

# Qualms about Dimension Reduction Theory: Central Subspace, Envelopes, Variable Selection

David J. Olive\*  
Southern Illinois University

December 16, 2025

## Abstract

Although envelopes, lasso, PLS, SDR, and variable selection are very useful, some qualms about the explanations and theory are given. Balms are needed.

**KEY WORDS:** Lasso; Model Selection; Outliers; PLS; SDR; Variable Selection

## 1 INTRODUCTION

Some notation for dimension reduction methods is needed. Let  $\mathbf{y} = (Y_1, \dots, Y_r)^T$  be a vector of  $r$  response variables and let  $\mathbf{x} = (x_1, \dots, x_p)^T$  be a vector of  $p$  predictor variables. For example, predict  $Y_1$  = mussel muscle mass and  $Y_2$  = mussel shell mass from  $x_1$  = height,  $x_2$  = width, and  $x_3$  = length of the mussel shell. The  $r = 2$  and  $p = 3$ .

Let  $\mathbf{A}$  be an  $k \times c$  matrix and let  $\mathcal{A} = \text{span}(\mathbf{A})$  be the subspace of  $\mathbb{R}^k$  spanned by the columns of  $\mathbf{A}$ . For a symmetric  $k \times k$  matrix  $\mathbf{C}$ , let  $\mathbf{C} > 0$  and  $\mathbf{C} \geq 0$  denote that  $\mathbf{C}$  is positive definite or positive semidefinite, respectively. Then the projection onto  $\mathcal{A}$  in the  $\mathbf{C}$  inner product is  $\mathbf{P}_{\mathcal{A}(\mathbf{C})} = \mathbf{P}_{\mathbf{A}(\mathbf{C})} = \mathbf{A}(\mathbf{A}^T \mathbf{C} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{C}$  provided  $\mathbf{A}^T \mathbf{C} \mathbf{A} > 0$ . Let  $\mathbf{Q}_{\mathcal{A}(\mathbf{C})} = \mathbf{I}_r - \mathbf{P}_{\mathcal{A}(\mathbf{C})}$  be the orthogonal projection. If  $\mathbf{C} = \mathbf{I}_r$ , let  $\mathbf{P}_{\mathcal{A}} = \mathbf{P}_{\mathbf{A}} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ .

Let the  $p \times p$  covariance matrix of  $\mathbf{x}$  be  $\text{Cov}(\mathbf{x}) = \mathbf{\Sigma}_{\mathbf{x}} = E[(\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))^T]$  and the  $p \times 1$  vector  $\text{Cov}(\mathbf{x}, Y) = \mathbf{\Sigma}_{\mathbf{x}Y} = E[(\mathbf{x} - E(\mathbf{x}))(Y - E(Y))]$  =  $(\text{Cov}(x_1, Y), \dots, \text{Cov}(x_p, Y))^T$ . Let the estimators be

$$\mathbf{S}_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \quad \text{and} \quad \mathbf{S}_{\mathbf{x}Y} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(Y_i - \bar{Y}).$$

Let  $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbf{\Sigma}_{\mathbf{x}\mathbf{y}}$  be estimated by  $\mathbf{S}_{\mathbf{x}\mathbf{y}}$ .

---

\*David J. Olive is Professor, School of Mathematical & Statistical Sciences, Southern Illinois University, Carbondale, IL 62901, USA.

