# Commentary: Measuring the shape of degree distributions

JENNIFER M. BADHAM

*School of Engineering and Information Technology, Australian Defence Force Academy,
Northcott Drive, Canberra ACT 2600, Australia*
(*e-mail:* `research@criticalconnections.com.au`)

## Abstract

Degree distribution is a fundamental property of networks. While mean degree provides a standard measure of scale, there are several commonly used shape measures. Widespread use of a single shape measure would enable comparisons between networks and facilitate investigations about the relationship between degree distribution properties and other network features. This paper describes five candidate measures of heterogeneity and recommends the Gini coefficient. It has theoretical advantages over many of the previously proposed measures, is meaningful for the broad range of distribution shapes seen in different types of networks, and has several accessible interpretations. While this paper focuses on degree, the distribution of other node-based network properties could also be described with Gini coefficients.

**Keywords:** *degree distribution, heterogeneity, centralization, Gini coefficient*

## 1 Introduction

Standard measures of network properties are important for many types of analysis. A detailed examination of a specific network that investigates the relationship between network location and influence is likely to report various network measures including size, existence of cliques, or other features of interest. While these measures are useful in understanding the specific network, it is not always clear how to compare values between networks so as to assess whether the network is unusual in some way. Researchers who consider patterns over multiple networks, such as the relationship between a particular structural feature of a network and behavior of the network as a whole, require measures that are valid for any network of interest and comparable over a broad variety of networks.

For degree, the scale of a network (of given size) can be described and compared using the degree mean or density. However, degree distributions with the same mean can have very different shapes. Presenting the cumulative degree distribution in full (such as displayed at Newman 2003, Figure 3.2) is informative for a single network or comparison between a small number of networks, but is unwieldy for comparing many networks. A common shape measure would facilitate network comparison, classification of networks into similar types, and description of the relationship between shape of the degree distribution and other properties of the network.

This paper assesses five commonly used statistics for the shape of distributions. As the standard statistical measure of central tendency, variance is particularly well known and its properties with respect to degree distributions have been examined

in Snijders (1981). Three other statistics are well established in network science: power law exponent (Barabási & Albert 1999), centralization (Freeman 1978), and hierarchization (Coleman 1964). However, these statistics are popular within specific disciplines; for example, power law exponent is popular with researchers working with strongly skewed, large, technologically supported information networks such as World Wide Web links or email address books (Newman 2003), while centralization is included as an available measure in social network analysis software such as UCINET and the sna R package. The final statistic assessed is the Gini coefficient. This is a standard heterogeneity measure for income and wealth distributions (which are typically strongly skewed) but has received only limited attention in network science, with the most theoretical analysis focusing on its relationship with the power law exponent (Hu & Wang 2008) and sporadic use for specific applications (such as in Lopes et al. 2011).

The first section of the paper provides the structure necessary to compare the measures. The main body of the paper then presents each measure and assesses its performance against desirable theoretical properties. The final section summarizes the theoretical analysis and compares the measures empirically, using both real world and constructed networks. The paper recommends the Gini coefficient as the most suitable shape measure for degree distributions, because it has desirable theoretical properties, is appropriate for any shaped distribution, and has several useful and intuitive interpretations.

## 2 Analysis framework

While there has not been a comprehensive examination of degree distribution shape measures and their properties, individual measures have been assessed on specific properties. For example, Snijders (1981) preferred degree variance over centralization because the latter focuses entirely on a single node and is insensitive to other high degree nodes. Similarly, Hu & Wang (2008, p. 3771) propose the Gini coefficient because it "*is superior to some other parameters in characterizing heterogeneity, such as variance or standard deviation, since … they demand that two networks studied should have the same average degree.*"

A systematic identification of desirable properties and comparison of measures has, however, been conducted for income inequality (Dalton 1920; Allison 1978; Cowell 2000). Three of these properties are equally relevant for comparing potential degree distribution shape measures: transfer, addition, and replication. These are considered in Section 2.1.

In addition, measures must be valid and sensitive over a broad range of distribution shapes. Several real-world networks taken from the literature and artificial networks constructed with well-established algorithms are used to demonstrate the measures empirically. These networks were selected for the diversity of their degree distribution shapes and are described in Section 2.2.

