

# A Family of Centrality Measures for Graph Data Based on Subgraphs

**Cristian Riveros**

Pontificia Universidad Católica de Chile, Santiago, Chile  
Millennium Institute for Foundational Research on Data, Santiago, Chile  
cristian.riveros@uc.cl

**Jorge Salas**

Pontificia Universidad Católica de Chile, Santiago, Chile  
Millennium Institute for Foundational Research on Data, Santiago, Chile  
jusalas@uc.cl

---

## Abstract

We present the theoretical foundations of a new approach in centrality measures for graph data. The main principle of our approach is very simple: the more relevant subgraphs around a vertex, the more central it is in the network. We formalize the notion of “relevant subgraphs” by choosing a family of subgraphs that, given a graph  $G$  and a vertex  $v$  in  $G$ , it assigns a subset of connected subgraphs of  $G$  that contains  $v$ . Any of such families defines a measure of centrality by counting the number of subgraphs assigned to the vertex, i.e., a vertex will be more important for the network if it belongs to more subgraphs in the family. We show many examples of this approach and, in particular, we propose the all-subgraphs centrality, a centrality measure that takes every subgraph into account. We study fundamental properties over families of subgraphs that guarantee desirable properties over the corresponding centrality measure. Interestingly, all-subgraphs centrality satisfies all these properties, showing its robustness as a notion for centrality. Finally, we study the computational complexity of counting certain families of subgraphs and show a polynomial time algorithm to compute the all-subgraphs centrality for graphs with bounded tree width.

**2012 ACM Subject Classification** Mathematics of computing → Graph theory; Information systems → Graph-based database models

**Keywords and phrases** Graph data, graph centrality, centrality measures

**Digital Object Identifier** 10.4230/LIPIcs.ICDT.2020.23

**Funding** C. Riveros and J. Salas were partially funded by the Millennium Institute for Foundational Research on Data.

## 1 Introduction

Which are the most important or “central” nodes in a network? This is a crucial question that has been asked in several areas like social science [21], biology [19], computer science [9] and essentially every area where graph data is relevant [25]. Given the graph structure of data one expects that more central nodes are more important for the network and they will be relevant in understanding its underlying structure. Several centrality measures have been proposed like closeness [4], betweenness [15], Page Rank [9], Katz index [20], among others [25], trying to give an answer or explanation to our first question.

Which centrality measure is the most meaningful for network analysis? This has been behind all proposals of centrality measures and it is an old question that has been discussed from the beginning of network analysis [7, 16, 27]. Over the years, some axioms or properties have been risen as crucial for a centrality measure and several centrality measures have been axiomatized [28, 29, 31]. However, as it was shown in [6] many commonly used centrality



© Cristian Riveros and Jorge Salas;

licensed under Creative Commons License CC-BY

23rd International Conference on Database Theory (ICDT 2020).

Editors: Carsten Lutz and Jean Christoph Jung; Article No. 23; pp. 23:1–23:18

Leibniz International Proceedings in Informatics



LIPIC Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

measures do not satisfy even a simple set of “desirable” axioms (i.e. properties). The question above then remains unanswered: how to naturally and formally define a centrality measure that has reasonable properties?

To motivate our approach that aims both questions, consider the following setting from graph data management. Suppose a graph database  $G$  and a query language  $\mathcal{L}$  for extracting patterns from  $G$ . Further, suppose  $Q$  is a query in  $\mathcal{L}$  such that the evaluation of  $Q$  over  $G$ , denoted by  $Q(G)$ , retrieves a set of nodes in  $G$ . How should we rank  $Q(G)$  in order to output the most meaningful outputs first? More specifically, suppose that  $G$  is a property graph and  $\mathcal{L}$  is a language of basic graph patterns [1]. Given that queries dynamically change over time [10, 22], one would expect that, if  $v \in Q(G)$  satisfies more patterns from  $\mathcal{L}$ , it will have more chances to appear later as an extension of  $Q$ . More general, one would expect that the more queries from  $\mathcal{L}$  where  $v$  is included, the more important is  $v$  on  $G$  with respect to  $\mathcal{L}$ . Furthermore, depending whether  $\mathcal{L}$  is designed to look for paths, trees, or maybe triangles [1] on  $G$ , maybe the user would like its measure of centrality to focus more on these patterns than in all basic graph patterns.

In this paper, we tackle the first question following the simple idea motivated from the graph data management setting: the more relevant subpatterns around a node, the more central it is in the network. Several proposals in the literature (e.g degree, betweenness [15], cross-clique [14]) already have considered relevant subpatterns like edges, paths, or cliques to define meaningful centrality measures. We generalize this approach by defining centrality measures based on families of subgraphs. Specifically, we formalize the notion of “relevant subgraphs” by choosing any family of subgraphs that, given a graph  $G$  and a vertex  $v$  in  $G$ , it assigns a subset of connected subgraphs from  $G$  that contains  $v$ . Any of such families defines a measure of centrality by counting the number of subgraphs assigned to the vertex, i.e., a vertex will be more important for the network if it belongs to more subgraphs in the family. We show several examples that can be derived by following this approach. In particular, a natural family of subgraphs is to consider all connected subgraphs around a vertex, that we called *all-subgraphs centrality*, and we show that it defines a well-behaved notion of centrality.

With a family of centrality measures at hand we embark on answering the second question. Generally speaking, we can consider any property on the family of subgraphs and see what “axiom” it implies in the respectively centrality notion that it defines. With this strategy, we no longer depend on comparing centrality measures of different nature (e.g. Page Rank vs Betweenness). Instead, we can understand all centrality notions proposed in this paper by just understanding the properties that satisfy the families of subgraphs. We consider simple axioms that has been proposed in the literature (e.g. monotonicity [27] or isolated vertex [16]). Then, look for simple properties in the family of subgraph that imply them. Interestingly, we can show natural examples of families of subgraphs that do not satisfy these properties and whose corresponding centrality notions do not satisfy the axioms. This allows to have a more deep understanding of why a centrality measure does not behave as expected and, moreover, to look into ways on how to “fix” it. Finally, the all-subgraphs centrality proposed in this paper satisfies all these properties and axioms, showing its robustness as a measure for centrality.

The general definition of centrality based on subgraphs allows us to easily extend the idea from vertices to sets of vertices, also called group centrality. We propose an approach to extend every centrality measure to groups, and prove a natural way to reduce the computation of all-subgraph centrality from groups to vertices. We show that this extension over sets allows to answer simple questions on the dynamic of graphs, like how to maximize the centrality of a vertex when an edge is added.

Towards the end of the paper, we study the computational complexity of counting certain families of subgraphs. Unfortunately, we show that the centrality measures defined from families of subgraphs like all subgraphs or trees lead to intractability. In terms of good news, we show that these centralities can be efficiently computed in acyclic graphs (i.e. trees). Moreover, we show that this result can be extended to more classes of graphs, by showing a polynomial time algorithm to compute the all-subgraphs centrality over all classes of graphs with bounded tree width.

**Related work.** Centrality measures have been extensively studied since the 50's [4, 21] and the subject is spread in different research areas. Moreover, the literature contains several alternative proposals that, given space restrictions, it will be impossible to cover all of them here (see [25]). Instead, we review here the work that is more closed in spirit to our proposal by stressing the main differences.

Centrality measures based on some relevant subgraphs have been studied before (e.g. betweenness [15], cross-clique [14]). The difference with our approach is that we take a step further and studied families of subgraphs in a more general setting. In particular, to the best of our knowledge all-subgraphs centrality and trees centrality (see Section 3 and 4) are new measures and have not been studied before.

There are several papers that have studied centrality measures in terms of properties [7, 16, 27]. Furthermore, in the last years there are several proposals to axiomatize standard centrality measures [5, 6, 28, 29, 31]. In this paper, we study properties and axioms in terms of families of subgraphs, which is a different goal compared to previous approaches.

Finally, a centrality measure called subgraph centrality was proposed in [13]. Although the name resemble our approach and the paper also motivates the use of subgraphs, subgraph centrality sums the number of closed-walks weighted by its length and not all the connected subgraphs that contains a nodes, as in our case.

## 2 Preliminaries

For a finite set  $V$ , we denote by  $\text{edges}(V) = \{\{u, v\} \subseteq V \mid u \neq v\}$  all subsets of  $V$  of size two. Sometimes, we consider a function  $f$  as a relation and write  $f' \subseteq f$  when  $f'$  is a (partial) function resulting to take a subset of the order pairs from  $f$ . In the sequel, all logarithms are in base 2 unless it is stated differently.

