

A primer on Partial Least Squares Regression

Dennis Cook

School of Statistics
University of Minnesota

Collaborating at times with
Liliana Forzani, Inge Helland, and Zhihua Su

PLS Regression \in PLS Family of Methods

Partial Least Squares (PLS), aka Projection onto Latent Structures, and PLS regression are not synonyms.

Herman Wold (1908–1992)

- PhD in 1938, Harald Cramér, advisor
- Contributions: Cramér-Wold theorem, time series, casual modeling, econometrics, . . . originator of PLS family of methods

Wold's algorithms leading to PLS

- **fixed-point algorithm**, early 1960's: iterative OLS regressions to estimate the parameters of multi-equation systems
- **NILES**, Nonlinear Iterative Least Squares: modification for principal components and canonical correlations
- **NIPALS**, late 1960's: Nonlinear Iterative Partial Least Squares: refinement of NILES.
- **PLS-PM**, 1970's: Path Modeling, aka Structural Equation Modeling, via PLS

PLS Regression Algorithm

Developed around 1980 by the Scandinavian chemometrics community – S. Wold, H. Martens & H. Wold – mainly for prediction in high-dimensional multi-response regressions, without requiring that the sample size n be larger than the number of predictors, p .

PLS regression often does better than OLS when $n > p$, as judged empirically by cross validation.

- 1 Cursory introduction to PLS regression
- 2 History of PLS regression
- 3 PLS regression formulation via envelopes, recent model-based dimension-reduction methodology that can greatly reduce estimative and predictive variation relative to standard methods.
- 4 PLS n, p -asymptotics

1. Introduction to PLS regression

Context

Consider the usual linear regression model

$$Y_i = \alpha + \beta^T \mathbf{X}_i + \varepsilon_i, \quad i = 1, \dots, n,$$

- Y univariate (Why?)
- \mathbf{X} : $p \times 1$, $\text{var}(\mathbf{X}) = \Sigma_{\mathbf{X}}$, sample version $\mathbf{S}_{\mathbf{X}}$.
- Errors ε have mean 0 and variance τ^2 , $\varepsilon \perp \mathbf{X}$
- β : $p \times 1$, unknown coefficients
- Goal: estimate β , prediction.
- \mathbf{B} = OLS estimator of β when $n > p$.

PLS regressions – algorithms for reducing the predictor dimension

- 1 Select a matrix $\hat{\Phi} \in \mathbb{R}^{p \times q}$, with $q \ll p$, $q =$ number of components, and then reduce $\mathbf{X} \mapsto \hat{\Phi}^T \mathbf{X}$.
- 2 Fit $Y = \alpha + \eta^T \{\hat{\Phi}^T \mathbf{X}\} + \varepsilon$ using OLS when $n > q$.
- 3 Estimator is of the form

$$\begin{aligned} \hat{\beta} &= \hat{\Phi} \hat{\eta} = \hat{\Phi} (\hat{\Phi}^T \mathbf{S}_X \hat{\Phi})^{-1} \hat{\Phi}^T \mathbf{S}_{XY} \\ &= \mathbf{P}_{\hat{\Phi}(\mathbf{S}_X)} \mathbf{B} \text{ when } n > p \end{aligned}$$

\mathbf{S}_X : Sample variance of \mathbf{X} ; \mathbf{S}_{XY} : sample covariance between \mathbf{X} and Y .

How do we get an estimator $\text{span}(\hat{\Phi})$? PLS regression algorithm.

SIMPLS algorithm (de Jong, 1993), $n > q$

Set $\hat{\boldsymbol{\varphi}}_0 = \mathbf{0}$ and $\hat{\boldsymbol{\Phi}}_0 = (\hat{\boldsymbol{\varphi}}_0)$. For $k = 0, \dots, q-1$, set