### 2.1 Desirable properties of a shape measure

The first property that should be held by a measure of heterogeneity is that it does indeed measure heterogeneity. That is, as distributions become more (or less)

unequal, the value of the measure should change in some consistent direction. This idea is captured by the transfer principle, which states that rewiring edges from a high degree node to a lower degree node should decrease heterogeneity, provided that the transfer does not lead to a ranking reversal of the two nodes. In the limit, this also implies that the minimum value is achieved when all nodes have the same degree (or as close as possible if mean degree is not an integer).

The addition principle considers the effect of increasing the degree of all nodes by the same absolute amount. There are two potentially appropriate consequences for the inequality measure. From the perspective that inequality refers to absolute differences (or variation about the mean), the measure should be unchanged. Alternatively, from the perspective that inequality is relative, the measure should decrease, on the basis that the differences between degrees are proportionally smaller if the same number of edges is added to each node.

Finally, the heterogeneity statistic should not be affected by replication. That is, the value of the shape statistic calculated over multiple instances of the same network should equal the value for a single instance. This principle extends to non-integer replications by considering the distribution based on the least common multiple of the two network sizes, generated by appropriate numbers of replications of each of the initial distributions.

### 2.2 Empirical and artificial degree distributions

In addition to the theoretical analysis, the heterogeneity measures are compared for four real-world networks taken from the literature, two networks constructed with well-established algorithms and a star network. The networks were selected to emphasize differences in the measures, with substantial diversity in size and in shape of degree distribution.

The four real-world distributions are:

- Friends: the number of friendship nominations received within a school study (Rapoport & Horvath 1961, Table 5), with 859 nodes, mean in-degree of 6.84 and maximum in-degree of 29.
- Yeast: the yeast protein interaction network described in Jeong et al. (2001), with 2,114 nodes, mean in-degree of 2.12 and maximum in-degree of 56.
- Collaborators: collaborations between authors on the condensed matter archive from January to March 2005 (updated version of Newman 2001), with 40,421 nodes, mean degree of 8.69 and maximum degree of 278.
- WWW: hyperlinks between domains in the World Wide Web (Albert et al. 1999), with 325,729 nodes, mean in-degree of 4.60 and maximum in-degree of 10,721.

The three artificial networks are of the same size and two have a common mean degree so as to emphasize the differences arising from shape. The random graph and preferential attachment algorithms are well established for generating artificial networks with very different degree distributions. The three generated distributions are:

- BA1000: a single instance generated with the preferential attachment algorithm described in Barabási & Albert (1999), with three edges per added node and
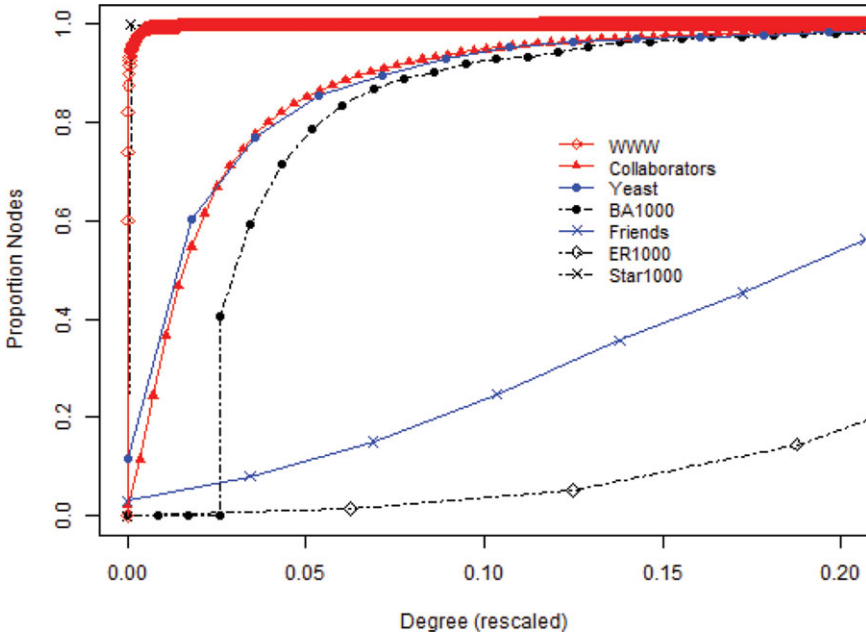
Fig. 1. Cumulative distribution for real-world and generated degree distributions, with degree displayed as proportion of maximum degree. Except for the Friends and ER1000 networks, the maximum degree is over five times the degree values for almost all other nodes. (color online)

a complete initial network with three nodes. This network has 1,000 nodes, mean degree of 5.99 and maximum degree of 116.

