

研究成果の刊行に関する一覧表

書籍

著者氏名	書籍全体の編集者名	書籍名	出版社名	出版地	出版年
奥野恭史ほか	石渡信一・桂勲・桐野豊・美宅成樹	生物物理学ハンドブック	朝倉書店	日本	2007
奥野恭史ほか	日本バイオインフォマティクス学会	バイオインフォマティクス事典	共立出版	日本	2006

雑誌

発表者氏名	論文タイトル名	発表誌名	巻号	ページ	出版年
Nijima, S. and Okuno, Y.	Laplacian Linear Discriminant Analysis Approach to Unsupervised Feature Selection.	IEEE/ACM Transactions on Computational Biology and Bioinformatics	-	-	in press
Okuno, Y., Tamon, A., Yabuchi, H., Nijima, S., Minowa, Y., Tonomura, K., Kunimoto, R. and Feng, C.	GLIDA: GPCR-Ligand Database for Chemical Genomics Drug Discovery - Database and Tools Update.	Nucleic Acids Research	36	D907-12	2008
Kitajima, M., Minowa, Y., Matsuda, H. and Okuno, Y.	Compound-Transporter Interaction Studies using Canonical Correlation Analysis.	Chem-Bio Informatics J.	7	24-34	2007
Yamamoto, H., Takematsu, H., Fujinawa, R., Naito, Y., Okuno, Y., Tsujimoto, G., Suzuki, A. and Kozutsumi, Y.	Correlation index-based responsible-enzyme gene screening (CIRE S), a novel DNA microarray-based method for glycan biosynthesis enzyme gene.	PLoS ONE	2	e1232	2007
Ikeda, A., Miyazaki, T., Kakizawa, S., Okuno, Y., Tsuchiya, S., Myomoto, A., Saito, SY., Yamamoto, T., Yamazaki, T., Iino, M., Tsujimoto, G., Watanabe, M. and Takeshima, H.	Abnormal features in mutant cerebellar Purkinje cells lacking junctophilins.	Biochem. Biophys. Res. Commun.	363	835-9	2007
Yamazaki, T., Sasaki, N., Nishi, M., Yamazaki, D., Ikeda, A., Okuno, Y., Komazaki, S., and Takeshima, H.	Augmentation of drug-induced cell death by ER protein BRI3BP.	Biochem. Biophys. Res. Commun.	362	971-5	2007

Naito, Y., Takematsu, H., Koyama, S., Miyake, S., Yamamoto, H., Fujinawa, R., Sugai, M., <u>Okuno, Y.</u> , Tsujimoto, G., Yamaji, T., Hashimoto, Y., Ito, S., Kawasaki, T., Suzuki, A., and Kozutsumi, Y.	Germinal center marker GL7 probes activation-dependent repression of N-glycolylneuraminic acid, a sialic acid species involved in the negative modulation of B cell activation.	Mol. Cell Biol.	27(8)	3008-22	2007
Zhu, S., <u>Okuno, Y.</u> , Tsujimoto, G., and Mamitsuka, H.	Application of a new probabilistic model for mining implicit associated cancer genes from OMIM and Medline	Cancer Informatics	2	361-71	2006
Osada, S., Naganawa, A., Misou, M., Tsuchiya, S., Tamura, S., <u>Okuno, Y.</u> , Nishikawa, J., Satoh, K., Imagawa, I., Tsujimoto, G., Sugimoto, Y., and Nishihara, T.	Altered gene expression of transcriptional regulatory factors in tumor marker-positive cells during chemically induced hepatocarcinogenesis.	Toxicology Letters	167	106-113	2006
Tsuchiya, S., <u>Okuno, Y.</u> , and Tsujimoto, G.	MicroRNA: biogenetic and functional mechanisms and involvement in cell differentiation and cancer.	J Pharmacol Sci	101(4)	267-70	2006
<u>Okuno, Y.</u> , Yang, J., Taneishi, K., Yabuuchi, H., and Tsujimoto, G.	GLIDA: GPCR-Ligand database for Chemical Genomic Drug Discovery	Nucleic Acids Research	34	D673-7	2006
<u>Okuno, Y.</u>	Chemical-genomic network for integrating bioinformatics and cheminformatics	Journal of Pharmacological Sciences.	97	29p-29p Suppl. 1	2005
Zhu, S., <u>Okuno, Y.</u> , Tsujimoto, G., and Mamitsuka, H.	A probabilistic model for mining implicit 'chemical compound-gene' relations from literature	Bioinformatics	21	ii245-5 1	2005
Sugimoto, Y., Fukada, Y., Mori, D., Tanaka, S., Yamane, H., <u>Okuno, Y.</u> , Deai, K., Tsuchiya, S., Tsujimoto, G., and Ichikawa, A.	Prostaglandin E2 Stimulates Granulocyte Colony-Stimulating Factor Production via the Prostanoid EP2 Receptor in Mouse Peritoneal Neutrophils	J. Immunol.	175	2606-12	2005
Adachi, T., <u>Okuno, Y.</u> , Takenaka, S., Matsuda, K., Ohta, N., Takashima, K., Yamazaki, K., Nishimura, D., Miyatake, K., Mori, C., and Tsujimoto, G.	Comprehensive analysis of the effect of phytoestrogen, daidzein, on a testicular cell line, using mRNA and protein expression profile.	Food Chem. Toxicol.	43	529-35	2005

Yamada, M., Katsuma, S., Adachi, T., Hirasawa, A., Shiojima, S., Kadowaki, T., Okuno, Y., Koshimizu, TA., Fujii, S., Sekiya, Y., Miyamoto, Y., Tamura, M., Yumura, Y., Nihiei, H., Kobayashi, M., and Tsujimoto, G.	Inhibition of protein kinase CK2 prevents the progression of glomerulonephritis.	Proc. Natl. Acad. Sci. U. S. A.	102	736-41	2005
---	--	---------------------------------	-----	--------	------

その他

著者氏名	執筆タイトル名	掲載誌名	巻号	ページ	出版年
奥野 恭史ほか	ケミカル・バイオ情報に基づく創薬インフォマティクス研究	Pharma VISION NEWS	No. 9	13-16	2007
薬事日報	2008. 3. 21	掲載記事			
薬事日報	2007. 5. 1	掲載記事			
日刊工業新聞	2007. 3. 26	掲載記事			
日刊工業新聞	2007. 3. 20	掲載記事			
日経産業新聞	2007. 3. 20	掲載記事			
日刊工業新聞	2006. 3. 28	掲載記事			

研究成果の刊行物・別刷

Laplacian Linear Discriminant Analysis Approach to Unsupervised Feature Selection

Satoshi Nijjima and Yasushi Okuno

Abstract—Until recently, numerous feature selection techniques have been proposed and found wide applications in genomics and proteomics. For instance, feature/gene selection has proven to be a powerful tool for biomarker discovery from microarray and mass spectrometry data. While supervised feature selection has been explored extensively, there are only a few unsupervised methods that can be applied to exploratory data analysis in which class information is unavailable. In this paper, we address the problem of unsupervised feature selection. First, we extend Laplacian linear discriminant analysis (LLDA) to unsupervised cases. Second, we propose an efficient algorithm for computing LLDA. Finally, a new unsupervised feature selection algorithm, called LLDA-based Recursive Feature Elimination (LLDA-RFE), is presented. We apply LLDA-RFE to several public datasets of cancer microarrays and compare its performance with those of state-of-the-art unsupervised methods, Laplacian score and SVD-entropy, and of a supervised filter method, Fisher score. Our results demonstrate that LLDA-RFE outperforms Laplacian score and shows favorable performance against SVD-entropy. It performs even better than Fisher score for some of the datasets, despite the fact that LLDA-RFE is fully unsupervised.

Index Terms—Feature selection, linear discriminant analysis, graph Laplacian, microarray data analysis.

I. INTRODUCTION

IN recent years, feature/gene selection methods have been widely used in genomics and proteomics to handle a deluge of data produced by high-throughput technologies such as microarray and mass spectrometry. In microarray studies, for instance, a small fraction of genes typically exhibit significant differential expression among tens of thousands of genes whose expression levels are measured simultaneously. Thus, it is of great importance to identify genes relevant to a biological phenomenon of interest and to characterize their expression profiles. Gene selection can be useful for multiple purposes: to save computational costs of subsequent analysis by reducing the number of genes; to improve the prediction performance of classifiers by using discriminative genes only; and to identify informative genes for further investigation of their biological relevance. Specifically, gene selection has proven to be a powerful tool for biomarker discovery, i.e. searching for potential marker genes contributing to classification of cancer subtypes or prediction of clinical outcomes, which leads to more reliable diagnosis and better treatments of cancer.

To date, numerous techniques for feature selection have been developed [12] and also been applied successfully to the analysis of biological data with many features. In contrast to supervised feature selection, however, unsupervised feature selection has not yet been explored extensively. Indeed, there have been only a few

unsupervised methods proposed until recently [7], [14], [28], [30]. Unsupervised feature selection is of great use particularly for class discovery, where class information is unavailable. For instance, clustering is usually performed to find clusters in microarray samples on the basis of the expression profiles of all genes, but the clusters so obtained can be obscured by the large number of irrelevant genes. Therefore, unsupervised feature selection is essential to the exploratory analysis of biological data. Moreover, even when class labels are provided by external knowledge, but may be unreliable or mislabeled, overfitting can be alleviated by performing feature selection in an unsupervised manner. It is obviously more challenging to identify features that reveal underlying cluster structures in the samples than to find those exhibiting similar patterns across all the samples.

To address this problem, we propose a new unsupervised feature selection method, called Laplacian linear discriminant analysis-based Recursive Feature Elimination (LLDA-RFE). LLDA-RFE is closely related to Laplacian score [15], which is also based on graph Laplacian and can be applied in an unsupervised manner. The major difference is that, whereas Laplacian score is a univariate approach, LLDA-RFE is multivariate, allowing for selecting features that contribute to discrimination in combination with other features. Recently, Wolf and Shashua [30] proposed the $Q - \alpha$ algorithm, taking advantage of the spectral properties of the graph Laplacian of features. While the $Q - \alpha$ algorithm has an interesting property that the sparsity of features naturally emerges, it does not scale well to the feature size. Also, it involves iterative computations on a matrix of the feature size in a least-squares optimization process to ensure a local maximum solution. In contrast, our proposed algorithm for LLDA-RFE is computationally tractable and has a global maximum solution. It is shown that LLDA includes the maximum margin criterion (MMC) [18] as a supervised case. Although LLDA-RFE is a natural extension of MMC-RFE, the proposed algorithm need not reduce dimensionality before applying LLDA, unlike the MMC-RFE algorithm proposed previously [23].

We compare the performance of LLDA-RFE with those of state-of-the-art unsupervised feature selection methods, Laplacian score and SVD-entropy [28], on seven public datasets of cancer microarrays. The performances of these methods are evaluated by their capability of identifying discriminative genes without using class information. We also compare the performance between LLDA-RFE and a supervised filter method, Fisher score [8], [15]. Experimental results demonstrate that LLDA-RFE outperforms Laplacian score and shows favorable performance against SVD-entropy. Despite the fact that LLDA-RFE is fully unsupervised, it performs even better than Fisher score for some of the datasets.

