

Expressions for Shrinkage Factors of PLS Estimator

Jong-Duk Kim¹⁾

Abstract

Partial least squares regression (PLS) is a biased, non-least squares regression method and is an alternative to the ordinary least squares regression (OLS) when predictors are highly collinear or predictors outnumber observations. One way to understand the properties of biased regression methods is to know how the estimators shrink the OLS estimator. In this paper, we introduce an expression for the shrinkage factor of PLS and develop a new shrinkage expression, and then prove the equivalence of the two representations. We use two near-infrared (NIR) data sets to show general behavior of the shrinkage and in particular for what eigendirections PLS expands the OLS coefficients.

Keywords : 고유방향, 부분최소제곱회귀, 축소인자

1. 서론

다음의 선형회귀모형을 고려한다.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1.1)$$

여기서 \mathbf{y} 는 $n \times 1$ 반응벡터, \mathbf{X} 는 $n \times p$ 설계행렬, $\boldsymbol{\beta}$ 는 $p \times 1$ 모수벡터, $\boldsymbol{\epsilon}$ 은 $n \times 1$ 오차(noise)벡터로 평균이 $\mathbf{0}$, 분산이 $\sigma^2 \mathbf{I}$ 라고 가정한다. \mathbf{X} 의 계수는 r 로서 \mathbf{X} 가 완전 열계수가 아닐 수 있는 일반적인 경우를 가정한다. \mathbf{X} 와 \mathbf{y} 는 중심화 또는 표준화된 것으로 간주한다.

비정칙값 분해(SVD, singular value decomposition)는 회귀추정량을 이해하는데 핵심 도구이며 이것을 요약하면 다음과 같다. 행렬 \mathbf{X} 의 SDV는 $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}'$ 인데 여기서 \mathbf{U} 는 $\mathbf{X}\mathbf{X}'$ 의 고유값으로 구성된 $n \times r$ 행렬, \mathbf{D} 는 대각원소가 \mathbf{X} 의 비정칙값이 크기순으로 나열된 $r \times r$ 대각행렬, \mathbf{V} 는 $\mathbf{X}'\mathbf{X}$ 의 고유값으로 구성된 $p \times r$ 행렬이다. 즉 $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_r)$, $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_r)$, $\mathbf{D} = \text{diag}(\delta_1, \dots, \delta_r)$ 로서 $\mathbf{U}'\mathbf{U} = \mathbf{I}_r$, $\mathbf{V}'\mathbf{V} = \mathbf{I}_r$ 이다. 행렬 \mathbf{U}

1) 부산시 남구 우암동 산 55-1, 부산외국어대학교 응용통계학과 교수
E-mail : jdkim@pufs.ac.kr

는 \mathbf{X} 의 주왼쪽비정칙벡터(principal left singular vectors)의 행렬, 행렬 \mathbf{V} 는 \mathbf{X} 의 주오른쪽비정칙벡터(principal right singular vectors)의 행렬이라 부른다. 그리고

$$\mathbf{X}'\mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}' = \sum_{i=1}^r \lambda_i \mathbf{v}_i \mathbf{v}_i' \text{이며 } \lambda_i = \delta_i^2 \text{이다.}$$

OLS 추정량은 선형이다. 회귀모형 (1.1)에서 \mathbf{X} 가 완전열계수이면 OLS 추정량은 유일하나 완전열계수가 되지 못하면 $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ 은 과대모수화되어 무수히 많은 $\hat{\boldsymbol{\beta}}$ 가 존재한다. 그렇지만 여전히 $\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|$ 의

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{ols} &= (\mathbf{X}'\mathbf{X})^+ \mathbf{X}'\mathbf{y} \\ \hat{\mathbf{y}}_{ols} &= \mathbf{X}(\mathbf{X}'\mathbf{X})^+ \mathbf{X}'\mathbf{y} \end{aligned}$$

을 사용할 수 있다. 여기서 첨자 $+$ 는 Moore-Penrose 일반화역행렬을 의미한다. $\hat{\boldsymbol{\beta}}_{ols}$ 은 유일하지 않지만 $\hat{\mathbf{y}}_{ols}$ 은 유일하다. OLS 추정량은 다음과 같이 고유벡터의 선형 결합으로 나타낼 수 있다.

$$\hat{\boldsymbol{\beta}}_{ols} = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}'\mathbf{y} = \sum_{j=1}^r \hat{\alpha}_j \mathbf{v}_j$$

여기서 $\hat{\alpha}_j = \mathbf{u}_j'\mathbf{y}/\delta_j$ 이다. MSE를 줄이기 위해 OLS 추정량에서 분산을 크게 만드는 고유방향을 축소시키는 방법을 생각해볼 수 있다. 이렇게 하면 당연히 편향이 발생한다. 따라서 우리는 편향의 증가량이 분산의 감소량에 비해 작기를 기대하면서 OLS 추정량의 축소를 시도한다.