**Undirected graphs.** We consider finite undirected graphs of the form  $G = (V, E)$  where  $V$  is a finite non-empty set and  $E \subseteq \text{edges}(V)$ . Given a graph  $G$ , we will denote by  $V(G)$  and  $E(G)$  the set of vertices and edges, respectively. We will usually use  $u$  and  $v$  for denoting vertices and  $e$  and  $f$  for edges. Furthermore, we will use edges as sets and write  $v \in e$  when  $e$  is an edge incident to  $v$ . We denote by  $N(v, G) = \{u \mid \{u, v\} \in E(G)\}$  the neighborhood of  $v$  in  $G$ . We say that a graph  $G'$  is a subgraph of  $G$ , denoted  $G' \subseteq G$ , if  $V(G') \subseteq V(G)$  and  $E(G') \subseteq E(G)$ . If two graphs  $G_1$  and  $G_2$  are isomorphic, we write  $G_1 \cong G_2$ . Furthermore, we write  $G_1, v_1 \cong G_2, v_2$  for  $v_1 \in V(G_1)$  and  $v_2 \in V(G_2)$  if  $G_1 \cong G_2$  and  $v_1$  is equivalent to  $v_2$  under the bijective function between  $G_1$  and  $G_2$ .

**Multigraphs.** We also work with graphs with multiple edges between vertices, called multigraphs. A multigraph  $M$  is a triple  $M = (V, E, r)$  such that  $V$  is a finite non-empty set,  $E$  is a finite set, and  $r : E \rightarrow \text{edges}(V)$  (i.e. the edge-assignment function). Intuitively,  $E$  is a set of identifiers for edges and  $r$  assigns identifiers to edges (i.e. there could be multiple edges

between two pair of vertices). Similar than for graphs, we denote by  $V(M)$ ,  $E(M)$ , and  $r(M)$  the corresponding set of vertices, edges, and edge-assignment of  $M$ , respectively. We say that a multigraph  $M'$  is a sub-multigraph of  $M$ , denoted by  $M' \subseteq M$ , if  $V(M') \subseteq V(M)$ ,  $E(M') \subseteq E(M)$ , and  $r(M') \subseteq r(M)$ .

Note that a simple graph is a multigraph  $M$  where  $r(M)$  is an injective function. For this reason, in the future we will not make distinction between graphs and multigraphs. Furthermore, we will usually work with graphs but all definitions and results also extend to multigraphs. When this is not the case, we will make the distinction explicitly.

**Connected graphs.** A path in a graph  $G$  is a sequence of nodes  $\pi = v_0, \dots, v_n$  such that  $\{v_i, v_{i+1}\} \in E(G)$  for every  $i < n$  and we say that the length of  $\pi$  is  $n$ . Note that  $v_0$  is the trivial path from  $v_0$  to itself of length 0. We say that  $G$  is connected if there exists a path between any pair of vertices. Furthermore, we say that  $G' \subseteq G$  is a connected component of  $G$  if  $G'$  is connected and its maximal element over all subgraphs of  $G$  under  $\subseteq$ . We denote by  $\text{ConnComp}(G)$  the set of all connected components of  $G$ . For  $u, v \in V(G)$  we say that  $u$  is at distance  $d$  of  $v$  if there exists a path from  $u$  to  $v$  of length  $d$  and every path from  $u$  to  $v$  is of length at least  $d$ . We denote the distance  $d$  from  $u$  and  $v$  in  $G$  by  $\text{dist}_G(u, v)$ . Given this distance, the diameter of  $G$  is defined as  $\max_{u, v \in V(G)} \text{dist}_G(u, v)$ .

**Families.** We consider several families of graphs through the paper to give examples or show some properties of our centrality measures. Given a vertex  $v$ , we denote by  $G_v$  the graph with one vertex  $v$  (i.e.  $V(G_v) = \{v\}$ ) and no-edges (i.e.  $E(G_v) = \emptyset$ ). Given an edge  $e = \{u, v\}$ , we denote by  $G_e$  the graph only containing  $e$  (i.e.  $V(G_e) = \{u, v\}$  and  $E(G_e) = \{e\}$ ). For any  $n \geq 1$ , we write  $S_n$  for the star with  $n + 1$  vertices such that  $V(S_n) = \{0, 1, \dots, n\}$  and all nodes are connected to 0, namely,  $E(S_n) = \{\{0, i\} \mid 0 < i \leq n\}$ . Similarly, we write  $L_n$  for the line with  $n$  vertices where  $V(L_n) = \{0, \dots, n - 1\}$  and  $E(L_n) = \{\{i, i + 1\} \mid 0 \leq i < n - 1\}$ . The circuit with  $n$  vertices is denoted by  $C_n$  with  $V(C_n) = \{0, \dots, n - 1\}$  and  $E(C_n) = \{\{i, (i + 1) \bmod n\} \mid 0 \leq i \leq n - 1\}$ . Finally, the clique of size  $n$  is denoted by  $K_n$  where  $V(K_n) = \{0, \dots, n - 1\}$  and  $E(K_n) = \text{edges}(V(K_n))$ .

**Operations.** Through the paper, we use several operations to create, modify, or combine graphs. Given  $v \in V(G)$ , we denote by  $G - v$  the result of removing  $v$  from  $G$  and all its incident edges, namely,  $V(G - v) = V(G) \setminus \{v\}$  and  $E(G - v) = \{e \in E(G) \mid v \notin e\}$ . Given  $e = \{u, v\}$ , we write  $G + e$  for the result of adding  $e$  into  $G$ , formally,  $V(G + e) = V(G) \cup e$  and  $E(G + e) = E(G) \cup \{e\}$  (i.e. if  $u$  or  $v$  are not in  $G$ , then they are included as new vertices). Instead, we write  $G - e$  for the result of removing all edges between  $u$  and  $v$ , namely,  $V(G - e) = V(G)$  and  $E(G - e) = E(G) \setminus \{e\}$ . Note that if  $G$  is a (simple) graph, then at most one edge is removed, but if  $G$  is a multigraph then all edges between  $u$  and  $v$  are removed. For  $G_1$  and  $G_2$ , we denote by  $G_1 \cup G_2$  the union of the two graphs, namely,  $V(G_1 \cup G_2) = V(G_1) \cup V(G_2)$  and  $E(G_1 \cup G_2) = E(G_1) \cup E(G_2)$ . In particular,  $G + e = G \cup G_e$ .

Let  $G$  be a graph and  $U \subseteq V(G)$ . We define the *set contraction* of  $U$  on  $G$  as the multigraph  $G/U$  by merging the vertices  $U$  to one vertex (called  $U$ ) and keeping multi-edges into  $U$ . Formally,  $G/U$  is the multigraph  $M$  such that  $V(M) = (V(G) \setminus U) \cup \{U\}$ ,  $E(M) = \{e \in E(G) \mid e \not\subseteq U\}$  and for every  $e \in E(M)$  either  $r(M)(e) = e$  whenever  $e \cap U = \emptyset$ , or  $r(M)(e) = \{v, U\}$  whenever  $e = \{v, u\}$  with  $v \notin U$  and  $u \in U$ . Note that we use  $U$  (i.e. the set) as the new vertex that represent the contraction in  $M/U$ . When  $G$  is a multigraph, the set contraction  $G/U$  easily follows from the above definition.

### 3 The all-subgraphs centrality

We start by introducing our first centrality measure based on all subgraphs, called the all-subgraphs centrality. In the next section, we generalize this idea to any family of subgraphs.

Fix a graph  $G$  and a vertex  $v \in V(G)$ . We denote by  $\mathcal{A}(v, G)$  the set of all connected subgraphs of  $G$  that contains  $v$ , formally,  $\mathcal{A}(v) = \{G' \subseteq G \mid G' \text{ is connected} \wedge v \in V(G')\}$ . The *all-subgraphs centrality* of  $v$  in  $G$  is defined as:

$$C_{\mathcal{A}}(v, G) := \log(|\mathcal{A}(v, G)|)$$

namely, the logarithm of the number of connected subgraphs of  $G$  that contains  $v$ . Intuitively, the all-subgraphs centrality of a node only considers connected graphs since it captures the importance of the node in the neighborhood that it belongs. We add more importance to a node if its neighborhood is richer in substructures. Furthermore, we consider connected subgraphs since there is no argument to say that a node has more centrality by counting another component that is not directly connected to it.

The function  $C_{\mathcal{A}}$  naturally induces a ranking between nodes: the higher the centrality  $C_{\mathcal{A}}(v, G)$ , the more important is  $v$  in  $G$ . We define the ranking  $<_{\mathcal{A}}$  over  $V(G)$  induced by  $C_{\mathcal{A}}$  (or just  $\mathcal{A}$ -ranking for short) such that  $u <_{\mathcal{A}} v$  if, and only if,  $C_{\mathcal{A}}(u, G) > C_{\mathcal{A}}(v, G)$ . Strictly speaking,  $<_{\mathcal{A}}$  is not an order in  $V(G)$ , given that there could exist vertices  $u$  and  $v$  such that  $C_{\mathcal{A}}(u, G) = C_{\mathcal{A}}(v, G)$  (e.g.  $u$  and  $v$  are isomorphic in  $G$ ). In this case, we write  $u =_{\mathcal{A}} v$ .

