

¹ Algorithms for spanning trees of ² unweighted networks

³ Lovro Šubelj

⁴ University of Ljubljana, Faculty of Computer and Information Science, Ljubljana,
⁵ Slovenia

⁶ Corresponding author:

⁷ Lovro Šubelj

⁸ Email address: lovro.subelj@fri.uni-lj.si

⁹ ABSTRACT

¹⁰ Spanning tree of a network or a graph is a subgraph connecting all the nodes with the minimum number
¹¹ of edges. Spanning tree retains the connectivity of a network and possibly other structural properties, and
¹² is one of the simplest techniques for network simplification or sampling, and for revealing its backbone or
¹³ skeleton. The Prim's algorithm and the Kruskal's algorithm are well known algorithms for computing a
¹⁴ spanning tree of a weighted network. In this paper, we study the performance of these algorithms on
¹⁵ unweighted networks, and compare them to different priority-first search algorithms. We show that the
¹⁶ distances between the nodes and the diameter of a network are best preserved by an algorithm based
¹⁷ on the breadth-first search node traversal. The algorithm computes a spanning tree with properties of
¹⁸ a balanced tree and a power-law node degree distribution. We support our results by experiments on
¹⁹ synthetic graphs and more than a thousand real networks, and demonstrate different practical applications
²⁰ of computed spanning trees. We conclude that, if a spanning tree is supposed to retain the distances
²¹ between the nodes or the diameter of an unweighted network, then the breadth-first search algorithm
²² should be the preferred choice.

²³ **Keywords** Unweighted networks, Spanning tree, Balanced tree, Breadth-first search, Network dis-
²⁴ tances, Network diameter

²⁵ INTRODUCTION

²⁶ Networks or graphs have become a popular tool for analysing complex real-world systems (Newman,
²⁷ 2018). Examples include predicting the spread of contagious viruses (Tizzoni et al., 2012), study of the
²⁸ interactome of species (Zitnik et al., 2019), understanding the structure of science (Fortunato et al., 2018)
²⁹ and outreach of social connections online (Backstrom et al., 2012). The sizes of today's networks are
³⁰ often in millions of nodes and edges, with the largest networks being the WWW with more than a trillion
³¹ web pages and human brain with close to a hundred of billions of neurons. In result, the size of real
³² networks can make many practical applications computationally very demanding.

³³ Techniques to alleviate this issue include network simplification or sampling (Leskovec and Faloutsos,
³⁴ 2006; Hamann et al., 2016; Blagus et al., 2017), and revealing the so-called network backbone or
³⁵ skeleton (Grady et al., 2012; Coscia and Neffke, 2017; Šubelj, 2018). These approaches try to reduce the
³⁶ size of a network in a way that the network still retains many of its structural properties. One of the most
³⁷ straightforward ways to simplify a network is to compute its spanning tree (Bollobás, 1998; Newman,
³⁸ 2018), which is a subgraph connecting all the nodes of a network with the minimum number of edges. A
³⁹ spanning tree retains the connectivity of a network and possibly other structural properties. In the case
⁴⁰ of weighted networks, one usually aims to compute the minimum spanning tree, which is a subgraph
⁴¹ connecting all the nodes with the minimum overall weight of the edges. In the case of unweighted
⁴² networks, any spanning tree is in fact a "minimum" spanning tree.

⁴³ The Prim's algorithm and the Kruskal's algorithm are well known algorithms for computing a
⁴⁴ minimum spanning tree of a weighted network (Bollobás, 1998; Newman, 2018). Although developed
⁴⁵ primarily for weighted networks, the algorithms can be readily applied to unweighted networks. However,
⁴⁶ the performance of these algorithms has not yet been properly explored for unweighted networks. In this

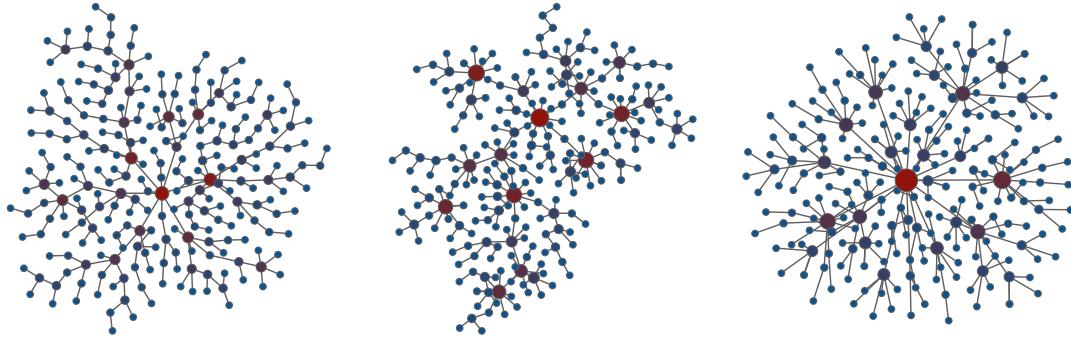


Figure 1. Wiring diagrams of spanning trees of a small random graph. The spanning trees were computed with the Prim's algorithm (left), the Kruskal's algorithm (middle) and the breadth-first search algorithm (right). The sizes of the nodes are proportional to their degrees, while the layouts were computed with the Large Graph Layout algorithm (Adai et al., 2004).

47 paper, we apply the algorithms to more than a thousand real networks and compare them to different
 48 priority-first search algorithms. We show that the structure of unweighted networks is best preserved by
 49 an algorithm using the breadth-first search node traversal. In particular, spanning trees computed with the
 50 breadth-first search algorithm well retain the distances between the nodes of the network¹ and exhibit
 51 power-law node degree distribution, which may be desirable in practical applications such as network
 52 visualization.

53 The rest of the paper is structured as follows. In the following section, we first describe different
 54 algorithms for computing a spanning tree of a network. Next, we analyze the structure of spanning trees of
 55 synthetic graphs and, afterwards, the structure of spanning trees of real networks. Finally, we demonstrate
 56 practical applications of spanning trees of networks. We conclude the paper with suggestions for future
 57 research.

58 COMPUTATION OF SPANNING TREES

59 Let a network be represented by an undirected connected graph $G = (V, E)$, where V denotes the set of
 60 nodes of G and E denotes the set of edges of G . (Where necessary to make explicit that these sets represent
 61 graph G , we write V_G and E_G .) The number of nodes equals $n = |V|$ and the number of edges equals
 62 $m = |E|$. We denote the average node degree as $\langle k \rangle = 2m/n$. Furthermore, let d_{ij} be the distance between
 63 the nodes $i, j \in V$, defined as the number of edges in the shortest paths between the nodes i and j . Since
 64 the graph is undirected and connected, $d_{ij} = d_{ji}$ and $d_{ij} < \infty$. Therefore, the average distance between the
 65 nodes equals $\langle d \rangle = \frac{2}{n(n-1)} \sum_{i < j} d_{ij}$ and the maximum distance or diameter equals $d_{\max} = \max_{i < j} d_{ij}$. To
 66 measure the variability of the distances between the nodes, we also define the coefficient of variation as
 67 $c_d = \sigma_d / \langle d \rangle$, where σ_d is the standard deviation of the distances.

