

 Open access • Posted Content • DOI:10.1101/149492

Selection of most relevant centrality measures: A systematic survey on protein-protein interaction networks — [Source link](#)

Minoo Ashtiani, Ali Salehzadeh-Yazdi, Zahra Razaghi-Moghadam, Holger Hennig ...+3 more authors

Institutions: Pasteur Institute of Iran, University of Rostock, University of Tehran, Tarbiat Modares University

Published on: 02 Oct 2017 - bioRxiv (Cold Spring Harbor Laboratory)

Topics: Random walk closeness centrality, Katz centrality, Centrality, Betweenness centrality and Network theory

Related papers:

- [A systematic survey of centrality measures for protein-protein interaction networks](#)
- [Revisiting the General Concept of Network Centralities: A Propose for Centrality Analysis in Network Science](#)
- [Graph Energy Based Centrality Measure to Identify Influential Nodes in Social Networks](#)
- [Comparative Analysis of Centrality Measures of Network Nodes based on Principal Component Analysis](#)
- [Centrality Measures in Complex Networks: A Survey.](#)

Share this paper:    


View more about this paper here: <https://typeset.io/papers/selection-of-most-relevant-centrality-measures-a-systematic-im0yiueyok>

RESEARCH ARTICLE

Open Access



A systematic survey of centrality measures for protein-protein interaction networks

Minoo Ashtiani^{1†}, Ali Salehzadeh-Yazdi^{2†}, Zahra Razaghi-Moghadam^{3,4}, Holger Hennig², Olaf Wolkenhauer², Mehdi Mirzaie^{5*} and Mohieddin Jafari^{1*} 

Abstract

Background: Numerous centrality measures have been introduced to identify “central” nodes in large networks. The availability of a wide range of measures for ranking influential nodes leaves the user to decide which measure may best suit the analysis of a given network. The choice of a suitable measure is furthermore complicated by the impact of the network topology on ranking influential nodes by centrality measures. To approach this problem systematically, we examined the centrality profile of nodes of yeast protein-protein interaction networks (PPINs) in order to detect which centrality measure is succeeding in predicting influential proteins. We studied how different topological network features are reflected in a large set of commonly used centrality measures.

Results: We used yeast PPINs to compare 27 common of centrality measures. The measures characterize and assort influential nodes of the networks. We applied principal component analysis (PCA) and hierarchical clustering and found that the most informative measures depend on the network’s topology. Interestingly, some measures had a high level of contribution in comparison to others in all PPINs, namely Latora closeness, Decay, Lin, Freeman closeness, Diffusion, Residual closeness and Average distance centralities.

Conclusions: The choice of a suitable set of centrality measures is crucial for inferring important functional properties of a network. We concluded that undertaking data reduction using unsupervised machine learning methods helps to choose appropriate variables (centrality measures). Hence, we proposed identifying the contribution proportions of the centrality measures with PCA as a prerequisite step of network analysis before inferring functional consequences, e.g., essentiality of a node.

Keywords: Network science, Centrality analysis, Protein-protein interaction network (PPIN), Clustering, Principal components analysis (PCA)

Background

Essential proteins play critical roles in cell processes such as development and survival. Deletion of essential proteins is more likely to be lethal than deletion of non-essential proteins [1]. Identifying essential proteins conventionally had been carried out with experimental methods which are time-consuming and expensive, and such experimental approaches are not always feasible.

Analyzing high-throughput data with computational methods promises to overcome these limitations. Various computational methods have been proposed to predict and prioritize influential nodes (e.g. proteins) among biological networks. Network-based ranking (i.e. centrality analysis) of biological components has been widely used to find influential nodes in large networks, with applications in biomarker discovery, drug design and drug repurposing [2–6]. Not only in molecular biology networks but also in all types of networks, finding the influential nodes is the chief question of centrality analysis [7]. Examples include predicting the details of information controlling or disease spreading within a specific network in order to delineate how to effectively implement target marketing or preventive healthcare [8–10]. Several centralities measures (mostly in the context of social network analyses)

* Correspondence: mirzaie@modares.ac.ir; mjafari@pasteur.ac.ir; <https://www.jafarilab.com>

[†]Minoo Ashtiani and Ali Salehzadeh-Yazdi contributed equally to this work.

⁵Department of Applied Mathematics, Faculty of Mathematical Sciences, Tarbiat Modares University, P.O. Box 14115-134, Tehran, Iran

¹Drug Design and Bioinformatics Unit, Medical Biotechnology Department, Biotechnology Research Center, Pasteur Institute of Iran, P.O. Box 13164, Tehran, Iran

Full list of author information is available at the end of the article



have been described [7] in the last decades. A comprehensive list of centrality measures and software resources can be found on the CentiServer [11].

The correlation of lethality and essentiality with different centrality measures has been subject of active research in biological areas, which has led to the centrality-lethality rule [1]. Typically, some classic centrality measures such as Degree, Closeness, and Betweenness centralities have been utilized to identify influential nodes in biological networks [9]. For example, in a pioneering work, the authors found that proteins with the high Degree centrality (hubs) among a yeast PPIN is likely to be associated with essential proteins [1]. In another study, this rule was re-examined in three distinct PPINs of three species which confirmed the essentiality of highly connected proteins for survival [12]. Similar results were reported for gene co-expression networks of three different species [13] and for metabolic network of *Escherichia coli* [14, 15]. Ernesto Estrada generalized this rule to six other centrality measures. He showed that the Subgraph centrality measure scored best compared to classic measures to find influential proteins, and generally using these measures performed significantly better than a random selection [16]. However, He and Zhang showed that the relationship between hub nodes and essentiality is not related to the network architecture [17]. Furthermore, regarding the modular structure of PPINs, Joy et al. concluded that the Betweenness centrality is more likely to be essential than the Degree centrality [18]. The predictive power of Betweenness as a topological characteristic was also mentioned in mammalian transcriptional regulatory networks which was clearly correlated to Degree [19]. Recently, it has been shown that presence of hubs, i.e. high Degree centralities, do not have a direct relationship with prognostic genes across cancer types [20].

On the other hand, Tew and Li demonstrated functional centrality and showed that it correlates more strongly than pure topological centrality [21]. More recently, localization-specific centrality measures had been introduced and claimed that their results is more likely essential in different species [22–25]. In the same way, some studies emphasized on the protein complex and topological structure of a sub-network to refine PPIN and identify central nodes [26–28]. Tang et al. integrated the gene co-expression data on PPIN as edge weights to realize the reliable prediction of essential proteins [24]. Khuri and Wuchty introduced minimum dominating sets of PPIN which are enriched by essential proteins. They described that there is a positive correlation between Degree of proteins in these sets and lethality [29]. In these studies, the solution of the controversy is ascribed to utilizing biological information.

Similar in methodology but different in the underlying physical system that the network represents, some other

studies attempted to quantify correlations between several classic centrality measures. In 2004, Koschützki and Schreiber compared five centrality measures in two biological networks and showed different patterns of correlations between centralities. They generally concluded that all Degree, Eccentricity, Closeness, random walk Betweenness and Bonacich's Eigenvector centralities should be considered to find central nodes and could be useful in various applications without explaining any preference among them [30]. Two years later, they re-expressed previous outcomes by explaining the independence behavior of centrality measures in a PPIN using 3D parallel coordinates, orbit-based and hierarchy-based comparison [31]. Valente et al. examined the correlation between the symmetric and directed versions of four measures which are commonly used by the network analysts. By comparing 58 different social networks, they concluded that network data collection methods change the correlation between the measures and these measures show distinct trends [32]. Batool and Niazi also studied three social, ecological and biological neural networks and they concluded the correlation between Closeness-Eccentricity and Degree-Eigenvector and insignificant pattern of Betweenness. They also demonstrated that Eccentricity and Eigenvector measures are better to identify influential nodes [33]. In 2015, Cong Li et al. further investigated the question of correlation between centrality measures and introduced a modified centrality measure called *m*th-order degree mass. They observed a strong linear correlation between the Degree, Betweenness and Leverage centrality measures within both real and random networks [34].

However, there is no benchmark for network biologists that provides insight, which of the centrality measures is suited best for the analysis of the given network. The result of the centrality analysis of a network may depend on the used centrality measure which can lead to inconsistent outcomes. Previously, a detailed study showed that the predictive power and shortcomings of centrality measures are not satisfactory in various studies [35]. While these centrality measures have proven to be essential in understanding of the roles of nodes which led to outstanding contributions to the analysis of biological networks, choosing the appropriate measure for given networks is still an open question. Which measure identifies best the centers of real networks? Do all measures independently highlight the central network elements and encompass independent information or are the measures correlated? Is the computation of all these measures meaningful in all different networks or does the best measure depend on the network topology and the logic of the network reconstruction? In this study, we used unsupervised machine learning to compare how well the most common centrality measures characterize nodes in networks. We comprehensively

compared 27 distinct centrality measures applied to 14 small to large biological and random networks. All biological networks were PPINs of the same set of proteins which are reconstructed using a variety of computational and experimental methods. We demonstrated how the ranking of nodes depends on the network structure (topology) and why this network concept i.e. centrality deserves renewed attention.

Methods

The workflow of this study was schematically presented in Fig. 1. Our workflow started by constructing and retrieving networks, followed by global network analysis. The centrality analysis and comparing them using machine learning methods were the next main steps. See basic definitions for more details.

