

Classification of Protein Interactions Based on Sparse Discriminant Analysis and Energetic Features

Katarzyna Stapor, Piotr Fabian

► **To cite this version:**

Katarzyna Stapor, Piotr Fabian. Classification of Protein Interactions Based on Sparse Discriminant Analysis and Energetic Features. 15th IFIP International Conference on Computer Information Systems and Industrial Management (CISIM), Sep 2016, Vilnius, Lithuania. pp.530-537, 10.1007/978-3-319-45378-1_47 . hal-01637469

HAL Id: hal-01637469

<https://hal.inria.fr/hal-01637469>

Submitted on 17 Nov 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Classification of protein interactions based on sparse discriminant analysis and energetic features

Katarzyna Stapor¹, Piotr Fabian¹

¹ Silesian Technical University, Faculty of Automatic Control,
Electronics and Computer Science, Akademicka 16, Gliwice, Poland

{Katarzyna.Stapor, Piotr.Fabian}@polsl.pl

Abstract. Prediction of protein-protein interaction (PPI) types is an important problem in life sciences because of fundamental role of PPIs in many biological processes. In this paper we propose a new classification approach based on the extended classical Fisher linear discriminant analysis (FLDA) to predict obligate and non-obligate protein-protein interactions. To characterize properties of the protein interaction, we proposed to use the binding free energies (total of 282 features). The obtained results are better than in the previous studies.

Keywords: sparse discriminant analysis, feature selection, protein-protein interaction.

1 Introduction

Prediction of protein-protein interaction (PPI) types is an important problem in life sciences because of fundamental role of PPIs in many biological processes. PPIs have been investigated in various ways, involving both experimental (in vivo or in vitro) and computational (in silico) approaches [2,10]. Experimental approaches tend to be costly, labor intensive and suffer from noise. Therefore, using computational approaches for prediction of PPIs is a good choice for many reasons.

PPIs can be divided into non-obligate and obligate complexes (binding components can/cannot form stable functional structures). Based on the duration and life time of the interactions, there are transient complexes and permanent ones. Although interfaces have been the main subject of study to predict protein-protein interactions, an accuracy of 70% has been independently achieved by several different groups ([9,10,13,14]). These approaches have been carried out by analyzing a wide range of parameters, including solvation energies, amino acid composition, conservation, electrostatic energies, and hydrophobicity.

These includes a method based on PCA combined with Chernoff extension of Fisher linear discriminant analysis [9]. PCA is necessary to reduce the dimensionality of the input feature space (i.e. to be less than the sample size). As a consequence some important information is lost.

In this paper, we propose a new classification approach based on sparse discriminant analysis [12] to predict obligate (permanent) and non-obligate (transient) protein-protein interactions. To characterize properties of protein interaction, we proposed to use the binding free energies.

2 Fisher and Sparse regularized linear discriminant analyses

Fisher Linear Discriminant analysis (FLDA) [5,11] is a multivariate technique which is concerned with the search for a linear transformation that reduces the dimension of a given p -dimensional statistical model to q ($q < p$) dimensions, while maximally preserving the discriminatory information for the several classes within the model.

Formally, suppose that there are k classes and let $x_{ij}, j = 1, \dots, n_i$ be vectors of observations from the i -th class, $i = 1, \dots, k$. Set $n = n_1 + \dots + n_k$ and let $X_{n \times p} = (x_{11}, \dots, x_{1n_1}, \dots, x_{k1}, \dots, x_{kn_k})^T$, where p is the dimensionality of the input space. FLDA determines a linear mapping L , i.e. a $q \times p$ matrix A , that maximizes the so-called Fisher criterion J_F :

$$J_F(A) = \text{tr}((AS_W A^T)^{-1}(AS_B A^T)) \quad (1)$$

where $S_B = \sum_{i=1}^k p_i(m_i - \bar{m})(m_i - \bar{m})^T$ and $S_W = \sum_{i=1}^k p_i S_i$ are the between-class and the average within-class scatter matrix, respectively;

$S_i = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (x_{ij} - m_i)(x_{ij} - m_i)^T$ is the within-class covariance matrix of class i ,

m_i is the mean vector of class i , p_i is its *a priori* probability, and $\bar{m} = \sum_{i=1}^k p_i m_i$ is the overall mean. FLDA maximizes the ratio of between-class scatter to average within-class scatter in the lower-dimensional space. Optimizing (1) comes down to determining an eigenvalue decomposition of $S_W^{-1} S_B$, and taking the rows of A to equal the q eigenvectors corresponding to the q largest eigenvalues. There are no more than $\min(p, k - 1)$ eigenvectors corresponding to nonzero eigenvalues.