68 Below we describe different algorithms for computing a spanning tree of an undirected connected
 69 graph. In the case of a disconnected graph consisting of more than one connected component, the
 70 algorithms should be applied to each of the connected components separately. For a more extensive
 71 discussion of the algorithms, reader is referred to classical graph theory literature (Bollobás, 1998) or
 72 network science literature (Barabási, 2016; Newman, 2018).

73 Prim's algorithm

74 The Prim's algorithm for computing a spanning tree T of an undirected connected graph G operates as
 75 follows (see Algorithm 1). First, the algorithm selects a random seed node $i \in V_G$ from graph G and adds
 76 it to an empty tree T (lines 2, 3). The node i serves as a starting point for computing the spanning tree T .
 77 Then, on each step of the algorithm (lines 4–8), a random edge $\{i, j\} \in E_G$ from graph G is selected that
 78 leads from a node $i \in V_T$ already in the tree T to a node $j \notin V_T$ not yet in the tree T (line 5). Both node j
 79 and edge $\{i, j\}$ are added to the tree T (lines 6, 7). Finally, when there is no further node $i \in V_G$ in graph

¹The breadth-first search algorithm is otherwise a standard approach for computing the distances between one selected node and all other nodes in an undirected network.

Algorithm 1 Prim's algorithm

Require: undirected graph G
Ensure: spanning tree T

```
1:  $T \leftarrow$  empty graph
2:  $i \leftarrow \text{RANDOM}(i \in V_G)$                                  $\triangleright$  Select seed node.
3: add node  $i$  to  $V_T$                                           $\triangleright$  Add selected seed node.
4: while  $\exists i \in V_G : i \notin V_T$  do                          $\triangleright$  There exists non-visited node?
5:    $\{i, j\} \leftarrow \text{RANDOM}(\{i, j\} \in E_G : i \in V_T \wedge j \notin V_T)$      $\triangleright$  Edge to non-visited node.
6:   add node  $j$  to  $V_T$                                           $\triangleright$  Add non-visited node.
7:   add edge  $\{i, j\}$  to  $E_T$                                       $\triangleright$  Add selected edge.
8: end while
9: return  $T$ 
```

80 G such that node $i \notin V_T$ is not already in the tree T , the algorithm stops (line 4). At this point, the tree T
81 is a spanning tree of graph G (line 9).

82 The Prim's algorithm is non-deterministic and can compute different spanning trees. The actual
83 spanning tree depends on random selection of the seed node (line 2 of Algorithm 1) and on random
84 selection of the edges to expand the tree (line 5). The latter is nontrivial and most efficiently implemented
85 by rejection sampling over an array list of edges to non-visited nodes. For simplicity, we do not make
86 these computations explicit in Algorithm 1.

87 Assume that the graph is represented with an adjacency list. In the case of weighted graphs, the time
88 complexity of the Prim's algorithm implemented with a Fibonacci heap is $\mathcal{O}(m + n \log n)$. For unweighted
89 graphs which we consider here, the heap can be replaced by a simple array list, which reduces the time
90 complexity to $\mathcal{O}(m)$. As an example, the left graph in Figure 1 shows a spanning tree computed with the
91 Prim's algorithm.

92 Kruskal's algorithm

93 The Kruskal's algorithm is conceptually different from the Prim's algorithm. Instead of starting with
94 a tree consisting of a seed node and then expanding it, the algorithm starts with a forest of trees, each
95 consisting of a single node. The trees are then incrementally merged into larger trees by adding edges
96 between them until only one tree remains. At this point the algorithm stops and the remaining tree is a
97 spanning tree of a graph.

98 The Kruskal's algorithm is non-deterministic and the actual spanning tree depends on random selection
99 of the edges to merge the trees on each step. The time complexity of the algorithm using a disjoint-set
100 data structure is $\mathcal{O}(m \log n)$, for either weighted or unweighted graphs. As an example, the middle graph
101 in Figure 1 shows a spanning tree computed with the Kruskal's algorithm.

102 Breadth-first search

103 The breadth-first search node traversal is very similar to the Prim's algorithm. The main difference is in
104 how the edges to non-visited nodes are processed. In contrast to the Prim's algorithm, where only one
105 such edge is processed on each step, the breadth-first search processes all edges from a selected node to
106 non-visited nodes in a single step.

107 The breadth-first search algorithm for computing a spanning tree T of an undirected connected graph
108 G operates as follows (see Algorithm 2). In contrast to before, we make all computations in Algorithm 2
109 explicit. First, the algorithm selects a random seed node $i \in V_G$ from graph G and adds it to an empty tree
110 T (lines 3, 4). The node i is also added to an empty queue $Q \subseteq V_T$ for further processing. Then, on each
111 step of the algorithm (lines 5-11), a node $i \in Q$ is removed from the beginning of the queue Q (line 6)
112 and all edges that lead from node $i \in V_T$ already in the tree T to nodes $j \notin V_T$ not yet in the tree T are
113 processed (lines 7-10). All nodes j and edges $\{i, j\}$ are added to the tree T (lines 8, 9), while nodes j are
114 also added to the queue Q for further processing. Finally, when there is no further node $i \in Q$ in the queue
115 Q , the algorithm stops (line 5). At this point the tree T is a spanning tree of graph G (line 12).

116 The breadth-first search algorithm is again non-deterministic, while the actual spanning tree depends
117 on random selection of the seed node (line 3 of Algorithm 2) and on the exact order in which the edges
118 to expand the tree are processed (line 7). The time complexity of the algorithm using a queue of non-
119 processed nodes is $\mathcal{O}(m)$, for either weighted or unweighted graphs. As an example, the right graph

Algorithm 2 Breadth-first search

Require: undirected graph G
Ensure: spanning tree T

```
1:  $T \leftarrow$  empty graph
2:  $Q \leftarrow$  empty queue
3:  $i \leftarrow \text{RANDOM}(i \in V_G)$ 
4: add node  $i$  to  $V_T$  and  $Q$ 
5: while  $\exists i \in Q$  do
6:    $i \leftarrow$  remove node from  $Q$ 
7:   for  $\{i, j\} \in E_G : j \notin V_T$  do
8:     add node  $j$  to  $V_T$  and  $Q$ 
9:     add edge  $\{i, j\}$  to  $E_T$ 
10:  end for
11: end while
12: return  $T$ 
```

120 in Figure 1 shows a spanning tree computed with the breadth-first search algorithm.