Reconstruction of the networks

In this study, a UniProtKB reviewed dataset [36] was used to retrieve proteins in *Saccharomyces cerevisiae* (6721 proteins). UniProtKB accessions were converted to STRING using the STRINGdb R package, which resulted in 6603 protein identifiers (3rd Sep 2016). Interactions among proteins were extracted based on the STRING IDs. In the 2017 edition of the STRING database the results of these interactions are structured in a way to provide maximum coverage; this is achieved by including indirect and predicted interactions on the top of the set. [37]. In this study, 13 evidence channels (related to the origin and type

of evidence) indicating PPIN of yeast were presented: co-expression, co-expression-transferred, co-occurrence, database, database-transferred, experiments, experiments-transferred, fusion, homology, neighborhood-transferred, textmining, textmining-transferred and combined-score (See Additional file 1). In the following, the name of the reconstructed network is basis of the corresponding channel name which made of. For the purpose of comparison with real network behavior, a null model network was generated. The null network is the Erdős–Rényi model [38] and was generated using the igraph R package [39]. The generated null network was created with a size similar to the yeast reconstructed PPIN in order to have a more fair comparison.

Fundamental network concepts analysis

To understand the network structure, we reviewed various network features using several R packages [40–42]. The network density, clustering coefficient, network heterogeneity, and network centralization properties of the network were calculated. The number of connected components and graph diameter for each network were also computed. Then, the power-law distribution was assessed by computing α values and r correlation coefficients. As most of centrality measures require a strongly connected component graph, the giant component of each PPINs and the null network were extracted. Moreover, for a general overview of the structure of the extracted giant components, some network features such as network density,

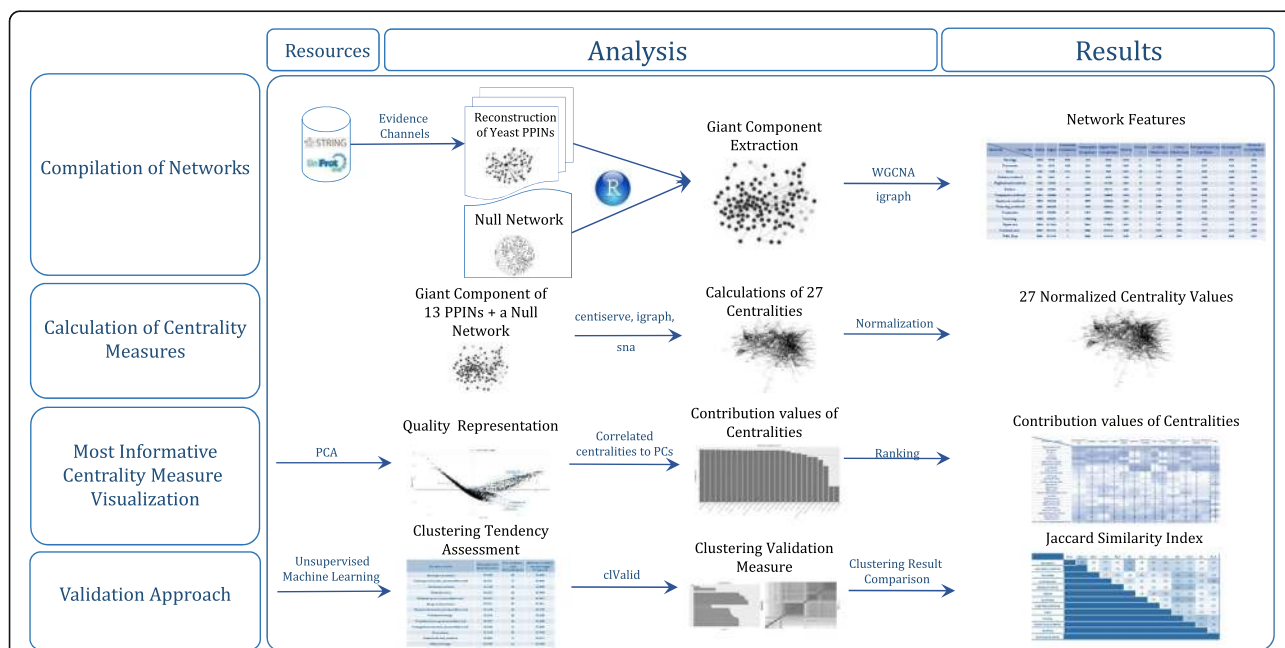


Fig. 1 Our workflow for studying the centrality measures. This was followed the reconstruction of the yeast PPIN relying on different kinds of evidence channels as well as the generation of a null network. The workflow contained a comparison of several centrality measures using machine learning methods such as principal components analysis and clustering procedures

clustering coefficient, network heterogeneity, and network centralization were calculated.

Centrality analysis

For this research study, we were only considered undirected, loop-free connected graphs according to the PPIN topology. For centrality analysis, the following 27 centrality measures were selected: Average Distance [43], Barycenter [44], Closeness (Freeman) [9], Closeness (Latora) [45], Residual closeness [46], ClusterRank [47], Decay [48], Diffusion degree [49], Density of Maximum Neighborhood Component (DMNC) [50], Geodesic K-Path [51, 52], Katz [53, 54], Laplacian [55], Leverage [56], Lin [57], Lobby [58], Markov [59], Maximum Neighborhood Component (MNC) [50], Radiality [60], Eigenvector [61], Subgraph scores [62], Shortest-Paths betweenness [9], Eccentricity [63], Degree, Kleinberg's authority scores [64], Kleinberg's hub scores [64], Harary graph [63] and Information [65]. All these measures are calculated for undirected networks in a reasonable time. These measures were calculated using the centiserve [11], igraph [39] and sna [66] R packages. Some of the centrality measures had a measurable factor to be specified which we used the default values. For a better visualization, We assorted the centrality measures into five distinct classes including Distance-, Degree-, Eigen-, Neighborhood-based and miscellaneous groups depend on their logic and formulas (Table 1).

Unsupervised machine learning analysis

Standard normalization (scaling and centering of matrix-like objects) has been undertaken on computed centrality values according to methodology explained in [67]. We

used PCA, a linear dimensionality reduction algorithm, [68] as a key step to understand which centrality measures better determine central nodes within a network. PCA was done on normalized computed centrality measures. To validate the PCA results in PPINs, we also examined whether the centrality measures in all networks can be clustered according to clustering tendency procedure. To do this, the Hopkins' statistic values and visualizing VAT (Visual Assessment of cluster Tendency) plots was calculated by factoextra R package [69]. We applied the clustering validation measures to access the most appropriate clustering method among hierarchical, k-means, and PAM (Partitioning Around Medoids) methods using cValid package [70]. This provides silhouette scores according to clustering measures which would be helpful for choosing the suitable method. After selection of the clustering technique, factoextra package was used to attain optimal number of clusters [69]. In order to measure the dissimilarity among clusters, we used Ward's minimum variance method. To compare the clustering results in aforementioned PPINs, the Jaccard similarity index was used relying on the similarity metrics of the clustering results within BiRewire package [71].

Results

Evaluation of network properties

By importing the same set of protein names, the 13 PPINs were extracted from the STRING database using different evidence channels. (Note: the PPI scores derived from the neighborhood channel of yeast were all zero). All these channels distinctly identify an interaction for each protein pair quantitatively. The dependency between evidence channels was also shown in Fig. 2 by a

Table 1 Centrality measures. The centrality measures were represented in five groups depending on their logic and formulae

Distance_based	Degree-based	Eigen-based	Neighborhood-based	Miscellaneous
Average Distance	Authority_score	Eigenvector centralities	ClusterRank	Geodesic K-Path Centrality
Barycenter	Degree Centrality	Katz Centrality (Katz Status Index)	Density of Maximum Neighborhood Component (DMNC)	Harary Graph Centrality
Closeness Centrality (Freeman)	Diffusion Degree	Laplacian Centrality	Maximum Neighborhood Component (MNC)	Information Centrality
Closeness centrality (Latora)	Kleinberg's hub centrality scores		Subgraph centrality scores	Markov Centrality
Decay Centrality	Leverage Centrality			Shortest-Paths Betweenness Centrality
Eccentricity of the vertices	Lobby Index (Centrality)			
Lin Centrality				
Radiality Centrality				
Residual Closeness Centrality				

Note that the first column (i.e. distance-based centralities) was specified according to the definition of distance between vertices in graph theory. The second one (i.e. degree-based centralities) was defined based on the number of immediate neighbors of each node within a given network. Eigen-values of adjacency matrix was the main idea to classify the Eigen-based centralities. Furthermore, the concept of subgraph or community structure was proposed in the neighborhood-based centralities. Others were collected in the miscellaneous group. Remind that this grouping was just applied to have better visualizations.