- ER1000: a single instance generated with the fixed number of edges algorithm described in Erdös & Rényi (1960), matched to the BA1000 network size and number of edges. This network has 1,000 nodes, mean degree of 5.99 and maximum degree of 16.
- Star1000: a star network with 1,000 nodes, one central node with a single edge to the other 999 nodes and no other edges. This network has mean degree of 2.00 and maximum degree of 999.

The degree distribution for each of these seven networks is shown in Figure 1, with degree rescaled to a proportion of its network specific maximum and truncated at 20%. The ER1000 and Friends networks have a clearly different degree distribution from the other five networks, with the latter group so skewed that almost all the nodes have been accounted for by 20% of the maximum degree. The WWW and Star1000 networks substantially overlap and are more extreme than the other three.

### 2.3 Notation and assumptions

Network size or number of nodes is denoted by $N$; and $k$ is used as a general indicator of degree, with $k_i$ for the degree of node $i$, $N_k$ for the number of nodes with degree $k$, $\mu_k$ for mean degree and $p_k = N_k/N$ for the proportion of nodes with degree $k$ (for empirical distributions) or probability of a given node having degree $k$

(for ideal distributions). The notation in equations taken from references is adapted accordingly.[1]

The networks are assumed to be simple, so self edges and multiple edges are not permitted and the maximum degree is given by $N - 1$. Three of the real-world networks (Friends, Yeast, and WWW) are directed, and the in-degree distribution is used for each. Despite the use of directed networks to generate example distributions, for those measures where there is a different calculation method for undirected and directed networks, this paper uses the undirected method in all cases. That is, these in-degree distributions are interpreted as degree distributions from undirected networks even if such a network is not realizable. This is for consistency, so as to provide a basis for comparison and more effectively demonstrate the properties of the candidate measures over different distribution shapes.

### 3 Candidate measures of distribution shape

Five broad shape measures are described; each normalized (where applicable) to facilitate comparability over networks of different sizes. The first two are hierarchization and centralization, which were developed specifically for social networks. The next two are variance and power law exponent, which apply standard statistical techniques to degree distributions. The final measure is the Gini coefficient, used predominantly to measure inequality of income or wealth, but applicable to any distribution.

#### 3.1 Coleman's hierarchization

Coleman (1964, pp. 434–441) developed two related hierarchization indices specifically for degree distributions. While they were defined only for directed networks, applying the same equations to undirected networks does not change their qualitative behavior. The first compared the network of interest to a multinomial distribution null model and is not discussed here. The second extends the concept of entropy as a measure of choice developed in Shannon (1948). In the context of networks, such choice is the degree values realized from edges "choosing" two nodes for its ends. The index $h_2$ is entropy normalized against the maximum possible conditional on the number of nodes and edges. It has a value between 0 and 1 and is given by

$$h_2 = \frac{S_{\max} - S}{S_{\max} - S_{\min}} \qquad \text{where } S \text{ is entropy} \tag{1a}$$

$$= \frac{\log_e N + \sum_{i=1}^{N} \frac{k_i}{N\mu_k} \log_e \frac{k_i}{N\mu_k}}{\log_e N - 0} \tag{1b}$$

$$= \frac{\sum_{k=0}^{\max k} N_k k \log_e k - N\mu_k \log_e \mu_k}{N\mu_k \log_e N} \tag{1c}$$

Entropy is not suitable because a transfer may introduce a new degree value in the distribution and potentially increase the value of entropy even where heterogeneity

---

[1] Notation adaptation may also include algebraic manipulation, such as where the original equation uses density rather than mean degree. Derivations are available from the author.

is reduced. Normalizing to maximum entropy corrects this problem and $h_2$ respects the transfer principle. However, this respect is at the expense of breaching the replication principle; because there is more choice available with additional nodes and edges, the maximum entropy is larger and the value of $h_2$ is reduced for the replicated distribution.

### 3.2 Freeman's centralization

A different approach was taken by Freeman (1978) in his classic study of different types of node centrality, including degree. He argued that the network measure of centralization should describe the extent to which a single node is more central than the others and dominates the network, normalized with respect to the maximum possible value for a network of the same size (which occurs for a star network).