$\boldsymbol{\beta}$ 의 축소추정량의 일반식으로 다음을 고려한다(Frank와 Friedman, 1993).

$$\hat{\boldsymbol{\beta}}_{shr} = \sum_{j=1}^r f_{j(\cdot)} \hat{\alpha}_j \mathbf{v}_j$$

$f_{j(\cdot)}$ 는 축소인자(shrinkage factor)라 부르기로 한다. 축소인자가 \mathbf{y} 에 의존하지 않으면 $\hat{\boldsymbol{\beta}}_{shr}$ 는 \mathbf{y} 에 선형이며, $f_{j(\cdot)} \neq 1$ 인 인자는 i 번째 성분의 편향을 증가시킨다. j 번째 성분의 분산은 $f_{j(\cdot)} < 1$ 이면 감소하고 $f_{j(\cdot)} > 1$ 이면 증가한다. 따라서 분산의 감소량에 비해 편향의 증가량이 작은 축소 방법이 바람직할 것이다. 이 축소추정량의 행렬 형은 $\hat{\boldsymbol{\beta}}_{shr} = \mathbf{V}\mathbf{F}^{(m)}\hat{\boldsymbol{\alpha}}$ 으로 쓸 수 있는데, 여기서 \mathbf{V} 는 \mathbf{X} 의 주오른쪽비정칙벡터의 행렬, $\mathbf{F}^{(m)} = \text{diag}(f_{1(\cdot)}, \dots, f_{r(\cdot)})$, $\hat{\boldsymbol{\alpha}} = (\hat{\alpha}_1, \dots, \hat{\alpha}_r)'$ 이다.

2. 부분최소제곱회귀

부분최소제곱회귀(partial least squares regression, PLS)는 H. Wold(1975)에 의해 개발되었으며 반복 알고리즘인 NIPALS 알고리즘으로 주어졌다. 알고리즘에서 PLS의 인자벡터 \mathbf{t}_k 는 각각 직교성을 유지하면서 순차적으로 한 개씩 계산되며 최대 r 개까지 구해질 수 있다. 인자수 m 개의 PLS 추정량은

$$\hat{\beta}_{pls}^{(m)} = \mathbf{W}_m (\mathbf{W}_m' \mathbf{X}' \mathbf{X} \mathbf{W}_m)^{-1} \mathbf{W}_m' \mathbf{X}' \mathbf{y}$$

또는

$$\hat{\beta}_{pls}^{(m)} = \mathbf{W}_m (\mathbf{P}_m' \mathbf{W}_m)^{-1} \mathbf{q}_m$$

으로 쓸 수 있는데, 여기서 $\mathbf{W}_m = (\mathbf{w}_1, \dots, \mathbf{w}_m)$ 은 가중값행렬, $\mathbf{P}_m = (\mathbf{p}_1, \dots, \mathbf{p}_m)$ 은 \mathbf{X} 의 적재행렬, $\mathbf{q}_m = (q_1, \dots, q_m)'$ 은 \mathbf{y} 의 적재벡터이다. 이와 같이 PLS의 해는 반복 알고리즘으로 주어지고 인자를 구할 때 \mathbf{X} 뿐만 아니라 \mathbf{y} 의 정보도 이용하기 때문에 복잡한 비선형이 되어 그 통계적 특성을 이해하기 어렵다. 이들 특성을 규명하기 위한 많은 노력이 이루어졌다.

PLS의 비반복적 해도 구해졌다(Helland, 1988). $m \geq 1$ 에서 다음의 Krylov 공간을 고려한다.

$$\mathbf{K}_m = \text{span}(\mathbf{X}'\mathbf{y}, (\mathbf{X}'\mathbf{X})\mathbf{X}'\mathbf{y}, \dots, (\mathbf{X}'\mathbf{X})^{m-1}\mathbf{X}'\mathbf{y})$$

그러면 위의 PLS 알고리즘에서의 가중값벡터 집합 $\{\mathbf{w}_1, \dots, \mathbf{w}_m\}$ 에 의해 생성되는 공간과 Krylov 공간은 동일하다. m 개 잠재인자의 PLS 추정량은 다음의 최적화의 해이다.

$$\begin{aligned} \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\| \\ \text{such that } \beta \in \mathbf{K}_m \end{aligned}$$

$p \times m$ 행렬 \mathbf{R}_m 을 열들이 부공간 \mathbf{K}_m 을 생성하는 한 행렬이라 하면 m 개 인자의 PLS 추정량은

$$\hat{\beta}_{pls}^{(m)} = \mathbf{R}_m (\mathbf{R}_m' \mathbf{X}' \mathbf{X} \mathbf{R}_m)^{-1} \mathbf{R}_m' \mathbf{X}' \mathbf{y}$$