When the number of variables exceeds the sample size, i.e., in the high-dimensional, low-sample size (HDLSS) settings, the within-class covariance matrix S_W is singular and the classical FLDA breaks down. Several extensions have been proposed to overcome this problem but all of them possess the data piling problem [8]. To ameliorate this problem, some sparse version of LDA have been proposed.

In our approach, to circumvent this problem, we adapt the sparse linear discriminant approach (slda) from [12] that incorporates feature selection in FLDA. The term ‘‘sparse’’ means that the discriminant vectors have only a small number of nonzero components. The underlying assumption is that, among the large number of variables there are many irrelevant or redundant variables for the purpose of classification. This method is based on the connection of FLDA and a generalized eigenvalue problem, stated formally by the following theorem:

Theorem [12]:

Suppose S_w is a positive definite matrix and denote its Cholesky decomposition as $S_w = R_w^T R_w$ (R_w is an upper triangular matrix). Let H_b be $k \times p$ matrix, V_1, \dots, V_q ($q < \min(p, k-1)$) denote the eigenvectors of $S_w^{-1} S_B$ corresponding to the q largest eigenvalues $\lambda_1 \geq \dots \geq \lambda_q$, $A = [\alpha_1, \dots, \alpha_q]$, $B = [\beta_1, \dots, \beta_q]$. For $\lambda > 0$ let \hat{A}, \hat{B} be the solution to the following problem:

$$\min_{A, B} \sum_{i=1}^k \left\| R_w^{-T} H_{b,i} - AB^T H_{b,i} \right\|^2 + \lambda \sum_{j=1}^q \beta_j^T (S_w) \beta_j \quad \text{subject to } A^T A = I_{p \times q},$$

where:

$$H_{b,i} = \sqrt{n_i} (\bar{x}_i - \bar{x})^T \quad \text{is the } i\text{-th row of the matrix}$$

$$H_b = \left(\sqrt{n_1} (\bar{x}_1 - \bar{x}), \dots, \sqrt{n_k} (\bar{x}_k - \bar{x}) \right)^T,$$

e^{n_i} is a vector of ones with length n_i ,

Then $\hat{\beta}_j, j = 1, \dots, q$, span the same linear space as $V_j, j = 1, \dots, q$.

The following method of regularization is applied in [12] to circumvent the singularity problem and to obtain the sparse linear discriminants: i.e. the first q sparse discriminant directions β_1, \dots, β_q are defined as the solutions to the following optimization problem:

$$\min_{A, B} \sum_{i=1}^k \left\| R_w^{-T} H_{b,i} - AB^T H_{b,i} \right\|^2 + \lambda \sum_{j=1}^q \beta_j^T \left(S_w + \gamma \frac{\text{tr}(S_w)}{p} I \right) \beta_j + \sum_{j=1}^q \lambda_{1,j} \|\beta_j\|_1$$

subject to $A^T A = I_{p \times q}$, where $B = [\beta_1, \dots, \beta_q]$, $\|\beta_j\|_1$ is the 1-norm of the vector β_j , the same λ is used for all q directions, different $\lambda_{1,j}$'s are allowed to penalize different discriminant directions.

Our empirical study suggests that the solution is very stable when λ varies in a wide range, for example in (0.01, 10000). We can use K-fold cross validation (CV) [11] to select the optimal parameters $\lambda_{1,j}$, but when the dimension of the input data is very large, the numerical algorithm becomes time consuming and we can let $\lambda_{1,1} = \dots = \lambda_{1,q}$. The tuning parameter γ controls the strength of the regularization of the matrix S_w , the large values will bias too much S_w towards identity matrix (high degree of regularization). In our empirical studies, we find that the results are not sensitive to the choice of γ if a small value that is less than 0.1 is used, in our studies we set $\gamma = 0.05$. More careful studies of choice of γ are left for future research.

