

**The Role of  
Singular Value Decomposition in  
Gene Expression Microarrays**

**Amir Niknejad**  
University of Illinois at Chicago

UVW, April 3, 2006

# 1 Outline of the talk

1. Estimation of missing values in given matrix data using the inverse eigenvalue problems techniques, and their applications to DNA microarrays and image processing.
2. A joint SVD decomposition of two or more matrices to compare several biological processes.

Most of the results can be found in the following recent papers, which are available at

<http://www.math.uic.edu/~friedlan/research.html>

## Singular Value Decomposition

$$G^{(n \times m)} = U^{(n \times n)} \Sigma^{(m \times m)} V^{(m \times n)T}, n \gg m$$

$$\Sigma = \text{DIAG}(\sigma_1, \sigma_2, \dots, \sigma_m)$$

$$\sigma_i = \sqrt{\lambda_i} \text{ eigenvalues of } GG^T \text{ or } G^T G$$

Columns of  $U$ , eigenvectors of  $GG^T$ .

Columns of  $V$ , eigenvectors of  $G^T G$ .

Let  $\text{rank}(G) = r$ . Then

- $\sigma_1 \geq \sigma_2 \geq \dots \sigma_r > \sigma_{r+1} = \dots \sigma_m = 0$
- $N(G) = \text{span}\{v_{r+1}, \dots, v_n\}$  and  
 $R(G) = \text{span}\{u_1, \dots, u_r\}$
- $G = \sum_{i=1}^r u_i v_i v_i^T$
- Norms:  $\|G\|_F^2 = \sigma_1^2 + \dots + \sigma_r^2$ , and  
 $\|G\|_2^2 = \sigma_1^2$

## Principal Component Analysis (PCA)

- Let  $G = (g_1, \dots, g_n)^T$ , genes through different experiments. Let  $\mu_G = E(G)$  be the mean and  $C_G = (G - \mu_G)(G - \mu_G)^T$  be the covariance matrix.
- Then compute the eigenvalues  $\lambda_i$  and corresponding eigenvector  $v_i, i = 1 \dots n$ .
- Sort the eigenvalues and corresponding eigenvectors in descending order.
- Construct orthogonal basis, by having the first eigenvector the direction of the largest variance.
- Represent the data with the few basis vector (corresponding to the largest eigen value).
- Let  $E_k$  be the matrix having the first  $k$  eigenvectors.
- Project the original data on the axis in dimension  $k$ . This way we could minimize MSE, by this representation of data. This way we lose a little bit of information.

## Low Rank Approximation

Let  $\mathcal{R}(n, m, k)$  denote the set of  $n \times m$  matrices of at most rank  $k$  ( $m \geq k$ ). Then for each  $k, k \leq r$ , the SVD of  $E$  gives the solution to the following approximation problem:

$$\begin{aligned} \min_{E \in \mathcal{R}(n, m, k)} \|G - E\|_{\mathcal{F}} &= \min_{E \in \mathcal{R}(n, m, k)} \left\| G - \sum_{q=1}^k \sigma_q \mathbf{u}_q \mathbf{v}_q^T \right\|_{\mathcal{F}} \\ &= \sqrt{\sum_{q=k+1}^r \sigma_q^2} \end{aligned}$$

Transcriptional Response:

$$\mathbf{g}_i = \sum_{k=1}^r u_{ik} \sigma_k \mathbf{v}_k$$

Array expression profile:

$$\mathbf{c}_j = \sum_{k=1}^r v_{jk} \sigma_k \mathbf{u}_k$$

## Computing Eigen-genes

Gene matrix  $G^{n \times m}$ , # of significant  $\sigma_i \leq \frac{m}{2}$

$u_1, \dots, u_n$  eigenarrays,  $v_1, \dots, v_n$  eigengenes,  
 $\sigma_1 \dots, \sigma_l$  eigenexpressions.

- $\text{Hous}(G) = (I - 2hh^t)G = A. \|h\|_2 = 1$

- Compute  $A^T A$

- $QR(A^T A) \rightarrow Q \begin{bmatrix} R \\ 0 \end{bmatrix} \rightarrow$

positive  $\sigma_1^2, \dots, \sigma_l^2$  and  $u_1, \dots, u_n$

## 2 SVD in inner product spaces

$U_i$  is  $m_i$ -dimensional IPS over  $\mathbb{C}$ , with  $\langle \cdot, \cdot \rangle_i, i = 1, 2$ .