이다.

m 개 인자의 PLS의 비반복적 해는 다음의 형태도 가능하다(Kim, 2003; Kim과 Moon, 2005).

$$\hat{\beta}_{pls}^{(m)} = \mathbf{V} \mathbf{D}^{-1} \mathbf{G}_{(m)} \mathbf{U}' \mathbf{y}$$

여기서 $G_{(m)} = Q_m(Q_m'Q_m)^{-1}Q_m'$, $Q_m = (D^2U'y, D^4U'y, \dots, D^{2m}U'y)$ 이고, V , D , U 는 비정칙값 분해 $X = UDV'$ 에서 나온 것이다.

3. PLS 추정량의 축소인자

Frank와 Friedman(1993)은 $n \times p$ 행렬 X 가 완전열계수인 경우에 PLS 해의 축소 인자의 식을 증명 없이 보였다. 여기서는 X 가 불완전열계수일 수 있는 $\text{rank}(X) = r$ 인 경우의 축소인자의 식으로 바꾸고 증명을 첨부한다.

정리 1. m 개 인자의 PLS 추정량의 축소인자는 다음과 같다.

$$f_{j(pls)}^{(m)} = \sum_{k=1}^m \delta_j^{2k} b_k, \quad j = 1, 2, \dots, r \quad (3.1)$$

그리고 벡터 $b = (b_1, \dots, b_m)'$ 은 $b = S^{-1}s$ 인데 여기서 S 는 (k, l) 번째 원소가 $S_{kl} = \sum_{j=1}^r \delta_j^{2(k+l+1)} \hat{\alpha}_j^2$ 인 $m \times m$ 행렬이고 s 는 k 번째 원소가 $s_k = \sum_{j=1}^r \delta_j^{2(k+1)} \hat{\alpha}_j^2$ 인 $m \times 1$ 벡터이다.

증명. m 개 인자의 PLS 추정량

$$\hat{\beta}_{pls}^{(m)} = W_m(W_m'X'XW_m)^{-1}W_m'X'y$$

에서 가중값 행렬 W_m 은 Krylov 행렬 K_m 의 QR-분해인 $K_m = W_mR$ 으로 얻어지며 따라서 $W_m = K_mR^{-1}$ 이다(Di Ruscio, 2000). 그리고 행렬 X 가 불완전열계수일 때 $VV' \neq I$ 이지만 $VV'W_m = W_m$ 은 성립하는데, 이것은 $K_m = VD^{-1}Q_m$ 을 이용하면

$$VV'W_m = VV'K_mR^{-1} = VV'VD^{-1}Q_mR^{-1} = K_mR^{-1} = W_m$$

으로 바로 확인된다($Q_m = (D^2U'y, D^4U'y, \dots, D^{2m}U'y)$). 이상의 두 결과에 의해 PLS의 해벡터는

$$\hat{\beta}_{pls}^{(m)} = VV'K_m(K_m'X'XK_m)^{-1}K_m'X'y$$

으로 쓸 수 있다. 이제 $S \equiv K_m'X'XK_m$, $s \equiv K_m'X'y$ 로 두고 또한 $b = S^{-1}s$ 로 두면

$$\hat{\beta}_{pls}^{(m)} = VV'K_m b$$

이다. 여기서 $V'K_m b = (DU'y, D^3U'y, \dots, D^{2m-1}U'y)b$ 이며 따라서 이것의 j 번째 원소는 $(\delta_j, \delta_j^3, \dots, \delta_j^{2m-1})u_j'yb$ 이며 이것은 또한 $(\delta_j^2, \delta_j^4, \dots, \delta_j^{2m})bu_j'y/\delta_j$ 이다. 따라서

$r \times 1$ 벡터 $V'K_m b$ 는 j 번째 원소가 $(\delta_j^2, \delta_j^4, \dots, \delta_j^{2m})b$ 인 $r \times r$ 대각행렬과 원소가 $\hat{\alpha}_j = \mathbf{u}_j' \mathbf{y} / \delta_j$ 인 $r \times 1$ 인 벡터의 곱으로 나타낼 수 있다. 이 대각행렬과 벡터를 각각 $F^{(m)}$ 와 $\hat{\alpha}$ 으로 표시하면 PLS 추정량은

$$\hat{\beta}_{pls}^{(m)} = V F^{(m)} \hat{\alpha}$$

으로 표현되며 따라서 $\mathbf{b} = (b_1, b_2, \dots, b_m)'$ 로 두면 축소인자의 식은

$$f_j^{(m)} = \sum_{k=1}^m \delta_j^{2k} b_k, \quad j = 1, 2, \dots, r$$

이다. $\mathbf{b} = S^{-1} \mathbf{s}$ 에서 행렬 S 의 (k, l) 번째 원소는 다음과 같이 구해진다.