The rest of this paper is organized as follows: In Section II, we give outlines of LDA and the MMC. We then introduce LLDA and extend it to unsupervised cases in Section III. An efficient algorithm for LLDA is also proposed. We present the LLDA-RFE

Manuscript received XXX XX, 2007; revised XXX XX, 2007.

S. Nijjima and Y. Okuno are with the Department of PharmacoInformatics, Frontier Education Center, Graduate School of Pharmaceutical Sciences, Kyoto University, 46-29 Yoshida Shimoadachi-cho, Sakyo-ku, Kyoto 606-8501, Japan.

algorithm for feature selection in Section IV. Section V describes related work on unsupervised feature selection. Experimental results on seven microarray datasets are presented and discussed in Section VI. Finally, we give concluding remarks in Section VII.

II. LDA AND MMC

Linear discriminant analysis (LDA) aims to find a set of projection vectors that maximize the between-class scatter and simultaneously minimize the within-class scatter, thereby achieving maximum discrimination [9].

Let $X \in \mathbb{R}^{p \times n}$ be a sample matrix containing $x_i, i = 1, \dots, n$ as columns, where n is the number of samples and p is the number of features. The between-class scatter matrix S_b and the within-class scatter matrix S_w are defined as:

$$S_b = \frac{1}{n} \sum_{k=1}^c n_k (m^{(k)} - m)(m^{(k)} - m)^T,$$

$$S_w = \frac{1}{n} \sum_{k=1}^c \sum_{j=1}^{n_k} (x_j^{(k)} - m^{(k)})(x_j^{(k)} - m^{(k)})^T,$$

where c is the number of classes; n_k is the number of samples in class k ; $x_j^{(k)}$ is the j th sample in class k ; $m^{(k)}$ and m are the mean vector of class k and the total mean vector, respectively. Then, classical LDA finds the projection matrix W by maximizing the Fisher criterion

$$J_{LDA}(W) = \text{trace}((W^T S_w W)^{-1} (W^T S_b W)). \quad (1)$$

By solving a generalized eigenvalue problem, W can be found as the eigenvectors of $S_w^{-1} S_b$ corresponding to the largest eigenvalues. However, when the dimensionality of samples is larger than the sample size, i.e. $p > n$, S_w becomes singular and we cannot compute $S_w^{-1} S_b$, which is a major drawback of classical LDA. This is known as the singularity problem or the small sample size problem.

To overcome this problem, Li *et al.* [18] recently proposed to use the maximum margin criterion (MMC) instead of (1) to find the projection vectors. The MMC is defined as:

$$J_{MMC}(W) = \text{trace}(W^T (S_b - S_w) W). \quad (2)$$

In this case, the projection matrix W that maximize the criterion (2) can be found as the eigenvectors of $S_b - S_w$ corresponding to the largest eigenvalues. Li *et al.* proposed an efficient algorithm to compute the projection matrix of the MMC under the constraint that $W^T S_t W = I$, where S_t is the total scatter matrix. This is found to be the same as the uncorrelated LDA (ULDA) algorithm in [32]. Also, an efficient algorithm for the MMC subject to the orthogonality constraint on W , i.e. $W^T W = I$, was presented in [23]. In both cases, we need not compute the inverse of S_w , hence the singularity problem can be easily avoided.

It should be noted that the MMC is not equivalent to the Fisher criterion. As Loog [21] disproved the theorem in [26], the discriminant vectors obtained by maximizing (2) are not generally the same as those obtained by maximizing (1). More precisely, although ULDA can be considered as an extension of classical LDA to small sample size cases [32], the MMC with the orthogonality constraint does not necessarily yield projection vectors that are optimal for discrimination. In practice, a better

discrimination can be achieved by balancing the between-class and within-class scatters using the following criterion as in [20]:

$$J_{MMC}(W) = \text{trace}(W^T (S_b - \mu S_w) W), \quad (3)$$

where μ is a non-negative constant. It is clear that (2) is a special case of (3). In the following sections, we focus on the MMC defined by (2) with the orthogonality constraint.

III. UNSUPERVISED LLDA

A. Extension of LLDA to unsupervised cases

We can write the total and within-class scatter matrices as follows:

$$S_t = \frac{1}{n} X (I - \frac{1}{n} e e^T) X^T$$

$$= \frac{1}{n} X (I - W_g) X^T,$$

$$S_w = \frac{1}{n} X (I - \sum_{k=1}^c \frac{1}{n_k} e^{(k)} e^{(k)T}) X^T$$

$$= \frac{1}{n} X (I - W_\ell) X^T,$$

where I is the identity matrix, $e = (1, 1, \dots, 1)^T$ is an n -dimensional vector, and $e^{(k)}$ is an n -dimensional vector with $e_i^{(k)} = 1$ if x_i belongs to class k , and 0 otherwise. In terms of graph Laplacians [6], $I - W_g$ can be viewed as the global Laplacian of a graph such that all vertices are connected with a constant weight of $1/n$, and $I - W_\ell$ as the local Laplacian of a graph such that vertices are connected with a constant weight of $1/n_k$ only when both belong to the k th class.

From the relationship [9]

$$S_t = S_b + S_w,$$

it follows that

$$S_b - S_w = S_t - 2S_w$$

$$= \frac{1}{n} X ((I - W_g) - 2(I - W_\ell)) X^T. \quad (4)$$

The MMC represented in this form is referred to as Laplacian linear discriminant analysis (LLDA) in [26] and was applied to extract discriminant features in supervised scenarios.

In this study, we extend (4) to unsupervised cases. We first define the global similarity matrix K_g and the local similarity matrix K_ℓ as:

$$[K_g]_{ij} = \begin{cases} k(x_i, x_j), & \text{if } i \neq j \\ 0, & \text{otherwise,} \end{cases}$$

$$[K_\ell]_{ij} = \begin{cases} k(x_i, x_j), & \text{if } x_i \text{ is among } k_\ell \text{ nearest neighbors of } x_j \\ & \text{or } x_j \text{ is among } k_\ell \text{ nearest neighbors of } x_i \\ 0, & \text{otherwise.} \end{cases}$$

Here, $k(\cdot, \cdot)$ represents the similarity between each pair of samples, and the standard measures include heat kernel (Gaussian kernel), inner product and Euclidean distance. In supervised cases, prior class information can also be reflected to guide the graph construction [31]. Let L_g and L_ℓ be the normalized global and local Laplacian matrices, respectively as:

$$L_g = I - D_g^{-\frac{1}{2}} K_g D_g^{-\frac{1}{2}},$$

$$L_\ell = I - D_\ell^{-\frac{1}{2}} K_\ell D_\ell^{-\frac{1}{2}},$$

where D_g and D_ℓ are diagonal matrices such that $[D_g]_{ii} = \sum_j [K_g]_{ji}$ and $[D_\ell]_{ii} = \sum_j [K_\ell]_{ji}$. Then, we seek to find a set of projection vectors W that maximize the following criterion:

$$J_{LLDA}(W) = \text{trace}(W^T(S_g - 2S_\ell)W), \quad (5)$$

where S_g and S_ℓ are the global and local scatter matrices defined as:

$$S_g = \frac{1}{n}XL_gX^T, \\ S_\ell = \frac{1}{n}XL_\ell X^T.$$

It is easy to check that, when we set $[K_g]_{ij} = 1/n$ for all i, j , and $[K_\ell]_{ij} = 1/n_k$ if x_i and x_j are both in the k th class, and 0 otherwise, L_g and L_ℓ respectively become $I - W_g$ and $I - W_\ell$, hence (5) includes the MMC as a special case.

In general, (5) does not require class information and can be used in an unsupervised manner. The construction of the local scatter matrix is based on the assumption that, if x_i and x_j are close, they are likely to belong to the same cluster. Under the condition that class labels are unavailable, we cannot explicitly consider the separability of different clusters, which is represented by the between-class scatter in classical LDA and the MMC. In the objective function (5), it is implicitly represented by the difference between the global scatter and the local scatter. Therefore, discriminative features can be extracted even in unsupervised scenarios. In this paper, we refer to unsupervised LLDA simply as LLDA.

Note that the reason for using the normalized graph Laplacians is that the criterion (5) without normalization may be affected by the scale of the similarity measure or by the choice of the number of nearest neighbors, since (5) is defined as the difference rather than the ratio of the global scatter to the local scatter. Also, the use of normalized graph Laplacian is known to be effective in spectral clustering (e.g. [24]).

B. Efficient algorithm for LLDA

Similarly to the case of (2), the projection matrix W that maximize the criterion (5) subject to the orthogonality constraint can be found as the eigenvectors of $S_g - 2S_\ell$ corresponding to the largest eigenvalues. When p , the number of features, is very large as in microarray data, however, it is computationally demanding to directly perform the eigenvalue decomposition (EVD) of $S_g - 2S_\ell$, which is of size $p \times p$. In [26], two approaches for computing LLDA have been presented. The first one directly computes the eigenvalues and eigenvectors, hence demands expensive computational costs. The other approach achieves this via the spectral decomposition of Laplacian matrix, but it still needs to compute the eigenvalues and eigenvectors of a $p \times p$ matrix. Even worse, the eigenvectors corresponding to the non-positive eigenvalues are discarded, thus it does not provide the exact solution to the maximization problem and results in losing discriminatory information.

Here, we propose a novel algorithm for computing W , which is particularly efficient when the feature size is much larger than the sample size, i.e. $p \gg n$, as is often the case with microarray data. The proposed algorithm is based on the following theorem (see the Appendix for the proof).

Theorem 1: Let PAQ^T be the reduced SVD [11] of $X \in \mathbb{R}^{p \times n}$, where $P \in \mathbb{R}^{p \times n}$ and $Q \in \mathbb{R}^{n \times n}$ are orthonormal

matrices and $\Lambda \in \mathbb{R}^{n \times n}$ is a diagonal matrix. Further, let $V\Delta V^T$ be the EVD of a symmetric matrix $\Lambda Q^T(L_g - 2L_\ell)QA$, where $V \in \mathbb{R}^{n \times n}$ is an orthonormal matrix and $\Delta \in \mathbb{R}^{n \times n}$ is a diagonal matrix. Then, W is constituted by the eigenvectors in PV corresponding to the largest eigenvalues in Δ .

Note that the main computation of the algorithm consists of the SVD of a $p \times n$ matrix and the EVD of an $n \times n$ matrix. Thus, it is very efficient in the case of $p \gg n$. The previous study [23] first removed the null space of the total scatter matrix via the SVD, thereby reducing the dimensionality of the data to $n - 1$, and then applied the MMC in the reduced space, where the rank of the mean subtracted matrix of X was implicitly assumed to be $n - 1$. However, in a more general case where the samples show multi-colinearity, the rank degenerates to less than $n - 1$ and needs to be estimated appropriately. In contrast, the proposed algorithm does not involve the dimension reduction before applying LLDA, allowing to deal with the degenerate case.

In this way, the graph Laplacian representation of the MMC enables both the extension to unsupervised LLDA and the efficiency of the algorithm.

IV. LLDA-RFE: FEATURE SELECTION BASED ON LLDA

The proposed algorithm for LLDA can be used in both supervised and unsupervised cases to extract discriminant features from high-dimensional data often encountered in e.g. face recognition [16], [18], [31], [32], text categorization [5], [32], and microarray cancer classification [18], [32], [33]. In the context of microarray data analysis, the features so extracted correspond to *metagenes*, a linear combination of multiple genes, but we are rather interested in identifying discriminative genes themselves.

