

Proceedings

Open Access

A kernel-based approach for detecting outliers of high-dimensional biological data

Jung Hun Oh and Jean Gao*

Address: Department of Computer Science and Engineering, The University of Texas, Arlington, Texas, USA

Email: Jung Hun Oh - jung.oh@uta.edu; Jean Gao* - gao@uta.edu

* Corresponding author

from IEEE International Conference on Bioinformatics and Biomedicine (BIBM) 2008
Philadelphia, PA, USA. 3–5 November 2008

Published: 29 April 2009

BMC Bioinformatics 2009, 10(Suppl 4):S7 doi:10.1186/1471-2105-10-S4-S7

This article is available from: <http://www.biomedcentral.com/1471-2105/10/S4/S7>

© 2009 Oh and Gao; licensee BioMed Central Ltd.

This is an open access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

Background: In many cases biomedical data sets contain outliers that make it difficult to achieve reliable knowledge discovery. Data analysis without removing outliers could lead to wrong results and provide misleading information.

Results: We propose a new outlier detection method based on Kullback-Leibler (KL) divergence. The original concept of KL divergence was designed as a measure of distance between two distributions. Stemming from that, we extend it to biological sample outlier detection by forming sample sets composed of nearest neighbors. KL divergence is defined between two sample sets with and without the test sample. To handle the non-linearity of sample distribution, original data is mapped into a higher feature space. We address the singularity problem due to small sample size during KL divergence calculation. Kernel functions are applied to avoid direct use of mapping functions. The performance of the proposed method is demonstrated on a synthetic data set, two public microarray data sets, and a mass spectrometry data set for liver cancer study. Comparative studies with Mahalanobis distance based method and one-class support vector machine (SVM) are reported showing that the proposed method performs better in finding outliers.

Conclusion: Our idea was derived from Markov blanket algorithm that is a feature selection method based on KL divergence. That is, while Markov blanket algorithm removes redundant and irrelevant features, our proposed method detects outliers. Compared to other algorithms, our proposed method shows better or comparable performance for small sample and high-dimensional biological data. This indicates that the proposed method can be used to detect outliers in biological data sets.

Background

Outlier detection is an active research area that has many applications such as network intrusion detection [1], fraud detection [2] and biomedical data analysis [3]. In particular, outliers caused from instrument error or

human error in the biomedical data analysis such as biomarker selection and disease diagnosis could deeply degrade the performance of the data analysis. Therefore, prior to the analysis, during preprocessing it is imperative to remove outliers to prevent wrong results. To detect such

anomalous observations from normal ones, data mining techniques are widely used.

Outlier detection has been studied by researchers using a diversity of approaches. Statistical methods often view objects that are located relatively far from the center of the data distribution as outliers. Several distance measures were implemented. The Mahalanobis distance is the most commonly used multivariate outlier criterion. Based on Akaike's Information Criterion (AIC), Kadota *et al.* developed a method for detecting outliers, which is free from a significance level [4]. Knorr and Ng introduced a distance-based approach in which outliers are those objects for which there are less than k points within a given threshold in the input data set [5,6]. Angiulli *et al.* proposed a distance-based outlier detection method which finds the top outliers and provides a subset of the data set, called outlier detection solving set, that can be used to predict if new unseen objects are outliers [7]. Distance-based strategies are advantageous since model learning is not required. As an alternative, clustering algorithms can be used for outlier detection in which objects that do not belong to any cluster are regarded as outliers. Wang and Chiang proposed an effective cluster validity measure with outlier detection and cluster merging strategies for support vector clustering (SVC) [8]. The validity measure is capable of finding suitable values for the kernel parameter and soft margin constant. Based on these parameters, SVC algorithm can identify the ideal cluster number and increase robustness to outliers and noises. Schölkopf proposed a method of adapting support vector machine (SVM) to one-class classification problems [9]. Manevitz and Yousef presented two versions using the one-class SVM, both of which can identify outliers: Schölkopf's method and their proposed suggestion [10]. In such methods, after mapping the original samples into a feature space using an appropriate kernel function, the origin is referred to as the second class. In the feature space, samples close to the origin or lying on the standard subspace such as axes are regarded as outliers. Bandyopadhyay and Santra applied a genetic algorithm to the outlier detection problem in a lower dimensional space of a given data set, dividing these spaces into grids and efficiently computing the sparsity factor of the grid [11]. Aggarwal and Yu studied the problem of outlier detection for high-dimensional data, which works by finding lower dimensional projections [12]. Malossini *et al.* proposed two methods for detecting potential labeling errors: Classification-stability algorithm (CL-stability) and Leave-One-Out-Error-sensitivity algorithm (LOOE-sensitivity) [13]. In CL-stability, the stability of the classification of a sample is evaluated with a small perturbation of the other samples. LOOE-sensitivity was derived from the fact that if a sample is mislabeled, flipping the label of the sample should improve the prediction power.

