

# 연속회분슬러지반응기(SBR)의 품질 예측을 위한 통계학적 방법

## I. 부분 최소 자승법 이론

### 유창규

SBR 은 유기물은 물론 질소와 인을 효과적으로 동시에 제거할 수 있다는 점에서 크게 주목을 받고 있다. 이 생물학적 유기물, 질소, 인 제거 공정에서 제어하고자 하는 주된 변수는 생산물 속의 질소와 인의 농도인데 질소와 인 농도는 실시간으로 직접 측정 및 제어하는 데 어려움이 있다. 실시간 측정의 문제점을 극복하기 위한 시도로써 좀더 신속하고 쉽게 측정되며 추정하고자 하는 변수의 거동과 밀접한 연관을 갖는 다른 공정 변수로부터 간접적으로 그 변수를 추정하는 추정기에 대한 연구가 일찍부터 있어왔는데, 기존의 Principal Component Regression(PCR), Partial Least Squares(PLS) 등과 같은 변수들 간의 상관관계를 고려하는 다변량 통계적 회귀분석 방법이 많이 사용되고 있다. 이 중 PLS 방법이 경험적인 추론 모델을 만드는 보편적인 방법을 제공하며 공정의 변수가 매우 많아서 그들 간에 심한 상관관계가 있을 때도 다른 방법에 비해 좋은 성능을 보임이 보고되었다. 그러나 비선형성을 띠는 공정을 모델링하기에는 한계를 지닌다. 따라서 최근에 많은 관심을 받고 있는 RKHS(Reproducing Kernel Hilbert Space)에서 PLS의 장점을 도입한 KPLS(Kernel partial least square)가 SBR 같은 비선형 공정 모델링에 좋은 성능을 나타낸다. SBR 이 배치 공정이므로 Multiway/Multiblock KPCA/KPLS 를 개발하여 배치가 끝난 후 유기물, 질소, 인 농도를 예측하는 것을 목표로 한다. . 그 후에 배치 공정의 batch-to-batch variation 을 고려한 state-space model 에 기초한 모니터링 및 품질 예측에 관한 연구를 진행할 예정이다. 이에 따라 본 보고서에서는 기존의 부분 최소 자승법의 이론을 먼저 소개한다.

# 부분 최소 자승법 (Partial Least Squares)

## 1. 도입 배경

두 데이터 ( $X$ ,  $Y$ )간의 회귀 분석을 하기 위해 이전에는 MLR (multiple linear regression)으로도 쉽게 분석을 할 수 있었다. 하지만 데이터의 변수가 방대하게 늘어나고 변수간의 correlation 이 강하게 나타남에 따라 MLR 에서 singularity 문제가 나타나게 되어 정확한 분석이 힘들어졌다. Correlation 이 강하고 noise 가 많이 함유된 다차원 데이터를 uncorrelated 된 저 차원의 모델로 해석하기 위해 PCA 가 도입되었으며  $X$  데이터에 PCA 를 적용한 후 그 score 값  $T$  를 구해 이를  $Y$  에 regression 하는 PCR (principal component regression)이 도입되어 MLR 에서 발생했던 singularity 문제를 해결하게 되었으나  $X$  의 variance 만을 설명하는 score  $T$  를 구하여 이를 단순히  $Y$  와 regression 을 하였기 때문에 결과적으로  $X$  와  $Y$  의 관계를 나타내는데 한계가 있게 되었다. 따라서  $X$  의 variance 를 잘 설명하면서  $Y$  와의 correlation 도 최대화할 수 있는 방법을 생각하게 되었는데 이 방법이 바로 PLS 이다.

## 2. 역사적 배경

PLS 접근 방법은 1975 년 무렵부터 Herman Wold 가 복잡한 데이터를 path model 이라 부르는 일련의 행렬들로 모델링 함으로써 유래되었다. 모델에 매개변수를 예측하기 위해 간단하고 효율적인 방법을 적용하였는데 이 방법을 바로 NIPALS (Non-linear Iterative Partial Least Squares)라 불렀으며 PLS 라는 용어는 바로 NIPALS 용어에서 나온 것이다. NIPALS 라는 용어에서 iterative 는 매개변수를 iteration 하여 구하는 것을 말하고 partial 이란 말은 partial regression 을 가리킨다. 1980 년대부터 Svante Wold 와 Herald Martens 에 의해 공학적 데이터에 적합한 간단한 PLS 모델이 개발되었으며 H. Wold et al.은 PLS 에 보다 정확한 의미를 부여하기 위해 'Projection to Latent Structures (PLS)' 라는 말로도 사용하기 시작했다.