The *response envelope* for the regression of  $\mathbf{y}$  on  $\mathbf{x}$  is the predictor envelope for the regression of  $\mathbf{x}$  on  $\mathbf{y}$ . Let  $\mathcal{L}$  be a subspace of  $\mathbb{R}^p$ . Let  $\{\gamma_1, \dots, \gamma_q\}$  be a basis for  $\mathcal{L}$ , and let  $\{\gamma_1, \dots, \gamma_q, \gamma_{q+1}, \dots, \gamma_p\}$  be a basis for  $\mathbb{R}^p$ . Then several dimension reduction methods involve estimating  $q$  and  $\hat{\gamma}_i$  for  $i = 1, \dots, q$ . Let  $\mathbf{S} = [\gamma_1 \dots \gamma_q]$  be the basis matrix for  $\mathcal{L}$ . Let  $\mathbf{P}_{\mathbf{S}} = \mathbf{P}_{\mathcal{L}}$  and  $\mathbf{Q}_{\mathcal{L}} = \mathbf{Q}_{\mathbf{S}}$ . Let  $\mathbf{y} \perp \mathbf{x} | \mathbf{Ax}$  indicate that the response vector  $\mathbf{y}$  is independent of the predictors  $\mathbf{x}$  given  $\mathbf{Ax}$ . Then a subspace  $\mathcal{L} \subseteq \mathbb{R}^p$  that satisfies  $\mathbf{y} \perp \mathbf{x} | \mathbf{P}_{\mathcal{L}} \mathbf{x}$  is a *dimension reduction subspace* for the regression of  $\mathbf{y}$  on  $\mathbf{x}$ . Thus the reduced predictors  $\mathbf{w} = \mathbf{P}_{\mathcal{L}} \mathbf{x}$  hold all of the information that  $\mathbf{x}$  has about  $\mathbf{y}$ . If the intersection of all dimension reduction subspaces is a dimension reduction subspace, then that subspace is called the *central subspace*, denoted by  $\mathcal{L}_{\mathbf{y}|\mathbf{x}}$ .

For a predictor envelope  $\mathcal{L}$ , we want a)  $\mathbf{y} \perp \mathbf{x} | \mathbf{P}_{\mathcal{L}} \mathbf{x}$ , and b)  $\mathbf{P}_{\mathcal{L}} \mathbf{x} \perp \mathbf{Q}_{\mathcal{L}} \mathbf{x}$ . Hence  $\mathbf{w} = \mathbf{P}_{\mathcal{L}} \mathbf{x}$  is (wanted or) *material* for the regression of  $\mathbf{y}$  on  $\mathbf{x}$  while  $\mathbf{u} = \mathbf{Q}_{\mathcal{L}} \mathbf{x}$  is (unwanted or) *immaterial* for the regression of  $\mathbf{y}$  on  $\mathbf{x}$ . Then  $\mathcal{L} \subseteq \mathbb{R}^p$  reduces  $\Sigma \mathbf{x}$  if and only if  $\text{Cov}(\mathbf{P}_{\mathbf{S}} \mathbf{x}, \mathbf{Q}_{\mathbf{S}} \mathbf{x}) = \mathbf{0}$ . Assume  $\mathcal{L}_{\mathbf{y}|\mathbf{x}} \subseteq \text{span}(\Sigma \mathbf{x})$ . The *predictor envelope* (subspace)  $\mathcal{E}_{\mathbf{x}}$  for the regression of  $\mathbf{y}$  on  $\mathbf{x}$  is the intersection of all dimension reduction subspaces that reduce  $\Sigma \mathbf{x}$  and contain  $\mathcal{L}_{\mathbf{y}|\mathbf{x}}$ . Hence  $\mathcal{L}_{\mathbf{y}|\mathbf{x}} \subseteq \mathcal{E}_{\mathbf{x}}$ . The envelope subspace may be larger than the central subspace, but tends to handle high predictor collinearity better than sufficient dimension reduction (SDR) estimators of the central subspace.

