Collection of Biostatistics Research Archive COBRA Preprint Series

Year 2012

Paper 100

PLS-ROG: Partial least squares with rank order of groups

Hiroyuki Yamamoto*

*Human Metabolome Technologies, Inc., h.yama2396@gmail.com This working paper is hosted by The Berkeley Electronic Press (bepress) and may not be commercially reproduced without the permission of the copyright holder.

http://biostats.bepress.com/cobra/art100

Copyright ©2012 by the author.

PLS-ROG: Partial least squares with rank order of groups

Hiroyuki Yamamoto

Abstract

Partial least squares (PLS), which is an unsupervised dimensionality reduction method, has been widely used in metabolomics. PLS can separate score depend on groups in a low dimensional subspace. However, this cannot use the information about rank order of groups. This information is often provided in which concentration of administered drugs to animals is gradually varies. In this study, we proposed partial least squares for rank order of groups (PLS-ROG). PLS-ROG can consider both separation and rank order of groups.

Partial least squares for rank order of groups (PLS-ROG)

Consider a matrix Y, which is a mean-centered dummy matrix of group information whose elements are 0 or 1. The latent variables t and s are related to the data matrix X and the dummy matrix Y by $\mathbf{t} = \mathbf{X}\mathbf{w}_x$ and $\mathbf{s} = \mathbf{Y}\mathbf{w}_y$. PLS for discrimination [1] is formulated as the optimization problem of maximizing the covariance between the latent explanatory variable t and the latent response variable s:

$$\max \operatorname{cov}(\mathbf{t}, \mathbf{s})$$

subject to $\mathbf{w}_{x}' \mathbf{w}_{x} = 1, \mathbf{w}_{y}' \mathbf{Y}' \mathbf{Y} \mathbf{w}_{y} = 1$ (1)

The differential penalty of mean of the latent variable s in each group is added to the constraint condition of PLS. PLS-ROG is formulated as follows:

 $\max \operatorname{cov}(\mathbf{t}, \mathbf{s})$

subject to $\mathbf{w}_x \mathbf{w}_x = 1$, $\mathbf{w}_y \mathbf{Y} \mathbf{Y} \mathbf{w}_y + \kappa \mathbf{w}_y \mathbf{Y} \mathbf{P} \mathbf{D} \mathbf{D} \mathbf{P} \mathbf{Y} \mathbf{w}_y = 1$

This formulation is similar as our previous study [2]. The differential matrix \mathbf{D} and averaging matrix of groups \mathbf{P} are set for the g class classification problems as follows:

$P = \begin{bmatrix} 1/n_1 & \cdot \\ \cdot & \cdot \end{bmatrix}$	·· 1,	n_1	5	0 ·.	-	0	
$P = \begin{bmatrix} 1 & n_1 \\ 0 & 0 \end{bmatrix}$	0	n_1		••.		0	

By using the method of Lagrange multipliers, the problem in eq. (1) can be reformulated as the maximization of

$$J = \frac{1}{n-1} \mathbf{w}_{\mathbf{x}}' \mathbf{X}' \mathbf{Y} \mathbf{w}_{\mathbf{y}} + \lambda_{\mathbf{x}} (1 - \mathbf{w}_{\mathbf{x}}' \mathbf{w}_{\mathbf{x}}) + \lambda_{\mathbf{y}} (1 - \mathbf{w}_{\mathbf{y}}' \mathbf{Y}' \mathbf{Y} \mathbf{w}_{\mathbf{y}} - \kappa \mathbf{w}_{\mathbf{y}}' \mathbf{Y}' \mathbf{P}' \mathbf{D}' \mathbf{D} \mathbf{P} \mathbf{Y} \mathbf{w}_{\mathbf{y}}).$$
(2)

Partial differentiation of eq. (2) with respect to \mathbf{w}_x and \mathbf{w}_y , followed by a

transformation, yields the following two equations.

$$\frac{1}{n-1} \mathbf{X'} \mathbf{Y} \mathbf{w}_{y} = 2\lambda_{x} \mathbf{w}_{x} \qquad (3)$$
$$\frac{1}{n-1} \mathbf{Y'} \mathbf{X} \mathbf{w}_{x} = 2\lambda_{y} (\mathbf{Y'} \mathbf{Y} + \kappa \mathbf{Y'} \mathbf{P'} \mathbf{D'} \mathbf{D} \mathbf{P} \mathbf{Y}) \mathbf{w}_{y} \qquad (4)$$