To this end, the previous study [23] proposed to combine the MMC with recursive feature elimination (RFE). The MMC-RFE algorithm in [23] recursively removes features with the smallest absolute values of the discriminant vectors of the MMC. The RFE approach has recently proven to be effective with regression [19], [34] as well as with support vector machine (SVM) [13]. In the present study, we propose an unsupervised recursive feature selection method using the discriminant vectors of LLDA to identify features that potentially reveal clusters in the samples.

While the number of discriminant vectors extracted by classical LDA is limited to at most $c - 1$, the MMC and LLDA are capable of extracting more than $c - 1$ discriminant vectors. It can also be shown that the maximum value of $J_{LLDA}(W)$ with the obtained discriminant vectors is equal to the sum of the corresponding eigenvalues. Because the eigenvalues reflect the discrimination ability, we use only the discriminant vectors corresponding to the positive eigenvalues to calculate the weight of each feature. Let $\delta_1 \geq \delta_2 \geq \dots \geq \delta_n$ be the eigenvalues in Δ . Then, we define the weight of feature j as the sum of the absolute values of d discriminant vectors in W , i.e. $\sum_{i=1}^d \sqrt{\delta_i} |W]_{ji}|$, where d is the number of positive eigenvalues. Here, the discriminant vectors are weighted by the corresponding eigenvalues.

Our proposed algorithm, LLDA-RFE, can be summarized as follows:

Algorithm: LLDA-RFE**Input:** sample matrix $X \in \mathbb{R}^{p \times n}$ k_ℓ : the number of nearest neighbors**Output:** r top-ranked features0. Set $q \leftarrow p$,Repeat the following steps until $q = r$ 1. Construct the complete and k_ℓ -nearest neighbor graphs on X and compute K_g, K_ℓ, L_g and L_ℓ ;2. Perform the SVD of X as $X = PAQ^T$;3. Compute $Z = \Lambda Q^T(L_g - 2L_\ell)QA$;4. Perform the EVD of Z as $Z = V\Delta V^T$;5. Set W to the eigenvectors in PV corresponding to the positive eigenvalues in Δ ;6. Remove the j th feature with the smallest weight of $\sum_{i=1}^d \sqrt{\delta_i} |W_{ji}|$.7. Set $q \leftarrow q - 1$, form X and go to step 1.

V. RELATED WORK ON UNSUPERVISED FEATURE SELECTION

Data variance is one of the most common unsupervised feature selection criteria, and often used as a baseline method for comparison [15], [28]. Although variance ranking can be useful for selecting features that show large variation across all samples, it is not suited for selecting ones that contribute to characterize different clusters in the samples. Hastie *et al.* [14] developed gene shaving to select informative genes from microarray data. Gene shaving iteratively removes genes having lowest correlation with the leading principal component. Because the principal components are found so that they capture the directions of maximum variance in the data, gene shaving is also unsuitable for identifying genes that reveal different clusters. The assumption that discriminative genes exhibit large variance is not necessarily valid particularly for noisy microarray data, due to the large number of irrelevant genes. Indeed, recent studies [28], [30] have shown that variance ranking, principal component analysis and gene shaving are not effective for yielding distinctive patterns between different classes of samples.

The latest and probably more effective unsupervised methods include Laplacian score [15], the $Q - \alpha$ algorithm [30] and SVD-entropy [28]. Among these, Laplacian score and SVD-entropy are employed for comparison in this study. In the following, we give a brief overview of these two methods.

A. Laplacian score

The idea of Laplacian score is to evaluate each feature by its locality preserving power, which is similar in spirit to Locality Preserving Projection [16].

Let $f_r = (f_{r1}, \dots, f_{rn})^T$, $r = 1, \dots, p$, denote the r th feature for n samples. First, we construct a nearest neighbor graph in the same way as for the LLDA-RFE algorithm. Then, we compute the weight matrix K , the diagonal matrix $[D]_{ii} = \sum_j [K]_{ji}$, and the graph Laplacian matrix $L = D - K$. Finally, the Laplacian score L_r of the r th feature is computed as

$$L_r = \frac{\tilde{f}_r^T L \tilde{f}_r}{\tilde{f}_r^T D \tilde{f}_r},$$

where

$$\tilde{f}_r = f_r - \frac{f_r^T D e}{e^T D e} e.$$

It is worth noting that Fisher score [8] can be related to Laplacian score as shown in [15].

TABLE I

CHARACTERISTICS OF THE DATASETS USED IN THIS STUDY.

Dataset	# samples	# classes	# genes
Leukemia	38	2	7129
Colon cancer	62	2	2000
Medulloblastoma	60	2	7129
Breast cancer	76	2	4918
Lung adenocarcinoma	86	2	7129
MLL	57	3	12582
SRBCT	63	4	2308

B. SVD-entropy

Let us assume that $p > n$ for a given sample matrix $X \in \mathbb{R}^{p \times n}$. Denoting by s_j the singular values of X , an SVD-based entropy is defined as [2]:

$$E = -\frac{1}{\log(n_r)} \sum_{j=1}^{n_r} V_j \log(V_j),$$

where

$$V_j = s_j^2 / \sum_{k=1}^{n_r} s_k^2.$$

Here, $n_r \leq n$ is the number of positive singular values, which is equal to the rank of X . Then, the contribution of the i th feature to the entropy is defined as [28]:

$$CE_i = E(X_{[p \times n]}) - E(X_{[(p-1) \times n]}),$$

where $X_{[(p-1) \times n]}$ denotes the sample matrix with the i th feature being removed.

Varshavsky *et al.* [28] have proposed three feature selection strategies based on SVD-entropy: simple ranking (SR), forward selection (FS) and backward elimination (BE). SVD-entropy-based BE has high computational complexity in the case of a large number of features, hence impractical to apply to microarray datasets. This is due to the fact that CE_i is calculated on a leave-one-out basis. Indeed, they used only SR in their experiments on microarrays. Accordingly, we employ SR in this study; top-ranked features are those with the largest values of CE_i .

VI. EXPERIMENTAL RESULTS

A. Datasets and preprocessing

In the experiments, we used seven public datasets of cancer microarrays. Since binary classification is a typical and fundamental issue in diagnostic and prognostic prediction of cancer, the different methods were primarily compared using binary-class datasets: ALL versus AML for Leukemia [10], normal versus tumor for Colon cancer [1], outcome prediction on Medulloblastoma [25], Breast cancer [27], and Lung adenocarcinoma [4]. In addition, we used multi-class datasets on MLL [3] and SRBCT [17] to assess their performances. The characteristics of these datasets are summarized in Table I, and the details are given below:

- Leukemia [10]: This Affymetrix high-density oligonucleotide array dataset contains 38 samples from 2 classes of leukemia: 27 acute lymphoblastic leukemia (ALL) and 11 acute myeloid leukemia (AML). The dataset is publicly available at <http://www.broad.mit.edu/cgi-bin/cancer/datasets.cgi>.

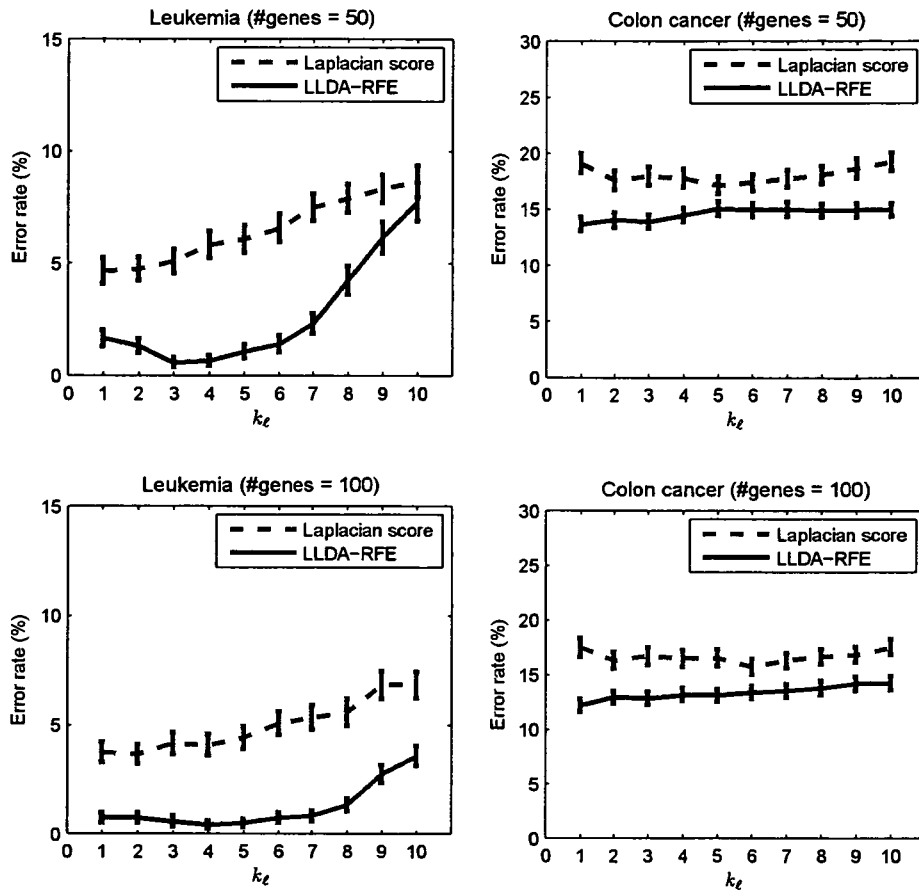


Fig. 1. Comparison between Laplacian score and LLDA-RFE using varying values of k_e for Leukemia and Colon cancer

- Colon cancer [1]: This Affymetrix high-density oligonucleotide array dataset contains 62 samples from 2 classes of colon-cancer patients: 40 normal healthy samples and 22 tumor samples. The dataset is publicly available at <http://microarray.princeton.edu/oncology/affydata/index.html>.
- Medulloblastoma dataset [25]: This Affymetrix high-density oligonucleotide array dataset contains 60 samples from 2 classes on patient survival with medulloblastoma: 21 treatment failures and 39 survivors. The dataset is publicly available at <http://www.broad.mit.edu/cgi-bin/cancer/datasets.cgi>.
- Breast cancer [27]: This cDNA microarray dataset contains 76 samples from 2 classes on five-year metastasis-free survival: 33 poor prognosis and 43 good prognosis. The dataset is publicly available at <http://www.rii.com/publications/2002/vantveer.html>.
- Lung adenocarcinoma [4]: This Affymetrix high-density oligonucleotide array dataset contains 86 samples from 2 classes on survival: an event of death for 34 and alive for 52. The dataset is publicly available at <http://dot.ped.med.umich.edu:2000/ourimage/pub/Lung/index.html>.
- MLL [3]: This Affymetrix high-density oligonucleotide array dataset contains 57 samples from 3 classes of leukemia: 20 acute lymphoblastic leukemia (ALL), 17 mixed-lineage leukemia (MLL), 20 acute myelogenous leukemia (AML). The dataset is publicly available at <http://www.broad.mit.edu/cgi-bin/cancer/datasets.cgi>.
- SRBCT [17]: This cDNA microarray dataset contains 63 samples from 4 classes of small round blue-cell tumors of childhood (SRBCT): 23 Ewing family of tumors, 20 rhabdomyosarcoma, 12 neuroblastoma, and 8 non-Hodgkin lymphoma. The dataset is publicly available at <http://research.nhgri.nih.gov/microarray/Supplement/>.

