

The Non-Negative Matrix Factorization Toolbox for Biological Data Mining

Yifeng Li^{*1} and Alioune Ngom¹

¹School of Computer Science, University of Windsor, Windsor, Ontario, Canada

Email: Yifeng Li* - li11112c@uwindsor.ca; Alioune Ngom - angom@cs.uwindsor.ca;

*Corresponding author

Abstract

Background: Non-negative matrix factorization (NMF) has been being an important approach to analyze biological data. Though there exists some packages implemented in R and other programming languages, they either only provide some optimization algorithms, or focus on a specific application field. There is no complete package for the bioinformatics community to perform various data mining on biological data.

Results: We provide a powerful but convenient MATLAB toolbox including both the implementations of various NMFs and a variety of NMF-based data mining approaches for analyzing biological data. Through this toolbox, data mining approaches such as clustering, biclustering, feature extraction, feature selection, classification, missing values imputation, visualization, and statistical comparison can be easily done.

Conclusions: A series of analysis such as molecular pattern discovery, biological process identification, dimension reduction, disease prediction, visualization, and statistical comparison can be conducted via this toolbox.

Keywords: non-negative matrix factorization, clustering, biclustering, feature extraction, feature selection, classification, missing values.

Background

Non-negative matrix factorization (NMF) is a matrix decomposition approach that factorizes a non-negative matrix into two low-rank non-negative matrices [1]. It has made great success in biological data mining. The following is some well-known examples. Ref. [2] and [3] used NMF as clustering method to discover metagenes and molecular patterns. Ref. [4] applied *non-smooth NMF* (nsNMF) for biclustering of gene expression data. *Least-squares NMF* (LS-NMF) was proposed to take into account of the uncertainty information in gene expression data [5]. Ref. [6] proposed kernel NMF for the dimension reduction of gene expression data.

Many authors indeed provide their implementation along with their publications. It is convenient for the

interested readers to verify and apply them to their specific fields. However, there exists at least three issues that prevent NMF from a better application in biological and medical studies. First, their codes are implemented in diverse programming languages, such as R, MATLAB, C++, and Java, and usually only one optimization algorithm is provided in a paper. It is inconvenient for some researchers who want to choose a suitable NMF for their data among many different implementations which are realized in different languages and termination criteria. Second, some papers only provide optimization algorithms at a basic level, rather than a data mining implementation at a higher level. If a biologist wants to do clustering, biclustering, and then visualizing the results, it might be hard for him or her to investigate these three

aspect of data mining and visualization. Third, the existing implementations often stick to a specific application. There is no systematic package for NMF-based biological data mining.

There are some existing NMF toolboxes which are introduced in the following. However, neither address the above three issues altogether. *NMFLAB* [7] is MATLAB toolbox for signal and image processing. *NMFLAB* provides a user-friendly interface to load data, processing signals, and save result. It includes a variety of algorithms such as multiplicative rules, exponentiated gradient, projected gradient, conjugate gradient, and Quasi-Newton. It also can visualize the signals and components. But it does not provide any data mining functionality, which is the second issue mentioned above. It is also lack of some important NMF concepts such as semi-NMF and kernel NMF. *NMF:DTU Toolbox* [8] is also a MATLAB toolbox including five NMF optimization algorithms such as multiplicative rules, projected gradient, probabilistic NMF, alternating least squares, and alternating least squares with *optimal brain surgeon* (OBS) method. It does not address the second and third issues. *NMFN: Non-negative Matrix Factorization* [9] is a R package that implemented several algorithms. It has the same issue as the DTU toolbox. *NMF: Algorithms and framework for Nonnegative Matrix Factorization* [10] is a R package. It implemented several algorithms and their main interface allows parallel computations. It does not solve the second issue. *Text to Matrix Generator (TMG)* is a MATLAB toolbox for text mining. It does not solve the third issue. Ref. [11] provided a NMF plug-in for BRB-ArrayTools. This plug-in only implemented the standard NMF and semi-NMF for clustering gene expression profiles. Therefore it does not address the second and third issues well. *Coordinated Gene Activity in Pattern Sets* (CoGAPS) is a new package implemented by C++ with R interface. *Bayesian decomposition* (DB, same as NMF) is implemented and statistical methods are provided for the inference of biological processes by CoGAPS. It has been reported that CoGAPS can obtain more precise result than some other NMF methods [12]. CoGAPS uses a Markov chain Monte Carlo (MCMC) scheme for estimating the DB model parameters, which may be slow than the NMFs implemented by the block-coordinate gradient descent scheme.

In order to address the three issues, we propose a NMF toolbox in MATLAB in this paper. The toolbox is implemented in two levels. The basic level is composed of the algorithms of NMF variants, and the advanced level consists of diverse applications of NMF in biological data mining. The contributions of our toolbox

are enumerated in the following:

1. The NMF algorithms are relatively complete, and all of them are implemented in MATLAB. This is to address the first issue. It is impossible and unnecessary to implement all NMF algorithms. We focus on the well-known NMF representatives. This repository of NMFs allows users to select the most suitable one in specific scenarios.
2. It includes comprehensive functionalities of data mining, such as clustering, biclustering, feature extraction, feature selection, and classification. This solves the second and third issues to a great extend.
3. It provides functions of biological data visualization. For instance, heat maps of a NMF can be plotted via our toolbox. It is pretty helpful to interpret result. It also provides statistical methods to compare the performance of multiple methods.

The rest of this paper is organized as below. The implementations of the basis level are given in the next section. After that, examples of applying the advanced level in biological data mining are demonstrated. The conclusion and future works are given finally.

Implementation

As mentioned above, this toolbox is implemented at two levels. The fundamental level is composed of various NMF algorithms. And the advanced level includes many data mining approaches based on the fundamental level. The critical issues of implementing these NMF variants are addressed in this section. Table 1 summarized all the NMF algorithms implemented in our toolbox. Potential users need to use the command `help nmfrule` in the command line, for example, to learn how to set the parameters of a function in our toolbox. The advanced level and some examples are given in the next section.

Standard-NMF

The *standard-NMF* decomposes a non-negative matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ into two non-negative factors $\mathbf{A} \in \mathbb{R}^{m \times k}$ and $\mathbf{Y} \in \mathbb{R}^{k \times n}$ (where $k < \min\{m, n\}$), that is

$$\mathbf{X}_+ = \mathbf{A}_+ \mathbf{Y}_+ + \mathbf{E}, \quad (1)$$

where \mathbf{E} is error (or residual), \mathbf{M}_+ indicates the matrix \mathbf{M} is non-negative. Its optimization in the Euclidean dis-

tance is formulated in the following equation

$$\min_{\mathbf{A}, \mathbf{Y}} \frac{1}{2} \|\mathbf{X} - \mathbf{A}\mathbf{Y}\|_F^2, \text{ s.t. } \mathbf{A}, \mathbf{Y} \geq 0. \quad (2)$$

Statistically speaking, this formulation is obtained from the log-likelihood function under the assumption of Gaussian error. If multivariate data points are arranged in the columns of \mathbf{X} , then \mathbf{A} is called *basis matrix*, each column of \mathbf{A} is thus called a *basis vector*, and \mathbf{Y} is called *coefficient matrix*. The interpretation is that each data point is a (sparse) non-negative linear combination of the basis vectors. It is well-known that the objective is non-convex. Block-coordinate descent is therefore the main prescription of such problem. Multiplicative update rules were provided in [13] to optimize Equation 2. Though simple to implement, this algorithm is even not guaranteed to converge to a stationary point [14]. Essentially the optimizations above with respect to \mathbf{A} and \mathbf{Y} , respectively, are *non-negative least squares* (NNLS). Therefore we implemented the alternating NNLS algorithm proposed in [14]. It can be proved that this algorithm converges to a stationary point. In our toolbox, functions `nmfrule` and `nmfnls` are the implementations of the two algorithms above.