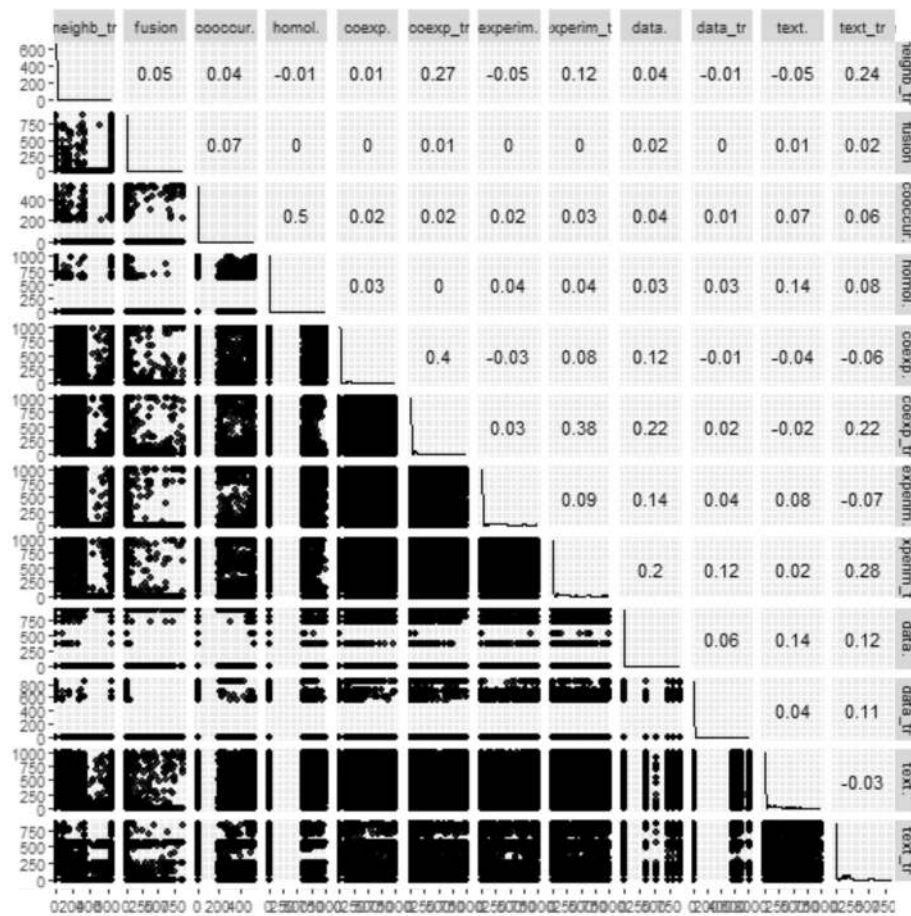


Fig. 2 Pairwise scatterplot between the evidence channel scores. The Pearson's r correlation coefficients between the evidence channels were shown in the upper triangle of the plot. The distributions of scores in each evidence were presented at the diameters of the figure

pairwise scatterplot and Pearson's r correlation coefficient. Most of the networks were not significantly correlated and correlation coefficients were around zero for all networks.

In the following, the 14 networks were utilized to undertake an examination of centrality measures. Note that the giant component of each network was accounted for computing several network properties (Table 2). The homology, fusion, co-occurrence and database networks contained high numbers of unconnected components. Except the homology network which had the smallest giant component, the densities of all networks were between 0.01–0.05, as was expected real network are typically sparse. The network diameter of the fusion, co-occurrence, database and co-expression were one order of magnitude greater than others. All of the PPINs except homology network were correlated to power-law distribution with high r correlation coefficients and diverse alpha power (see Additional file 2). The high value of the average clustering coefficients of the database and homology indicated the modular structure of these networks.

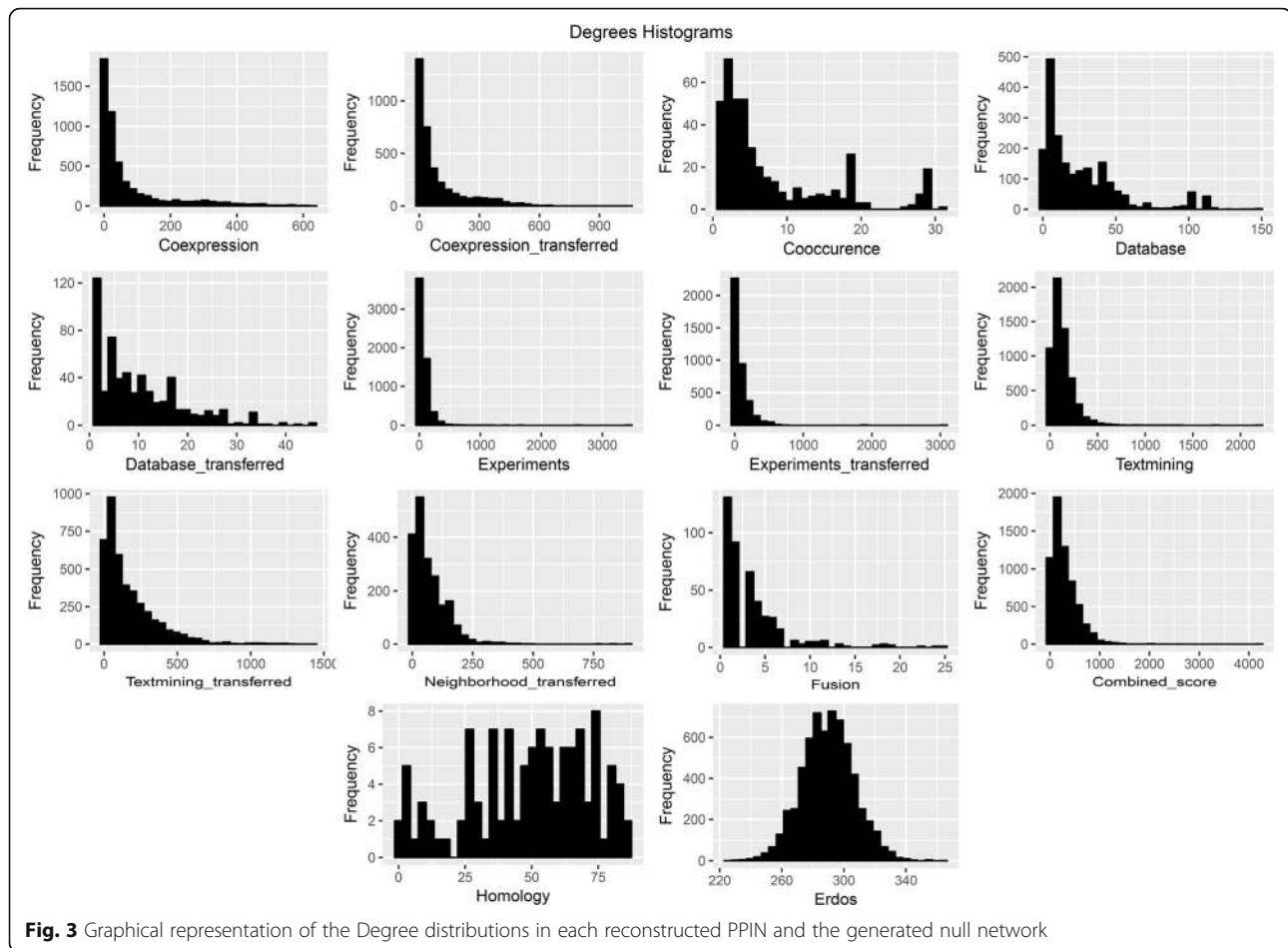
Compared with the null network, most of the PPINs had a high value of heterogeneity and network centralization. The Degree distribution and clustering coefficients for the networks were also plotted in Figs. 3 and 4 respectively. Except the homology network, all the Degree distributions were left-skewed similar to scale-free networks. The dependency of PPINs was further assessed and confirmed statistically by Wilcoxon rank sum test (Table 3).

Centrality analysis

In the next step, the 27 centrality measures of nodes were computed in all 14 networks. The distribution and pairwise scatter plots of the computed measures were represented in Fig. 5 to point out pairwise relationship between them. (For the other PPINs see Additional file 3). The r correlation coefficients were also shown in this figure in which some of the centrality measures displayed a clear correlation and the others revealed a vast diversity among all five centrality classes. This diversity especially enriched in Distance-, Neighborhood-based and miscellaneous classes for combined-score PPIN compared with Erdos-Renyi

Table 2 Network global properties of all PPINs and the null network

Networks Properties	Nodes	Edges	Connected Components	Nodes/giant component	Edges/Giant Component	Density	Diameter	α value (Power Law)	r value (Power Law)	Average Clustering Coefficient	Heterogeneity	Network Centralization
Homology	2479	7545	648	115	2813	0.43	5	0.01	0.00	0.64	0.47	0.34
Cooccurrence	1221	3275	209	425	1653	0.02	22	1.07	0.61	0.47	1.02	0.06
Fusion	1187	1408	222	437	789	0.01	26	1.76	0.91	0.07	1.00	0.05
Database_transferred	622	2975	10	583	2930	0.02	9	1.23	0.69	0.36	0.86	0.06
Neighborhood_transferred	2004	71,236	1	2004	71,236	0.04	6	0.91	0.72	0.26	1.04	0.41
Database	2496	27,766	100	2058	26,574	0.01	20	1.18	0.56	0.69	1.07	0.06
Coexpression_transferred	3614	168,368	4	3607	168,364	0.03	8	0.96	0.79	0.35	1.39	0.26
Experiments_transferred	3870	163,403	1	3870	163,403	0.02	6	1.03	0.81	0.35	1.59	0.77
Textmining_transferred	4207	364,816	1	4207	364,816	0.04	6	0.88	0.75	0.32	1.09	0.30
Coexpression	5310	195,676	27	5254	195,643	0.01	15	1.08	0.82	0.44	1.56	0.11
Textmining	5896	379,341	1	5896	379,341	0.02	5	1.01	0.66	0.30	0.92	0.35
Experiments	6026	211,613	2	6024	211,612	0.01	6	1.22	0.82	0.19	1.54	0.56
Combined_score	6294	911,414	2	6292	911,413	0.05	7	0.76	0.63	0.27	0.90	0.62
Erdős_Rényi	6292	911,413	1	6292	911,413	0.05	3	_1.095	0.01	0.05	0.06	0.01



network. Analogously, this special profile of centrality measures was repeated in all PPINs to some extent. Another remarkable distinction was the multimodality of distributions in the random network but not in real networks which was repeated for most of the Distance-based centrality measures. Furthermore, according to r correlation coefficients, the pairwise association of centrality measures were roughly higher in the null network than PPINs.