One weakness of this measure is that the feasible range of values for any network is conditional also on the mean degree, which makes it difficult to use to compare heterogeneity between different networks. This was addressed in Butts (2006, Equation (36)) by instead normalizing to the maximum feasible given $N$ and $\mu_k$, so that centralization has a potential range of 0 to 1 for any degree distribution:

$$C = \frac{k_{\max} - \lceil \mu_k \rceil}{\mu_k + \min\left[(N-1) - \mu_k, \mu_k\left(\frac{N}{2}-1\right)\right] - \lceil \mu_k \rceil} \tag{2}$$

A more significant weakness is that centralization is concerned only with the highest degree node compared to the average, whereas real-world networks can show more subtle degree heterogeneity with "*a vaguely outlined center, consisting of more than one point; or there are several centers; or just a gradual transition from more central to more peripheral points*" (Snijders 1981, p. 164). The consequence of this focus on a single node is that $C$ breaches the transfer principle, even with the normalization, because centralization is unaffected by redistribution of edges within a network provided the maximum degree is unaffected. $C$ also breaches the addition property; increasing the degree of every node can actually increase its value in some situations where one of the denominator terms $(N - 1 - \mu_k)$ decreases.

### 3.3 Variance

Variance ($\sigma_k^2$) and its square root, standard deviation, are well-established statistical measures of deviation from the average for any distribution. Despite its ubiquity, the application of variance to degree was not specifically considered until Snijders (1981) proposed it as a more sensitive measure than Freeman's centralization index as it takes into account the full distribution rather than only the maximum and mean. That paper also investigated properties such as maximum variance and expected value of degree variance for different types of networks.

Snijders (1981, p. 172) proposed the index $J$, constructed as the standard deviation of degree normalized to the maximum possible for a network of the given size and density.[2] There are two general network structures with maximum degree variance:

---

[2] Other indices proposed in Snijders (1981) are not pursued further in this paper as they measured difference from a null model assuming random distribution of the given number of edges between the given number of nodes rather than reference to a baseline of all nodes with equal degree.

the star network (or its complement), where one node has a single edge with each of the other nodes; and a network with all the edges concentrated into a complete subnetwork with leftover nodes as isolates. There is no single equation for $J$ because of the need to select between these two possible maxima and realizability issues. Instead, the maximum standard deviation of degree for a given number of nodes and edges is calculated using the algorithm (see Snijders 1981, pp. 167–168 for details):

- Let $T_1$ be the largest integer for which $T_1(T_1 - 1) \leqslant N\mu_k$ (that is, total degree);
- Let $E_1$ be the leftover edges $E_1 = [N\mu_k - T_1(T_1-1)]/2$;
- Let $T_2$ be the largest integer for which $T_2(T_2 - 1) \leqslant N(N - 1) - N\mu_k$;
- Let $E_2$ be the equivalent leftover edges $E_2 = [N(N-1) - N\mu_k - T_2(T_2-1)]/2$;
- Use Equation (3) to calculate standard deviation of degree for the two networks constructed with a complete subnetwork of size $T_i$ and $E_i$ leftover edges (for $i = 1, 2$) and the larger value is the maximum standard deviation of degree for the original network.

With this maximum calculated, $J$ follows:

$$\sigma_{\max} = \sqrt{\frac{T_i(T_i - 1)^2 + E_i(2T_i + E_i - 1)}{N} - \frac{(T_i(T_i - 1) + 2E_i)^2}{N^2}} \tag{3}$$

and then normalize:

$$J = \frac{\sigma_k}{\sigma_{\max}} \tag{4}$$

This measure complies with both the transfer and addition properties, with the value of $J$ changing in the appropriate direction. However, like $h_2$, normalizing against the maximum leads to a breach of the replication property because the maximum variance increases as there are more nodes and edges available and $J$ consequently decreases.

The simplest standardization of variance that meets the transfer, addition, and replication principles is the coefficient of variation, the ratio of the standard deviation to the mean:

$$V_k = \frac{\sigma_k}{\mu_k} = \frac{\sqrt{\frac{1}{N} \sum_{i=1}^{N} (k_i - \mu_k)^2}}{\mu_k} \tag{5}$$

Coefficient of variation has a range of 0 (homogeneous) to $\infty$ (though lower for realized networks). As degree is always non-negative, $V_k$ can be interpreted as relative variability.

