

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

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 (Chen & Guestrin 2016; Pedregosa et al. 2011).

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 Pedregosa et al. 2011).

129
130 **Extra Tree Classifier:** Extra Tree classifier is a new tree-based ensemble method for
131 supervised classification of feature selection (Geurts et al. 2006; Pedregosa et al.
132 2011).

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

139 INFORMATION REGARDING MIPS DATABASE

140
141 The PPIN of MIPS (Munich Information Center for Protein Sequences) (Mewes et
142 al. 2002) database of yeast has been used in this work for considering protein pair
143 along with its corresponding functions. This dataset is available at their website:

144 (<ftp://ftpmips.helmholtzmuellen.de/fungi/Saccharomycetes/CYGD/PPI/>)

145 The Munich Information Center for Protein Sequences (MIPS-GSF, Neuherberg,
146 Germany) (Mewes et al. 2002) continues to provide genome-related information in
147 a systematic way. MIPS supports both national and European sequencing and
148 functional analysis projects, develops and maintains automatically generated and
149 manually annotated genome-specific databases, develops systematic classification
150 schemes for the functional annotation of protein sequences, and provides tools for
151 the comprehensive analysis of protein sequences. The MIPS dataset of yeast
152 obtained from the link mentioned above contains protein pairs along with their
153 corresponding functions like Protein A|Protein B|DNA Repair (suppose for
154 example) i.e. when Protein A interacts with Protein B they perform the function
155 DNA Repair (annotated). Mewes et al. (Mewes et al. 2002) stated these as
156 “*Genomes that are being annotated and published by MIPS*”.

157 According to Mewes et al. (Mewes et al. 2002), as the amount of specialist yeast
158 related data continues to grow, they are exploring a model to integrate additional
159 data collections and knowledge into the Comprehensive Yeast Genome Database
160 (CYGD). CYGD is built upon collaboration with several yeast laboratories and
161 includes specialized databases. “20000 newly identified genes” from 13
162 hemiascomycetous yeasts, generated by the Genolevure project, have already been
163 integrated. He also classified some of the genes to be “*Unfinished and/or*
164 *unpublished genomic sequences*” in which he stated that “Gene prediction
165 conducted by ORPHEUS in a completely automatic fashion, usually allows large
166 overlaps between ORFs. This leads to many overpredicted ORFs, but ensures that
167 *fewer real ORFs are missed.*”

168 This lead to the development of “two types of interactions” like this in their MIPS
169 dataset:

- 170 1. Protein B|Protein C|unknown
- 171 2. Protein D|Protein E|missing

172 **Protein B|Protein C|unknown:** This signifies that Protein B interacts with
173 Protein C. This interaction already exists but their implementing automated
174 methodology for protein pair function prediction failed to predict the functions
175 when Protein B interacts with Protein C. That’s why they have given “unknown”
176 in the function field to signify that this is an existing pair whose function is yet to
177 be predicted.

178 **Protein D|Protein E|missing:** This signifies that Protein D interacts with Protein
179 E. But these are the newly identified proteins and interactions. While attempting to
180 predict functional annotations in automated fashion this interaction gets missed
181 due to excessive overlapping or overprediction of ORFs. That’s why they have
182 given “missing” in the function field to signify that this is a missing interaction
183 pair (newly predicted) whose function is yet to be predicted.

184 So both can be basically classified as “*unpredicted protein pair interactions*” i.e.
185 protein interactions whose functional annotations are not yet predicted.

186 INFORMATION REGARDING PROTEIN-PAIR FUNCTION 187 PREDICTION

189 A protein may perform various functions in isolation. But it does not perform all
190 the functions while reacting with another. It may perform some specific functions
191 while interacting with one protein while perform some other specific functions
192 while reacting with other proteins. So besides predicting protein function, protein
193 pair function also needs to be determined. So various researches have been
194 conducted in this field of study (Chatterjee et al. 2012; Shatsky et al. 2016). In
195 disease based PPIN, where function of one protein (say Protein A) is known but
196 the function of its interacting protein (responsible for causing disease)(say Protein

197 B) is not known then function of Protein B can be predicted from Protein A since
198 conventional approaches associate protein interaction with the sharing of
199 functions: “if proteins A and B belong to the same functional pathway, A is likely
200 to interact with B; therefore when A and B are observed to interact, they are likely
201 to share functions” (Chua et al. 2006).

202 **INFORMATION REGARDING METHODOLOGY**

203

204 FunPred 3 is broadly classified into two sections:

205 *First section* involves:

- 206 1. Selection of test set proteins (proteins considered as unannotated which are
207 annotated in real).
- 208 2. Prediction of functional annotations of test set by the proposed
209 methodology.
- 210 3. Computation of the effectiveness of the prediction of our proposed
211 methodology through the computation of precision, recall and F-Score.

212 *Second section* involves only the prediction of unknown/missing protein pair
213 function by our proposed methodology.

214 Hence Precision, recall and F-score have been computed in the first section. Since
215 originally the functions of the test proteins are annotated but we consider them to
216 be unannotated, so after the predicting the functions of test set proteins we can
217 match them with the original defined ones. If it matches then it is considered as
218 true positives. Similarly False Positives etc. are also computed. Suppose, for
219 example Protein A has originally DNA Repair function and it is included in our
220 test set proteins. So we consider that the function of Protein A is unannotated and
221 hence predict it by our proposed methodology. Now if our proposed methodology
222 predicts the function of Protein A as DNA Repair then we consider it as a match
223 with the original one (i.e. True Positive) else not.