In this paper, we propose a new outlier detection method based on KL divergence [14]. Due to the possible non-linearity of data structure, we deal with this problem in a higher feature space rather than the original space. Several issues arise after data mapping such as singularity because of small sample size versus high feature dimension. We address the computational issues and show the effectiveness of the proposed approach, KL divergence for outlier detection (KLOD).

Methods

Markov blanket

Markov blanket algorithm proposed by Koller and Sahami is a cross-entropy based technique to identify redundant and irrelevant features [15]. Let F be a full set of features and $M \subseteq F$ be a subset of features which does not contain feature F_i . Then, M is called a Markov blanket for F_i if F_i is conditionally independent of $F - M - \{F_i\}$ given M . Generally, the Markov blanket M_i of F_i is defined as a subset of features that consists of some features that have the highest Pearson correlation with F_i . To evaluate the closeness between F_i and its Markov blanket M_i , the following expected cross-entropy Δ is estimated:

$$\Delta(F_i | M_i) = \sum_{f_{M_i}, f_i} P(M_i = f_{M_i}, F_i = f_i) \times D(P(c | M_i = f_{M_i}, F_i = f_i) || P(c | M_i = f_{M_i})),$$

where f_{M_i} and f_i are feature values to M_i and F_i , respectively, c is the class label, and $D(.||.)$ represents the cross-entropy (a.k.a. Kullback-Leibler divergence). For each feature, Δ value is computed and a feature with the smallest Δ value is eliminated from the whole feature set. With the remaining features, the procedure is repeated until a predefined number of features remains.

Kullback-Leibler (KL) divergence

KL divergence, widely used in information theory, is adopted in Markov blanket as a core component. As shown in Markov blanket, KL divergence represents a measure of the distance between two probability distributions [16], i.e., for two probability densities $p(x)$ and $q(x)$, the KL-divergence is defined as

$$D_{KL}(p || q) = \int_x p(x) \log \frac{p(x)}{q(x)} dx.$$

Suppose that $\mathcal{N}(\mu, \Sigma)$ is a multivariate Gaussian distribution defined as

$$\mathcal{N}(\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right),$$

where $x \in \mathcal{R}^m$ and $|\Sigma|$ is the determinant of covariance matrix Σ . Given two different probability density

functions, $p(\mathbf{x}) = \mathcal{N}_1(\mu_1, \Sigma_1)$ and $q(\mathbf{x}) = \mathcal{N}_2(\mu_2, \Sigma_2)$, the KL divergence is defined as

$$D_{\text{KL}}(\mathcal{N}_1||\mathcal{N}_2) = \frac{1}{2} \{ (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) + \log \frac{|\Sigma_2|}{|\Sigma_1|} + \text{tr}[\Sigma_1 \Sigma_2^{-1} - \mathbf{I}_m] \}.$$

Concept of KL divergence for outlier detection (KLOD)

In Markov blanket, based on KL divergence, after calculating Δ value of Eq. (1) for each feature, a feature with the lowest Δ value is considered to be the most redundant. Using KL divergence, our new outlier detection method, called KLOD, employs similar strategy to the Markov blanket, i.e., while Markov blanket algorithm detects redundant and irrelevant features, our method identifies outliers. In KLOD, each sample \mathbf{x}_i has a sample set that consists of t samples close to the \mathbf{x}_i . To calculate the distance between samples, Euclidean metric is used. More specifically, we define two sample sets, i.e., S_1 and S_2 : S_2 is a sample set close to \mathbf{x}_i in Euclidean distance and the other set S_1 consists of \mathbf{x}_i and all samples in S_2 . The similarity, $D_{\mathbf{x}_i}(S_1||S_2)$, between S_1 and S_2 for each sample can be measured by using KL divergence, where $1 \leq i \leq n$ and n is the total number of samples in the data set. Intuitively, in our strategy, a sample \mathbf{x}_i with the largest D is regarded as an outlier.

$$o = \arg \max_{1 \leq i \leq n} D_{\mathbf{x}_i}$$

Given a data set with nonlinear data structure, if we model the linearity for the data set, it will cause our strategy to fail. Here, we focus on modeling the nonlinearity. Accordingly, with a mapping function ϕ , the original space is mapped into a higher dimensional feature space. Let S_1^ϕ and S_2^ϕ denote the two sample sets in the feature space in which we compute the similarity $D(S_1^\phi || S_2^\phi)$ between S_1^ϕ and S_2^ϕ . For each sample, its $D(S_1^\phi || S_2^\phi)$ is calculated. A sample which has the largest $D(S_1^\phi || S_2^\phi)$ is referred to as an outlier.

Please see an example in Figure 1. However, the calculation leads to several important issues to be considered, such as kernel trick, singularity problem, and calculation of KL divergence in the feature space. In the following sections, we will describe them.