For partial least squares (PLS), PLS1 algorithms are for regression with a univariate response  $Y$ , and while PLS2 algorithms are for a multivariate response  $\mathbf{y}$ , and are also PLS1 algorithms. These algorithms produce estimated basis vectors  $\hat{\gamma}_i$  sequentially, using the response variables. Let  $\hat{\mathbf{A}}_k$  be the matrix with  $i$ th row  $\hat{\gamma}_i^T$  for  $i = 1, \dots, k$ . Let the estimated coefficient matrix  $\hat{\mathbf{B}}_k = \hat{\mathbf{A}}_k^T (\hat{\mathbf{A}}_k \mathbf{S} \mathbf{x} \hat{\mathbf{A}}_k^T)^{-1} \hat{\mathbf{A}}_k \mathbf{S} \mathbf{x} \mathbf{y}$  for  $k = 1, \dots, J$  where  $k = q$  is especially important. Then  $\hat{\mathbf{B}}_k$  can also be obtained by the OLS multivariate linear regression of  $\mathbf{y}$  on  $\mathbf{w} = \hat{\mathbf{A}}_k \mathbf{x} = (W_1, \dots, W_k)^T$  where  $W_i = \hat{\gamma}_i^T \mathbf{x}$ .  $J \leq \min(n-1, p)$  needs to be small enough so that  $(\hat{\mathbf{A}}_k \mathbf{S} \mathbf{x} \hat{\mathbf{A}}_k^T)^{-1}$  exists. If  $\mathbf{S}^{-1}$  exists, then  $\hat{\mathbf{B}}_k = \mathbf{P}_{\hat{\mathbf{A}}_k^T(\mathbf{S} \mathbf{x})} \hat{\beta}_{OLS}$ , and if  $\hat{\mathbf{A}}_p^{-1}$  exists, then  $\hat{\mathbf{B}}_p = \hat{\mathbf{B}}_{OLS} = \mathbf{S}^{-1} \mathbf{S} \mathbf{x} \mathbf{y}$ . Replace  $\mathbf{B}$  by  $\beta$  if  $r = 1$  so  $Y$  is used instead of  $\mathbf{y}$ .

The NIPALS algorithm uses  $\hat{\gamma}_1 = 1$ st eigenvector of  $\mathbf{S} \mathbf{x} \mathbf{y} \mathbf{S}^T \mathbf{x} \mathbf{y}$ . For  $k < q$ , the SIMPLS algorithm uses  $\hat{\mathbf{A}}_k^T$  and  $\hat{\gamma}_{k+1} = \max_{\gamma \in \mathbb{R}^p} \gamma^T \mathbf{S} \mathbf{x} \mathbf{y} \mathbf{S}^T \mathbf{x} \mathbf{y} \gamma = \max_{\gamma \in \mathbb{R}^p} (\gamma^T \mathbf{S} \mathbf{x} \mathbf{y})^2$  subject to  $\gamma^T \mathbf{S} \mathbf{x} \hat{\mathbf{A}}_k^T = \mathbf{0}$  and  $\gamma^T \gamma = 1$ . From canonical correlation analysis, finding  $\max_{\gamma \neq \mathbf{0}} \text{Cor}(\gamma^T \mathbf{x}, Y)$  is equivalent to maximizing

$$\max_{\gamma \neq \mathbf{0}} \frac{\gamma^T \Sigma \mathbf{x} \mathbf{y} \Sigma_{\mathbf{x} \mathbf{y}}^T \gamma}{\gamma^T \Sigma \mathbf{x} \gamma}$$

which has a maximum at  $\gamma = \Sigma \mathbf{x}^{-1} \Sigma \mathbf{x} \mathbf{y} = \beta_{OLS}$ . This result suggests that for PLS1 algorithms, the  $[\text{Cor}(\hat{\gamma}_i^T \mathbf{x}, Y)]^2$  are fairly high for  $i$  near 1. See Olive (2025).

For a univariate response  $Y$ ,  $\text{span}(\mathbf{A}_k^T)$  and  $\text{span}(\hat{\mathbf{A}}_k^T)$  are the same for the SIMPLS, NIPALS, and HPLS PLS1 algorithms for  $k = 1, \dots, q$ . See Cook and Forzani (2024, p. 96). The Helland algorithm (HPLS) uses  $\gamma_i = \Sigma_{\mathbf{x}}^{i-1} \Sigma \mathbf{x} \mathbf{y}$  and  $\hat{\gamma}_i = \mathbf{S} \mathbf{x}^{i-1} \mathbf{S} \mathbf{x} \mathbf{y}$  with  $\mathbf{S}^0_{\mathbf{x}} = \Sigma_{\mathbf{x}}^0 = \mathbf{I}_p$ . Thus  $\hat{\gamma}_1 = \mathbf{S} \mathbf{x} \mathbf{y}$  and  $W_i = \hat{\gamma}_i^T \mathbf{x}$ .