121 **Depth-first search**

122 For completeness, we also describe the depth-first search node traversal. While the breadth-first search
123 algorithm processes nodes of a graph using a level order traversal, the depth-first search algorithm uses a
124 preorder traversal. This means that the only change required to the breadth-first search algorithm is to
125 replace the queue of non-processed nodes Q with a stack (line 2 in Algorithm 2). This is because the
126 queue Q operates on a first-in first-out basis, while a stack operates on a last-in first-out basis (lines 6, 8).
127 The time complexity of the depth-first search algorithm is again $\mathcal{O}(m)$.

128 **SPANNING TREES OF SYNTHETIC GRAPHS**

129 Consider an Erdős-Rényi random graph (Erdős and Rényi, 1959) with n nodes and the probability of an
130 edge between each pair of nodes $p = \langle k \rangle / (n - 1)$, where $\langle k \rangle$ is the expected node degree. A spanning tree
131 of any connected graph with n nodes consists of n nodes and $n - 1$ edges with the average node degree
132 $\langle k \rangle = 2 - 2/n$. Since it is a tree, the average clustering coefficient equals $\langle C \rangle = 0$ (Watts and Strogatz,
133 1998). We, therefore, here focus on other graph properties. In particular, we study the average distance
134 between the nodes $\langle d \rangle$ and the diameter d_{\max} . A theoretical estimate for the diameter d_{\max} of a random
135 graph equals $\log n / \log \langle k \rangle$ (Newman, 2018), which is $d_{\max} = 2.40$ for $n = 250$ and $\langle k \rangle = 10$. Due to the
136 sensitivity of the diameter d_{\max} for relatively small n and $\langle k \rangle$, this turns out to be a better estimate of the
137 average distance between the nodes $\langle d \rangle$. Indeed, the empirical estimates for the considered random graph
138 are $\langle d \rangle \approx 2.64$ and $d_{\max} \approx 4.39$.

139 Figure 1 shows particular realizations of spanning trees of a random graph with the above parameters
140 computed with the Prim's algorithm, the Kruskal's algorithm and the breadth-first search algorithm. The
141 diameter d_{\max} of the spanning trees equals 14, 17 and 6, respectively. While the diameters of the spanning
142 trees computed with Prim's algorithm and the Kruskal's algorithm are much higher than in the random
143 graph, the diameter of the spanning tree computed with the breadth-first search algorithm is very close
144 to the diameter of the random graph. These observations are closely related to the question whether the
145 computed spanning trees are balanced (Knuth, 2011).

146 The average distance $\langle d \rangle$ and the diameter d_{\max} of a balanced tree are in $\mathcal{O}(\log n)$ for any practical
147 definition of balance (Knuth, 2011). However, in the case of a random tree, both values are almost
148 certainly in $\mathcal{O}(\sqrt{n})$ (Rényi and Szekeres, 1967; Meir and Moon, 1970). Since these results only talk
149 about the scaling, they can not be directly employed to measure whether a particular spanning tree is
150 balanced or not, and to what extent. To the best of our knowledge, no such approach exists. One can,
151 however, study the scaling of the average distance $\langle d \rangle$ and the diameter d_{\max} of spanning trees of graphs
152 with an increasing number of nodes n and try to empirically estimate whether the values scale as $\mathcal{O}(\log n)$
153 or worse. Note that only in the case of the former the spanning trees can possibly retain short distances
154 between the nodes in random graphs and also real small-world networks (Watts and Strogatz, 1998).

155 Besides Erdős-Rényi random graphs (Erdős and Rényi, 1959), we also analyse triangular lattices and
156 Barabási-Albert scale-free graphs (Barabási and Albert, 1999). We vary the number of nodes n , while

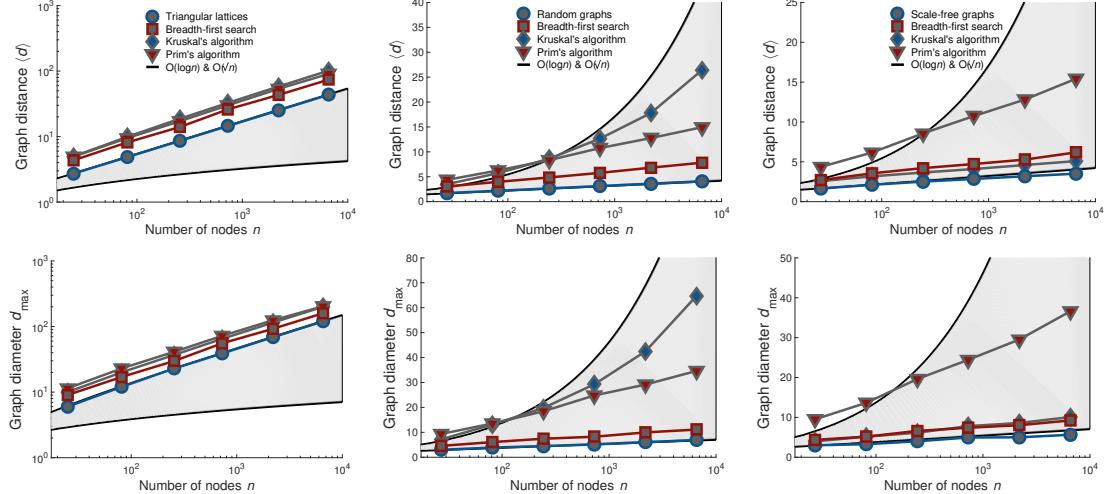


Figure 2. The average distance $\langle d \rangle$ and the diameter d_{\max} of triangular lattices (left), random graphs (middle) and scale-free graphs (right), and their spanning trees computed with different algorithms (legend). The plots show estimates over 100 realizations, while the shaded areas span between theoretical estimates for random graphs $\mathcal{O}(\log n)$ and two-dimensional lattices $\mathcal{O}(\sqrt{n})$, and are consistent between the plots.

157 we keep the average node degree fixed to $\langle k \rangle = 10$. Figure 2 shows the scaling of the average distance
158 $\langle d \rangle$ and the diameter d_{\max} for selected synthetic graphs and their spanning trees computed with different
159 algorithms.

160 We first consider triangular lattices, as these results serve as a baseline for further analyses. The
161 average distance $\langle d \rangle$ and the diameter d_{\max} of any two-dimensional lattice scale as $\mathcal{O}(\sqrt{n})$ (Newman,
162 2018). This can be observed as a straight line with slope 0.5 on double-logarithmic plots in the left
163 column of Figure 2. Notice also that all the spanning trees computed with different algorithms show
164 similar scaling $\mathcal{O}(\sqrt{n})$.

165 Next, we consider Erdős-Rényi random graphs (Erdős and Rényi, 1959) shown in the middle column
166 of Figure 2. The average distance $\langle d \rangle$ and the diameter d_{\max} of random graphs, and also small-world
167 networks (Watts and Strogatz, 1998), scale as $\mathcal{O}(\log n)$ (Newman, 2018). This can be observed as a straight
168 line on semi-logarithmic plots in Figure 2, whereas any upward concave function would imply scaling
169 faster than $\mathcal{O}(\log n)$. Notice that the spanning trees computed with the breadth-first search algorithm best
170 retain the distances in random graphs, while both the average distance $\langle d \rangle$ and the diameter d_{\max} appear
171 to scale as $\mathcal{O}(\log n)$, at least for moderate number of nodes $n \leq 10^4$. In contrast, the distances between
172 the nodes of the spanning trees computed with the Kruskal's algorithm scale faster than $\mathcal{O}(\log n)$ (see
173 also Figure 3 and discussion below).