Semi-NMF

The standard NMF only works for non-negative data, which limits its applications. Ref. [15] extended it to *semi-NMF* which removes the non-negative constraints on the data \mathbf{X} and basis matrix \mathbf{A} . It can be expressed in the following equation:

$$\min_{\mathbf{A}, \mathbf{Y}} \frac{1}{2} \|\mathbf{X} - \mathbf{A}\mathbf{Y}\|_F^2, \text{ s.t. } \mathbf{Y} \geq 0. \quad (3)$$

Semi-NMF can be applied to the matrix of mixed signs, therefore it expands NMF to many fields. However, the gradient-descent-based update rule proposed in [15] is slow to converge (implemented in function `seminmfrule` in our toolbox). Keeping \mathbf{Y} fixed, updating \mathbf{A} is a least squares problem which has an analytical solution

$$\mathbf{A} = \mathbf{X}\mathbf{Y}^T(\mathbf{Y}\mathbf{Y}^T)^{-1} = \mathbf{X}\mathbf{Y}^\dagger, \quad (4)$$

where $\mathbf{Y}^\dagger = \mathbf{Y}^T(\mathbf{Y}\mathbf{Y}^T)^{-1}$ is Moore-Penrose pseudoinverse. Updating \mathbf{Y} while fixing \mathbf{A} is a NNLS problem essentially as above. Therefore we implemented the fast NNLS based algorithm to optimize semi-NMF in function `seminmfnls`.

Sparse-NMF

The standard NMF and semi-NMF have the issues of scale-invariance and non-unique solution, which imply that the non-negativity constrained on the least squares is insufficient in some cases. Sparsity is a popular regularization principle in statistical modeling [16]. It has been applied to reduce non-uniqueness and enhance interpretation of NMF. The *sparse-NMF* proposed in [3] is expressed in the following equation

$$\min_{\mathbf{A}, \mathbf{Y}} \frac{1}{2} \|\mathbf{X} - \mathbf{A}\mathbf{Y}\|_F^2 + \frac{\eta}{2} \|\mathbf{A}\|_F^2 + \frac{\lambda}{2} \sum_{i=1}^n \|\mathbf{y}_i\|_1 \quad (5)$$

s.t. $\mathbf{A}, \mathbf{Y} \geq 0$,

where \mathbf{y}_i is the i -th column of \mathbf{Y} . From the Bayesian viewpoint, This formulation is obtained from the log-posterior probability under the assumptions that Gaussian error, Gaussian distributed basis vector, and Laplace distributed coefficient. Keeping one matrix fixed, updating another matrix can be formulated into a NNLS problem. In order to improve the interpretation of basis vectors and speed up the algorithm, we implemented the following model instead:

$$\min_{\mathbf{A}, \mathbf{Y}} \frac{1}{2} \|\mathbf{X} - \mathbf{A}\mathbf{Y}\|_F^2 + \lambda \sum_{i=1}^n \|\mathbf{y}_i\|_1 \quad (6)$$

s.t. $\mathbf{A}, \mathbf{Y} \geq 0$,

$\|\mathbf{a}_i\|_2^2 = 1, \quad i = 1, \dots, k.$

We optimize this using three alternating steps in each iteration. First of all, the following task is optimized:

$$\min_{\mathbf{Y}} \frac{1}{2} \|\mathbf{X} - \mathbf{A}\mathbf{Y}\|_F^2 + \lambda \sum_{i=1}^n \|\mathbf{y}_i\|_1 \quad (7)$$

s.t. $\mathbf{Y} \geq 0$.

Second, \mathbf{A} is updated in the following equation:

$$\min_{\mathbf{A}} \frac{1}{2} \|\mathbf{X} - \mathbf{A}\mathbf{Y}\|_F^2 \quad (8)$$

s.t. $\mathbf{A} \geq 0$.

After that, the columns of \mathbf{A} are normalized to have unit l_2 norm. The first and second step can be solved by *non-negative quadratic programming* (NNQP), whose general formulation is

$$\min_{\mathbf{Z}} \sum_{i=1}^n \frac{1}{2} \mathbf{z}_i^T \mathbf{H} \mathbf{z}_i + \mathbf{g}_i^T \mathbf{z}_i + c_i \quad (9)$$

s.t. $\mathbf{Z} \geq 0$,

where z_i is the i -th column of matrix variable \mathbf{Z} . It is easy to prove that NNLS is a specific problem of NNQP. For example, Equation 7 can be rewritten as

$$\begin{aligned} \min_{\mathbf{Y}} \sum_{i=1}^n \frac{1}{2} \mathbf{y}_i^T (\mathbf{A}^T \mathbf{A}) \mathbf{y}_i + (\lambda - \mathbf{A}^T \mathbf{x}_i)^T \mathbf{y}_i + \mathbf{x}_i^T \mathbf{x}_i \\ \text{s.t. } \mathbf{Y} \geq 0. \end{aligned} \quad (10)$$

The implementations of the method in [3] and our method are given in functions `sparsenmfnnls` and `sparseNMFNNQP`, respectively. We also implemented the sparse semi-NMF in function `sparseseminmfnnls`.

Versatile Sparse Matrix Factorization

When training data \mathbf{X} is of mixed signs, \mathbf{A} should not be restrained by non-negativity. However, \mathbf{A} is not sparse any more. In order to obtain sparse basis vectors \mathbf{A} for some analysis, we may use l_1 norm on them to induce sparsity. The drawback of l_1 norm is that correlated variables may not be simultaneously non-zero in the l_1 induced sparse result. This is because l_1 norm is able to produce sparse but non-smooth result. It is known that l_2 norm is able to obtain smooth but not sparse result. Combining both norms has been proven that correlated variables can be selected or removed simultaneously [17]. When smoothness is needed on \mathbf{Y} , we can also enforce l_2 norm on it. We thus generalize the aforementioned NMF models into a versatile form as expressed below

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{Y}} f(\mathbf{A}, \mathbf{Y}) = \frac{1}{2} \|\mathbf{X} - \mathbf{A}\mathbf{Y}\|_F^2 + \sum_{i=1}^k \left(\frac{\alpha_2}{2} \|\mathbf{a}_i\|_2^2 \right. \\ \left. + \alpha_1 \|\mathbf{a}_i\|_1 \right) + \sum_{i=1}^n \left(\frac{\lambda_2}{2} \|\mathbf{y}_i\|_2^2 + \lambda_1 \|\mathbf{y}_i\|_1 \right) \quad (11) \\ \text{s.t. } \begin{cases} \text{if } t_1 = 1 & \mathbf{A} \geq 0 \\ \text{if } t_2 = 1 & \mathbf{Y} \geq 0 \end{cases}, \end{aligned}$$

where parameter $\alpha_1 \geq 0$ controls the sparsity of the basis vectors, parameter $\alpha_2 \geq 0$ controls the smoothness and scale of the basis vectors, parameter $\lambda_1 \geq 0$ controls the sparsity of the coefficient vectors, parameter $\lambda_2 \geq 0$ controls the smoothness of coefficient vectors, parameters t_1 and t_2 are boolean variables (0: false, 1: true) that indicate if non-negativity needs to be enforced on \mathbf{A} and \mathbf{Y} , respectively. We can call this model *versatile sparse matrix factorization* (VSMF). It can be easily seen

that the standard NMF, semi-NMF, and sparse-NMFs are special cases of VSMF.