Kernel function

Suppose that $\{\mathbf{x}_1, \mathbf{x}_2, \cup \mathbf{x}_n\}$ are the given samples in the original space. After mapping the samples into a higher

feature space by a nonlinear mapping function ϕ , the samples in the feature space are observed as $\Phi_{m \times n} = [\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \cup, \phi(\mathbf{x}_n)]$ where m is the number of features. Denote \mathbf{K} as follows:

$$\mathbf{K} = \Phi^T \Phi.$$

The calculation can be performed using kernel trick, i.e., the ij th element, $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$, of the \mathbf{K} matrix can be computed as a kernel function $k(\mathbf{x}_i, \mathbf{x}_j)$. In literatures, the polynomial kernel and the Gaussian kernel are the most widely used kernel functions. In this study, the Gaussian kernel function is used:

$$k(\mathbf{x}, \mathbf{y}) = \exp \left(- \frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2} \right),$$

where σ controls the kernel width. Similar to Eq. (6), we define \mathbf{K}_{ij} as follows:

$$\mathbf{K}_{ij} = \Phi_i^T \Phi_j,$$

where if $i \neq j$, Φ_i and Φ_j are different sample sets in the feature space; if $i = j$, \mathbf{K}_{ij} is equivalent to the definition of \mathbf{K} . Indeed, the feature space and the mapping function may not be explicitly known. However, once the kernel function is known, we can easily deal with the nonlinear mapping problem by replacing the mapping functions by the kernel functions.

KL divergence equation is composed of mean and covariance components. The mean and the covariance matrix in the feature space are estimated as

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) = \Phi \mathbf{s},$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\phi(\mathbf{x}_i) - \mu)(\phi(\mathbf{x}_i) - \mu)^T = \Phi \mathbf{J} \mathbf{J}^T \Phi^T,$$

where $\mathbf{s}_{n \times 1} = \frac{1}{n} \bar{\mathbf{1}}^T$, $\mathbf{J} = \frac{1}{\sqrt{n}} (\mathbf{I}_n - \mathbf{s} \bar{\mathbf{1}}^T)$ and $\bar{\mathbf{1}} = [1, 1, \cup, 1]$.

Then, an $m \times n$ matrix \mathbf{W} is denoted as

$$\mathbf{W} = \Phi \mathbf{J} = \frac{1}{\sqrt{n}} [(\phi(\mathbf{x}_1) - \mu), \dots, (\phi(\mathbf{x}_n) - \mu)].$$

Singularity problem

The covariance matrix in Eq. (10) is rank-deficient due to the small number of samples against the number of features. This problem, called singularity problem, makes it impossible to calculate the inverse of the covariance matrix. To overcome the problem, several methods have