174 Last, we consider Barabási-Albert scale-free graphs (Barabási and Albert, 1999) shown in the
175 right column of Figure 2. The average distance $\langle d \rangle$ and the diameter d_{\max} of scale-free graphs scale
176 as $\mathcal{O}(\log n / \log \log n)$ (Cohen and Havlin, 2003), while such graphs are usually called ultra small-
177 world (Barabási, 2016). Note that $\mathcal{O}(\log n / \log \log n)$ is indistinguishable from $\mathcal{O}(\log n)$ for $n \leq 10^4$,
178 thus this scaling can again be observed as a straight line on semi-logarithmic plots in Figure 2. The
179 spanning trees computed with both the breadth-first search algorithm and the Kruskal's algorithm well
180 retain the distances between the nodes of scale-free graphs and appear to scale as $\mathcal{O}(\log n)$. On the other
181 hand, the distances between the nodes of the spanning trees computed with the Prim's algorithm can be
182 more than five times larger than the distances in scale-free graphs, e.g., $\langle d \rangle = 15.42$ and $d_{\max} = 36.6$
183 compared to 3.51 and 5.64 for graphs with $n = 6561$ nodes.

184 Above observations are confirmed in Figure 3, where we show the coefficient of variation of the
185 distances between the nodes c_d . This is a standard measure of dispersion of a probability distribution,
186 while the distributions with $c < 1$ are considered low-variance distributions and those with $c > 1$ are
187 high-variance distributions. Note that the distributions of the distances between the nodes of random
188 and scale-free graphs, and real small-world networks, are low-variance with $c_d \ll 1$ (Watts and Strogatz,

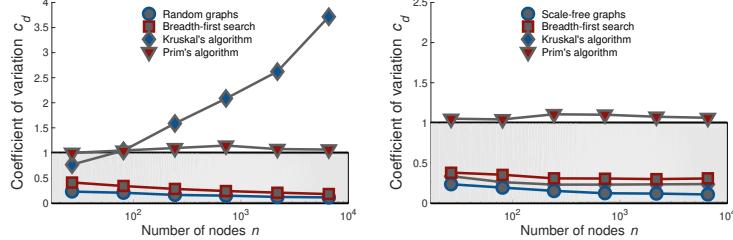


Figure 3. The coefficient of variation c_d for random (left) and scale-free graphs (right), and their spanning trees computed with different algorithms (legend). The plots show estimates over 100 realizations, while the error bars are smaller than the symbol sizes.

189 1998; Barabási, 2016). As one can see in Figure 3, all distributions of the distances in the spanning trees
 190 computed with the breadth-first search algorithm are low-variance $c_d \ll 1$. The distributions for spanning
 191 trees computed with the Kruskal's algorithm are high-variance $c_d \gg 1$ for random graphs with $n > 100$,
 192 while the results for the Prim's algorithm are inconclusive $c_d \approx 1$.

193 In summary, if a spanning tree is supposed to retain the distances between the nodes of a graph, then
 194 the breadth-first search algorithm should be the preferred choice, at least for random and scale-free graphs.
 195 In the following section, we consider also real networks.

196 SPANNING TREES OF REAL NETWORKS

197 Table 1 shows statistics of collections of over a thousand real networks analyzed in the paper. These
 198 represent citations between the papers published in the Physical Review E journal (American Physical
 199 Society, 2015), paper collaborations between the authors extracted from the SICRIS database (Institute
 200 of Information Science, 2010), protein interactions of different species collected from the BioGRID
 201 repository (Stark et al., 2006; Biological General Repository for Interaction Datasets, 2016), interactions
 202 between the users at the stack exchange web site MathOverflow (Paranjape et al., 2017; Leskovec and
 203 Krevl, 2014), Facebook friendships between the students at different US universities (Traud et al., 2012;
 204 Rossi and Ahmed, 2015) and links between autonomous systems extracted by the Oregon Route Views
 205 project (Leskovec et al., 2005; Leskovec and Krevl, 2014). Some collections represent temporal networks
 206 that grow through time (*e.g.*, paper citations and author collaborations), while other represent similar
 207 networks of different size (*e.g.*, protein interactions and online friendships). All networks were reduced to
 208 a simple graph of their largest connected component.

Table 1. Statistics of collections of real networks. These are the number of networks in the collection N ,
 and the number of nodes n and edges m , the average node degree $\langle k \rangle$ and the average distance between
 the nodes $\langle d \rangle$.

Networks	N	n	m	$\langle k \rangle$	$\langle d \rangle$
Paper citations	46	[3,37511]	[2,135260]	[1.3, 7.2]	[1.33, 21.79]
Author collaborations	25	[18,1735]	[42,6710]	[4.1, 7.7]	[1.85, 8.75]
Protein interactions	40	[5,19961]	[4,238886]	[1.6, 83.1]	[1.47, 6.06]
User interactions	75	[2,20969]	[1,86137]	[1.0, 10.1]	[1.00, 3.80]
Online friendships	97	[762,41536]	[16651,1590651]	[39.1, 116.2]	[2.24, 3.21]
Autonomous systems	733	[103,6474]	[239,12572]	[3.4, 4.7]	[2.65, 3.98]

209 Figure 4 shows the average distance between the nodes $\langle d \rangle$ of real networks and their spanning trees,
 210 where we have used semi-logarithmic axes as in Figure 2. First, we consider the networks. As expected
 211 for small-world networks (Watts and Strogatz, 1998), the average distance $\langle d \rangle$ increases with the number
 212 of nodes n and appears to scale no faster than $\mathcal{O}(\log n)$ in all network collections but two. In the case of
 213 temporal networks representing paper citations and author collaborations in the first two plots of Figure 4,
 214 the average distance $\langle d \rangle$ actually starts to decrease when the number of nodes exceeds $n \approx 500$. This is
 215 a consequence of network densification known as shrinking diameter (Leskovec et al., 2005). Figure 8
 216 shows also the diameter d_{\max} of real networks, where the interpretation is exactly the same.