► **Example 1.** Let  $v$  be a vertex. Recall that  $G_v$  is the trivial graph with one vertex  $v$  and no edges. Then one can easily check that  $C_{\mathcal{A}}(v, G_v) = 0$  given that  $\mathcal{A}(v, G_v) = \{G_v\}$  and then  $\log(|\mathcal{A}(v, G_v)|) = \log(1) = 0$ . Note that this is the only vertex and graph (up to isomorphism) where the centrality is equal to 0. This follows the intuition that an isolated vertex must have 0 centrality since no one is connected to him.

► **Example 2.** Recall that  $S_n$  denotes the star graph with  $n + 1$  vertices. Note that every connected subgraph of  $S_n$  corresponds to a subset of  $E(S_n)$ , and there are  $2^n$  subsets of  $E(S_n)$ . Therefore, the centrality of the center of the star (i.e. the 0 vertex) is  $C_{\mathcal{A}}(0, S_n) = n$ . Interestingly, the all-subgraphs centrality of the center of a star coincides with its degree-centrality [25], following the intuition of what should be the centrality in this case. One can easily show that, for any  $i \neq 0$ ,  $C_{\mathcal{A}}(i, S_n) = n - 1 + \epsilon$  with  $\epsilon \in o(1)$ . Thus, in terms of ranking we have that  $0 <_{\mathcal{A}} i$  and  $i =_{\mathcal{A}} j$  for every  $i, j > 0$ .

The all-subgraphs centrality is measuring the worst-case entropy [12, 24] of the set  $\mathcal{A}(v, G)$ , namely, the minimum number of bits that are required to represent the set  $\mathcal{A}(v, G)$  with bit-codes. Of course, using the size of  $|\mathcal{A}(v, G)|$  will give the same ranking of centrality over the vertex of  $G$ . Nevertheless, the log-function gives a better interpretation of the centrality in terms of information theory. Moreover, it normalizes the value  $|\mathcal{A}(v, G)|$  in a scale that is in correspondence with the intuition of a centrality notion, e.g. Examples 1 and 2 above.

The next lemma is another result that validates the use of worst-case entropy and it will be useful for computing the all-subgraphs centrality over simple graphs. Recall that a vertex  $v \in V(G)$  is a *cut vertex* of  $G$  if  $|\text{ConnComp}(G - v)| < |\text{ConnComp}(G)|$ , namely, whose removal increases the number of connected components of  $G$ .

► **Lemma 3.** *Let  $v$  be a cut-vertex of graph  $G$  and  $G_1, \dots, G_n$  are all the subgraphs that partition  $G$  and whose pairwise intersection is  $v$ , that is,  $V(G) = \cup_{i=1}^n V(G_i)$ ,  $E(G) = \cup_{i=1}^n E(G_i)$ , and  $V(G_i) \cap V(G_j) = \{v\}$  for  $i \neq j$ . Then*

$$C_{\mathcal{A}}(v, G) = \sum_{i=1}^n C_{\mathcal{A}}(v, G_i).$$

*Namely, the centrality of  $v$  in  $G$  is the sum of its centrality in all the components  $G_i$ .*

This property is usually known in the literature as cut-vertex additivity [29]. Since not every centrality measure satisfies it, this can be seen as the first distinction between all-subgraphs centrality and commonly used centrality measures (e.g. pagerank, betweenness).

► **Example 4.** Let  $G$  be any graph,  $u \in V(G)$ , and  $v$  be a new vertex not in  $G$ . For  $e = \{u, v\}$ , recall that  $G_e$  is the graph only containing  $e$ . Then one can easily see that  $C_{\mathcal{A}}(u, G_e) = 1$ . Since  $G + e = G \cup G_e$  and  $u$  is a cut-vertex of  $G + e$ , by Lemma 3 we get:

$$C_{\mathcal{A}}(u, G + e) = C_{\mathcal{A}}(u, G) + C_{\mathcal{A}}(u, G_e) = C_{\mathcal{A}}(u, G) + 1$$

Thus, by connecting one new vertex directly to  $u$  its centrality grows exactly in one unit. This property is very appealing for a centrality measure and follows verbatim the intuition of the score-monotonicity axiom in [6] (see Section 5 for more discussion). On the other hand, one can check that the new vertex  $v$  in  $G + e$  absorbs part of the centrality of  $u$  in  $G$ . Specifically, one can easily see that  $|\mathcal{A}(v, G + e)| = |\mathcal{A}(u, G)| + 1$  and then  $C_{\mathcal{A}}(v, G + e) = \log(|\mathcal{A}(u, G)| + 1) = C_{\mathcal{A}}(u, G) + \epsilon$ , where  $\epsilon$  is a negligible factor.

► **Example 5.** For  $n \geq 1$ , recall that  $L_n$  is the line with  $n$  nodes starting from 0 and ending in  $n - 1$ . For the 0-vertex in  $L_n$  there are  $n$ -different subgraphs, one for each vertex, and then  $C_{\mathcal{A}}(0, L_n) = \log(n)$ . The line graph is the most sparse graph with  $n$  vertices and 0 is the most extreme vertex in the graph. As one could expect, the centrality of 0 grows very slow, logarithmic in the number of vertex.

For the  $i$ -vertex in  $L_n$ , we can easily compute its centrality by using Lemma 3. Indeed, the centrality for  $i$  is the composition of two lines with  $i + 1$  and  $n - i$  vertices, respectively. Therefore, by Lemma 3:

$$C_{\mathcal{A}}(i, L_n) = C_{\mathcal{A}}(0, L_{i+1}) + C_{\mathcal{A}}(0, L_{n-i}) = \log(i + 1) + \log(n - i).$$

If  $n$  is odd, the vertex with maximum centrality is reached by the middle node  $\frac{n-1}{2}$  and  $C_{\mathcal{A}}(\frac{n-1}{2}, L_n) = 2(\log(n + 1) - 1)$ . Thus, the middle point of a line doubles the centrality of the extreme vertices, nevertheless, the grow of its centrality is still logarithmic in  $n$ . Finally, note that the centrality is maximized in the middle node and the ranking decreases towards the extremes (i.e.  $i <_{\mathcal{A}} i + 1$  for every  $i < \frac{n-1}{2}$ ).

A natural question at this point is to think in lower and upper bounds of the centrality with respect to the number of edges of a graph. Indeed, the number of subgraphs  $\mathcal{A}(v, G)$  could be exponential in  $G$  but its entropy is bounded by the number of edges as follows.

► **Proposition 6.** *For any connected graph  $G$  and  $v \in V$ , it holds that:*

$$\log(|E(G)| + 1) \leq C_{\mathcal{A}}(v, G) \leq |E(G)|.$$

From Example 2 above, we can infer that the upper bound is reached by the central vertex of a star. This follows the intuition that the central vertex of a star must be the most central vertex regarding the number of edges (i.e. all edges are pointing to him). Furthermore, in Example 5 we show that the extreme vertex of a line  $L_n$  has centrality  $\log(n) = \log(|E| + 1)$ . That is, the minimum centrality is reached in the extreme points of a line, agreeing with the intuition that the line graph is the most sparsest graph over all undirected graphs.

#### 4 A family of centralities based on subgraphs

The idea of measuring the centrality of a vertex based on relevant substructures is not new [14, 15]. For example, the degree centrality counts how many edges are incident to a vertex and the betweenness centrality [15] counts how many geodesic paths passed through a vertex. In our case, all-subgraphs centrality measures all connected subgraphs including  $v$ , but maybe for an expert not all subgraphs are equally important and he will be interested in counting some of them. In this section we generalize the notion of all-subgraphs centrality to propose a framework of centrality notions based on measuring the worst-case entropy of relevant substructures surrounding a vertex.

A family of substructures is a function  $\mathcal{F}$  that, given a graph  $G$  and a vertex  $v \in V(G)$ , it assigns a non-empty subset of connected subgraphs in  $G$  that contains  $v$ . Formally,  $\mathcal{F}$  is a function such that  $\mathcal{F}(v, G) \subseteq \mathcal{A}(v, G)$  and  $\mathcal{F}(v, G) \neq \emptyset$ . We also assume that  $\mathcal{F}$  is closed under isomorphism, namely, if  $G_1, v_1 \cong G_2, v_2$  then  $\mathcal{F}(v_1, G_1)$  is isomorphic to  $\mathcal{F}(v_2, G_2)$ , by extending the isomorphism between  $G_1, v_1$  and  $G_2, v_2$  to subgraphs. For example,  $\mathcal{A}$  is a family of substructures where  $\mathcal{A}(v, G)$  contains all connected subgraphs in  $G$  containing  $v$  and is closed under isomorphism. Given a family of substructures we define the  $\mathcal{F}$ -subgraph centrality (denoted by  $C_{\mathcal{F}}(v, G)$ ) as:

$$C_{\mathcal{F}}(v, G) := \log(|\mathcal{F}(v, G)|)$$

for any graph  $G$  and vertex  $v \in V(G)$ . In other words, following the idea of all-subgraphs centrality it measures the worst-case entropy of the substructures  $\mathcal{F}(v, G)$ . We could have left the framework open to any monotone positive function over  $\mathcal{F}(v, G)$  instead of the logarithm, leading to the same ranking of centrality between vertices. Of course, this will derive in a more complex and enriched theory, however, for the purpose of this paper we will keep the simplicity of the logarithm as it still give place to novel results.

