 0.2127 0.0570 0.0481 0.5534 0.7345 0.9608 0.3488 0.6625 0.7112 0.1489 0.8013 0.8157 0.1949 0.4330 0.6883 0.1325 0.1727 0.0966 0.7839 0.9036 0.4127 0.9668 0.8633 0.4193 0.3203 0.7087 0.6845 0.4385 0.7336 0.6205 0.6192 0.4785 0.4806 0.6953 0.6688 0.7502 0.4919 0.6218 0.9048 0.5265 0.3456 0.8864 0.8660 0.3778 0.9136 0.4313 0.8129 0.0885 0.8174 0.4761 0.4204 0.2255 0.8537 0.8913 0.2520 0.6595 0.2562 0.9785 0.8553 0.4041 0.2658 0.7795 0.7146 0.0832 0.7320 0.3238 0.8888 0.2955 0.2398 0.1183 0.2938 0.2185 0.7987 0.8182 0.2311 0.8126 0.1276 0.4324 0.7959 0.6746 0.7437 0.4448 0.6157 0.0861 0.2889 0.2595 0.6244 0.6891 0.7091 0.7434 0.1345 0.1222 0.4363 0.4279 0.7151 0.6998 0.7234 0.8915 0.4348 0.3832 0.3280 0.8491 0.3135 0.6284 0.4668 0.4424 0.1462 0.4433 0.9492 0.7469 0.4715 0.6155 0.1713 0.5995 0.4480 0.0319 0.0498 0.1424 0.7526 0.6726 0.0430 0.9596 0.8844 0.8127 0.9041 0.9458 0.2208 0.4243 0.8260 0.1939 0.0189 0.3805 0.5153 0.4134 0.9489 0.5548 0.2810 0.1878 0.3875 0.0624 0.2313 0.2484 0.3966 0.5436 0.4840 0.1725 0.4667 0.0513 0.3517 0.1666 0.2706 0.2089 0.6938 0.8814 0.0500 0.3889 0.3863 0.4236 0.4438 0.2886 0.2650 0.5937 0.2040 0.6808 0.4045 0.1480 0.2775 0.2464 0.7388 0.6426 0.6884 0.4641 0.7389 0.2842 0.7885 0.9789 0.8126 0.9314 0.3849 0.2500 0.9458 0.0005 0.5238 0.8651 0.7640 0.6918 0.6341 0.1806 0.8203 0.3724 0.6476 0.5529 0.3825 0.0008 0.6431 0.7708 0.3846 0.8809 0.8467 0.4321 0.4165 0.2953 0.9108 0.1148 0.3389 0.0128 0.3429 0.8917 0.8340 0.3487 0.9839 0.9331 0.5234 0.2907 0.3156 0.2338 0.9109 0.3021 0.2537 0.9999 0.7875 0.1380 0.4026 0.2953 0.8454 0.6252 0.9310 0.8847 0.8835 0.2581 0.1540 0.1048 0.3675 0.6487 0.7020 0.6586 0.8987 0.8038 0.7029 0.2084 0.4806 0.2185 0.9962 0.9381 0.3302 0.3623 0.7209 0.5342 0.4340 0.2009 0.8548 0.5294 