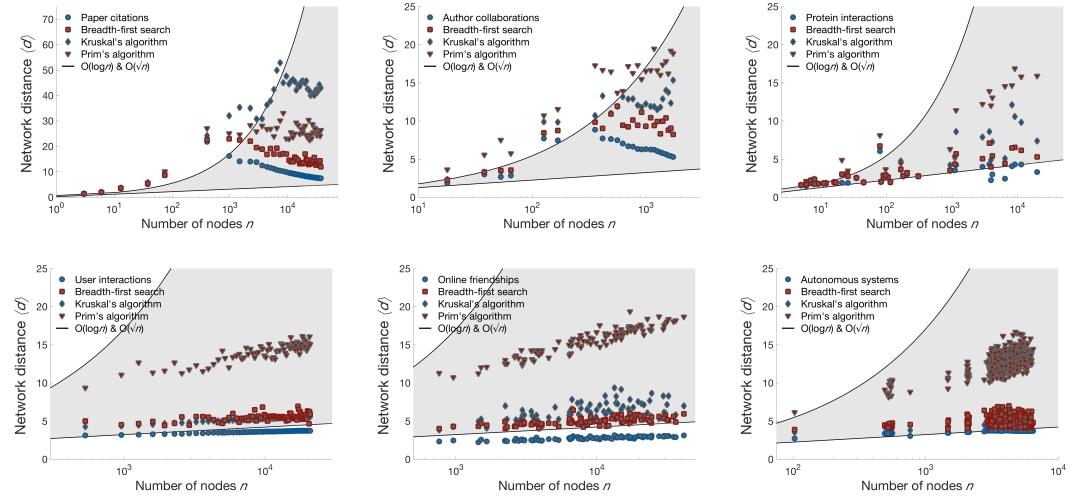


Figure 4. The average distance $\langle d \rangle$ of real networks and their spanning trees computed with different algorithms (legend), while the shaded areas are the same as in Figure 2.

217 Next, we consider spanning trees of these networks computed with different algorithms. Consistent
 218 with the results for synthetic graphs, the spanning trees computed with the breadth-first search algorithm
 219 best preserve the average distance between the nodes $\langle d \rangle$ in all network collections but two. In the case
 220 of networks representing user interactions and autonomous systems in the bottom row of Figure 4, the
 221 Kruskal's algorithm performs similarly well. Furthermore, in non-temporal networks that are not subject
 222 to densification law (Leskovec et al., 2005), the average distance $\langle d \rangle$ of the spanning trees computed with
 223 the breadth-first search algorithm, and also other algorithms, appears to scale no faster than $\mathcal{O}(\log n)$.
 224 In other networks, the scaling of the average distance $\langle d \rangle$ closely follows the scaling for real networks.
 225 Again, Figure 8 shows also the diameter d_{\max} of spanning trees, where the interpretation is exactly the
 226 same.

227 Above observations are confirmed in Figure 5, where we show the coefficient of variation of the
 228 distances between the nodes c_d . Notice that all distributions of the distances in real networks and
 229 spanning trees computed with the breadth-first search algorithm are low-variance with $c_d < 1$, as long
 230 as the networks are large enough $n \geq 10^4$. In contrast, this holds neither for the Kruskal's algorithm
 231 nor the Prim's algorithm where even $c_d \gg 1$ for some network collections (see the first and the last
 232 plot of Figure 5). It is further interesting that the node degree distribution p_k of the spanning trees
 233 computed with the breadth-first search algorithm often follows a power-law $p_k \sim k^{-\gamma}$ (Clauset et al.,
 234 2009), regardless of whether the network is scale-free or not (Barabási and Albert, 1999; Broido and
 235 Clauset, 2019) (see Figure 9 and discussion alongside).

236 To summarize, we again conclude that, if a spanning tree should retain the average distance between
 237 the nodes $\langle d \rangle$ and the diameter d_{\max} of a network, then the breadth-first search algorithm should be used.
 238 Whether preserving the distances is actually desired or favorable depends on a particular application,
 239 which we consider in the following section.

240 APPLICATIONS OF SPANNING TREES

241 A spanning tree can be seen as a technique for revealing a network backbone or skeleton (Coscia and
 242 Neffke, 2017; Šubelj, 2018), with applications in network visualization and link prediction. Furthermore, a
 243 spanning is one of the simplest approaches for network simplification or sampling (Leskovec and Faloutsos,
 244 2006; Blagus et al., 2017). Any computation that can be well approximated from a spanning tree of a
 245 network, without the need of applying the algorithms to the entire network, can provide computational
 246 benefits. In particular, many applications on networks require super-linear $\mathcal{O}(m \log n)$ or even quadratic
 247 algorithms $\mathcal{O}(mn)$, where n and m are the number of nodes and edges. These include community detection
 248 algorithms (Fortunato and Hric, 2016) and revealing node importance or similarity (Newman, 2018).
 249 In contrast, the computation of a spanning tree with the breadth-first search algorithm has linear time

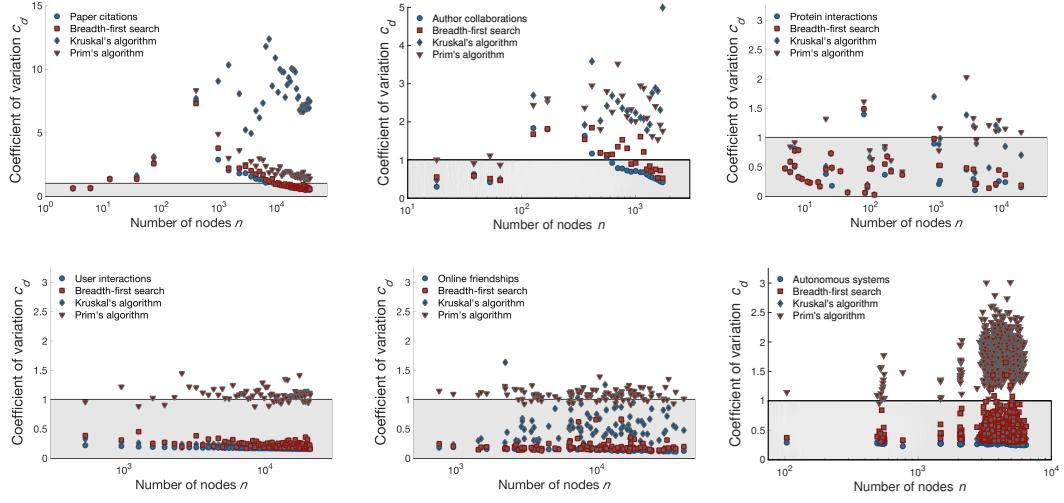


Figure 5. The coefficient of variation c_d for real networks and their spanning trees computed with different algorithms (legend), while other details are the same as in Figure 3.

complexity $\mathcal{O}(m)$ and thus does not contribute to the overall time complexity.

In this section, we consider two applications of spanning trees, where it is desired to preserve the distances between the nodes of a network.

Node importance

Revealing the importance of nodes in a network is a classical application of network science with different use cases. There exist many node measures or indices (Schoch and Brandes, 2016), known as measures of node position or centrality in the social networks analysis literature (Freeman, 1977, 1979). The measure based on the distances between the nodes is called closeness centrality, which measures the extent to which the node appears to be in the “center” of a network. The closeness centrality of a node $i \in V$ of graph $G(V, E)$ is defined as $\frac{1}{n-1} \sum_{j \neq i} d_{ij}^{-1}$ (Freeman, 1979), where $n = |V|$ is the number of nodes and d_{ij} is the distance between the nodes $i, j \in V$. The time complexity of computing closeness centrality of all nodes in a network is $\mathcal{O}(nm)$ and no more efficient algorithm exists (Knuth, 2011).