Suppose  $\Sigma_{\mathbf{x}}$  has eigenvalue eigenvector pairs  $(\lambda_1, \mathbf{d}_1), \dots, (\lambda_p, \mathbf{d}_p)$  where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ . Let the eigenvalue eigenvector pairs of  $\mathbf{S}_{\mathbf{x}}$  be  $(\hat{\lambda}_1, \hat{\mathbf{d}}_1), \dots, (\hat{\lambda}_p, \hat{\mathbf{d}}_p)$  where  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$ . Then a common principal component analysis (PCA) uses  $\hat{\gamma}_i = \hat{\mathbf{d}}_i$ .

Although envelopes, lasso, PLS, SDR, and variable selection are very useful, there are some qualms about the explanations and theory.

## 2 Qualms

### 2.1 Predictor Envelope – PCA Relationship

A very interesting property of the predictor envelope is that if the eigenvalues of  $\Sigma_{\mathbf{x}}$  are unique, and if  $\dim(\mathcal{E}_{\mathbf{x}}) = q$ , then  $\mathcal{E}_{\mathbf{x}}$  is spanned by  $q$  eigenvectors of  $\Sigma_{\mathbf{x}}$ . See Cook and Forzani (2024, p. 38). Hence the predictor envelope and principal component analysis (PCA) have some similarities. Then PCA replaces the variables  $x_1, \dots, x_p$  by  $w_1, \dots, w_J$  where  $w_i = \hat{\mathbf{d}}_i^T \mathbf{x}$  for  $i = 1, \dots, J$  where  $J$  is chosen using a scree plot or some other method. The estimated predictor envelope replaces the variables  $x_1, \dots, x_p$  by  $w_1, \dots, w_q$  where  $w_i = \hat{\mathbf{d}}_i^T \mathbf{x}$  for  $i = 1, \dots, q$  where  $q$  is chosen in some way using the response variables  $Y_1, \dots, Y_r$ . Hence  $q$  is often much smaller than  $J$ . In particular, let  $r = 1$  so there is a single response variable  $Y$ . If  $\Sigma_{\mathbf{x}} = \text{diag}(1, \dots, p)$ , then the eigenvectors of  $\Sigma_{\mathbf{x}}$  are the columns of the identity matrix  $\mathbf{I}_p$  (up to spans). Suppose the MLR model  $Y = \alpha + \beta^T \mathbf{x} + e$  holds with  $\beta = (1, \dots, 1)^T = \mathbf{1}$  where the zero mean errors are iid with variance  $V(e_i) = \sigma^2$ . Then the central subspace =  $\text{span}(\beta)$ , but PCA needs all  $p$  eigenvectors. This result appears to imply that  $q = p$  for the predictor envelope. For PCA, the assumption that  $\beta \in \mathbb{R}^p$  is a much weaker assumption than  $\beta \in \mathbb{R}^q$  where  $1 \leq q < p$ . As  $r$  increases, the regularity conditions get stronger.

Intuitively, since PCA does not use the response variables to form the  $\hat{\gamma}_i = \hat{\mathbf{d}}_i$ , then SDR methods, response envelopes, and PLS should outperform PCA at least if  $n \geq 10p$  and  $r$  is small.

Let  $W_1, \dots, W_p$  be the linear combinations  $\hat{\mathbf{d}}_k^T \mathbf{x}$  ordered with respect to the highest squared correlations  $r_1^2 \geq r_2^2 \geq \dots \geq r_p^2$  where the sample correlation  $r_{i,Y} = \text{cor}(x_i, Y)$ . Then for the predictor envelope, does  $w_i = W_i$  for  $i = 1, \dots, q$ ?

The above ordering is marginal maximum likelihood estimator (MMLE) variable selection. See Fan and Lv (2008) and Fan and Song (2010), who give some drawbacks of the method.

