

1. PLS*-Regression: $Y \sim X$

Here Y is the given response vector ($N \times 1$) and X the given matrix ($N \times p$) of describing variables.

2. Computation of VIP^{**} for each variable X_1, \dots, X_p

3. Order variables X_1, \dots, X_p with respect to VIP (descending). $\rightarrow X_{1:p}, \dots, X_{p:p}$

4. For $j = 1, \dots, p$ do PLS*-Regression: $Y \sim (X_{1:p}, \dots, X_{j:p})$

5. Select number j_{opt} such that the corresponding model $Y \sim (X_{1:p}, \dots, X_{j_{opt}:p})$ has highest R^2 in CV. $\rightarrow X_{new} := (X_{1:p}, \dots, X_{j_{opt}:p})$

6. PLS*-Regression: $Y \sim X_{new}$

7. Generate 5000 different permutations of response Y . $\rightarrow Y_{perm}[1], \dots, Y_{perm}[5000]$

8. Test significance of the correlation $R = \text{cor}(Y, Y_{cv})$ by comparing to $\text{cor}(Y_{perm}[k], Y_{cv})$ for $k = 1, \dots, 5000$:

Compute mean and s of $\text{cor}(Y_{perm}[k], Y_{cv})$, $k = 1, \dots, 5000$. Compute distance of R and mean as measured by s ; estimate P-value, assuming $N(\text{mean}, s^2)$.

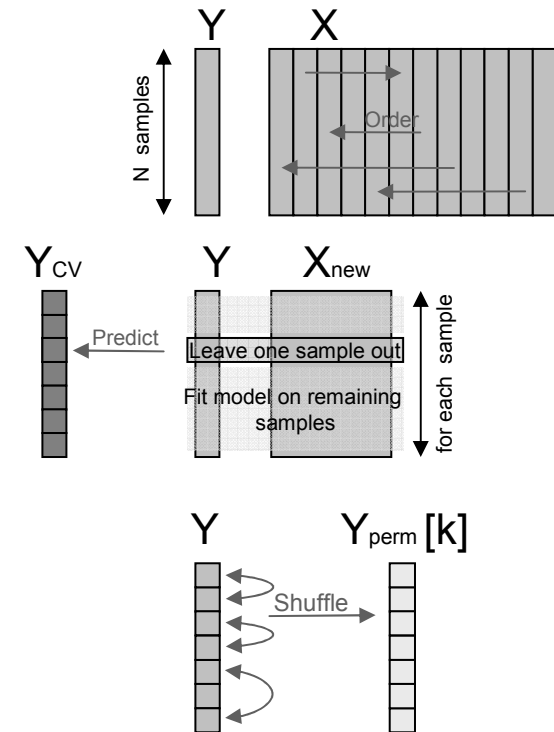
* Partial Least Squares Regression: Projecting describing variables to latent variables T_1, \dots, T_m , such that R^2 is maximized in CV of the OLS-model $Y \sim T$ with respect to m

** Variable Importance in the Projection of X_j : Contribution to Y of those latent variables with high weights in X_j

Fit $Y = c_0 + c_1 \cdot T_1 + \dots + c_m \cdot T_m$

with $T_i = w_{i1}X_1 + \dots + w_{ip}X_p$ for $i = 1, \dots, m$

$$VIP(X_j) = \sqrt{\frac{\sum_{i=1}^m (w_{ij} \cdot \text{cor}(Y, T_i))^2}{\sum_{i=1}^m (\text{cor}(Y, T_i))^2}} \cdot p$$



Abbreviations:

- CV : Cross Validation (Leave-one-out)
- Y_{cv} : Response predicted in CV
- R^2 : Squared correlation of Y and Y_{cv}
- s : Sample standard deviation
- OLS: Ordinary Least Squares