224 Since our methodology FunPred 3 chooses only the essential proteins as test set
225 proteins (by the application of node and edge weight) in the entire PPIN of yeast
226 so it has been observed that these essential proteins belong to near about 155
227 diversified functional groups which is extensively large when compared to its
228 predecessors FunPred 1 and FunPred 2.

229 Both GO and the MIPS functional catalogues are hierarchical. But MIPS contain
230 certain common GO functions. Moreover we have not used FunCat id (like
231 11.06.03.01, 16 etc.) for function prediction. Instead we have used direct functions
232 (like mRNA editing, transcription) of proteins i.e. if protein A is originally
233 annotated to “mRNA editing” in MIPS dataset and our prediction model annotates
234 it as “transcription” then it is not considered as true positive. True positive is

235 considered only when our prediction model annotates it as “mRNA editing”.
236 FunCat id is considered as one of our future work which is already in progress.

237 REFERENCES

238

- 239 Bjellqvist B, Basse B, Olsen E, and Celis JE. 1994. Reference points for comparisons of two-
240 dimensional maps of proteins from different human cell types defined in a pH scale where
241 isoelectric points correlate with polypeptide compositions. *ELECTROPHORESIS*
242 15(1):529-539. 10.1002/elps.1150150171.
- 243
- 244 Breiman L. 2001. Random Forests. *Machine Learning* 45(1):5-32. 10.1023/a:1010933404324.
- 245
- 246 Chatterjee T, Chatterjee P, Basu S, Kundu M, and Nasipuri M. 2012. Protein function by minimum
247 distance classifier from protein interaction network. In: 2012 International Conference on
248 Communications, Devices and Intelligent Systems (CODIS).588-591.
- 249
- 250 Chen T, and Guestrin C. 2016. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the
251 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data
252 Mining. San Francisco, California, USA: ACM.785-794.
- 253
- 254 Chua HN, Sung W-K, and Wong L. 2006. Exploiting indirect neighbours and topological weight to
255 predict protein function from protein–protein interactions. *Bioinformatics* 22(13):1623-
256 1630. 10.1093/bioinformatics/btl145.
- 257
- 258 Geurts P, Ernst D, and Wehenkel L. 2006. Extremely randomized trees. *Machine Learning* 63(1):3-
259 42. 10.1007/s10994-006-6226-1.
- 260
- 261 Guruprasad K, Reddy BVB, and Pandit MW. 1990. Correlation between stability of a protein and
262 its dipeptide composition: a novel approach for predicting in vivo stability of a protein
263 from its primary sequence. *Protein Engineering, Design and Selection* 4(2):155-161.
264 10.1093/protein/4.2.155.
- 265
- 266 Kyte J, and Doolittle RF. 1982. A simple method for displaying the hydropathic character of a
267 protein. *Journal of Molecular Biology* 157(1):105-132. [https://doi.org/10.1016/0022-
268 2836\(82\)90515-0](https://doi.org/10.1016/0022-2836(82)90515-0).
- 269
- 270 Lobry JR, and Gautier C. 1994. Hydrophobicity, expressivity and aromaticity are the major trends
271 of amino-acid usage in 999 Escherichia coli chromosome-encoded genes. *Nucleic Acids
272 Research* 22(15):3174-3180.
- 273
- 274 Mewes HW, Frishman D, Güldener U, Mannhaupt G, Mayer K, Mokrejs M, Morgenstern B,
275 Münsterkötter M, Rudd S, and Weil B. 2002. MIPS: a database for genomes and protein
276 sequences. *Nucleic Acids Research* 30(1):31-34.
- 277
- 278 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer
279 P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M,
280 and Duchesnay E. 2011. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*
281 12(2825-2830).
- 282
- 283 Saha S, and Chatterjee P. 2014. Protein Function Prediction from Protein Interaction Network
284 using Physico-Chemical Properties of Amino Acids. *International journal of pharmacy
285 and biological sciences* 4(2):55-65.
- 286

287 Shatsky M, Allen S, Gold BL, Liu NL, Juba TR, Reveco SA, Elias DA, Prathapam R, He J, Yang
288 W, Szakal ED, Liu H, Singer ME, Geller JT, Lam BR, Saini A, Trotter VV, Hall SC,
289 Fisher SJ, Brenner SE, Chhabra SR, Hazen TC, Wall JD, Witkowska HE, Biggin MD,
290 Chandonia J-M, and Butland G. 2016. Bacterial Interactomes: Interacting Protein Partners
291 Share Similar Function and Are Validated in Independent Assays More Frequently Than
292 Previously Reported. *Molecular & Cellular Proteomics* : MCP 15(5):1539-1555.
293 10.1074/mcp.M115.054692.
294
295 Singh M, Wadhwa PK, and Kaur S. 2008. Predicting Protein Function using Decision Tree. *World*
296 *Academy of Science, Engineering and Technology* 2(3):300-303.
297
298 Wang S, and Wu F. 2013. Detecting overlapping protein complexes in PPI networks based on
299 robustness. *Proteome Science* 11(Suppl 1):S18-S18. 10.1186/1477-5956-11-S1-S18.

300