Note that  $\mathcal{F}(v, G)$  is non-empty and, therefore,  $C_{\mathcal{F}}(v, G)$  is always well-defined. Similar to all-subgraphs centrality, the centrality measures induced a ranking between nodes: we define the  $\mathcal{F}$ -ranking  $<_{\mathcal{F}}$  over  $V(G)$  such that  $u <_{\mathcal{F}} v$  if, and only if,  $C_{\mathcal{F}}(u, G) < C_{\mathcal{F}}(v, G)$ .

► **Example 7.** Given a graph  $G$  and  $v \in V(G)$ , denote by  $\mathcal{T}(v, G)$  all subgraphs  $T \in \mathcal{A}(v, G)$  such that  $T$  is a tree. Note that an isolated vertex is defined as a trivial tree, so  $\mathcal{T}(v, G)$  is always non-empty. Furthermore, the family  $\mathcal{T}$  is closed under isomorphism. Then  $C_{\mathcal{T}}$  measures the centrality of a vertex based on trees and we call it *the trees centrality*. For example, if  $L_n$  is a line graph with  $n$  vertices (see Example 5) then we have that  $C_{\mathcal{A}}(v, G) = C_{\mathcal{T}}(v, G)$ . Indeed, if  $T$  is a tree, then  $C_{\mathcal{A}}(v, G) = C_{\mathcal{T}}(v, G)$  for every  $v \in V(G)$ . However, this is not always the case if  $G$  has cycles and one can find examples where the two measures give different values and ranking.

The motivation behind trees centrality is to considered substructures defined by acyclic graphs like trees or paths. For example, path queries [1] are at the core of graph queries languages and they are used to find path substructures between pair of nodes. Also, basic graph patterns that are acyclic (e.g. tree-shaped queries) forms a well-behaved core of graph query languages that can be evaluated efficiently [18]. Therefore, if the query languages mostly uses queries that are acyclic, maybe it makes sense to rank the results by a centrality notion based on trees.

The generalization of all-subgraphs centrality to any family of subgraphs opens the possibilities of defining any centrality notion based on a particular group of relevant subgraphs. In the next section, we use this framework to understand which properties in the family leads to desirable properties in the corresponding centrality measure. This will help to guide the design of a centrality notion based on subgraphs and, moreover, to have a better understanding of this framework and all-subgraphs centrality.

## 5 What families of subgraphs define good centrality measures?

Several attempts have been taken to define which properties a centrality measure should satisfy and how to axiomatize them [3, 28, 29]. In our framework, each family of subgraphs defines a new centrality measure, so it is not our purpose here to axiomatize them. In some sense, each family of subgraphs captures the know-how of an expert who knows what are the relevant subpatterns around a vertex. From this point of view, it does not make sense to prefer one notion of centrality over the other. Instead, we study here which properties over the family of subgraphs lead to desirable properties on the corresponding centrality notion. We hope that these properties will guide experts on the design of a centrality based on subgraphs and they will help to understand the benefits and problems of choosing one family over the other. Towards this goal, we consider several axioms of centrality that has been proposed in the literature and study which natural property on the family of subgraphs is enough to satisfy it. We also give several examples for showing what happens when a property is not satisfied.

In the sequel, a centrality measure is any function  $C$  that given a graph  $G$  and  $v \in V(G)$ , it outputs a non-negative value, i.e.,  $C(v, G) \geq 0$ .

**Default axioms.** We start our discussion by showing three natural axioms proposed in the literature that any  $C_{\mathcal{F}}$  satisfies, for any family of substructures  $\mathcal{F}$ . We discuss these three axioms briefly and show that they are naturally satisfied by definition.

In [28] they present the so-called *locality axiom*, which says that the centrality of a vertex should only depend on the connected component it belongs. In other words, after removing components that are not connected to a vertex  $v$  the centrality of  $v$  should not change. This natural axiom is satisfied by any centrality measure based on subgraphs because we define a family of substructures as a subset of connected subgraphs. This might be seen as an irrelevant detail but it is an important design decision of our approach. In second place, an axiom called *anonymity* is introduced in [27]. This is the same as saying that a centrality measure is closed under isomorphism. In our definition we explicitly say that any feasible substructure must be closed under isomorphism, which means that the centrality measure as defined will satisfy this axiom. Finally, in [17] the authors propose a minimum value for any centrality. More specifically, the centrality of an isolated vertex is the minimum possible and it should be 0. These two properties are called isolated minimization and isolated zero, respectively. In our case, these axioms are satisfied by definition, because the set of substructures associated to a vertex must be always non-empty, which means that  $|\mathcal{F}(v, G_v)| = 1$  for any family of substructures. Therefore, the minimum possible value for any centrality measure defined in this way is 0.

**Monotonicity.** The monotonicity axiom is probably the property that more people [6, 27, 28] agree that any centrality notion should satisfy. In [6], the definition of this axiom says that if an edge is added to the graph, then the centrality of the vertex that is incident with the new edge should not decrease. Clearly, a vertex is more central the more edges it has and, thus, a new edge should help to increase its relevance in the graph. A more general definition of this axiom was introduced in [27] where the effect of adding any new edge in the graph should not decrease the centrality of every vertex.

► **Axiom 1 (Monotonicity).** *A centrality measure  $C$  satisfies the monotonicity axiom if for every graph  $G$ ,  $v \in V(G)$  and  $e \notin E(G)$ , it holds that  $C(v, G) \leq C(v, G + e)$ .*



Note that the axiom implies that if  $G_1$  is a subgraph of  $G_2$  and  $v \in V(G_1)$ , then  $C(v, G_1) \leq C(v, G_2)$ . This coincides with the intuition that  $v$  in  $G_2$  has the same or more connections than in  $G_1$  and, thus, its relevance in  $G_2$  should be at least the one in  $G_1$ .

What property should a family of subgraphs  $\mathcal{F}$  satisfy in order that  $C_{\mathcal{F}}$  satisfy Axiom 1? Intuitively, when edge  $e$  is added to  $G$  we have  $G \subseteq G + e$  and all subgraphs that are relevant for  $v$  in  $G$  should also be relevant for  $v$  in  $G + e$ . Moreover, if a subgraph  $S$  is relevant for  $v$  in  $G + e$  but  $S$  is a subgraph of  $v$  in  $G$ , then it should also be a relevant subgraph of  $v$  in  $G$ . That is, all subgraphs of  $G$  that are relevant should also be relevant in  $G + e$  and vice versa. We call this the containment property.

► **Property 1 (Containment).** *A family of subgraphs  $\mathcal{F}$  satisfies the containment property if for every graphs  $G_1$  and  $G_2$  such that  $G_1 \subseteq G_2$  and for every  $v \in V(G_1)$  and  $S \in \mathcal{A}(v, G_1)$ , it holds that  $S \in \mathcal{F}(v, G_1)$  if, and only if,  $S \in \mathcal{F}(v, G_2)$ .*

In particular, the containment property implies that  $\mathcal{F}(v, G_1) \subseteq \mathcal{F}(v, G_2)$  whenever  $G_1 \subseteq G_2$ . As one could expect, the containment property is enough to satisfy the monotonicity axiom.

► **Theorem 8.** *If a family of subgraphs  $\mathcal{F}$  satisfies the containment property, then the corresponding centrality measure  $C_{\mathcal{F}}$  satisfies the monotonicity axiom.*

One can easily see that the family of all-subgraphs and trees satisfies the containment property and, therefore, the all-subgraphs centrality and trees centrality satisfy monotonicity as expected. Next, we show that this is not always the case.

► **Example 9.** Given a graph  $G$  and  $v \in V(G)$ , denote by  $\mathcal{W}(v, G)$  all subgraphs  $P \in \mathcal{A}(v, G)$  such that  $P = v_0, \dots, v_n$  is a geodesic path in  $G$ , namely, it is a path of minimal distance between  $v_0$  and  $v_n$ . We assume here that the isolated vertex  $v$  is the only geodesic path from  $v$  to  $v$ . In [8],  $|\mathcal{W}(v, G)|$  is defined as the stress centrality of vertex  $v$ . Then we define log-stress centrality of  $v$  in  $G$  as  $C_{\mathcal{W}}(v, G)$ . Of course,  $C_{\mathcal{W}}(v, G)$  is not equivalent to  $\text{Betweenness}(v, G)$  as a value and in how we aggregate the number of geodesic paths. Nevertheless, it will be useful below to understand  $\text{Betweenness}$  in the context of counting subgraphs.