The above problem can be numerically solved by alternating optimization over A and B [12] and the resulting algorithm is summarized below.

Regularized sparse LDA (rSLDA) algorithm (based on [12]):

1. Form the matrices from the input data:

$$H_w = X - \begin{pmatrix} e^{n_1} (\bar{x}_1)^T \\ \vdots \\ e^{n_k} (\bar{x}_k)^T \end{pmatrix}$$

$$H_b = \left(\sqrt{n_1} (\bar{x}_1 - \bar{x}), \dots, \sqrt{n_k} (\bar{x}_k - \bar{x}) \right)^T$$

2. Compute upper triangular matrix R_w from the Cholesky decomposition of:

$$\left(S_w + \gamma \frac{\text{tr}(S_w)}{p} I \right) \text{ such that } \left(S_w + \gamma \frac{\text{tr}(S_w)}{p} I \right) = R_w^T R_w$$

3. Solve the q independent optimization problems

$$\min_{\beta_j} \beta_j^T (\tilde{W}^T \tilde{W}) \beta_j - 2 \tilde{y}^T \tilde{W} \beta_j + \lambda_1 \|\beta_j\|_1, j = 1, \dots, q$$

where

$$\tilde{W}_{(n+p) \times p} = \begin{pmatrix} H_b \\ \sqrt{\lambda} \cdot R_w \end{pmatrix} \quad \tilde{y}_{(n+p) \times 1} = \begin{pmatrix} H_b R_w^{-1} \alpha_j \\ 0 \end{pmatrix}$$

4. Compute SVD:

$$R_w^{-T} (H_B^T H_B) B = U D V^T \text{ and let } A = U V^T$$

5. Repeat steps 3 and 4 until converges.

3 Protein-protein interaction classification method

To characterize properties of protein interaction, we proposed to use the binding free energies. These were computed using FastContact [4], which obtains their fast estimates. FastContact delivers the electrostatic energy, solvation free energy, and the top 20 maximum and minimum values for:

1. residues contributing to the binding free energy,
2. ligand residues contributing to the solvation free energy,
3. ligand residues contributing to the electrostatic energy,
4. receptor residues contributing to the solvation free energy,
5. receptor residues contributing to the electrostatic energy,
6. receptor-ligand residue solvation constants,
7. receptor-ligand residue electrostatic constants.

Thus, all these values and the total solvation and electrostatic energy values compose a total of 282 features characterizing interaction.

To create a dataset for classification, we used the pre-classified dataset from previous study [9] containing 62 transient and 75 obligate complexes as two different classes for classification. Each complex is listed in the form of chains for ligand and receptor respectively. The relevant data about the structure of each complex was obtained from the Protein Data Bank (PDB) [1] and then obtaining the 282 features by invoking FastContact.

Due to the fact that the number of features (282) is greater than the number of samples in a dataset (137), we have HDLSS setting, so we apply sparse regularized linear discriminant analysis for the calculation of discriminant directions, i.e. the algorithm sparsed rLDA described above.

For the classification of the samples in the new discriminant space, we applied the nearest centroid method [11] as the classification algorithm.

4 Rapid estimation of contact and binding free energies

The estimation of contact and binding free energies may be in general a time consuming job. One of components of the binding energy is electrostatic energy. This term applies to a system of charges and is defined as the work necessary to move all the electric charges from infinity to positions occupied in the analyzed system. This work does not depend on the path traveled by the charges and is one of properties of a static arrangement of charges in space. Electrostatic interaction works on relatively long distances [7]. For proteins, it refers to the interaction between electrically charged atoms in different proteins or interactions between charges in the surface of the protein and charges in the environment. The exact computation of this energy for each possible conformation would be time consuming.