Figure 6 shows the Pearson correlation coefficient between node closeness centrality in real networks and their spanning trees computed with the breadth-first search algorithm, the Kruskal’s algorithm and the Prim’s algorithm. The coefficients for the breadth-first search algorithm are shown in the center of the heatmaps in the left column of Figure 6. These correlations are all ≥ 0.70 , which means strong linear correlation. On the other hand, the correlation coefficients for the spanning trees computed with the Kruskal’s algorithm in the middle column of Figure 6 and the Prim’s algorithm in the right column are ≤ 0.61 and ≤ 0.37 , respectively. Therefore, consistent with the previous results, the breadth-first search algorithm best preserves the distances between the nodes of real networks, even on the level of individual nodes.

For comparison, Figure 6 also shows the Pearson correlation coefficient for node degree centrality and betweenness centrality (Freeman, 1977). The latter measures the extent to which a node appears to serve as a “bridge” in a network and is defined as the fraction of the shortest paths between all pairs of nodes that go through a node. The time complexity of computing betweenness centrality of all nodes in a network is again $\mathcal{O}(nm)$ and no more efficient algorithm exists (Brandes, 2001).

Network visualization

Besides providing computational benefits, a network simplification technique such as a spanning tree can be useful in network visualization. Any visualization with a wiring diagram or some other approach is limited in the size of a network it can represent and in the structural characteristics of a network it can reveal (Ma and Mueler, 2013; Gibson et al., 2013). Classical algorithms for computing a layout of a network include spring embedding or force-directed algorithms (Eades, 1984; Fruchterman and Reingold, 1991) and algorithms that embed the nodes in an Euclidean plane thus their Euclidean distance matches

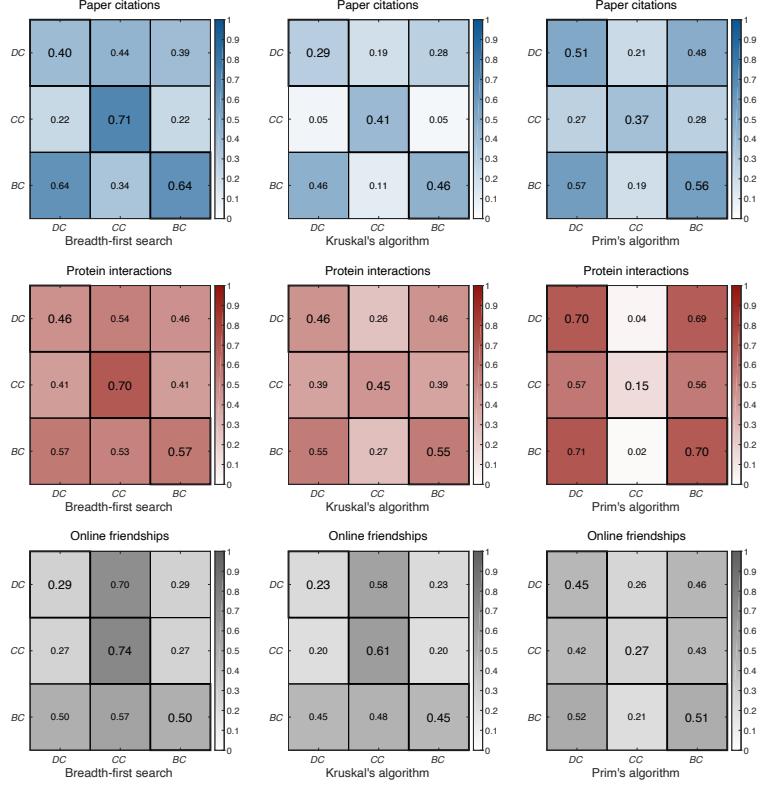


Figure 6. Pearson correlation coefficient between the measures of node centrality in real networks and their spanning trees computed with the breadth-first search algorithm (left), the Kruskal's algorithm (middle) and the Prim's algorithm (right). The measures include node degree centrality DC , closeness centrality CC and betweenness centrality BC (Freeman, 1977), while the values are estimates over 25 realizations.

their network distance as best as possible (Kamada and Kawai, 1989). It is, therefore, important that a spanning tree preserves the distances between the nodes of a network if it is to be used in network visualization.

Figure 7 shows a wiring diagram of a spanning tree of the largest connected component of the SICRIS author collaboration network (Institute of Information Science, 2010). The spanning tree was computed with the breadth-first search algorithm. The wiring diagram illustrates how authors from the same discipline cluster in certain regions and how authors from different disciplines collaborate. In contrast, the visualization of the entire network is much more involved (Adai et al., 2004), while it provides limited insight into the patterns of author collaboration. Since there exists no objective measure of the quality of a network visualization, we refrain from providing any further subjective interpretation.

CONCLUSIONS

A spanning tree is one of the most straightforward ways to network simplification or sampling, and for revealing its backbone or skeleton (Blagus et al., 2017; Coscia and Neffke, 2017). Well known algorithms for computing a spanning tree of a weighted network are the Prim's algorithm and the Kruskal's algorithm (Bollobás, 1998; Newman, 2018). However, when applied to unweighted networks, these algorithms do not necessarily capture the structure of a network such as short distances between the nodes and small network diameter (Watts and Strogatz, 1998). On the other hand, an algorithm based on the breadth-first search node traversal well retains the distances between the nodes in synthetic graphs and real networks. As we demonstrate, this can provide computational benefits and is also important in practical applications like network visualization. Thus, if a spanning tree of an unweighted network is supposed to retain the distances between the nodes, then the breadth-first search algorithm should be preferred to other algorithms considered here.