One can easily show that the family  $\mathcal{W}$  does not satisfy the containment condition. Consider just a line  $L_3 = \text{---}$ . Then if we connect the black vertices and make a triangle  $K_3 = \triangle$ , then the geodesic path  $\text{---}$  is in  $\mathcal{W}(1, L_3)$  but  $\text{---}$  is not in  $\mathcal{W}(1, K_3)$ . Coincidentally, log-stress centrality (and betweenness centrality as well) do not satisfy the monotonicity axiom. Actually, one can show pathological examples where monotonicity does not hold [16]. For example, if one compares the circuit  $C_n$  with the clique  $K_n$  one can see that  $C_n \ll K_n$  but  $C_{\mathcal{W}}(0, C_n) > C_{\mathcal{W}}(0, K_n)$ , and  $\text{Betweenness}(0, C_n) > \text{Betweenness}(0, K_n)$  as well.

It is important to note that, for some axiomatic approaches [28], it is desirable that the center of a star  $S_{n-1}$  is the most central node in a graph with  $n$  vertices, namely, a centrality measure  $C$  satisfies this axiom if, for any  $n$  and for any graph  $G$  with  $|V(G)| = n$ , it holds that  $C(v, G) \leq C(0, S_{n-1})$  for every  $v \in V(G)$ . Unfortunately, this assumption contradicts the idea behind the monotonicity axiom, since we can add edges to  $S_{n-1}$  but the centrality of the center will never increase. Thus, given that this axiom contradicts the monotonicity axiom, we do not consider it in our analysis.

**Rank monotonicity.** Another axiom that has been remarked as important in the literature is rank monotonicity [5, 6, 11, 27]. Similar than for monotonicity, this axiom says that if  $v$  is more central than  $u$  in  $G$ , then when we add a new edge  $e$  to  $v$  the ranking between  $u$  and  $v$  is preserved. In particular, if  $v$  is the most central vertex in  $G$ , then it will be the most central vertex in  $G + e$  as well. We generalize this intuition as follows.

## 23:10 A Family of Centrality Measures for Graph Data Based on Subgraphs

► **Axiom 2** (Rank monotonicity). *A centrality measure  $C$  satisfies the rank monotonicity axiom if for every graph  $G$ ,  $u, v \in V(G)$  and  $e \notin E(G)$  with  $v \in e$ , then  $C(u, G) \leq C(v, G)$  implies that  $C(u, G + e) \leq C(v, G + e)$ .*

Note that with  $e = \{u, v\}$  it could happen that the increment in centrality for  $u$  is bigger than the increment on  $v$ , but the axiom says that the centrality of  $v$  will be still bigger than the centrality of  $u$ . In other words, if I meet Donald Trump, my centrality will rise more than his centrality, however, Donald Trump will still be the president of US.

It is important to say that in [6] an axiom called *density axiom* was proposed, which is a special case of rank monotonicity. Specifically, take a clique  $K_n$ , a circuit  $C_n$ , and vertices  $u \in V(K_n)$  and  $v \in V(C_n)$ . Then the density axiom says that if we connect  $u$  and  $v$  with an edge  $e = \{u, v\}$ , then  $G = (K_n \cup C_n) + e$  satisfies  $C(u, G) > C(v, G)$  for a centrality measure  $C$ . Intuitively, given that the neighborhood of  $u$  is more dense than in  $v$ , then its centrality should be bigger. One can see that if  $C$  satisfies monotonicity (i.e. vertices in  $K_n$  has more centrality than in  $C_n$ ), then rank monotonicity implies the density axiom [6]. Therefore, we can see rank monotonicity as a generalization of the density axiom in [6].

The containment property is useful to imply rank monotonicity but it is not enough. One can easily find centrality measures that satisfies Axiom 1 but it does not satisfy Axiom 2 (see Example 11 below). For this, one needs a notion of “fairness” in the family of subgraphs. Intuitively, if  $S$  is a relevant subgraph for  $v$  in  $G$  and  $S$  contains a vertex  $u$ , then  $S$  should also be relevant for  $u$  in  $S$ .

► **Property 2** (Fairness). *A family of subgraphs  $\mathcal{F}$  satisfies the fairness property if for every graph  $G$ ,  $u, v \in V(G)$  and  $S \subseteq G$  with  $u, v \in V(S)$  it holds that  $S \in \mathcal{F}(u, G)$  iff  $S \in \mathcal{F}(v, G)$ .*

As we show next, fairness is what you need if you want to preserve the ranking between vertices in a graph.

► **Theorem 10.** *If a family of subgraphs  $\mathcal{F}$  satisfies the containment property and fairness, then the corresponding centrality measure  $C_{\mathcal{F}}$  satisfies the rank monotonicity axiom.*

The family of all-subgraphs, trees and even betweenness (i.e. geodesic paths) satisfy fairness. Given that all-subgraphs and trees also satisfy the containment property, we conclude that both satisfy the rank monotonicity axiom. Next we show a natural family that satisfy the containment property but does not satisfy fairness.

► **Example 11.** A natural approach to define a family of subgraphs is to consider subpatterns on a neighborhood of bounded size around a vertex. Intuitively, an expert would not care if a vertex  $v$  can reach a far vertex  $u$  as long as there are many other substructures close to  $v$ . To formalize this, let  $k \geq 1$ . For a graph  $G$ , fix a vertex  $v$  and let  $N_k$  be the induced subgraph of all vertices at distance at most  $k$  of  $v$ , i.e.,  $V(N_k) = \{u \in V(G) \mid \text{dist}_G(u, v) \leq k\}$  and  $E(N_k) = \{e \in E(G) \mid e \subseteq V(N_k)\}$ . We define the family of subgraphs  $\mathcal{N}_k$  such that  $\mathcal{N}_k(v, G) = \mathcal{A}(v, N_k)$ , that is, all subgraphs in the neighborhood of  $v$  with radius  $k$ . Then we define the  $k$ -neighborhood centrality of  $v$  on  $G$  as  $C_{\mathcal{N}_k}(v, G)$ . Note that if the diameter of the graph is less than  $k$  then  $\mathcal{N}_k(v, G)$  and  $\mathcal{A}(v, G)$  coincide.

The family of  $k$ -neighborhood satisfies monotonicity but it does not satisfy fairness. Moreover, it does not satisfy the rank monotonicity axiom. To see this, consider the family  $\mathcal{N}_2$  and  $G_1 = \begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array}$ . By counting, one can check that the left white vertex, called  $u$ , and the right white vertex, called  $v$ , satisfy  $|\mathcal{N}_2(u, G_1)| = 8$  and  $|\mathcal{N}_2(v, G_1)| = 5$ , respectively. Then  $C_{\mathcal{N}_2}(u, G_1) > C_{\mathcal{N}_2}(v, G_1)$ . However, if we add an edge  $e$  between the two and create the graph  $G_1 + e = \begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \\ | \quad | \\ \bullet \quad \bullet \end{array}$ , then one can check that  $\mathcal{N}_2$  does not satisfy fairness. For instance,

the whole graph  $G_1 + e \in \mathcal{N}_2(v, G_1 + e)$ , contains  $u$  and  $v$ , but  $G_1 + e \notin \mathcal{N}_2(u, G_1 + e)$ . One can also check by counting that  $|\mathcal{N}_2(u, G_1 + e)| = 24$  and  $|\mathcal{N}_2(v, G_1 + e)| = 45$ . Thus,  $C_{\mathcal{N}_2}(u, G_1 + e) < C_{\mathcal{N}_2}(v, G_1 + e)$  and 2-neighborhood does not satisfy the rank monotonicity axiom as well.

The previous example shows that, if we want to approximate all-subgraphs centrality by only counting subgraphs up to a certain radius, one will have to lose some natural properties, like rank monotonicity.

**Line minimization.** Everyone would agree that any reasonable notion for centrality should assign 0 centrality to an isolated vertex [3, 28]. Basically, there is nothing less central to a community than the vertex that is not connected to any other vertex. One can generalize this idea by considering, what is the most sparse connected graph with  $n$  vertices. Clearly, the line  $L_n$  should be this graph: it is the only graph with  $n$  vertices that maximizes the diameter. Then the vertices that minimize the centrality in the line  $L_n$  are its extreme points, 0 and  $n - 1$ , and one would expect that this should be the vertices that have less centrality over all connected graphs with  $n$ -vertices.

