

Feature Extraction from Degree Distribution for Comparison and Analysis of Complex Networks

Sadegh Aliakbary, Jafar Habibi, and Ali Movaghar

Department of Computer Engineering, Sharif University of Technology, Tehran, Iran
aliakbary@ce.sharif.edu, jhabibi@sharif.edu, movaghar@sharif.edu

July 27, 2021

Abstract

The degree distribution is an important characteristic of complex networks. In many data analysis applications, the networks should be represented as fixed-length feature vectors and therefore the feature extraction from the degree distribution is a necessary step. Moreover, many applications need a similarity function for comparison of complex networks based on their degree distributions. Such a similarity measure has many applications including classification and clustering of network instances, evaluation of network sampling methods, anomaly detection, and study of epidemic dynamics. The existing methods are unable to effectively capture the similarity of degree distributions, particularly when the corresponding networks have different sizes. Based on our observations about the structure of the degree distributions in networks over time, we propose a feature extraction and a similarity function for the degree distributions in complex networks. We propose to calculate the feature values based on the mean and standard deviation of the node degrees in order to decrease the effect of the network size on the extracted features. The proposed method is evaluated using different artificial and real network datasets, and it outperforms the state of the art methods with respect to the accuracy of the distance function and the effectiveness of the extracted features.

1 Introduction

Degree distribution is an important and informative characteristic of a complex network [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17]. Although the degree distribution does not capture all aspects of the topology of a network [18], it reflects the overall pattern of connections [14] and is an important determinant of network properties [15]. Degree distribution is also a sign of link formation process in the network. The degree distribution of complex networks often follows a heavy-tailed distribution [6, 7, 8, 9, 19], but networks still show different characteristics in their degree distributions. Hence, we frequently need an appropriate similarity function in order to compare the degree distribution of the complex networks. Such a similarity function plays an important role in many network analysis applications, such as evaluation of network generation models [1, 2, 20, 5, 7, 16], evaluation of sampling methods [9, 21, 11, 22, 18], generative model selection [23, 24, 25], classification or clustering of network instances [3, 12, 24, 25, 23], anomaly detection [26, 27], and study of epidemic dynamics [28, 29, 30]. For example, in evaluation of sampling algorithms, the given network instance is compared with its sampled counterpart in order to ensure that the structure of the degree distribution is preserved [9, 21, 11, 22, 18, 12].

In addition to the need for network comparison, representing the network as a fixed-size feature vector is also an important step in every data analysis process [23, 24, 25]. In order to employ the degree distribution in such applications, a procedure is needed for extracting

a feature vector from the degree distribution. The extracted feature vector is also useful in developing a distance function for comparing two degree distributions. Although there exist accepted quantification methods for many network features (e.g., average clustering coefficient, modularity, and average shortest path length), the quantification and comparison of the degree distribution is not a trivial task. We will show that the existing methods have major weaknesses in feature extraction and comparison of networks, particularly when the considered networks have different sizes. Even if no network comparison is required, feature extraction from the degree distributions has independent applications. For example, data-analysis algorithms require the network features, including the degree distribution, to be represented as some real numbers in the form of a fixed-length feature vector [23, 12, 24, 25].

According to the mentioned applications for the demanded similarity function, we regard two degree distributions similar if their networks follow similar structure, similar connection patterns and similar link formation processes, even if the networks are of completely different scales and sizes. The state of the art approaches for comparing degree distributions are eyeballing the distribution diagrams (usually, to satisfy a heavy-tailed distribution) [2, 7, 13], Kolmogorov-Smirnov (KS) test [6, 8, 9, 31, 32, 33], comparison based on fitted power-law exponent [13, 34], and comparison based on distribution percentiles [24]. Eyeballing is obviously an inaccurate, error-prone and manual task. Comparison based on power-law exponent is based on the assumption that the degree distributions obeys a power-law, which is invalid for many complex networks [10, 35, 36, 37]. KS-test is based on a point-to-point comparison of the cumulative distribution function, which is not a good approach for comparing networks with different ranges of node degrees. Percentile method is too sensitive to the outlier values of node degrees. As a result, the existing methods are actually inappropriate for comparing the degree distribution of networks.

In order to reveal the limitations of the existing methods and the main ideas of our proposed method, Figure 1 and Figure 2 provide intuitive examples of networks with different sizes over time. Figure 1 illustrates the degree distribution of two citation networks and two collaboration (co-authorship) networks which are extracted from CiteSeerX digital library [38]. The networks represent two snapshots in the years 1992 and 2010. For example, *Citation_1992* represents the graph of citations in the papers of the CiteSeerX repository which are published before 1992. The figure shows two similar distributions for the two citation networks (resembling a power-law distribution) and two similar distributions for the two collaboration networks (similar to log-normal model). Obviously, the degree distributions of the citation networks are dissimilar to those of the collaboration networks. The existing similarity functions are unable to capture such similarity and dissimilarity in the shape of the degree distributions appropriately. For example, the Kolmogorov-Smirnov test will return the least similarity between the two citation networks, the power-law exponents are uninformatively different for the four networks, and the percentile quantification method [24] returns an equal vector for all the four networks. Although the shape of the degree distributions in this example clearly reveals the similarity of the networks, in many situations the similarity of the degree distributions are not that obvious. Therefore, we need an automatic method for feature extraction and/or quantified network comparison using the degree distributions. It is worth noting that many real networks do not follow a specific distribution model such as the power-law or log-normal models. As a result, fitting the degree distribution to some predefined distribution models is not a good approach for feature extraction or comparison of the networks.