We have used a method called FastContact, developed by Camacho et al. [4,3]. The binding energy is computed as a sum of the electrostatic potential and the desolvation free energy in proteins: $G_{\text{binding}} = E_{\text{electrostatic}} + G_{\text{desolvation}}$. In this formula, $E_{\text{electrostatic}}$ is the standard intermolecular Coulomb electrostatic potential with relative permittivity varying with the distance r and equal to $4r$. The term $G_{\text{desolvation}}$ includes basic features of the desolvation free energy in proteins: hydrophobic interactions, self-energy change resulting from desolvating charge on polar atom groups and side-chain entropy loss. The $G_{\text{desolvation}}$ term is calculated as a sum of atomic contact potentials (ACP) between all pairs of atoms, where the first atom belongs to the receptor, the second to the ligand. Each term of this sum is additionally multiplied by a function $g(r)$ depending on the distance r between involved atoms. For $r > 7 \text{ \AA}$ the value is 0, for $r < 5 \text{ \AA}$ is 1 and between 5 \AA and 7 \AA $g(r)$ is a smooth function. These ACPs are defined for 18 atom types. The function $g(r)$ makes possible faster computation by zeroing interactions between distant atoms.

5 Experimental results

In our experiments we have used the dataset of 137 protein complexes described in [13]. 75 samples in this dataset belong to the first class (i.e. “obligate interactions”) and 62 samples to the second class (i.e. “non-obligate interactions”). This dataset is randomly divided into a “training set” and “testing set” in a ratio of 4:1.

As we have only two classes ($k=2$), there is only one discriminant direction β_1 ($q=1$). Using all variables in constructing the discriminant vector β_1 might cause the overfitting of the training data, resulting in high testing error rate. Moreover it is computationally demanding, so sparsification would be a good choice.

Denote the number of significant variables involved in specifying the discriminant direction β_1 (i.e. giving the best prediction), to be m . To find these most significant variables we have performed the experiment with varying values of m . For a given value of m , only the m maximum values of the coordinates of the vector β_1 (so called beta values) are left, the rest is zeroed.

Fig. 2 shows the components of vector β_1 obtained by the SLDA algorithm in one of experiments converted to the absolute values and sorted in the ascending order.

We leave only m biggest values, zeroing all others. We keep track of indices of these biggest values and modify the original β_1 leaving only m biggest values. These values are used to cast the original 282-dimensional vector onto a one-dimensional space. The projection of the samples from the protein dataset uses only these m non-zero coefficients.

Then, classification is performed in such new discriminant space by the nearest centroid method.

The results are shown in Fig. 1. We can observe that the error rate of the nearest mean classifier grows rapidly and then decreases with the rise of m , up to 28 (error = $\sim 25\% \pm 5$ measured on the testing set). Then, for bigger values of m , almost a constant error rate was observed.

From the plot it is clear that if we specify $m=28$ as the number of component variables in discriminant vector β_1 - sparse LDA algorithm can discriminate the two classes fairly well (the classifier performance = $\sim 75\% \pm 5$).

These 28 input features (“selected” by the rslda algorithm) are the most significant for classification. These are the following (corresponding to the ascending order of the absolute value of the coefficients composing vector β_1):

202 198 281 200 48 42 243 203 47 133 128 121 161 160
157 132 49 156 46 134 241 131 155 158 127 119 135 41

Among these 28 features – 13 are from the receptor residues contributing to the desolvation free energy, but these are not from the beginning of the above list. It can be observed that in each of the 7 groups of energetic features – only features with extreme (min or max) contribution to the energy are always selected. The features from the beginning of the list are those from the receptor residues contributing to the electrostatics energy. One may conclude that electrostatic energy is the most important in the prediction of obligate/non-obligate protein-protein interactions. Electro-

static energy involves a long-range interaction and occur between charged atoms of two interacting proteins.

Thus, the rslda algorithm does suggest which constituents are the most important in the classification of interactions.