► **Axiom 3** (Line minimization). *A centrality measure  $C$  satisfies the line minimization axiom if for every  $n$  and every connected graph  $G$  with  $|V(G)| = n$  it holds that  $C(0, L_n) \leq C(v, G)$  for every  $v \in V(G)$ .*

All centralities that we consider in this paper satisfy the line minimization axiom. Of course, one can manage to find unnatural families of subgraphs that produce centrality measures not satisfying this axiom. Still, one would like to find under which circumstances a centrality measure defined from a family of subgraphs satisfies it. For this, we need to introduce the following property.

► **Property 3** (Inclusion). *A family of subgraphs  $\mathcal{F}$  satisfies the inclusion property if for every graph  $G$ ,  $v \in V(G)$ , and  $S \in \mathcal{F}(v, G)$ , if  $S' \subseteq S$  and  $v \in V(S')$ , then  $S' \in \mathcal{F}(v, G)$ .*

Intuitively, this property is saying that every subgraph of a relevant subgraph should also be relevant for the family. Actually, this property is satisfied by all families of subgraphs proposed so far.

► **Theorem 12.** *If a family of subgraphs  $\mathcal{F}$  satisfies the containment and inclusion properties, then the corresponding centrality measure  $C_{\mathcal{F}}$  satisfies the line minimization axiom.*

**Continuity.** The inclusion property plus the containment property actually imply a natural property over centrality measures defined by family of subgraphs. Given that all subgraphs of a relevant subgraph are also included, it gives a sense of “continuity” in the centrality notion. Specifically, each time that we add a set of edges that rises the centrality of a vertex, there exists a way to add them, one at a time, in such a way that the centrality of the vertex always increases. We formalize this intuition as follows.

► **Axiom 4** (Continuity). *A centrality measure  $C$  satisfies the continuity axiom if for every graphs  $G$  and  $F$ , and  $v \in V(G)$ , if  $C(v, G) < C(v, G \cup F)$ , then there exists edges  $e_1, \dots, e_k \in E(F)$  such that:  $C(v, G) < C(v, G + e_1) < \dots < C(v, G + e_1 + \dots + e_k) = C(v, G \cup F)$ .*

To the best of our knowledge, the continuity axiom has not been proposed before in the literature. Furthermore, the inclusion and containment property implies the continuity axiom over the corresponding centrality measure.

► **Theorem 13.** *If a family of subgraphs  $\mathcal{F}$  satisfies the inclusion and containment properties, then the corresponding centrality measure  $C_{\mathcal{F}}$  satisfies the continuity axiom.*

All families of subgraphs so far satisfy the inclusion property and their corresponding centrality measures satisfy the continuity axiom as well. We give below a centrality measure based on cliques as a counter-example of this theorem.

► **Example 14.** Cliques are relevant substructure in network analysis and they are usually used to measure the importance of vertices [25]. In [14], this idea has been taken a step further by counting the number of cliques that a vertex belongs, which is called the cross-clique centrality. We can define this centrality with families of subgraphs as follows. Define the family  $\mathcal{K}$  such that  $\mathcal{K}(v, G)$  contains all subgraphs  $K \in \mathcal{A}(v, G)$  such that  $K$  is a clique of size 1 (i.e.  $v$ ) or size greater than 2 for every graph  $G$  and  $v \in V(G)$ . Then the clique centrality of  $v$  on  $G$  is defined as  $C_{\mathcal{K}}(v, G)$ . Note that  $C_{\mathcal{K}}(v, G) = \log(\text{Cross-Clique}(v, G) + 1)$  and, thus, we can use  $C_{\mathcal{K}}$  as a proxy to understand cross-clique centrality.

Cliques  $\mathcal{K}$  is a family that does not satisfy the inclusion property. Indeed, any subgraph of a clique is not necessarily a clique. One can also check that its centrality  $C_{\mathcal{K}}$  also does not satisfy the continuity property. For example, consider a single edge  $G = \bullet\circ$  where the white vertex  $v$  has clique centrality  $C_{\mathcal{K}}(v, G) = 0$ . Then, if a triangle  $F = \blacklozenge$  is added to  $G$ , producing the graph  $G + F = \bullet\blacklozenge$  with  $C_{\mathcal{K}}(v, G + F) = 1$ , there is no way to rise the centrality of  $v$  from 0 to 1 by adding the edges of the triangle one-by-one.

**Size.** The last axiom that we study here is the one proposed in [6] about size. This was formalized as follows: for any  $n > 0$  if we consider clique  $K_n$  and a circuit  $C_n$ , for a centrality measure  $C$  one would expect that  $C(0, K_n) > C(0, C_n)$ . Then no matter how big is  $C(0, K_n)$ , there should exist a value  $m > n$  where the centrality of the cycle  $C_m$  passes the centrality of the clique  $K_n$ , namely,  $C(0, K_n) < C(0, C_m)$ . This argument is related to the size of graphs in the sense that no matter how slow the centrality of  $C_m$  grows, at some point it should beat the clique of size  $n$ . We propose a generalization of this axiom as follows.

► **Axiom 5 (Size).** *A centrality measure  $C$  satisfies the size axiom if for every infinite sequence  $\{G_n\}_{0 \leq n}$  of connected graphs with  $V(G_n) = \{0, \dots, n\}$  and for every value  $N$  there exists  $m$  such that  $C(0, G_m) \geq N$ .*

Here the sequence  $\{G_n\}_{0 \leq n}$  is playing the role of the circuits and  $N$  the role of the centrality in the clique. Thus, if a centrality measure satisfies Axiom 5 then it satisfies the size axiom of [6], but the converse of course is not true.

This axiom is clearly satisfied by all-subgraphs and trees centrality. Indeed, by Proposition 6 we know that  $C_{\mathcal{A}}(v, G)$  is always bounded below by  $\log(n)$  and thus the all-subgraphs satisfy the axiom (similar argument can be given for trees centrality). Typical centrality measures that do not satisfy the size axiom are “local measures” that only consider subgraphs of bounded size, i.e., degree or  $k$ -neighborhood centrality. However, there are families of subgraphs of unbounded size that also do not satisfy this axiom, i.e., clique centrality. In both cases, if we consider the sequence of lines  $\{L_n\}_{0 \leq n}$ , we can see that the centrality on the vertex 0 is not growing and, thus, for a reasonable  $N$  the axiom does not hold. Actually, the next theorem shows that this counter-example is enough to show whether a centrality measure satisfy the size axiom or not.

► **Theorem 15.** *Let  $\mathcal{F}$  be a family of subgraphs that satisfies the line minimization axiom. Then  $C_{\mathcal{F}}$  satisfies the size axiom if, and only if,  $\lim_{n \rightarrow \infty} |\mathcal{F}(0, L_n)| = \infty$ .*

We remark that all centrality measures considered in this paper satisfy the line minimization axiom. Therefore, it is enough to check whether the family of subgraphs grows on the line to see whether the centrality notion is “local” or not.

We want to end this section by pointing out that in [6] it was shown that all standard notions of centrality in the literature (like closeness [4], betweenness [15], Page Rank [9], Katz index [20], etc) do not satisfy at least one of its axioms and, therefore, do not satisfy at least one of the general axioms stated above. This shows that all the standard notions for centrality studied in the literature are different with all-subgraphs centrality.

## 6 Extension to group centrality measures

Any natural centrality measure should come with a simple extension to measure the centrality of sets of vertices (also called *group centrality*). Although this is a desirable property, it is not always clear how to do it (i.e. not many centrality measures in the literature have a standard extension to group centrality). In this section, we embark on extending our families of centralities from vertices to sets and give a natural characterization for all-subgraphs group centrality. Towards the end, we show an application of this notion regarding the centrality maximization of a vertex.

Given an arbitrary family of subgraphs  $\mathcal{F}$ , what should be its extension to groups? A first approach is to consider all connected subgraphs in  $\mathcal{F}$  that contains all elements in the group. Formally, given  $U \subseteq V(G)$  one could consider the family of relevant subgraphs:

$$\mathcal{F}^*(U, G) = \{ S \subseteq G \mid U \subseteq V(S) \wedge \exists v \in U. S \in \mathcal{F}(v, G) \}.$$

In other words, all relevant subgraphs of vertices in  $U$  that cover  $U$ . Although this is the direct extension for connected subgraphs, this definition rises two issues. First, some local families (e.g.  $k$ -neighborhood) could not keep the restriction of having all vertices in  $U$  inside a subgraph (i.e.  $U \subseteq V(S)$ ). Moreover, if the size of  $U$  grows then there will be less subgraphs satisfying such restriction, making the definition impractical for some families of subgraphs. Second, sets that have more relevant subgraphs under this definition are likely to be closer in the graph. For example, if we look at the extension of all-subgraphs  $\mathcal{A}^*(U, G)$ , then in a circuit  $C_n$  a set  $U$  of  $k$ -vertices that has maximum centrality will be any set of  $k$  contiguous vertices. Clearly, if one looks for a central group of  $k$ -vertices in  $C_n$ , one would prefer a set of  $k$ -vertices that are equidistant in  $C_n$  because they cover more relevant structures of the graph as a group.