$$S = K_m' X' X K_m = (K_m' V D)(D V' K_m)$$

에서 $K_m' V D$ 의 k 번째 행은 $(\delta_1^{2k+1} \hat{\alpha}_1, \dots, \delta_r^{2k+1} \hat{\alpha}_r)$ 이고 따라서 S 의 (k, l) 번째 원소는

$$(\delta_1^{2k+1} \hat{\alpha}_1)(\delta_1^{2l+1} \hat{\alpha}_1) + \dots + (\delta_r^{2k+1} \hat{\alpha}_r)(\delta_r^{2l+1} \hat{\alpha}_r) = \sum_{j=1}^r \delta_j^{2(k+l+1)} \hat{\alpha}_j^2$$

이다. 행렬 s 의 k 번째 원소는 다음과 같이 구해진다.

$$\mathbf{s} = K_m' X' \mathbf{y} = K_m' V D U' \mathbf{y} = K_m' V D^2 \hat{\alpha}$$

이고, $K_m' V D^2$ 은 앞서 나온 $K_m' V D$ 와 대각행렬 D 의 곱이므로 이것의 k 번째 행은 $(\delta_1^{2(k+1)} \hat{\alpha}_1, \dots, \delta_r^{2(k+1)} \hat{\alpha}_r)$ 이고 따라서 $K_m' V D^2 \hat{\alpha}$ 의 k 번째 값은

$$\delta_1^{2(k+1)} \hat{\alpha}_1^2 + \dots + \delta_r^{2(k+1)} \hat{\alpha}_r^2 = \sum_{j=1}^r \delta_j^{2(k+1)} \hat{\alpha}_j^2$$

이다. ■

Frank와 Friedman(1993)은 축소인자 $f_{j(pls)} > 1$ 인 것도 있음을 언급하였다. 즉 어떤 고유방향으로는 팽창하는 것도 있다는 것이다. PLS의 축소인자는 다음과 같이 쓸 수도 있다.

정리 2. m 개 인자의 PLS의 축소인자는 다음과 같다.

$$f_{j(pls)}^{(m)} = \frac{\sum_{l=1}^r g_{jl} \mathbf{u}_l' \mathbf{y}}{\mathbf{u}_j' \mathbf{y}}, \quad j = 1, 2, \dots, r \quad (3.2)$$

단 g_{jl} 은 $\mathbf{G}^{(m)}$ 의 (j,l) 번째 원소, \mathbf{u}_j 는 \mathbf{U} 의 i 번째 열, $\mathbf{G}_{(m)} = \mathbf{Q}_m (\mathbf{Q}_m' \mathbf{Q}_m)^{-1} \mathbf{Q}_m'$, $\mathbf{Q}_m = (D^2 \mathbf{U}' \mathbf{y}, D^4 \mathbf{U}' \mathbf{y}, \dots, D^{2m} \mathbf{U}' \mathbf{y})$ 이다.

증명. 이것의 증명은 2절에서 언급된 PLS의 해벡터 $\hat{\beta}_{pls}^{(m)} = \mathbf{V} \mathbf{D}^{-1} \mathbf{G}_{(m)} \mathbf{U}' \mathbf{y}$, $\mathbf{G}_{(m)} = \mathbf{Q}_m (\mathbf{Q}_m' \mathbf{Q}_m)^{-1} \mathbf{Q}_m'$, $\mathbf{Q}_m = (D^2 \mathbf{U}' \mathbf{y}, D^4 \mathbf{U}' \mathbf{y}, \dots, D^{2m} \mathbf{U}' \mathbf{y})$ 에서 시작한다. 이 해벡터 식에서 $\mathbf{G}_{(m)}$ 의 j 번째 행벡터를 \mathbf{g}_j' 이라 하면 $\mathbf{D}^{-1} \mathbf{G}_{(m)}$ 은 j 번째 행벡터가 \mathbf{g}_j' / δ_j 인 $r \times r$ 인 행렬이고, $\mathbf{D}^{-1} \mathbf{G}_{(m)} \mathbf{U}' \mathbf{y}$ 는 j 번째 원소가 $\mathbf{g}_j' \mathbf{U}' \mathbf{y} / \delta_j$ 인 $r \times 1$ 인 벡터이다. 이것은 j 번째 대각원소가 $\mathbf{g}_j' \mathbf{U}' \mathbf{y} / \delta_j \mathbf{u}_j' \mathbf{y}$ 인 $r \times r$ 대각행렬과 j 번째 원소가 $\mathbf{u}_j' \mathbf{y}$ 인 $r \times 1$ 인 벡터로 분해되며, 전자의 대각행렬은 다시 j 번째 대각원소가 $\mathbf{g}_j' \mathbf{U}' \mathbf{y} / \mathbf{u}_j' \mathbf{y}$ 인 $r \times r$ 대각행렬과 j 번째 원소가 $1/\delta_j$ 인 $r \times r$ 대각행렬로 분해된다. 여기서

$$\mathbf{F}^{(m)} = \text{diag}(f_1^{(m)}, \dots, f_r^{(m)}), \quad \text{단 } f_j^{(m)} = \frac{\mathbf{g}_j' \mathbf{U}' \mathbf{y}}{\mathbf{u}_j' \mathbf{y}}$$