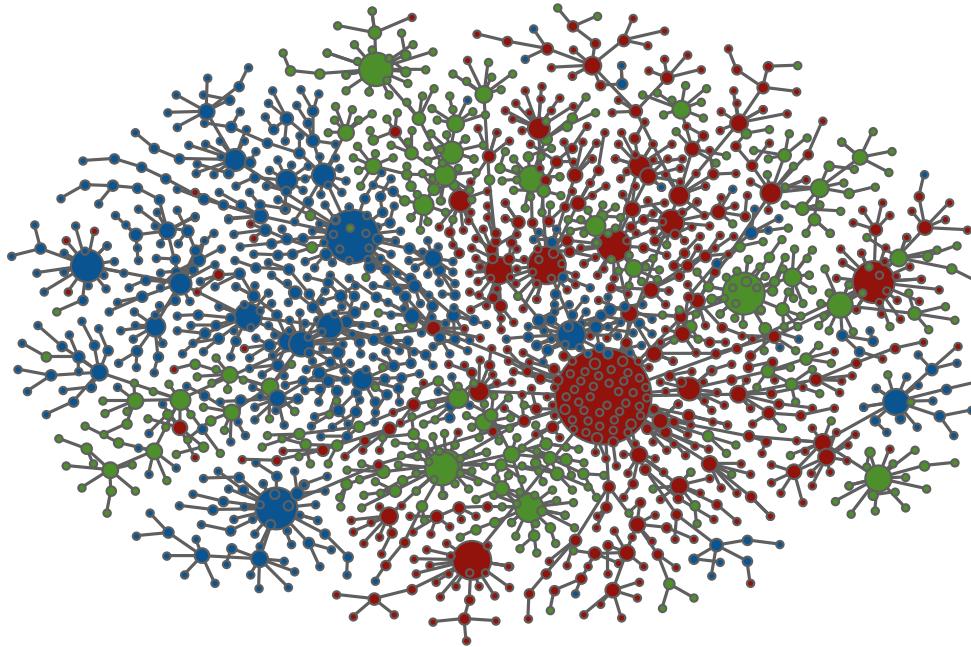


Figure 7. Spanning tree of the SICRIS author collaboration network computed with the breadth-first search algorithm. The sizes of the nodes are proportional to their degrees, while the colors represent primary author disciplines consisting of natural sciences (red), engineering (green), medical sciences (blue) and other. The layout was computed with the Large Graph Layout algorithm (Adai et al., 2004).

305 In the language of the theory of computation, the breadth-first search algorithm computes a spanning
 306 tree with properties of a balanced tree (Knuth, 2011). Due to the lack of a formal definition of some sort
 307 of approximate balance that would *already* imply short distances between the nodes, our results here are
 308 merely empirical. In future research, we plan to develop such a definition that could support the results
 309 also analytically.

310 ACKNOWLEDGMENTS

311 The authors thank Luka Kronegger for sharing the SICRIS data. This work has been supported by the
 312 Slovenian Research Agency ARRS under the program P5-0168.

313 DIAMETER AND DEGREE DISTRIBUTION

314 Figure 8 shows the diameter d_{\max} of real networks and their spanning trees computed with the Prim's
 315 algorithm, the Kruskal's algorithm and the breadth-first search algorithm. The interpretation of these
 316 results is the same as for the average distance $\langle d \rangle$ in Figure 4 and therefore omitted here.

317 Figure 9 shows the node degree distribution p_k of the largest networks in Table 1 and their spanning
 318 trees computed with the breadth-first search algorithm. Under the goodness-of-fit test at p -value =
 319 0.1 (Clauset et al., 2009), a power-law $p_k \sim k^{-\gamma}$ is a plausible fit of the degree distribution p_k for the
 320 protein interactions and autonomous systems networks in the right column of Figure 9. On the other hand,
 321 the degree distribution p_k of the spanning trees follows a power-law in all cases but the online friendships
 322 network. The maximum likelihood estimates of the power-law exponents γ are shown with solid lines
 323 in Figure 9.

324 REFERENCES

- 325 Adai, A. T., Date, S. V., Wieland, S., and Marcotte, E. M. (2004). LGL: Creating a map of protein
 326 function with an algorithm for visualizing very large biological networks. *Journal of Molecular Biology*,
 327 340(1):179–190.

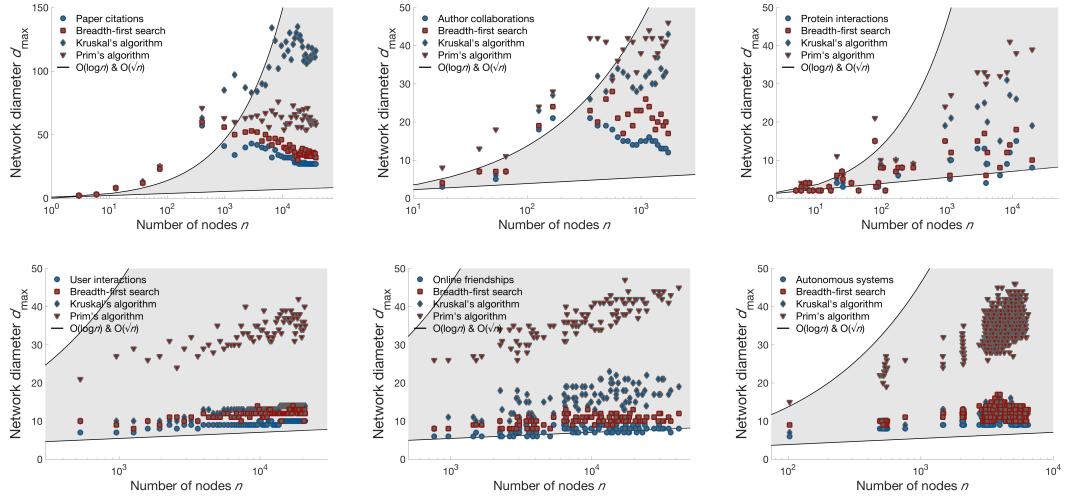


Figure 8. The diameter d_{\max} of real networks and their spanning trees computed with different algorithms (legend), while the shaded areas are the same as in Figure 4.

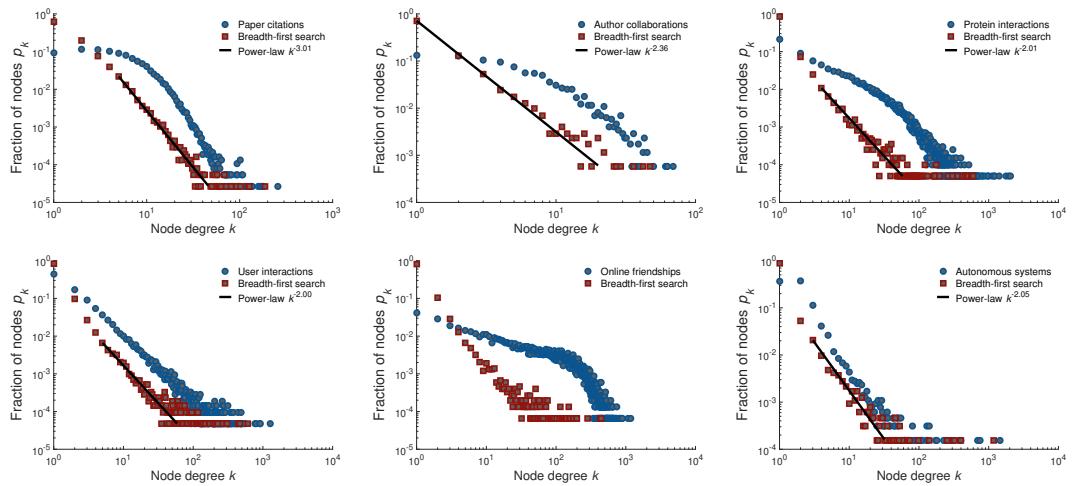


Figure 9. Node degree distribution p_k of real networks and their spanning trees computed with the breadth-first search algorithm. The power-law distributions are maximum likelihood estimates at p -value = 0.1 (Clauset et al., 2009).