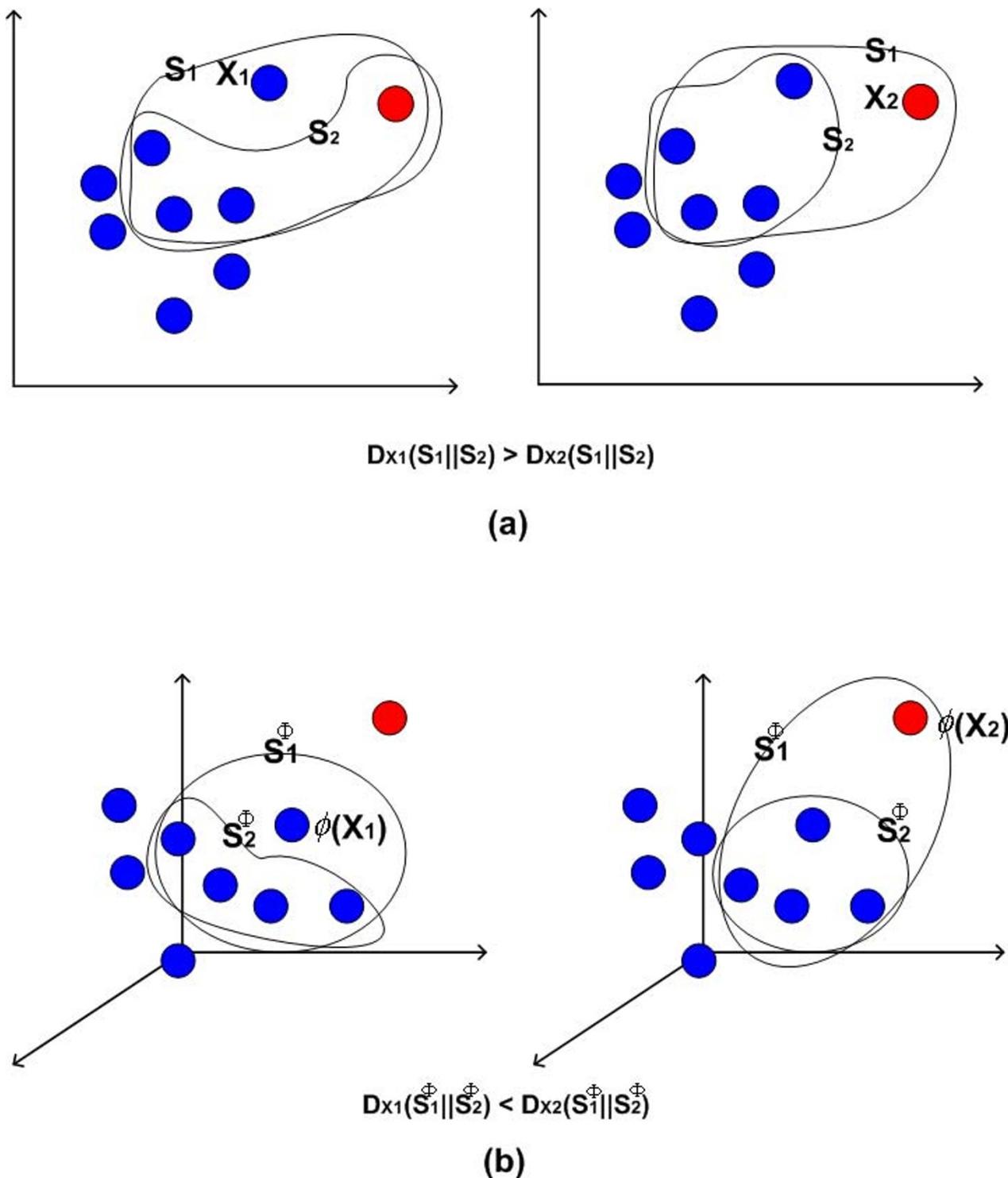


Figure 1
Outlier detection in a high feature space. Suppose that the red dot is a real outlier which is the farthest one from the majority of data. (a) in the original space, x_1 is regarded as an outlier. (b) in the higher feature space, x_2 is correctly detected as an outlier.

been proposed. In this study, we make use of a simple regularized approximation in which some positive constant values are added to the diagonal elements of the covariance matrix [17]. Therefore, the modified covariance matrix is of full rank, hence nonsingular. Let C denote

$$\begin{aligned} C &= \Phi J J^T \Phi^T + \rho I_m, \\ &= W W^T + \rho I_m, \\ &= \Phi R \Phi^T + \rho I_m, \end{aligned}$$

where $R = J J^T$, $\rho > 0$, and I_m is an identity matrix. In this study, $\rho = 1$ is used. Then, the inversion of the matrix C can be computed by using *Woodbury formula*:

$$\begin{aligned} C^{-1} &= (\rho I_m + \Phi J J^T \Phi^T)^{-1}, \\ &= (\rho I_m + W W^T)^{-1}, \\ &= \rho^{-1} (I_m - \rho^{-1} W (I_n + \rho^{-1} W^T W)^{-1} W^T), \\ &= \rho^{-1} (I_m - W (\rho I_n + W^T W)^{-1} W^T), \\ &= \rho^{-1} (I_m - \Phi J M^{-1} J^T \Phi^T), \\ &= \rho^{-1} (I_m - \Phi B \Phi^T), \end{aligned}$$

where $B = J M^{-1} J^T$ and $M = \rho I_n + W^T W = \rho I_n + J^T \Phi^T \Phi J = \rho I_n + J^T K J$.

Definition (Woodbury formula): Let A be a square $r \times r$ invertible matrix, where U and V are two $r \times k$ matrices with $k \leq r$. Assume that a $k \times k$ matrix $\Sigma = I_k + \beta V^T A^{-1} U$, in which I_k denotes a $k \times k$ identity matrix and β is an arbitrary scalar, is invertible. Then

$$(A + \beta U V^T)^{-1} = A^{-1} - \beta A^{-1} U \Sigma^{-1} V^T A^{-1}.$$

Calculation of KL divergence

Suppose that S_1^Φ and S_2^Φ are two sample sets in the feature space as mentioned in section. We know that the covariance matrices for both sets are singular. Let C_1 and C_2 denote the approximated covariance matrices for S_1^Φ and S_2^Φ , respectively, where the size of S_1^Φ is one larger than that of S_2^Φ . Also, let μ_1 and μ_2 be mean matrices for S_1^Φ and S_2^Φ , respectively. Therefore, KL divergence for S_1^Φ and S_2^Φ is expressed as follows:

$$2D_{KL}(\mathcal{N}_1 | \mathcal{N}_2) = (\mu_1 - \mu_2)^T C_2^{-1} (\mu_1 - \mu_2) + \log \frac{|C_2|}{|C_1|} + \text{tr}[C_1 C_2^{-1}] - m.$$

The KL divergence above is composed of three terms, i.e.,

$$\left\{ \begin{aligned} &(\mu_1 - \mu_2)^T C_2^{-1} (\mu_1 - \mu_2) \\ &\log \frac{|C_2|}{|C_1|} \\ &\text{tr}[C_1 C_2^{-1}]. \end{aligned} \right.$$

It should be noted that as shown in Eq. (9), Eq. (12) and Eq. (13), μ_i , C_i and C_i^{-1} ($i = 1$ or 2) have mapping functions rather than kernel functions.