Given the previous discussion, we define the group extension of  $\mathcal{F}$  to sets of vertices  $U$  on  $G$ , denoted as  $\mathcal{F}(U, G)$ , as follows:

$$\mathcal{F}(U, G) = \{ S \subseteq G \mid U \subseteq V(S) \wedge \forall H \in \text{ConnComp}(S). \exists v \in U. H \in \mathcal{F}(v, G) \}.$$

Note that this extension is similar to the one discussed above (i.e.  $\mathcal{F}^*(U, G)$ ), but we asked that each connected component from  $S$  comes from a relevant subgraph of a vertex in  $U$ . This allows to use disconnected subgraphs to cover  $U$  and, at the same time, each connected component comes from connected subgraphs in  $\mathcal{F}$ . Unlike our first extension, this definition is not local anymore and gives meaningful results for any set  $U$ . In particular, when  $U = \{v\}$  this definition generalizes the family of subgraphs for vertices given that  $\mathcal{F}(U, G) = \mathcal{F}(v, G)$ .

With a family of subgraphs for sets of vertices, it is natural to develop its corresponding group centrality. Similar than for vertices, given a set  $U \subseteq V(G)$  from a graph  $G$ , we define the  $\mathcal{F}$ -group centrality measure of  $U$  in  $G$  as the worst-case entropy of  $\mathcal{F}(U, G)$ , namely:

$$C_{\mathcal{F}}(U, G) = \log(|\mathcal{F}(U, G)|).$$

## 23:14 A Family of Centrality Measures for Graph Data Based on Subgraphs

All families introduced in previous sections have a corresponding group centrality measure. From now, we restrict our analysis to the all-subgraphs family and its centrality over groups, and leave the understanding of other families for future work.

► **Example 16.** Let  $C_n$  be a circuit of length  $n \geq 3$  and consider all sets  $U \subseteq V(C_n)$  of two vertices. Then one can check that the set  $U$  that maximizes  $C_{\mathcal{A}}(U, C_n)$  is any pair of vertices that are at distance  $\frac{n}{2}$  (assuming  $n$  even). Furthermore, if  $U$  are sets of  $k$  vertices with  $k$  a factor of  $n$ , then  $C_{\mathcal{A}}(U, C_n)$  is maximized when all vertices in  $U$  are distributed in  $C_n$  with equal distance. Intuitively, this is the best way of covering a circuit  $C_n$  with  $k$  vertices.

Next we show that all-subgraphs group centrality over  $U$  can be reduced to computing the centrality of a vertex. Recall that we denote by  $G/U$  the set contraction of  $U$  on  $G$ , namely, to merge the vertices  $U$  to one vertex and keeping multi-edges into  $U$  (see Section 2). In particular, recall that  $U$  is a vertex in the multigraph  $G/U$ .

► **Theorem 17.** *Let  $G$  be a graph and  $U \subseteq V(G)$ . Then:*

$$C_{\mathcal{A}}(U, G) = C_{\mathcal{A}}(U, G/U) + |\{e \in E(G) \mid e \subseteq U\}|$$

The all-subgraphs group centrality of a set  $U$  in  $G$  is then reduced to the centrality of  $U$  (i.e. as a vertex) in the set-contraction of  $U$  on  $G$  plus the number of edges between vertices in  $U$ . Note that, in particular, this shows that if we look for  $k$ -sets of high centrality, then the all-subgraphs centrality is balancing between the number of edges of the set (i.e. how similar is the set to a clique) versus how central it is if we contract it into a vertex.

This connection between both definitions (i.e. vertices and sets) for all-subgraphs centrality is strictly related to the properties of the family. Given two subgraphs  $G_1$  and  $G_2$  of  $G$  with  $V(G_1) \cap V(G_2) \neq \emptyset$ , we can generate a new subgraph  $G_1 \cup G_2$  by merging the nodes they share. Unfortunately, this is not possible for all families like the family of trees  $\mathcal{T}$ , that is, the union of two trees is not necessarily in  $\mathcal{T}$ . This means that Theorem 17 cannot be directly extended for families like trees, in particular, for trees centrality.

To end this section, we show an example how the all-subgraphs group centrality allows us to study simple questions regarding the maximization of the centrality of a vertex. Given a graph  $G$  and a vertex  $v \in V(G)$ , with whom should we connect  $v$  in  $G$  in order to maximize its centrality? In other words, if I am in a social network, with whom should I connect in order to maximize my centrality? A naive answer to this question is to connect  $v$  to the most central vertex in  $G$ . Actually, from the perspective of all-subgraphs centrality this is not the right answer: connecting to the most central node will rise its centrality but maybe the centrality of the most central vertex is highly dependent of  $v$ 's centrality. Instead, all-subgraphs centrality says that  $v$  must be connected to the vertex  $u$  where  $\{u, v\}$  (as a group) is more central in  $G$ .

► **Theorem 18.** *Given  $G$  and  $v \in V(G)$  with  $\{u \in V(G) \mid \{u, v\} \notin E(G)\} \neq \emptyset$ , it holds that:*

$$\arg \max_{u \in V(G)} C_{\mathcal{A}}(v, G + \{v, u\}) = \arg \max_{\{u, v\} \notin E(G)} C_{\mathcal{A}}(\{v, u\}, G)$$

## 7 On computing centrality measures based on subgraphs

We study here the problem of computing centrality measures based on subgraphs. In particular, we study the problem of computing the all-subgraphs centrality. We state the problem as follows: given a family of subgraphs  $\mathcal{F}$ , consider the problem

<b>Problem:</b>	$\text{COUNT}(\mathcal{F})$
<b>Input:</b>	A graph $G$ and a vertex $v \in V(G)$
<b>Output:</b>	$ \mathcal{F}(v, G) $

Furthermore, given a class of graphs  $\mathcal{G}$  we write  $\text{COUNT}(\mathcal{F})[\mathcal{G}]$  for the parametrized version of  $\text{COUNT}(\mathcal{F})$  when input graph  $G$  is restricted to  $\mathcal{G}$ . Of course, given a family  $\mathcal{F}$  computing its centrality  $C_{\mathcal{F}}$  requires also taking the logarithm to the output of  $\text{COUNT}(\mathcal{F})$ . Although these are not the same problems, the conclusions obtained here sheds light on the pitfalls of computing a centrality based on a family  $\mathcal{F}$ .

We start by giving an algorithm for computing  $\text{COUNT}$  over all-subgraphs  $\mathcal{A}$ . Algorithm 1 shows a simple recursive algorithm for counting all connected subgraphs that contains a vertex  $v \in V(G)$  in a (multi)graph  $G$ . The main idea is indeed very simple. Recall that  $N(v, G)$  denotes the neighborhood of  $v$  in  $G$  (see Section 2). If  $N(v, G) = \emptyset$ , the vertex  $v$  is an isolated vertex and there is exactly one subgraph. Otherwise,  $v$  is connected to at least one vertex, called it  $u \in N(v, G)$ , and by some edge  $e = \{u, v\}$ . Then we can partition the set of connected subgraphs  $\mathcal{A}(v, G)$  into those that  $u$  and  $v$  are directly connected by some edge, and those that are not. For the former, we can compute the exact number recursively as  $\text{COUNTALL}(G - e, v)$  (recall here that  $G - e$  contains no edges between  $u$  and  $v$ ). For the latter, let  $w(e)$  be the number of edges between  $u$  and  $v$  in  $G$  (recall that  $G$  could be a multigraph). Then all connected subgraphs where  $u$  and  $v$  are directly connected by some edge can be formed by choosing a non-empty set of edges between  $u$  and  $v$  (i.e.  $2^{w(e)} - 1$  many possibilities) plus a connected subgraph from  $\mathcal{A}(e, G/e)$  where  $G/e$  is the set contraction of  $e$  on  $G$  (i.e.  $\text{COUNTALL}(G/e, e)$  many possibilities). Therefore, we can compute  $\text{COUNTALL}(G, v)$  by recursively computing  $\text{COUNTALL}(G - e, v)$  and  $\text{COUNTALL}(G/e, e)$ . In both cases, the number of edges or the number of vertices is reduced, and  $\text{COUNTALL}$  will eventually finish.

Although Algorithm 1 is easy to implement, it could take exponential time in the number of edges. Actually, this is the best that one can hope as we show in the next result. Recall that  $\#\text{P}$  is the class of counting problems that can be defined as counting the number of accepting runs of a polynomial-time non-deterministic Turing machine. Further, a counting problem is  $\#\text{P}$ -complete if it is in  $\#\text{P}$  and all counting problems in  $\#\text{P}$  can be reduced to it [30]. It is known that a polynomial-time algorithm for solving a  $\#\text{P}$ -complete problem, if it existed, would imply that  $\text{P} = \text{NP}$ . For this reason,  $\#\text{P}$ -complete is a class of counting problems considered as hard [2].