Dimensionality reduction and clustering analysis

In the next step, PCA-based dimensionality reduction was used to reveal which centrality measures contain the most relevant information in order to effectively identify important or influential nodes in networks. As illustrated in Fig. 6, the profile of the distance to the center of the plot and their directions were mostly consonant except for the homology which was similar to the random network. The rank of contribution values of each centrality measure were shown in Table 4, depend on their corresponding principal components. The percentage of contribution of variables (i.e. centrality measures) in a given PC were computed as $(\text{variable.Cos2} \times 100) / (\text{total Cos2 of the component})$. A similar profile of the contribution of

centrality measures was observed among all biological networks even in homology network opposed to the random null network (See Additional file 4). On average, Latora closeness centrality was the major contributor of the principal components in PPINs. In contrast, other well-known centralities i.e. Betweenness and Eccentricity revealed a low contribution value in all PPINs. Analogous to the null network, their values were lower than random threshold depicted in Fig. 8 and Additional file 4. On the contrary, the Degree displayed moderate levels of contribution in all real networks whilst it was the fourth rank of random network contributors. Although the profile of contributions were similar, each PPIN exhibited a special fingerprint of the centrality ranking. Finally, by performing unsupervised categorization, we aimed to cluster centrality values computed in the networks. First, we performed a clustering tendency procedure. We found that the centrality values are clusterable in each network as all values in the Hopkins statistics were more than the cutoff (0.05). The results are shown in the first column of Table 5 and Additional file 5. Then, by calculating silhouette scores, three methods (i.e. hierarchical, k-means, and PAM) were evaluated in clustering the data sets (Additional files 6 and 7). The output

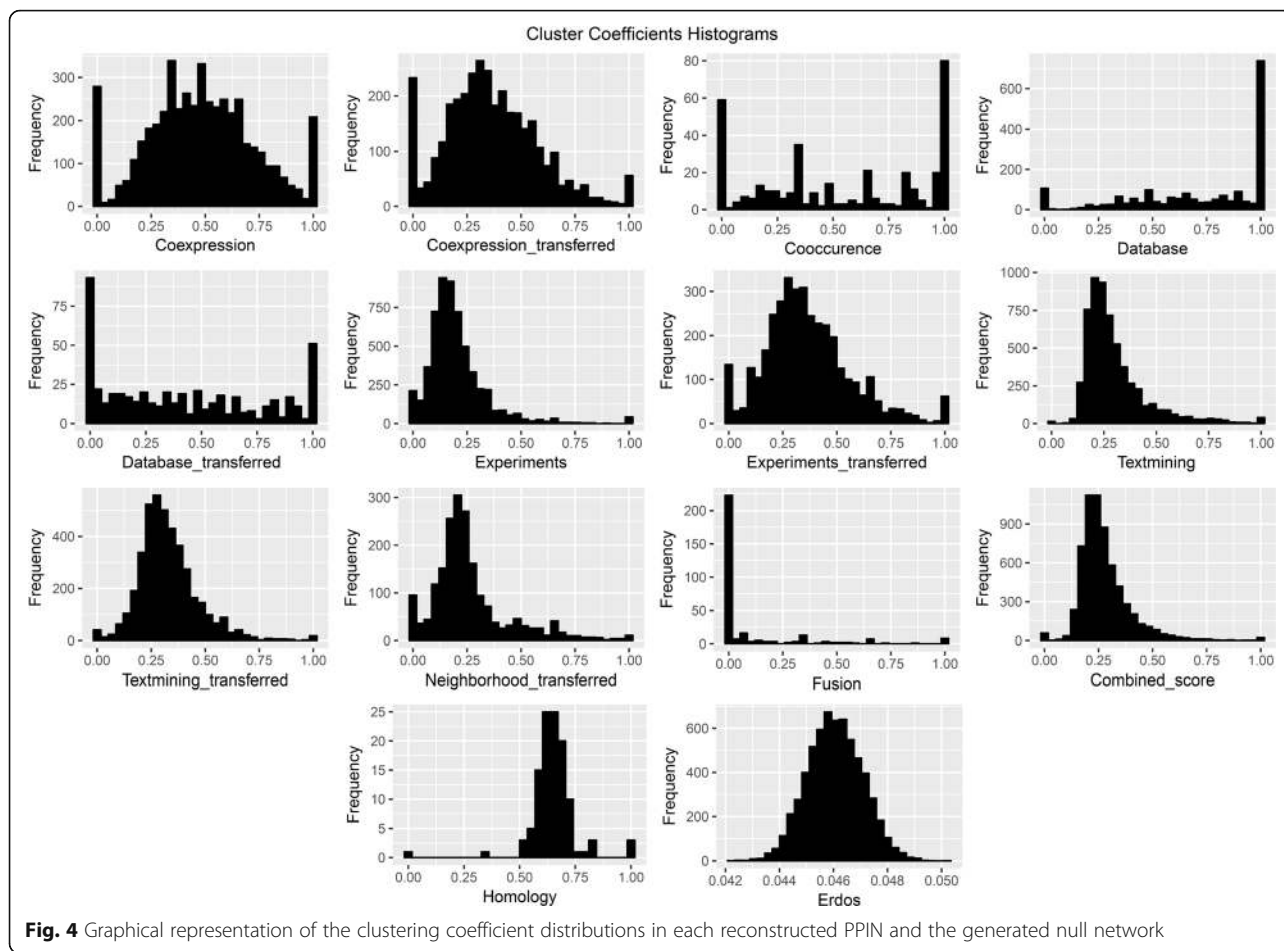


Fig. 4 Graphical representation of the clustering coefficient distributions in each reconstructed PPIN and the generated null network

of applying these algorithms and the corresponding number of clusters were also shown in Table 5 and Additional file 8. Using the hierarchical algorithm based on Ward’s method [72], the centrality measures were clustered in each PPINs (Fig. 7). Number of clusters, distance between centrality measures and centrality composition in all 13 PPINs indicated that each centrality ranks nodes within a given network distinctly. For a better comparison, we provided Table 6 containing pairwise Jaccard similarity indices for each network pair. The lowest values were related to the homology, neighborhood-transferred and co-occurrence PPINs while among these genome context prediction methods, fusion PPIN was more associated to the other networks. The high similarity between co-expression and co-expression-transferred was expected however the similar clusters of the database derived PPIN with both aforementioned PPINs and also combined-score with textmining-transferred are noteworthy.

Discussion

Interestingly, silhouette scores of centrality measures were closely related to corresponding contribution value of the measures (Fig. 8). Where there was a high

silhouette value, a high contribution value was observed, however, a high contribution value did not always mean a high silhouette value. The relationship between the silhouette scores and contribution values of each centrality measure was also examined by regression analysis. Latora closeness, Radiality, Residual, Decay, Lin, Leverage, Freeman closeness and Barycenter centrality measures were present together in the same cluster where the corresponding silhouette scores were all at a high level except the Leverage’s score (Fig. 8a). The average silhouette score was around 0.66 in this cluster. On the other hand, the Leverage’s contribution value was below the threshold line and placed in the group with the least amount of contribution (Fig. 8b). The centrality measures namely Lobby index, ClusterRank, Laplacian, MNC, Degree, Markov, Diffusion degree, Kleinberg’s hub, Eigen vector, Authority score, Katz group together where the mean of their silhouette scores (i.e. 0.61) was higher than the overall average and in the same way, their corresponding contribution values were high, too. On the other hand, we observed that Shortest path Betweenness (which was in a separated cluster) and Geodesic k path, Subgraph and DMNC (which are all in one cluster) showed the low

