Non-OLS Linear SEL Prediction: Partial Least Squares Regression

Stephen Vardeman

Analytics Iowa LLC

ISU Statistics and IMSE

Derived inputs and PLS

Here we continue to assume that the columns of X have been standardized and Y has been centered. We discuss a second prediction methodology involving some number M < p of predictors derived from the original inputs x_j , so-called Partial Least Squares (PLS) Regression. The shrinking methods treated thus far take no account of Y in determining directions or amounts of shrinkage. PLS employs Y.

First PLS component

Let

$$\mathbf{z}_1 = \sum_{j=1}^{p} \langle \mathbf{Y}, \mathbf{x}_j \rangle \, \mathbf{x}_j = \mathbf{X} \mathbf{X}' \mathbf{Y}$$

For $\mathbf{w}_1 = \mathbf{X}'\mathbf{Y} / \|\mathbf{X}'\mathbf{Y}\|$, $\mathbf{X}\mathbf{w}_1 = \mathbf{z}_1 / \|\mathbf{X}'\mathbf{Y}\|$ is a linear combination of the columns of \mathbf{X} maximizing, subject to the constraint that $\|\mathbf{w}\| = 1$, the quantity

 $|\langle \mathbf{Y}, \mathbf{X} \mathbf{w} \rangle|$

which is essentially the absolute sample covariance between the variables y and $\mathbf{x}'\mathbf{w}$. (The first principal component of \mathbf{X} maximizes $|\langle \mathbf{X}\mathbf{w}, \mathbf{X}\mathbf{w} \rangle|$ subject to the same constraint.)

Second PLS component

Then define X^1 by orthogonalizing the columns of X with respect to z_1 . That is, define the *j*th column of X^1 by

$$\mathbf{x}_{j}^{1} = \mathbf{x}_{j} - rac{\langle \mathbf{x}_{j}, \mathbf{z}_{1}
angle}{\langle \mathbf{z}_{1}, \mathbf{z}_{1}
angle} \mathbf{z}_{1}$$

and take

$$\mathbf{z}_2 = \sum_{j=1}^p \left< \mathbf{Y}, \mathbf{x}_j^1 \right> \mathbf{x}_j^1 = \mathbf{X}^1 \mathbf{X}^{1\prime} \mathbf{Y}$$

(For $\mathbf{w}_2 = \mathbf{X}^{1'}\mathbf{Y} / \|\mathbf{X}^{1'}\mathbf{Y}\|$, $\mathbf{X}^1\mathbf{w}_2 = \mathbf{z}_2 / \|\mathbf{X}^{1'}\mathbf{Y}\|$ is the linear combination of the columns of \mathbf{X}^1 maximizing

 $\left|\left< \mathbf{Y}, \mathbf{X}^1 \mathbf{w} \right> \right|$

subject to the constraint that $\| {f w} \| = 1.)$

PLS regression

Then for l > 1, define \mathbf{X}^{l} by orthogonalizing the columns of \mathbf{X}^{l-1} with respect to \mathbf{z}_{l} . That is, define the *j*th column of \mathbf{X}^{l} by

$$\mathbf{x}_{j}^{\prime} = \mathbf{x}_{j}^{\prime-1} - rac{\left\langle \mathbf{x}_{j}^{\prime-1}, \mathbf{z}_{l}
ight
angle}{\left\langle \mathbf{z}_{l}, \mathbf{z}_{l}
ight
angle} \mathbf{z}_{l}$$

and let

$$\mathbf{z}_{l+1} = \sum_{j=1}^{p} \left\langle \mathbf{Y}, \mathbf{x}_{j}^{l} \right\rangle \mathbf{x}_{j}^{l} = \mathbf{X}^{l} \mathbf{X}^{l\prime} \mathbf{Y}$$

Partial least squares regression uses the first M of these variables z_j as input variables.

PLS coefficient vector and predictor The PLS predictors \mathbf{z}_j are orthogonal by construction. Using the first Mof these as regressors, one has the vector of fitted output values $\widehat{\mathbf{Y}}^{PLS} = \sum_{j=1}^{M} \frac{\langle \mathbf{Y}, \mathbf{z}_j \rangle}{\langle \mathbf{z}_j, \mathbf{z}_j \rangle} \mathbf{z}_j$ Since the PLS predictors are (albeit recursively-computed data-dependent) linear combinations of columns of \mathbf{X} , it is possible to find a *p*-vector $\widehat{\boldsymbol{\beta}}_M^{PLS}$ (namely $(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \widehat{\mathbf{Y}}^{PLS}$) such that

$$\widehat{\mathbf{Y}}^{\mathsf{PLS}} = \mathbf{X}\widehat{\boldsymbol{\beta}}_{M}^{\mathsf{PLS}}$$

and thus produce the corresponding linear prediction rule

$$\widehat{f}(\mathbf{x}) = \mathbf{x}' \widehat{\boldsymbol{\beta}}_{M}^{\mathsf{PLS}}$$
(1)

PLS predictor complexity

It seems like in (1), the number of components, M, should function as a complexity parameter. But then again there is the following. When the \mathbf{x}_j are orthogonal, it's fairly easy to see that \mathbf{z}_1 is a multiple of $\widehat{\mathbf{Y}}^{\text{OLS}}$, $\widehat{\boldsymbol{\beta}}_1^{\text{PLS}} = \widehat{\boldsymbol{\beta}}_2^{\text{PLS}} = \cdots = \widehat{\boldsymbol{\beta}}_p^{\text{PLS}} = \widehat{\boldsymbol{\beta}}^{\text{OLS}}$, and all steps of partial least squares after the first are simply providing a basis for the orthogonal complement of the 1-dimensional subspace of $C(\mathbf{X})$ generated by $\widehat{\mathbf{Y}}^{\text{OLS}}$ (without improving fitting at all). That is, here changing M doesn't change flexibility of the fit at all.

