

1 **Complex networks for streamflow dynamics**

2
3 **B. Sivakumar^{1,2} and F. M. Woldemeskel¹**

4 [1]{School of Civil and Environmental Engineering, The University of New South Wales,
5 Sydney, Australia}

6 [2]{Department of Land, Air and Water Resources, University of California, Davis, USA}

7 Correspondence to: B. Sivakumar (s.bellie@unsw.edu.au)

8 9 **Abstract**

10 Streamflow modeling is an enormously challenging problem, due to the complex and
11 nonlinear interactions between climate inputs and landscape characteristics over a wide range
12 of spatial and temporal scales. A basic idea in streamflow studies is to establish connections
13 that generally exist, but attempts to identify such connections are largely dictated by the
14 problem at hand and the system components in place. While numerous approaches have been
15 proposed in the literature, our understanding of these connections remains far from adequate.
16 The present study introduces the *theory of networks*, and in particular *complex networks*, to
17 examine the connections in streamflow dynamics, with a particular focus on spatial
18 connections. Monthly streamflow data observed over a period of 52 years from a large
19 network of 639 monitoring stations in the contiguous United States are studied. The
20 connections in this streamflow network are examined primarily using the concept of
21 *clustering coefficient*, which is a measure of local density and quantifies the network's
22 tendency to cluster. The clustering coefficient analysis is performed with several different
23 threshold levels, which are based on correlations in streamflow data between the stations. The
24 clustering coefficient values of the 639 stations are used to obtain important information
25 about the connections in the network and their extent, similarity and differences between
26 stations/regions, and the influence of thresholds. The relationship of the clustering coefficient
27 with the number of links/actual links in the network and the number of neighbors is also
28 addressed. The results clearly indicate the usefulness of the network-based approach for
29 examining connections in streamflow, with important implications for interpolation and
30 extrapolation, classification of catchments, and predictions in ungaged basins.

1 **1 Introduction**

2 Streamflow forms an important input for a wide range of applications in hydrology, water
3 resources, environment, and ecosystem. However, its estimation or prediction is an
4 enormously challenging problem, since streamflow arises as a result of complex and
5 nonlinear interactions between climate inputs (external factors) and landscape characteristics
6 (internal factors) that occur over a wide range of spatial and temporal scales. For instance,
7 streamflow is governed not only by the distribution of rainfall (in both space and time) but
8 also by the nature and state of the catchment (e.g. topography, vegetation, soil, geology); see
9 Beven (2006) for a compilation of, and stimulating insight into, some early ‘benchmark’
10 studies (1933–1984) on streamflow generation processes. Attempts to monitor, model, and
11 predict streamflow have been a central topic in hydrology during the last century or so; see,
12 for example, Salas et al. (1995), Grayson and Blöschl (2000), Duan et al. (2003), Mishra and
13 Coulibaly (2009), and Hrachowitz et al. (2013) for comprehensive accounts on streamflow
14 monitoring, modeling, and prediction.

15 Despite their efforts and contributions, studies on streamflow have and continue to encounter
16 at least two major challenges: (1) determination of the locations, number, and density of
17 streamflow gaging stations for monitoring data and representation of process variability; and
18 (2) identification of the appropriate scientific concepts and mathematical techniques/models
19 for a more solid conceptual understanding of the catchment systems, proper analysis of the
20 data, and reliable interpretation of the outcomes. It is true that recent developments in
21 measurement technology, computational power, and mathematical sophistication have
22 generally played an important role in overcoming these challenges to a certain extent. It can
23 also not be denied, however, that the same developments have, at times, played an indirect
24 role in creating imbalance and hindering true progress, as they have contributed to the perhaps
25 unnecessary complexification in models (rather than simplification), highly specialized
26 conceptual notions that are often suitable only for specific situations (rather than
27 generalization frameworks that suit all conditions), difficult-to-bridge gaps between theory
28 and practice, and lack of communication among researchers as well as between researchers
29 and practitioners; see, for example, Perrin et al. (2001), Beven (2002), Kirchner (2006),
30 Sivakumar (2008), and Young and Ratto (2009) for some details.

31 It is important to recognize that a fundamental idea in streamflow (and other hydrologic)
32 studies is to establish connections that generally exist between the different elements or items

1 (known or assumed) of the underlying system. Depending upon the situation (e.g. catchment,
2 purpose, problem), these elements include hydroclimatic variables, catchment characteristics,
3 model parameters, and others (and their combinations), and their connections are often
4 different with respect to space, time, and space-time. Unraveling the nature and extent of
5 these connections has always been a great challenge, not to mention the challenge in the
6 identification of all the relevant elements in the first place. Thus far, a plethora of concepts
7 and methods has been proposed and applied for studying the connections associated with
8 streamflow, including those based on time, distance, correlation, variability, scale, patterns,
9 and many other properties/measures as well as their combinations and variants, in both single-
10 variable and multi-variable perspectives; see, for example, Gupta et al. (1986), Salas et al.
11 (1995), Grayson and Blöschl (2000), Yang et al. (2004), Archfield and Vogel (2010), and Li
12 et al. (2012) for some details. Despite the progress made through these concepts and methods,
13 our understanding of the connections in streamflow is still far from adequate.

14 In view of this, there is indeed a need to greatly advance our studies on streamflow
15 connections. Some important current and foreseeable future problems, including our ever-
16 increasing demands for water, the potential impacts of climate change on water security and
17 hydroclimatic disasters, and the numerous issues associated with the management of our
18 environment and ecosystems, further reflect the urgency to this need. A greater understanding
19 of streamflow connections will also enhance our recent and current efforts in the estimation of
20 data at ungaged locations (e.g. predictions in ungaged basins – PUB) (see Hrachowitz et al.,
21 2013) and development of a generalization framework for hydrologic modeling (e.g.
22 catchment classification) (see Sivakumar et al., 2014), among others. The question, however,
23 remains on the identification of a suitable theory that can help bring advancement to studies
24 on streamflow connections. In this regard, recent developments in the field of complex
25 systems science can offer some crucial clues. The present study introduces the theory of
26 *complex networks*, or simply *networks*, for studying connections in streamflow. In particular,
27 the study focuses on spatial connections in streamflow.

28 The origin of the concept of networks can be traced back to the works of Leonhard Euler,
29 during the first half of the eighteenth century, on the Seven Bridges of Königsberg (Euler,
30 1741), which laid the foundations of what would become popularly known as *graph theory*.
31 Graph theory witnessed several important theoretical developments in the nineteenth century,
32 including *topology* (originally introduced as *topologie* in German) (Listing, 1848) and *trees*

1 (Cayley, 1857). Further significant advances were made during the twentieth century,
2 especially with the development of *random graph theory* by Erdős and Rényi (Erdős and
3 Rényi, 1960). The concepts of graph theory, and random graph theory in particular, have
4 found a wide variety of applications in numerous fields, including linguistics, physics,
5 chemistry, biology, sociology, engineering, economics, and ecology; see, for example, Berge
6 (1962), Bondy and Murty (1976), and Bollobás (1998) for extensive reviews.