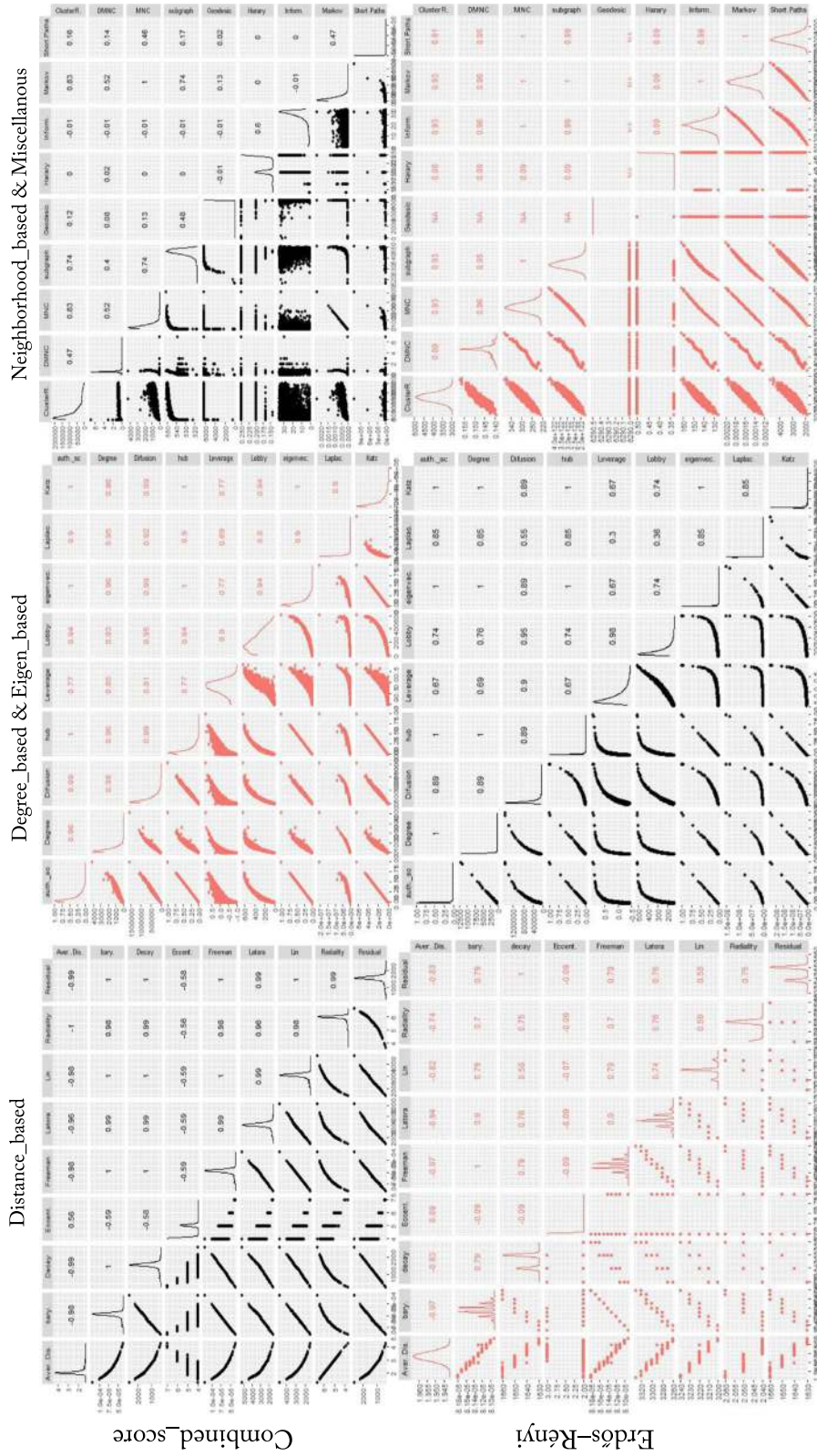
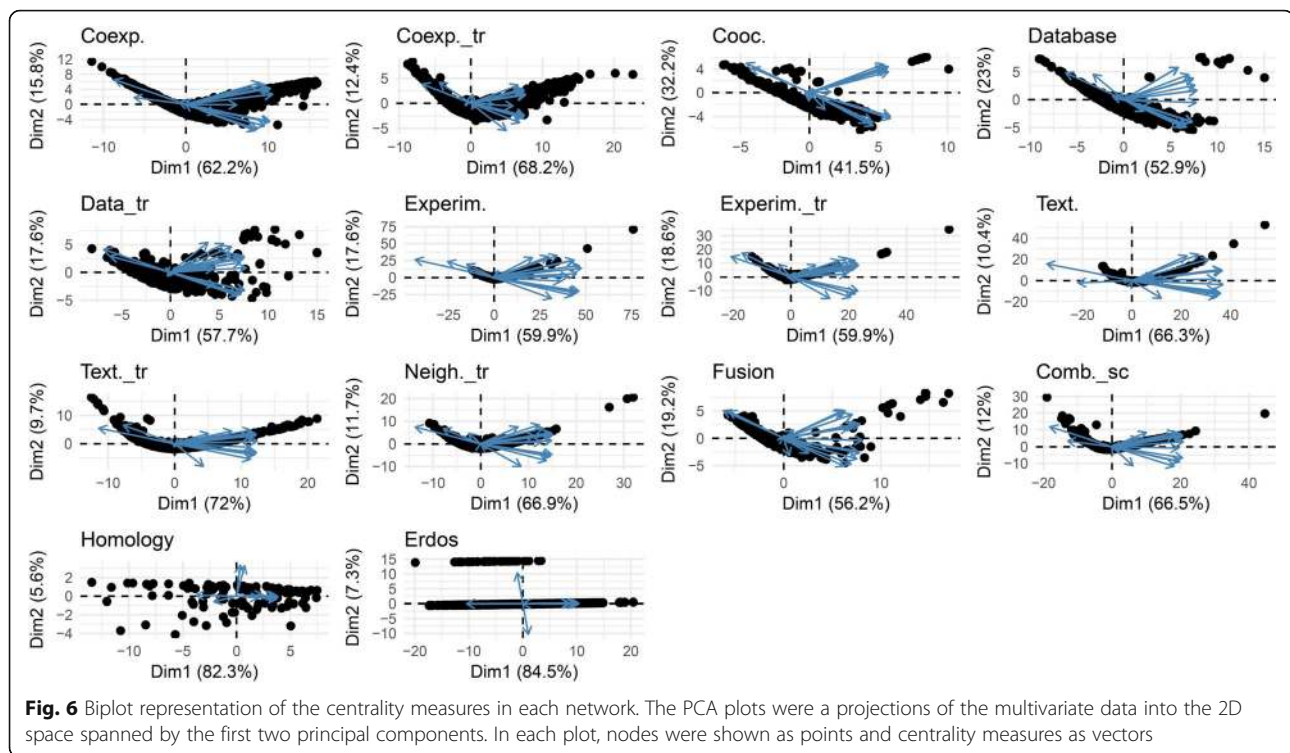


Fig. 5 Pairwise scatterplot between the centrality measures. This figure contains combined-score PPIN and the null network. In this figure, the r Pearson correlation coefficients between centralities beside the centralities distribution were also presented in both networks. For better representation, red and black colors were used and the scatterplot was divided into three parts corresponding to Table 1 groups. For the scatterplot visualizations of all PPINs see Additional file 2



silhouette value mean (i.e. 0.03) much lower than the average. In all other PPINs, the same relationship between silhouette scores and contribution values was observed as shown in Additional files 4 and 7.

Our results demonstrated that a unique profile of centrality measures including Latora closeness, Barycenter, Diffusion degree, Freeman closeness, Residual, Average distance, Radiality centralities, was the most significant indicator in ranking PPIN nodes. We inferred that the rationale and logic of network reconstruction dictates which centrality measures should be chosen. Also, we demonstrated the relationship between contribution value derived from PCA and silhouette width as a cluster validity index. Regarding to the robustness issue, we first reasserted that the architecture and global properties of a network impact on the centrality analysis results [73–75]. Therefore, the center of a network would be different, depending on the network's inherent topology. In other words, we addressed this issue whether a given centrality measure has enough information via-a-vis and it demonstrates a same behavior in some other networks.

Conclusion

Network-based methods have been introduced as an emergent approach for simplification, reconstruction, analysis, and comprehension of complex behavior in biological systems. Network-based ranking methods (i.e. centrality analysis) have been found widespread use for predicting essential proteins, proposing drug targets

candidates in treatment of cancer, biomarker discovery, human disease genes identification and creation a cell with the minimal genome [76]. However, there is no consensus pipeline for centrality analysis regarding aforementioned applications among network analysts.

In this study, we worked on yeast PPINs which were built using 13 evidence channels in the STRING database. Subsequently, 27 centrality measures were used for the prioritization of the nodes in all PPINs. We illustrated that data reduction and low-dimensional projection help to extract relevant features (i.e. centrality measures) and corresponding relationships. Thus, to quantify connectivity in biological networks, we recommend that before arbitrary picking centrality measures to pinpoint important nodes, PCA (as an example of data projection methods) conduce how to use these measures. In the other word, the analysis of principal components clarifies which measures have the highest contribution values, i.e., which measures comprise much more information about centrality. Freshly, the application of these approach for discovering essential proteins was assayed in a polypharmacology study to prevent epithelial-mesenchymal transition in cancer [77].

Basic definitions

- **Giant component of a graph** defines the largest connected component of a graph in which there is a path between each pair of nodes [78].

Table 5 Clustering information values for PPINs. The Hopkin's statistics threshold for clusterability was 0.05

Network	Hopkins Statistic	Number of Clusters	Silhouette Average Value
Coexpression	0.25	6	0.36
Coexpression_transferred	0.21	7	0.33
Cooccurrence	0.18	6	0.55
Database	0.24	6	0.33
Database_transferred	0.20	9	0.32
Experiments	0.21	9	0.31
Experiments_transferred	0.16	6	0.43
Textmining	0.24	8	0.28
Textmining_transferred	0.20	6	0.35
Neighborhood_transferred	0.26	2	0.39
Fussion	0.16	5	0.48
Combined_score	0.30	7	0.27
Homology	0.23	2	0.46

- **Network density** is a representation of the number of interactions to the number of possible interactions among a given network [79].
- **Network centralization** refers to a topological spectrum from star to grid topologies (where each node has a same number of links) of a graph varies from 1 to 0 [79].
- The **network heterogeneity** measure describes as the coefficient of variation of connectivity distribution. A high heterogeneous network implies that the network is exhibited approximate scale-free topology [79, 80].
- The **clustering coefficient** of a node is the number of triangles (3-loops) that pass through it, relative to the maximum number of 3-loops that could pass through the node. The network clustering coefficient