For the Leukemia, Medulloblastoma, Lung adenocarcinoma and MLL datasets, expression values were first thresholded with a floor of 100 and a ceiling of 16000, followed by a base 10 logarithmic transform. Then, each sample was standardized to zero mean and unit variance across genes. For the Colon cancer dataset, after a base 10 logarithmic transform, each sample was standardized. For the Breast cancer dataset, after the filtering of genes following [27], each sample was standardized. For the SRBCT dataset, the expression profiles already preprocessed following [17] were used.

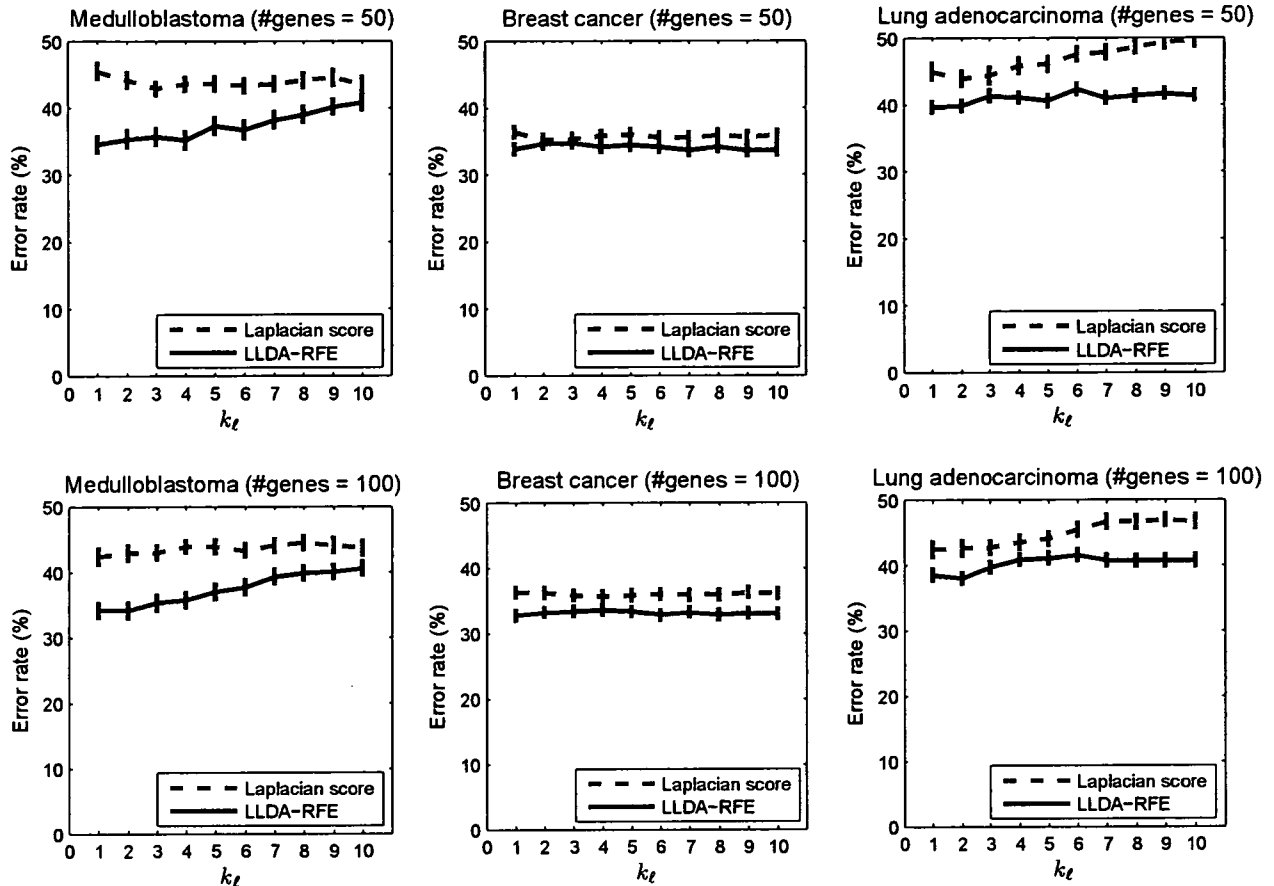


Fig. 2. Comparison between Laplacian score and LLDA-RFE using varying values of k_ℓ for Medulloblastoma, Breast cancer and Lung adenocarcinoma

B. Performance evaluation and experimental settings

We compare the performance of LLDA-RFE with those of state-of-the-art unsupervised feature selection methods, Laplacian score and SVD-entropy. The performances of the unsupervised methods are evaluated by their capability of identifying discriminative genes without using class information. Varshavsky *et al.* [28] employed the Jaccard score of clustering algorithms such as K-means, showing how clusters can be discovered by using a smaller number of genes selected from several thousand or more genes in the same samples. Wolf and Shashua [30] measured the performances by the classification accuracy of a linear SVM classifier using leave-one-out cross-validation; gene selection was performed in an unsupervised setting, but classification in a supervised setting using selected genes only. Because we also compare the performance between LLDA-RFE and a supervised gene selection method, Fisher score [8], we employed the nearest mean classifier (NMC) and measured the performances by its classification accuracy. It is known that NMC is highly effective for cancer classification despite its simplicity [29]. Note that since Fisher score is a supervised filter, it is generally expected to perform better than unsupervised methods.

We assessed the performance of each gene selection method with NMC by repeated random splitting as in [22]; the samples

were partitioned randomly in a class proportional manner into a training set consisting of two-thirds of the whole samples and a test set consisting of the held-out one-third of the samples. To avoid selection bias, gene selection was performed using only the training set, and the classification error rate of the learnt classifier was obtained using the test set. This splitting was repeated 100 times. The error rates averaged over the 100 trials and the corresponding standard error rates are reported here.

To save computational time of RFE, we removed half of the genes until less than 500, and then a single gene at a time. For the computation of graph Laplacians, we used the Euclidean distance for nearest neighbor search and a simple 0-1 weighting as the similarity of the graph, i.e. $k(x_i, x_j) = 1$ if x_i and x_j are connected, and 0 otherwise.

C. Results and discussion

1) *Effect of k_ℓ* : We first compare the performance between LLDA-RFE and Laplacian score by varying the number of nearest neighbors, on the binary-class datasets: Leukemia, Colon cancer, Medulloblastoma, Breast cancer and Lung adenocarcinoma. Figs. 1 and 2 show the average error and standard error rates for $k_\ell = 1, \dots, 10$. For LLDA-RFE, k_ℓ was fixed to the same value

during elimination. The number of genes selected and used for classification is 50 and 100.

It is clear that LLDA-RFE consistently achieves better performance than Laplacian score. This can be attributed to the difference that while Laplacian score is univariate, LLDA-RFE is multivariate and gene subsets are refined by the recursive elimination.

It may be difficult to set an appropriate value of k_ℓ in fully unsupervised settings, because we cannot rely on cross-validation unless class labels are provided and the value can also be largely dependent on the sample size of each dataset and on the potential number of clusters therein. Although an adaptive setting of the value might be preferable during elimination, our results suggest that $k_\ell = 1-3$ is a reasonable choice when applying LLDA-RFE to microarray datasets with small sample size.

2) *Comparison on binary-class datasets:* Table II shows the average error and standard error rates of NMC with four gene selection methods on the binary-class datasets. Fig. 3 plots the average error rates as a function of the number of genes from 1 to 100. The number of nearest neighbors for Laplacian score was set as follows: $k_\ell = 2$ for Leukemia, $k_\ell = 6$ for Colon cancer, $k_\ell = 3$ for Medulloblastoma and Breast cancer, and $k_\ell = 1$ for Lung adenocarcinoma. For LLDA-RFE, $k_\ell = 3$ was used for Leukemia and $k_\ell = 1$ for the other datasets.

We can observe that LLDA-RFE outperforms Laplacian score for a wide range of gene sizes. In comparison with SVD-entropy, LLDA-RFE yields lower error rates for Leukemia, Medulloblastoma and Breast cancer. Although SVD-entropy performs better for Colon cancer and Lung adenocarcinoma, LLDA-RFE consistently shows satisfactory performances for all the datasets. Also, note that LLDA-RFE performs better than Fisher score for Leukemia, Medulloblastoma and Breast cancer, despite the fact that LLDA-RFE is fully unsupervised. However, this does not imply that unsupervised gene selection is preferred to supervised one for these datasets. In fact, Fisher score, which can be viewed as a supervised version of Laplacian score, improves the performance of Laplacian score by using class information. Likewise, we can expect further improvement when using LLDA-RFE in a supervised manner.

3) *Comparison on multi-class datasets:* Table III shows the average error and standard error rates for the MLL and SRBCT datasets. Fig. 4 plots the average error rates as a function of the number of genes from 1 to 100. For Laplacian score, $k_\ell = 1$ was used for both datasets, and for LLDA-RFE, $k_\ell = 4$ and 3 were used for MLL and SRBCT, respectively.

It can be seen that LLDA-RFE reaches smaller error rates with a smaller number of genes, showing superior performance to Laplacian score and SVD-entropy. Notably, LLDA-RFE achieves even better performance than Fisher score. These results indicate that LLDA-RFE can also be useful for filtering genes from microarray samples potentially comprising multiple clusters.

In summary, our comparison using several microarray datasets has demonstrated that LLDA-RFE is effective for identifying genes that contribute to characterize different clusters in the samples. Although we used 0-1 weighting as the similarity measure, the performance could be improved by using other data-dependent similarity measures. Also, more discriminative features can be found by balancing the global and local scatters as in (3).

TABLE II
COMPARISON ON BINARY-CLASS DATASETS. BEST RESULTS IN BOLD
FACE.

# genes	Fisher score	SVD-entropy	Laplacian score	LLDA-RFE
Leukemia				
20	4.9 ± 0.6	2.6 ± 0.4	9.6 ± 1.0	3.6 ± 0.5
50	3.9 ± 0.4	1.6 ± 0.4	4.8 ± 0.5	0.6 ± 0.2
100	2.9 ± 0.4	1.6 ± 0.4	3.7 ± 0.4	0.6 ± 0.2
Colon cancer				
20	12.4 ± 0.6	14.9 ± 0.6	18.2 ± 0.8	15.9 ± 0.6
50	13.0 ± 0.6	11.5 ± 0.6	17.4 ± 0.7	13.7 ± 0.6
100	12.8 ± 0.5	11.7 ± 0.6	15.7 ± 0.7	12.2 ± 0.6
Medulloblastoma				
20	38.8 ± 0.9	36.8 ± 1.1	43.7 ± 1.0	34.1 ± 1.1
50	38.7 ± 1.0	38.4 ± 1.0	43.0 ± 0.9	34.6 ± 1.1
100	38.5 ± 1.0	37.1 ± 1.0	43.1 ± 1.0	34.2 ± 1.1
Breast cancer				
20	35.2 ± 0.9	42.6 ± 0.9	35.1 ± 0.8	33.3 ± 0.8
50	36.0 ± 0.8	42.0 ± 0.8	35.4 ± 0.8	33.8 ± 0.7
100	36.3 ± 0.8	42.4 ± 0.8	35.9 ± 0.7	32.8 ± 0.7
Lung adenocarcinoma				
20	37.8 ± 0.8	40.5 ± 0.8	45.7 ± 1.2	42.6 ± 0.7
50	36.3 ± 0.8	40.0 ± 0.7	45.0 ± 1.1	39.7 ± 0.8
100	35.1 ± 0.8	38.3 ± 0.8	42.4 ± 1.1	38.6 ± 0.8

TABLE III
COMPARISON ON MULTI-CLASS DATASETS. BEST RESULTS IN BOLD
FACE.