As another example, Figure 2 shows the degree distribution of the citation networks, extracted from the same CiteSeerX repository for different snapshots from 1990 to 2010. These networks have become bigger over years, with 191,443 nodes (papers) in 1990 and 1,039,952 nodes in 2010. Although all the networks show similar degree distributions, the bigger networks

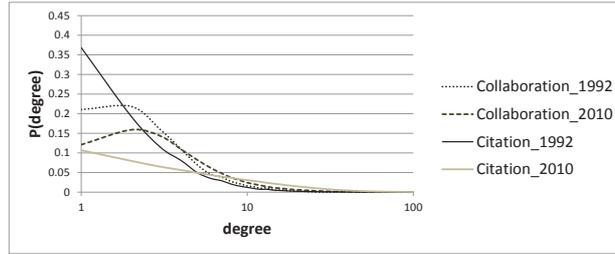


Figure 1: Degree distribution of two citation networks and two collaboration networks

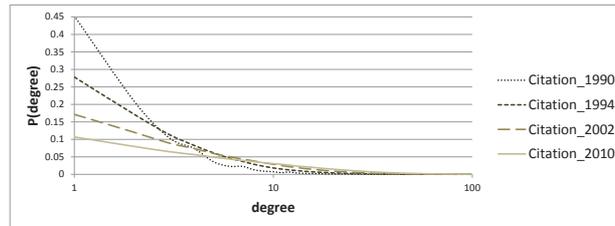


Figure 2: Degree distribution of a citation network over different times

contain a wider range of node degrees. As the network size increases, maximum node degree, average degree, and standard deviation of the node degrees also increase. As a result, it seems that a degree distribution should be normalized according to its mean and deviation so that the comparison of the networks with different scales become valid. We follow this idea in our proposed methods of feature extraction and comparison of the network degree distributions. In our proposed “feature extraction” method, a fixed-length feature vector is extracted from the degree distribution, which can be used in data analysis applications, data-mining algorithms and comparison of degree distributions. In the “comparison method”, we propose a distance function that computes the distance (amount of dissimilarity) between two given network degree distributions. With such a distance function, we can figure out how similar the given networks are, according to their degree distributions.

Although our proposed approach is of general nature and equally applicable to other network types, in this paper we focus on simple undirected networks. In the rest of this paper, our proposed method is called “Degree Distribution Quantification and Comparison (DDQC)”. The rest of this paper is organized as follows: In section 2, we briefly overview the related works. In section 3, we propose a new method for degree distribution quantification and comparison. In section 4, we evaluate the proposed method and we compare it with baseline methods. Finally, we conclude the paper in section 5.

2 Related Works

The degree distribution of many real-world networks are heavy tailed [6, 7, 8, 9, 19], and the power-law distribution is the most suggested model for complex networks[4, 6, 8]. In power-law degree distribution the number of nodes with degree d is proportional to $d^{-\gamma}$ ($N_d \propto d^{-\gamma}$) where γ is a positive number called “the power-law exponent”. The value of γ is typically in the range $2 < \gamma < 3$ [2, 39, 36]. The fitted power-law exponent can be used to characterize graphs [34]. A common approach for feature extraction from the degree distribution is to fit it on a power-law model and estimate the power-law exponent (γ). As a result, it will be possible to compare networks according to their fitted power-law exponents. One of the drawbacks of this

approach is that power-law exponent is too limited to represent a whole degree distribution. This approach also follows the assumption that the degree distribution is power-law, which is not always valid, because many networks follow other degree distribution models such as log-normal distribution [10, 35, 36, 37]. In addition, the power-law exponent does not reflect the deviation of the degree distribution from the fitted power-law distribution. As a result, two completely different distributions may have similar quantified feature (fitted power-law exponent).

An alternative approach for comparing the degree distributions is to utilize statistical methods of probability distribution comparison. Degree distribution is a kind of probability distribution and there are a variety of measures for calculating the distance between two probability distributions. In this context, the most common method is the Kolmogorov-Smirnov (KS) test, which is defined as the maximum distance between the cumulative distribution functions (CDF) of the two probability distributions [6]. KS-test is used for comparing two degree distributions (two-sample KS test) [9, 31] and also for comparing a degree distribution with a baseline (usually the power-law) distribution [6, 10, 40, 41]. The KS distance of two distributions is calculated according to Equation 1, in which $S_1(d)$ and $S_2(d)$ are the CDFs of the two degree distributions, and d indicates the node degree. KS-test is largely utilized in the literature for comparing degree distribution of complex networks [6, 8, 9, 31, 32, 33]. KS-test is a method for comparing the degree distributions and calculating their distance, and it does not provide feature extraction mechanism. As a result, we should maintain the CDF of the degree distributions so that we can compare them according to KS-test. This is a drawback of KS-test since in other investigated approaches the degree distribution is summarized in a small fixed-length feature vector. Additionally, the KS-test is not applicable in data analysis applications that rely on feature vector representation of networks. KS-test is also sensitive to the scale and size of the networks, since it performs a point-to-point comparison of CDFs. Therefore, for two networks with different ranges of node degrees, the KS-test may return a large value as their distance even if the overall views of the degree distributions are similar (refer to Figure 1).

$$distance_{KS}(S_1, S_2) = \max_d |S_1(d) - S_2(d)| \quad (1)$$

Janssen et. al., [24] propose an alternative method for feature extraction from the degree distribution. In this method, the degree distribution is divided into eight equal-size regions and the sum of degree probabilities in each region is extracted as distribution percentiles. This method is sensitive to the range of node degrees and also to outlier values of degrees. We recall this technique as ‘‘Percentiles’’ and we include it in baseline methods, along with ‘‘KS-test’’ and ‘‘Power-law’’ (the power-law exponent) in order to evaluate our proposed distance metric which is called ‘‘DDQC’’.

3 Proposed Method

The degree distribution of a network is described in Equation 2 as a probability distribution function. The equation shows the probability of node degrees in the given graph G , in which $D(v)$ is the degree of node v , and $P_G(d)$ is the probability that the degree of a node is equal to d . Based on our observations in networks of different sizes, we propose a new method for feature extraction and comparison of the degree distributions. In this method, a vector of eight real numbers is extracted from the degree distribution. A distance function is also suggested for comparing the feature vectors. In order to reduce the impact of the network size on the extracted features, we considered the mean and standard deviation of the degree distribution in the feature extraction procedure. Equation 3 and Equation 4 show the mean and standard deviation of the degree distribution respectively.