### 3.4 Fitted function parameter

The other common statistical approach to distribution characterization is to fit the coefficients or parameters of some function that represents the idealized probability density. In the literature concerning very large social and information networks such as email address books, citations, and web page links, degree is treated as continuous and the functional form fitted to the degree distribution is typically the power law

(Brinkmeier & Schank 2005). That is, the degree for each node is given by the probability:

$$p_k = Ck^{-\alpha} \quad \text{with } k > 0 \tag{6}$$

where $\alpha$ is the parameter of the distribution and referred to as the power exponent. The constant $C$ is completely determined by $\alpha$ and the requirement that probabilities over all degrees $k$ sum to 1.

In practice, the power law is fitted to only the higher degree part of the distribution because the function diverges as $k \rightarrow 0$ and would substantially overestimate the number of low degree nodes if applied to all degrees. Identifying an appropriate minimum $k$ from which to apply the power law is somewhat arbitrary but has little impact on the value of the power exponent. A more serious flaw is that rigorous tests suggest that at least some of the network degree (and other) distributions identified in the literature as following power laws may be better described with some other skewed functional form such as the lognormal distribution (Newman 2005; Clauset et al. 2009).

Fitting the correct function limits the usefulness of this approach because coefficients can be compared only between distributions with the same functional form. Nevertheless, the derived power law exponent allows degree distributions of other shapes to be compared if the power law functional form is considered as simply an approximation, though there is no formal standard for acceptability of a fit. In addition, care must be taken to use a robust estimator as the "obvious" logarithmic transformation and least squares regression can lead to substantial error (Clauset et al. 2009). Typical values for $\alpha$ for degree distributions range between 1.5 and 3.0 (Newman 2003, Figure 3.2), with a larger value leading to a faster reduction in the probability of higher degree nodes and therefore less heterogeneity.

Assuming that a power law function can be fitted to all of the degree distributions required, the power law exponent satisfies the transfer, addition, and replication properties. Thus, it can provide meaningful comparisons between the heterogeneity of such distributions.

### 3.5 Gini coefficient

The Gini coefficient is a widely used inequality (or heterogeneity) measure developed for the skewed distributions of income and wealth, with an extensive body of theoretical support and over 100 years of use (Allison 1978; Cowell 2000). However, it has received limited attention in network science with the exception of Hu & Wang (2008), who proposed its use for degree distributions and examined its relationship with the power law exponent.

Adapted to degree distributions, the Gini coefficient is formally defined as the normalized expected difference in degree between two randomly selected nodes, given by (Gini 1912; Dalton 1920):

$$G = \frac{1}{\mu_k} \frac{1}{2N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} |k_i - k_j| \tag{7}$$

An alternative definition relies on the Lorenz curve, which provides a visual measure of inequality (Lorenz 1905). For degree distribution, the Lorenz curve plots
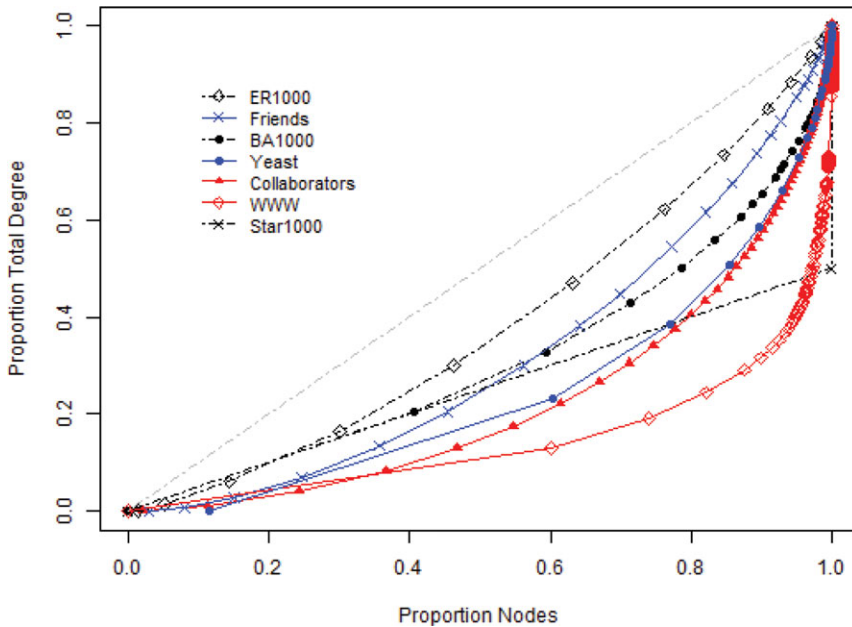
Fig. 2. Lorenz curves of degree for real-world degree distributions (Friends, Yeast, Collaborators, WWW) and networks generated with algorithms (ER1000, BA1000, Star1000). The WWW and ER1000 distributions are the most and least extreme, respectively. (color online)

the cumulative proportion of the nodes ordered by degree against the cumulative proportion of the degree held by those nodes and also includes a (diagonal) reference curve that indicates the Lorenz curve for a distribution where all nodes have the same degree. A greater "bend" away from the reference curve indicates greater inequality. Examples of Lorenz curves for the seven degree distributions described in Section 2.2 are displayed in Figure 2.