7 Despite the above-mentioned developments and applications, studies on graph theory,
8 including random graph theory, had some major deficiencies. First, the studies largely
9 focused on networks that are regular, simple, small, and static. As a result, they are generally
10 unsuitable for examining real networks, as such networks are often highly irregular, complex,
11 large, and dynamically evolving in time. Second, even while examining complex and large-
12 scale networks, they assumed that such networks are wired randomly together (Erdős and
13 Rényi, 1960). Such an assumption, however, is not necessarily valid for real networks, since
14 order and determinism are inherent in real systems and networks. Indeed, real networks are
15 neither completely ordered nor completely random, but generally exhibit important properties
16 of both. These observations motivated a renewed and fresh look of random graph theory
17 towards the end of the last century (e.g. Watts and Strogatz, 1998; Barabási and Albert, 1999),
18 and gave birth to a new movement of interest and research in studying real and complex
19 networks, under the umbrella of the *new science of networks*. They also led to new
20 discoveries about complex networks, including *small-world networks* (Watts and Strogatz,
21 1998), *scale-free networks* (Barabási and Albert, 1999), *network motifs* (Milo et al., 2002), as
22 well as other notable advances, such as a new method for identifying *community structure*
23 (Girvan and Newman, 2002). Since then, the science of networks has found applications in
24 many different fields, including natural and physical sciences, social sciences, medical
25 sciences, economics, and engineering and technology (e.g. Albert et al., 1999; Bouchaud and
26 Mézard, 2000; Newman, 2001; Liljeros et al., 2001; Tsonis and Roebber, 2004; Davis et al.,
27 2013). In hydrology, applications of networks are just starting to emerge, and so far include
28 river networks, virtual water trade, precipitation, and agricultural pollution due to
29 international trade, among others (Rinaldo et al., 2006; Suweis et al., 2011; Dalin et al., 2012;
30 Boers et al., 2013; Scarsoglio et al., 2013). In a very recent study, Sivakumar (2014) has
31 argued that networks can be useful for studying all types of connections in hydrology and,
32 hence, can provide a generic theory for hydrology.

1 With the encouraging results reported by the above studies, the present study explores the
2 usefulness of the theory of networks for studying connections in streamflow, especially the
3 spatial connections. To this end, monthly streamflow data observed over a period of 52 years
4 (1951–2002) from each of 639 gaging stations in the contiguous United States are studied.
5 The connections are examined primarily using the concept of *clustering coefficient*. The
6 clustering coefficient is a measure of local density and, hence, quantifies the tendency of a
7 network to cluster. The implications of the clustering coefficient results for
8 interpolation/extrapolation of streamflow data as well as for classification of catchments are
9 also discussed. To put the clustering coefficient analysis in a proper perspective, traditional
10 linear correlation analysis (Pearson correlation coefficient) and another simple network-based
11 analysis (degree centrality) are also performed.

12 The rest of this paper is organized as follows. Section 2 introduces the concept of networks
13 and describes the procedure for calculation of degree centrality and clustering coefficient in a
14 network. Section 3 presents details of the study area and streamflow data considered. Section
15 4 reports the results, first from the traditional linear correlation analysis and then from the
16 network-based degree centrality and clustering coefficient analysis. Section 5 highlights the
17 implications of the results.

18 **2 Network and clustering coefficient**

19 **2.1 Network**

20 A *network* or a *graph* is a set of points connected together by a set of lines, as shown in
21 Figure 1. The points are referred to as *vertices* or *nodes* and the lines are referred to as *edges*
22 or *links*; here, the term *nodes* are used for points and the term *links* are used for lines.
23 Mathematically, a network can be represented as $G = \{P, E\}$, where P is a set of N nodes
24 (P_1, P_2, \dots, P_N) and E is a set of n links. The network shown in Figure 1 has $N = 7$ (nodes) and n
25 = 8 (links), with $P = \{1, 2, 3, 4, 5, 6, 7\}$ and $E =$
26 $\{\{1, 7\}, \{2, 3\}, \{2, 5\}, \{2, 7\}, \{3, 7\}, \{4, 7\}, \{5, 6\}, \{6, 7\}\}$.

27 Figure 1 is perhaps the simplest form of network, i.e. one with a set of identical nodes
28 connected by identical links. There are, however, many ways in which networks may be more
29 complex. For instance, a network: (1) may have more than one different type of node and/or
30 link; (2) may contain nodes and links with a variety of properties, such as different weights
31 for different nodes and links depending on the strength of nodes and connections; (3) may
32 have links that can be directed (pointing in only one direction), with either cyclic (i.e.

1 containing closed loops of links) or acyclic form; (4) may have multilinks (i.e. repeated links
2 between the same pair of nodes), self-links (i.e. links connecting a node to itself), and
3 hyperlinks (i.e. links connecting more than two nodes together); and (5) may be bipartite, i.e.
4 containing nodes of two distinct types, with links running only between unlike types.

5 There are many different ways and measures to study the characteristics of networks. In the
6 context of the modern theory of *complex networks* (which also include random graphs),
7 degree centrality, clustering coefficient, small-world networks, and degree distribution are
8 some of the prominent concepts. As the present study uses the concepts of degree centrality
9 and clustering coefficient for studying streamflow connections, they are described next.

10 **2.2 Degree centrality**

11 Centrality is one of the most basic and intuitive measures of a network; see Freeman (1979)
12 for an early comprehensive review. The idea behind the use of centrality as a network
13 measure is that it identifies whether a given node, say, i in a network is more central or more
14 influential than another node in the network. The degree centrality of node i in a network of N
15 nodes is defined as the number of first neighbors (or simply *neighbors*) of node i divided by
16 the total number of possible neighbors ($N - 1$) in the network.

17 Let us consider a selected node i in a network of N nodes, having k_i links which connect it to
18 k_i other nodes. For illustration, Figure 2 presents a network consisting of nine nodes (i.e. $N =$
19 9), with the node i having four links (i.e. with four other nodes) (see Figure 2, left). In this
20 case, the four nodes corresponding to the four links are the *neighbors* of node i , which are
21 identified based on some conditions (e.g. correlation between node i and other nodes in the
22 network), while the total number of possible neighbors for node i is eight (i.e. $N - 1$).

23 **2.3 Clustering coefficient**

24 The clustering coefficient quantifies the tendency of a network to cluster, which is one of the
25 most fundamental properties of networks (Watts and Strogatz, 1998). The clustering
26 coefficient of a network is basically a measure of local density. The concept of clustering has
27 its origin in sociology, under the name “fraction of transitive triples” (Wasserman and Faust,
28 1994). The procedure for calculating the clustering coefficient is as follows.