Here, we will show how each term can be expressed by kernel functions instead of mapping functions. The first term consists of four sub-terms,

$$(\mu_1 - \mu_2)^T C_2^{-1} (\mu_1 - \mu_2) = \mu_1^T C_2^{-1} \mu_1 + \mu_2^T C_2^{-1} \mu_2 - \mu_1^T C_2^{-1} \mu_2 - \mu_2^T C_2^{-1} \mu_1.$$

Substituting Eq. (9) and Eq. (13) into each sub-term $\mu_i^T C_j^{-1} \mu_k$, we have

$$\begin{aligned} \mu_i^T C_j^{-1} \mu_k &= s_i^T \Phi_i^T \rho^{-1} (I_m - \Phi_j B_j \Phi_j^T) \Phi_k s_k, \\ &= \rho^{-1} (s_i^T K_{ik} s_k - s_i^T K_{ij} B_j K_{jk} s_k), \\ &= \rho^{-1} \theta_{ijk}. \end{aligned}$$

As a result of the effort, all mapping functions in the first term are replaced with kernel functions. Before dealing with the second term, we want to introduce the following three properties of determinant that are essential in the calculation of the second term.

Properties of determinant

- (a) If A is an r -by- r matrix, $\det|dA| = \det|dI_r A| = d^r \det|A|$.
- (b) If A and B are k -by- r matrices, $\det|I_k + AB^T| = \det|I_r + B^T A|$.
- (c) If A is invertible, $\det|A^{-1}| = 1/\det|A|$.

In the second term, we should compute the determinant of $C(C_1$ or $C_2)$. Instead of directly calculating the determinant of C , we try to obtain it through the determinant of C^{-1} . That is,

$$\begin{aligned} |C^{-1}| &= |\rho^{-1} (I_m - \Phi B \Phi^T)|, \\ &= \rho^{-m} |I_m - \Phi B \Phi^T|, \quad \text{by property (a)} \\ &= \rho^{-m} |I_m - Q \Phi^T|, \\ &= \rho^{-m} |I_n - \Phi^T Q|, \quad \text{by property (b)} \\ &= \rho^{-m} |I_n - \Phi^T \Phi B|, \\ &= \rho^{-m} |I_n - K B|, \end{aligned}$$

where $\mathbf{Q} = \Phi\mathbf{B}$. Here, by property (c), we can easily calculate $|\mathbf{C}|$, i.e.,

$$|\mathbf{C}| = \frac{1}{|\mathbf{C}^{-1}|} = \frac{\rho^m}{|\mathbf{I}_n - \mathbf{KB}|}.$$

By taking logarithm of $|\mathbf{C}|$, we have

$$\log |\mathbf{C}| = \log \frac{\rho^m}{|\mathbf{I}_n - \mathbf{KB}|} = m \log \rho - \log |\mathbf{I}_n - \mathbf{KB}|.$$

Note that the size of \mathbf{S}_1^Φ is one larger than that of \mathbf{S}_2^Φ . If the size of \mathbf{S}_2^Φ is k , the size of \mathbf{S}_1^Φ becomes $k + 1$.

Now we have the second term composed of kernel functions:

$$\begin{aligned} \log \frac{|\mathbf{C}_2|}{|\mathbf{C}_1|} &= \log |\mathbf{C}_2| - \log |\mathbf{C}_1|, \\ &= \log |\mathbf{I}_{k+1} - \mathbf{K}_{11}\mathbf{B}_1| - \log |\mathbf{I}_k - \mathbf{K}_{22}\mathbf{B}_2|. \end{aligned}$$

The third term can be replaced with kernel functions using properties of trace:

$$\begin{aligned} \text{tr}[\mathbf{C}_1\mathbf{C}_2^{-1}] &= \text{tr}[(\Phi_1\mathbf{R}_1\Phi_1^T + \rho\mathbf{I}_m)\rho^{-1}(\mathbf{I}_m - \Phi_2\mathbf{B}_2\Phi_2^T)], \\ &= \rho^{-1}\text{tr}[\Phi_1\mathbf{R}_1\Phi_1^T] - \rho^{-1}\text{tr}[\Phi_1\mathbf{R}_1\Phi_1^T\Phi_2\mathbf{B}_2\Phi_2^T] + m - \text{tr}[\Phi_2\mathbf{B}_2\Phi_2^T], \\ &= \rho^{-1}\text{tr}[\mathbf{R}_1\mathbf{K}_{11}] - \rho^{-1}\text{tr}[\mathbf{R}_1\mathbf{K}_{12}\mathbf{B}_2\mathbf{K}_{21}] + m - \text{tr}[\mathbf{B}_2\mathbf{K}_{22}]. \end{aligned}$$

Successfully, we substitute all mapping functions in the three terms of KL divergence by kernel functions so that we can calculate KL divergence between two sample sets in the feature space.

Results and discussion

To evaluate the performance of KLOD method, we performed several experiments using a synthetic data, two gene expression data sets, and a high-resolution mass spectrometry data. To obtain unbiased results, all experiments were repeated 30 times with 10-fold cross validation (CV) and the performance was averaged. The performance of KLOD was compared with one-class SVM and Mahalanobis distance based outlier detection methods. Given n samples, the Mahalanobis distance for each multivariate sample \mathbf{x}_i is as follows:

$$D_i = \sqrt{(\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})}$$

where $\boldsymbol{\Sigma}$ and $\boldsymbol{\mu}$ are the sample covariance matrix and sample mean vector, respectively. Samples with a large Mahalanobis distance are regarded as outliers.