According to Equation 5, for any network G , we divide the range of node degrees into four regions ($R_G(r), r = 1..4$). The regions are defined based on the mean, standard deviation, minimum, and maximum of node degrees. Equation 6 shows the length of each region ($|R_G(r)|$), in which $left(R)$ indicates the lower-bound (minimum degree) in region R and $right(R)$ is its upper-bound. As Equation 7 shows, each region is further divided into two equal-size intervals ($I_G(i), i = 1..8$). Although it is possible to consider more than two intervals in each region, our experiments showed that considering more than two intervals per region results in no considerable improvements in the accuracy of the distance metric. Equation 8 defines the “interval degree probability” (IDP_G) which shows the probability that the specified interval I includes the degree of a randomly chosen node. We have defined eight intervals in the degree distribution (four regions and two intervals per region), and the degree distribution can be quantified based on the IDP of the eight intervals. Equation 9 shows the final quantification (feature-vector) of the degree distribution. The proposed feature extraction method can be utilized in network analysis applications which rely on feature extraction and/or comparison of network degree distributions.

$$P_G(d) = P(D(v) = d); v \in V(G) \quad (2)$$

$$\mu_G = \sum_{d=\min_G(D(v))}^{\max_G(D(v))} d \times P_G(d) \quad (3)$$

$$\sigma_G = \sqrt{\sum_{d=\min_G(D(v))}^{\max_G(D(v))} P_G(d) \times (d - \mu_G)^2} \quad (4)$$

$$R_G(r) = \begin{cases} [\min_G(D(v)), \mu_G - \sigma_G] & r = 1 \\ [\mu_G - \sigma_G, \mu_G] & r = 2 \\ [\mu_G, \mu_G + \sigma_G] & r = 3 \\ [\mu_G + \sigma_G, \max_G(D(v))] & r = 4. \end{cases} \quad (5)$$

$$|R_G(r)| = \max(right(R_G(r)) - left(R_G(r)), 0) \quad (6)$$

$$I_G(i) = \begin{cases} [left(R_G(\lceil \frac{i}{2} \rceil)), left(R_G(\lceil \frac{i}{2} \rceil)) + \frac{|R_G(\lceil \frac{i}{2} \rceil)|}{2}] & i \text{ is odd} \\ [left(R_G(\lceil \frac{i}{2} \rceil)) + \frac{|R_G(\lceil \frac{i}{2} \rceil)|}{2}, right(R_G(\lceil \frac{i}{2} \rceil))] & i \text{ is even} \end{cases} \quad (7)$$

$$IDP_G(I) = P(left(I) \leq D(v) < right(I)); v \in V(G) \quad (8)$$

$$Q(G) = \langle IDP_G(I_G(i)) \rangle_{i=1..8} \quad (9)$$

After the feature extraction phase, we can compare the degree distribution of two networks G_1 and G_2 according to their quantified feature vectors. We propose the Equation 10 for comparing two degree distributions. This equation compares two networks based on their corresponding feature vectors (Q_i values). Intuitively, $d(G_1, G_2)$ compares the corresponding interval degree probabilities of the two networks and sums their differences. Equation 10 is a distance function for degree distribution of networks, and it is the result of a comprehensive study of different real, artificial and temporal networks.

$$d(G_1, G_2) = distance(G_1, G_2) = \sum_{i=1}^8 |Q_i(G_1) - Q_i(G_2)| \quad (10)$$

4 Evaluation

In this section, we evaluate our proposed method. In subsection 4.1 we describe different network datasets which are used in our evaluations and in subsection 4.2 we compare our proposed method with baseline methods.

4.1 Datasets

In our problem setting, we aim a distance function that given the degree distribution of two networks, calculates how similar they are. But what does this “similarity” mean for degree distributions? What benchmark is available for evaluating such a distance function? For evaluating different distance metrics, an approved dataset of networks with known distances of its instances is sufficient. Although there is no such an accepted benchmark of networks with known “distance values”, there exist some similarity witnesses among the networks. For evaluating different distance metrics, we have prepared two network datasets with admissible similarity witnesses among the networks of these datasets:

- **Artificial Networks.** We generated a dataset of 6,000 labeled artificial networks using six different generative models. Generative models are network generation methods for synthesizing artificial graphs that resemble the topological properties of real-world complex networks. We considered six generative models in the evaluations: Barabási-Albert model [5], Erdős-Rényi [42], Forest Fire [7], Kronecker model [16], random power-law [17], and Small-world (Watts-Strogatz) model [20]. Each evaluation scenario consists of 100 iterations of network generation and the average of the evaluation results of the 100 iterations are reported. In each iteration 60 networks are generated using the six generative models (10 networks per model). As a result, the evaluations on the artificial networks includes the generation of a total of 6,000 network instances using different parameters. The number of nodes in generated networks ranges from 1,000 to 5,000 nodes with the average of 2,936.34 nodes in each network instance. The average number of edges is 13,714.75. In this dataset, the generative models (generation methods) are the witnesses of the similarity: The networks generated by the same model follow identical link formation rules, and their degree distributions are considered similar. The networks of this data-set are further described in the A, along with an overview of the selected generative models.
- **Real-world Networks.** We have collected a dataset of 33 real-world networks of different types. The networks are selected from six different network classes: Friendship networks, communication networks, collaboration networks, citation networks, peer to peer networks and graph of linked web pages. The category of networks is a sign of similarity: networks of the same type usually follow similar link formation procedures and produce similar degree distributions. So, when comparing two network instances, we expect the distance metric to return small distances (in average) for networks of the same type and relatively larger distances for networks with different types. The “real-world networks” data-set is described in the A, along with the basic properties and the source of its networks.

We assume that the category of the networks in the artificial and real datasets is a sign of their similarity. Networks of the same type follow similar link formation procedures and produce networks with similar structures. Although it is possible for two different-class networks to be more similar than two same-class networks, we assume that the overall “expected similarity” among networks of the same class is more than the expected similarity of different-class networks. This definition of the network similarity which is based on the network types is frequently utilized in the literature [23, 24, 25, 43, 44, 45, 46].

