# Gene subset selection in kernel-induced feature space

Satoshi Niijima [a,*], Satoru Kuhara [b]

[a] *Department of Bioinformatics, Graduate School of Systems Life Sciences, Kyushu University, Hakozaki 6-10-1, Higashi-ku, Fukuoka 812-8581, Japan*
[b] *Faculty of Agriculture, Kyushu University, Hakozaki 6-10-1, Higashi-ku, Fukuoka 812-8581, Japan*

## Abstract

This paper proposes a new filter approach to gene subset selection for kernel-based classifiers. We derive kernel forms of several well-known class separability criteria, and gene subset selection based on the kernelized criteria is applied to microarray cancer classification problems. The performance of our proposed strategy is compared in experiments with those of the conventional filter approach as well as gene ranking methods.
© 2006 Elsevier B.V. All rights reserved.

## 1. Introduction

Microarray technology allows us to measure the expression levels of thousands of genes simultaneously, producing a vast amount of data. Microarrays pose a great challenge on the data analysis, because the number of genes often exceeds tens of thousands, whereas the number of samples available is at most a few hundred.

In microarray data analysis, *gene selection* has been a central issue in recent years (see Stolovitzky (2003) and Cuperlovic-Culf et al. (2005) for reviews). Gene selection is often used to identify genes most relevant to a specific classification task, for example, those differentiate between normal and cancerous tissue samples. Gene selection plays several important roles in classification tasks (Krishnapuram et al., 2004; Hochreiter and Obermayer, 2004). It improves the prediction accuracy of classifiers by using only discriminative genes. It also saves computational costs by reducing dimensionality. More importantly, if it is possible to identify a small subset of biologically relevant genes, it

may provide insights into understanding the underlying mechanism of a certain biological phenomenon. Also, such information can be useful for designing less expensive experiments by targeting only a small number of genes.

The most common gene selection approach used in practice is so-called *gene ranking*. It is a univariate approach in the sense that each gene is evaluated individually with respect to a certain criterion that represents class discrimination ability, and is ranked according to the assigned score. The criteria frequently used are *t*-statistics, the signal-to-noise ratio (Golub et al., 1999), the Fisher's criterion (Dudoit et al., 2002), information gain and $\chi^2$ statistics (Liu et al., 2002) among others. Usually, top-ranked genes are selected for subsequent classification tasks. Although they are simple to implement, those genes that function in combination with other genes can be lowly ranked, if a single gene is not discriminative by itself. On the other hand, *gene subset selection* evaluates the discrimination ability of a subset of genes, hence a multivariate approach.

In general, gene selection can be performed in different manners: the *filter* and *wrapper* approaches. The filter approach employs some separability measure as an evaluation criterion, while the wrapper approach evaluates a subset of genes based on the performance of a specific

---

* Corresponding author. Tel.: +81 92 642 3043; fax: +81 92 642 3043.
  *E-mail addresses:* niijima@grt.kyushu-u.ac.jp (S. Niijima), kuhara@grt.kyushu-u.ac.jp (S. Kuhara).

classifier. Although the wrapper approach often gives satisfactory classification accuracy (Inza et al., 2004; Xiong et al., 2001), evaluating the classification performance for each subset of genes demands expensive computational costs. Moreover, the selected genes may overfit the classifier used, hence there is no guarantee that they also give high performance for other classifiers.

Among many classification methods previously proposed, recent studies have demonstrated the effectiveness of kernel methods such as support vector machines (SVMs) in various classification problems of bioinformatics (Schölkopf et al., 2004). In this paper, we propose a new filter approach to gene subset selection for kernel methods. Our major goal is to enhance the performance of kernel-based classifiers without resorting to the wrapper approach. This can be realized by directly performing gene subset selection in a kernel-induced feature space, and also performing classification in the same feature space. It is different from the conventional filter approach combined with a kernel-based classifier, where gene subset selection is performed in input space.

We show that several well-known class separability criteria can be kernelized. Gene subset selection is performed based on the kernelized criteria. The performance of our strategy is assessed by combining it with three kernel-based classifiers, from simple to advanced ones. We apply it to cancer classification problems on acute lymphoblastic leukemia (ALL) (Yeoh et al., 2002) and compare the performance of our proposed strategy with those of gene ranking methods and the conventional filter approach to gene subset selection.