$T : U_1 \rightarrow U_2$  linear operator.  $T^* : U_2 \rightarrow U_1$  the adjoint operator:  $\langle T\mathbf{x}, \mathbf{y} \rangle_2 = \langle \mathbf{x}, T^*\mathbf{y} \rangle_1$ .

$S_1 := T^*T : U_1 \rightarrow U_1$ ,

$S_2 := TT^* : U_2 \rightarrow U_2$ .

$S_1, S_2$  self-adjoint:  $S_1^* = S_1, S_2^* = S_2$  and nonnegative definite:  $\langle S_i \mathbf{x}_i, \mathbf{x}_i \rangle_i \geq 0$ .

$\sigma_1^2 \geq \dots \geq \sigma_r^2 > 0$  positive eigenvalues of  $S_1$  and  $S_2$  and  $r = \text{rank } T = \text{rank } T^*$ . Let

$S_1 \mathbf{v}_i = \sigma_i^2 \mathbf{v}_i, \langle \mathbf{v}_i, \mathbf{v}_j \rangle_1 = \delta_{ij}, i, j = 1, \dots, r$ .

Define  $\mathbf{u}_i := \sigma_i^{-1} T \mathbf{v}_i, i = 1, \dots, r$ . Then

$\langle \mathbf{u}_i, \mathbf{u}_j \rangle_2 = \delta_{ij}, i, j = 1, \dots, r$ .

Complete  $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$  and  $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$  to orthonormal bases  $[\mathbf{v}_1, \dots, \mathbf{v}_{m_1}]$  and  $[\mathbf{u}_1, \dots, \mathbf{u}_{m_2}]$  in  $U_1$  and  $U_2$ .

### 3 Matrix SVD

Let  $A \in \mathbb{C}^{m \times n}$ . Then  $A : \mathbb{C}^n \rightarrow \mathbb{C}^m$ . Assume  $\mathbb{C}^n, \mathbb{C}^m$  equipped with standard inner product  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{y}^* \mathbf{x}$ .

Then  $A = U \Sigma V^*$ , where  $U \in U(m), V \in U(n)$ ,  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min(m,n)}) \in \mathbb{R}_+^{m \times n}$ .

$U, V$  transition matrices from  $[\mathbf{u}_1, \dots, \mathbf{u}_m], [\mathbf{v}_1, \dots, \mathbf{v}_n]$  to the standard bases in  $\mathbb{C}^m, \mathbb{C}^n$  respectively.

For  $k \leq r$  let  $\Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{k \times k}$ , and  $U_k \in U(m, k), V_k \in U(n, k)$  having the first  $k$  columns of  $U, V$  respectively. Then  $A_k := U_k \Sigma_k V_k^*$  the best rank  $k$  approximation in Frobenius and operator norm of  $A$ :

$$\min_{B \in \mathcal{R}(m,n,k)} \|A - B\| = \|A - A_k\|.$$

$A = U_r \Sigma_r V_r^*$  is Reduced SVD

$(r \geq) \nu$  numerical rank of  $A$  if  $\frac{\sigma_{\nu+1}}{\sigma_\nu} \approx 0$ .

$A_\nu$  is a noise reduction of  $A$ .

Noise reduction has many applications in image processing, DNA-Microarrays analysis, data compression.



## 4 MISSING ENTRIES PROBLEM

In DNA Microarrays experiments one measure thousands of genes  $i = 1, \dots, m$  in  $n$  different conditions, typically  $n \in [3, 20]$ .

$0 \leq a_{ij}$  measures the intensity of gene  $i$  in  $j$  – th experiment. The results are recorded in the matrix

$$A = (a_{ij}) \in \mathbb{R}^{m \times n}.$$

Sometimes the entries  $a_{ij}$  are missing (corrupted, up to 20%).

Let  $\mathcal{T} \subset \{1, \dots, n\} \times \{1, \dots, m\}$  missing entries set.

Set  $a_{ij} = 0$  if  $(i, j) \in \mathcal{T}$ .

Let  $\mathcal{X}$  be all  $X = (x_{ij}) \in \mathbb{R}^{m \times n}$  where  $x_{ij} = 0$  if  $(i, j) \notin \mathcal{T}$ .