29 Let us consider first a selected node i in the network, having k_i links which connect it to k_i
30 other nodes, as shown in Figure 2 (left). If the neighbors of the original node (i) were part of a
31 cluster, there would be $k_i(k_i - 1)/2$ links between them. As shown in Figure 2 (right), there are

1 $4(4 - 1)/2 = 6$ links in the *cluster* of node i . The clustering coefficient of node i is then given
2 by the ratio between the number E_i of links that actually exist between these k_i nodes (shown
3 as solid lines on Figure 2, right) and the total number $k_i(k_i - 1)/2$ (i.e. all lines on Figure 2,
4 right),

$$5 \quad C_i = \frac{2E_i}{k_i(k_i-1)} \quad (1)$$

6 The clustering coefficient of the whole network C is the average of the clustering coefficients
7 C_i 's of all the individual nodes.

8 The clustering coefficient of a random graph is $C = p$ (where p is the probability of two nodes
9 being connected), since the links in a random graph are distributed randomly. However, the
10 clustering coefficient of real networks is generally much larger than that of a comparable
11 random network (i.e. having the same number of nodes and links as the real network).
12 Therefore, the clustering coefficient analysis offers useful information about the nature of the
13 network and, hence, the appropriate model (e.g. level of complexity), among others.

14 **3 Study area and data**

15 In the present study, streamflow data from the United States are studied to explore the
16 usefulness of the theory of networks for identifying connections in streamflow, with a focus
17 on spatial connections. Monthly data from an extensive network of 639 streamflow gaging
18 stations in the contiguous US are studied. The locations of these 639 stations are shown in
19 Figure 3. The above streamflow data are obtained from the US Geological Survey database, in
20 particular from the Hydro-Climatic Data Network (HCDN), originally developed by Slack
21 and Landwehr (1992) and subsequently updated at different times, with the last update in
22 2009; see Lins (2012) for details (<http://water.usgs.gov/osw/hcdn-2009/>). The HCDN is a
23 subset of all USGS streamgages for which the streamflow primarily reflects prevailing
24 meteorological conditions for specified years; see Kiang et al. (2013) for the latest and
25 comprehensive account of USGS streamflow gages across the entire United States. The
26 HCDN streamgage stations were screened to exclude sites where human activities, such as
27 artificial diversions, storage, and other activities in the drainage basin or the stream channel,
28 affect the natural flow of the watercourse.

29 Streamflow data in the US are commonly expressed in “water years,” which commence in
30 October. The data used in this study are those observed over a period of 52 years (1951–
31 2002), obtained from an earlier version of HCDN. The data are average monthly values (not

1 anomalies). During the past few decades, a large number of studies have investigated the
2 above streamflow dataset (or a part or variant of it) in many different contexts (e.g. Slack and
3 Landwehr, 1992; Kahya and Dracup, 1993; Vogel and Sankarasubramanian, 2000;
4 Sivakumar, 2003; Tootle and Piechota, 2006; Patil and Stieglitz, 2012; Sivakumar and Singh,
5 2012; Kiang et al., 2013). Some of these studies have explicitly addressed the connections of
6 streamflow between the stations, including in the context of data correlations, catchment
7 similarities, and other measures; see, Patil and Stieglitz (2012) and Kiang et al. (2013) for
8 some recent studies. Many studies have explored the connections of streamflow with large-
9 scale climatic patterns and relevant indices, including El-Niño, La-Niña, Southern Oscillation
10 Index (SOI), Pacific North America (PNA) Index, and Pacific Decadal Oscillation (PDO).
11 However, within the specific context of the network analysis for connections among
12 streamflow stations presented here, as well as in the broader context of complex systems
13 science for streamflow analysis, the studies by Sivakumar (2003) and Sivakumar and Singh
14 (2012) are worth mentioning, as they have addressed the aspects of streamflow variability,
15 nonlinearity, and dominant governing mechanisms, especially for studies on model
16 simplification, data interpolation/extrapolation, and catchment classification framework.

17 The above 639 streamflow stations and the observed streamflow data exhibit tremendous
18 variations in their characteristics, often by about four orders of magnitude. For instance: (1)
19 basin drainage area ranges from 10.62 km² (4.1 mi²) to 35224 km² (13600 mi²) (2) station
20 elevation ranges from 0 m to 2996 m (9830 ft); (3) mean flow ranges from 0.0549 m³/s (1.94
21 ft³/s) to 381.59 m³/s (13476 ft³/s); (4) maximum flow ranges from 0.878 m³/s (31 ft³/s) to
22 2489 m³/s (87900 ft³/s); and (5) number of zero-flow months ranges from none to 424. Figure
23 3, for instance, presents the variations in the mean (Figure 3a), standard deviation (Figure 3b),
24 and coefficient of variation (Figure 3c) of flow values in all the 639 stations. The significant
25 differences in catchment and flow characteristics can play important roles in the nature and
26 extent of connections in streamflow between the different stations. While studying their
27 influences is clearly important, the present study does not specifically attempt to address this.
28 Rather, the focus of the present study is in identifying the extent of connections among the
29 stations based on streamflow data alone.

30 **4 Analysis and results**

31 The usefulness of the theory of networks for studying connections in streamflow is examined
32 primarily through the clustering coefficient analysis on the monthly streamflow data from the

1 above 639 stations in the United States. To put the clustering coefficient analysis in a proper
2 perspective, however, linear correlation and degree centrality analyses are also performed.

3 **4.1 Linear correlation analysis**

4 A common approach to examine connections between streamflow observed at different
5 stations is through a simple linear cross correlation analysis, where the correlation for any
6 given station is given by the average of its correlation with all the other stations. Several
7 variants of this procedure are also usually considered. These include: *nearest* neighbors – for
8 example, *number of nearby stations based on distance* or stations within a pre-defined *region*
9 *of geographic promixity* or *neighborhood*, with equal or unequal weightage (e.g. inverse
10 distance); and *similar* stations – stations with *similar* properties (e.g. in terms of climate,
11 rainfall, basin characteristics, land use), which may or may not include nearest stations. These
12 and many other *correlation-based* procedures (e.g. spline fitting) are routinely employed for
13 interpolation and extrapolation of streamflow and other hydrologic data.