로 두면 바로

$$\hat{\beta}_{pls}^{(m)} = \mathbf{V} \mathbf{F}^{(m)} \mathbf{D}^{-1} \mathbf{U}' \mathbf{y}$$

을 얻는다. g_{jl} 을 $\mathbf{G}_{(m)}$ 의 (j,l) 번째 원소라 하면 $\mathbf{g}_j' \mathbf{U}' \mathbf{y} = \sum_{l=1}^r g_{jl} \mathbf{u}_l' \mathbf{y}$ 이고 따라서 증명된다. ■

식 (3.2)의 축소인자와 식 (3.1)의 축소인자의 동일성은 다음과 같이 보일 수 있다. j 번째 축소인자는

$$f_j^{(m)} = \frac{\sum_{l=1}^r g_{jl} \mathbf{u}_l' \mathbf{y}}{\mathbf{u}_j' \mathbf{y}} = \frac{\sum_{l=1}^r g_{jl} \delta_l \hat{\alpha}_l}{\delta_j \hat{\alpha}_j}$$

으로 쓸 수 있다. 그리고 \mathbf{q}_j' 을 \mathbf{Q}_m 의 j 번째 행이라 하면

$$g_{jl} = \mathbf{q}_j' (\mathbf{Q}_m' \mathbf{Q}_m)^{-1} \mathbf{q}_l$$

이다. 따라서

$$\begin{aligned} f_j^{(m)} &= \frac{1}{\delta_j \hat{\alpha}_j} \sum_{l=1}^r \mathbf{q}_j' (\mathbf{Q}_m' \mathbf{Q}_m)^{-1} \mathbf{q}_l \delta_l \hat{\alpha}_l \\ &= \frac{1}{\delta_j \hat{\alpha}_j} \mathbf{q}_j' (\mathbf{Q}_m' \mathbf{Q}_m)^{-1} \sum_{l=1}^r \mathbf{q}_l \mathbf{u}_l' \mathbf{y} \end{aligned}$$

이다. 여기서 $\mathbf{Q}_m = (D^2 \mathbf{U}' \mathbf{y}, D^4 \mathbf{U}' \mathbf{y}, \dots, D^{2m} \mathbf{U}' \mathbf{y}) = D \mathbf{V}' \mathbf{K}_m$ 이므로 $\mathbf{Q}_m' \mathbf{Q}_m$ 은 정리 1에서의 \mathbf{S} 와 동일하고, $\sum_{l=1}^r \mathbf{q}_l \mathbf{u}_l' \mathbf{y} = \mathbf{Q}_m' \mathbf{U}' \mathbf{y} = \mathbf{K}_m' \mathbf{V} D \mathbf{U}' \mathbf{y} = \mathbf{K}_m' \mathbf{X}' \mathbf{y}$ 는 정리 1에서의 \mathbf{s} 와 동일하다. 그리고 남은 앞부분은

$$\frac{1}{\delta_j \hat{\alpha}_j} \mathbf{q}_j' = \frac{1}{\delta_j \hat{\alpha}_j} (\delta_j^3 \hat{\alpha}_j, \delta_j^5 \hat{\alpha}_j, \dots, \delta_j^{2m+1} \hat{\alpha}_j) = (\delta_j^2, \delta_j^4, \dots, \delta_j^{2m})$$

으로 정리 1의 $\mathbf{V}' \mathbf{K}_m$ 와 동일하다. 따라서 이 축소인자와 정리 1의 축소인자는 동일함이 확인된다.

Frank와 Friedman(1993)은 m 개 인자의 PLS는 m 번째 고유값과 근처에서 계수가 팽창한다고 결론을 내렸다. 그러나 Butler과 Denham(2000)은 다른 예를 통해 항상 그렇게 되는 것은 아니라고 하였다.