Presumably, when the \mathbf{x}_j are nearly orthogonal, something similar happens, and one might thus expect PLS to be most effective as a shrinkage method where there are substantial correlations among columns of \mathbf{X} .

PLS complexity and effective df

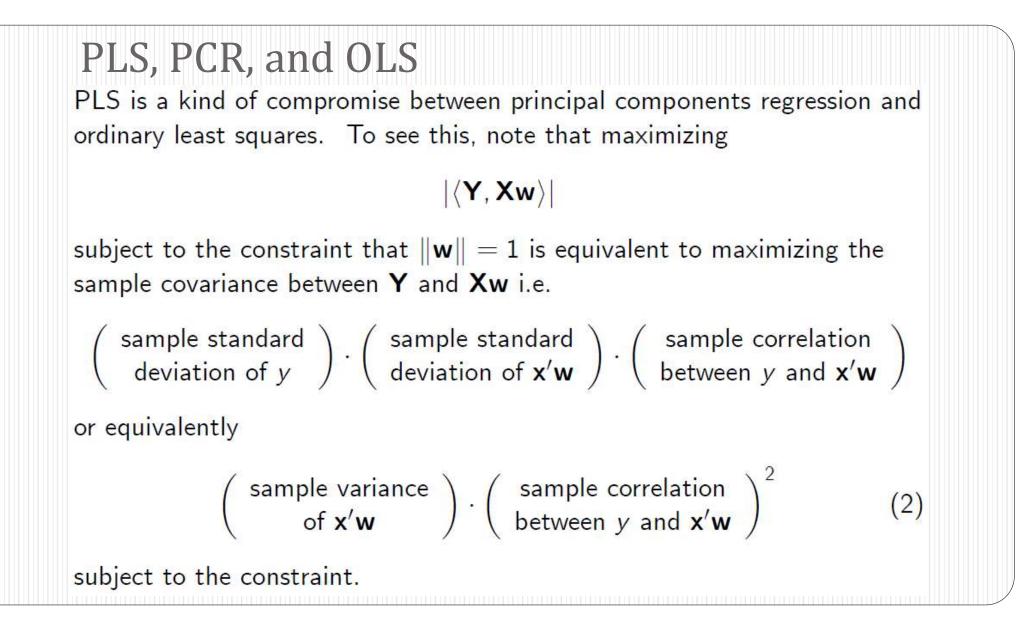
This observation about PLS in cases where predictors are orthogonal has another simple implication. That is that there will be no naive form for effective degrees of freedom for PLS. Since with z_j the *j*th principal component of **X** and, say,

$$\mathbf{Z}^{M} = (\mathbf{z}_{1}, \mathbf{z}_{2}, \dots, \mathbf{z}_{M})$$

we have

$$\widehat{\mathbf{Y}}^{\mathsf{PCR}} = \mathbf{Z}^{M} \left(\left(\mathbf{Z}^{M} \right)^{\prime} \mathbf{Z}^{M} \right)^{-1} \left(\mathbf{Z}^{M} \right)^{\prime} \mathbf{Y}$$

principal components regression on M components has effective degrees of freedom M. But the fact that the " Z^{M} " matrix corresponding to PLS depends upon Y makes PLS *nonlinear in* Y. And the "orthogonal X" argument shows that a PLS predictor with M = 1 can have effective degrees of freedom as large as rank (X).



PLS, PCR, and OLS cont.

Now if only the first term (the sample variance of $\mathbf{x}'\mathbf{w}$) were involved in (2), a first principal component direction would be an optimizing \mathbf{w}_1 , and $\mathbf{z}_1 = \|\mathbf{X}'\mathbf{Y}\| \mathbf{X}\mathbf{w}_1$ a multiple of the first principal component of \mathbf{X} . On the other hand, if only the second term were involved, $\hat{\boldsymbol{\beta}}^{\text{OLS}} / \|\hat{\boldsymbol{\beta}}^{\text{OLS}}\|$ would be an optimizing \mathbf{w}_1 , and $\mathbf{z}_1 = \hat{\mathbf{Y}}^{\text{OLS}} \|\mathbf{X}'\mathbf{Y}\| / \|\hat{\boldsymbol{\beta}}^{\text{OLS}}\|$ a multiple of the vector of ordinary least squares fitted values. The use of the product of two terms can be expected to produce a compromise between these two.

This logic applied at later steps in the PLS algorithm then produces for \mathbf{z}_l a compromise between a first principal component of \mathbf{X}^{l-1} and a suitably constrained multiple of the vector of least squares fitted values for \mathbf{Y} based on the matrix of inputs \mathbf{X}^{l-1} . \mathbf{X}^l has columns that are the corresponding columns of \mathbf{X} minus their projections onto the span of $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_l\}$ and $\mathcal{C}(\mathbf{X}) \supset \mathcal{C}(\mathbf{X}^1) \supset \mathcal{C}(\mathbf{X}^2) \cdots$