14 In this study, two of the above-mentioned procedures are employed for examining the
15 monthly streamflow from the 639 stations: (1) for each station, the correlation is the average
16 of its correlation with all the other 638 stations; and (2) for each station, the correlation is the
17 average of correlations for a *certain number of nearest neighbors* – 30, 15, and 5 neighbors.
18 The neighbors are selected based on the geographical distance from the reference station. For
19 the three different number of neighbors (i.e. 30, 15, and 5) considered in the latter, the mean
20 distances are 111, 73, and 41 km, respectively, and the standard deviations are 94, 63, and 37
21 km, respectively. The correlation considered here is the Pearson correlation coefficient, and
22 the streamflow values themselves (rather than their logarithms) are used for computation. The
23 Pearson correlation coefficient can be sensitive to outliers in the data. However, the impact of
24 this sensitivity is minimal for monthly streamflow (when compared to streamflow at shorter
25 timescales, e.g. daily), as the monthly data assumes approximately normal distribution from
26 additive errors at finer timescales through the central limit theorem (Anderson, 2010).

27 When all the 638 stations are considered, the correlation values are generally very low, as
28 expected, with only 0.5% of the stations exceeding a value of 0.4 (see Figure 4a). This is
29 mainly due to the consideration of a very large region, with the stations coming from different
30 climatic, catchment, land use, and other characteristics. When the number of stations is
31 reduced, the results get generally better – see Figure 4b (30 neighbors), Figure 4c (15
32 neighbors), and Figure 4d (5 neighbors). Among the three neighborhood cases, the best

1 correlation results are obtained when the neighborhood is the smallest, i.e. 5 neighbors
2 (Figure 4d), with a large number of stations having correlations above 0.7.

3 While one can study a large number of combinations in terms of the *neighborhood*, what is
4 evident from even the very few cases presented here is that there are obvious *regional* patterns
5 in terms of correlations, regardless of the number of neighbors. These *regional* patterns are
6 considered to have important implications for a wide range of studies in hydrology and water
7 resources, as they are commonly used as a basis for interpolation and extrapolation of
8 streamflow and, subsequently, for water resources assessment, planning, and management.
9 However, as Sivakumar and Singh (2012) point out, through their nonlinear dynamic study on
10 streamflow data from the western United States, the use of *regional* patterns as basis for
11 streamflow studies may be misleading, as such patterns are not necessarily a true
12 representation of the actual connections between the stations but may just be spurious. The
13 obvious question, therefore, is: how to identify if the connections are actual or spurious? This
14 is where the ideas from the theory of networks can be particularly useful.

15 **4.2 Degree centrality analysis**

16 The degree centrality is calculated for the monthly streamflow data from the network of 639
17 stations in the United States, according to the procedure described in Section 2.2. The essence
18 of the procedure for the streamflow data is as follows. For a given streamflow station or node
19 i , the nearest neighbors k_i in the network of 639 stations (more specifically, the remaining 638
20 stations) are identified based on a (pre-specified) threshold value (T). To define the threshold
21 value, the correlations in streamflow data between different stations are considered as a
22 reasonable measure. With this, if, for example, the correlation between station i and any other
23 station(s) in the entire network of 639 stations exceeds the threshold value, then that station(s)
24 is considered as a *neighbor(s)*, k_i , for station i . The degree centrality of station i is then given
25 by the ratio of the number of neighbors to the total number of possible neighbors (i.e. 638).

26 In this study, several different threshold values are considered for calculation of the degree
27 centrality. Although there are no definitive guidelines for selection of the threshold values for
28 streamflow (and other hydrologic) data, our experience in streamflow studies, especially
29 spatial and temporal correlations, offers some useful clues. For instance, streamflow data
30 generally exhibit high spatial correlations (when compared to rainfall values, for example),
31 especially at the monthly scale. With this knowledge, and also with the condition that $-1 < T$

1 < 1.0, closer intervals of values are considered at the higher end of correlations and vice-
2 versa. In addition, very low values (say, $T < 0.30$) and very high values (say, $T > 0.85$) do not
3 offer much help in the analysis; for instance, $T < 0.30$ normally results in a very large number
4 of neighbors, while $T > 0.85$ results in a very small number. Considering all these, eight
5 threshold values are used for analysis: 0.30, 0.40, 0.50, 0.60, 0.70, 0.75, 0.80, and 0.85.

6 Figure 5a–d, for example, shows the results from the degree centrality analysis for the 639
7 stations for threshold values of 0.70, 0.75, 0.80, and 0.85. The results offer some interesting
8 observations. For instance, only a very small number of streamflow stations (blue circles)
9 have connections with more than 10% of the other stations in the network of 639 stations,
10 while a large number of stations (cyan circles) have connections to less than just 1% of the
11 other stations. Indeed, for thresholds of 0.70, 0.75, 0.80, and 0.85, the number of stations
12 having connections with more than 10% of the stations is only 39, 0, 0, and 0, respectively,
13 while the number of stations having connections with less 1% of the stations is 118, 160, 257,
14 and 429, respectively. This clearly suggests that only a small proportion of stations has
15 considerable influence in the network, while a large proportion of stations has only very little
16 or almost no influence. This result has significant implications, for example, in interpolation
17 and extrapolation, especially from the viewpoint of *dominant* stations (as is the case of 39
18 stations for $T = 0.70$; Figure 5a). It is also important to note, however, that not all of the
19 stations (i.e. *neighbors*) that a given station has connection with (see Figure 5a–d) are the
20 *geographic* neighbors, and some are over long geographic distances (see Section 4.3 for
21 further details on this). These observations seem to suggest that the streamflow network of
22 639 stations is neither a completely ordered network nor a random graph, but some other.

23 **4.3 Clustering coefficient analysis**

24 Following the description in Section 2.3, the procedure for the calculation of the clustering
25 coefficient for the monthly streamflow data from the network of 639 stations in the United
26 States is as follows. For a given streamflow station or node i , the nearest neighbors k_i in the
27 network of 639 stations are identified based on a (pre-specified) threshold value (T), as
28 explained above. The *cluster* of these k_i neighbors then forms the basis for identifying the
29 *actual connections*. Therefore, the *actual connections* are those links in the *cluster* of stations
30 (not just *nearest* stations) having correlations among themselves exceeding the threshold
31 value. Similar to the degree centrality analysis above, eight threshold values are considered in
32 the cluster coefficient analysis as well: 0.30, 0.40, 0.50, 0.60, 0.70, 0.75, 0.80, and 0.85.