Assume that the completed matrix of the experiment should have the numerical rank  $\nu$ . Then we complete the entries by solving the problem:

$$(1) \quad \min_{X \in \mathcal{X}} \sum_{i=\nu+1}^n \sigma_i^2(A + X) = \min_{X \in \mathcal{X}} \sum_{i=\nu+1}^n \lambda_i((A + X)^T(A + X))$$

## 5 FRAA

### Fixed Rank Approximation Algorithm: [8]

Let  $G_p \in \mathcal{X}$  be the  $p^{\text{th}}$  approximation to a solution of optimization problem (1). Let

$B_p := (A + G_p)^T (A + G_p)$  and find an orthonormal set of eigenvectors for  $B_p$ ,  $v_{p,1}, \dots, v_{p,m}$ . Then  $G_{p+1}$  is a solution to the following minimum of a convex nonnegative quadratic function  $\min_{X \in \mathcal{X}} \sum_{q=l+1}^m ((A + X)v_{p,q})^T ((A + X)v_{p,q})$ .

Flow chart of the algorithm:

#### Fixed Rank Approximation Algorithm (FRAA)

**Input:** integers  $m, n, L, iter$ , the locations of non-missing entries  $\mathcal{S}$ , initial approximation  $G_0$  of  $n \times m$  matrix  $G$ .

**Output:** an approximation  $G_{iter}$  of  $G$ .

**for**  $p = 0$  **to**  $iter - 1$

- Compute  $B_p := (A + G_p)^T (A + G_p)$  and find an orthonormal set of eigenvectors for  $B_p$ ,  $v_{p,1}, \dots, v_{p,m}$ .

-  $G_{p+1}$  is a solution to the minimum problem (1) with  $\nu = L - 1 = l$ .

Let  $f_l(\mathbf{X}) := \sum_{i=\nu+1}^n \sigma_i^2(\mathbf{A} + \mathbf{X})$ . In each step of the algorithm  $f_l(\mathbf{G}_p) \geq f_l(\mathbf{G}_{p+1})$ .  $\mathbf{G}_p, p = 1, \dots$  converges to a critical point  $\tilde{\mathbf{G}}$ . FRAA gives a good approximation of  $\tilde{\mathbf{G}}$ . In many simulations  $\tilde{\mathbf{G}} = \mathbf{G}^*$ .

FRAA is an adaptation of an algo for IEP:

**Inverse Eigenvalue Problem:** Find the values of the missing entries of  $\mathbf{G}$  such that the nonnegative definite matrix  $\mathbf{G}^T \mathbf{G}$  will have  $m - l$  smallest eigenvalues equal to zero.

IEP appear often in engineering. See [9] for examples of IEP and a number of good algorithms to solve these problems.

FRAA is a robust algorithm which performs good, but not as well as KNNimpute, BPCA and LSSimpute.

All other algo reconstruct the missing values of each gene from similar genes.

## More about FRAA

### Proposition

Let  $G'$  be the  $n' \times m'$  matrix constructed by deleting  $i$  and  $j$  columns, the  $\sigma_q(G) \geq \sigma_q(G')$  for  $q = 1 \dots m$ .

### Ky-Fan Characterization

$$\begin{aligned} \sum_{q=l}^m \lambda_q(A) &= \sum_{q=l}^m z_q^T A z_q \\ &= \min_{\{y_l, \dots, y_m\} \in \Omega_{m-l+1}} \sum_{q=l}^m y_q^T A y_q \end{aligned}$$

## Computational Aspect of (FRAA)

The algorithm tries to minimize  $\sum_{l=L}^M \sigma_{l+1}^2(G)$ , where  $G$  ranges over all the completion of the missing data. We assume that the ideal completed gene matrix should have rank ( $k$ ).