## 2. Gene subset selection

In classification problems using microarray data, we are given $n$ samples consisting of gene expression levels of $p$ genes, i.e. $x_i \in \mathbb{R}^p (i = 1, \ldots, n)$, and the known class labels $y_i \in \{1, \ldots, C\}$ $(i = 1, \ldots, n)$, where $C$ is the number of classes. The problem is to predict the class label of a given test sample with unknown class.

### 2.1. Kernelization of class separability criteria

In order to select an appropriate subset of genes for classification, an evaluation criterion that measures class separability is needed. In our study, we employ several common measures based on the within-class and total scatter matrices (Fukunaga, 1990; Theodoridis and Koutroumbas, 1999).

For simplicity, we treat binary cases, i.e. $C = 2$. As mentioned later, however, the proposed strategy can be used also in multiclass cases. Let $X \in \mathbb{R}^{n \times p}$ be a sample matrix containing $x_i$ $(i = 1, \ldots, n)$ as rows. Then, the total scatter matrix $S_T$ can be expressed as

$$S_T = \frac{1}{n} X^T X - X^T \vec{1}_n \vec{1}_n^T X, \tag{1}$$

where $(\cdot)^T$ is the transpose operator, and $\vec{1}_n$ denotes an $n$-dimensional vector with all the components equal to $1/n$.

We shall further denote a sample matrix for class $k$ ($k = 1, 2$) by $X_k \in \mathbb{R}^{n_k \times p}$, where $n_k$ is the number of samples for class $k$. Then, the pooled within-class scatter matrix $S_W$ can be expressed as

$$S_W = \frac{n_1}{n} S_1 + \frac{n_2}{n} S_2, \tag{2}$$

where

$$S_k = \frac{1}{n_k} X_k^T X_k - X_k^T \vec{1}_{n_k} \vec{1}_{n_k}^T X_k.$$

For obtaining good class separability, larger between-class scatter and smaller within-class scatter are preferable. Such criteria used in our study are

$$J_1 = \frac{\text{Tr}(S_T)}{\text{Tr}(S_W)}, \tag{3}$$

$$J_2 = \text{Tr}(S_W^{-1} S_T), \tag{4}$$

and the Mahalanobis distance defined as

$$J_3 = (\mu_1 - \mu_2)^T S_W^{-1} (\mu_1 - \mu_2), \tag{5}$$

where $\text{Tr}(\cdot)$ and $(\cdot)^{-1}$ denote the trace and inverse of a matrix, respectively. $\mu_k$ is a sample mean vector for class $k$, which can be written as

$$\mu_k = X_k^T \vec{1}_{n_k}.$$

We may find gene subsets that maximize these criteria.

In terms of kernel methods, selecting a subset of genes based on the criteria corresponds to performing gene subset selection in *input space*. It turns out that these criteria can be rewritten by using the symmetric matrix

$$K = \{x_i \cdot x_j\}_{i,j} = XX^T, \tag{6}$$

which represents the inner products of samples.

Following Ruiz and López-de-Teruel (2001), let us define:

$$Z = \left( \frac{1}{n} I_n - \vec{1}_n \vec{1}_n^T \right)^{1/2}, \tag{7}$$

$$Z_k = \left( \frac{1}{n_k} I_{n_k} - \vec{1}_{n_k} \vec{1}_{n_k}^T \right)^{1/2}, \tag{8}$$

where $I_n \in \mathbb{R}^{n \times n}$ and $I_{n_k} \in \mathbb{R}^{n_k \times n_k}$ denote the identity matrices. Then, we have

$$J_1 = \frac{n \text{Tr}(Z^2 K)}{n_1 \text{Tr}(Z_1^2 K_{11}) + n_2 \text{Tr}(Z_2^2 K_{22})}, \tag{9}$$

where

$$K_{kk} = X_k X_k^T.$$

From the Sherman–Morrison–Woodbury formula (Meyer, 2000) and Corollary 2 in (Ruiz and López-de-Teruel, 2001), we also have

$$J_2 = \text{Tr}(K_1 V (I_{n_1} - K_{12}((W_2^2)^{-1} + K_{12}^T V K_{12})^{-1} K_{12}^T V) K_1^T Z^2), \tag{10}$$