



Available online at www.sciencedirect.com

ScienceDirect



Procedia Computer Science 60 (2015) 850 - 859

19th International Conference on Knowledge Based and Intelligent Information and Engineering Systems

A variable selection method considering cluster loading for labeled high dimension low sample size data

Jiaxin Chen^a, Mika Sato-Ilic^{b,*}

^aGraduate School of Systems and Information Engineering, University of Tsukuba, Tsukuba, Ibaraki 305-8573, Japan
^bFaculty of Engineering, Information and Systems, University of Tsukuba, Tsukuba, Ibaraki 305-8573, Japan

Abstract

As the information society rapidly develops, there is an increased importance placed of dealing with high dimension low sample size (HDLSS) data, whose number of variables is much larger than the number of objects. Moreover, the selection of effective variables for HDLSS data is becoming more crucial. In this paper, a variable selection method considering cluster loading for labeled HDLSS data is proposed. Related to cluster loading, the conventional model considering principal component analysis has been proposed. However, the model can not be used for HDLSS data. Therefore, we propose a cluster loading model using a clustering result. By using the obtained cluster loading, we can select variables which belong to clusters unrelated with the given discrimination information represented by the labels of objects. Several numerical examples show a better performance of the proposed method.

© 2015 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/).

Peer-review under responsibility of KES International

Keywords: HDLSS data; variable selection; cluster loading; clustering

1. Introduction

With the rapid development of the information society, data analysis of HDLSS data ^{1,2,3,4} is getting more and more valued, especially in the specific application areas such as microarray gene data and image data. In HDLSS data, the number of variables is much larger than the number of objects, so how to select the effective variables plays an important role in analyzing HDLSS data. In this paper, a variable selection method considering cluster loading for labeled HDLSS data is proposed.

The cluster loading is defined as the relationship between given labels of objects and obtained clusters of variables. Related to the cluster loading, a method for constrained principal component analysis (CPCA)⁵ has been proposed, which is based on principal component analysis with external information on both objects and variables. In this model, the relationship between objects and variables is estimated by using least squares method. As a mathematically similar model to the CPCA, a model for interpreting principle components by using discrimination information⁶ has

^{*} Corresponding author. Tel.: +81-29-853-5006; fax: +81-29-853-5006. E-mail address: mika@risk.tsukuba.ac.jp

been proposed. In this model, the discrimination information is represented as the labels of objects and relationship between given labels and obtained principal components is estimated. The cluster loading is defined by using the framework of this model mathematically. However, the previously proposed model for the interpretation of principal components can not be used for HDLSS data. Because the result from principal component analysis is associated with the eigenvalues which are calculated through variance-covariance matrix. In HDLSS data, these obtained eigenvalues of samples are not consistent with the eigenvalues of population. As a result, this model is not applicable to HDLSS data

By using the obtained cluster loading, we can select variables which belong to clusters unrelated with the given labels. This means that we can select variables which do not relate to the discrimination of objects. The selected variables are the targets to be deleted variables.

The construction of this paper is as follows: In Section 2, a method for constrained principal component analysis (CPCA) is referred. In Section 3, a method for interpreting principal components using discrimination information is described. In Section 4, a cluster loading model for HDLSS data is proposed, and section 5 proposes a variable selection method using the cluster loading. Section 6 shows numerical examples of the proposed method. In Section 7, conclusions are stated.

2. Constrained principal component analysis

Let **X** be a data matrix consisted of *n* objects and *p* variables, where n > p. Suppose **G** a $n \times Q$ object information matrix and **H** a $p \times K$ variable information matrix. A model of constrained principal component analysis⁵ has been defined as follows:

$$\mathbf{X} = \mathbf{G}\mathbf{M}\mathbf{H}^{\mathrm{T}} + \mathbf{B}\mathbf{H}^{\mathrm{T}} + \mathbf{G}\mathbf{C} + \mathbf{E},\tag{1}$$

where M, B and C are the matrices of coefficients to be estimated, and E is a matrix of errors. In equation (1), X is defined by three terms, where the first term can be explained by both G and H, the second term is only defined by H, the third term is defined by G. That is, M shows the relationship between summarized objects and summarized variables, B shows relationship of objects through summarized variables and C shows relationship of variables through summarized objects.

Estimates of **M**, **B** and **C** are obtained by minimizing $SS(\mathbf{E}) = tr(\mathbf{E^T E})$, where $SS(\mathbf{E})$ shows the sum of square of errors. From equation (1), we consider the first term of the model and define a model as follows:

$$\mathbf{X} = \mathbf{G}\mathbf{M}\mathbf{H}^{\mathsf{T}} + \mathbf{E}_{\mathbf{1}}.\tag{2}$$

From equation (2),

$$SS(\mathbf{E}_{1}) = tr(\mathbf{E}_{1}^{\mathsf{T}}\mathbf{E}_{1})$$

$$= tr((\mathbf{X} - \mathbf{G}\mathbf{M}\mathbf{H}^{\mathsf{T}})^{\mathsf{T}}(\mathbf{X} - \mathbf{G}\mathbf{M}\mathbf{H}^{\mathsf{T}}))$$

$$= tr(\mathbf{X}^{\mathsf{T}}\mathbf{X}) - tr(\mathbf{X}^{\mathsf{T}}\mathbf{G}\mathbf{M}\mathbf{H}^{\mathsf{T}}) - tr(\mathbf{H}\mathbf{M}^{\mathsf{T}}\mathbf{G}^{\mathsf{T}}\mathbf{X}) + tr(\mathbf{H}\mathbf{M}^{\mathsf{T}}\mathbf{G}^{\mathsf{T}}\mathbf{G}\mathbf{M}\mathbf{H}^{\mathsf{T}}).$$
(3)

From equation (3) and

$$\begin{split} \frac{\partial tr(\mathbf{E_1^T E_1})}{\partial \mathbf{M}} &= \frac{\partial tr(\mathbf{X^T X})}{\partial \mathbf{M}} - \frac{\partial tr(\mathbf{X^T G M H^T})}{\partial \mathbf{M}} - \frac{\partial tr(\mathbf{HM^T G^T X})}{\partial \mathbf{M}} + \frac{\partial tr(\mathbf{HM^T G^T G M H^T})}{\partial \mathbf{M}} \\ &= -2(\mathbf{G^T X H}) + 2(\mathbf{G^T G M H^T H}) \\ &= 0. \end{split}$$

we obtain $\hat{\mathbf{M}}$ which is the estimate of \mathbf{M} as follows:

$$\hat{\mathbf{M}} = (\mathbf{G}^{\mathsf{T}}\mathbf{G})^{\mathsf{T}}\mathbf{X}\mathbf{H}(\mathbf{H}^{\mathsf{T}}\mathbf{H})^{\mathsf{T}},\tag{4}$$

where $(\mathbf{G}^{\mathsf{T}}\mathbf{G})^{\mathsf{-}}$ and $(\mathbf{H}^{\mathsf{T}}\mathbf{H})^{\mathsf{-}}$ are g-inverses of $\mathbf{G}^{\mathsf{T}}\mathbf{G}$ and $\mathbf{H}^{\mathsf{T}}\mathbf{H}$, respectively.

Download English Version:

https://daneshyari.com/en/article/489615

Download Persian Version:

https://daneshyari.com/article/489615

<u>Daneshyari.com</u>