- We start with  $G_0$  using SVD. Suppose we have  $G_p$ .  $G_p$  is the missing reconstruction data.
- Then algorithm for constructing  $G_{p+1}$ ,

$$A_p := G_p^T G_p \text{ an } n \times m \text{ matrix.}$$

- We use MatLab to find

$$\lambda_1(\mathbf{G}_p) \geq \dots \geq \lambda_m(\mathbf{G}_p),$$

and their corresponding eigenvector  $\mathbf{u}_{p,1} \dots \mathbf{u}_{p,m}$ .

$$\mathbf{G} = \mathbf{G}_p + \mathbf{X}.$$

$\mathbf{X}$  has all zero entries for the non missing data. Now we consider the quadratic form,

$$\begin{aligned} \sum_{i=1}^M \|\mathbf{G}_p + \mathbf{X}\mathbf{u}_{p,i}\|^2 = \\ \sum_{i=1}^M \mathbf{u}_{p,i}^T (\mathbf{G}_p + \mathbf{X})^T (\mathbf{G}_p + \mathbf{X}) \mathbf{u}_{p,i} \\ \mathbf{x}^T \mathbf{B}_p \mathbf{x} + 2\mathbf{w}_p^T \mathbf{x} + \sum_{i=1}^M \sigma_i^2(\mathbf{G}_p) = f(\mathbf{x}) \end{aligned}$$

and we want to minimize  $f(\mathbf{x})$ .

## Algorithm for $G_p$

- In the original matrix we substitute  $0$  for missing value or write down the row average
- Compute the singular value of  $G_p^T G_p$ .
- Construct  $B_p$ , with entries  $b_p(s, t)$ , in  $(s, t)$  places.

$$b_p(s, t) = \frac{1}{2} \sum_{i=l}^m v_{p,i}^T (F(i_s, j_s)^T F(i_t, j_t) + F(i_t, j_t)^T F(i_s, j_s)) v_{p,i}.$$

- The matrix  $B_p$  gives the exact solution of  $G_{p+1}$ .
- Solve for  $\mathbf{x}$  in  $B_p \cdot \mathbf{x}_{p+1} = -\mathbf{w}_p$ .

$$G_{p+1} := G_p + X_{p+1} \text{ until}$$

$$\frac{\sum_{i=L}^M \sigma_i^2(G_{p+1})}{\sum_{i=1}^M \sigma_i^2(G_{p+1})} \text{ is very small } (L = \ell + 1).$$

## 6 IMPROVED FRAA (IFRAA)

Improved Fixed Rank Approximation Algorithm [7].

First use FRAA to find a completion  $G$ .

Then use a cluster algorithm,

(We used K-means repeating & refining cluster size),

to find a reasonable number of clusters of similar genes, each cluster is a relatively smaller matrix having an effective low rank.

For each cluster of genes apply FRAA separately to recover the missing entries in this cluster.

*These results suggest that IFRAA has a potential for being an effective algorithm to recover blurred spots in digital images.*



## 7 Clustering

Given a metric space  $X$ ,  $d : X \times X \rightarrow \mathbb{R}_+$  and  $\mathcal{X} := \{x_1, \dots, x_n\} \subset X$  are  $n$  distinct points, associate  $M := (d(x_i, x_j))_{i,j=1}^n \in \mathbb{R}_+^{n \times n}$ .

**Problem:** Partition  $\mathcal{X}$  to clusters  $\mathcal{X} = \cup_{j=1}^m \mathcal{X}_j$  using  $M$ .

There are many different approaches to solve this problem.

## **K-Means Clustering**

- **Divide genes randomly in K-clusters**
- **Compute centroid**
- **Compute genes distance to clusters**
- **Pick closest cluster**
- **Recompute new centroid**
- **Repeat until convergence**

### **IFRAA(FRAA +K-Means Clustering)**

- **Impute  $G$  via FRAA**
- **Find clusters  $c_i$  using K-means.**
- **Compare  $c_i$  with  $G$ , remove data points corresponding to missing data.**
- **Impute applying FRAA to  $c_i$ .**