Version I

$$\hat{\boldsymbol{\varphi}}_{k+1} = \arg \max_{\mathbf{w}} \mathbf{w}^T \mathbf{S}_{XY} \mathbf{S}_{XY}^T \mathbf{w}, \text{ subject to}$$

$$\mathbf{w}^T \mathbf{S}_X \hat{\boldsymbol{\Phi}}_k = 0 \text{ and } \mathbf{w}^T \mathbf{w} = 1$$

$$\hat{\boldsymbol{\Phi}}_{k+1} = (\hat{\boldsymbol{\varphi}}_0, \dots, \hat{\boldsymbol{\varphi}}_k, \hat{\boldsymbol{\varphi}}_{k+1})$$

Version II With $\mathbf{Q}_{S_k} = \mathbf{I} - \mathbf{P}_{S_k}$,

$$S_k = \text{span}(\mathbf{S}_X \hat{\boldsymbol{\Phi}}_k)$$

$$\hat{\boldsymbol{\varphi}}_{k+1} = \mathbf{Q}_{S_k} \mathbf{S}_{XY} / \|\mathbf{Q}_{S_k} \mathbf{S}_{XY}\|$$

$$\hat{\boldsymbol{\Phi}}_{k+1} = (\hat{\boldsymbol{\varphi}}_0, \dots, \hat{\boldsymbol{\varphi}}_k, \hat{\boldsymbol{\varphi}}_{k+1}).$$

At termination $\hat{\boldsymbol{\Phi}}_{\text{pls}} = \hat{\boldsymbol{\Phi}}_q$.

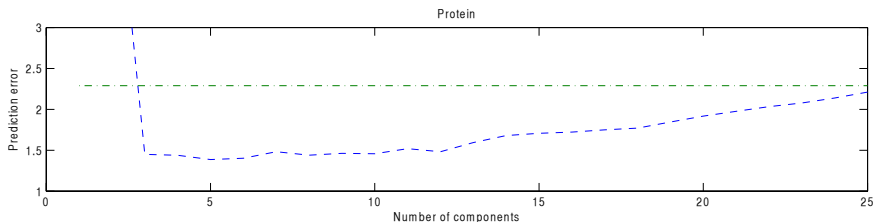
Version III

$$\text{span}(\hat{\boldsymbol{\Phi}}_{\text{pls}}) = \text{span}(\mathbf{S}_{XY}, \mathbf{S}_X \mathbf{S}_{XY}, \dots, \mathbf{S}_X^{q-1} \mathbf{S}_{XY})$$

Beef protein



Predict protein content (Y) of beef based on spectral measurements at $p = 50$ wave lengths, $n = 103$. 5-fold cross validation used for root mean squared prediction error. **Blue, PLS**; **Green, OLS**.



2. Brief History of PLS regression, Mixed with Personal Encounters

- S. Wold, H. Martens and H. Wold (1983).
- Los Alamos National Labs, ~1980.
- Gordon Research Conferences on Statistics in Chemistry and Chemical Engineering, 1980–mid 1990's
- Statistics: Helland (1990, 1992, 2001), Næs and Helland (1993), Frank and Friedman (1993), Garthwaite (1994)
- Chemometrics: de Jong's (1993) SIMPLS algorithm
- Chemometrics conference early 2000's
- Outside of Chemometrics: Micro-array data, FMRI data, biomedical analyses, tumor classification, bioprocesses, forecasting, characteristics of craft beer,
Selling point – serviceable when $n < p$.

- Chun and Keleş (2010): PLS regression is consistent for β only if $p/n \rightarrow 0$ as n & $p \rightarrow \infty$. That is, $\hat{\beta}_{\text{pls}} \xrightarrow{pr} \beta$ only if only if $p/n \rightarrow 0$. (same context)

1**2****3****4**

- Chun and Keleş (2010): PLS regression is consistent for β only if $p/n \rightarrow 0$ as n & $p \rightarrow \infty$. That is, $\hat{\beta}_{\text{pls}} \xrightarrow{p} \beta$ only if only if $p/n \rightarrow 0$. (same context)
 - 1 Consistency doesn't always signal the usefulness of a method in practice.
 - 2
 - 3
 - 4

- Chun and Keleş (2010): PLS regression is consistent for β only if $p/n \rightarrow 0$ as n & $p \rightarrow \infty$. That is, $\hat{\beta}_{\text{pls}} \xrightarrow{p} \beta$ only if only if $p/n \rightarrow 0$. (same context)
 - 1 Consistency doesn't always signal the usefulness of a method in practice.
 - 2 Operationally easy way to fit $n < p$ regressions that performed well against the competing methods of the day.
 - 3
 - 4

