

FunPred 3.0: Improved Protein function prediction using protein interaction network: Supplementary Document

Sovan Saha¹, Piyali Chatterjee², Subhadip Basu³, Mita Nasipuri³, Dariusz Plewczynski^{4,5}

¹ Department of Computer Science and Engineering, Dr. Sudhir Chandra Sur Degree Engineering College, DumDum, Kolkata, India

² Department of Computer Science and Engineering, Netaji Subhash Engineering College, Garia, Kolkata, India

³ Department of Computer Science and Engineering, Jadavpur University, Kolkata, West Bengal, India

⁴ Centre of New Technologies, University of Warsaw, Warsaw, Poland

⁵ Faculty of Mathematics and Information Science, Warsaw University of Technology, Warsaw, Poland

Corresponding Author:

Dariusz Plewczynski^{4,5}

Banacha 2c Street, 02-097 Warsaw, Poland

Email address: d.plewczynski@cent.uw.edu.pl

1 DEFINITIONS AND NOTATIONS

2

3 Before proceeding into the main section of our work, it is important to discuss the
4 graphical properties as well as other relevant terms associated with our work.

5

6 **Protein-protein interaction network:** Protein-protein interactions occur when
7 two or more proteins bind together, often to carry out their biological function.
8 These protein interactions form a network like structure which is known as a
9 protein interaction network. Protein interaction network is generally represented as
10 a graph consisting of a set of nodes connected by edges or links. Proteins are
11 represented as nodes in the graph and the edges signify interactions between two
12 proteins. Here protein interaction network is represented as a graph G_p which
13 consists of a set of vertex V (nodes) connected by edges E (links). Thus $G =$
14 (V, E) .

15 **Protein complex/cluster:** It can be defined as group of proteins (usually in close
16 proximity to one another) interconnected through a network to work as one
17 centralized data processing resource. Here it is defined by C_i where i represent
18 cluster number.

19 **Sub graph:** A graph G'_p is a sub graph of G_p if the vertex set of G'_p is a subset of
20 the vertex set of G_p and if the edge set of G'_p is a subset of the edge set of G_p .
21 That is, if $G'_p = (V', E')$ and $G_p = (V, E)$, then G'_p is called as sub graph of G_p
22 if $V' \subseteq V$ and $E' \subseteq E$. G'_p may be defined as a set of $\{P_U | P_A\}$ where P_U represents
23 the set of un-annotated proteins while P_A represents the set of annotated protein.

24 **Level – 1 Neighbors:** For any vertex v in G'_p , all those vertices in G'_p that are
25 connected with v through an edge are deemed Level – 1 neighbors of v .

26 **Edge weight (W_{uv}):** The weight W_{uv} of edge (u, v) (Wang & Wu 2013) is defined
27 as the similarity between u and v . It is obvious that two nodes with an edge
28 between them belong to the same cluster if they have high similarity. The
29 similarity between u and v is measured by Jaccard's coefficient. Jaccard's
30 coefficient adopts the proportion of common neighbors of two nodes in all distinct
31 neighbors of these nodes to measure node similarity in complex networks.
32 Obviously, the more common neighbors two nodes share, the higher similarity
33 these nodes have. Therefore, the edge weight W_{uv} is represented by

$$34 \quad w_{uv} = (\Gamma(u) \cap \Gamma(v)) / (\Gamma(u) \cup \Gamma(v)) \quad (1)$$

35 where, $\Gamma(u)$ and $\Gamma(v)$ are neighbors of u and v respectively. $\Gamma(u) \cap \Gamma(v)$ represents
36 all common neighbors of u and v , and $\Gamma(u) \cup \Gamma(v)$ represents all distinct neighbors
37 of u and v . In our algorithm, edge weight is used to guarantee that in the same

38 cluster every pair of nodes with an edge between them should have relatively high
39 similarity.

40 **Neighborhood graph (G_v):** The neighborhood graph of $v \in V$ consists of v , all its
41 neighbors and the edges among them. It is defined as $G_v = (V', E')$, in which $V' =$
42 $\{v\} \cup \{u | u \in V, (u, v) \in E\}$, and $E' = \{(u_i, u_j) | (u_i, u_j) \in E, u_i, u_j \in V'\}$.