0.8851 0.4286 0.8887 0.3520 0.8913 0.7310 0.8814 0.6287 0.9831 0.6520 0.0040 0.7756 0.3048 0.7423 0.6250 0.0542 0.8078 0.4747 0.3061 0.1804 0.7103 0.3887 0.9717 0.0310 0.1740 0.7540 0.2578 0.1489 0.5246 0.1911 0.9772 0.4566 0.7259 0.6586 0.8204 0.0740 0.0207 0.6917 0.7512 0.1619 0.3211 0.3807 0.4021 0.9789 0.1514 0.5388 0.5161 0.0869 0.8348 0.1111 0.8428 0.3777 0.6434 0.9931 0.8249 0.7143 0.2100 0.8005 0.8748 0.2847 0.3088 0.3286 0.5305 0.4925 0.0195 0.0165 0.6189 0.9447 0.6844 0.7995 0.9942 0.7869 0.9559 0.1289 0.1209 0.3274 0.5401 0.5138 0.7725 0.1788 0.4615 0.6883 0.8772 0.6167 0.2401 0.2165 0.1115 0.2803 0.7658 0.4587 0.0361 0.9538 0.8953 0.4838 0.6333 0.8882 0.2105 0.0114 0.8722 0.4888 0.2730 0.8863 0.5214 0.7411 0.0438 0.9988 0.7180 0.4886 0.8682 0.2855 0.2481 0.1433 0.1808 0.6793 0.7134 0.8893 0.5542 0.8886 0.0755 0.4548 0.2209 0.8583 0.2490 0.9611 0.4182 0.1406 0.1575 0.1386 0.2001 0.0315 0.1482 0.8584 0.5009 0.1125 0.8843 0.2882 0.7654 0.3232 0.7809 0.8653 0.5512 0.2409 0.9003 0.9833 0.2151 0.7562 0.8952 0.1236 0.8025 0.7385 0.8748 0.4575 0.2863 0.5591 0.9220 0.8020 0.3245 0.0340 0.4400 0.3347 0.6895 0.0042 0.7167 0.2449 0.8888 0.7038 0.8884 0.1787 0.3695 0.7615 0.0566 0.9679 0.6693 0.7828 0.1136 0.3210 0.0574 0.9770 0.7103 </pre>	<pre> 0.95966 0.356713 0.486145 0.638897 0.797838 0.678419 0.559 0.42675 0.336552 0.277129 0.891221 0.893862 0.814337 0.136546 0.229739 0.475745 0.488544 0.414508 0.697255 0.862532 0.979597 0.132789 0.4318 0.815738 0.021419 0.676203 0.430445 0.612759 0.920089 0.206251 0.89003 0.321661 0.248884 0.387838 0.380958 0.98113 0.827221 0.388784 0.483895 0.124121 0.201202 0.21384 0.47244 0.221774 0.251334 0.581154 0.386232 0.951575 0.437263 0.669892 0.636684 0.071696 0.356715 0.691153 0.025848 0.046264 </pre>