## 3. PLS 개요

PCA 에 이론적 바탕을 두고  $X$  와  $Y$  의 관계를 짓는 방법이 바로 PLS 이다.  $X$  와

Y를 mapping 하는 transfer matrix 를 구성하는 방법에는 PLS 외에도 MultipleLinear Regression(MLR)와 Principal Component Regression (PCR) 이 있다. 그러나 가장 널리 응용되고 있는 MLR 은 centering 되고 scaling 된 자료 행렬에 대해 collinearity 와 singularity(Wold, S. et al. , 1984)의 문제를 안고 있고 PCR 은 X의 score vectors(or PCs)에 대해 Y 의 변수들 각각을 regression 시킴으로 collinearity 와 singularity 문제는 발생하지 않지만 X 를 가장 잘 설명하는 첫 번째 PC 에 대한 score vector 가 Y 를 가장 잘 설명하리라는 보장은 없으므로 PCR 은 한계가 있다. 다시 말하면 Y 를 구성하는 데이터들이 크게 상관 관계를 가진다면 PCR 에서처럼 X 에 대한 데이터 공간의 상관 관계를 가장 잘 설명한다고 해서(첫번째 PC) Y 에 대한 데이터 공간의 상관 관계까지도 가장 잘 설명한다고 볼 수는 없다. 오히려 X 에 대한 첫번째 PC 가 가장 못 설명할 수도 있다.

반면 Partial Least Squares(PLS)는 위에서 언급한 문제들을 가장 잘 해결할 수 있는 방법이다. 먼저 X 와 Y 공간 상에서 각각에 대해 PCA 를 적용해서 아래 그림과 같이 분해한다. 이런 관계들을 각각에 대한 outer relation 이라 한다.

$$\begin{array}{c}
 \begin{array}{c} \boxed{X} \\ n \times m \end{array} = \begin{array}{c} \boxed{T} \\ n \times a \end{array} \begin{array}{c} \boxed{P} \\ a \times m \end{array} + \begin{array}{c} \boxed{E} \\ n \times m \end{array} \\
 \\
 \begin{array}{c} \boxed{Y} \\ n \times p \end{array} = \begin{array}{c} \boxed{U} \\ n \times a \end{array} \begin{array}{c} \boxed{Q} \\ a \times p \end{array} + \begin{array}{c} \boxed{F^*} \\ n \times p \end{array}
 \end{array}$$

위 그림에서 n 은 sample 수를, m 은 X의 변수 개수를, p 는 Y의 변수 개수를, a 는 principal component 수를 나타낸다.

PLS 가 PCR 보다 나은 점은 X 의 score vector( $t_h$ )와 Y 의 score vector( $u_h$ ) 사이에  $u_h = b_h t_h$  와 같은 inner relation 을 만들어 서로에 대한 정보를 공유한다는 점이다. 이 inner relation 은 X 의 score vector 들을 Y 의 데이터 공간을 잘 설명할 수 있도록 회전시키는 것을 의미하므로 X 공간에서의 평면이 Y 를 더 잘 예측할 수 있도록 기울어진다는 것을 의미하기도 한다. 그런 다음 X 의 score vector 들에 대해 weight( $w_h$ )를 주어 그들 각각이 Y 의 데이터 공간을 설명하는 정도를 정해준다. 이런 PLS 를 이용하면 score vector 들로 regression 하므로 collinearity 와 singularity 문제도 발생하지 않을 뿐더러 X 와 Y 에게 서로 정보를

공유하고 contribution 의 정도를 weight 를 주어 조정하므로 PCR 에서 생겼던 문제도 발생하지 않는다. 수학적으로 PCA 에서의 loading vector 들은 covariance 행렬( $\mathbf{S}=\mathbf{X}^T\mathbf{X}$ )의 고유 벡터들이듯이 PLS 의 loading vector 들은  $(\mathbf{X}^T\mathbf{Y})(\mathbf{Y}^T\mathbf{X})$ 의 고유 벡터들이다.