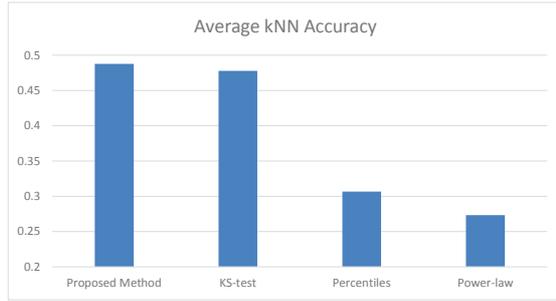


Figure 3: Average kNN accuracy in real networks dataset for $K=1..10$

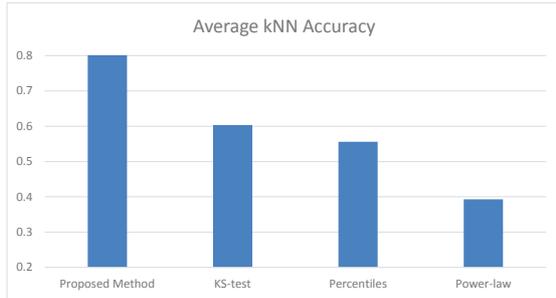


Figure 4: Average kNN accuracy in artificial networks dataset for $K=1..10$

4.2 Evaluation Results

In the section 4.1, we described our two network datasets and we introduced different signs and witnesses of similarities among networks of these datasets. We can consider these witnesses in the evaluation of the proposed method. In this subsection, we evaluate our proposed distance function and we compare it with baseline methods based on their consistency to mentioned witnesses of the similarity.

We start the evaluations by evaluating the precision of k-Nearest-Neighbor classifier based on different distance functions. The k-Nearest-Neighbor rule (kNN) [47] is a common method for classification. It categorizes an unlabeled example by the majority label of its k-nearest neighbors in the training set. The performance of kNN is essentially dependent on the way that similarities are computed between different examples. Hence, better distance metrics result in better classification accuracy of kNN. In order to evaluate the accuracy of different distance functions, we employ them in kNN classification and we test the accuracy of this classifier. This evaluation is performed for both labeled datasets of real-world and artificial networks. As Equation 11 shows, in a *dataset* of labeled instances, the KNN-accuracy of a distance metric d is the probability that the predicted class of an instance is equal to its actual class, when the distance metric d is used in the KNN classifier. Figure 3 and Figure 4 illustrate this evaluation in the real and artificial network datasets respectively. In these evaluations, kNN rule is performed with different values of k from 1 to 10 and the average kNN accuracy is computed. As the figures show, the proposed method results in the best kNN accuracy among the baseline methods.

$$KNN\text{-Accuracy}(d) = P(KNN\text{-Classify}_d(x) = class(x)), x \in dataset \quad (11)$$

As an alternative evaluation criterion, we can consider the Precision-at-K ($P@K$) for assessing the proposed method. As Equation 12 shows, $P@K$ indicates the percentage of classmates in the K nearest neighbors. In other words, $P@K$ is equal to the average number of classmates in the K nearest neighbors of an instance, divided by k . $P@K$ is dependent on the distance metric

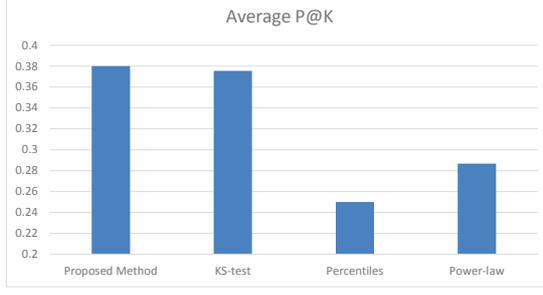


Figure 5: Average P@K in real networks dataset for K=1..10

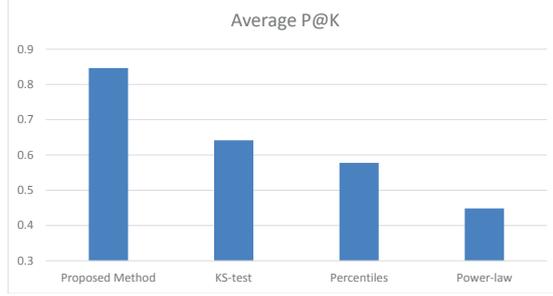


Figure 6: Average P@K in artificial networks dataset for K=1..10

d that is utilized for computing the distances among the dataset instances. Figure 5 and Figure 6 show the average P@K for K=1..10 in the networks of the real and artificial network datasets respectively, according to different distance metrics. As both the figures show, the proposed method outperforms all the baseline methods with respect to the average P@K measure.

$$P@K(d) = \frac{E(c)}{k}; c = \text{count}(m), m \in KNN_d(x) \text{ and } \text{class}(m) = \text{class}(x), x \in \text{dataset} \quad (12)$$

An appropriate distance metric should return smaller distances for networks of the same class. In this context, Dunn Index [48] is an appropriate measure for comparing the inter/intra class distances. For any partition setting U , in which the set of instances are clustered or classified into c groups ($U \leftrightarrow X = X_1 \cup X_2 \cup \dots \cup X_c$), Dunn defined the *separation index* of U as described in Equation 13 [48]. Dunn index investigates the ratio between the average distance of the two nearest classes (Equation 14) and the average distance between the members of the most extended class (Equation 15). Figure 7 and Figure 8 illustrate the Dunn index for different network distance functions in the real and artificial network datasets respectively. The proposed method shows the best (the largest) Dunn index among the baseline methods.