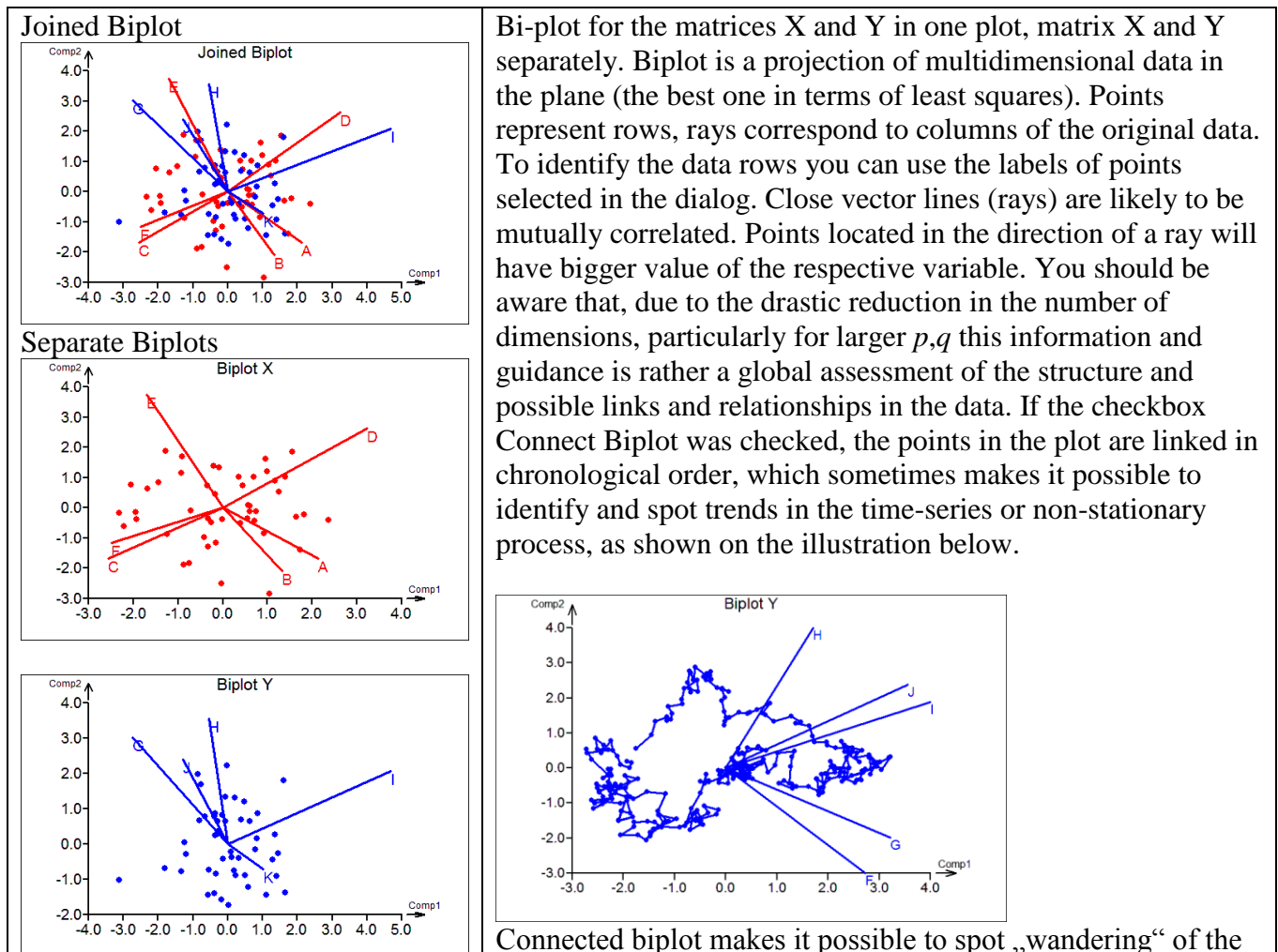
Fig. 2 Typical data for PLS regression

Protocol

Input data	
No of rows	Number of valid rows
No of columns	Number of columns of X a Y matrices.
Columns	Column names of both input matrices.
Chosen dimension	Dimension k for the PLS model chosen in the input dialog window. The dimension must be less or equal to $\min(p, q)$. Scree plot may be used as an

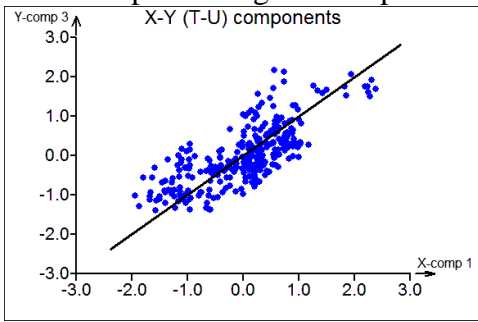
	aid to select suitable k if required.
PLS - coefficients, B	Diagonal elements of the matrix B .
Explained sum of squares	Table of the squares sum of residuals with growing dimension of the model, $i = 1, \dots, k$, these values are used for constructing scree plot.
No of components	Number of components (dimensions) used for the squares sum.
RSS	Residual square sum value, for 0 components the RSS is the total squares sum without a model.
Percent Explained %	% of the RSS (100 – %RSS).
Loadings X, P	Loadings matrix P .
Loadings Y, Q	Loadings matrix Q .
Regression coefficients, A	Matrix of regression coefficients a_{ij} formally similar to those in the separate classical multiple linear regression models $\mathbf{Y} = \mathbf{XA}$, or $\mathbf{y}_j = \sum a_{ij}\mathbf{x}_i$. The coefficient values are generally different from the classical coefficients, since they are based on the orthogonal component regression and therefore they are biased, shortened (with lower standard deviations) and more stable.
Prediction	Predicted values for the data selected in the field X-Prediction in the PLS dialog box. This part of output is not generated unless the checkbox is checked.