1 Figure 6a–d, for instance, presents the clustering coefficient values for the 639 stations for
2 threshold values of 0.70, 0.75, 0.80, and 0.85. Table 1 presents the number of stations falling
3 under different ranges of clustering coefficient values. For better illustration and discussion,
4 the clustering coefficient values are grouped into six different ranges. In Figure 6 and Table 1,
5 a clustering coefficient of 0.0 indicates that there are *no actual connections*, while ‘NA’
6 indicates there are *no nearest neighbors*. From an overall perspective, the clustering
7 coefficient results indicate certain similarity at some stations/regions but significant
8 differences at others. They also offer some specific observations:

- 9 • Even *nearest* stations have significantly different characteristics (e.g. connections), as
10 part of a network. Some stations have very strong connections, while others have
11 almost no or only very weak connections. For instance, the few geographically closer
12 stations in Florida in the southeast region are an excellent example. These few stations
13 have clustering coefficient values varying anywhere from 0 to 1.0, especially for $T =$
14 0.7 and 0.75 (Figures 6a and 6b).
- 15 • Even *distant* stations have significantly similar characteristics, i.e. they have very
16 strong (or very weak or even no) connections as part of a network. The similar (very
17 high or very low) clustering coefficient values obtained for a number of stations all
18 across the United States, regardless of their geographic promixity, offer evidence to
19 this; for example, regardless of the threshold value, the green circles (see Figure 6a–d),
20 representing the clustering coefficient range 0.76–1.0, are present all over the United
21 States, northwest to southwest to midwest to northeast to southeast. Similar
22 observations are made also for other clustering coefficient ranges, for one or more
23 threshold values; see the deep pink circles ($C_i = 0.51–0.75$) and blue circles ($C_i = \text{NA}$);
- 24 • There are significant changes in characteristics with respect to the threshold values.
25 For instance, as can be seen from Figure 6 and Table 1, for threshold values of 0.7 and
26 0.85, the number of stations falling within the clustering coefficient range of 0.51–0.75
27 is 348 and 197, respectively. Indeed, in some cases, further breakdown in the range of
28 clustering coefficient values indicate an even wider difference in the (percentage)
29 number of stations for these thresholds;
- 30 • Although there are changes in the number of stations having similar clustering
31 coefficient values with respect to thresholds, there is no consistency in the trend of
32 changes (see, for example, the number of stations falling within the clustering
33 coefficient range of 0.51–0.75).

1 While the usefulness of the clustering coefficient values in assessing connections between
2 streamflow stations and identifying regions having similarity/differences is abundantly clear,
3 the *actual links* in the network would certainly offer more specific details as to where and
4 how connections exist. To facilitate this, Figure 7 shows the *actual links* for four selected
5 streamflow stations (red circles) for threshold values of 0.75 (Figure 7a), 0.80 (Figure 7b),
6 and 0.85 (Figure 7c); the nodes and links for $T = 0.70$ are too many, and so do not offer a
7 good visualization. In each of these plots, for the station of interest (red circle), a green circle
8 indicates a station that has a correlation coefficient value exceeding the threshold, and a black
9 circle indicates a station that has a correlation coefficient value smaller than the threshold.
10 The lines are the *actual links* among all the links available for the *cluster of neighbors* (green
11 circles only). The plots clearly indicate which stations are *actually* connected to which other.
12 The plots make it abundantly clear that *geographic proximity* does not always result in greater
13 correlation, and the *actual links* can go for large distances. Among the various observations
14 that can be made, the ones for the two stations in the northwest are certainly interesting.
15 Despite being in the same region, the two stations exhibit significantly different connectivity
16 characteristics, for example, for threshold level 0.85 (Figure 7c), with one showing all the
17 actual connections within a small neighborhood (see the enlarged plot on the top left) while
18 the other showing no clear neighborhood for connectivity (see the enlarged plot on the bottom
19 left). The latter station (see bottom left) is an even more curious case, as most of the
20 neighbors of this station seem to be beyond its (perceived) *circle of geographic influence*. The
21 actual links observed for the other threshold values also support the above observations.

22 These observations clearly suggest that our usual approach with consideration of geographic
23 proximity, nearest neighbors, regional patterns, and linear correlation-based techniques for
24 studying connections in streamflow may have serious limitations. Clustering coefficient, and
25 other network-based techniques, offers a better means to examine streamflow connections. In
26 what follows, we explore the clustering coefficient results even further.

27 As the clustering coefficient of a network is based on the *actual links* among *all links* in the
28 *cluster* of neighbors of a node (rather than just the links between a node and its neighbors), it
29 would be interesting to see how it changes with respect to *all links* and *actual links*. To this
30 end, Figure 8a–d shows the clustering coefficient values against the *number of all links* (red
31 circles) and the *number of actual links* (blue circles) for threshold values of 0.70, 0.75, 0.80,

1 and 0.85 for the monthly streamflow data from the United States. The results lead to the
2 following major observations:

- 3 • In general, regardless of the threshold value, there is an inverse relationship between
4 the clustering coefficient and number of links (both for *all links* and *actual links*), i.e.
5 higher clustering coefficient for smaller number of links and vice-versa;
- 6 • The inverse relationship between the clustering coefficient and number of links is
7 generally more evident for lower thresholds (see Figure 8a and b) when compared to
8 higher thresholds (see Figure 8c and d). When the threshold is very high ($T = 0.85$),
9 this relationship seems to cease to exist;
- 10 • The clustering coefficient is generally far more sensitive when the number of links is
11 smaller (see the significant larger spread of circles on the Y-axis), but has only very
12 little or almost no sensitivity for a larger number of links (see the very narrow spread
13 followed by a tapering towards a fixed value – especially in Figure 8a and b). Further,
14 larger numbers of links almost always give lower clustering coefficients;
- 15 • For a given number of links, the clustering coefficient for a lower threshold is
16 generally higher than that for a higher threshold.

17 Another useful way to look at the clustering coefficient of a network is its relationship with
18 the number of neighbors (k_i), which is defined by the threshold value and dictates the (number
19 of) links and actual links. Figure 9a–d shows the relationship between the clustering
20 coefficient values and the number of neighbors for threshold values of 0.70, 0.75, 0.80, and
21 0.85 for the monthly streamflow data. The results generally indicate an inverse relationship
22 between the clustering coefficient and number of neighbors, but such a relationship is far
23 more evident for lower threshold values (see Figure 9a and b) than that for higher threshold
24 values (see Figure 9c and d). Again, the clustering coefficient is generally far more sensitive
25 when the number of neighbors is smaller (see the larger spread towards the left), but becomes
26 less sensitive for a larger number of neighbors (see the narrow spread towards the right).
27 These observations are somewhat consistent with those made in regard to the number of links
28 (Figure 8). It is important to recall, however, that the neighbors are not necessarily geographic
29 but defined by the threshold values (as shown in Figure 7).

30 While these results and observations are still preliminary in nature, they seem to suggest that
31 there is a particular threshold value or range beyond which the inverse relationship between
32 the clustering coefficient and number of neighbors/links/actual links in the streamflow

1 network may not hold well for monthly streamflow data from the United States, and
2 streamflow data in general.