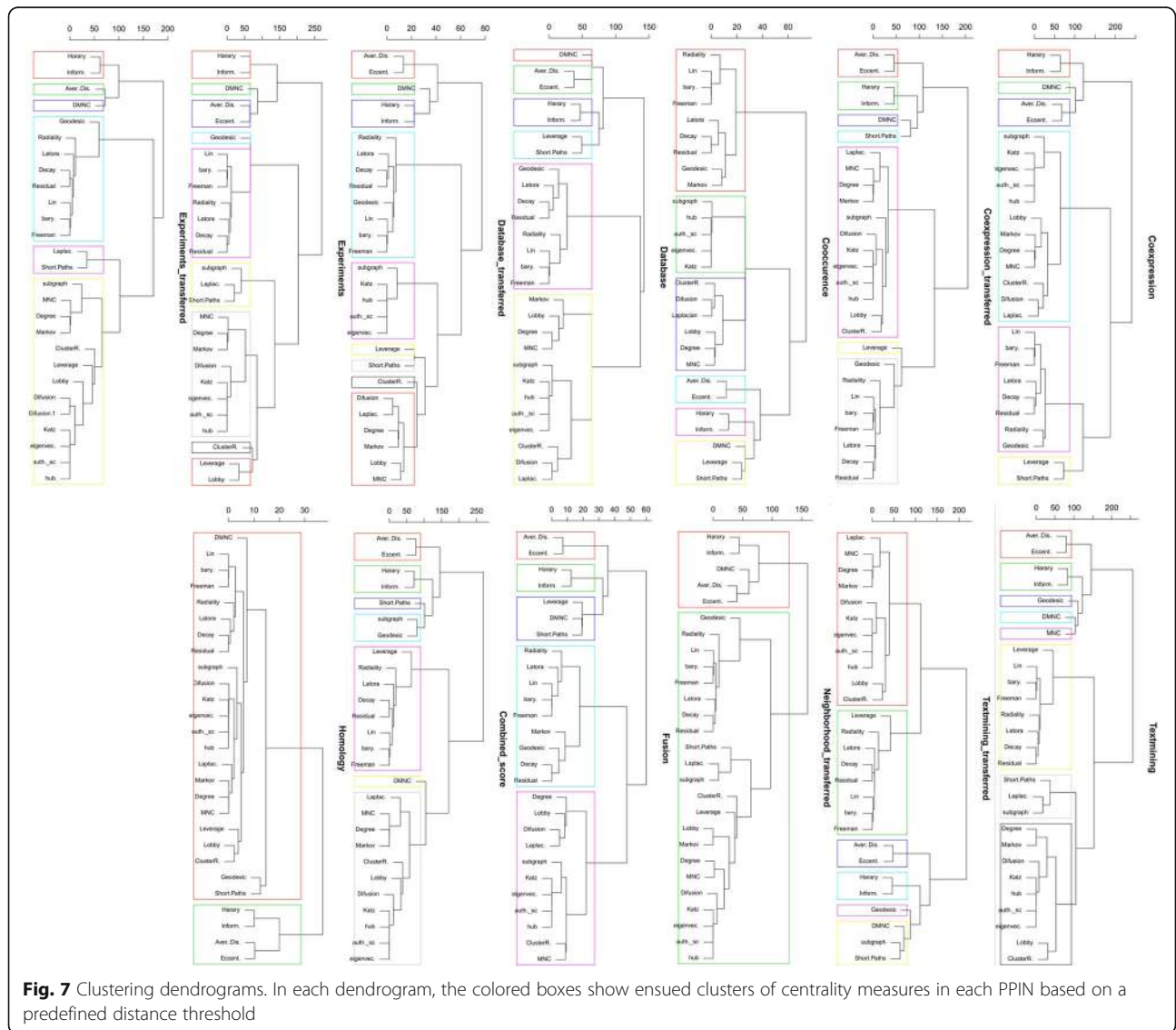


Fig. 7 Clustering dendrograms. In each dendrogram, the colored boxes show ensued clusters of centrality measures in each PPIN based on a predefined distance threshold

Table 6 Jaccard index coefficient values for PPINs. The values represent how similar the networks are, in terms of their clustering results. A value of 1 indicates an exact match while values equal to 0 show dissimilarity

	coexp.	coexp_tr	coocc.	comb.	dat_tr	dat.	exp.	exp_tr	fus.	hom.	nei_tr	tex.	tex_tr
coexpression		0.99	0.58	0.77	0.62	1.00	0.58	0.80	0.83	0.41	0.43	0.62	0.76
coexpression_transferred			0.57	0.78	0.63	0.99	0.58	0.81	0.82	0.40	0.43	0.62	0.77
cooccurrence				0.47	0.75	0.58	0.44	0.50	0.73	0.29	0.30	0.43	0.48
combined_score					0.52	0.77	0.62	0.63	0.64	0.37	0.39	0.78	0.96
database_transferred						0.62	0.55	0.55	0.55	0.25	0.27	0.47	0.51
database							0.58	0.80	0.83	0.41	0.43	0.62	0.76
experiments								0.59	0.49	0.25	0.27	0.67	0.63
experiments_transferred									0.67	0.41	0.43	0.62	0.62
fussion										0.40	0.42	0.52	0.64
homology											0.91	0.30	0.37
neighborhood_transferred												0.32	0.39
textmining													0.78
textmining_transferred													

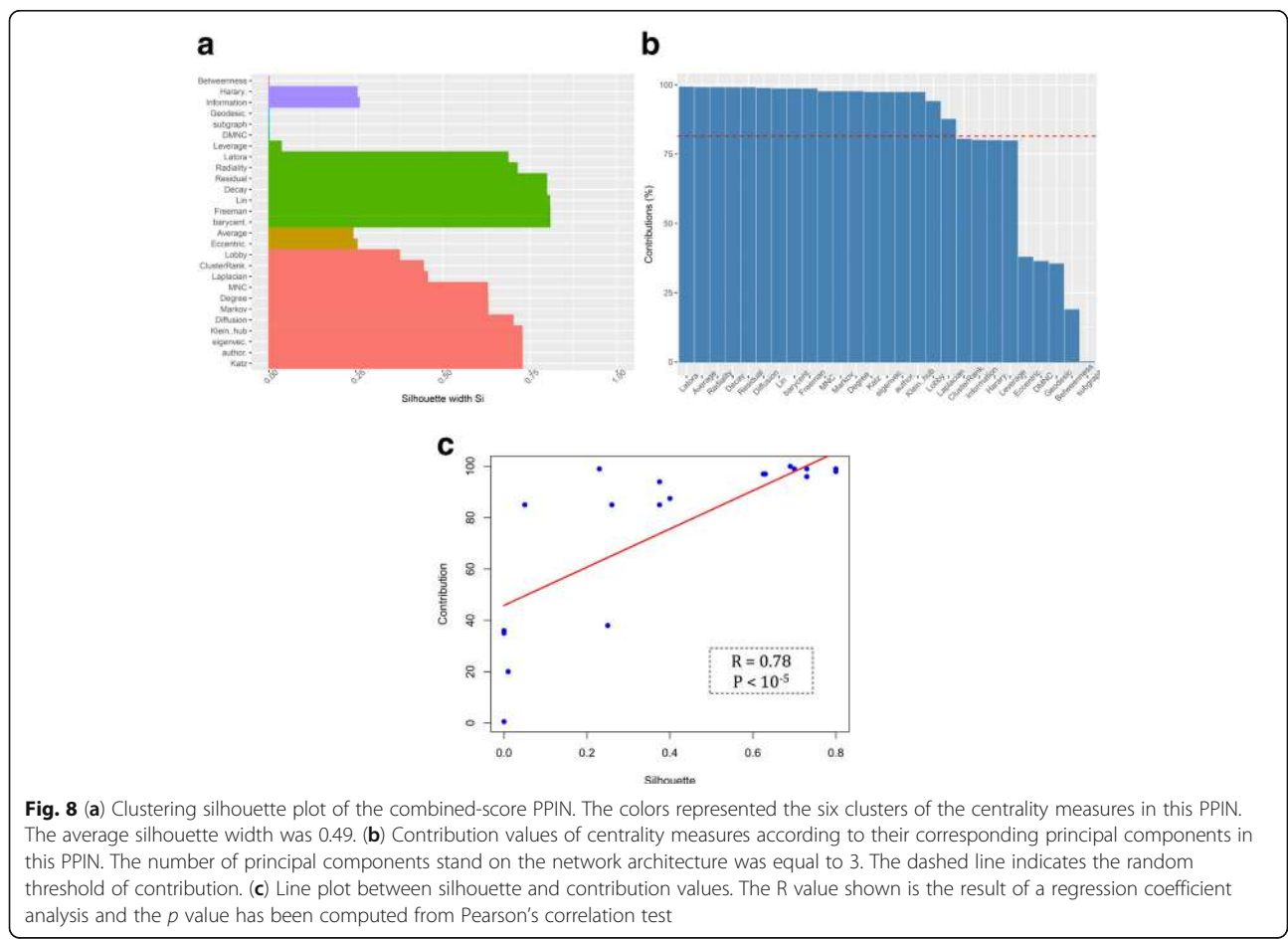


Fig. 8 (a) Clustering silhouette plot of the combined-score PPIN. The colors represented the six clusters of the centrality measures in this PPIN. The average silhouette width was 0.49. (b) Contribution values of centrality measures according to their corresponding principal components in this PPIN. The number of principal components stand on the network architecture was equal to 3. The dashed line indicates the random threshold of contribution. (c) Line plot between silhouette and contribution values. The R value shown is the result of a regression coefficient analysis and the p value has been computed from Pearson's correlation test

defines as the mean of the clustering coefficients for all nodes in the network [81, 82].