축소인자식 (3.2)에서 j 번째 항은 $f_{j(pls)}^{(m)} = g_{jj} + \sum_{l \neq j}^r g_{jl} \mathbf{u}_l' \mathbf{y} / \mathbf{u}_j' \mathbf{y}$ 로 쓸 수 있다. 그리고 $\mathbf{G}_{(m)}$ 은 대칭 멱등행렬이므로 (i) $0 \leq g_{jj} \leq 1$ (ii) $g_{jj} = \sum_{l=1}^r g_{jl}^2$ 이 성립된다. j 번째 항 $f_{j(pls)}^{(m)}$ 의 값이 1보다 크면 그 고유방향으로 팽창된다.

4. 수치 예

두 개의 실제 데이터 셋으로 PLS의 축소인자의 값을 구해보도록 한다. 이 값의 계산은 본 논문에서 유도된 식 (3.2)를 이용하였으며, 모든 계산은 SAS/IML로 수행되었다.

데이터 셋 1 (Fearn, 1983): x 값은 NIR 반사의 여섯 파장(설명변수)에서 나온 것이고 y 값은 밀에 포함된 단백질 비율(%)이다. 총 표본 수는 24이다. 모든 변수는 사전 처리로서 단위길이표준화(상관변환)를 하였다. 이 변환된 $\mathbf{X}' \mathbf{X}$ 의 고유값은

$$(5.8679606, 0.1007471, 0.0186959, 0.0122150, 0.0002301, 0.0001513)$$

으로 매우 높은 공선성의 데이터임을 알 수 있다. <표 1>은 PLS의 각 인자 수에서 고유방향에 따른 축소인자의 값이다.

<표 1> Fearn 데이터에 적용된 PLS의 축소인자 $f_{j(pls)}^{(m)}$ 의 값 (볼드체로 한 것은 1보다 큰 값)

고유방향 j	인자수 m					
	1	2	3	4	5	6
1	1.0218501	0.9999097	1.0000000	1.0000000	1.0000000	1.0000000
2	0.0175442	1.4820444	0.9997565	1.0000011	1.0000000	1.0000000
3	0.0032557	0.2788941	1.0049619	0.9998805	1.0000001	1.0000000
4	0.0021271	0.1824156	0.6995192	1.0112845	0.9999923	1.0000000
5	0.0000401	0.0034425	0.0146740	0.0338496	1.2184397	1.0000000
6	0.0000264	0.0022647	0.0096600	0.0223391	0.8107309	1.0000000

데이터 셋 2 (Brereton, 2003): 각 샘플은 27개의 파장(220nm에서 350nm 사이의 파장을 5nm 간격)과 화합물에서의 Pyrene의 농도(10 PAHs)를 측정된 것으로 구성되어 있다. 총 샘플(EAS 스펙트럼)수는 25이다. 여기서도 모든 변수는 사전 처리로서 상관변환을 하였다. 이 변환된 데이터의 $\mathbf{X}'\mathbf{X}$ 의 양의 고유값은

(20.953006, 2.5443417, 1.3252614, 1.0724578, 0.6091839, 0.2301074, 0.0833184, 0.0651115, 0.0374254, 0.0324147, 0.0201557, 0.0134545, 0.0061830, 0.0022937, 0.0018606, 0.0011580, 0.0009120, 0.0005756, 0.0002793, 0.0002632, 0.0001596, 0.0000396, 0.0000229, 0.0000144)

이다. 따라서 크기 25×27 의 \mathbf{X} 의 계수(rank)는 24로 불완전계수이다. 모든 경우의 축소인자의 값은 <표 2>에 주어져 있다.

<표 2> Brereton의 데이터에 적용된 PLS의 축소인자 $f_{j(pls)}^{(m)}$ 의 값 (볼드체로 한 것은 1보다 큰 값)