We design the following multiplicative update rules for VSMF model in the case of $t_1 = t_2 = 1$ (implemented in function `sparsenmf2rule`):

$$\begin{cases} \mathbf{A} = \mathbf{A} * \frac{\mathbf{X}\mathbf{Y}^T}{\mathbf{A}\mathbf{Y}\mathbf{Y}^T + \alpha_2 \mathbf{A} + \alpha_1} \\ \mathbf{Y} = \mathbf{Y} * \frac{\mathbf{A}^T \mathbf{X}}{\mathbf{A}^T \mathbf{A} \mathbf{Y} + \lambda_2 \mathbf{Y} + \lambda_1} \end{cases}, \quad (12)$$

where $\mathbf{A} * \mathbf{B}$ and $\frac{\mathbf{A}}{\mathbf{B}}$ are element-wise multiplication and division between matrix \mathbf{A} and \mathbf{B} , respectively. Alternatively, we devise active-set algorithm for VSMF (implemented in function `vsmf`). When t_1 (or t_2) = 1, \mathbf{A} (or \mathbf{Y}) can be updated by NNQP (this case is also implemented in `sparsenmf2nnqp`). When t_1 (or t_2) = 0, \mathbf{A} (or \mathbf{Y}) can be updated by l_1 QP.

Kernel-NMF

Two features of a kernel approach are that i) it can represent complex patterns, and ii) the optimization of the model is dimension-free. We now show that NMF can also be kernelized. The difficulty is that the basis matrix is dimension-related, and therefore it is impossible to represent it in a very high (even infinite) dimensional fashion. We notice that the NNLS optimization of updating \mathbf{Y} in Equation 10 only needs inner products $\mathbf{A}^T \mathbf{A}$, $\mathbf{A}^T \mathbf{X}$, and $\mathbf{X}^T \mathbf{X}$. From Equation 4, we obtain that $\mathbf{A}^T \mathbf{A} = (\mathbf{Y}^\dagger)^T \mathbf{X}^T \mathbf{X} \mathbf{Y}^\dagger$, $\mathbf{A}^T \mathbf{X} = (\mathbf{Y}^\dagger)^T \mathbf{X}^T \mathbf{X}$. Therefore, we can see that only inner product $\mathbf{X}^T \mathbf{X}$ is needed in the optimization of NMF! It is hence very convenient to obtain *kernel-NMF* by replacing the inner product $\mathbf{X}^T \mathbf{X}$ by a kernel matrix $K(\mathbf{X}, \mathbf{X})$. For further detail, interested readers are referred to our recent paper [6]. Based on the above derivation, we implemented the kernel semi-NMF using multiplicative update rule (in `kernelseminmfrule`) and using NNLS (in `kernelseminmfnnls`). The sparse kernel semi-NMFs are implemented in functions `kernelsparseseminmfnnls` and `kernelSparseNMFNNQP` which are equivalent to each other. The kernel method of decomposing kernel matrix proposed in [18] is implemented in `kernelnmfdecom`.

Other Variants

Ref. [15] proposed *Convex-NMF*, where the columns of \mathbf{A} are constrained to be the convex combinations of data points in \mathbf{X} . It is thus formulated as $\mathbf{X}_\pm = \mathbf{X}_\pm \mathbf{W}_+ \mathbf{Y}_+ + \mathbf{E}$, where \mathbf{M}_\pm indicates that matrix \mathbf{M}

is of mixed signs. $\mathbf{XW} = \mathbf{A}$ and each column of \mathbf{W} contains the convex coefficients of all the data points to get the corresponding column of \mathbf{A} . It has been verified that columns of \mathbf{A} obtained by convex-NMF are close to the real cluster centroids. Convex-NMF can be kernelized as well [15]. We implemented convex-NMF and its kernel version in `convexnmfrule` and `kernelconvexnmf`, respectively.

The basis vectors obtained by all the above NMF are non-orthogonal. Alternatively, *orthogonal NMF* (ortho-NMF) imposes orthogonality to enhance sparsity [19]. Its formulation is

$$\begin{aligned} \mathbf{X} &= \mathbf{A}\mathbf{S}\mathbf{Y} + \mathbf{E} \\ \text{s.t. } \mathbf{A}^T\mathbf{A} &= \mathbf{I}, \quad \mathbf{Y}\mathbf{Y}^T = \mathbf{I}, \quad \mathbf{A}, \mathbf{S}, \mathbf{Y} \geq 0, \end{aligned} \quad (13)$$

where the input \mathbf{X} is non-negative, \mathbf{S} absorbs the magnitude due to the normalization of \mathbf{A} and \mathbf{Y} . Function `orthnmfrule` is its implementation in our toolbox. Ortho-NMF is very similar with the *non-negative sparse PCA* (NSPCA) proposed in [20]. The disjoint property on ortho-NMF may be too restrictive for many applications, therefore this property is relaxed in NSPCA. Ortho-NMF does not guarantee the maximum-variance property which is also relaxed in NSPCA. However NSPCA only enforces non-negativity on the basis vectors, even when the training data have negative values. As plan to design a model where the disjoint property, maximum-variance property, non-negativity, and sparsity can be nicely controlled on both basis vectors and coefficient vectors.

If we apply NMF on data subject to missing values, there are two fast ways. First, the missing values can be estimated prior to running NMF. Alternatively, *weighted-NMF* [21] can be directly applied to decompose it. Weighted-NMF puts a zero weight on the missing elements and hence only the present data contribute to the final result. An expectation-maximization (EM) based missing value estimation during the running of NMF might not a wise choice due to computational concern. The weighted-NMF is given in our toolbox in function `wnmfrule`.

Results and discussion

Based on the various NMFs implemented, a series of data mining can be conducted via our toolbox. Table 2 lists the functions in this level. These applications are described along with examples as below.