For  $r > 1$ , let  $W_1, \dots, W_p$  be the linear combinations  $\hat{\gamma}_k^T \mathbf{x}$  ordered with respect  $r_1^2 \geq r_2^2 \geq \dots \geq r_p^2$  where  $r_{i,\mathbf{y}}^2 = \frac{1}{r} \sum_{j=1}^r [\text{cor}(x_i, Y_j)]^2$  for  $i = 1, \dots, p$ . Let  $U_1, \dots, U_p$  be the linear combinations  $\hat{\gamma}_k^T \mathbf{x}$  ordered with respect to  $m_1^2 \geq m_2^2 \geq \dots \geq m_p^2$  where  $m_{i,\mathbf{y}}^2 = \max_{j=1, \dots, r} [\text{cor}(x_i, Y_j)]^2$ . Hence  $m_i^2 = m_{(i),\mathbf{y}}^2$ , the order statistics of the  $m_{i,\mathbf{y}}^2$ . These two methods of ordering variables appear to be new, although the  $r_i^2$  can be motivated by marginal maximum likelihood estimator (MMLE) variable selection. The Olive (2025) SC scree plot of  $i$  versus  $r_i^2$  behaves like a scree plot of  $i$  versus the eigenvalues. Hence quantities like  $\sum_{i=1}^j r_i^2 / \sum_{i=1}^p r_i^2$  are of interest for  $j = 1, \dots, p$ , and scree plot techniques could be adapted to choose  $q$  for PCA.

## 2.2 Envelopes and PLS

Cook and Forzani (2024, pp. 83,92,111): to establish the link between envelopes and PLS, and for the theory for  $\hat{\beta}_{1PLS}$ , a very strong regularity condition for a univariate response is that  $\Sigma_{\mathbf{x}Y}$  is an eigenvector of  $\Sigma_{\mathbf{x}}$ , e.g.,  $\Sigma_{\mathbf{x}} = \tau^2 \mathbf{I}_p$ . The Cook and Forzani (2024, equation (4.1)) model is  $Y|\mathbf{x} = \alpha_1 + \beta_{1PLS}^T \mathbf{x} + e$ . OLS is a consistent estimator for such models, hence (4.1) forces  $\beta_{OLS} = \Sigma_{\mathbf{x}}^{-1} \Sigma_{\mathbf{x}Y} = \lambda \Sigma_{\mathbf{x}Y} = \beta_{1PLS}$ . If  $\Sigma_{\mathbf{x}} = \text{diag}(1, \dots, p)$ , then the eigenvectors of  $\Sigma_{\mathbf{x}}$  are the columns of the identity matrix  $\mathbf{I}_p$  (up to spans). Hence  $\beta_{OLS}$  and  $\Sigma_{\mathbf{x}Y}$  have  $p - 1$  zero elements and one nonzero element under (4.1).

Note that if  $\Sigma_{\mathbf{x}Y}$  is an eigenvector of  $\Sigma_{\mathbf{x}}$ , then for HPLS,  $\text{span}(\gamma_1, \dots, \gamma_p) = \text{span}(\Sigma_{\mathbf{x}Y})$  which is equal to  $\text{span}(\beta_{OLS})$  if  $\Sigma_{\mathbf{x}}^{-1}$  exists. For  $r = 1$ , SIMPLS, NIPALS, and HPLS, to get linearly independent population basis vectors, a necessary condition is that  $\Sigma_{\mathbf{x}Y}$  is not an eigenvector of  $\Sigma_{\mathbf{x}}$ .

## 2.3 What Is Variable Selection Actually Doing?

Currently, this discussion is from Olive (2025b).

Some simple results hold for model selection estimators when the number of predictors  $p$  is fixed. Several methods, including PLS, use  $p$  linear combinations  $\eta_1^T \mathbf{x}, \dots, \eta_p^T \mathbf{x}$ . Performing the ordinary least squares (OLS) regression of  $Y$  on  $(\hat{\eta}_1^T \mathbf{x}, \hat{\eta}_2^T \mathbf{x}, \dots, \hat{\eta}_k^T \mathbf{x})$  and a constant gives the  $k$ -component estimator.