## 8 SIMULATIONS 1

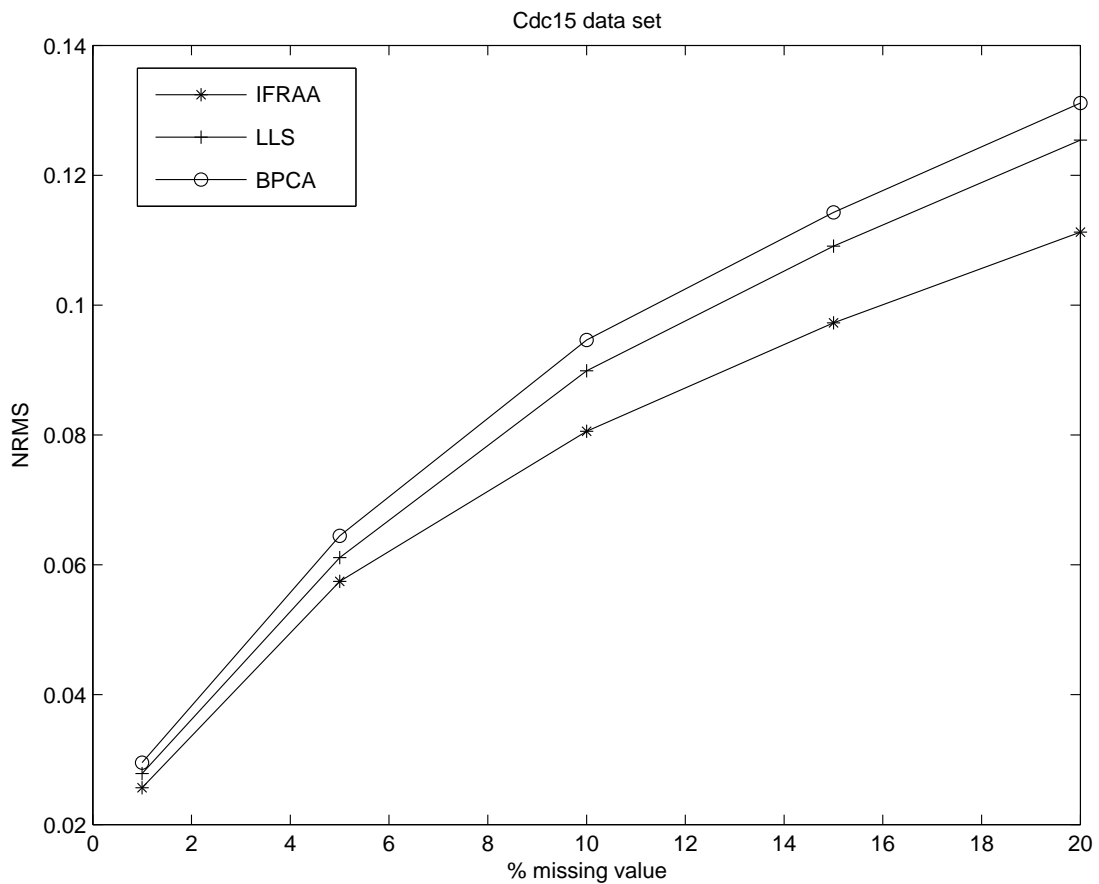


Figure 1: Comparison of NRMSE against percent of missing entries for three methods: IFRAA, BPCA and LLS. Cdc15 data set in [17] with 24 samples.