Results on synthetic data

First, using a synthetic data, we evaluated KLOD to see the ability in detecting outliers. The synthetic data consists of

100 samples, denoted as \mathbf{N} , each of which has 100 features generated from a mixture of Gaussian $\mathcal{N}(0, \mathbf{I})$. In addition, two sample sets called quasi-outlier set \mathbf{Q} and perfect outlier set \mathbf{P} were produced, each of which has 10 samples with 100 features, which were generated from a mixture of Gaussian $\mathcal{N}(0, \mathbf{I})$ and $\mathcal{N}(2, \mathbf{I})$, respectively. It is noted that \mathbf{Q} was created from the same distribution as \mathbf{N} . Here, we corrupted \mathbf{Q} by changing the values in some features. To do so, some features from each sample in \mathbf{P} were randomly selected. The values of the selected features replaced those of features randomly selected from the corresponding sample in \mathbf{Q} . Finally, we merged \mathbf{N} and \mathbf{Q} , which were used as a synthetic data. Figure 2 illustrates an example of generating the synthetic data. In this experiment, we tested KLOD changing the number of corrupted features from 10 to 30 increasing by 2 and the size of a set, denoted as t , that consists of close samples of each sample from 5 to 20 increasing by 5. With the synthetic data, we measured how accurate our method is in identifying outliers in a way that the number of real outliers is counted out of the first 10 samples detected by KLOD.

Figure 3 shows the experimental results. When the number of noisy features increases, the accuracy shows a tendency to increase as well. It should be noted that for all set sizes, when the number of noisy features is 18, an accuracy of over 90% was obtained. Particularly, for $t = 10, 15$ and 20, when the number of noisy features is 30, an accuracy of 100% was achieved.

Performance evaluation after outlier removal

Before introducing the outlier removal for real biomedical data, we first introduce the performance evaluation metric we will use which is PCA (principal component analysis) + LDA (linear discriminant analysis). LDA maps the data into a very low dimensionality of $c - 1$, where c is the number of classes. In the reduced space, a simple matching procedure is used for classification. However, in order to guarantee a non-degenerate result from LDA, before the LDA task, the dimensionality of the data must be reduced to at most $n - c$ where n is the number of samples. Principal component analysis (PCA) is often used in the analysis of high dimensional data set. PCA performs a transformation of the original space into a lower dimensional space with little or no information loss while maximally preserving variance.

Lilien *et al.* used the PCA+LDA method in the analysis of mass spectrometry data sets [18]. In this framework, the PCA dimensionality-reduced samples are projected by LDA onto a hyperplane in the way of maximizing the between-class variance and minimizing the within-class

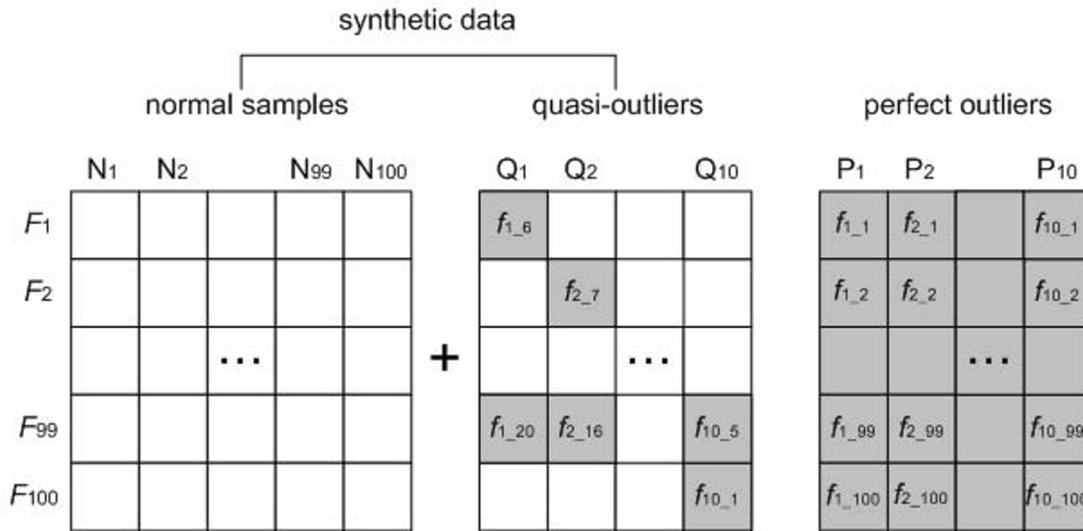


Figure 2
Generation of a synthetic data. This example shows a way used in this study to generate a synthetic data.

variance of the projected samples. To evaluate the performance after outlier removal in our experiments, we employed the PCA+LDA strategy.

Results on gene expression data sets

In this study, two public microarray data sets were used.