Consider the OLS regression of  $Y$  on a constant and  $\mathbf{w} = (W_1, \dots, W_p)^T$  where, e.g.,  $W_j = x_j$  or  $W_j = \hat{\eta}_j^T \mathbf{x}$ . Let  $I$  index the variables in the model so  $I = \{1, 2, 4\}$  means that  $\mathbf{w}_I = (W_1, W_2, W_4)^T$  was selected. The full model  $I = F$  uses all  $p$  predictors and the constant with  $\beta_F = \beta = \beta_{OLS}$ . Let  $r$  be the residuals from the full OLS model and let  $r_I$  be the residuals from model  $I$  that uses  $\hat{\beta}_I$  with  $k$  predictors including a constant where  $2 \leq k \leq p + 1$ . Olive (2025a) noted that if the Mallows (1973) criterion  $C_p(I) \leq 2k$ , then the sample correlation

$$\text{cor}(r, r_I) \geq \sqrt{1 - \frac{p+1}{n}}. \quad (1.1)$$

Since the correlation gets arbitrarily close to 1 as  $n \rightarrow \infty$ , the model selection estimator and full OLS estimator are estimating the same population parameter  $\beta$ . This result holds if  $\hat{\beta}_{OLS}$  is a consistent estimator of  $\beta$ : heterogeneity is allowed and the cases do not need to be independent and identically distributed (iid). For moderate sample size  $n$ , “weak” predictors  $W_j$  will often be omitted as long as  $\text{cor}(r_I, r)$  stays high.

The PLS literature often assumes (a1):  $Y|\mathbf{x} = \alpha + \mathbf{x}^T \beta_{kPLS} + e$  for some  $k$ . If  $Y|\mathbf{x} = \alpha + \mathbf{x}^T \beta + e$ , then under mild regularity conditions,  $\beta = \beta_{OLS}$ . Hence assumption (a1) forces  $\beta_{kPLS} = \beta_{OLS}$ . For  $k = 1$ , (a1) forces  $\beta_{OLS} = \beta_{1PLS}$  = an eigenvector of the covariance matrix  $\text{Cov}(\mathbf{x}) = \Sigma_{\mathbf{x}}$ . Assume instead that the cases  $(\mathbf{x}_i, Y_i)$  are iid with  $E(Y) = \mu_Y$  and  $E(\mathbf{x}) = \mu_{\mathbf{x}}$ . Let  $\text{Cov}(\mathbf{x}, Y) = \Sigma_{\mathbf{x}Y}$ . Then  $\beta = \beta_{OLS} = \Sigma_{\mathbf{x}}^{-1} \Sigma_{\mathbf{x}Y}$ ,  $\beta_{1PLS} = \theta \Sigma_{\mathbf{x}Y}$  where  $\theta$  is a constant, and  $\Sigma_{\mathbf{x}Y} = \Sigma_{\mathbf{x}} \beta$ , even when heterogeneity is present. Since  $\Sigma_{\mathbf{x}}$  is positive definite,  $\beta = \mathbf{0}$  iff  $\Sigma_{\mathbf{x}Y} = \mathbf{0}$ . This hypothesis can be tested by applying a one sample test to  $\mathbf{v}_i = (\mathbf{x}_i - \bar{\mathbf{x}})(Y_i - \bar{Y})$  for  $i = 1, \dots, n$ . See Olive and Zhang (2025), Olive et al. (2025), and Abid, Quaye, and Olive (2025).

Let an OLS working model for  $\hat{\beta}_{kPLS}$  be  $Y_i = \alpha_k + \theta_{k1}W_{1i} + \dots + \theta_{kk}W_{ki} + e_{ki}$  where  $W_{ji} = \mathbf{x}_i^T \hat{\Sigma}_{\mathbf{x}}^{j-1} \hat{\Sigma}_{\mathbf{x}Y}$  with  $\hat{\Sigma}_{\mathbf{x}}^0 = \Sigma_{\mathbf{x}}^0 = \mathbf{I}_p$ . Then  $\hat{\beta}_{kPLS} = (\sum_{j=1}^k \hat{\theta}_{kj} \hat{\Sigma}_{\mathbf{x}}^{j-1}) \hat{\Sigma}_{\mathbf{x}Y}$  and  $\beta_{kPLS} = (\sum_{j=1}^k \theta_{kj} \Sigma_{\mathbf{x}}^{j-1}) \Sigma_{\mathbf{x}Y}$  (under iid cases) using  $Y = \alpha_{kPLS} + \mathbf{x}^T \beta_{kPLS} + e_k$ . This result suggests that the  $\beta_{kPLS}$  are typically different for each  $k = 1, \dots, p$ , but  $\beta_{kPLS} = \beta_{qPLS}$  for  $k \leq q \leq p$  if  $\theta_{pj} = 0$  for  $k+1 \leq j \leq p$ .