## 9 SIMULATIONS 2

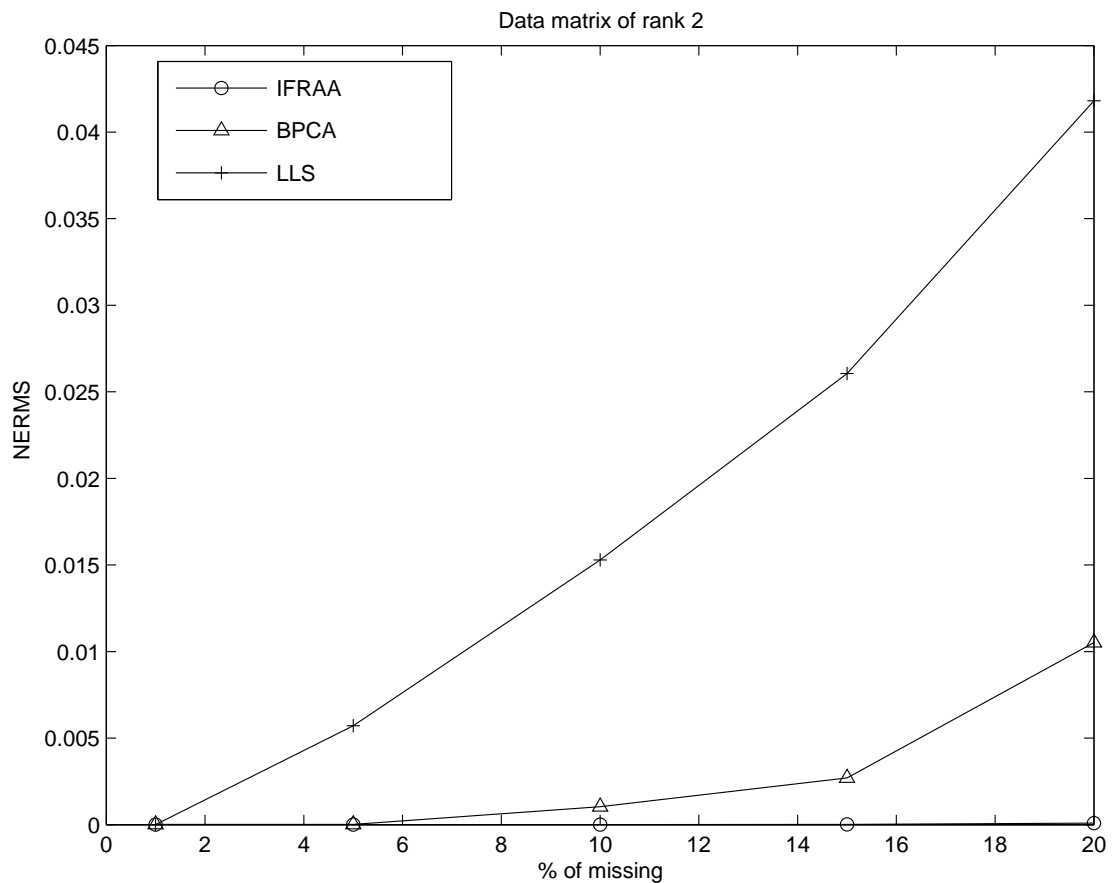


Figure 2: Comparison of NRMSE against percent of missing entries for three methods: IFRAA, BPCA and LLS. Data set was a  $2000 \times 20$  randomly generated matrix of rank 2.

**Bayesian principal component analysis-BPCA** [15]: A global method consisting of three components. First, principal component regression, which is basically a low rank approximation of the data set is performed. Second, Bayesian estimation, which assumes that the residual error and the projection of each gene on principal components behave as normal independent random variables with unknown parameters, is carried out. Third, Bayesian estimation follows by iterations based on the expectation-maximization (EM) of the unknown Bayesian parameters.

**Local least squares imputation method LLS** [14]: A local methods, which use similarity structure of the data to impute the missing values. **LLS** has two versions to find similar genes whose expressions are not corrupted: the  $L_2$ -norm and the Pearson's correlation coefficients. After a group of similar genes  $C$  are identified, the missing values of the gene are obtained using least squares applied to the group  $C$ . The recovery of missing data is done independently, i.e. the estimation of each missing entry does not influence the estimation of the other missing entries.

## 10 TABLE

The performance of the BCPA, IFRAA and LLS algorithms depends on the unknown distribution of missing position of the entries.

Table 1: Comparison of NRMSE for three methods: IFRAA, LLS and BPCA for actual missing values distribution for three gene expression data sets with different percentage of missing values.

Data sets	IFRAA	LLS	BPCA
Cdc15 data set %0.81 missing	0.0175	0.0200	0.0216
Evolution data set %9.16	0.0703	0.0969	0.1247
Calcineurin data set %3.68	0.0421	0.0445	0.0453

## 11 Generalized SVD

Let  $A \in \mathbb{C}^{m \times n}$ ,  $B \in \mathbb{C}^{l \times n}$ . Then Van Loan 70s:

$$A = F\Gamma R, \quad B = G\Delta R,$$

$$F \in U(m), \quad G \in U(l), \quad R \in GL(n, \mathbb{C}),$$

$$\Gamma \in \mathbb{R}_+^{m \times n}, \quad \Delta \in \mathbb{R}_+^{l \times n} \text{ diagonal matrices.}$$

Numerical computations of GSVD are very unstable.

Thm ([5]). Let  $P := A^*A + B^*B$  and  $r := \text{rank } P$ .