$$DI(U) = \min_{1 \leq i \leq c} \left\{ \min_{\substack{1 \leq j \leq c \\ j \neq i}} \left\{ \frac{\delta(X_i, X_j)}{\max_{1 \leq k \leq c} \{\Delta(X_k)\}} \right\} \right\} \quad (13)$$

$$\delta(S, T) = \delta_{avg}(S, T) = \frac{1}{|S||T|} \sum_{x \in S, y \in T} d(x, y) \quad (14)$$

$$\Delta(S) = \frac{1}{|S| \cdot (|S| - 1)} \sum_{x, y \in S, x \neq y} d(x, y) \quad (15)$$

In the last experiment, we investigate the integration of the degree distribution features with other network features, such as clustering coefficient and assortativity. In this experiment,

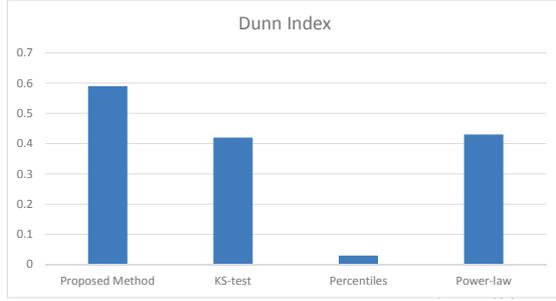


Figure 7: Dunn index for different distance metrics in real networks dataset.

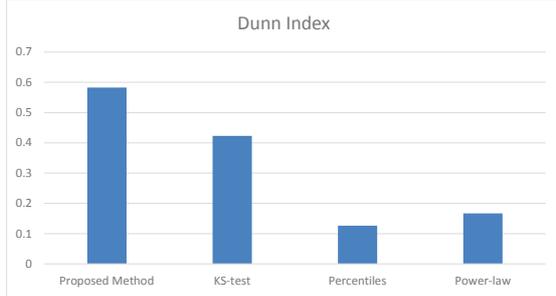


Figure 8: Dunn index for different distance metrics in artificial networks dataset.

we represent the networks with feature vectors that not only include the degree distribution features, but also some other important features that reflect the structural characteristics of complex networks. Then, we utilize supervised machine learning algorithms in order to classify the networks of the artificial and real network datasets using the integrated feature vectors. The aim of this experiment is to find the degree distribution features that best improves the accuracy of the classifier. Along with the degree distribution features, we consider four well-known network features: 1- “Average clustering coefficient” [20] reflects the transitivity of connections. 2- “Average path length” [3] shows the average shortest path between any pair of nodes. 3- The “assortativity” measure [49] shows the degree correlation between pairs of linked nodes. 4- “Modularity” [50] is a measure for quantifying community structure of a network. When a network is represented with a feature vector of these four properties, we call the feature vector “Features”. We consider three other feature vector representations in which the degree distribution features are also considered: “Features+Powerlaw” adds the fitted power-law exponent, “Features+Percentiles” includes the eight percentile features and “Features+DDQC” adds our proposed features of degree distribution. Since Kolmogorov-Smirnov (KS) does not provide a feature extraction method, we do not consider it in this experiment. We used Support Vector Machines (SVM) [51] as the classification method with 10-fold cross-validation. SVM performs a classification by mapping the inputs into a high-dimensional feature space and constructing hyperplanes to categorize the data instances. We utilized Sequential Minimal Optimization (SMO) [51] which is a common method for solving the optimization problem. Figure 9 shows the accuracy of the classifier based on the described four versions of the feature vectors in the real networks dataset. Figure 10 shows the result of the same experiment for the artificial networks dataset. As both the figures show, among the feature extraction methods for degree distribution, our proposed method results in the best classifier. In other words, our proposed method extracts the most informative features from the degree distribution that can improve the accuracy of network classification.

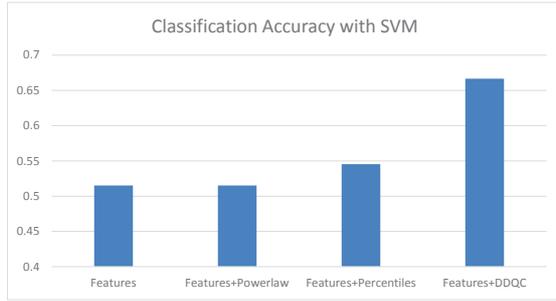


Figure 9: SVM classification accuracy in real networks dataset.

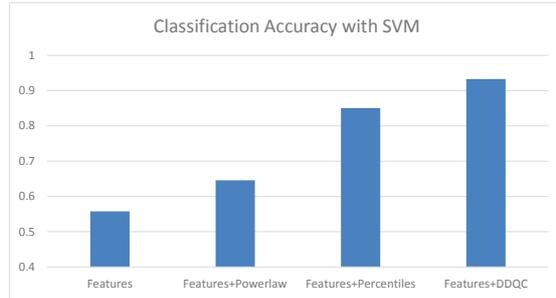


Figure 10: SVM classification accuracy in artificial networks dataset.

5 Conclusion

In this paper, we proposed a novel method for quantification (feature extraction) and comparison of network degree distributions. The aim of “quantification” is extracting a fixed-length feature vector from the degree distribution. Such a feature vector is used in network analysis applications such as network comparison. The “network comparison” is performed by returning a real number as the distance between two degree distributions. The distance is the counterpart of “similarity” and larger distances indicate less similarity. The degree distribution is an indicator of the link formation process in the network, which reflects the overall pattern of connections [14]. Similarly evolving networks have analogous degree distributions and we derive the similarity of degree distributions according to the similarity of link formation process in the networks. For deriving the amount of similarity of networks, we introduced admissible witnesses for network similarity: Similarity among the networks in two categories (same-type real networks and same-model artificial networks). We assume the networks in each of these categories have similar degree distributions. This assumption is the base of our evaluations for different degree-distribution distance metrics. Our proposed method, named DDQC, outperforms the baseline methods with regard to its accuracy in various evaluation criteria. The evaluations are performed based on different criteria such as the ability of the similarity metric to classify networks and comparison of inter/intra class distances. Although the integration of other network features (such as assortativity and modularity) improves the accuracy of the network similarity function, a similarity metric which is only based on the degree distribution has independent and important applications.