Graphs



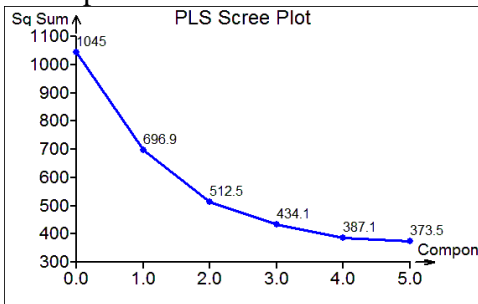
process in time.

X-Y Components agreement plot



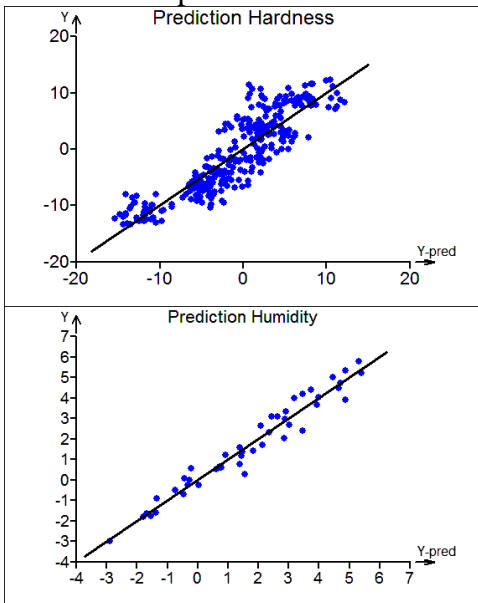
Plot of agreement between the columns of **T** and **U**. This graph shows the global success of a PLS model fitting. The closer the points to the line, the more successful a PLS model is.

Scree plot

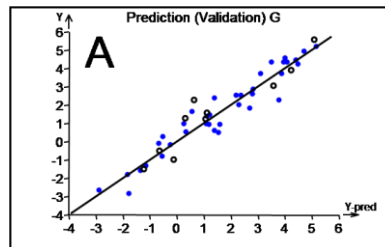


The effectiveness of the model expressed by reduction of unexplained (residual) sum of squares, depending on the number of factors included (columns of matrix **T** and **U**).

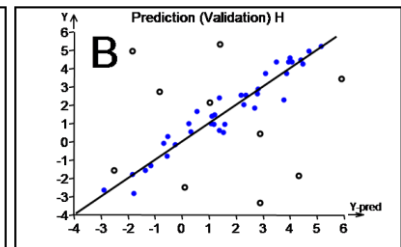
Y-Prediction plot



This plot expresses compliance of dependent variables and model prediction. The closer are the points to the line, the better the fit. This plot is created for each dependent variable. Some variables can be predicted better, others worse. If the plot does not show a visible trend, the suitable model for this variable was probably not found, a model is not able to predict dependent variable. If Validation checkbox was selected, the validation (test) points in the plot are marked in red (in the example below marked empty circles). In the Fig A below both the training and the validation data are fitted well, showing the model is reliable and the dependence is real. However, if the validation data strongly disagree with other data, as in the Fig B below, the PLS model may be over fitted, describing only the training data, and probably is not usable for the prediction of new values. It is advisable to try to reduce the dimension of the model by entering numbers less than $\min(p, q)$ into the field dimension.



Good prediction



Poor prediction

Validation residuals plot

Plot Validation is used to assess the quality of prediction of validation data. Unique very remote points may represent outlying measurements. On the Y-axis are Euclidean distances of the data from the model.