- 328 American Physical Society (2015). APS Dataset. <http://journals.aps.org/datasets>.
- 329 Backstrom, L., Boldi, P., Rosa, M., Ugander, J., and Vigna, S. (2012). Four degrees of separation. In
Proceedings of the ACM International Conference on Web Science, pages 45–54, Evanston, IL, USA.
- 330 Barabási, A.-L. (2016). *Network Science*. Cambridge University Press, Cambridge.
- 331 Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks. *Science*, 286(5439):509–
512.
- 332 Biological General Repository for Interaction Datasets (2016). BioGRID Database. <http://thebiogrid.org>.
- 333 Blagus, N., Šubelj, L., and Bajec, M. (2017). Empirical comparison of network sampling: How to choose
the most appropriate method? *Physica A: Statistical Mechanics and its Applications*, 477:136–148.
- 334 Bollobás, B. (1998). *Modern Graph Theory*. Springer, Heidelberg.
- 335 Brandes, U. (2001). A faster algorithm for betweenness centrality. *Journal of Mathematical Sociology*,
25(2):163–177.
- 336 Broido, A. D. and Clauset, A. (2019). Scale-free networks are rare. *Nature Communications*, 10(1):1017.
- 337 Clauset, A., Shalizi, C. R., and Newman, M. E. J. (2009). Power-law distributions in empirical data. *SIAM
Review*, 51(4):661–703.
- 338 Cohen, R. and Havlin, S. (2003). Scale-free networks are ultrasmall. *Physical Review Letters*,
90(5):058701.
- 339 Coscia, M. and Neffke, F. (2017). Network backboning with noisy data. In *Proceedings of the IEEE
International Conference on Data Engineering*, pages 425–436, San Diego, CA, USA.
- 340 Eades, P. (1984). A heuristic for graph drawing. *Congressus Numerantium*, 42:149–160.
- 341 Erdős, P. and Rényi, A. (1959). On random graphs I. *Publicationes Mathematicae Debrecen*, 6:290–297.
- 342 Fortunato, S., Bergstrom, C. T., Börner, K., Evans, J. A., Helbing, D., Milojević, S., Petersen, A. M.,
Radicchi, F., Sinatra, R., Uzzi, B., Vespiagnani, A., Waltman, L., Wang, D., and Barabási, A.-L. (2018).
Science of science. *Science*, 359(6379):eaao0185.
- 343 Fortunato, S. and Hric, D. (2016). Community detection in networks: A user guide. *Physics Reports*,
659:1–44.
- 344 Freeman, L. (1977). A set of measures of centrality based on betweenness. *Sociometry*, 40(1):35–41.
- 345 Freeman, L. C. (1979). Centrality in social networks: Conceptual clarification. *Social Networks*, 1(3):215–
239.
- 346 Fruchterman, T. M. J. and Reingold, E. M. (1991). Graph drawing by force-directed placement. *Software:
Practice and Experience*, 21(11):1129–1164.
- 347 Gibson, H., Faith, J., and Vickers, P. (2013). A survey of two-dimensional graph layout techniques for
information visualisation. *Information Visualization*, 12(3-4):324–357.
- 348 Grady, D., Thiemann, C., and Brockmann, D. (2012). Robust classification of salient links in complex
networks. *Nature Communications*, 3:864.
- 349 Hamann, M., Lindner, G., Meyerhenke, H., Staudt, C. L., and Wagner, D. (2016). Structure-preserving
sparsification methods for social networks. *Social Network Analysis and Mining*, 6(1):22.
- 350 Institute of Information Science (2010). SICRIS Database. <http://www.sicris.si>.
- 351 Kamada, T. and Kawai, S. (1989). An algorithm for drawing general undirected graphs. *Information
Processing Letters*, 31(1):7–15.
- 352 Knuth, D. E. (2011). *The Art of Computer Programming*. Addison-Wesley Professional, Amsterdam.
- 353 Leskovec, J. and Faloutsos, C. (2006). Sampling from large graphs. In *Proceedings of the ACM SIGKDD
International Conference on Knowledge Discovery and Data Mining*, pages 631–636, Philadelphia, PA,
USA.
- 354 Leskovec, J., Kleinberg, J., and Faloutsos, C. (2005). Graphs over time: Densification laws, shrinking
diameters and possible explanations. In *Proceedings of the ACM SIGKDD International Conference on
Knowledge Discovery and Data Mining*, pages 177–187, Chicago, IL, USA.
- 355 Leskovec, J. and Krevl, A. (2014). SNAP Datasets. <http://snap.stanford.edu/data>.
- 356 Ma, K.-L. and Muelter, C. W. (2013). Large-scale graph visualization and analytics. *Computer*, 46(7):39–
46.
- 357 Meir, A. and Moon, J. W. (1970). The distance between points in random trees. *Journal of Combinatorial
Theory*, 8(1):99–103.
- 358 Newman, M. E. J. (2018). *Networks*. Oxford University Press, Oxford, 2nd edition.
- 359 Paranjape, A., Benson, A. R., and Leskovec, J. (2017). Motifs in temporal networks. In *Proceedings of*

- the ACM International Conference on Web Search and Data Mining, pages 601–610, Cambridge, UK.

Rényi, A. and Szekeres, G. (1967). On the height of trees. *Journal of the Australian Mathematical Society*, 7(4):497–507.

Rossi, R. A. and Ahmed, N. K. (2015). Network data repository. <http://networkrepository.com>.

Schoch, D. and Brandes, U. (2016). Re-conceptualizing centrality in social networks. *European Journal of Applied Mathematics*, 27(6):971–985.

Stark, C., Breitkreutz, B.-J., Reguly, T., Boucher, L., Breitkreutz, A., and Tyers, M. (2006). BioGRID: A general repository for interaction datasets. *Nucleic Acids Research*, 34(1):535–539.

Šubelj, L. (2018). Convex skeletons of complex networks. *Journal of the Royal Society Interface*, 15(145):20180422.

Tizzoni, M., Bajardi, P., Poletto, C., Ramasco, J. J., Balcan, D., Gonçalves, B., Perra, N., Colizza, V., and Vespignani, A. (2012). Real-time numerical forecast of global epidemic spreading: case study of 2009 A/H1N1pdm. *BMC Medicine*, 10(1):165.

Traud, A. L., Mucha, P. J., and Porter, M. A. (2012). Social structure of Facebook networks. *Physica A: Statistical Mechanics and its Applications*, 391(16):4165–4180.

Watts, D. J. and Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks. *Nature*, 393(6684):440–442.

Zitnik, M., Sosić, R., Feldman, M. W., and Leskovec, J. (2019). Evolution of resilience in protein interactomes across the tree of life. *Proceedings of the National Academy of Sciences of the United States of America*, 116(10):4426–4433.