- The leukemia data set covers two types of acute leukemia: 47 acute lymphoblastic leukemia (ALL) samples and

25 acute myeloid leukemia (AML) samples with 7,129 genes. The data set is publicly available at <http://www.broad.mit.edu/cgi-bin/cancer/datasets.cgi/> [19].

- The colon data set contains 40 tumor and 22 normal colon tissues with 2,000 genes. The data set is available at <http://microarray.princeton.edu/oncology/> [20].

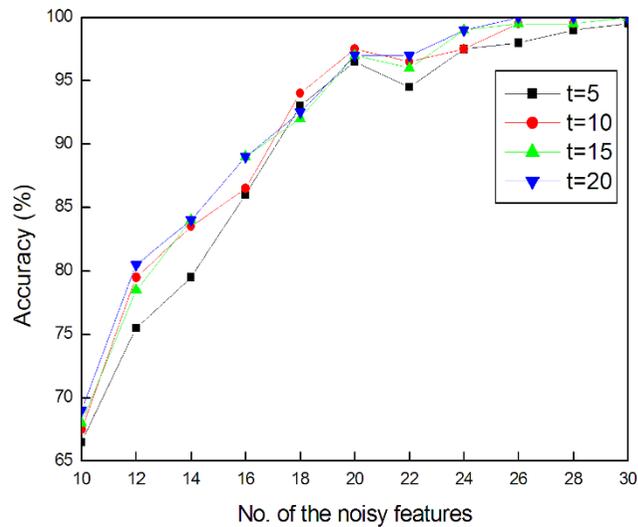


Figure 3
Accuracy of detecting outliers on a synthetic data. The data consists of 100 normal samples and 10 outliers, each having 100 features.

In experiments with the two microarray data sets, specificity, sensitivity, and accuracy were measured using PCA+LDA classification strategy after removing outliers detected by KLOD with $t = 10$, Mahalanobis distance based method, and one-class SVM. We define the specificity as the ratio of correctly classified negatives to the actual number of negatives. For leukemia and colon microarray data sets, negatives are ALL and normal samples, respectively. For KLOD and Mahalanobis distance based method, the performance was measured after removing a sample having the largest distance from each class at each iteration. If the prediction rate (specificity or sensitivity) decreases more than a threshold γ compared to the prediction rate before the outlier removal, we stop the outlier detection in the corresponding class. In this study, we used $\gamma = 0.5\%$. In contrast, for one-class SVM, after excluding all samples regarded as outliers in each class, the performance was assessed.

Table 1 shows the experimental results obtained using leukemia and colon microarray data sets. For the leukemia data set, KLOD achieved the best accuracy with 9 outliers (2 ALL and 7 AML samples).

Table 1: Performance after outlier detection in leukemia and colon data sets.

Data set	Measurements	Without outlier removal	After outlier removal		
			KLOD	Mahalanobis	One-class SVM
Leukemia	Specificity (%)	96.17	99.00	97.37	100
	Sensitivity (%)	95.60	99.44	100	95.24
	Accuracy (%)	95.97	99.13	98.28	98.33
	No. of the outliers	ALL	2	9	8
		AML	7	5	4
Colon	Specificity (%)	82.50	85.95	83.25	85.26
	Sensitivity (%)	88.25	94.43	85.90	94.17
	Accuracy (%)	86.21	91.25	85.00	91.09
	No. of the outliers	normal	1	2	3
		tumor	5	1	4

Mahalanobis distance based method and one-class SVM found 14 and 12 outliers, respectively. For the colon data set, KLOD found 6 outliers (1 normal and 5 tumor samples) with 84.95% specificity, 94.43% sensitivity, and 91.25% accuracy. It should be noted that the performance of Mahalanobis distance based method was degraded in terms of sensitivity and accuracy compared to the performance obtained using all samples without outlier removal, suggesting that outliers detected by Mahalanobis distance based method are unlikely to be real ones.

Results on mass spectrometry data

To evaluate the effectiveness of KLOD, we also used a public mass spectrometry data for liver cancer study that consists of 201 spectra containing hepatocellular carcinoma (HCC) (78), cirrhosis (51), and health (72) [3]. From <http://microarray.georgetown.edu/ressomlab/>, we downloaded the binned spectra that have 23,846 peaks for each spectrum. To test outlier detection methods, only cirrhosis and HCC spectra were used as in [3]. By using t-test with the significance level of 0.05 in cirrhosis and HCC spectra, we selected 10,682 peaks. That is, the top 10,682 peaks selected by t-test with cirrhosis and HCC spectra were used in outlier detection methods. The same way as performed

with the microarray data sets was employed. Here cirrhosis samples are regarded as negatives. As shown in Table 2, KLOD obtained slightly higher performance with the smallest number of outliers than Mahalanobis distance based method and one-class SVM. From the results in experiments using mass spectrometry and microarray data sets, it seems that one-class SVM detects more outliers than KLOD and Mahalanobis distance based method.