It can be shown that $G$ is the area bounded by the Lorenz and reference curves divided by one half (reported at Dalton 1920, p. 354). The half is a normalization factor, representing the size of the area for a distribution of maximum inequality. This approach allows $G$ to be calculated for a theoretical degree distribution provided the mean is finite (Dorfman 1979).

The use of the Lorenz curve also provides an efficient computation method for $G$ for the degree distribution of an empirical network (Brown 1994):

$$G = 1 - \sum_{g=\min k}^{\max k} \left( \sum_{k=1}^{g} p_k - \sum_{k=1}^{g-1} p_k \right) \left( \sum_{k=1}^{g} k p_k + \sum_{k=1}^{g-1} k p_k \right) \qquad (8)$$

where $g$ and $k$ iterate only through those degrees that exist in the network.[3]

Compliance with the properties of transfer, addition, and replication is a key reason for the popularity of the Gini coefficient for describing income and wealth distributions. It was developed specifically to compare these distributions in different locations and over time.

---

[3] There is no mathematical difference to iterate through all degree values, including those that are not realized; it is simply less efficient.

Table 1. *Compliance of potential measures with principles.*

| Measure<br>desirable | Transfer<br>decrease | Addition<br>not increase | Replication<br>no change |
|---|---|---|---|
| Normalized Hierarchization ($h_2$) | Decrease | Decrease | Decrease |
| Normalized Centralization ($C$) | Varies | Varies | Decrease |
| Coefficient of variation ($V_k$) | Decrease | Decrease | No change |
| Normalized (by max) deviation ($J$) | Decrease | Decrease | Decrease |
| Power law exponent ($\alpha$) | Decrease* | Decrease | No change |
| Gini coefficient ($G$) | Decrease | Decrease | No change |

* The numerical value of $\alpha$ increases with a transfer of degree away from high degree nodes, and this indicates a decrease in heterogeneity.

## 4 Comparison of proposed measures

As discussed in the presentation of each measure, only coefficient of variation ($V_k$) and the Gini coefficient ($G$) are valid for all distributions and have the desired properties of consistency over changes in distributions arising from transfer, addition, or replication (see Table 1). The power law exponent is effective for those distributions that can be reasonably approximated by the power law functional form, but this is a strong restriction and there are other measures available. The two measures normalized by their theoretical maximum ($h_2$ and $J$) breach the replication requirements because replication increases the number of nodes and the edges that can be distributed between them thereby increasing maximum heterogeneity.

There is no theoretical reason to prefer coefficient of variation ($V_k$) or the Gini coefficient ($G$). However, there may be normative reasons and further insight can be gained by examining the measure values empirically (see Table 2).

The theoretical weakness of $h_2$, $C$, and $J$ is clearly demonstrated by their values for the larger networks. For all three, the combination of substantial heterogeneity and moderate size for BA1000 results in values that are much higher than the most heterogeneous WWW network, because the moderate size limits the potential maximum that is being used for normalization. These three candidates are therefore insufficiently independent of the network scale to be useful in comparing shapes of degree distributions from different networks.

If the Lorenz curves do not intersect, ordering will be consistent (Allison 1978). For the example networks, the curve for Star1000 crosses many of the other distributions and the BA1000 and Friends pair also intersect. Table 2 is ordered by decreasing values of $V_k$ and $G$, except for Star1000 which ranks inconsistently between the two measures.

Even where the two measures provide the same ranking, they have different patterns of sensitivity. Using similarity of values in Table 2, $V_k$ groups the BA1000 network with Yeast and Collaborators, while $G$ groups it with the Friends network. From Figure 1, it is apparent that the former is the more natural grouping and this supports a preference for $V_k$. However, the actual degrees of the highest degree influence $V_k$ much more than $G$, which could lead to difficulties in comparing distributions that are very heterogeneous as these high degree values may dominate changes in the shape of the bulk of the distribution.