43 **Node weight (W_v):** In G_v , there are some nodes with degree 1 that only have
44 connections with v and the connections among these nodes are often false positive
45 according to topological reliability measures (Wang & Wu 2013). So nodes with
46 degree 1 and corresponding edges are removed from G_v . The remaining sub graph
47 of G_v is marked as G'_v . The node weight w_v of node $v \in V$ in PPI networks is the
48 average degree of all nodes in G'_v . It is represented by

$$49 \quad w_v = \sum_{u \in V''} \deg(u) / |V''| \quad (2)$$

50 where, V'' is the set of nodes in G'_v . $|V''|$ is the number of nodes in G'_v . And $\deg(u)$
51 is the degree of a node $u \in V''$ in W_v . In our algorithm, the weight W_v of a node v
52 $\in V$ is used in the step of seed chosen. Higher value of W_v of a graph indicates a
53 collection of nodes with maximum interactions among them and hence the graph
54 is densely connected region.

55 **Physico-Chemical Properties (PCP):** Physico-Chemical Properties (Saha &
56 Chatterjee 2014; Singh et al. 2008) of amino acids are the various features of
57 protein which are used to predict protein class. These properties are very
58 important in protein class prediction. The various Physico-Chemical Properties
59 used in this work are as given below:

60 **1. Extinction Coefficient (E_{protein}):** Extinction Coefficient (Singh et al. 2008) is a
61 protein parameter that is commonly used in the laboratory for determining the
62 protein concentration in a solution by spectrophotometry. It describes to what
63 extent light is absorbed by the protein and depends upon the protein size and
64 composition as well as the wavelength of the light. For proteins measured in water
65 at wavelength of 280nm, the value of the Extinction coefficient can be determined
66 from the composition of Tyrosine, Tryptophan and Cystine.

67 Mathematically it can be defined as:

$$68 \quad E_{\text{protein}} = (N_{\text{tyr}} \times E_{\text{tyr}}) + (N_{\text{trp}} \times E_{\text{trp}}) + (N_{\text{cys}} \times E_{\text{cys}}) \quad (3)$$

69 where $E_{\text{tyr}} = 1490$, $E_{\text{trp}} = 5500$, $E_{\text{cys}} = 125$ are the Extinction coefficients of the
70 individual amino acid residues.

71 **2. Absorbance (Optical Density):** For proteins measured in water at wavelength
72 of 280nm the absorbance can be determined by the ratio of Extinction coefficient
73 and the molecular weight of the protein. It is a representation of a material's light
74 blocking ability (Singh et al. 2008).

75 Mathematically absorbance is defined as:

76
$$\text{Absorbance} = E_{\text{protein}} / \text{Molecular Weight} \quad (4)$$

77 **3.Number of Negatively Charged Residues (N_{neg}):** This can be calculated from
78 the composition of Aspartic acid and Glutamic acid (Singh et al. 2008).

79

80 **4.Number of Positively Charged Residues (N_{pos}):** This can be calculated from
81 the composition of Arginine and Lysine (Singh et al. 2008).

82

83 **5.Aliphatic Index (AI):** The aliphatic index of a protein is defined as the relative
84 volume occupied by aliphatic side chains (alanine, valine, isoleucine, and leucine).
85 It may be regarded as a positive factor for the increase of thermo stability of
86 globular (Singh et al. 2008).

87 Mathematically aliphatic index is defined as:

88
$$AI = X_{\text{ala}} + a \times X_{\text{val}} + b \times (X_{\text{ile}} + X_{\text{leu}}) \quad (5)$$

89 where X_{ala} , X_{val} , X_{ile} and X_{leu} are the mole percentages of alanine, valine,
90 isoleucine and leucine respectively. Coefficients a and b are the relative volume of
91 valine side chain and side chains to the side chain of alanine i.e. a = 2.9 and b =
92 3.9.

93 **6.Compute IP/Mol weight:** It calculates the isoelectric point by molecular weight
94 (Singh et al. 2008) of the input amino acid sequence. IP stands for isoelectric point
95 of the input amino acid sequence. Mol weight stands for molecular weight of the
96 input amino acid sequence.

