

Non-OLS Linear SEL Prediction: Principal Components Regression

Stephen Vardeman
Analytics Iowa LLC
ISU Statistics and IMSE

Linear predictors based on derived inputs

Another possible approach to the problem of finding an appropriate level of complexity in a fitted linear SEL prediction rule is to consider regression on some number $M < p$ of predictors derived from the original inputs x_j . Two such methods are those of Principal Components Regression and Partial Least Squares Regression. Here we continue to assume that **the columns of X have been standardized and Y has been centered.**

Principal components regression

Here the p columns of predictors in \mathbf{X} are replaced with the first M principal components of \mathbf{X}

$$\mathbf{z}_j = \mathbf{X}\mathbf{v}_j = d_j\mathbf{u}_j$$

The vector of fitted predictions is thus

$$\hat{\mathbf{Y}}^{\text{PCR}} = \sum_{j=1}^M \frac{\langle \mathbf{Y}, \mathbf{z}_j \rangle}{\langle \mathbf{z}_j, \mathbf{z}_j \rangle} \mathbf{z}_j = \sum_{j=1}^M \langle \mathbf{Y}, \mathbf{u}_j \rangle \mathbf{u}_j$$

Comparing this to ridge regression, we see that ridge regression shrinks the coefficients of the principal components \mathbf{u}_j according to their importance in making up \mathbf{X} , while principal components regression "zeros out" those least important in making up \mathbf{X} . Since the \mathbf{u}_j constitute an orthonormal basis for $C(\mathbf{X})$, for $\text{rank}(\mathbf{X}) = r$,

$$\left\| \hat{\mathbf{Y}}^{\text{PCR}} \right\|^2 = \sum_{j=1}^M \langle \mathbf{Y}, \mathbf{u}_j \rangle^2 \leq \sum_{j=1}^r \langle \mathbf{Y}, \mathbf{u}_j \rangle^2 = \left\| \hat{\mathbf{Y}}^{\text{OLS}} \right\|^2$$

PCR and shrinkage

Further, $\hat{\mathbf{Y}}^{\text{PCR}}$ can be written in terms of the original inputs as

$$\hat{\mathbf{Y}}^{\text{PCR}} = \sum_{j=1}^M \langle \mathbf{Y}, \mathbf{u}_j \rangle \frac{1}{d_j} \mathbf{X} \mathbf{v}_j = \mathbf{X} \left(\sum_{j=1}^M \frac{1}{d_j^2} \langle \mathbf{Y}, \mathbf{X} \mathbf{v}_j \rangle \mathbf{v}_j \right)$$

so that

$$\hat{\boldsymbol{\beta}}^{\text{PCR}} = \sum_{j=1}^M \frac{1}{d_j^2} \langle \mathbf{Y}, \mathbf{X} \mathbf{v}_j \rangle \mathbf{v}_j$$

$\hat{\boldsymbol{\beta}}^{\text{OLS}}$ is the $M = \text{rank}(\mathbf{X})$ version of $\hat{\boldsymbol{\beta}}^{\text{PCR}}$. As the \mathbf{v}_j are orthonormal, it is clear (as on the previous slide for $\hat{\mathbf{Y}}$) that

$$\|\hat{\boldsymbol{\beta}}^{\text{PCR}}\| \leq \|\hat{\boldsymbol{\beta}}^{\text{OLS}}\|$$

So principal components regression shrinks both $\hat{\mathbf{Y}}^{\text{OLS}}$ toward $\mathbf{0}$ in \mathbb{R}^N and $\hat{\boldsymbol{\beta}}^{\text{OLS}}$ toward $\mathbf{0}$ in \mathbb{R}^p .

Choice of M and a caveat

The number of components used in PCR is a complexity parameter. The larger is M , the more complex is the predictor. So rational choice of M is in practice based on cross-validation.

There is, of course, nothing guaranteeing that directions in \mathbb{R}^p of large variation in training set inputs are directions in which y has large rates of change with respect to inputs. That said, PCR is a popular form of "dimension-reduction" for inputs that is sometimes highly effective in sorting out a low-dimensional "signal" in a "large p "-dimensional input.