Our proposed method enables the data analysis applications and data mining algorithms to employ the degree distribution as a fixed-length feature vector. Hence, it is now possible to represent a network instance with a record of features (including clustering coefficient, average path length and the quantified degree distribution) and use such records in data analysis applications. As the future works, we will combine different network features along with the

quantified degree distribution in an integrated distance metric for complex networks. Such an integrated distance metric will be the main building block of our future researches in evaluation and selection of network generative models and sampling methods.

Acknowledgements

We appreciate Sadegh Motallebi, Sina Rashidian and Javad Gharechamani for their helps in implementation and evaluation of the methods and in preparation of network datasets. We also thank Masoud Asadpour, Mahdiah Soleymani Baghshah and Hossein Rahmani for their valuable comments.

A Overview of Artificial and Real-world Network Datasets

In this appendix, we briefly describe the utilized network datasets of this research. The “artificial networks” dataset consists of a total of 6,000 networks, which are synthesized using six generative models. For each generative model, 1000 network instances are generated using different parameters. The number of nodes in generated networks is configured from 1,000 to 5,000 nodes, with the average of 2,936.34 nodes and 13,714.75 edges in each network instance. The selected generative models are some of the important and widely used network generation methods which cover a wide range of degree distribution structures. The generative models are described in the following, along with their configuration parameters in generation of “artificial networks” dataset. The values of the model parameters are selected according to the hints and recommendations of the cited original papers, along with a concern of keeping the number of network edges balanced for different models. The datasets are available upon request.

- **Barabási-Albert model (BA)**. This is the classical preferential attachment model which generates scale free networks with power-law degree distributions [5]. In this model, new nodes are incrementally added to the graph, one at a time. Each new node is randomly connected to k existing nodes with a probability that is proportional to the degree of the available nodes. In the artificial networks dataset, k is randomly selected as an integer number from the range $1 \leq k \leq 10$.
- **Erdős-Rényi (ER)**. This model generates completely random graphs with a specified density [42]. Network density is defined as the ratio of the existing edges to potential edges. The density of the ER networks in the artificial networks dataset is randomly selected from the range $0.002 \leq density \leq 0.005$.
- **Forest Fire (FF)**. This model, in which edge creation is similar to fire-spreading process, supports shrinking diameter and densification properties along with heavy-tailed in-degrees and community structure [7]. This model is configured by two main parameters: Forward burning probability (p) and backward burning probability (p_b). For generating artificial networks dataset, we fixed $p_b = 0.32$ and selected p randomly from the range $0 \leq p \leq 0.3$.
- **Kronecker graphs (KG)**. This model generates realistic synthetic networks by applying a matrix operation (the kronecker product) on a small initiator matrix [16]. The model is mathematically tractable and supports many network features including small path lengths, heavy tail degree distribution, heavy tails for eigenvalues and eigenvectors, densification, and shrinking diameters over time. The KG networks of the artificial networks dataset are generated using a 2×2 initiator matrix. The four elements of the initiator

matrix are randomly selected from the ranges: $0.7 \leq P_{1,1} \leq 0.9, 0.5 \leq P_{1,2} \leq 0.7, 0.4 \leq P_{2,1} \leq 0.6, 0.2 \leq P_{2,2} \leq 0.4$.

- **Random power-law (RP)**. This model follows a variation of ER model and generates synthetic networks with power law degree distribution [17]. This model is configured by the power-law degree exponent (γ). In our parameter setting for generating artificial networks dataset, γ is randomly selected from the range $2.5 < \gamma < 3$.
- **Watts-Strogatz model (WS)**. The classical Watts-Strogatz small-world model synthesizes networks with small path lengths and high clustering [20]. It starts with a regular lattice, in which each node is connected to k neighbors, and then randomly rewires some edges of the network with rewiring probability β . In WS networks of the artificial networks dataset, β is fixed as $\beta = 0.5$, and k is randomly selected from the integer numbers between 2 and 10 ($2 \leq k \leq 10$).

Table 1 describes the graphs of the “real-world networks” dataset, along with the category, number of nodes and edges, and the source of these graphs. Most of these networks are publicly available datasets. Two temporal networks (Cit_CiteSeerX and Collab_CiteSeerX) are extracted from CiteSeerx digital library [38], using a web crawler software tool.

References

References

- [1] R. Albert, A.-L. Barabási, Statistical mechanics of complex networks, *Reviews of modern physics* 74 (1) (2002) 47–97.
- [2] M. E. Newman, The structure and function of complex networks, *SIAM review* 45 (2) (2003) 167–256.
- [3] L. d. F. Costa, F. A. Rodrigues, G. Travieso, P. Villas Boas, Characterization of complex networks: A survey of measurements, *Advances in Physics* 56 (1) (2007) 167–242.
- [4] D. Easley, J. Kleinberg, *Networks, Crowds, and Markets: Reasoning about a highly connected world*, Cambridge university press, 2010.
- [5] A.-L. Barabási, R. Albert, Emergence of scaling in random networks, *Science* 286 (5439) (1999) 509–512.
- [6] A. Clauset, C. R. Shalizi, M. E. Newman, Power-law distributions in empirical data, *SIAM review* 51 (4) (2009) 661–703.
- [7] J. Leskovec, J. Kleinberg, C. Faloutsos, Graphs over time: densification laws, shrinking diameters and possible explanations, in: *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, ACM, 2005, pp. 177–187.
- [8] L. Muchnik, S. Pei, L. C. Parra, S. D. Reis, J. S. Andrade Jr, S. Havlin, H. A. Makse, Origins of power-law degree distribution in the heterogeneity of human activity in social networks, *Scientific reports* 3.
- [9] J. Leskovec, C. Faloutsos, Sampling from large graphs, in: *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, 2006, pp. 631–636.