# genes	Fisher score	SVD-entropy	Laplacian score	LLDA-RFE
MLL				
20	7.2 ± 0.5	26.9 ± 0.9	10.2 ± 0.8	6.1 ± 0.5
50	6.6 ± 0.5	8.1 ± 0.6	9.4 ± 0.6	5.1 ± 0.5
100	5.9 ± 0.5	4.8 ± 0.4	9.1 ± 0.6	3.8 ± 0.4
SRBCT				
20	3.6 ± 0.5	22.6 ± 1.0	17.4 ± 1.2	16.2 ± 1.0
50	2.6 ± 0.4	17.4 ± 0.8	13.4 ± 1.0	12.0 ± 0.7
100	4.6 ± 0.4	11.8 ± 0.7	11.3 ± 0.9	11.4 ± 0.7

VII. CONCLUSIONS

In this paper, we have proposed a new unsupervised feature selection method based on Laplacian linear discriminant analysis (LLDA). In particular, we have extended LLDA to unsupervised cases and proposed an efficient algorithm for computing the discriminant vectors of LLDA. The LLDA-based Recursive Feature Elimination (LLDA-RFE) algorithm was applied to several microarray datasets to identify discriminative genes without using class labels.

Our comparison with other state-of-the-art unsupervised feature selection methods and with a supervised filter method has demonstrated the feasibility and effectiveness of the proposed algorithm. LLDA-RFE is capable of identifying discriminative features that contribute to reveal underlying class structures, providing a useful tool for the exploratory analysis of biological data.

A possible application of interest is the use of LLDA-RFE in semi-supervised scenarios; when labels are partially given, we

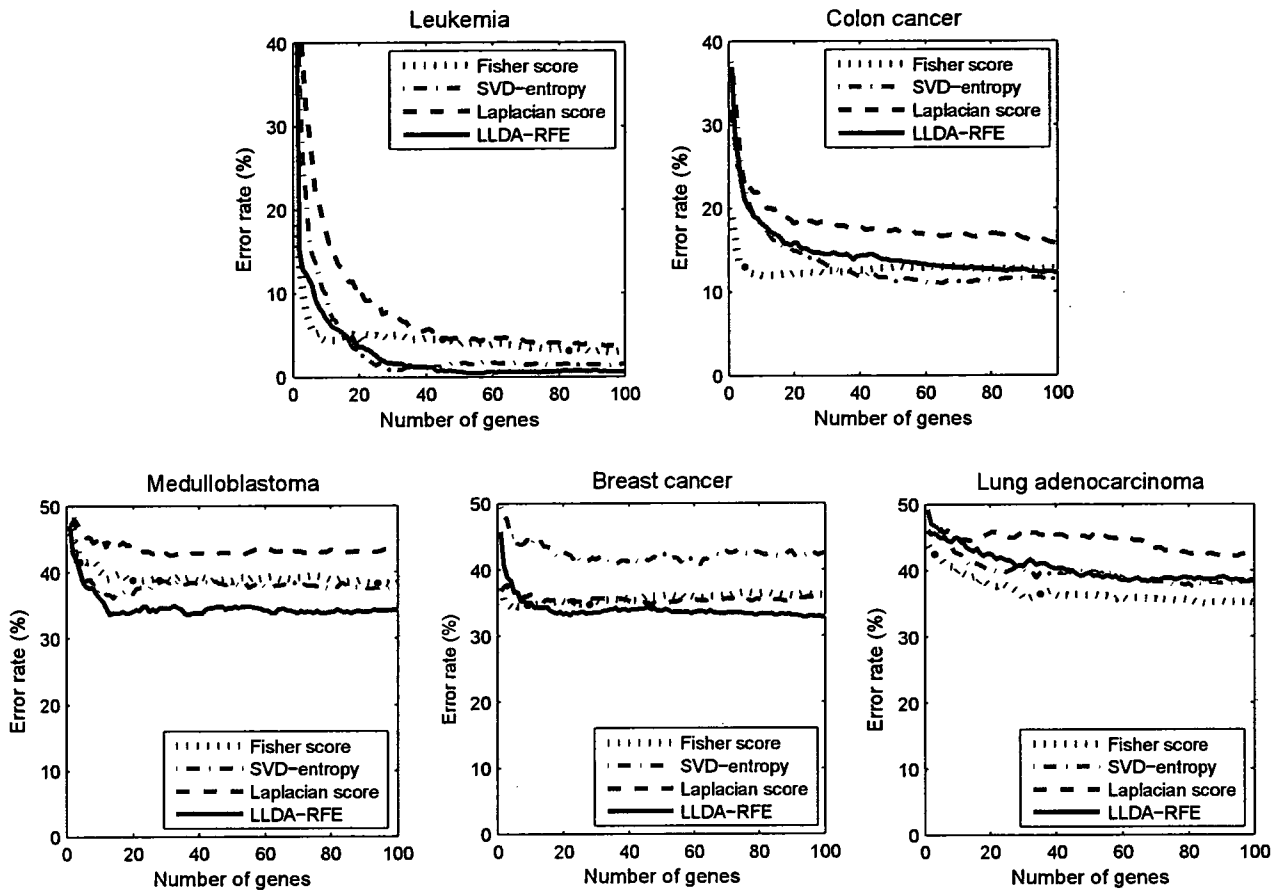


Fig. 3. Average error rates for Fisher score, SVD-entropy, Laplacian score and LLDA-RFE on binary-class datasets.

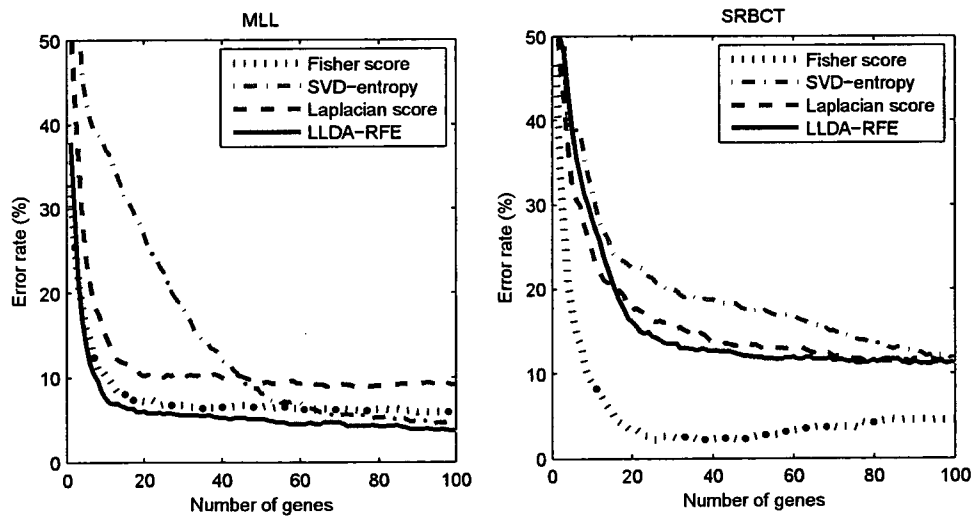


Fig. 4. Average error rates for Fisher score, SVD-entropy, Laplacian score and LLDA-RFE on multi-class datasets.

construct a graph such that samples with the same class are always connected, while those with different classes disconnected, and those with no class labels are adaptively connected or disconnected depending on the nearest neighbors.

APPENDIX PROOF OF THEOREM 1.

It is straightforward to verify that $S_g - 2S_\ell$ can be decomposed as follows:

$$\begin{aligned} S_g - 2S_\ell &= \frac{1}{n} X(L_g - 2L_\ell)X^T \\ &= \frac{1}{n} PAQ^T(L_g - 2L_\ell)(PAQ^T)^T \\ &= \frac{1}{n} PAQ^T(L_g - 2L_\ell)QAP^T \\ &= \frac{1}{n} PV\Delta V^T P^T \\ &= \frac{1}{n} (PV)\Delta(PV)^T. \end{aligned}$$

Since both P and Q are orthonormal,

$$(PV)^T(PV) = V^T P^T P V = V^T V = I.$$

Hence, the theorem holds. \square

ACKNOWLEDGMENTS

This work was supported by grant from the 21st Century COE program "Knowledge Information Infrastructure for Genome Science". A part of this work was done while S. Nijima was in the Graduate School of Systems Life Sciences at Kyushu University, and also supported in part by KAKENHI (Grant-in-Aid for Scientific Research) on Priority Areas "Comparative Genomics" from the Ministry of Education, Culture, Sports, Science and Technology of Japan (to Prof. S. Kuhara).