고유방향 j	인자수 m							
	1	2	3	4	5	6	7	8
1	1.9923154	0.9966199	1.0000510	0.9999996	1.0000000	1.0000000	1.0000000	1.0000000
2	0.2419286	1.0346828	0.9942856	1.0006016	0.9999287	1.0000094	0.9999999	1.0000000
3	0.1260124	0.5704465	1.2793579	0.8342482	1.0552426	0.9820824	1.0003673	0.9999981
4	0.1019746	0.4669183	1.1723000	1.2627035	0.8496363	1.0698040	0.9982125	1.0000115
5	0.0579242	0.2707269	0.8131543	1.6260574	1.2531818	0.4276831	1.0269797	0.9996865
6	0.0218797	0.1039634	0.3543623	1.0085407	1.1819695	1.1777862	0.9730427	1.0009072
7	0.0079223	0.0378821	0.1350767	0.4308210	0.5842390	0.8254315	1.2171041	0.9676150
8	0.0061911	0.0296272	0.1062199	0.3433620	0.4737220	0.6951742	1.0864045	1.2294962
9	0.0035586	0.0170496	0.0616325	0.2032826	0.2877938	0.4465271	0.7585403	1.2990434
10	0.0030821	0.0147701	0.0534716	0.1770031	0.2517504	0.3944651	0.6798825	1.2408421
11	0.0019165	0.0091890	0.0333873	0.1114949	0.1603729	0.2573277	0.4591870	0.9682218
12	0.0012793	0.0061356	0.0223374	0.0749513	0.1084702	0.1762897	0.3204679	0.7268335
13	0.0005879	0.0028205	0.0102903	0.0347071	0.0505613	0.0833118	0.1544751	0.3775991
14	0.0002181	0.0010465	0.0038224	0.0129277	0.0188994	0.0313693	0.0587751	0.1492925
15	0.0001769	0.0008489	0.0031010	0.0104912	0.0153434	0.0254877	0.0478103	0.1219533
16	0.0001101	0.0005284	0.0019305	0.0065344	0.0095626	0.0159058	0.0298924	0.0767686
17	0.0000867	0.0004161	0.0015205	0.0051475	0.0075347	0.0125384	0.0235793	0.0606994
18	0.0000547	0.0002626	0.0009598	0.0032500	0.0047586	0.0079238	0.0149146	0.0385186
19	0.0000266	0.0001274	0.0004657	0.0015774	0.0023103	0.0038490	0.0072506	0.0187787
20	0.0000250	0.0001201	0.0004389	0.0014865	0.0021772	0.0036274	0.0068334	0.0177010
21	0.0000152	0.0000728	0.0002661	0.0009014	0.0013203	0.0022001	0.0041458	0.0107498
22	3.7695E-6	0.0000181	0.0000661	0.0002240	0.0003281	0.0005469	0.0010308	0.0026758
23	2.1751E-6	0.0000104	0.0000382	0.0001292	0.0001893	0.0003156	0.0005949	0.0015445
24	1.3682E-6	6.5655E-6	0.0000240	0.0000813	0.0001191	0.0001985	0.0003742	0.0009716

<표 2> (연속)

고유방향 j	인자수 m							
	9	10	11	12	13	14	15	16
1	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
2	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
3	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
4	0.9999998	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
5	1.0000077	0.9999998	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
6	0.9999310	1.0000052	0.9999999	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
7	1.0153703	0.9960846	1.0003526	0.9999750	1.0000040	0.9999999	1.0000000	1.0000000
8	0.5718795	1.1554998	0.9801095	1.0021593	0.9995240	1.0000160	0.9999999	1.0000000
9	0.9440552	0.8567916	1.5867934	0.2623218	1.3752734	0.9768609	1.0001560	0.9999975
10	1.0516061	0.9184716	1.0194313	1.0153412	0.9897570	1.0007440	0.9999942	1.0000001
11	1.1599374	1.1420655	0.5666853	0.9726728	1.0805571	0.9896328	1.0001324	0.9999959
12	1.0287419	1.1462505	1.1239025	0.9904235	0.9703944	1.0065972	0.9998709	1.0000063
13	0.6294814	0.8086885	1.5402409	1.6016306	1.4284747	0.5262981	1.0220557	0.9972744
14	0.2698295	0.3735544	0.9346358	1.1834574	1.2117964	1.1760879	0.9577953	1.0336192
15	0.2223497	0.3103555	0.7981199	1.0326740	1.0761210	1.1416687	1.0824059	0.7858279
16	0.1419477	0.2007686	0.5392025	0.7222548	0.7745486	0.9395089	1.1260308	0.9495891
17	0.1127851	0.1602589	0.4368573	0.5922269	0.6415391	0.8131933	1.0518961	1.0723384
18	0.0720490	0.1030221	0.2865294	0.3948116	0.4336212	0.5824223	0.8317614	1.1114234
19	0.0353311	0.0507994	0.1437722	0.2009436	0.2233895	0.3151033	0.4885267	0.8163537
20	0.0333140	0.0479136	0.1357323	0.1898528	0.2111988	0.2986851	0.4650820	0.7862053
21	0.0202728	0.0292134	0.0832590	0.1170343	0.1307458	0.1880095	0.3009302	0.5471239
22	0.0050582	0.0073052	0.0209655	0.0296391	0.0332742	0.0487651	0.0805295	0.1587408
23	0.0029205	0.0042192	0.0121206	0.0171487	0.0192651	0.0283084	0.0469497	0.0935731
24	0.0018375	0.0026551	0.0076311	0.0108011	0.0121383	0.0178599	0.0296852	0.0594933

<표 2> (연속)