3 Finally, the question arises as to the type of network. As mentioned previously, the clustering
4 coefficient of a whole network (C) is the average of the clustering coefficients C_i 's of all the
5 individual nodes. The clustering coefficient of the eight different networks of the above 639
6 streamflow stations corresponding to threshold values of 0.30, 0.40, 0.50, 0.60, 0.70, 0.75,
7 0.80, and 0.85 is 0.79, 0.76, 0.71, 0.68, 0.65, 0.63, 0.58, and 0.51 (see Table 1). These
8 generally high clustering coefficient values seem to suggest that the streamflow monitoring
9 network of 639 stations is not a random graph, since a (comparable) random graph, where the
10 links are distributed randomly, will have a typically very low clustering coefficient, i.e. $C = p$,
11 where p is the probability of two nodes being connected. As (natural) streamflow dynamics
12 are neither completely random (there are inherent deterministic patterns) nor completely
13 ordered (there are inherent stochastic components) (see Sivakumar, 2011; Sivakumar and
14 Singh, 2012 for some details), it is also reasonable to assume that streamflow networks are
15 not random graphs, but networks of some other nature. Whether they are *small-world* or
16 *scale-free* or other types of networks remains to be seen. Studies in this direction are currently
17 underway, details of which will be reported in the future.

18 **5 Study implications**

19 One of the basic requirements in studying streamflow dynamics is to identify connections in
20 space or time or space-time, depending upon the purpose. Although a wide variety of
21 approaches have been developed and applied to identify connections in streamflow dynamics,
22 there is no question that significant improvements are still needed. In this regard, modern
23 developments in the field of network theory, especially complex networks, offer new avenues,
24 both for their generality about systems and for their holistic perspective about connections.

25 The present study has made an initial attempt to apply the ideas developed in the field of
26 complex networks to examine connections in streamflow dynamics, with particular focus on
27 spatial connections. Application of the concepts of clustering coefficient, which is a measure
28 of local density and quantifies the tendency of a network to cluster to monthly streamflow
29 data from a large network of 639 monitoring stations in the contiguous United States has
30 offered some very interesting results. The clustering coefficient values for the 639 stations
31 suggest that: (1) even nearest stations can have significantly different connections and distant
32 stations can have significantly similar connections; (2) connections can be significantly

1 different for different threshold levels; (3) there is generally an inverse relationship between
2 the clustering coefficient and number of neighbors, number of all links, and actual links (in
3 the cluster of neighbors); (4) the clustering coefficient is far more sensitive when the number
4 of neighbors/number of links is smaller, but has only little or no sensitivity when the latter is
5 larger; and (5) the high clustering coefficient value obtained for the entire network is not
6 consistent with the one expected for a random graph, suggesting that the streamflow network
7 is likely to be small-world or scale-free or some other type. The results from the degree
8 centrality analysis suggest that a very small number of streamflow stations have more
9 influence in the network of 639 stations with connections to more than 10% of the other
10 stations, while a large number of stations have very little influence with connections to just
11 less than 1% of the other stations. These observations seem to further eliminate the possibility
12 of random nature of the streamflow monitoring network.

13 Although the present results are preliminary, they offer important information about the
14 connections that possibly exist in the streamflow network, and especially their extent. The
15 clustering coefficient values, and the *actual links*, are particularly useful in the identification
16 of the specific regions where interpolation and extrapolation of streamflow data may be more
17 effective and also of the specific stations whose data can be more reliable for such purposes.
18 For instance, regions consisting of stations with high clustering coefficient values would
19 generally provide a more accurate estimation of streamflow when interpolation and
20 extrapolation schemes are employed. It is also important to emphasize, however, that such a
21 region is identified based on *cluster of actual connections*, rather than based on our traditional
22 way of geographic proximity, nearest neighbors, regional patterns, and linear correlations.
23 The clustering coefficient values can also offer important clues and guidelines as to the setting
24 up/removal of streamflow monitoring stations in a region. For instance, if a region consists of
25 stations with very high clustering coefficients, then installing additional monitoring stations
26 will not offer any significant benefits. Indeed, one or more monitoring stations from such a
27 region may be removed and the resources can be used in regions where additional stations
28 might offer greater benefits (e.g. in regions where the clustering coefficient values are low).
29 The identification of stations that play more influential roles in the network, as is reflected by
30 the degree centrality results, can also be useful in identifying stations/regions around which
31 interpolation/extrapolation might work better.

1 Finally, the present study and the results obtained have important implications for a wide
2 range of issues and associated efforts in streamflow modeling, and hydrologic modeling in
3 general. Among these are: (1) predictions in ungaged basins (PUB), where approaches based
4 on nearest neighbors, regionalization, similarity, and other concepts are commonly adopted;
5 (2) formulation of a catchment classification framework, for simplification and generalization
6 in our modeling paradigm and better communication among/between researchers and
7 practitioners; and (3) development of an integrated framework for water planning and
8 management, including in studies on climate change impacts on water resources, that involves
9 proper consideration and inclusion of stakeholders and concepts from a vast number of
10 disciplines, including climate, hydrology, engineering, environment, ecology, social sciences,
11 political sciences, economics, and psychology. In view of these, ideas gained from the
12 modern theory of complex networks, and network theory at large, seem to have immense
13 potential in hydrology and water resources.

14 **Acknowledgments**

15 Support for this work was provided by the Australian Research Council (ARC). Bellie
16 Sivakumar acknowledges the financial support from ARC through the Future Fellowship
17 grant (FT110100328). We thank the two reviewers, Stefania Scarsoglio and Stacey Archfield,
18 for their constructive comments/suggestions, which have helped improve the quality and
19 presentation of our work further.

20

21 **References**

- 22 Albert, R., Jeong, H., and Barabasi A.-L.: Internet: Diameter of the world wide web, *Nature*,
23 401, 130–131, 1999.
- 24 Anderson, C. J.: Central Limit Theorem, the Corsini Encyclopedia of Psychology, John Wiley
25 & Sons, 2010.
- 26 Archfield, S. A. and Vogel, R. M.: Map correlation method: Selection of a reference
27 streamgauge to estimate daily streamflow at ungaged catchments, *Water Resour. Res.*, 46,
28 W10513, doi:10.1029/2009WR008481.
- 29 Barabási, A.-L. and Albert, R.: Emergence of scaling in random networks, *Science*, 286, 509–
30 512, 1999.