Table 2. *Degree heterogeneity measures for selected networks.*

| Network | Size | $h_2$ | $C$ | $V_k$ | $J$ | $\alpha^*$ | $G$ |
|---|---|---|---|---|---|---|---|
| WWW (in) | $N = 325{,}729$ $\mu_k = 4.6$ | 0.132 | 0.03 | 8.5 | 0.047 | 2.1 | 0.71 |
| Collaborators | $N = 40{,}421$ $\mu_k = 8.7$ | 0.055 | 0.01 | 1.5 | 0.031 | – | 0.55 |
| Yeast (in) | $N = 2{,}114$ $\mu_k = 2.1$ | 0.072 | 0.03 | 1.4 | 0.066 | 2.4 | 0.51 |
| BA1000 | $N = 1{,}000$ $\mu_k = 6.0$ | 0.357 | 0.11 | 1.2 | 0.209 | 3.0 | 0.37 |
| Friends (in) | $N = 859$ $\mu_k = 6.8$ | 0.034 | 0.03 | 0.7 | 0.090 | – | 0.37 |
| ER1000 | $N = 1{,}000$ $\mu_k = 6.0$ | 0.013 | 0.01 | 0.4 | 0.046 | – | 0.23 |
| Star1000 | $N = 1{,}000$ $\mu_k = 2.0$ | 0.400 | 1.00 | 0.5 | 1.000 | – | 0.50 |

\* A missing value for $\alpha$ indicates that it is not available in the literature, which may occur because the power law functional form is inappropriate or because it may be appropriate but was not reported. Unlike other measures in the table, a higher value indicates lower heterogeneity.

The interpretation of the measure supports a preference for $G$. The interpretation of $G$ as the (normalized) expected difference in degree between two randomly selected nodes is natural in the network context, where edges are connecting two nodes. In contrast, $V_k$ describes deviation from the center of the distribution and many degree distributions have no meaningful center.

There are many other heterogeneity measures not considered in this paper. However, they have more fundamental flaws than those included. For example, an intuitive approach is to measure the distance between the cumulative degree distribution and the (homogeneous) delta distribution, with the Kolmogorov–Smirnov statistic providing a distance measure (Massey 1951). However, the cumulative probability distribution of the delta distribution is a step function at mean degree and the maximum distance between the cumulative distributions will occur at the mean degree. Hence, assuming the typical positive skew and that the network is sufficiently large that degree can be considered continuous, the Kolmogorov–Smirnov statistic is simply the proportion of nodes with degree less than or equal to the mean degree. Another approach is to use diversity measures such as the Herfindahl–Hirschman Index (Hirschman 1964). For degree distribution, this index is the sum of the squared contributions from each node to total degree. However, it is sensitive to the network size, as more nodes dilute the effect of each node's degree.

## 5 Conclusions

A measure that describes the shape of the distribution, regardless of the functional form of the distribution, would facilitate research to examine how the shape of the distribution is related to other properties of the network or processes occurring over the network. Properties of a simple network that could reasonably be expected to vary with some measure of degree heterogeneity include maximum degree assortativity (Hakimi 1962), expected degree assortativity (Newman 2002), and size of an epidemic occurring on the network (Diekmann et al. 1990).

Both the coefficient of variation ($V_k$) and the Gini coefficient ($G$) are suitable measures as they respond to redistribution in ways that facilitate valid comparisons between networks of different sizes and mean degree. For both, a value of 0 indicates

that all the nodes have equal degree and larger values indicate greater heterogeneity. However, the coefficient of variation is difficult to interpret in the context of a highly skewed distribution. In particular, its conceptual source is the width of the peak in a distribution and such a peak may not exist for some networks. In contrast, the Gini coefficient represents the difference between degrees for pairs of nodes, rather than comparing a node's degree to the mean degree. Thus, it can be easily interpreted for distributions that are highly non-normal as well as normal distributions. The Gini coefficient is therefore proposed as the most suitable shape measure for degree distributions.

While this paper has focused on degree distribution, there are other network properties that are calculated by node (or node pairs) but typically reported only as mean over all nodes. These properties include clustering coefficient, betweenness, and shortest path lengths (Newman 2003). The shape of these distributions could also be described with Gini coefficients (with discretization) and values of this measure could reasonably be expected to be linked to other network properties.

## Acknowledgments

## References

Albert, R., Jeong, H., & Barabási, A.-L. (1999). Diameter of the world wide web. *Nature (London)*, **401**(6749), 130–131.