REFERENCES

- [1] U. Alon, N. Barkai, D.A. Notterman, K. Gish, S. Ybarra, D. Mack, and A.J. Levine, "Broad Patterns of Gene Expression Revealed by Clustering Analysis of Tumor and Normal Colon Tissues Probed by Oligonucleotide Arrays," *Proc. Nat'l Academy of Sciences USA*, vol. 96, pp. 6745-6750, 1999.
- [2] O. Alter, P.O. Brown, and D. Botstein, "Singular Value Decomposition for Genome-wide Expression Data Processing and Modeling," *Proc. Nat'l Academy of Sciences USA*, vol. 97, pp. 10101-10106, 2000.
- [3] S.A. Armstrong, J.E. Staunton, L.B. Silverman, R. Pieters, M.L. den Boer, M.D. Minden, S.E. Sallan, E.S. Lander, T.R. Golub, and S.J. Korsmeyer, "MLL Translocations Specify a Distinct Gene Expression Profile That Distinguishes a Unique Leukemia," *Nature Genetics*, vol. 30, pp. 41-47, 2002.
- [4] D.G. Beer, S.L.R. Kardia, C.-C. Huang, T.J. Giordano, A.M. Levin, D.E. Misek, L. Lin, G. Chen, T.G. Gharib, D.G. Thomas, M.L. Lizyness, R. Kuick, S. Hayasaka, J.M.G. Taylor, M.D. Iannettoni, M.B. Orringer, and S. Hanash, "Gene-expression Profiles Predict Survival of Patients with Lung Adenocarcinoma," *Nature Medicine*, vol. 8, no. 8, pp. 816-824, 2002.
- [5] D. Cai, X. He, and J. Han, "Document Clustering Using Locality Preserving Indexing," *IEEE Trans. Knowledge and Data Engineering*, vol. 17, no. 12, pp. 1624-1637, 2005.
- [6] F.R.K. Chung, *Spectral Graph Theory*, Regional Conference Series in Mathematics, no. 92, 1997.
- [7] C.H.Q. Ding, "Unsupervised Feature Selection via Two-way Ordering in Gene Expression Analysis," *Bioinformatics*, vol. 19, no. 10, pp. 1259-1266, 2003.
- [8] S. Dudoit, J. Fridlyand, and T. Speed, "Comparison of Discrimination Methods for the Classification of Tumors Using Gene Expression Data," *J. Amer. Statist. Assoc.*, vol. 97, pp. 77-87, 2002.
- [9] K. Fukunaga, *Introduction to Statistical Pattern Recognition*, second ed. Boston, MA: Academic Press, 1990.
- [10] T.R. Golub, D.K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J.P. Mesirov, H. Coller, M.L. Loh, J.R. Downing, M.A. Caligiuri, C.D. Bloomfield, and E.S. Lander, "Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring," *Science*, vol. 286, pp. 531-537, 1999.
- [11] G.H. Golub and C.F. Van Loan, *Matrix Computations*, third ed. Baltimore, MD: Johns Hopkins University Press, 1996.
- [12] I. Guyon and A. Elisseeff, "An Introduction to Variable and Feature Selection," *J. Machine Learning Research*, vol. 3, pp. 1157-1182, 2003.
- [13] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene Selection for Cancer Classification Using Support Vector Machines," *Machine Learning*, vol. 46, pp. 389-422, 2002.
- [14] T. Hastie, R. Tibshirani, M.B. Eisen, A. Alizadeh, R. Levy, L. Staudt, W.C. Chan, D. Botstein, and P.O. Brown, "'Gene Shaving' as a Method for Identifying Distinct Sets of Genes with Similar Expression Patterns," *Genome Biology*, vol. 1, no. 2, research0003, 2000.
- [15] X. He, D. Cai, and P. Niyogi, "Laplacian Score for Feature Selection," In Y. Weiss, B. Schölkopf and J. Platt (eds.), *Advances in Neural Information Processing Systems 18*, pp. 507-514, Cambridge, MA: MIT Press, 2006.
- [16] X. He, S. Yan, Y. Hu, P. Niyogi, and H.-J. Zhang, "Face Recognition Using Laplacianfaces," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 3, pp. 328-340, 2005.
- [17] J. Khan, J.S. Wei, M. Ringnér, L.H. Saal, M. Ladanyi, F. Westermann, F. Berthold, M. Schwab, C. Antonescu, C. Peterson, and P.S. Meltzer, "Classification and Diagnostic Prediction of Cancers Using Gene Expression Profiling and Artificial Neural Networks," *Nature Medicine*, vol. 7, no. 6, pp. 673-679, 2001.
- [18] H. Li, T. Jiang, and K. Zhang, "Efficient and Robust Feature Extraction by Maximum Margin Criterion," *IEEE Trans. Neural Networks*, vol. 17, no. 1, pp. 157-165, 2006.
- [19] F. Li and Y. Yang, "Analysis of Recursive Gene Selection Approaches from Microarray Data," *Bioinformatics*, vol. 21, no. 19, pp. 3741-3747, 2005.
- [20] Q. Liu, X. Tang, H. Lu, and S. Ma, "Face Recognition Using Kernel Scatter-difference-based Discriminant Analysis," *IEEE Trans. Neural Networks*, vol. 17, no. 4, pp. 1081-1085, 2006.
- [21] M. Loog, "On an Alternative Formulation of the Fisher Criterion That Overcomes the Small Sample Problem," *Pattern Recognition*, vol. 40, pp. 1753-1755, 2007.
- [22] S. Michiels, S. Koscielny, and C. Hill, "Prediction of Cancer Outcome with Microarrays: a Multiple Random Validation Strategy," *Lancet*, vol. 365, pp. 488-492, 2005.
- [23] S. Nijima and S. Kuhara, "Recursive Gene Selection Based on Maximum Margin Criterion: a Comparison with SVM-RFE," *BMC Bioinformatics*, vol. 7, 543, 2006.
- [24] A.Y. Ng, M.I. Jordan, and Y. Weiss, "On Spectral Clustering: Analysis and an Algorithm," In T. Dietterich, S. Becker and Z. Ghahramani(eds.), *Advances in Neural Information Processing Systems 14*, pp. 849-856, Cambridge, MA: MIT Press, 2002.
- [25] S.L. Pomeroy, P. Tamayo, M. Gaasenbeek, L.M. Sturla, M. Angelo, M.E. McLaughlin, J.Y.H. Kim, L.C. Goumnerova, P.M. Black, C. Lau, J.C. Allen, D. Zagzag, J.M. Olson, T. Curran, C. Wetmore, J.A. Biegel, T. Poggio, S. Mukherjee, R. Rifkin, A. Califano, G. Stolovitzky, D.N. Louis, J.P. Mesirov, E.S. Lander, and T.R. Golub, "Prediction of Central Nervous System Embryonal Tumor Outcome Based on Gene Expression," *Nature*, vol. 415, pp. 436-442, 2002.
- [26] H. Tang, T. Fang, and P.-F. Shi, "Laplacian Linear Discriminant Analysis," *Pattern Recognition*, vol. 39, pp. 136-139, 2006.
- [27] L.J. van 't Veer, H. Dai, M.J. van de Vijver, Y.D. He, A.A.M. Hart, M. Mao, H.L. Peterse, K. van der Kooy, M.J. Marton, A.T. Witteveen, G.J. Schreiber, R.M. Kerkhoven, C. Roberts, P.S. Linsley, R. Bernards, and S.H. Friend, "Gene Expression Profiling Predicts Clinical Outcome of Breast Cancer," *Nature*, vol. 415, pp. 530-536, 2002.
- [28] R. Varshavsky, A. Gottlieb, M. Linal, and D. Horn, "Novel Unsupervised Feature Filtering of Biological Data," *Bioinformatics*, vol. 22, no. 14, pp. e507-e513, 2006.
- [29] L.F.A. Wessels, M.J.T. Reinders, A.A.M. Hart, C.J. Veenman, H. Dai, Y.D. He, L.J. van't Veer, "A Protocol for Building and Evaluating Predictors of Disease State Based on Microarray Data," *Bioinformatics*, vol. 21, no. 19, pp. 3755-3762, 2005.
- [30] L. Wolf and A. Shashua, "Feature Selection for Unsupervised and Supervised Inference: the Emergence of Sparsity in a Weight-based Approach," *J. Machine Learning Research*, vol. 6, pp. 1855-1887, 2005.
- [31] S. Yan, D. Xu, B. Zhang, H.-J. Zhang, Q. Yang, and S. Lin, "Graph Embedding and Extensions: a General Framework for Dimensionality Reduction," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 29, no. 1, pp. 40-51, 2007.

- [32] J. Ye, "Characterization of a Family of Algorithms for Generalized Discriminant Analysis on Undersampled Problems," *J. Machine Learning Research*, vol. 6, pp. 483-502, 2005.
- [33] J. Ye, T. Li, T. Xiong, and R. Janardan, "Using Uncorrelated Discriminant Analysis for Tissue Classification with Gene Expression Data," *IEEE/ACM Trans. Comput. Biol. Bioinform.*, vol. 1, no. 4, pp. 181-190, 2004.
- [34] J. Zhu and T. Hastie, "Classification of Gene Microarrays by Penalized Logistic Regression," *Biostatistics*, vol. 5, no. 3, pp. 427-443, 2004.

GLIDA: GPCR—ligand database for chemical genomics drug discovery—database and tools update

Yasushi Okuno^{1,*}, Akiko Tamon², Hiroaki Yabuuchi¹, Satoshi Niijima¹,
Yohsuke Minowa¹, Koichiro Tonomura¹, Ryo Kunimoto¹ and Chunlai Feng¹

¹Department of Pharmacoinformatics, Center for Integrative Education of Pharmacy Frontier, Graduate School of Pharmaceutical Sciences, Kyoto University and ²Bio Science Group, IT Solution Div.1, Industry Solution Business Unit, Mitsui Knowledge Industry, Osaka city, Japan

Received September 6, 2007; Revised October 14, 2007; Accepted October 15, 2007

ABSTRACT

G-protein coupled receptors (GPCRs) represent one of the most important families of drug targets in pharmaceutical development. GLIDA is a public GPCR-related Chemical Genomics database that is primarily focused on the integration of information between GPCRs and their ligands. It provides interaction data between GPCRs and their ligands, along with chemical information on the ligands, as well as biological information regarding GPCRs. These data are connected with each other in a relational database, allowing users in the field of Chemical Genomics research to easily retrieve such information from either biological or chemical starting points. GLIDA includes a variety of similarity search functions for the GPCRs and for their ligands. Thus, GLIDA can provide correlation maps linking the searched homologous GPCRs (or ligands) with their ligands (or GPCRs). By analyzing the correlation patterns between GPCRs and ligands, we can gain more detailed knowledge about their conserved molecular recognition patterns and improve drug design efforts by focusing on inferred candidates for GPCR-specific drugs. This article provides a summary of the GLIDA database and user facilities, and describes recent improvements to database design, data contents, ligand classification programs, similarity search options and graphical interfaces. GLIDA is publicly available at <http://pharminfo.pharm.kyoto-u.ac.jp/services/glida/>. We hope that it will prove very useful for Chemical Genomics research and GPCR-related drug discovery.

INTRODUCTION

The family of G-protein coupled receptors (GPCRs) represents one of the most important classes of pharmaceutical targets (1). Among the more than 1000 GPCRs encoded in the human genome, more than 400 are of potential therapeutic interest (2). Currently the drugs available on the market address only 30 GPCRs, which represent a small fraction of the GPCR target family. A large majority of human-derived GPCRs still remain promising drug targets, and thus a key goal of GPCR research related to drug design is to identify new ligands for such target GPCRs.

With the unprecedented accumulation of genomic information, databases and bioinformatics have become essential tools to guide GPCR research (3). The GPCRDB (2) and IUPHAR receptor database (IUPHAR-RD) (4) are representatives of widely used public databases covering GPCRs. These databases, which provide substantial data on the GPCR proteins and pharmacological information on receptor proteins containing GPCRs, are mainly focused on biological aspects of the GPCR gene products or proteins. In spite of the significance of ligand compounds as drug leads, the relationships between GPCRs and their ligands and/or chemical information on the ligands themselves are not yet fully covered.

On the other hand, there is increasing interest in publicly collecting and applying chemical as well as biological information in the post-genome era (5–8). This new trend is called ‘Chemical Genomics’, and it aims to identify all possible chemical ligands and drugs for all targets families (9,10). There is a vast amount of information on the interactions between small molecules and proteins/genes. However, compound–protein interactions have not yet been analyzed on a large scale, and there are no effective methods to extract meaningful

*To whom correspondence should be addressed. Tel: +81 75 753 4559; Fax: +81 75 753 4544; Email: okuno@pharm.kyoto-u.ac.jp

© 2007 The Author(s)

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/2.0/uk/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

information from the data in a comprehensive manner. Therefore, we need to integrate chemoinformatics and bioinformatics into a common computational platform for mining of Chemical Genomics data (11).

GLIDA (GPCR-Ligand DAtabase) is a public GPCR-related Chemical Genomics database designed to simultaneously mine biological information on GPCRs and chemical information on their ligands. It provides various analytical data regarding GPCR–ligand correlations by incorporating bioinformatics and chemoinformatics techniques, and thus it should prove very useful for GPCR-related drug discovery from the viewpoint of Chemical Genomics research. There have been several major improvements to GLIDA since it was last described in Ref. (12): (i) there are more increments in the entries of the ligands and the corresponding ligand–GPCR pairs; (ii) the ligands are originally classified using a new strategy; (iii) additional options are available within the similarity search program for the GPCRs and ligands and (iv) the graphical interface to display the correlation maps between GPCRs and ligands has been enhanced.