Table 1: Dataset of real-world networks

Category	ID	Vertices	Edges	Source
Citation Network	Cit-HepPh	34,546	420,899	SNAP [52]
	Cit-HepTh	27,770	352,304	SNAP [52]
	dblp_cite	475,886	2,284,694	DBLP [53]
	Cit_CiteSeerX	1,106,431	11,791,228	CiteSeerX [38]
Collaboration Network	CA-AstroPh	18,772	198,080	SNAP [52]
	CA-CondMat	23,133	93,465	SNAP [52]
	CA-HepTh	9,877	25,985	SNAP [52]
	Collab_CiteSeerX	1,260,292	5,313,101	CiteSeerX [38]
	com-dblp	317,080	1,049,866	SNAP [52]
	dblp_collab	975,044	3,489,572	DBLP [53]
	dblp20080824	511,163	1,871,070	Sommer [54]
	IMDB-09	4,155	16,679	Rossetti [55]
	CA-GrQc	5,242	14,490	SNAP [52]
CA-HepPh	12,008	118,505	SNAP [52]	
Communication Network	EmailURV	1,133	5,451	Aarenas [56]
	Email-Enron	36,692	183,831	SNAP [52, 57]
	Email-EuAll	265,214	365,025	Konect [58]
	WikiTalk	2,394,385	4,659,565	SNAP [52]
Friendship Network	Dolphins	62	159	NetData [59]
	facebook-links	63,731	817,090	MaxPlanck [60]
	Slashdot0811	77,360	507,833	SNAP [52]
	Slashdot0902	82,168	543,381	SNAP [52]
	soc-Epinions1	75,879	405,740	SNAP [52]
	Twitter-Richmond	2,566	8,593	Rossetti [55]
Graph of Web Pages	youtube-d-growth	1,138,499	2,990,443	MaxPlanck [60]
	web-BerkStan	685,230	6,649,470	SNAP [52]
	web-Google	875,713	4,322,051	SNAP [52]
	web-NotreDame	325,729	1,103,835	SNAP [52]
P2P Network	web-Stanford	281,903	1,992,636	SNAP [52]
	p2p-Gnutella04	10,876	39,994	SNAP [52]
	p2p-Gnutella05	8,846	31,839	SNAP [52]
	p2p-Gnutella06	8,717	31,525	SNAP [52]
	p2p-Gnutella08	6,301	20,777	SNAP [52]

- [10] V. Gómez, A. Kaltenbrunner, V. López, Statistical analysis of the social network and discussion threads in slashdot, in: Proceedings of the 17th international conference on World Wide Web, ACM, 2008, pp. 645–654.
- [11] M. P. Stumpf, C. Wiuf, R. M. May, Subnets of scale-free networks are not scale-free: sampling properties of networks, Proceedings of the National Academy of Sciences of the United States of America 102 (12) (2005) 4221–4224.
- [12] E. M. Airoidi, X. Bai, K. M. Carley, Network sampling and classification: An investigation of network model representations, Decision support systems 51 (3) (2011) 506–518.
- [13] A. Sala, L. Cao, C. Wilson, R. Zablit, H. Zheng, B. Y. Zhao, Measurement-calibrated graph models for social network experiments, in: Proceedings of the 19th international conference on World wide web, ACM, 2010, pp. 861–870.
- [14] V. Boginski, S. Butenko, P. M. Pardalos, Mining market data: a network approach, Computers & Operations Research 33 (11) (2006) 3171–3184.
- [15] C. J. Stam, J. C. Reijneveld, Graph theoretical analysis of complex networks in the brain, Nonlinear biomedical physics 1 (1) (2007) 3.
- [16] J. Leskovec, D. Chakrabarti, J. Kleinberg, C. Faloutsos, Z. Ghahramani, Kronecker graphs: An approach to modeling networks, The Journal of Machine Learning Research 11 (2010) 985–1042.
- [17] D. Volchenkov, P. Blanchard, An algorithm generating random graphs with power law degree distributions, Physica A: Statistical Mechanics and its Applications 315 (3) (2002) 677–690.
- [18] J.-D. J. Han, D. Dupuy, N. Bertin, M. E. Cusick, M. Vidal, Effect of sampling on topology predictions of protein-protein interaction networks, Nature biotechnology 23 (7) (2005) 839–844.
- [19] J. Leskovec, K. J. Lang, M. Mahoney, Empirical comparison of algorithms for network community detection, in: Proceedings of the 19th international conference on World wide web, ACM, 2010, pp. 631–640.
- [20] D. J. Watts, S. H. Strogatz, Collective dynamics of 'small-world' networks, nature 393 (6684) (1998) 440–442.
- [21] M. P. Stumpf, C. Wiuf, Sampling properties of random graphs: the degree distribution, Physical Review E 72 (3) (2005) 036118.
- [22] S. H. Lee, P.-J. Kim, H. Jeong, Statistical properties of sampled networks, Physical Review E 73 (1) (2006) 016102.
- [23] S. Motallebi, S. Aliakbary, J. Habibi, Generative model selection using a scalable and size-independent complex network classifier, Chaos: An Interdisciplinary Journal of Nonlinear Science 23 (4) (2013) 043127.
- [24] J. Janssen, M. Hurshman, N. Kalyaniwalla, Model selection for social networks using graphlets, Internet Mathematics 8 (4) (2012) 338–363.