Note that  $\beta \in \mathbb{R}^p$  is a much weaker assumption than  $\beta \in \mathbb{R}^m$  where  $1 \leq m < p$ . The model selection result of Equation (1) implies that  $\beta \in \mathbb{R}^m$  approximately true in that  $\theta_{pj} \approx 0$  for  $k^* + 1 \leq j \leq p$  for some  $m = k^* < p$ . Hence the predictors  $W_j$  are “weak” or “almost immaterial” for  $m+1 \leq j \leq p$ .

## 2.4 Oracle Property

## 2.5 Bigger Model

## 2.6 MLR with Heterogeneity

## 2.7 High Dimensional Multivariate Linear Regression

Let the multivariate linear regression model be  $\mathbf{y} = \alpha + \mathbf{B}^T \mathbf{x} + \epsilon$  where  $\mathbf{B} = [\beta_1 \beta_2 \dots \beta_r]$ . As before, let  $\hat{\mathbf{A}}_k \mathbf{x} = \mathbf{w} = (W_1, \dots, W_k)^T$  where  $\hat{\mathbf{A}}_k$  is the matrix with  $i$ th row  $\hat{\gamma}_i^T$  for  $i = 1, \dots, k$ . Let  $\mathbf{D}_k = [\theta_1 \theta_2 \dots \theta_k]$ . Let, for example,  $\hat{\Sigma}_{\mathbf{x}} = n \mathbf{S}_{\mathbf{x}} / (n-1)$ . Fit the working model

$$\mathbf{y} = \alpha + \mathbf{D}_k^T \mathbf{w} + \epsilon = \alpha + \mathbf{D}_k^T \hat{\mathbf{A}}_k \mathbf{x} + \epsilon$$

with  $\hat{\mathbf{B}}_k^T = \hat{\mathbf{D}}_k^T \hat{\mathbf{A}}_k$ . Then  $\hat{\mathbf{D}}_k = \hat{\Sigma}_{\mathbf{w}}^{-1} \hat{\Sigma}_{\mathbf{w}y}$ , and  $\hat{\mathbf{B}}_k = \hat{\mathbf{A}}_k^T \hat{\mathbf{D}}_k = \hat{\mathbf{A}}_k^T (\hat{\mathbf{A}}_k \hat{\Sigma}_{\mathbf{x}} \hat{\mathbf{A}}_k^T)^{-1} \hat{\mathbf{A}}_k \hat{\Sigma}_{\mathbf{x}y}$ . PLS2 algorithms can be used to get  $\hat{\mathbf{B}}_k$ .

Alternatively, let  $r$  univariate MLR's be fit to get  $\hat{\mathbf{B}}_U = \hat{\mathbf{A}}_H^T = [\hat{\eta}_1 \hat{\eta}_2 \dots \hat{\eta}_r]$  and

$$\hat{\mathbf{y}} = \hat{\alpha} + \hat{\mathbf{B}}_U^T \mathbf{x}.$$

Then  $\hat{\mathbf{A}}_H \mathbf{x} = \mathbf{w} = (W_1, \dots, W_r)^T$  where  $\hat{\mathbf{A}}_H$  is the matrix with  $i$ th row  $\hat{\eta}_i^T$  for  $i = 1, \dots, r$ . Let  $\mathbf{D}_H = [\theta_1 \theta_2 \dots \theta_r]$ , and fit the working model

$$\mathbf{y} = \alpha + \mathbf{D}_H^T \mathbf{w} + \epsilon = \alpha + \mathbf{D}_H^T \hat{\mathbf{A}}_H \mathbf{x} + \epsilon$$

to get  $\hat{\mathbf{D}}_H = \hat{\Sigma}_{\mathbf{w}}^{-1} \hat{\Sigma}_{\mathbf{w}y}$ , and  $\hat{\mathbf{B}}_H = \hat{\mathbf{A}}_H^T \hat{\mathbf{D}}_H = \hat{\mathbf{A}}_H^T (\hat{\mathbf{A}}_H \hat{\Sigma}_{\mathbf{x}} \hat{\mathbf{A}}_H^T)^{-1} \hat{\mathbf{A}}_H \hat{\Sigma}_{\mathbf{x}y}$ . PLS1 algorithms could be used to get the  $\hat{\eta}_i$ .