Clustering and biclustering

NMF has been applied for clustering. Given data \mathbf{X} with multivariate data points in the columns, the idea is that, after applying NMF on \mathbf{X} , a multivariate data point, say \mathbf{x}_i is a non-negative linear combination of the columns of \mathbf{A} , that is $\mathbf{x}_i \approx \mathbf{A}\mathbf{y}_i = y_{1i}\mathbf{a}_1 + \cdots + y_{ki}\mathbf{a}_k$. The greatest coefficient in the i th column of \mathbf{Y} indicates the cluster this data point belongs to. The reason is that if the data points are mainly composed by the same basis vector, they should be in the same group. A basis vector is usually viewed as a cluster centroid or prototype. It has been applied for clustering microarray data for the discovery of tumor subtypes in [2]. We implemented function `NMFCluster`, through which various NMF algorithms can be chosen. An example is provided in `exampleCluster` file in the folder of our toolbox.

If we interpret both of the basis matrix and coefficient matrix, NMF can be applied as biclustering approach. Interested readers are referred to [22] for a survey of biclustering and to [4] for a biclustering method by NMF. We implemented a biclustering approach based on NMF in `biCluster` function. The biclusters can be visualized via function `NMFBicHeatMap`. We applied NMF to simultaneously grouping the genes and samples of a leukemia dataset [2] which includes tumor samples of three subtypes. The goal is to find strongly correlated genes over a subset of samples. A subset of such genes and a subset of such samples form a bicluster. The heat map is shown in Figure 1. Readers can find the script in `exampleBiCluster` file of our toolbox.

Basis Vector Analysis for Biological Process Discovery

we can obtain more interpretation via the basis vectors. When applying NMF on microarray data, the basis vectors are interpreted as potential biological processes [23] [3] [12]. In the following, we give one example for finding biological factors on gene-sample data, and two examples on time-series data. Please note they only serve as simple examples. Fine tuning of the parameters of NMF is necessary for accurate results.

First, we ran our VSMF on the ALLAML gene-sample data [2] with the setting $k = 3$, $\alpha_1 = 0.01$, $\alpha_2 = 0.01$, $\lambda_1 = 0$, $\lambda_2 = 0.01$, $t_1 = 1$, and $t_2 = 1$. After that, we obtain 81, 37, and 448 genes for the three factors, respectively. As in [3], we then conducted gene enrichment analysis by applying Onto-Express [24]. Please see file `exampleBioProcessGS` for details. The gene set analysis can also be conducted via other tools such as MIPS [25], GOTermFinder [26], and DAVID [27] [28].

Second, we used NMF to cluster a time-series data of yeast metabolic cycle in [29]. Figure 2 shows the heat map of NMF clustering, and Figure 3 illustrates the three basis vectors. We used `nmfnls` function to decompose the data and `NMFHeatMap` to plot the heat map. The detailed script is given in the `exampleBioProcessTSYeast` file in the toolbox. We can clearly see that the three periodical biological processes corresponds exactly to the Ox (oxidative), R/B (reductive, building), and R/C (reductive, charging) processes discovered in [29].

Third, we used NMF to factorize a breast cancer time-series data set which include wide type MYCN cell lines and mutant MYCN cell lines [30]. The purpose of example is to show that NMF is a potential tool to fining drivers of cancers. One basic methodology is in the following. First, basis vectors are produced applying NMF on a time-series data. Then factor-specific genes are identified by computational or statistical methods. Finally, the regulators of the factor-specific genes are identified by biological knowledge. This data set has 8 time points (0, 2, 4, 8, 12, 24, 36, 48 hr.). The zero time point is untreated. Samples were collected at its following time points after treatment of 4-hydroxytamoxifen (4-OHT). In our computational experiment, we use our VSMF implementation (function `vsmf`). we set $k = 2$. Because this data set has negative values we set $t_1 = 0$ and $t_2 = 1$. We set $\alpha_1 = 0.01$, $\alpha_2 = 0$, $\lambda_1 = 0$, $\lambda_2 = 0.01$. The basis vectors of both wide-type and mutant data are compared in Figure 4. From the wide-type time-series data, we can successfully identify two patterns. The rising pattern corresponds to the induced signature and the falling pattern corresponding to the repressed signature in [30]. It is reported, in [30], that both patterns are contributed by MYC target genes. From the mutant time-series, we can obtain two flat processes, which is plausible. The source code of this example can be found in `exampleBioProcessMYC`. Readers are also suggested to look at the NMF based method, proposed in [12] and [31], for identifying signaling pathways.

Basis Vector Analysis for Variable Selection

The columns in \mathbf{A} for the gene expression data is called *metasample* in [2]. They can be interpreted as biological processes, because its values imply the activation and silence of all the genes. Gene selection aims to find marker genes to aid the disease prediction and understand pathways. Rather than selecting genes on the original data, the novel idea is to conduct gene selec-

tion on the metasamples. The reason is that the discovered biological process via NMF are biologically meaningful for the class discrimination, and the genes, expressing differently across these processes, should contribute to the classification. In Figure 3, for example, three biological processes are discovered and only selected genes are shown. We have implemented the information-entropy-based gene selection approach proposed in [3] in function `featureFilterNMF`. We give an example on how to call this function in file `exampleFeatureSelection`. It has been reported that it can select meaningful genes, which has been verified by gene ontology analysis. Feature selection based on supervised NMF will also be implemented.

Feature Extraction

Microarray data and mass spectrometry data have tens of thousands of features but only tens or hundreds of samples. This leads to the notorious curses of dimensionality. For example, it is impossible to estimate the parameters of some statistical models as the number of their parameters grow exponentially as the dimension increases. Another issue is that biological data is normally subject to much noise, which crucially affects the analysis. In cancer study, a common hypothesis is that only several biological factors (for example oncogenes) play a crucial role in the development of a cancer. When we generate data from control and unhealthy patients, the huge amount of data (due to high dimensions) therefore contain lots of irrelevant and redundant information. Orthogonal factors obtained by *principal component analysis* (PCA) or *independent component analysis* (ICA) are not appropriate in many cases. Since NMF generates non-orthogonal (non-negative) factors, therefore it is much plausible to extract new features from such data using NMF. As mentioned above, training data $\mathbf{X}_{m \times n}$, with m features and n samples, can be decomposed into k metasamples $\mathbf{A}_{m \times k}$ and $\mathbf{Y}_{k \times n}$, that is

$$\mathbf{X} \approx \mathbf{A}\mathbf{Y}_{\text{tr}}, \text{ s.t. } \mathbf{A}, \mathbf{Y}_{\text{tr}} \geq 0, \quad (14)$$

where \mathbf{Y}_{tr} means the \mathbf{Y} is obtained from the training data. The k columns of \mathbf{A} span the k -dimensional *feature space* and each column of \mathbf{Y}_{tr} is the representation of the corresponding original training sample in the feature space. In order to project the p future unknown samples $\mathbf{S}_{m \times p}$ into this feature space, we have to solve the following non-negative least squares problem:

$$\mathbf{S} \approx \mathbf{A}\mathbf{Y}_{\text{uk}}, \text{ s.t. } \mathbf{Y}_{\text{uk}} \geq 0, \quad (15)$$