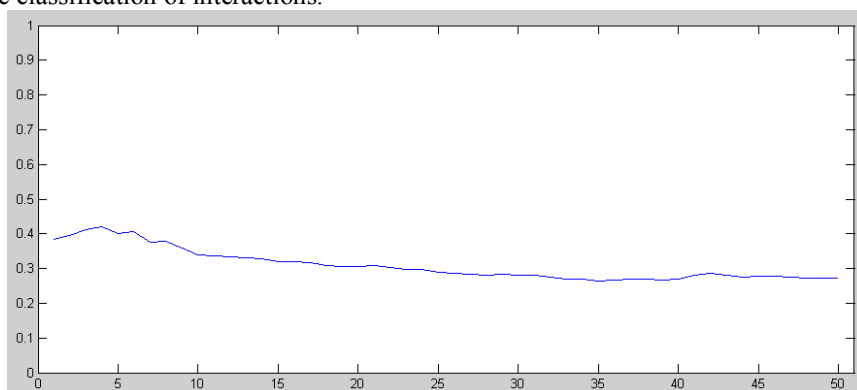


Fig. 1. The average classification error rate as a function of the number of variables using nearest centroid method [10] on the projected data (based on 5 random partitions of the dataset into training and test) – the local minimum is at 28

6 Conclusions

We have proposed a classification approach for obligate/non-obligate (transient) protein-protein complexes. We have used regularized version of sparse linear discriminant analysis algorithm [12] for feature extraction as well as for input variable selection. To discriminate between two types of protein interactions: obligate and non-obligate, we have used the “energetic features”. These are based on the binding free energy defined as the sum of the desolvation and electrostatic energies. These were computed effectively using the package FastContact [4]. The results on the protein-protein interactions dataset showed that using only 28 from 282 input variables enables the classification of the mentioned two types of interactions with the performance of 75%. Among the most important features are those from residues contributing to the electrostatic energy.

The hypothesis on the importance of the electrostatic energy in the prediction of obligate/non-obligate protein-protein interactions should be confirmed by the additional experiments on bigger protein datasets. This will be the subject of our future research.

References

1. Berman H. et al. (2000) *The Protein Data Bank*. Nucleid Acid Research 28, 235-242.
2. Bordner A., Abagyan R. (2005) *Statistical analysis and prediction of protein-protein interfaces*. Proteins 60(3), 353-366.

3. Camacho C. J., Gatchell D. W., Kimura S. R., Vajda S. (2000) *Scoring Docked Conformations Generated by Rigid-Body Protein-Protein Docking*. *PROTEINS: Structure, Function, and Genetics* 40, 525–537.
4. Camacho C., Zhang C. (2005) *FastContact: rapid estimate of contact and binding free energies*. *Bioinformatics* 21(10), 2534-2536.
5. Fukunaga K. (1990) *Introduction to statistical pattern recognition*. New York: Academic Press.
6. Jones S., Thornton J.M. (1996) *Principles of protein-protein interactions*. *Proc. Natl. Acad. Sci. USA* 93(1), 13-20.
7. Maleki M., Vasudev G., Rueda L. (2013) *The role of electrostatic energy in prediction of obligate protein-protein interactions*. *Proteome Sci.*, 11 (Suppl. 1) (2013), p. S11.
8. Marron J. et al. (2007). *Distance-weighted discrimination*. *Journal of American Statistical Association*, 102, 1267-1273.
9. Rueda L. et al. (2010) *Biological protein-protein interaction prediction using binding free energies and linear dimensionality reduction*. In: Dijkstra T., et al. (eds): *PRIB 2010*, LNBI 6282, 383-394, Springer Berlin.
10. Skrabanek L. et al (2008) *Computational prediction of protein-protein interactions*. *Molecular Biotechnology*, 38(1), 1-17.
11. Stapor K. (2011) *Classification methods in computer vision*. PWN Warszawa (in Polish).
12. Qiao Z., Zhou L., Huang J. (2009) *Sparse linear discriminant analysis with applications to high dimensional low sample size data*. *IAENG Int. Journal of Applied Mathematics*, 39, 1.
13. Zhou H., Shan Y. (2001) *Prediction of protein-protein interaction sites from sequence profile and residue neighbor list*. *Proteins* 44(3), 336-343.
14. Zhu H., et al. (2006) *NoxClass: prediction of protein-protein interaction types*. *BMC Bioinformatics* 7(27).