Eqs. (3) and (4) can be rewritten as the eigenvalue problems

$$\frac{1}{(n-1)^2} \mathbf{X'} \mathbf{Y} (\mathbf{Y'} \mathbf{Y} + \kappa \mathbf{Y'} \mathbf{P'} \mathbf{D'} \mathbf{DPY})^{-1} \mathbf{Y'} \mathbf{X} \mathbf{w}_x = \lambda \mathbf{w}_x \qquad (5)$$
$$\frac{1}{(n-1)^2} \mathbf{Y'} \mathbf{X} \mathbf{X'} \mathbf{Y} \mathbf{w}_y = \lambda (\mathbf{Y'} \mathbf{Y} + \kappa \mathbf{Y'} \mathbf{P'} \mathbf{D'} \mathbf{DPY}) \mathbf{w}_y. \qquad (6)$$

where $\lambda = 4\lambda_x\lambda_y$. These eigenvalue problems can be computed by using singular value decomposition.

Factor loading in PLS-ROG

We now describe the statistical properties of the eigenvector \mathbf{w}_x that can be used for factor loading in PLS-ROG. The correlation coefficient between the latent response variables **s** and the *p*-th explanatory variable \mathbf{x}_p can be written as

$$\operatorname{corr}\left(\mathbf{s}, \mathbf{x}_{p}\right) = \frac{\mathbf{s}' \mathbf{x}_{p} / n - 1}{\sqrt{\operatorname{var}(\mathbf{s})} \sqrt{\operatorname{var}(\mathbf{x}_{p})}}$$
(7)

Substituting $\mathbf{s} = \mathbf{Y}\mathbf{w}_y$ and $\mathbf{x}_p = \mathbf{X}\mathbf{c}$, where **c** the column vector in which the *p*-th element is 1 and the others are 0, yields

$$\operatorname{corr}(\mathbf{s}, \mathbf{x}_{p}) = \frac{\mathbf{w}_{y} \, \mathbf{Y} \, \mathbf{X} \, \mathbf{c} \, / \, n - 1}{\sqrt{\operatorname{var}(\mathbf{s})} \sqrt{\operatorname{var}(\mathbf{x}_{p})}} \,. \tag{8}$$

Transposing eq. (3) gives $\mathbf{w}_{y}'\mathbf{Y}'\mathbf{X}/n-1=2\lambda_{x}\mathbf{w}_{x}'$ which can be substituted in eq.

(8), giving

$$\operatorname{corr}(\mathbf{s}, \mathbf{x}_{p}) = \frac{2\lambda_{x} \mathbf{w}_{x}' \mathbf{c}}{\sqrt{\operatorname{var}(\mathbf{s})} \sqrt{\operatorname{var}(\mathbf{x}_{p})}}$$
(9)
Collection of Biostatistics
Research Archive

Now

$$\operatorname{var}(\mathbf{s}) = \frac{1}{n-1} \mathbf{w}_{y} \mathbf{Y} \mathbf{Y}_{y} \quad (10)$$

which can be substituted into eq. (9), giving

$$\operatorname{corr}(\mathbf{s}, \mathbf{x}_{\mathbf{p}}) = \frac{2\lambda_{x} \mathbf{w}_{x}' \mathbf{c}}{\sqrt{\operatorname{var}(\mathbf{s})} \sqrt{\operatorname{var}(\mathbf{x}_{p})}} = \frac{2\lambda_{x} w_{x,p}}{\sqrt{\mathbf{w}_{y}' \mathbf{Y}' \mathbf{Y} \mathbf{w}_{y}} / (n-1)} \sigma_{p}$$

$$= \frac{\sqrt{(n-1)\lambda} w_{x,p}}{\sqrt{\mathbf{w}_{y}' \mathbf{Y}' \mathbf{Y} \mathbf{w}_{y}} \sigma_{p}}$$
(11)

The scalar $\mathbf{w_y'Y'Yw_y}$ can be assumed to be constant because it does not depend on the *p*-th variable $\mathbf{x_p}$. With autoscaling of data, the *p*-th component of the eigenvector \mathbf{w}_x is proportional to the correlation coefficient between the latent response variable **s** and the *p*-th explanatory variable $\mathbf{x_p}$. For this reason, factor loadings in PLS can be defined by correlation coefficient in eq. (11). Using this definition we can perform a statistical test for factor loading.

Reference

- M. Barker, W. Rayens, Partial least squares for discrimination, J. Chemom. 17 (2001) 166-173.
- [2] H.Yamamoto, H.Yamaji, Y.Abe, K.Harada, D.Waluyo, E.Fukusaki, A.Kondo, H. Ohno, H.Fukuda, Dimensionality reduction for metabolome data using PCA, PLS, OPLS, and RFDA with differential penalties to latent variables, Chemom. Intell. Lab. Syst., 98 (2009) 136-142