where Y_{uk} means the Y is obtained from the unknown samples. After obtaining Y_{tr} and Y_{uk} , the learning and prediction can be quickly done in the k -dimensional feature space instead of the m -dimensional original space. A classifier can learn over Y_{tr} , and then predict the class labels of the representations of unknown samples, that is Y_{uk} . From the aspect of interpretation, the advantage of NMF over PCA and ICA is that the metasamples are meaningful to understand the underlining physical processes, as mentioned above. We implemented a pair of functions `featureExtractionTrain` and `featureExtractionTest` including many linear and kernel NMF algorithms. The basis matrix (or inner product of basis matrix in the kernel case) is learned from training data via the former function, and unknown samples can be projected into the feature space via the later function. We give the examples of how to use these functions in files `exampleFeatureExtraction` and `exampleFeatureExtractionKernel`. Figure 5 shows the classification performance of no dimension reduction, linear NMF, kernel NMF using *radial basis function* (rbf) kernel, and PCA on two datasets, SRBCT [32] and Breast [33]. Since ICA is computationally very slow, we did not compare with it. The bars represent the averaged 4-fold cross-validation accuracies using *support vector machine* (SVM) as classifier over 20 runs. We can see that NMF is comparable with PCA on SRBCT, and is slightly better than PCA on Breast data. Also, with only few factors, the performance after dimension reduction using NMF is similar with, even better than, that without any dimension reduction. As future work, supervised NMF will be investigated and implemented in order to extract discriminative features.

Classification

If we make a reasonable assumption that every unknown sample is a sparse non-negative linear combination of the training samples, then we can directly derive a classifier from NMF. Indeed, this is a specific case of NMF where the training samples are the basis vectors. Since the optimization is actually NNLS problem, we dub this approach as *NNLS classifier* [34]. A NNLS problem is essentially quadratic programming problem as formulated in Equation 9, therefore, only inner products are needed for the optimization. We hence can naturally extend the NNLS classifier to kernel version. Two glaring features of this approach are that i) the sparsity regularization can avoid overfitting the model; and ii) kernelization allows dimension-free optimization and linearizes complex patterns. The implementation of NNLS classifier

is in file `nnlsClassifier`. Also, we provide many other classification approaches, for example SVM, in our toolbox. Please see file `exampleClassification` for demonstration. In our experiment of 4-fold cross-validation, accuracies of 0.7865 and 0.7804 are, respectively, obtained by linear and kernel (rbf kernel) NNLS classifier on Breast dataset. They achieved accuracies of 0.9762 and 0.9785, respectively, over SRBCT data.

Biological data are usually very noisy and sometimes suffer from the severe issue of missing values. Two key strengths of NNLS classifier are that it is very robust to noise and missing values, which make NNLS classifier very suitable for classifying biological data. In order to show its robustness to noise, we added Gaussian noise of mean 0 and variance from 0 to 4 by step 0.5 on SRBCT. Figure 6 illustrates the results of NNLS, SVM, and *1-nearest neighbor* classifiers using the noisy data. It can be seen that as the noise increases, NNLS outperforms SVM and 1-NN significantly. Biological data sometimes suffer from missing values. To deal with this problem, three strategies, removal, imputation, and disregard, are usually used. The first one will remove much other useful information, particularly when there is a large percent of missing values. The second one has a high risk of introducing false data. The third one is to avoid using missing values during analysis, which is the best way. Our idea in the following belongs to the last case. We base our idea of coping with missing values on the fact that NNLS optimization only needs inner products of pairs of samples. When computing the inner product of two samples, say x_i and x_j , we normalize them to have unit l_2 norm using only the features present in both samples, and then take the inner product. In order to impress the readers, we removed 10% to 70% of the data in STBCT. Using such data, we compared our method with 0-imputation method (the more sophisticated $k - NN$ method [35] would fail in high missing rate because no complete vector exists) in Figure 7. We can see that the NNLS classifier using our idea outperforms the imputation method in the case of large missing rate.

Statistical Comparison

In order to evaluate the performance of a method statistically, we need to conduct statistical tests. In this toolbox, we provide two methods for statistical evaluation. The first is a two-stage method proposed in [36]. The importance of this method is that it can estimate the data-size requirement for significant accuracy and extrapolate the performance based on current available data. Generating biological data is usually very expensive. This

method helps researchers to evaluate the necessity of producing more data. At the first stage, the minimum data size required for significant accuracy is estimated. This is implemented in function `significantAcc`. The second stage is to fit the learning curve using error rates of large data sizes. It is implemented in function `learnCurve`. In our experiment, we found that the NNLS classifier usually requires few number of samples for significant accuracy. For example on SRBCT data, NNLS requires only 4 training samples, while SVM needs 19 training samples for significant accuracy. The fitted learning curves of NNLS and SVM classifiers are shown in Figure 8. We provide an example of how to plot this figure in file `exampleFitLearnCurve`.

The second method is the nonparametric Friedman test coupled with post-hoc Nemenyi test to compare multiple classifiers over multiple datasets [37]. It is difficult to draw overall conclusion, if we compare multiple approaches in a pairwise fashion. Friedman test is recommended in [37] because it is simple, safe, and robust, compared with parametric tests. It is implemented in function `FriedmanTest`. The result can be presented graphically by CD diagram as implemented in function `plotNemenyiTest`. Figure 9 is an example of the result Nemenyi test which was conducted to compare 8 classifiers over 13 high dimensional biological datasets. This example can be found in file `exampleFriedmanTest`. In this figure, CD stands for crucial difference. If the distance of two methods is greater than CD, then we conclude that they differ significantly.

Conclusions

In order to address the issues of existing NMF implementations, we propose the powerful NMF MATLAB Toolbox which includes basic algorithmic level and advanced data mining level. It enable users to analyze biological data via NMF-based data mining approaches, such as clustering, biclustering, feature extraction, feature selection, and classification.

The following are the future works to improve the toolbox. First, we will include more NMF algorithms such as nsNMF, LS-NMF, and supervised NMF. Second, we are very interested in implementing and speeding up Bayesian decomposition which is actually a probabilistic NMF invented independently in the same period as the standard NMF. Third, we will implement more statistical comparison methods. Furthermore, we will investigate the performance of NMF on denoising and data compres-

sion.

Availability and requirements

Project name: The NMF Toolbox in MATLAB

Project home page: <https://sites.google.com/site/nmftool> and <http://cs.uwindsor.ca/~li11112c/nmf>

Operating system(s): Platform independent

Programming language: MATLAB

Other requirements: MATLAB 7.11 or higher

License: GNU GPL Version 3

Any restrictions to use by non-academics: Licence needed

Acknowledgements

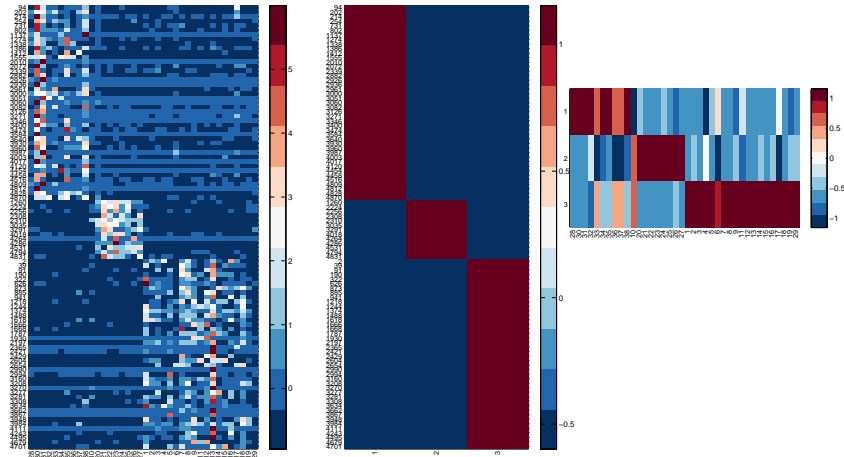
This research has been partially supported by IEEE CIS Walter Karplus Summer Research Grant 2010, Ontario Graduate Scholarship 2011-2013, and Canadian NSERC Grants #RGPIN228117-2011.