Then  $A = U\Phi V^*$ ,  $U \in U(m, r)$ ,  $V \in \mathbb{C}^{n \times r}$ ,

$$B = W\Psi V^*, \quad W \in U(l, r),$$

$$\Phi = \text{diag}(\phi_1, \dots, \phi_r),$$

$$\Psi = \text{diag}(\psi_1, \dots, \psi_r) \in \mathbb{R}_+^{r \times r} \text{ and}$$

$$\Phi^2 + \Psi^2 = I_r.$$

Hence  $P = VV^*$  and the columns of  $V$  form an orthonormal basis of the subspace  $\mathbf{X}$ , spanned by the columns of  $A^*$ ,  $B^*$  with respect to the inner product  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{y}^* P \mathbf{x}$  on  $\mathbf{V}$ .

Reason:  $T_A^* T_A + T_B^* T_B = I|_{\mathbf{X}} \Rightarrow$

$$(T_A^* T_A)(T_B^* T_B) = (T_B^* T_B)(T_A^* T_A)$$



## References

- [1] O. Alter, P.O. Brown and D. Botstein, Singular value decomposition for genome-wide expression data processing and modelling, *Proc. Nat. Acad. Sci. USA* 97 (2000), 10101-10106.
- [2] O. Alter, P.O. Brown and D. Botstein, Generalized singular decomposition for comparative analysis of genome-scale expression data sets of two different organisms, *Proc. Nat. Acad. Sci. USA* 100 (2003), 3351-3356.
- [3] O. Alter, G.H. Golub, P.O. Brown and D. Botstein, Novel genome-scale correlation between DNA replication and RNA transcription during the cell cycle in yeast is predicted by data-driven models, 2004 *Miami Nature Winter Symposium*, Jan. 31 - Feb. 4, 2004, to appear.
- [4] A. Deshpande, L. Rademacher, S. Vemapala and G. Wang, Matrix Approximation and Projective Clustering via Volume Sampling, *SODA*, 2006.
- [5] S. Friedland, A New Approach to Generalized Singular Value Decomposition, *SIMAX* 27 (2005), 434-444.

- [6] S. Friedland, M. Kaveh, A. Niknejad and H. Zare, Fast Monte-Carlo Low Rank Approximations for Matrices, Proc. IEEE SoSE, 2006.
- [7] S. Friedland, M. Kaveh, A. Niknejad and H. Zare, An Algorithm for Missing Value Estimation for DNA Microarray Data, Proc. ICASSP, 2006
- [8] S. Friedland, A. Niknejad and L. Chihara, A simultaneous reconstruction of missing data in DNA microarrays, to appear in *Linear Algebra and Its Applications*.
- [9] S. Friedland, J. Nocedal and M. Overton, The formulation and analysis of numerical methods for inverse eigenvalue problems, SIAM J. Numer. Anal. 24 (1987), 634-667.
- [10] A. Frieze, R. Kannan and S. Vempala, Fast Monte-Carlo algorithms for finding low rank approximations, *Proceedings of the 39th Annual Symposium on Foundation of Computer Science*, 1998.
- [11] G.H. Golub and C.F. Van Loan, *Matrix Computation*, John Hopkins Univ. Press, 3rd Ed., 1996.

- [12] R.A. Horn and C.R. Johnson, *Matrix Analysis*, Cambridge Univ. Press, 1987.
- [13] R.A. Johnson, D. W. Wichern, *Applied Multivariate Statistical Analysis*, Prentice Hall, New Jersey, 4th edition (1998).
- [14] H. Kim, G.H. Golub and H. Park, Missing value estimation for DNA microarray gene expression data: local least squares imputation, *Bioinformatics* 21 (2005), 187-198.
- [15] S. Oba, M. Sato, I. Takemasa, M. Monden, K. Matsubara and S. Ishii, A Bayesian missing value estimation method for gene expression profile data, *Bioinformatics* 19 (2003), 2088-2096.
- [16] C.C. Paige and M. A. Saunders, Towards a generalized singular value decomposition, *SIAM J. Numer. Anal.* 18 (1981), 398-405.
- [17] P.T. Spellman, G. Sherlock, M.Q. Zhang, V.R. Iyer, K. Anders, M.B. Eisen, P.O. Brown, D. Botstein and B. Futcher, Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces*

*cerevisiae* by microarray hybridization, *Mol. Biol. Cell*, **9** (1998), 3273-3297.

- [18] G.W. Stewart, A method for computing the generalized singular value decomposition, *Matrix Pencils*, B. Kagström and A. Ruhe, *Lecture Notes in Mathematics*, 973 (1982), 207-220.