#### 4. PLS algorithm

PLS 를 실제 적용할 때는 PCA 에서 사용했던 NIPALS 알고리즘과 비슷한 다음과 같은 알고리즘을 이용한다. 먼저 다음과 같은 outer relation 을 PCA 를 적용해 구한다.

$$\mathbf{X} = \mathbf{TP}^T + \mathbf{E} = \sum_{h=1}^a \mathbf{t}_h \mathbf{p}_h^T$$

$$\mathbf{Y} = \mathbf{UQ}^T + \mathbf{F}^* = \sum_{h=1}^a \mathbf{u}_h \mathbf{q}_h^T + \mathbf{F}^*$$

여기서  $\|\mathbf{F}^*\|$  를 가능한 한 작게 하면서 Y 를 설명할 수 있도록 해야 하며 동시에 X 와 Y 의 유용한 관계를 얻어내는 것이 목표이다. Inner relation 은 각 PC 마다 Y block 의 score  $\mathbf{u}$  와 X block 의 score  $\mathbf{t}$  의 관계를 맺어주는 것이다.

$$\hat{\mathbf{u}}_h = b_h \mathbf{t}_h$$

여기서  $b_h = \frac{\mathbf{u}_h^T \mathbf{t}_h}{\mathbf{t}_h^T \mathbf{t}_h}$  로 구할 수 있으며 이는 MLR 과 PCR 모델에서의 regression coefficient 역할과 같다. 하지만 이렇게 한다고 해서 이 모델이 가장 좋은 solution 을 준다고 볼 수 없다. 왜냐하면 PC 가 각각 block 에서 따로따로 계산되므로 서로 weak relation 을 갖는다. 따라서 X, Y 서로에게 정보를 주어 약간 회전된 component 가 regression line 에 근접하도록 해준다. 이를 고려한 PLS algorithm 은 다음과 같다.

- 1) 일단 X 와 Y 데이터를 mean-centering 하고 scaling 한다.

For each component:

- 2) set  $\mathbf{u}$  equal to a column of Y

In the X block:

3) X의 각 변수를  $\mathbf{u}$ 에 투영시켜  $\mathbf{w}^T$ 를 구한다.

$$\mathbf{w}^T = \mathbf{u}^T \mathbf{X} / \mathbf{u}^T \mathbf{u} \text{ (regress columns of X on u)}$$

4) Normalize  $\mathbf{w}$  to unit length

5) X의 각 sample을  $\mathbf{w}$ 에 투영시켜 score  $\mathbf{t}$ 를 구한다.

$$\mathbf{t} = \mathbf{X} \mathbf{w} / \mathbf{w}^T \mathbf{w}$$

In the Y block:

6) Y의 각 변수를 X block에서 구한  $\mathbf{t}$ 에 투영시켜  $\mathbf{q}^T$ 를 구한다.

$$\mathbf{q}^T = \mathbf{t}^T \mathbf{Y} / \mathbf{t}^T \mathbf{t} \text{ (regress columns of Y on t)}$$

7) Normalize  $\mathbf{q}$  to unit length

8) Y의 각 sample을  $\mathbf{q}$ 에 투영시켜 score  $\mathbf{u}$ 를 구한다.

$$\mathbf{u} = \mathbf{Y} \mathbf{q} / \mathbf{q}^T \mathbf{q}$$

Check convergence:

9) step 5)에서 구한  $\mathbf{t}$ 와 이전 iteration에서 구했던  $\mathbf{t}$ 를 비교하여 같으면 step 10)으로 넘어가고 그렇지 않으면 step 3)으로 간다. 만약 Y block이 단변수이면  $\mathbf{q}=1$ 이라 두고 step 6)부터 9)까지는 생략될 수 있으며 더 이상의 iteration이 필요하지 않는다.

위에서 구한  $\mathbf{t}$ 는 orthogonal하지 않기 때문에 추가적인 algorithm이 필요하다.

10) X loadings:  $\mathbf{p} = \mathbf{X}^T \mathbf{t} / \mathbf{t}^T \mathbf{t}$

11) Normalize  $\mathbf{p}$  to unit length

$$12) \mathbf{t}_{new} = \mathbf{t} \|\mathbf{p}^T\|$$

$$13) \mathbf{w}_{new}^T = \mathbf{w}^T \|\mathbf{p}^T\|$$

여기서  $\mathbf{p}^T$ ,  $\mathbf{q}^T$ ,  $\mathbf{w}^T$ 는 prediction을 위해 저장되어야 한다.

14) inner relation을 위해 regression coefficient  $\mathbf{b}$ 를 구한다.

$$\mathbf{b} = \mathbf{u}^T \mathbf{t} / \mathbf{t}^T \mathbf{t}$$

**Calculate residual matrices:** 각 component  $h$  에 대하여

$$15) \mathbf{E}_h = \mathbf{E}_{h-1} - \mathbf{t}_h \mathbf{p}_h^T, \quad \mathbf{X} = \mathbf{E}_0$$

$$16) \mathbf{F}_h = \mathbf{F}_{h-1} - b_h \mathbf{t}_h \mathbf{q}_h^T, \quad \mathbf{Y} = \mathbf{F}_0$$