- 1 Berge, C.: *The Theory of Graphs and Its Applications*, Matheun, Ann Arbor, MI, USA, 1962.
- 2 Beven, K. J.: Uncertainty and the detection of structural change in models of environmental
3 systems, in: *Environmental Foresight and Models: a Manifesto*, edited by: Beck, M. B.,
4 Elsevier Science Ltd, Oxford, UK, 227–250, 2002.
- 5 Beven, K. J.: *Benchmark papers in Streamflow Generation Processes*, IAHS Press,
6 Wallingford, UK, 2006.
- 7 Boers, N., Bookhagen, B., Marwan, N., Kurths, J., and Marengo, J.: Complex networks
8 identify spatial patterns of extreme rainfall events of the South American Monsoon System,
9 *Geophys. Res. Lett.*, 40, 4386–4392, doi:10.1002/grl.50681, 2013.
- 10 Bollobás, B.: *Modern Graph Theory*, Springer, New York, USA, 1998.
- 11 Bondy, J. A. and Murty, U. S. R.: *Graph Theory with Applications*, Elsevier Science Ltd,
12 New York, USA, 1976.
- 13 Bouchaud, J.-P. and Mézard, M.: Wealth condensation in a simple model of economy,
14 *Physica A*, 282, 536–540, 2000.
- 15 Cayley, A.: On the theory of the analytical forms called trees, *Philos. Mag.*, 13, 172–176,
16 1857.
- 17 Dalin, C., Konar, M., Hanasaki, N., Rinaldo, A., and Rodriguez-Iturbe, I.: Evolution of the
18 global virtual water trade network, *Proc. Natl. Acad. Sci. USA*, 109, 5989–5994, 2012.
- 19 Davis, K. F., D’Odorico, P., Laio, F., and Ridolfi, L.: Global spatio-temporal patterns in
20 human migration: A complex network perspective, *PLoS ONE*, 8, e53723, doi:
21 10.1371/journal.pone.0053723, 2013.
- 22 Duan, Q., Gupta, H. V., Sorooshian, S., Rousseau, A. N., Turcotte, R.: *Calibration of*
23 *Watershed Models*, Water Science and Application Series, vol. 6, American Geophysical
24 Union, Washington, DC, USA, 2003.
- 25 Erdős, P. and Rényi, A.: On the evolution of random graphs, *Publ. Math. Inst. Hung. Acad.*
26 *Sci.*, 5, 17–61, 1960.
- 27 Euler, L.: *Solutio problematis ad geometriam situs pertinentis*, *Comment. Acad. Sci.*
28 *Petropolitanae*, 8, 128–140, 1741.

- 1 Freeman, L. C.: Centrality in social networks: conceptual clarification. *Social Networks*, 1,
2 215–239, 1978/79.
- 3 Girvan, M., and Newman, M. E. J.: Community structure in social and biological networks,
4 *Proc. Natl. Acad. Sci. USA*, 99, 7821–7826, 2002.
- 5 Grayson, R. B. and Blöschl, G.: *Spatial Patterns in Catchment Hydrology: Observations and*
6 *Modeling*, Cambridge University Press, Cambridge, UK, 2000.
- 7 Gupta, V. K., Rodriguez-Iturbe, I., and Wood, E. F.: *Scale Problems in Hydrology: Runoff*
8 *Generation and Basin Response*, Water Science and Technology Library Series, Springer,
9 Dordrecht, Holland, 1986.
- 10 Hrachowitz, M., Savenije, H. H. G., Blöschl, G., McDonnell, J. J., Sivapalan, M., Pomeroy, J.
11 W., Arheimer, B., Blume, T., Clark, M. P., Ehret, U., Fenicia, F., Freer, J. E., Gelfan, A.,
12 Gupta, H. V., Hughes, D. A., Hut, R. W., Montanari, A., Pande, S., Tetzlaff, D., Troch, P. A.,
13 Uhlenbrook, S., Wagener, T., Winsemius, H. C., Woods, R. A., Zehe, E., and Cudennec, C.:
14 A decade of predictions in ungaged basins (PUB) – a review, *Hydrol. Sci. J.*, 58, 1198–1255,
15 2013.
- 16 Kahya, E. and Dracup, J. A.: U.S. streamflow patterns in relation to the El Niño/Southern
17 Oscillation, *Water Resour. Res.*, 29, 2491–2503, 1993.
- 18 Kiang, J. E., Stewart, D. W., Archfield, S. A., Osborne, E. B. and Eng., K.: A national
19 streamflow network gap analysis, US Geological Survey Scientific Investigations Report
20 2013–5013, Reston, Virginia, USA, 2013.
- 21 Kirchner, J. W.: Getting the right answers for the right reasons: Linking measurements,
22 analyses, and models to advance the science of hydrology, *Water Resour. Res.*, 42, W03S04,
23 doi:10.1029/2005WR004362.
- 24 Konar, M. and Caylor, K. K.: Virtual water trade and development in Africa, *Hydrol. Earth*
25 *Syst. Sci.*, 17, 3969–3982, 2013.
- 26 Li, C., Singh, V. P., and Mishra, A. K.: Entropy theory-based criterion for hydrometric
27 network evaluation and design: maximum information minimum redundancy, *Water Resour.*
28 *Res.*, 48, W05521, doi:10.1029/2011WR011251, 2012.
- 29 Liljeros, F., Edling, C., Amaral, L. N., Stanley, H. E., and Åberg, Y.: The web of human
30 sexual contacts, *Nature*, 411, 907–908, 2001.

1 Lins, H. F.: USGS Hydro-climatic data network 2009 (HCDN–2009): U.S. Geological Survey
2 Fact Sheet 2012–3047,

3 Listing, J. B.: Vorstudien zur Topologie. Vandenhoeck und Ruprecht: Göttingen, Germany,
4 pp. 811–875, 1848.

5 Milo, R., Shen-Orr, S., Itzkovitz, S., Kashtan, N., Chklovskii, D., and Alon, U.: Network
6 motifs: simple building blocks of complex networks, *Science*, 298, 824–827, 2002.

7 Mishra, A. K. and Coulibaly, P.: Developments in hydrometric network design: a review,
8 *Rev. Geophys.*, 47, RG2001, doi:2007RG000243, 2009.

9 Newman, M. E. J.: The structure of scientific collaboration networks, *Proc. Nat. Acad. Sci.*
10 *USA*, 98, 404–409, 2001.

11 Patil, S. and Stieglitz, M.: Controls on hydrologic similarity: role of nearby gauged
12 catchments for prediction at an ungauged catchment, *Hydrol. Earth Syst. Sci.*, 16, 551–562,
13 2012.

14 Perrin, C., Michel, C., and Andréassian, V.: Does a large number of parameters enhance
15 model performance? Comparative assessment of common catchment model structures on 429
16 catchments, *J. Hydrol.*, 242, 275–301, 2001.

17 Rinaldo, A., Banavar, J. R., and Maritan, A.: Trees, networks, and hydrology, *Water Resour.*
18 *Res.*, 42, W06D07, doi:10.1029/2005WR004108, 2006.

19 Salas, J. D., Delleur, J. W., Yevjevich, V., and Lane, W. L.: *Applied Modeling of Hydrologic*
20 *Time Series*, Water Resources Publications, Littleton, Colorado, USA, 1995.