97

98 **7.Grand average of hydropathicity (GRAVY):** The GRAVY value for a protein
99 or a peptide (Kyte & Doolittle 1982) is calculated by adding the hydropathy values
100 of each amino acid residues and dividing by the number of residues in the
101 sequence or length of the sequence. Increasing positive score indicates a greater
102 hydrophobicity.

103

104 **8.Instability index:** The instability index (Guruprasad et al. 1990) provides an
105 estimate of the stability of your protein in a test tube. A protein whose instability
106 index is smaller than 40 is predicted as stable, a value above 40 predicts that the
107 protein may be unstable.

108

109 **9.Aromaticity:** Aromaticity (Lobry & Gautier 1994) is simply the relative
110 frequency of phenylalanine (Phe), tryptophan (Trp), tyrosine (Tyr).

111

112 **10. Isoelectric point:** The isoelectric point (Bjellqvist et al. 1994) is the pH at
113 which a molecule or surface carries no net electrical charge.

114

115 11. **PCP_{score}** : PCP_{score} is defined as scaling of the mean value obtained from
116 top-ranked physico-chemical properties among the properties mentioned above
117 which are obtained by the execution of four classifiers: XGBoost classifier,
118 Random Forest classifier, Extra Tree classifier and Recursive feature elimination
119 classifier.

120

121 **XGBoost Classifier:** XGBoost is a scalable end to end tree boosting system
122 which proves to be highly effective and widely used machine learning method for
123 feature selection .

124

125 **Random forest Classifier:** A random forest is defined to be a meta-estimator that fits
126 a number of decision tree classifiers on various sub-samples of the dataset and use
127 averaging to improve the predictive accuracy and control over-fitting (Breiman 2001).

128

129 **Extra Tree Classifier:** Extra Tree classifier is a new tree-based ensemble method for
130 supervised classification of feature selection (Geurts et al. 2006; Pedregosa et al.
131 2011).

132

133 **Recursive feature elimination (RFE) Classifier:** RFE is used to select features by
134 recursively considering smaller and smaller sets of features (Pedregosa et al. 2011).
135 First, the estimator is trained on the initial set of features and the importance of each
136 feature is obtained either through a *coef_* attribute or through a *feature_importances_*
137 attribute. Then, the least important features are pruned from current set of features.

138 REFERENCES

139

140 Bjellqvist B, Basse B, Olsen E, and Celis JE. 1994. Reference points for comparisons of two-
141 dimensional maps of proteins from different human cell types defined in a pH scale where
142 isoelectric points correlate with polypeptide compositions. *ELECTROPHORESIS*
143 15(1):529-539.

144

145 Breiman L. 2001. Random Forests. *Machine Learning* 45(1):5-32.

146

147 Geurts P, Ernst D, and Wehenkel L. 2006. Extremely randomized trees. *Machine Learning* 63(1):3-
148 42.

149

150 Guruprasad K, Reddy BVB, and Pandit MW. 1990. Correlation between stability of a protein and
151 its dipeptide composition: a novel approach for predicting in vivo stability of a protein
152 from its primary sequence. *Protein Engineering, Design and Selection* 4(2):155-161.

153

154 Kyte J, and Doolittle RF. 1982. A simple method for displaying the hydropathic character of a
155 protein. *Journal of Molecular Biology* 157(1):105-132.

156

157 Lobry JR, and Gautier C. 1994. Hydrophobicity, expressivity and aromaticity are the major trends
158 of amino-acid usage in 999 Escherichia coli chromosome-encoded genes. *Nucleic Acids*
159 *Research* 22(15):3174-3180.

160

161

162 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer
163 P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M,
164 and Duchesnay E. 2011. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*
165 12:2825-2830.
166
167 Saha S, and Chatterjee P. 2014. Protein Function Prediction from Protein Interaction Network
168 using Physico-Chemical Properties of Amino Acids. *International journal of pharmacy*
169 *and biological sciences* 4(2):55-65.
170
171 Singh M, Wadhwa PK, and Kaur S. 2008. Predicting Protein Function using Decision Tree. *World*
172 *Academy of Science, Engineering and Technology* 2(3):300-303.
173
174 Wang S, and Wu F. 2013. Detecting overlapping protein complexes in PPI networks based on
175 robustness. *Proteome Science* 11(Suppl 1):S18-S18.
176