고유방향 j	인자수 m							
	17	18	19	20	21	22	23	24
1	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
2	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
3	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
4	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
5	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
6	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
7	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
8	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
9	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
10	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
11	1.0000001	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
12	0.9999999	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
13	1.0001276	0.9999896	1.0000006	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
14	0.9940270	1.0016872	0.9997091	1.0000071	1.0000000	1.0000000	1.0000000	1.0000000
15	1.0576994	0.9774143	1.0051769	0.9998403	1.0000007	1.0000000	1.0000000	1.0000000
16	0.8954743	1.1192965	0.9410151	1.0031443	0.9999758	1.0000026	1.0000000	1.0000000
17	1.0070854	0.9591509	1.0341812	0.9975541	1.0000250	0.9999965	1.0000000	1.0000000
18	1.2314917	1.0170466	0.5046608	1.0667690	0.9987404	1.0003233	0.9999972	1.0000000
19	1.1098764	1.2008956	1.1287361	0.8792057	1.0289581	0.9721995	1.0005145	1.0000000
20	1.0810069	1.1922335	1.1904082	0.9715945	0.9444149	1.0627281	0.9987612	1.0000000
21	0.8083214	1.0100389	1.4399625	1.6983135	0.9947972	0.9564494	1.0015162	1.0000000
22	0.2542736	0.3663714	0.7302764	1.2517745	1.3727341	1.2592333	0.9230022	1.0000000
23	0.1515541	0.2226464	0.4623342	0.8285242	0.9918333	1.0038498	1.0215410	1.0000000
24	0.0968961	0.1437455	0.3045902	0.5578935	0.6970128	0.7417211	0.8691119	1.0000000

PLS의 $f_{j(pls)}$ 는 복잡한 비선형이어서 주어진 데이터에 따라 형태가 달라지며 그 성질을 알기가 어렵다. 그러나 <표 1>과 <표 2>의 축소인자 $f_{j(pls)}$ 값을 보면 PLS의 특이한 축소구조가 어느 정도 드러난다. m 개 인자의 모형에서, m 번째 고유방향과 그 전의 몇 고유방향에서 확대와 축소를 반복하는 패턴이 있고 m 번째와 그 이후 한두 개의 고유방향은 연속적으로 확대하고 그 이후는 급속히 축소하는 경향을 보였다. m 개 인자의 모형에서 m 번째 고유방향으로는 뚜렷하게 확대하는 패턴을 보였으나, 일부 그렇지 않는 경우도 나타났으며 데이터 셋 2에서 인자수 $m=9, 10, 11, 12, 16, 20, 21$ 에서는 각각의 9, 10, 11, 12, 16, 20, 21 번째 고유방향으로 확대하지 않고 오히려 축소하였다. 전체적으로, 각 고유방향으로 팽창하는 경우 1 보다 약간 큰 경우가 대부분이나 팽창이 상대적으로 매우 큰 것도 일부 나타났다. <표 1>에서는 $m=2$ 의 두 번째 고유방향의 $f_{j(pls)}$ 값은 1.482이고 <표 2>에서는 $m=1$

의 첫 번째 고유방향과 $m = 20$ 의 21 번째 고유방향의 $f_{j(pls)}$ 값이 각각 1.992와 1.698로 상당한 편차를 보였다. 인자수가 행렬 \mathbf{X} 의 계수인 24일 때는 모든 $f_{j(pls)}$ 값이 1이 되어 OLS가 됨이 확인된다.

참고문헌

1. Brereton, R. (2003). *Chemometrics: Data Analysis for the Laboratory and Chemical Plant*, John Wiley & Sons.
2. Butler, N. and Denham, M. (2000). The Peculiar Shrinkage Properties of Partial Least Squares Regression, *Journal of Royal Statistical Society*, B, 62, Part 3, 585-593.
3. Di Ruscio, D. (2000). A Weighted View on the Partial Least Squares Algorithm, *Automatica*, 36, 831-850.
4. Fearn, T. (1983). A Misuse of Ridge Regression in the Calibration of a Near Infrared Reflectance Instrument, *Journal of Applied Statistics*, 32, 73-79.
5. Frank, I. E. & Friedman, J. H. (1993). A Statistical View of Some Chemometrics Regression Tools (with Discussion), *Technometrics*, 35, 109-148.
6. Helland, I. S. (1988). On the structure of partial least squares regression, *Communications in Statistics - Simulation and Computation*, 17, 581-607.
7. Kim, J. D. (2003). Unified Non-iterative algorithm for Principal Component Regression, Partial Least Squares and Ordinary Least Squares, *Journal of Korean Data & Information Science Society*, 14, 355-366.
8. Kim, J. D. and S. Moon (2005). Connecting Partial Least Squares and Generalized Ridge Regression, *Journal of the Korean Data Analysis Society*, 7, 61-71.
9. Wold, H. (1975). Soft Modelling by Latent Variables; the Nonlinear Iterative Partial Least Squares Approach, in *Perspectives in Probability and Statistics*. Papers in Honour of M. S. Bartlett (ed. J. Gani). Academic Press, New York.

[2006년 9월 접수, 2006년 11월 채택]