- **Influential nodes** which is generally used in social networks analysis point as nodes with good spreading properties in networks [83]. Different centrality measures are used to find influential nodes.
- **Centrality-lethality rule** explains nodes with high centrality values in which maintain the integrity of the network structure, are more related to the survival of the biological system [84].
- The **silhouette criterion** defines how similar a centrality is to its own cluster compared to other clusters. It ranges from -1 to 1 , where a high value infers that the centrality is well matched to its own cluster and poorly matched to neighboring clusters. If most centralities have a high value, then the clustering configuration is proper. If they have low or negative values, then the clustering configuration may have too many or too few clusters [5, 85].

In order to see definitions of all used centrality measures, see <http://www.centiserver.org>.

Additional files

Additional file 1: Evidence channel dataset. The contents of 13 evidence channels illustrating the yeast PPIN from STRING database. (downloaded in 3rd Sep 2016) are provided. (TXT 28629 kb)

Additional file 2: Fitted power law distribution. The Degree distribution of each network has been compared to the power law distribution in order to visualize the scale free property in the structure of each network. (PDF 203 kb)

Additional file 3: Scatterplots between groups of centralities. Each panel indicates scatterplots between centralities groups of two networks. (PPTX 1963 kb)

Additional file 4: Contribution values of centralities in each network. These values were computed based on the principal components. The red line shows the threshold used for identifying effective centralities. (PDF 441 kb)

Additional file 5: Visual assessment of cluster tendency plots. Each rectangular represents the clusters of the calculated results of the centrality measures. (PDF 313 kb)

Additional file 6: Clustering properties results. These properties include connectivity, Dunn and Silhouette scores. These scores suggest the sufficient clustering method by a specific number of clusters. (DOCX 16 kb)

Additional file 7: Clusters silhouette plots. Each color represents a cluster and each bar with specific color indicates a centrality. (PDF 417 kb)

Additional file 8: Optimal number of clusters. The suitable number of clusters for hierarchical clustering method was computed using the average silhouette values. (PDF 321 kb)

Abbreviations

DMNC: Density of Maximum Neighborhood Component; MNC: Maximum Neighborhood Component; PAM: Partitioning Around Medoids; PCA: Principal Component Analysis; PPIN: Protein-protein interaction network; VAT: Visual Assessment of cluster Tendency

Authors' contributions

MM and MJ participated in the design of the study. MA carried out acquisition, computational analysis and interpretation of the data and drafted the manuscript. MJ helped computational analysis and in the writing of the draft.

ASY, ZRM, HH and OW conceived of the study and participated in revising the manuscript critically. All authors read and approved the final manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Drug Design and Bioinformatics Unit, Medical Biotechnology Department, Biotechnology Research Center, Pasteur Institute of Iran, P.O. Box 13164, Tehran, Iran. ²Department of Systems Biology and Bioinformatics, University of Rostock, P.O. Box 18051, Rostock, Germany. ³Faculty of New Sciences and Technologies, University of Tehran, P.O. Box 143995-71, Tehran, Iran. ⁴Max-Planck Institute of Molecular Plant Physiology, Potsdam, Germany. ⁵Department of Applied Mathematics, Faculty of Mathematical Sciences, Tarbiat Modares University, P.O. Box 14115-134, Tehran, Iran.

Received: 23 November 2017 Accepted: 22 June 2018

Published online: 31 July 2018

References

1. Jeong H, Mason SP, Barabasi A-L, Oltvai ZN. Lethality and centrality in protein networks. 2001. arXiv preprint cond-mat/0105306.
2. Csérmelyi P, Korcsmáros T, Kiss HJ, London G, Nussinov R. Structure and dynamics of molecular networks: a novel paradigm of drug discovery: a comprehensive review. *Pharmacol Ther.* 2013;138(3):333–408.
3. Azimzadeh Jamalkandi S, Mozghani S-H, Gholami Pourbadie H, Mirzaie M, Noorbakhsh F, Vaziri B, Gholami A, Ansari-Pour N, Jafari M. Systems biomedicine of rabies delineates the affected signaling pathways. *Front Microbiol.* 2016;7:1688.
4. Azimzadeh S, Mirzaie M, Jafari M, Mehrani H, Shariati P, Khodabandeh M. Signaling network of lipids as a comprehensive scaffold for omics data integration in sputum of COPD patients. *Biochim Biophys Acta, Mol Cell Biol Lipids.* 2015;1851(10):1383–93.
5. Jafari M, Mirzaie M, Sadeghi M, Marashi S-A, Rezaei-Tavirani M. Exploring biological processes involved in embryonic stem cell differentiation by analyzing proteomic data. *Biochim Biophys Acta, Proteins Proteomics.* 2013;1834(6):1063–9.
6. Rezadoost H, Karimi M, Jafari M. Proteomics of hot-wet and cold-dry temperaments proposed in Iranian traditional medicine: a network-based study. *Sci Rep.* 2016;6(130):30133.
7. Freeman LC. Going the wrong way on a one-way street: centrality in physics and biology. *J Soc Structure.* 2008;9(2):1–15.
8. Landherr A, Friedl B, Heidemann J. A critical review of centrality measures in social networks. *Bus Inform Syst Eng.* 2010;2(6):371–85.
9. Freeman LC. Centrality in social networks conceptual clarification. *Soc Networks.* 1978;1(3):215–39.
10. Buttner K, Scheffler K, Czycholl I, Krieter J. Social network analysis - centrality parameters and individual network positions of agonistic behavior in pigs over three different age levels. *Springerplus.* 2015;4:185.
11. Jalili M, Salehzadeh-Yazdi A, Asgari Y, Arab SS, Yaghmaie M, Ghavamzadeh A, Alimoghaddam K. CentiServer: a comprehensive resource, web-based application and R package for centrality analysis. *PLoS One.* 2015;10(11):e0143111.
12. Hahn MW, Kern AD. Comparative genomics of centrality and essentiality in three eukaryotic protein-interaction networks. *Mol Biol Evol.* 2005;22(4):803–6.
13. Bergmann S, Ihmels J, Barkai N. Similarities and differences in genome-wide expression data of six organisms. *PLoS Biol.* 2004;2(1):E9.
14. Wagner A, Fell DA. The small world inside large metabolic networks. *Proc Biol Sci / R Soc.* 2001;268(1478):1803–10.
15. Ma HW, Zeng AP. The connectivity structure, giant strong component and centrality of metabolic networks. *Bioinformatics.* 2003;19(11):1423–30.