Allison, P. D. (1978). Measures of inequality. *American Sociological Review*, **43**(6), 865–880.

Barabási, A.-L., & Albert, R. (1999). Emergence of scaling in random networks. *Science*, **286**, 509–512.

Brinkmeier, M., & Schank, T. (2005). Network statistics. In U. Brandes & T. Erlebach (Eds.), *Network Analysis* (pp. 293–317). Springer-Verlag.

Brown, M. C. (1994). Using Gini style indices to evaluate the spatial patterns of health practitioners: Theoretical considerations and an application based on Alberta data. *Social Science and Medicine*, **38**(9), 1243–1256.

Butts, C. T. (2006). Exact bounds for degree centralization. *Social Networks*, **28**, 283–296.

Clauset, A., Shalizi, C. R., & Newman, M. E. J. (2009). Power-law distributions in empirical data. *SIAM Review*, **51**(4), 661–703.

Coleman, J. S. (1964). *Introduction to Mathematical Sociology*. Free Press (MacMillan).

Cowell, F. A. (2000). Measurement of inequality. In A. B. Atkinson, & F. Bourguignon (Eds.), *Handbook of Income Distribution* (pp. 87–166). Elsevier.

Dalton, H. (1920). The measurement of the inequality of incomes. *The Economic Journal*, **30**(119), 348–361.

Diekmann, O., Heesterbeek, J. A. P., & Metz, J. A. J. (1990). On the definition and the computation of the basic reproduction ratio $R_0$ in models for infectious diseases in heterogeneous populations. *Journal of Mathematical Biology*, **28**, 365–382.

Dorfman, R. (1979). A formula for the Gini coefficient. *The Review of Economics and Statistics*, **61**(1), 146–149.

Erdös, P., & Rényi, A. (1960). On the evolution of random graphs. *Publications of the Institute of Mathematics, Hungarian Academy of Science*, **5**, 17–60.

Freeman, L. C. (1978). Centrality in social networks: Conceptual clarification. *Social Networks*, **1**, 215–239.

Gini, C. (1912). Variabilità e mutabilità. *Studi Economico-Giuridici dell'Università di Cagliari*, **3**, 1–158.

Hakimi, S. L. (1962). On realizability of a set of integers as degrees of the vertices of a linear graph. *Journal of the Society for Industrial and Applied Mathematics*, **10**(3), 496–506.

Hirschman, A. O. (1964). The paternity of an index. *American Economic Review*, **54**(5), 761–762.

Hu, H. B., & Wang, X. F. (2008). Unified index to quantifying heterogeneity of complex networks. *Physica A: Statistical Mechanics and Its Applications*, **387**(14), 3769–3780.

Jeong, H., Mason, S. P., Barabási, A.-L., & Oltvai, Z. N. (2001). Lethality and centrality in protein networks. *Nature (London)*, **411**(6833), 41–42.

Lopes, G. R., da Silva, R., & de Oliveira, J. P. M. (2011). Applying Gini coefficient to quantify scientific collaboration in researchers network. *Proceedings of the International Conference on Web Intelligence, Mining and Semantics, WIMS '11* (pp. 68:1–68:6). New York: ACM.

Lorenz, M. O. (1905). Methods of measuring the concentration of wealth. *Publications of the American Statistical Association*, **9**(70), 209–219.

Massey Jr, F. J. (1951). The Kolmogorov-Smirnov test for goodness of fit. *Journal of the American Statistical Association*, **46**(253), 68–78.

Newman, M. E. J. (2001). The structure of scientific collaboration networks. *Proceedings of the National Academy of Sciences of the United States of America*, **98**, 404–409.

Newman, M. E. J. (2002). Assortative mixing in networks. *Physical Review Letters*, **89**, 208701.

Newman, M. E. J. (2003). The structure and function of complex networks. *SIAM Review*, **45**(2), 167–256.

Newman, M. E. J. (2005). Power laws, Pareto distributions and Zipf's law. *Contemporary Physics*, **46**(5), 323–351.

Rapoport, A., & Horvath, W. J. (1961). A study of a large sociogram. *Behavioral Science*, **6**(4), 279–291.

Shannon, C. (1948). A mathematical theory of communication. *The Bell System Technical Journal*, **27**, 379–423.

Snijders, T. A. B. (1981). The degree variance: An index of graph heterogeneity. *Social Networks*, **3**(3), 163–174.