DATA CONTENTS

GLIDA contains three types of primary data: biological information on GPCRs, chemical information on their ligands and information on binding of the GPCR–ligand pairs. The GPCR entries were acquired from human, mouse and rat entries deposited in the GPCRDB because these three species include sufficient information regarding ligands, and rats and mice are representative model animals used in drug discovery research. The ligand-binding information was manually collected and curated using various public web sites and commercial databases such as the IUPHAR-RD, PubMed (5), PubChem (5), DrugBank (13), Ki Database (14) and MDL ISIS/Base 2.5. Table 1 indicates the size and scope of the GLIDA database. In particular, we have dramatically expanded the entry number of ligands and the corresponding ligand–GPCR pairs. The latest GLIDA version includes 24077 ligand entries and 39140 GPCR–ligand pair entries, representing nearly 35-fold and 20-fold increases, respectively, since the last publication of GLIDA in 2006. The total number of GPCR entries remains unchanged, but entries with associated ligand information have increased slightly, suggesting that it is difficult to de-orphan the GPCRs whose ligands have not yet been identified (15).

GPCR and ligand data

The database lists general information on GPCR and ligand data, respectively. The general information table listing GPCRs contains gene names, family names, protein sequences (in fasta format) and links to other biological databases, such as GPCRDB, UniProt (16), IUPHAR-RD, Entrez Gene (17) and KEGG (18). The ligand result page provides a general information table containing names, molecular structures, CAS registry numbers, formulas, molecular weights, structure files and links to

Table 1. The current numbers of GLIDA ligands and GPCRs and their respective links.

Information item	Number of entries
GPCR entries	3738
Links to Entrez Gene	3073
Links to GPCRDB	3738
Links to UniProt	3738
Links to IUPHAR	446
Links to KEGG	595
Ligand entries	24077
Cas registry number	2425
Molecular structure	23 216 ^a
Links to PubChem	1821
Links to ChEBI	103
Links to KEGG	664
Links to DrugBank	479
Cluster number	300 ^b
GPCR–ligand pair entries	39 140
GPCR entries	410
Ligand entries	24 077
Activity	
Agonist	8305
Full Agonist	2325
Partial Agonist	262
Antagonist	28 132
Inverse Agonist	116

^aMolecular structures consist of MDL MOL files and original files converted into KEGG atom types. The numbers of MDL MOL files and KEGG-type files are 23 216 and 23 214, respectively. PCA calculation was performed for 23 214 KEGG-type files.

^bThis cluster number (300) is different from the number of the selected principal components (314). No compounds were assigned to 14 principal components.

PubChem, KEGG, ChEBI (8) and DrugBank that are in publicly available chemical databases.

Information on binding of GPCR–ligand pairs

The interaction information relating GPCRs to particular ligands, a key issue for GPCR-related drug discovery, is deposited in a relational database. GLIDA allows users to retrieve GPCR–ligand-binding information dynamically and continuously. When users retrieve a GPCR (or ligand) entry, its result page displays all entries showing the corresponding ligands (or GPCR) entries with their binding activity types, as well as references. The references are hyperlinked with the corresponding PubMed literature. The activity types include agonist, antagonist and full, partial or inverse agonist (Table 1). Here the detail annotations such as full, partial or inverse agonist are not finished yet. The ligands classified as agonists are possible full agonists or partial agonists. Inverse agonists can be also contained among the antagonists.

WEB INTERFACE AND APPLICATION

GLIDA is available at <http://pharminfo.pharm.kyoto-u.ac.jp/services/glida/>. The web interface of GLIDA includes a GPCR search page (Figure 1a) and a ligand search page (Figure 1b). Each page consists of a classification menu and a keyword search box. The users can search a GPCR (or ligand) manually using the classification tool,

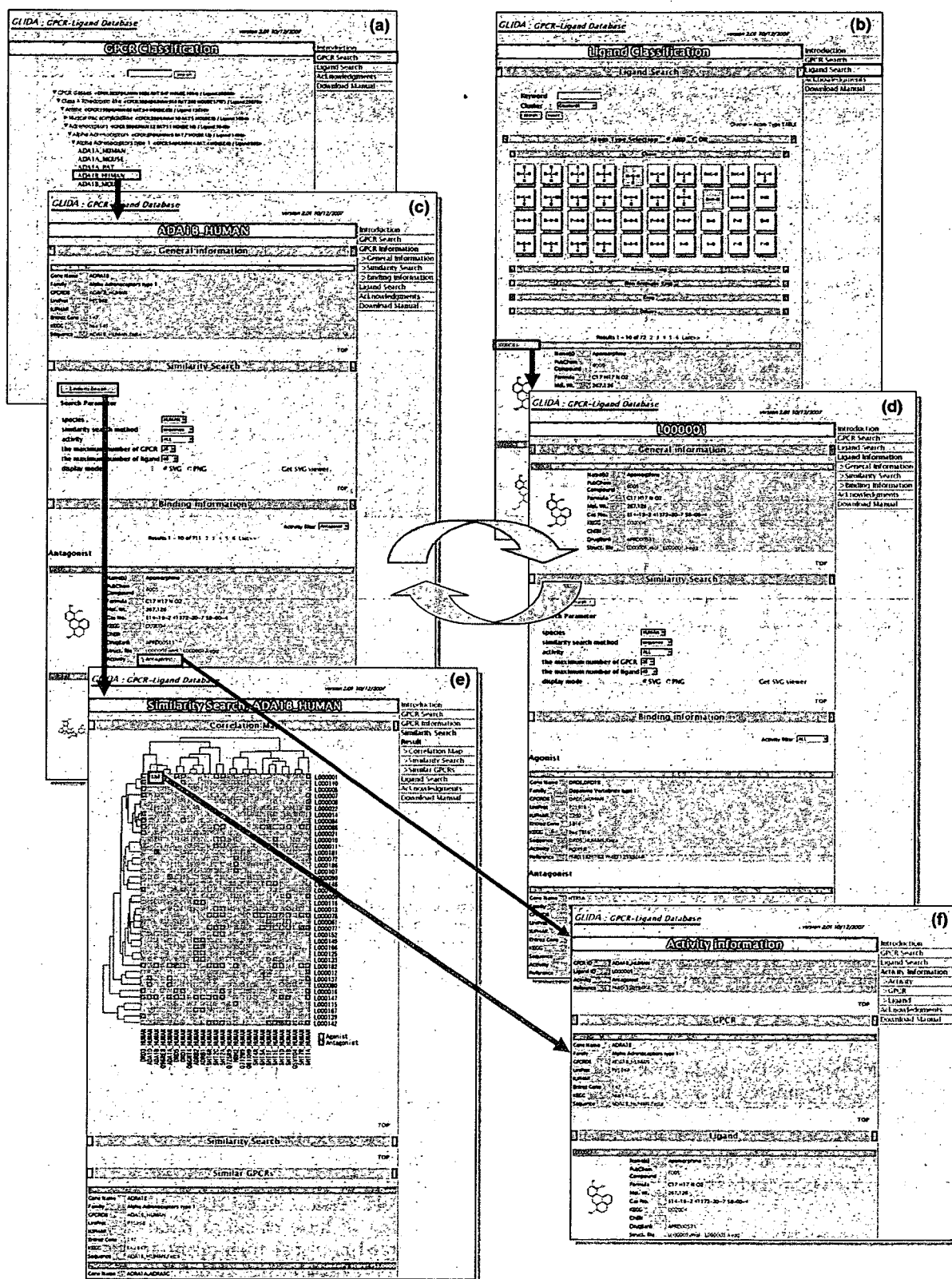


Figure 1. A screenshot of GLIDA showing linked relations among search pages (a and b), result pages (c and d), an analytical report page (e), and a binding information page (f). The analytical report page consists of a correlation map and a list resulting from a similarity search. Red and blue colors of the spots on the correlation map indicate the ligand activities of antagonists including inverse agonist and agonists including full/partial agonist, respectively.

GLIDA : GPCR-Ligand Database
version 2.01 10/12/2007

Ligand Classification

Introduction
GPCR Search
Ligand Search
Acknowledgments
Download Manual

Ligand Search

Keyword 2) Select Cluster

Cluster Cluster87
Select Cluster
Cluster53
Cluster85
Cluster86
Cluster87
Cluster90
Cluster91
Cluster98
Cluster113
Cluster121
Cluster130
Cluster141
Cluster154
Cluster160
Cluster170
Cluster174
Cluster177
Cluster181
Cluster187
Cluster194
Cluster196
Cluster199
Cluster205
Cluster210
Cluster214
Cluster216
Cluster227
Cluster228
Cluster235
Cluster237

3) Click Search

1) Select Atom Type

Atom Type selection CAND: G OR

Chain

Aromatic Ring

Non-Aromatic Ring

Cluster - Atom Type TABLE

Cluster	Atom Type	Number of Ligands
Cluster 1	R-C-H, R-C-R, R-N-R, P-O-R, P-R, S-O	1058 ligands
Cluster 2	R-C-H, R-C-R, R-N-R, P-O-R, P-R, S-O	219 ligands
Cluster 3	R-C-H, R-C-R, R-N-R, P-O-R, P-R, S-O	89 ligands
Cluster 4	R-C-H, R-C-R, R-N-R, P-O-R, P-R, S-O	313 ligands
Cluster 5	R-C-H, R-C-R, R-N-R, P-O-R, P-R, S-O	183 ligands

Results 1 - 10 of 34 2 3 4

1002247	<p>Name(s) DP 092</p> <p>PubChem Compound 5574513</p> <p>Formula C23 H31 N3 O2 S</p> <p>Mol. Wt. 413.213</p> <p>Cas No. 56384-09-7</p> <p>KEGG CHEBI</p> <p>DrugBank</p> <p>Struct. file L002247.mol L002247.fegg</p> <p>Binding Antagonist(2)</p>
1002278	<p>Name(s) L002278</p> <p>PubChem Compound</p> <p>Formula C23 H31 N3 O2 S</p> <p>Mol. Wt. 413.213</p> <p>Cas No.</p> <p>KEGG CHEBI</p> <p>DrugBank</p> <p>Struct. file L002278.mol L002278.fegg</p> <p>Binding Antagonist(1)</p>

Figure 2. A screenshot of the ligand search process on the ligand classification page. Users can search the ligands from two starting points: keyword search and cluster selection. If they have a chemical structure of their query compound, the ligand search is performed using the cluster selection tool as follows. Selecting a set of atom types (step 1) that the query compound contains, the pull-down menu of cluster selection displays the list of the only clusters that include selected atom types as the principal components (step 2). By selecting a cluster from the list, users can check the principal component's atoms on the upper right section of the page. Finally, upon clicking the search button, GLIDA displays the list of all ligands classified in the selected cluster (step 3). The 'Atom Type TABLE' button links the user to the page showing the cluster size and representative atom types for each cluster.

or automatically by using the keyword search function. Every GPCR (or ligand) has its own results page (Figure 1c or d) containing a general information table regarding a GPCR (or ligand), a table of its correlated ligands (or GPCRs) and a menu button to carry out a similarity search and correlation analysis.