16. Estrada E. Virtual identification of essential proteins within the protein interaction network of yeast. *Proteomics*. 2006;6(1):35–40.
17. He X, Zhang J. Why do hubs tend to be essential in protein networks? *PLoS Genet*. 2006;2(6):e88.
18. Joy MP, Brock A, Ingber DE, Huang S. High-betweenness proteins in the yeast protein interaction network. *J Biomed Biotechnol*. 2005;2005(2):96–103.
19. Potapov AP, Voss N, Sasse N, Wingender E. Topology of mammalian transcription networks. *Genome Inform Int Confer Genome Inform*. 2005; 16(2):270–8.
20. Yang Y, Han L, Yuan Y, Li J, Hei N, Liang H. Gene co-expression network analysis reveals common system-level properties of prognostic genes across cancer types. *Nat Commun*. 2014;5:3231. <https://doi.org/10.1038/ncomms4231>.
21. Tew KL, Li XL, Tan SH. Functional centrality: detecting lethality of proteins in protein interaction networks. *Genome Inform Ser*. 2007;19:166–77.
22. Peng X, Wang J, Wang J, Wu FX, Pan Y. Rechecking the centrality-lethality rule in the scope of protein subcellular localization interaction networks. *PLoS One*. 2015;10(6):e0130743.
23. Luo J, Qi Y. Identification of essential proteins based on a new combination of local interaction density and protein complexes. *PLoS One*. 2015;10(6): e0131418.
24. Tang X, Wang J, Zhong J, Pan Y. Predicting essential proteins based on weighted degree centrality. *IEEE/ACM Trans Comput Biol Bioinform*. 2014; 11(2):407–18.
25. Li M, Ni P, Chen X, Wang J, Wu F, Pan Y. Construction of refined protein interaction network for predicting essential proteins. *IEEE/ACM Trans Comput Biol Bioinform*. 2017; <https://doi.org/10.1109/TCBB.2017.2665482>.
26. Li M, Li W, Wu F-X, Pan Y, Wang J. Identifying essential proteins based on sub-network partition and prioritization by integrating subcellular localization information. *J Theor Biol*. 2018;447:65–73.
27. Li M, Lu Y, Niu Z, Wu F-X. United complex centrality for identification of essential proteins from PPI networks. *IEEE/ACM Trans Comput Biol Bioinform*. 2017;14(2):370–80.
28. Tang Y, Li M, Wang J, Pan Y, Wu F-X. CytoNCA: a cytoscape plugin for centrality analysis and evaluation of protein interaction networks. *Biosystems*. 2015;127:67–72.
29. Khuri S, Wuchty S. Essentiality and centrality in protein interaction networks revisited. *BMC Bioinformatics*. 2015;16(1):1.
30. Koschützki D, Schreiber F. Comparison of centralities for biological networks, in 'Proceedings of the German Conference on Bioinformatics 2004', Vol. P-53 of Lecture Notes in Informatics, Springer; 2004. pp. 199–206.
31. Dwyer T, Hong SH, Koschützki D, Schreiber F, Xu K. Visual analysis of network centralities. In Proceedings of the 2006 Asia-Pacific Symposium on Information Visualisation-Volume 60 (pp. 189-197). Australian Computer Society, Inc. Proceeding APVis '06 Proceedings of the 2006 Asia-Pacific Symposium on Information Visualisation - Volume 60 Pages 189-197, Tokyo, Japan, Australian Computer Society, Inc. Darlinghurst, Australia, Australia; 2006. ISBN: 1-920682-41-4.
32. Valente TW, Coronges K, Lakon C, Costenbader E. How correlated are network centrality measures? *Connections* (Toronto, Ont). 2008;28(1):16.
33. Batool K, Niazi MA. Towards a methodology for validation of centrality measures in complex networks. *PLoS one*. 2014;9(4):e90283. <https://doi.org/10.1371/journal.pone.0090283>.
34. Li C, Li Q, Van Mieghem P, Stanley HE, Wang H. Correlation between centrality metrics and their application to the opinion model. *Eur Phys J B*. 2015;88(3):65.
35. Jalili M, Salehzadeh-Yazdi A, Gupta S, Wolkenhauer O, Yaghmaie M, Resendis-Antonio O, Alimoghaddam K. Evolution of centrality measurements for the detection of essential proteins in biological networks. *Front Physiol*. 2016;7:375.
36. Boutet E, et al. UniProtKB/Swiss-Prot, the manually annotated section of the UniProt KnowledgeBase: how to use the entry view. In: Edwards D, editor. *Plant Bioinformatics. Methods in Molecular Biology*. New York: Humana Press; 2016;1374.
37. Szklarczyk D, Morris JH, Cook H, Kuhn M, Wyder S, Simonovic M, Santos A, Doncheva NT, Roth A, Bork P, Jensen LJ, von Mering C. The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res*. 2017;45(D1):D362–8. <https://doi.org/10.1093/nar/gkw937>.
38. Erdos P, Rényi A. On the evolution of random graphs. *Publ Math Inst Hung Acad Sci*. 1960;5(1):17–60.
39. Csardi G, Nepusz T. The igraph software package for complex network research. *Int J Complex Sys*. 2006;1695(5):1–9.
40. Butts CT. Network: a package for managing relational data in R. *J Stat Softw*. 2008;24(2):1–36.
41. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinform*. 2008;9(1):1.
42. *Weighted Network Analysis: Applications in Genomics and Systems Biology*, Steve Horvath, Springer Science & Business Media, Ordibehesht 10, 1390 AP - Science. p. 421.
43. del Rio G, Koschützki D, Coello G. How to identify essential genes from molecular networks? *BMC Syst Biol*. 2009;3(1):1.
44. Viswanath M. Ontology-based automatic text summarization (Doctoral dissertation, uga). 2009.
45. Latora V, Marchiori M. Efficient behavior of small-world networks. *Phys Rev Lett*. 2001;87(19):198701.
46. Dangalchev C. Residual closeness in networks. *Physica A: Stat Mechanics Appl*. 2006;365(2):556–64.
47. Chen D-B, Gao H, Lü L, Zhou T. Identifying influential nodes in large-scale directed networks: the role of clustering. *PLoS One*. 2013;8(10):e77455.
48. Jackson MO. Representing and Measuring Networks. *Social and economic networks*; 2008. pp. 37-43.
49. Kundu S, Murthy CA, Pal SK. A new centrality measure for influence maximization in social networks. *Lect Notes Comput Sc*. 2011;6744:242–7.
50. Lin C-Y, Chin C-H, Wu H-H, Chen S-H, Ho C-W, Ko M-T. Hubba: hub objects analyzer—a framework of interactome hubs identification for network biology. *Nucleic Acids Res*. 2008;36(suppl 2):W438–43.
51. Borgatti SP, Everett MG. A graph-theoretic perspective on centrality. *Soc Networks*. 2006;28(4):466–84.
52. De Meo P, Ferrara E, Fiumara G, Ricciardello A. A novel measure of edge centrality in social networks. *Knowl-Based Syst*. 2012;30:136–50.
53. Grassler J, Koschützki D, Schreiber F. CentiLib: comprehensive analysis and exploration of network centralities. *Bioinformatics*. 2012;28(8):1178–9.
54. Junker BH, Koschützki D, Schreiber F. Exploration of biological network centralities with CentiBIN. *BMC Bioinformatics*. 2006;7(1):1.
55. Qi X, Fuller E, Wu Q, Wu Y, Zhang C-Q. Laplacian centrality: a new centrality measure for weighted networks. *Inf Sci*. 2012;194:240–53.
56. Joyce KE, Laurienti PJ, Burdette JH, Hayasaka S. A new measure of centrality for brain networks. *PLoS One*. 2010;5(8):e12200.
57. Hoffman AN, Stearns TM, Shrader CB. Structure, context, and centrality in interorganizational networks. *J Bus Res*. 1990;20(4):333–47.
58. Korn A, Schubert A, Telcs A. Lobby index in networks. *Physica A: Stat Mechanics Appl*. 2009;388(11):2221–6.
59. White S, Smyth P. Algorithms for estimating relative importance in networks. In: Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM. Published in: Proceeding KDD '03 Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining. p. 266-275. Washington, D.C. — August 24 - 27, 2003 ACM New York, NY, USA; 2003. ISBN: 1-58113-737-0.
60. Zotenko E, Mestre J, O'Leary DP, Przytycka TM. Why do hubs in the yeast protein interaction network tend to be essential: reexamining the connection between the network topology and essentiality. *PLoS Comput Biol*. 2008;4(8):e1000140.
61. Bonacich P. Power and centrality: a family of measures. *Am J Sociol*. 1987; 92(5):1170–82.
62. Estrada E, Rodriguez-Velazquez JA. Subgraph centrality in complex networks. *Phys Rev E*. 2005;71(5):056103.
63. Hage P, Harary F. Eccentricity and centrality in networks. *Soc Networks*. 1995;17(1):57–63.
64. Kleinberg JM. Authoritative sources in a hyperlinked environment. *JACM*. 1999;46(5):604–32.
65. Stephenson K, Zelen M. Rethinking centrality: methods and examples. *Soc Networks*. 1989;11(1):1–37.
66. Butts CT. sna: Tools for Social Network Analysis. R package version 2.2-0. 2010.
67. Becker RA, Chambers JM, Wilks AR. The new S language, vol. 1988. Pacific Grove: Brooks; 1988.
68. Abdi H, Williams LJ. Principal component analysis. *Wiley Interdiscip Rev Comput Stat*. 2010;2(4):433–59.
69. Kassambara A. Factoextra: visualization of the outputs of a multivariate analysis. R Package version. 2015;1(1):1–75.

70. Brock G, Pihur V, Datta S, Datta S. cValid, an R package for cluster validation. *J Stat Software* (Brock et al., March 2008). 2011.
71. Gobbi A, Albanese D, Iorio F: Package 'BiRewire'. 2016.
72. Ward JH Jr. Hierarchical grouping to optimize an objective function. *J Am Stat Assoc.* 1963;58(301):236–44.
73. Tsugawa S, Matsumoto Y, Ohsaki H. On the robustness of centrality measures against link weight quantization in social networks. *Computat Math Org Theory.* 2015;21(3):318–39.
74. Niu QK, Zeng A, Fan Y, Di ZR. Robustness of centrality measures against network manipulation. *Physica A.* 2015;438:124–31.
75. Tsugawa S, Matsumoto Y, Ohsaki H. On the robustness of centrality measures against link weight quantization in social networks. *Comput Math Organ Th.* 2015;21(3):318–39.
76. Glass JI, Hutchison CA 3rd, Smith HO, Venter JC. A systems biology tour de force for a near-minimal bacterium. *Mol Syst Biol.* 2009;5:330.
77. Barneh F, Mirzaie M, Nickchi P, Tan TZ, Thiery JP, Piran M, Salimi M, Goshadrou F, Aref AR, Jafari M. Integrated use of bioinformatic resources reveals that co-targeting of histone deacetylases, IKBK and SRC inhibits epithelial-mesenchymal transition in cancer. *Brief Bioinform.* bby030. <https://doi.org/10.1093/bib/bby030>.
78. Barneh F, Jafari M, Mirzaie M. Updates on drug-target network; facilitating polypharmacology and data integration by growth of DrugBank database. *Brief Bioinform.* 2016;17(6):1070–80.
79. Horvath S. *Weighted network analysis.* New York: Springer New York; 2011.
80. Jafari M, Sadeghi M, Mirzaie M, Marashi S-a, Rezaei-Tavirani M. Evolutionarily conserved motifs and modules in mitochondrial protein-protein interaction networks. *Mitochondrion.* 2013;13(6):668–75.
81. Junker BH, Schreiber F. *Analysis of biological networks.* Wiley; 2008.
82. Jafari M, Mirzaie M, Sadeghi M. Interlog protein network: an evolutionary benchmark of protein interaction networks for the evaluation of clustering algorithms. *BMC Bioinformatics.* 2015;16(1):319.
83. Malliaros FD, Rossi M-EG, Vazirgiannis M. Locating influential nodes in complex networks. *Sci Rep.* 2016;6(1):19307.
84. Jeong H, Mason SP, Barabasi AL, Oltvai ZN. Lethality and centrality in protein networks. *Nature.* 2001;411(6833):41–2.
85. Rousseeuw PJ. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *J Comput Appl Math.* 1987;20:53–65.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