References

1. Lee DD, Seung S: **Learning the Parts of Objects by Non-Negative Matrix Factorization**. *Nature* 1999, **401**:788–791.
2. Brunet J, Tamayo P, Golub T, Mesirov J: **Metagenes and Molecular Pattern Discovery Using Matrix Factorization**. *PNAS* 2004, **101**(12):4164–4169.
3. Kim H, Park H: **Sparse Non-Negative Matrix Factorization via Alternating Non-Negativity-Constrained Least Squares for Microarray Data Analysis**. *SIAM J. Matrix Analysis and Applications* 2007, **23**(12):1495–1502.
4. Carmona-Saez P, Pascual-Marqui RD, Tirado F, Carazo JM, Pascual-Montano A: **Biclustering of Gene Expression Data by Non-Smooth Non-Negative Matrix Factorization**. *BMC Bioinformatics* 2006, **7**:78.
5. Wang G, Kossenkova A, Ochs M: **LS-NMF: A Modified Non-Negative Matrix Factorization Algorithm Utilizing Uncertainty Estimates**. *BMC Bioinformatics* 2006, **7**:175.
6. Li Y, Ngom A: **A New Kernel Non-Negative Matrix Factorization and Its Application in Microarray Data Analysis**. In *CIBCB*, IEEE CIS Society, Piscataway: IEEE Press 2012:371–378.
7. Cichocki A, Zdunek R: **NMFLAB - MATLAB Toolbox for Non-Negative Matrix Factorization**. Tech. rep. 2006, [http://www.bsp.brain.riken.jp/ICALAB/nmflab.html].
8. **The NMF: DTU Toolbox**. Tech. rep., Technical University of Denmark [http://www.bsp.brain.riken.jp/ICALAB/nmflab.html].
9. Liu S: **NMFN: Non-negative Matrix Factorization**. Tech. rep., Duke University 2011, [http://cran.r-project.org/web/packages/NMFN].
10. Gaujoux R, Seoighe C: **A Flexible R Package for Nonnegative Matrix Factorization**. *BMC Bioinformatics* 2010, **11**:367, [http://cran.r-project.org/web/packages/NMF].
11. Qi Q, Zhao Y, Li M, Simon R: **Non-Negative Matrix Factorization of Gene Expression Profiles: A Plug-in for BRB-ArrayTools**. *Bioinformatics* 2009, **25**(4):545–547.
12. Ochs M, Fertig E: **Matrix Factorization for Transcriptional Regulatory Network Inference**. In *CIBCB*, IEEE CIS Society, Piscataway: IEEE Press 2012:387–396.
13. Lee D, Seung S: **Algorithms for Non-Negative Matrix Factorization**. In *Advances in Neural Information Processing Systems*, MIT Press 2001:556–562.
14. Kim H, Park H: **Nonnegative Matrix Factorization based on Alternating Nonnegativity Constrained Least Squares and Active Set Method**. *SIAM J. Matrix Analysis and Applications* 2008, **30**(2):713–730.
15. Ding C, Li T, Jordan MI: **Convex and Semi-Nonnegative Matrix Factorizations**. *TPAMI* 2010, **32**:45–55.
16. Tibshirani R: **Regression Shrinkage and Selection via The Lasso**. *Journal of the Royal Statistical Society. Series B (Methodological)* 1996, **58**:267–288.
17. Zou H, Hastie T: **Regularization and Variable Selection via The Elastic Net**. *Journal of the Royal Statistical Society - Series B: Statistical Methodology* 2005, **67**(2):301–320.
18. Zhang D, Zhou Z, Chen S: **Non-Negative Matrix Factorization on Kernels**. *LNCS* 2006, **4099**:404–412.
19. Ding C, Li T, Peng W, Park H: **Orthogonal Nonnegative Matrix Tri-Factorizations for Clustering**. In *KDD*, ACM, New York: ACM 2006:126–135.
20. Zass R, Shashua A: **Non-Negative Sparse PCA**. In *NIPS*, NIPS, MIT Press 2006.
21. Ho N: **Nonnegative Matrix Factorization Algorithms and Applications**. *PhD thesis*, Louvain-la-Neuve, Belgium 2008.
22. Madeira S, Oliveira A: **Biclustering Algorithms for Biological Data Analysis: A Survey**. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 2004, **1**:24–45.
23. Kim P, Tidor B: **Subsystem Identification Through Dimensionality Reduction of Large-Scale Gene Expression Data**. *Genome Research* 2003, **13**:1706–1718.
24. Draghici S, Khatri P, Bhavsar P, Shah A, Krawetz S, Tainsky M: **Onto-Tools, The Toolkit of The Modern Biologist: Onto-Express, Onto-Compare, Onto-Design and Onto-Translate**. *Nucleic Acids Res.* 2003, **31**(13):3775–3781.
25. Mewes H, Frishman D, Gruber C, Geier B, Haase D, Kaps A, Lemcke K, Mannhaupt G, Pfeiffer F, Schuller C, Stocker S, Weil B: **MIPS: A Database for Genomes and Protein Sequences**. *Nucleic Acids Res.* 2000, **28**:37–40.
26. Boyle E, Weng S, Gollub J, Jin H, Botstein D, Cherry J, Sherlock G: **GO::TermFinder – Open Source Software for Accessing Gene Ontology Information and Finding Significantly Enriched Gene Ontology Terms Associated with a List of Genes**. *Bioinformatics* 2004, **20**:3710–3715.
27. Huang D, Sherman B, Lempicki R: **Systematic and Integrative Analysis of Large Gene Lists Using DAVID Bioinformatics Resources**. *Nature Protoc.* 2009, **4**:44–57.
28. Huang D, Sherman B, Lempicki R: **Bioinformatics Enrichment Tools: Paths Toward the Comprehensive Functional Analysis of Large Gene Lists**. *Nucleic Acids Res.* 2009, **37**:1–13.
29. Tu B, Kudlicki A, Rowicka M, McKnight S: **Logic of the Yeast Metabolic Cycle: Temporal Compartmentalization of Cellular Processes**. *Science* 2005, **310**:1152–1158.
30. Chandriani S, Frengen E, Cowling V, Pendergrass S, Perou C, Whitfield M, Cole M: **A Core MYC Gene Expression Signature is Prominent in Basal-Like Breast Cancer but only Partially Overlaps The Core Serum Response**. *PLoS ONE* 2009, **4**(5):e6693.
31. Ochs M, Rink L, Tam C, Mburu S, Taguchi T, Eisenberg B, Godwin A: **Detection of Treatment-Induced Changes in Signaling Pathways in Gastrointestinal Stromal Tumors Using Transcriptomic Data**. *Cancer Res.* 2009, **69**(23):9125–9132.
32. Khan J: **Classification and Diagnostic Prediction of Cancers Using Gene Expression Profiling and Artificial Neural Networks**. *Nature Medicine* 2001, **7**(6):673–679.
33. Hu Z: **The Molecular Portraits of Breast Tumors are Conserved Across Microarray Platforms**. *BMC Genomics* 2006, **7**:96.
34. Li Y, Ngom A: **Classification approach based on non-negative least squares**. *Neurocomputing* 2013, **in press**.
35. Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R, Botstein D, Altman R: **Missing Value Estimation Methods for DNA Microarrays**. *Bioinformatics* 2001, **17**(6):520–525.
36. Mukherjee S, Tamayo P, Rogers S, Rifkin R, Engle A, Campbell C, Golub T, Mesirov J: **Estimating Dataset Size Requirements for Classifying DNA Microarray Data**. *Journal of Computational Biology* 2003, **10**(2):119–142.
37. Demsar J: **Statistical Comparisons of Classifiers over Multiple Data Sets**. *Journal of Machine Learning Research* 2006, **7**:1–30.