Conclusion

We proposed a new outlier detection method based on KL divergence called KLOD. Our idea was derived from Markov blanket algorithm where redundant and irrelevant features are removed based on KL divergence. We tackled the outlier detection problem in a higher feature space after mapping the original data. The mapping leads to several issues. In particular, we showed how to calculate KL divergence in the higher feature space by using the properties of determinant and trace of matrix. To assess the usefulness of KLOD, we used a synthetic data and real life data sets. Compared to Mahalanobis distance based method and one-class SVM, KLOD achieved higher or comparable performance.

Table 2: Performance after outlier detection in liver cancer mass spectrometry data.

Measurements	Without outlier removal	After outlier removal		
		KLOD	Mahalanobis	One-class SVM
Specificity (%)	93.63	94.69	94.29	94.35
Sensitivity (%)	92.82	93.95	93.51	93.89
Accuracy (%)	93.14	94.23	93.82	94.07
No. of the outliers	Cirrhosis	3	2	5
	HCC	2	4	6

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

JHO performed data analysis and wrote the manuscript. JG supervised the project and edited the paper.

Acknowledgements

This work was supported in part by NSF under grants IIS-0612152 and IIS-0612214.

This article has been published as part of *BMC Bioinformatics* Volume 10 Supplement 4, 2009: Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine (BIBM) 2008. The full contents of the supplement are available online at <http://www.biomedcentral.com/1471-2105/10?issue=S4>.

References

- Lee W, Stolfo S, Mok K: **Mining audit data to build intrusion detection models.** *Proc Int Conf Knowledge Discovery and Data Mining (KDD 1998)* 1998:66-72.
- Fawcett T, Provost F: **Adaptive fraud detection.** *Data Mining and Knowledge Discovery* 1997, **1**:291-316.
- Ressom H, Varghese R, Drake S, Hortin G, Abdel-Hamid M, et al.: **Peak selection from MALDI-TOF mass spectra using ant colony optimization.** *Bioinformatics* 2007, **23**:619-626.
- Kadota K, Tominaga D, Akiyama Y, Takahashi K: **Detecting outlying samples in microarray data: A critical assessment of the effect of outliers on sample classification.** *Chem-Bio Informatics Journal* 2003, **3**:30-45.
- Knorr E, Ng R: **Algorithms for mining distance-based outliers in large datasets.** *Proc Int Conf Very Large Databases (VLDB 1998)* 1998:392-403.
- Knorr E, Ng R, Tucakov V: **Distance-based outlier: algorithms and applications.** *Proc Int Conf Very Large Databases (VLDB 2000)* 2000:237-253.
- Angiulli F, Basta S, Pizzuti C: **Distance-based detection and prediction of outliers.** *IEEE Trans on Knowledge and Data Engineering* 2006, **18**:145-160.
- Wang JS, Chiang JC: **A cluster validity measure with outlier detection for support vector clustering.** *IEEE Trans on Systems, Man, and Cybernetics, Part B* 2008, **38**:78-89.
- Schölkopf B, Platt J, Shawe-Taylor J, Smola A, Williamson R: **Estimating the support of a high-dimensional distribution.** *Neural Computation* 2001, **13**:1443-1471.
- Manevitz L, Yousef M: **One-class SVMs for document classification.** *Journal of Machine Learning Research* 2001, **2**:139-154.
- Bandyopadhyay S, Santra S: **A genetic approach for efficient outlier detection in projected space.** *Pattern Recognition* 2008, **41**:1338-1349.
- Aggarwal C, Yu P: **Outlier detection for high dimensional data.** *Proc ACM SIGMOD* 2001:37-46.
- Malossini A, Blanzieri E, Ng R: **Detecting potential labeling errors in microarrays by data perturbation.** *Bioinformatics* 2006, **22**:2114-2121.
- Oh J, Gao J, Rosenblatt K: **Biological data outlier detection based on Kullback-Leibler divergence.** *Proc IEEE Int Conf on Bioinformatics and Biomedicine (BIBM 2008)* 2008:249-254.
- Koller D, Sahami M: **Toward optimal feature selection.** *Proc Int Conf on Machine Learnin* 1996.
- Tumminello M, Lillo F, Mantegna R: **Kullback-Leibler distance as a measure of the information filtered from multivariate data.** *Physical Review E* 2007, **76**:256-67.
- Zhou S, Chellappa R: **From sample similarity to ensemble similarity: probabilistic distance measures in reproducing kernel Hilbert space.** *IEEE Trans on Pattern Analysis and Machine Intelligence* 2006, **28**:917-929.
- Lilien R, Farid H, Donald B: **Probabilistic disease classification of expression-dependent proteomic data from mass spectrometry of human serum.** *Journal of Computational Biology* 2003, **10**:925-946.
- Golub T, Slonim D, Tamayo P, Huard C, Gaasenbeek M, et al.: **Molecular classification of cancer: class discovery and class prediction by gene expression monitoring.** *Science* 1999, **286**:531-537.
- Alon U, Barkai N, Notterman D, Gish K, Ybarra S, et al.: **Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays.** *Proc Natl Acad Sci U S A* 1999, **96**:6745-6750.

Publish with **BioMed Central** and every scientist can read your work free of charge

"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:
http://www.biomedcentral.com/info/publishing_adv.asp