21 Scarsoglio, S., Laio, F., and Ridolfi, L.: Climate dynamics: A network-based approach for the
22 analysis of global precipitation, *PLoS ONE*, 8, e71129, doi:10.1371/journal.pone.0071129,
23 2013.

24 Sivakumar, B.: Forecasting monthly streamflow dynamics in the western United States: a
25 nonlinear dynamical approach, *Environ. Modell. Softw.*, 18, 721–728, 2003.

26 Sivakumar, B.: Dominant processes concept, model simplification and classification
27 framework in catchment hydrology, *Stoch. Environ. Res. Risk Assess.*, 22, 737–748, 2008.

- 1 Sivakumar, B.: Chaos theory for modeling environmental systems: Philosophy and
2 pragmatism, edited by: Wang, L. and Garnier, H., System Identification, Environmental
3 Modelling, and Control System Design, Springer-Verlag, London, 533–555, 2011c.
- 4 Sivakumar, B.: Networks: a generic theory for hydrology?, *Stoch. Environ. Res. Risk A.*,
5 doi:10.1007/s00477-014-0902-7, in press, 2014.
- 6 Sivakumar, B., and Singh, V. P.: Hydrologic system complexity and nonlinear dynamic
7 concepts for a catchment classification framework, *Hydrol. Earth Syst. Sci.*, 16, 4119–4131,
8 2012.
- 9 Sivakumar, B., Singh, V. P., Berndtsson, R., and Khan, S. K.: Catchment classification
10 framework in hydrology: challenges and directions, *J. Hydrol. Eng.*, A4014002,
11 doi:10.1061/(ASCE)E.1943-5584.0000837, 2014.
- 12 Slack, J. R. and Landwehr, V. M.: Hydro-climatic data network (HCDN): a US Geological
13 Survey streamflow data set for the United States for the study of climate variations, 1847–
14 1988, US Geological Survey Open File Report 92–129, Reston, Virginia, USA, 1992.
- 15 Suweis, S., Konar, M., Dalin, C., Hanasaki, N., Rinaldo, A., and Rodriguez-Iturbe, I.:
16 Structure and controls of the global virtual water trade network, *Geophys. Res. Lett.*, 38,
17 L10403, doi:10.1029/2011GL046837, 2011.
- 18 Tootle, G. A. and Piechota, T. C.: Relationships between Pacific and Atlantic ocean sea
19 surface temperatures and U.S. streamflow variability, *Water Resour. Res.*, 42, W07411,
20 doi:10.1029/2005WR004184, 2006.
- 21 Tsonis, A. A. and Roebber, P. J.: The architecture of the climate network, *Physica A*, 333,
22 497–504, 2004.
- 23 Vogel, R. M. and Sankarasubramanian, A.: Spatial scaling properties of annual streamflow in
24 the United States, *Hydrol. Sci. J.*, 45(3), 465–476, 2000.
- 25 Wasserman, S. and Faust, K.: *Social Network Analysis*, Cambridge University Press,
26 Cambridge, UK, 1994.
- 27 Watts, D. J. and Strogatz, S. H.: Collective dynamics of ‘small-world’ networks, *Nature*, 393,
28 440–442, 1998.
- 29 Yang, D., Li, C., Hu, H., Lei, Z., Yang, S., Kusuda, T., Koike, T., and Musiake, K.: Analysis

1 of water resources variability in the Yellow River of China during the last half century using
2 historical data, *Water Resour. Res.*, 40, W06502, doi:10.1029/2003WR002763, 2004.

3 Young, P. C. and Ratto, M.: A unified approach to environmental systems modeling, *Stoch.*
4 *Environ. Res. Risk Assess.*, 23, 1037–1057, 2009.

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

1 **Table 1.** Clustering coefficient values for monthly streamflow data from the United States

| Clustering coefficient range | Number of stations within each clustering coefficient range for threshold (T) | | | | | | | |
|------------------------------|---|------|------|------|------|------|------|------|
| | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.75 | 0.8 | 0.85 |
| 0.76–1.0 | 415 | 327 | 242 | 192 | 190 | 186 | 179 | 171 |
| 0.51–0.75 | 219 | 291 | 372 | 398 | 348 | 323 | 274 | 197 |
| 0.26–0.50 | 3 | 12 | 10 | 30 | 62 | 72 | 82 | 114 |
| 0.01–0.25 | 0 | 1 | 2 | 0 | 2 | 1 | 1 | 4 |
| 0 | 0 | 2 | 3 | 4 | 2 | 4 | 10 | 18 |
| NA | 2 | 6 | 10 | 15 | 35 | 53 | 93 | 135 |
| Entire Network | 0.79 | 0.76 | 0.71 | 0.68 | 0.65 | 0.63 | 0.58 | 0.51 |

2
3
4
5
6
7
8
9
10
11
12
13

1 **Figure Captions**

2

3 **Figure 1.** Network in its simplest form, i.e. an undirected network with only a single type of
4 node and a single type of link.

5 **Figure 2.** Connections in networks and calculation of clustering coefficient: nearest
6 neighbors and actual connections.

7 **Figure 3.** Characteristics of monthly streamflow observed at 639 stations in the United
8 States: **(a)** mean; **(b)** standard deviation; and **(c)** coefficient of variation.

9 **Figure 4.** Linear correlation for streamflow: average of correlation with **(a)** all the 638
10 stations; **(b)** nearest 30 neighbors; **(c)** nearest 15 neighbors; and **(d)** nearest 5 neighbors

11 Figure 5. Degree centrality for four correlation thresholds: **(a)** 0.70; **(b)** 0.75; **(c)** 0.80; and
12 **(d)** 0.85.

13 Figure 6. Clustering coefficients for four correlation thresholds: **(a)** 0.70; **(b)** 0.75; **(c)** 0.80;
14 and **(d)** 0.85. The six ranges are chosen for better visualization of results.

15 **Figure 7a.** Links in streamflow network for threshold $T = 0.75$. Four nodes (stations) are
16 chosen for better visualization.

17 **Figure 7b.** Links in streamflow network for threshold $T = 0.80$. Four nodes (stations) are
18 chosen for better visualization.

19 **Figure 7c.** Links in streamflow network for threshold $T = 0.85$. Four nodes (stations) are
20 chosen for better visualization.

21 **Figure 8.** Relationship between clustering coefficient and number of links: **(a)** $T = 0.70$; **(b)**
22 $T = 0.75$; **(c)** $T = 0.80$; and **(d)** $T = 0.85$. Both *all links* (red circles) and *actual links* (blue
23 circles) are presented.

24 **Figure 9.** Relationship between clustering coefficient and number of nearest neighbors: **(a)** T
25 $= 0.70$; **(b)** $T = 0.75$; **(c)** $T = 0.80$; and **(d)** $T = 0.85$.

26

27

28