Classification of GPCRs and ligands

The GPCR classification table on the search page was adapted from the phylogenetic tree within the GPCRDB information system (<http://www.gpcr.org/7tm/phylo/phylo.html>). The GPCR classification table displays the entries of the corresponding GPCRs at the tree branches,

and these are hyperlinked with the corresponding result pages (Figure 1a). GLIDA also provides an original ligand classification (Figures 1b and 2). With the great increase in ligand entries, we have to improve our method of classifying all the ligands in GLIDA. Hierarchical clustering and its tree representation, which were used in the old version of GLIDA, are unsuitable for the data mining of huge chemical databases. We therefore have adopted principal component analysis (PCA) for clustering of 23214 ligand structures in this new version, as follows. We generated frequency profiles of the atoms and the bonds converted into the KEGG atom types from MDL MOL files of ligand entries (19). The KEGG-type profile for each ligand is shown in 'Struct. file' item of

general information table of GLIDA. PCA was applied to the data matrix consisting of 700 KEGG-type features' columns and 23214 ligand entries' rows. The resulting principal components (PCs) constitute a new set of linearly independent, orthogonal axes that capture the directions of maximum variance in the data. The samples (chemical compounds) were then projected onto these PC axes. Herein, we used the top 314 PCs as seeds of clusters that account for >80% (cumulative proportion) of the total variance. Finally, each compound was assigned to the PC cluster having the maximum score among the 314 PCs. In order to annotate the features of each cluster (PC), we selected for each PC the atom types and their bonds corresponding to the top 10 loadings having the largest magnitude. The ligand classification page displays a table of all the atom types selected by PCA (Figure 2). By clicking on some of the atoms in this table, users can search clusters that include the selected atom types. Consequently, the ligands relevant to users' interests are included in the retrieved cluster.

Similarity search and correlation map for GPCRs and ligands

The fact that similar proteins bind similar ligands is the underlying principle of the Chemical Genomics approach to drug discovery (11). GLIDA has a variety of similarity search functions for GPCRs and ligands, respectively, on its result pages (Figure 1c or d). Alignment scores for protein sequences generated by the BLAST algorithm provide similarity measures for GPCRs. In addition to sequence similarity, gene expression patterns in tissue origins and developmental stages were used as similarity measures. The expression data for each GPCR was generated from the EST sequences in different libraries served from NCBI/Unigene (<http://www.ncbi.nlm.nih.gov/UniGene/ddd.cgi>). We can thereby retrieve the GPCRs that present tissue-/stage-specific distribution similar to a query GPCR. For example, co-expression information on specific GPCRs enables us to speculate about GPCR-heterodimerization that might have an effect on their activity (1). Ligand similarity is defined by the dissimilarity (distance) of frequency profile patterns generated from the constitutive atoms and bonds of the chemical structure, using the KEGG atom types (19,20). From the similarity search, the most similar GPCRs (or ligands) within the users' selected parameters are retrieved and listed with their similarity scores on an analytical report page (Figure 1e). In the latest GLIDA version, various parameters have been added as search options, such as selections of species, ligand activities, displayed number of GPCRs/ligands and map graphical mode. As another result of similarity search calculations, GLIDA illustrates the correlation map (Figure 1e) showing homologous GPCRs (or ligands) and their ligands (or GPCRs) that are retrieved. This map shows spots that match the GPCRs and their ligands in a 2D matrix. The ordering along the *x*-axis and the *y*-axis are calculated respectively by two-way clustering of the GPCRs and the ligands, based on their similarities. In particular, the ordering along the *x*- and *y*-axes allows users to evaluate

the sequence similarities among GPCRs and the correlation coefficients among ligands simultaneously. By analyzing the correlation patterns between GPCRs and ligands that are illustrated by these maps, we can gain detailed knowledge about their interactions. We can then utilize this information to infer possible candidates for development of GPCR-specific drugs. Furthermore, we have enhanced a graphical interface to display the correlation map between GPCRs and ligands. Graphics are an important tool to aid visualization and interpretation of high-dimensional data. The old version of GLIDA used only the PNG (Portable Network Graphics) format to display a GPCR–ligand correlation map. Due to the great increase in entries, the latest GLIDA version introduces the SVG (Scalable Vector Graphics) format, which is adaptable to an enormous correlation map size. The SVG vector image can be scaled indefinitely without loss of image quality, while the PNG bitmap image cannot. Users must install the free plug-in software on their computer in advance to use the SVG format (<http://www.adobe.com/svg/viewer/install/>). In the case of uninstalled devices, PNG representation should be selected as a graphical mode. Figure 1 shows an example of the GPCR–ligand search and analysis process starting from a GPCR query using GLIDA.

DISCUSSION AND FUTURE DIRECTIONS

GLIDA provides a unique database useful for GPCR-related Chemical Genomics research and drug discovery. GLIDA is distinct from other public Chemical Genomics databases because it contains original, GPCR-specific chemical entries and offers a common mining platform of bioinformatics and cheminformatics. GLIDA provides several advantages over other databases, in that a search can be started either from a GPCR or from a ligand. Thus, searches can be carried out in a dynamic and user-friendly way. GLIDA's coverage of chemical and biological information simultaneously also provides an advantage to users by saving them the time and labor required to search multiple databases. The ligand search page is another distinct characteristic of GLIDA, in that it displays the structural distribution of ligands. It thereby facilitates research on GPCR-related drugs by incorporating structural aspects of the ligand compounds into the search. The analytical report pages resulting from the calculated structural similarities of GPCRs and ligands can give the user deep insights into the GPCR–ligand relationships. The lists of neighboring ligands (or GPCRs) and the correlation maps are useful visualization tools for analyzing correlations among the structural features and the GPCR–ligand-binding properties. Because this database system can be applied to proteins other than the GPCR family, it may also be considered as a promising database for other types of Chemical Genomics research. One critical issue is how to define similarity metrics for proteins and ligands, because the underlying principle of GLIDA is that similar receptors bind similar ligands. For example, ligand similarity can be defined by manifold representations such as graph, fingerprint and descriptors.

Protein similarity can be also measured in different ways such as overall sequence homology (phylogenetic relationships), consensus motifs, common binding sites, 3D structures and reported functional annotations. Therefore we will add new menus incorporating these various similarity metrics for GPCRs and ligands. GLIDA will be updated continuously. In particular, we are now planning to add the drawing tool of chemical structures and to expand the ligand-searching function for an arbitrary chemical query.

ACKNOWLEDGEMENTS

This work was supported by grants from the Ministry of Education, Culture, Sports, Science and Technology of Japan, from the JSPS, KAKENHI, Grant-in-Aid for Publication of Scientific Research Results and from the Ministry of Health, Labour and Welfare of Japan. Financial support from the SUNTORY INSTITUTE FOR BIOORGANIC RESEARCH, the TATEISI SCIENCE AND TECHNOLOGY FOUNDATION and the Okawa Foundation for Information and Telecommunications is gratefully acknowledged. Funding to pay the Open Access publication charges for this article was provided by the Ministry of Education, Culture, Sports, Science and Technology of Japan.

Conflict of interest statement. None declared.

REFERENCES

- George, S.R., O'Dowd, B.F. and Lee, S.P. (2002) G-protein-coupled receptor oligomerization and its potential for drug discovery. *Nature Rev. Drug Discov.*, **1**, 808–820.
- Horn, F., Bettler, E., Oliveira, L., Campagne, F., Cohen, F.E. and Vriend, G. (2003) GPCRDB information system for G protein-coupled receptors. *Nucleic Acids Res.*, **31**, 294–297.
- Strachan, R., Ferrara, G. and Roth, B.L. (2006) Screening the receptorome: an efficient approach for drug discovery and target validation. *Drug Discov. Today*, **11**, 708–716.
- Foord, S.M., Bonner, T.I., Neubig, R.R., Rosser, E.M., Pin, J.P., Davenport, A.P., Spedding, M. and Harman, A.J. (2005) International Union of Pharmacology. XLVI. G Protein-Coupled Receptor List. *Pharmacol. Rev.*, **57**, 279–288.
- Wheeler, D.L., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Chetvermin, V., Church, D.M., DiCuccio, M., Edgar, R. *et al.* (2007) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.*, **35**, 12.
- Schreiber, S.L. (2004) Stuart Schreiber: biology from a chemist's perspective. Interview by Joanna Owens. *Drug Discov. Today*, **9**, 299–303.
- Goto, S., Okuno, Y., Hattori, M., Nishioka, T. and Kanehisa, M. (2002) LIGAND: database of chemical compounds and reactions in biological pathways. *Nucleic Acids Res.*, **30**, D402–D404.
- Brooksbank, C., Cameron, G. and Thornton, J. (2005) The European Bioinformatics Institute's data resources: towards systems biology. *Nucleic Acids Res.*, **33**, D46–D53.
- Zerhouni, E. (2003) The NIH Roadmap. *Science*, **302**, 63–72.
- Dobson, C.M. (2004) Chemical space and biology. *Nature*, **432**, 824–828.
- Klabunde, T. (2007) Chemogenomic approaches to drug discovery: similar receptors bind similar ligands. *Br. J. Pharmacol.*, **152**, 5–7.
- Okuno, Y., Yang, J., Taneishi, K., Yabuuchi, H. and Tsujimoto, G. (2006) GLIDA: GPCR-ligand database for chemical genomic drug discovery. *Nucleic Acids Res.*, **34**, D673–D677.
- Wishart, D.S., Knox, C., Guo, A.C., Shrivastava, S., Hassanali, M., Stothard, P., Chang, Z. and Woolsey, J. (2006) DrugBank: a comprehensive resource for in silico drug discovery and exploration. *Nucleic Acids Res.*, **34**, D668–D672.
- Roth, B.L., Lopez, E., Beischel, S., Westkaemper, R.B. and Evans, J.M. (2004) Screening the receptorome to discover the molecular targets for plant-derived psychoactive compounds: a novel approach for CNS drug discovery. *Pharmacol. Ther.*, **102**, 99–110.
- Civelli, O. (2005) GPCR deorphanizations: the novel, the known and the unexpected transmitters. *Trends Pharmacol. Sci.*, **26**, 15–19.
- The UniProt Consortium (2007) The Universal Protein Resource (UniProt). *Nucleic Acids Research*, **35**, D193–D197.
- Maglott, D., Ostell, J., Pruitt, K.D. and Tatusova, T. (2007) Entrez Gene: gene-centered information at NCBI. *Nucleic Acids Res.*, **35**, D26–D31.
- Kanehisa, M., Goto, S., Kawashima, S., Okuno, Y. and Hattori, M. (2004) The KEGG resource for deciphering the genome. *Nucleic Acids Res.*, **32**, D277–D280.
- Hattori, M., Okuno, Y., Goto, S. and Kanehisa, M. (2003) Development of a Chemical Structure Comparison Method for Integrated Analysis of Chemical and Genomic Information in the Metabolic Pathways. *J. Am. Chem. Soc.*, **125**, 11853–11865.
- Kotera, M., Okuno, Y., Hattori, M., Goto, S. and Kanehisa, M. (2004) Computational assignment of the EC numbers for genomic-scale analysis of enzymatic reactions. *J. Am. Chem. Soc.*, **126**, 16487–16498.