- Chun and Keleş (2010): PLS regression is consistent for β only if $p/n \rightarrow 0$ as n & $p \rightarrow \infty$. That is, $\hat{\beta}_{\text{pls}} \xrightarrow{p} \beta$ only if only if $p/n \rightarrow 0$. (same context)
 - 1 Consistency doesn't always signal the usefulness of a method in practice.
 - 2 Operationally easy way to fit $n < p$ regressions that performed well against the competing methods of the day.
 - 3 The technical restrictions (eg. Σ_X bdd as $p \rightarrow \infty$) of Chun and Keleş ruled out important practical contexts.
 - 4

- Chun and Keleş (2010): PLS regression is consistent for β only if $p/n \rightarrow 0$ as n & $p \rightarrow \infty$. That is, $\hat{\beta}_{\text{pls}} \xrightarrow{p} \beta$ only if only if $p/n \rightarrow 0$. (same context)
 - 1 Consistency doesn't always signal the usefulness of a method in practice.
 - 2 Operationally easy way to fit $n < p$ regressions that performed well against the competing methods of the day.
 - 3 The technical restrictions (eg. Σ_X bdd as $p \rightarrow \infty$) of Chun and Keleş ruled out important practical contexts.
 - 4 Chun and Keleş used their result to motivate a sparse version of PLS, but H. Wold, et al. (circa 1988) & S. Wold, et al. (1996) argued against sparse versions of PLS.

- Inge Helland's 2011 visit: Envelopes & PLS regression algorithms
- Cook, Helland and Su (2013):
 - gave first firm statistical model for PLS regression
 - proved that PLS regression provides a root- n consistent moment-based envelope method (p fixed)
 - showed that likelihood-based envelope methods dominate PLS when $n \gg p$ (p fixed)
- Cook and Forzani (2017, 2018) studied PLS in high-dimensional ($n < p$) regressions.

3. PLS regression via Envelopes (Cook, Helland & Su, 2013)

Motivation

Driving question: Are there linear combinations $\Phi^T \mathbf{X}$ of \mathbf{X} that carry all the information that \mathbf{X} has about Y ? So, changing \mathbf{X} while holding $\Phi^T \mathbf{X}$ fixed has no impact on the distribution of Y .

Motivation

Driving question: Are there linear combinations $\Phi^T \mathbf{X}$ of \mathbf{X} that carry all the information that \mathbf{X} has about Y ? So, changing \mathbf{X} while holding $\Phi^T \mathbf{X}$ fixed has no impact on the distribution of Y .

Let $\mathcal{E} = \text{span}(\Phi)$. Formally, we seek the smallest subspace $\mathcal{E} \subseteq \mathbb{R}^p$ so that

$$(Y, \mathbf{P}_{\mathcal{E}} \mathbf{X}) \perp\!\!\!\perp \mathbf{Q}_{\mathcal{E}} \mathbf{X}.$$

Consequently, \mathbf{X} affects Y only via $\mathbf{P}_{\mathcal{E}} \mathbf{X}$, the material part of \mathbf{X} .

The condition $(Y, \mathbf{P}_\varepsilon \mathbf{X}) \perp \mathbf{Q}_\varepsilon \mathbf{X}$ holds iff

- 1 $\mathcal{B} := \text{span}(\boldsymbol{\beta}) \subseteq \varepsilon$, so ε **envelops** \mathcal{B}
 - 2 $\boldsymbol{\Sigma}_X = \mathbf{P}_\varepsilon \boldsymbol{\Sigma}_X \mathbf{P}_\varepsilon + \mathbf{Q}_\varepsilon \boldsymbol{\Sigma}_X \mathbf{Q}_\varepsilon$ so ε is a **reducing subspace** of $\boldsymbol{\Sigma}_X$
- The intersection of all such subspaces ε is called the $\boldsymbol{\Sigma}_X$ -envelope of \mathcal{B} , $\mathcal{E}_{\boldsymbol{\Sigma}_X}(\mathcal{B})$ with semi-orthogonal basis $\boldsymbol{\Phi}$
 - Model:

$$\begin{aligned} Y &= \alpha + \boldsymbol{\eta}^T (\boldsymbol{\Phi}^T \mathbf{X}) + \varepsilon \\ \boldsymbol{\Sigma}_X &= \boldsymbol{\Phi} \boldsymbol{\Delta} \boldsymbol{\Phi}^T + \boldsymbol{\Phi}_0 \boldsymbol{\Delta}_0 \boldsymbol{\Phi}_0^T \end{aligned}$$

$\boldsymbol{\beta} = \boldsymbol{\Phi} \boldsymbol{\eta}$ and $\boldsymbol{\Sigma}_X$ are still of primary interest. \mathbf{X} not ancillary.

- SIMPLS return a basis for $\mathcal{E}_{\boldsymbol{\Sigma}_X}(\mathcal{B})$ in the population and $q = \dim(\mathcal{E}_{\boldsymbol{\Sigma}_X}(\mathcal{B}))$, the *number of components*

How do we get an estimator of $\mathcal{E}_{\Sigma_X}(\mathcal{B})$? Two options for now:

- Use MLE when $n \gg p$, giving $\hat{\Phi}_{\text{mle}}$
- Use SIMPLS since $\hat{\Phi}_{\text{pls}}$ is a root- n consistent estimator of a basis for $\mathcal{E}_{\Sigma_X}(\mathcal{B})$ when q is known.

In either case, as we saw, we estimate

$$\hat{\beta} = \hat{\Phi} \hat{\eta} = \hat{\Phi} (\hat{\Phi}^T \mathbf{S}_X \hat{\Phi})^{-1} \hat{\Phi}^T \mathbf{S}_{XY}$$

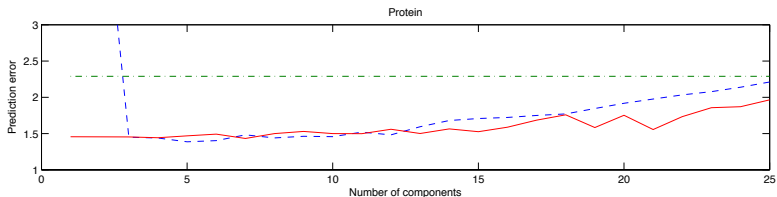
with $\hat{\Phi} = \hat{\Phi}_{\text{pls}}$ or $\hat{\Phi} = \hat{\Phi}_{\text{mle}}$, both of which provide root- n consistent estimators of β when q is assumed known.

Asymptotically, $\hat{\beta}_{\text{mle}}$ is better than $\hat{\beta}_{\text{pls}}$, but perhaps not when $n < p$ or when $n \not\gg p$.

Beef protein



Predict protein content (Y) of beef based on spectral measurements at $p = 50$ wave lengths, $n = 103$. 5-fold cross validation used for prediction errors. **Blue, PLS;** **Green, OLS,** **Red, Env.**



Comparison of PLS and Envelopes

Rimal, Almøy and Sæbø(2019). Comparison of Multi-response Prediction Methods. *Chemometrics and Intelligent Laboratory Systems*, to appear.

Analysis using both simulated data and real data has shown that the envelope methods are more stable, less influenced by . . . [predictor collinearity] and in general, performed better than PCR and PLS methods. These methods are also found to be less dependent on the number of components.

When $n < p$ they used PCA to reduce the predictors prior to applying envelope methodology.

The results from this study will . . . encourage [researchers] to use newly developed methods such as the envelopes.

4. PLS n, p Asymptotics

(Cook & Forzani, 2017, 2018)

Same context with

- $q = \dim(\mathcal{E}_{\Sigma_X}(\mathcal{B}))$ fixed. Φ semi-orthogonal basis for $\mathcal{E}_{\Sigma_X}(\mathcal{B})$. $\Phi^T \mathbf{X}$ true active predictors.
- $\tau^2 = \text{var}(Y | \mathbf{X})$ bounded away from 0
- $\beta \neq 0$, so $q \geq 1$

Gauges: Fitted value \hat{Y}_N at a new independent \mathbf{X}_N :

$\hat{Y}_N - E(Y|\mathbf{X}_N)$ & Estimation: $\|\hat{\beta}_{\text{pls}} - \beta\|_{\Sigma_X}$

Governing Quantities

Let $\Phi^T \mathbf{X}$ and $\Phi_0^T \mathbf{X}$ be true active and inactive predictors, (Φ, Φ_0) orthogonal.