Figure 1 - Heat map of NMF biclustering

The left is the gene expression data where each column corresponds to a sample, the middle is the basis matrix, and the right is the coefficient matrix.



Figures

Figure 2 - Heat map of NMF clustering on a yeast metabolic cycle time-series data

The left is the gene expression time-series data where each column corresponds to a gene, the middle is the basis matrix, and the right is the coefficient matrix.

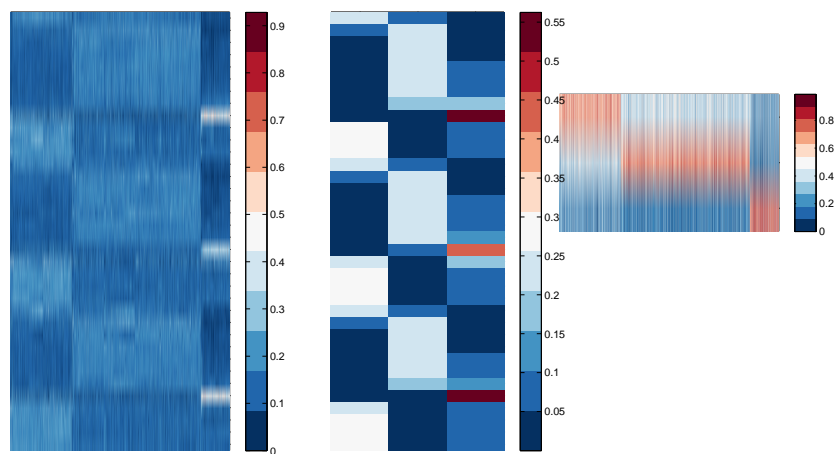


Figure 3 - Biological processes discovered by NMF on a yeast metabolic cycle time-series data

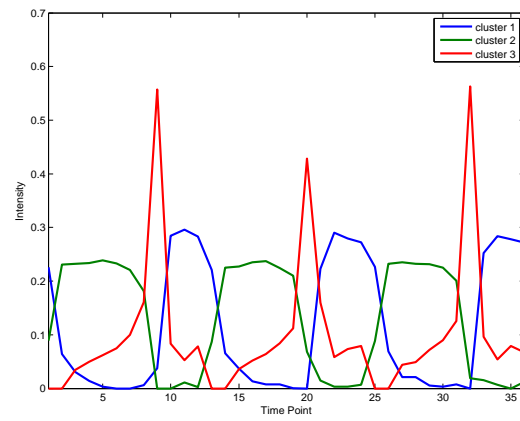


Figure 4 - Biological processes discovered by NMF on breast cancer time-series data

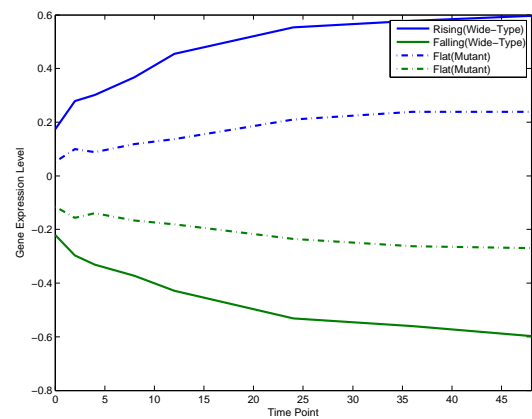


Figure 4 - The mean accuracy and standard deviation of NMF-based feature extraction on SRBCT data

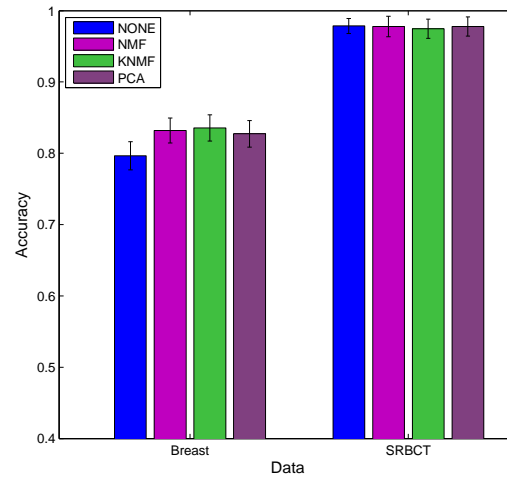


Figure 5 - The mean accuracy of NNLS classifier for different missing rate on SRBCT data

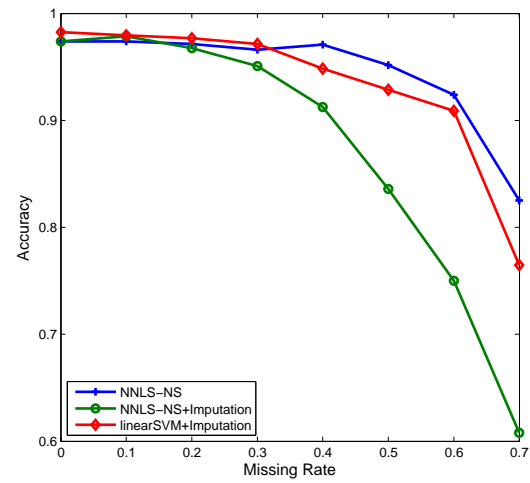


Figure 6 - The mean accuracy of NNLS classifier for different amount of noise on SRBCT data