17) 다음 component 를 구하기 위해 step2)로 간다. 첫번째 component 를 구한 뒤 앞에서의  $\mathbf{X}$  와  $\mathbf{Y}$  는  $\mathbf{E}_h$  와  $\mathbf{F}_h$  로 교체된다.

step 10) 이전의 과정에서의  $t$  는 서로 정보를 교환함으로써 PCA 의 NIPALS 에서와 달리 직교하지 않는  $t$  이지만 step 10)~12) 과정에서의  $t$  는 이들을 다시 직교하도록 만들어 준 것이다. 직교하는  $t$  는 꼭 있어야 하는 것은 아니지만 PCR 과의 비교를 용이하게 하며 PCA 와의 연결이 쉽다는 장점이 있으므로 이 과정을 알고리즘에 포함시키는 것이 일반적이다. 또한 위 알고리즘은 선형 회귀법(step 14))을 사용한 것이므로 원래 변수들을 변형시키거나 step 14) 과정에 알고 있는 비선형 관계를 집어 넣음으로써 PLS 모델에 비선형성을 줄 수도 있다.

이러한 PLS 알고리즘은 입출력 데이터의 상관성이 가장 큰 방향으로부터 순차적으로  $a$  개의 PC 를 구하는 알고리즘으로써 수렴이 빠르고 각각의 PC 가 직교하는 원리를 이용하였기 때문에 공정이 선형성을 가지고 있으면 다변량 데이터의 차원 감소 및 측정 잡음(measurement noise)제거를 통하여 뛰어난 추정능력을 발휘할 수 있다.

PLS 도 PCA 에서처럼  $a$  개의 PLS 성분을 결정해야 하는데 이때 많이 사용하는 방법은 cross validation 과 PRESS 같은 통계적 방법이다.

#### 4. Properties of the PLS factors

$$1) \|\mathbf{p}_h\| = \|\mathbf{q}_h\| = 1$$

$$2) \mathbf{w}_h^T \text{ 는 orthogonal}$$

$$3) \mathbf{t}_h \text{ 는 orthogonal}$$

4) 결과적으로 앞의 algorithm 에서 구한  $\mathbf{u}$ ,  $\mathbf{q}$ ,  $\mathbf{t}$ ,  $\mathbf{w}$  는 각각  $\mathbf{Y}^T \mathbf{X} \mathbf{X}^T \mathbf{u}$ ,  $\mathbf{Y}^T \mathbf{X} \mathbf{X}^T \mathbf{Y} \mathbf{q}$ ,  $\mathbf{X} \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{t}$ ,  $\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w}$  를 eigenvalue decomposition 했을 때 가장 큰 eigenvalue 에 해당하는 eigenvector 와 동일하다.

## 5. Prediction

모델링에서 구한  $\mathbf{p}^T, \mathbf{q}^T, \mathbf{w}^T, b$  를 이용한다.

새로운  $X$  에 대하여 예측값  $Y$  는 다음과 같이 구한다.

$$\hat{\mathbf{t}}_h = \mathbf{E}_{h-1} \mathbf{w}_h$$

$$\mathbf{E}_h = \mathbf{E}_{h-1} \mathbf{t}_h \mathbf{p}_h^T$$

$$\mathbf{Y} = \mathbf{F}_h = \sum b_h \mathbf{t}_h \mathbf{q}_h^T$$

## 커널 부분 최소 자승법 (Kernel Partial Least Squares in RKHS)

KPLS(Rosipal and Trejo, 2001)는 RKHS 하에서 가능한 high-dimensional feature space 에서 PLS 의 방법을 비선형 회귀모델에 적용방법이다. 다음 보고서에서는 이에 대한 이론 정리와 실제 배치 공정데이터의 품질예측 결과를 기존의 부분최소자승법과 비교한다.

참고문헌)

1. Geladi, P. and Kowalski, B. R., PARTIAL LEAST-SQUARES REGRESSION: A TUTORIAL. *Analytica Chimica Acta*, **185**, 1-17 (1986).
2. Hoskuldsson, A., PLS REGRESSION METHODS, *Journal of Chemometrics*, **2**, 211-228 (1988).
3. Wold, S., Sjostrom, M., Eriksson, L., PLS-regression: a basic tool of chemometrics, *Chemometrics and Intelligent Laboratory Systems*, **58**, 109-130 (2001).
4. Scholkopf, B., Smola A. and Muller, K., Nonlinear component analysis as a kernel eigenvalue problem, *Neural computation*, **10**, 1299-1319 (1998).
5. Rosipal, R. and Trejo, L., Kernel partial least squares regression in reproducing kernel Hilbert space, *Journal of Machine learning research*, **2**, 97-123 (2001).