- [25] M. Middendorf, E. Ziv, C. H. Wiggins, Inferring network mechanisms: the drosophila melanogaster protein interaction network, *Proceedings of the National Academy of Sciences of the United States of America* 102 (9) (2005) 3192–3197.
- [26] K. Juszczyszyn, N. T. Nguyen, G. Kolaczek, A. Grzech, A. Pieczynska, R. Katarzyniak, Agent-based approach for distributed intrusion detection system design, in: *Proceedings of the 6th international conference on Computational Science*, Springer, 2006, pp. 224–231.
- [27] P. Papadimitriou, A. Dasdan, H. Garcia-Molina, Web graph similarity for anomaly detection, *Journal of Internet Services and Applications* 1 (1) (2010) 19–30.
- [28] R. Pastor-Satorras, A. Vespignani, Epidemic dynamics in finite size scale-free networks, *Physical Review E* 65 (3) (2002) 035108.
- [29] A. Montanari, A. Saberi, The spread of innovations in social networks, *Proceedings of the National Academy of Sciences of the United States of America* 107 (47) (2010) 20196–20201.
- [30] L. Briesemeister, P. Lincoln, P. Porras, Epidemic profiles and defense of scale-free networks, in: *Proceedings of the 2003 ACM workshop on Rapid malcode*, ACM, 2003, pp. 67–75.
- [31] G. Kossinets, D. J. Watts, Empirical analysis of an evolving social network, *Science* 311 (5757) (2006) 88–90.
- [32] M. L. Goldstein, S. A. Morris, G. G. Yen, Problems with fitting to the power-law distribution, *The European Physical Journal B-Condensed Matter and Complex Systems* 41 (2) (2004) 255–258.
- [33] W. Deng, W. Li, X. Cai, Q. A. Wang, The exponential degree distribution in complex networks: Non-equilibrium network theory, numerical simulation and empirical data, *Physica A: Statistical Mechanics and its Applications* 390 (8) (2011) 1481–1485.
- [34] M. Faloutsos, P. Faloutsos, C. Faloutsos, On power-law relationships of the internet topology, in: *ACM SIGCOMM Computer Communication Review*, Vol. 29, ACM, 1999, pp. 251–262.
- [35] N. Z. Gong, W. Xu, L. Huang, P. Mittal, E. Stefanov, V. Sekar, D. Song, Evolution of social-attribute networks: measurements, modeling, and implications using google+, in: *Proceedings of the 2012 ACM conference on Internet measurement conference*, ACM, 2012, pp. 131–144.
- [36] H. Kwak, C. Lee, H. Park, S. Moon, What is twitter, a social network or a news media?, in: *Proceedings of the 19th international conference on World wide web*, ACM, 2010, pp. 591–600.
- [37] M. Kim, J. Leskovec, Multiplicative attribute graph model of real-world networks, *Internet Mathematics* 8 (1-2) (2012) 113–160.
- [38] Citeseerx digital library, <http://citeseerx.ist.psu.edu>
URL <http://citeseerx.ist.psu.edu>
- [39] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, U. Alon, Network motifs: simple building blocks of complex networks, *Science* 298 (5594) (2002) 824–827.

- [40] A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, B. Bhattacharjee, Measurement and analysis of online social networks, in: Proceedings of the 7th ACM SIGCOMM conference on Internet measurement, ACM, 2007, pp. 29–42.
- [41] D. Corlette, F. Shipman, Capturing on-line social network link dynamics using event-driven sampling, in: Proceedings of the International Conference on Computational Science and Engineering, IEEE, 2009, pp. 284–291.
- [42] P. Erdős, A. Rényi, On the central limit theorem for samples from a finite population, Publications of the Mathematical Institute of the Hungarian Academy 4 (1959) 49–61.
- [43] J. P. Bagrow, E. M. Boltt, J. D. Skufca, D. Ben-Avraham, Portraits of complex networks, EPL (Europhysics Letters) 81 (6) (2008) 68004.
- [44] M. Berlingerio, D. Koutra, T. Eliassi-Rad, C. Faloutsos, Network similarity via multiple social theories, in: Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM '13, ACM, New York, NY, USA, 2013, pp. 1439–1440. doi:10.1145/2492517.2492582. URL <http://doi.acm.org/10.1145/2492517.2492582>
- [45] A. Mehler, Structural similarities of complex networks: A computational model by example of wiki graphs, Applied Artificial Intelligence 22 (7-8) (2008) 619–683.
- [46] J.-P. Onnela, D. J. Fenn, S. Reid, M. A. Porter, P. J. Mucha, M. D. Fricker, N. S. Jones, Taxonomies of networks from community structure, Physical Review E 86 (3) (2012) 036104.
- [47] T. Cover, P. Hart, Nearest neighbor pattern classification, IEEE Transactions on Information Theory 13 (1) (1967) 21–27.
- [48] J. C. Bezdek, N. R. Pal, Some new indexes of cluster validity, Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on 28 (3) (1998) 301–315.
- [49] M. E. Newman, Assortative mixing in networks, Physical review letters 89 (20) (2002) 208701.
- [50] M. E. Newman, Modularity and community structure in networks, Proceedings of the National Academy of Sciences of the United States of America 103 (23) (2006) 8577–8582.
- [51] J. C. Platt, Fast training of support vector machines using sequential minimal optimization, in: Advances in kernel methods, MIT Press, 1999, pp. 185–208.
- [52] Stanford large network dataset collection, <http://snap.stanford.edu/data/>. URL <http://snap.stanford.edu/data/>
- [53] Xml repository of dblp library, <http://dblp.uni-trier.de/xml/>. URL <http://dblp.uni-trier.de/xml/>
- [54] Christian sommer's graph datasets, <http://www.sommer.jp/graphs/>. URL <http://www.sommer.jp/graphs/>
- [55] Giulio rossetti networks dataset, <http://giuliorossetti.net/about/ongoing-works/datasets>. URL <http://giuliorossetti.net/about/ongoing-works/datasets>

- [56] Alex arenas's network datasets, <http://deim.urv.cat/~aarenas/data/welcome.htm>.
URL <http://deim.urv.cat/~aarenas/data/welcome.htm>
- [57] B. Klimt, Y. Yang, The enron corpus: A new dataset for email classification research, in: Proceedings of the 15th European Conference on Machine Learning (ECML 2004), Springer, 2004, pp. 217–226.
- [58] The koblenz network collection, <http://konect.uni-koblenz.de/>.
URL <http://konect.uni-koblenz.de/>
- [59] Newman's netdata collection, <http://www-personal.umich.edu/~mejn/netdata/>.
URL <http://www-personal.umich.edu/~mejn/netdata/>
- [60] Network datasets at max planck, <http://socialnetworks.mpi-sws.org>.
URL <http://socialnetworks.mpi-sws.org>