The new estimator  $\hat{\mathbf{B}}_H$  may be competitive with  $\hat{\mathbf{B}}_q$  if  $n \geq Jr$  with  $J \geq 5$ ,  $r \geq 2$ , and if the  $\hat{\eta}_i$  are not obtained from a working model for a multivariate linear regression.  $J$  much larger than 5 may be needed. Univariate lasso estimators could be used to give a competitor for the `glmnet` lasso estimator with “mgaussian.”

For a comparison of using  $r$  MLR's (e.g. PLS1) and multivariate linear regression of  $\mathbf{y}$  on  $\mathbf{w}$  using  $\hat{\gamma}_i \mathbf{x} = W_i$  (possibly with  $\hat{\gamma}_i = \gamma_j = j$ th column of the identity matrix  $\mathbf{I}_p$  so that  $W_i = x_j$ ), see, for example, Cook and Forzani (2024, pp. 106-107).

## 2.8 Outlier Resistance

To make outlier resistant analogs for many statistical techniques, including envelopes, lasso, PCA, and PLS, suppose the cases are  $\mathbf{w}_1, \dots, \mathbf{w}_n$ . Find the Euclidean distances  $D_i$  of the  $\mathbf{w}_i$  from the coordinatewise median of the  $n$  cases. Find the  $n_R$  cases that satisfy  $D_j \leq \text{MED}(D_1, \dots, D_n) + 5\text{MAD}(D_1, \dots, D_n)$ , and apply the statistical technique on the  $n_R$  cases. For variants,  $R$  code, and more explanation, see Olive (2025).

## 3 CONCLUSIONS

Outliers are important, and outlier resistant methods should be used for low and high dimensional statistics.

### References

- Abid, A.M., Quaye, P.A., and Olive, D.J. (2025), “A High Dimensional Omnibus Regression Test,” *Stats*, 8, 107.
- Basa, J., Cook, R.D., Forzani, L., and Marcos, M. (2024), “Asymptotic Distribution of One-Component Partial Least Squares Regression Estimators in High Dimensions,” *The Canadian Journal of Statistics*, 52, 118-130.
- Cook, R.D., and Forzani, L. (2024), *Partial Least Squares Regression: and Related Dimension Reduction Methods*, Chapman and Hall/CRC, Boca Raton, FL.
- Cook, R.D., Helland, I.S., and Su, Z. (2013), “Envelopes and Partial Least Squares Regression,” *Journal of the Royal Statistical Society, B*, 75, 851-877.
- Fan, J., and Lv, J. (2008), “Sure Independence Screening for Ultrahigh Dimensional Feature Space,” *Journal of the Royal Statistical Society, B*, 70, 849-911.
- Fan, J., and Song, R. (2010), “Sure Independence Screening in Generalized Linear Models with np-Dimensionality,” *The Annals of Statistics*, 38, 3217-3841.
- Mallows, C. (1973), “Some Comments on  $C_p$ ,” *Technometrics*, 15, 661-676.
- Olive, D.J. (2025a), “Some Useful Techniques for High Dimensional Statistics,” *Stats*, 8, 60.
- Olive, D.J. (2025b), “David J. Olive’s Contribution to the Discussion of “On Optimal Linear Prediction” by I. Helland,” *Scandinavian Journal of Statistics*, to appear. <http://doi.org/10.1111/sjos.70037>
- Olive, D.J., Alshammari, A.A., Pathiranage, K.G., and Hettige, L.A.W. (2025), “Testing with the One Component Partial Least Squares and the Marginal Maximum Likelihood Estimators,” *Communications in Statistics: Theory and Methods*, to appear.
- Olive, D.J., and Zhang, L. (2025), “One Component Partial Least Squares, High Dimensional Regression, Data Splitting, and the Multitude of Models,” *Communications in Statistics: Theory and Methods*, 54, 130-145.