- Colinearity: $\rho(p) = \text{sum of population VIFs for the reg. of } Y \text{ on } \Phi^T \mathbf{X}$.
- Signal: $\eta(p) \asymp \text{trace}(\Phi^T \Sigma_{\mathbf{X}} \Phi) = \text{trace}\{\text{var}(\Phi^T \mathbf{X})\}$
- Noise: $\kappa(p) \asymp \text{trace}(\Phi_0^T \Sigma_{\mathbf{X}} \Phi_0) = \text{trace}\{\text{var}(\Phi_0^T \mathbf{X})\}$

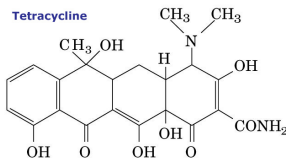
Special cases for reference:

- Abundance: If, as $p \rightarrow \infty$, $\|\Sigma_{\mathbf{X}Y}\|^2 \rightarrow \infty$ then $\eta \rightarrow \infty$
- Sparsity: If the regression is sparse (only q predictors are material) then η is bounded.

Conclusions, all assuming bounded $\rho(p)$ and noise $\kappa(p) = \text{trace}\{\text{var}(\Phi_0^T \mathbf{X})\} \asymp p$

- If the signal $\eta(p) = \text{trace}\{\text{var}(\Phi^T \mathbf{X})\}$ is bounded then PLS is not consistent and may not work well in high dimensional regressions unless $n \gg p$. (sparsity, as argued against by H. Wold, S. Wold et al.)
- If the signal $\eta(p)$ is unbounded then PLS is consistent and should work well in high dimensional regressions, even when $n < p$. (abundance, as argued for by H. Wold, S. Wold et al.)
- Chun-Keleş' assumption that the eigenvalues of $\text{var}(\mathbf{X})$ are bounded away from 0 and ∞ implies sparsity, which motivated their sparse methodology.

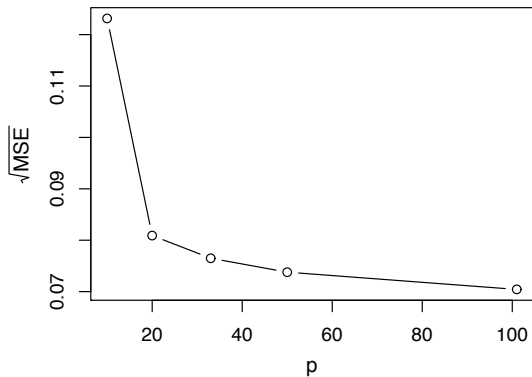
Tetracycline



Goal: Predict Tetracycline concentration in blood. $n = 50$ training samples constructed by spiking blank sera. $m = 57$ validation samples constructed in the same way.

Predictors: fluorescence intensity at $p = 101$ equally spaced points in 450 – 550 nm. Original PLS analysis indicated that $q = 4$ (Goicoechea and Olivieri, 2009) .

We repeated the analysis by selecting $p = 10, 20, 33, 50, 101$ equally spaced points, computing the validation root-MSE prediction in each case.



Overall conclusions

- Envelope methodology and SIMPLS estimate the same population quantity, $\mathcal{E}_{\Sigma_X}(\mathcal{B})$.
- Envelope methodology generally gives better results than SIMPLS, even when $n < p$ (Rimal, et al.).
- SIMPLS estimates are generally easier to compute.
- SIMPLS remains a serviceable method in abundant regressions.

Computing: z.umn.edu/envelopes

Papers: www.stat.umn.edu/~dennis

Book: Cook, R. D. (2018). *An introduction to Envelopes*. Wiley

Thank you!