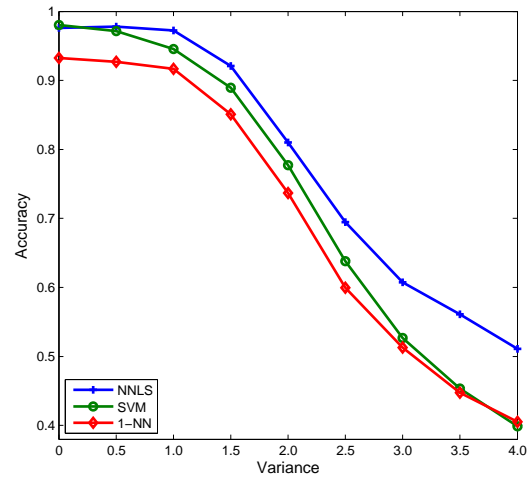


Figure 7 - The fitted Learning curves of NNLS and SVM classifiers on SRBCT data

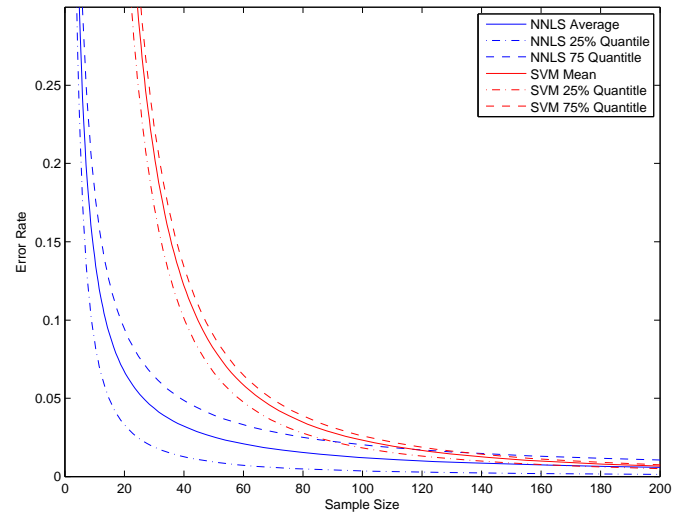
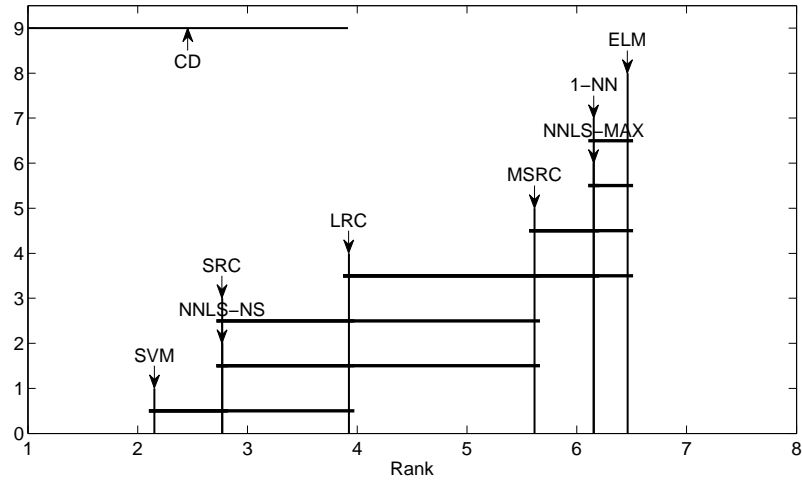


Figure 8 - Graphical representation of Nemenyi test for comparing 8 classifiers over 13 high dimensional biological data



Tables

Table 1 - Algorithms of NMF variants

Function	Description
nmfrule	The standard NMF optimized by gradient-descent-based multiplicative rules.
nmfnls	The standard NMF optimized by NNLS active-set algorithm.
seminmfrule	Semi-NMF optimized by multiplicative rules.
seminmfnls	Semi-NMF optimized by NNLS.
sparsenmfnls	Sparse-NMF optimized by NNLS.
sparsenmfNNQP	Sparse-NMF optimized by NNQP.
sparseseminmfnls	Sparse semi-NMF optimized by NNLS.
kernelnmfdecom	Kernel NMF through decomposing the kernel matrix of input data.
kernelseminmfrule	Kernel semi-NMF optimized by multiplicative rule.
kernelseminmfnls	Kernel semi-NMF optimized by NNLS.
kernelsparseseminmfnls	Kernel sparse semi-NMF optimized by NNLS.
kernelSparseNMFNNQP	Kernel sparse semi-NMF optimized by NNQP.
convexnmfrule	Convex-NMF optimized by multiplicative rules.
kernelconvexnmf	Kernel convex-NMF optimized by multiplicative rules.
orthnmfrule	Orth-NMF optimized by multiplicative rules.
wnmfrule	Weighted-NMF optimized by multiplicative rules.
sparsenmf2rule	Sparse-NMF on both factors optimized by multiplicative rules.
sparsenmf2nnqp	Sparse-NMF on both factors optimized by NNQP.
vsmf	Versatile sparse matrix factorization optimized by NNQP and l_1 QP.
nmf	The omnibus of the above algorithms.
computeKernelMatrix	Compute the kernel matrix $k(A,B)$ given a kernel function.

Table 2 - NMF-based data mining approaches

Function	Description
NMFCluster	Take the coefficient matrix produced by a NMF algorithm, and output the clustering result.
chooseBestk	Search the best number of clusters based on dispersion Coefficients.
biCluster	The biclustering method using one of the NMF algorithms.
featureExtractionTrain	General interface. Using training data, generate the bases of the NMF feature space.
featureExtractionTest	General interface. Map the test/unknown data into the feature space.
featureFilterNMF	On training data, select features by various NMFs.
featSel	feature selection methods.
nnlsClassifier	The NNLS classifier.
perform	Evaluate the classifier performance.
changeClassLabels01	Change the class labels to be in $\{0, 1, 2, \dots, C - 1\}$ for C -class problem.
gridSearchUniverse	a framework to do line or grid search.
classificationTrain	Train a classifier, many classifiers are included.
classificationPredict	Predict the class labels of unknown samples via the model learned by classificationTrain.
multiClassifiers	Run multiple classifiers on the same training data.
cvExperiment	Conduct experiment of k-fold cross-validation on a dataset.
significantAcc	Check if the given data size can obtain significant accuracy.
learnCurve	Fit the learning curve.
FriedmanTest	Friedman test with post-hoc Nemenyi test to compare multiple classifiers on multiple datasets.
plotNemenyiTest	Plot the CD diagram of Nemenyi test.
NMFHeatMap	Draw and save the heat maps of NMF clustering.
NMFBicHeatMap	Draw and save the heat maps of NMF biclustering.
plotBarError	Plot Bars with STD.
writeGeneList	write the gene list into a .txt file.
normmean0std1	Normalization to have mean 0 and STD 1.
sparsity	Calculate the sparsity of a matrix.
MAT2DAT	Write a dataset from MATLAB into .dat format in order to be readable by other languages.

Table 3 - Gene enrichment analysis using Onto-Express for the factor specific genes identified by NMF

Factor 1		Factor 2		Factor 3	
biological process	p-value	biological process	p-value	biological process	p-value
reproduction (5)	0	response to stimulus (15)	0.035	regulation of bio. proc. (226)	0.009
metabolic process (41)	0	biological regulation(14)	0.048	multi-organism proc. (39)	0.005
cellular process (58)	0			biological regulation (237)	0.026
death (5)	0				
developmental process (19)	0				
regulation of biological process (19)	0				