► **Theorem 19.**  *$\text{COUNT}(\mathcal{A})$  and  $\text{COUNT}(\mathcal{T})$  are  $\#\text{P}$ -complete.*

This is a negative result for using all-subgraphs centrality or trees centrality in practice. Nevertheless, we believe that this should not overshadow the impact that both measures can have in defining good centrality notions. As we show in Section 5, both notions behaved well as centrality measures and, although they are difficult to compute, they can still be used, for example, to guide the definition of new centrality measures or to design new efficient algorithms for computing the most relevant vertices in a graph.

Given that computing all-subgraphs over any graph is a difficult problem, our next step is to consider classes of graphs  $\mathcal{G}$  where  $\text{COUNT}(\mathcal{F})[\mathcal{G}]$  can be solved efficiently. A natural class to start here are trees. Indeed, when  $G$  is a tree every internal vertex is a cut-vertex and we can use the ideas of Lemma 3 for computing  $|\mathcal{A}(v, G)|$  efficiently. More specific, from Lemma 3 one can show that if  $v$  is a cut-vertex of a graph  $G$  and  $G_1, \dots, G_n$  are subgraphs that partitions  $G$  on  $v$  (i.e.  $V(G) = \cup_{i=1}^n V(G_i)$ ,  $E(G) = \cup_{i=1}^n E(G_i)$ , and

**Algorithm 1** All-subgraphs counting.

---

```

1: Require: A graph  $G$  and vertex  $v \in V(G)$ 
2: procedure COUNTALL( $G, v$ )
3:   if  $N(v, G) = \emptyset$  then
4:     return 1
5:   else
6:     let  $u \in N(v, G)$ 
7:      $e \leftarrow \{u, v\}$ 
8:     return COUNTALL( $G - e, v$ ) +
9:        $(2^{w(e)} - 1) \cdot \text{COUNTALL}(G/e, e)$ 

```

---

**Algorithm 2** All-subgraphs on trees.

---

```

1: Require: A tree  $T$  and vertex  $v \in V(T)$ 
2: procedure COUNTTREES( $T, v$ )
3:   if  $N(v, T) = \emptyset$  then
4:     return 1
5:   else
6:     let  $u \in N(v, T)$ 
7:      $e \leftarrow \{u, v\}$ 
8:     return COUNTTREES( $T - e, v$ ) ·
9:       (COUNTTREES( $T - e, u$ ) + 1)

```

---

$V(G_i) \cap V(G_j) = \{v\}$  for  $i \neq j$ ), then  $\mathcal{A}(v, G) = \prod_{i=1}^n \mathcal{A}(v, G_i)$ . We can exploit this in a tree by considering all subtrees  $T_1, \dots, T_n$  hanging from  $v$  and computing  $\mathcal{A}(v, G)$  as the product of  $\mathcal{A}(v, T_i)$ .

In the procedure COUNTTREES of Algorithm 2 we use the previous idea for computing  $|\mathcal{A}(v, T)|$  when  $T$  is a tree and  $v \in V(T)$ . It follows a similar approach to that in Algorithm 1. First, if  $v$  is an isolated vertex (i.e.  $N(v, T) = \emptyset$ ), then it outputs 1. Otherwise, it takes a vertex  $u \in N(v, T)$ , defines the edge  $e = \{u, v\}$ , and decompose  $T$  in two subtrees by removing  $e$  from the graph. Notice that, if we remove  $e$  from  $T$ , we create two connected components  $T_v$  and  $T_u$ , where  $T_v$  and  $T_u$  contains  $v$  and  $u$ , respectively. One can easily check that  $T_v$  and  $T_u + e$  partitions  $T$  on  $v$  and we have  $|\mathcal{A}(v, T)| = |\mathcal{A}(v, T_v)| \cdot |\mathcal{A}(v, T_u + e)|$  by the previous discussion above. Furthermore, it is straightforward to check that  $|\mathcal{A}(v, T_v)| = |\mathcal{A}(v, T - e)|$  and  $|\mathcal{A}(v, T_u + e)| = |\mathcal{A}(u, T - e) + 1|$ . Thus, we can compute  $|\mathcal{A}(v, T)|$  by recursively computing COUNTTREES( $T - e, v$ ) multiplied by COUNTTREES( $T - e, u$ ) + 1.

In contrast to Algorithm 1, the recursion in COUNTTREES separates the graph in two disjoint subtrees. This implies that the recursion eventually finishes and, moreover, it takes linear time in the size of the tree. Interestingly, we can extend this idea to any graph of bounded tree-width. To formalize the notion of bounded tree-width, we need to introduce some notation. Given a graph  $G$ , a tree decomposition  $T$  of  $G$  is a tree such that  $V(T)$  are sets of  $V(G)$  (i.e.  $X \subseteq V(G)$  for every  $X \in V(T)$ ) and satisfies the following three properties: (1)  $V(G) = \bigcup_{X \in V(T)} X$ , (2) if  $v \in X \cap Y$  for  $X, Y \in V(T)$ , then  $v \in Z$  for all  $Z \in V(T)$  in the simple path from  $X$  to  $Y$  in  $T$ , and (3) for every  $e \in E(G)$ , there exists  $X \in V(T)$  such that  $e \subseteq X$ . The width of a tree decomposition  $T$  is equal to  $\max_{X \in V(T)} |X| - 1$  and the tree-width  $\text{tw}(G)$  of  $G$  is the minimum width among all possible tree decompositions of  $G$  [26]. A class  $\mathcal{G}$  of graphs has bounded tree width if there exists a uniform bound  $k$  such that  $\text{tw}(G) \leq k$  for every  $G \in \mathcal{G}$ . For example, all trees is a class that has tree-width bounded by 1.

► **Theorem 20.** *If  $\mathcal{G}$  has bounded tree-width, then  $\text{COUNT}(\mathcal{A})[\mathcal{G}]$  can be solved in PTIME.*

The previous result shows that the problem becomes tractable when graphs has bounded tree-width. Despite that graphs have high tree-width in practice [23], this result gives some clues on how to tackle the problem of computing the all-subgraphs centrality.

## 8 Future work

This work arises several research opportunities regarding centrality measures based on subgraphs, which are briefly discussed here. One of the most important question is whether all-subgraphs centrality can be approximated efficiently, or even if the rank order given by this measure can be approximated. Another interesting question is to consider when a family



of graphs can approximate another family over some particular class of graphs (e.g. plain graphs). For the sake of simplification, we only considered undirected graphs but another relevant question is to study how to extend these results to directed graphs or to hypergraphs. Furthermore, the initial motivation of our approach came from centrality measures for graph query languages, but in order to incorporate this approach, several properties must be understood like, for example, how to mix the centrality measures to the output of a query. Finally, it would be interesting to consider a randomized version of our approach where not all subgraphs have the same chances to appear. Instead of considering the worst-case entropy, one could study the entropy of a family given a particular distribution and study their properties. We leave this and other questions for future work.

---

## References

- 1 Renzo Angles, Marcelo Arenas, Pablo Barceló, Aidan Hogan, Juan L. Reutter, and Domagoj Vrgoc. Foundations of Modern Query Languages for Graph Databases. *ACM Comput. Surv.*, 50(5):68:1–68:40, 2017.
- 2 Sanjeev Arora and Boaz Barak. *Computational complexity: a modern approach*. Cambridge University Press, 2009.
- 3 Sambaran Bandyopadhyay, Ramasuri Narayanam, and M Narasimha Murty. A Generic Axiomatic Characterization for Measuring Influence in Social Networks. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 2606–2611. IEEE, 2018.
- 4 Alex Bavelas. Communication patterns in task-oriented groups. *The Journal of the Acoustical Society of America*, 22(6):725–730, 1950.
- 5 Paolo Boldi, Alessandro Luongo, and Sebastiano Vigna. Rank monotonicity in centrality measures. *Network Science*, 5(4):529–550, 2017.
- 6 Paolo Boldi and Sebastiano Vigna. Axioms for centrality. *Internet Mathematics*, 10(3-4):222–262, 2014.
- 7 Stephen P Borgatti and Martin G Everett. A graph-theoretic perspective on centrality. *Social networks*, 28(4):466–484, 2006.
- 8 Ulrik Brandes. *Network analysis: methodological foundations*, volume 3418. Springer Science & Business Media, 2005.
- 9 Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1-7):107–117, 1998.
- 10 Carlos Buil-Aranda, Martin Ugarte, Marcelo Arenas, and Michel Dumontier. A preliminary investigation into SPARQL query complexity and federation in Bio2RDF. In *Alberto Mendelzon International Workshop on Foundations of Data Management*, page 196, 2015.
- 11 Steve Chien, Cynthia Dwork, Ravi Kumar, Daniel R Simon, and D Sivakumar. Link evolution: Analysis and algorithms. *Internet mathematics*, 1(3):277–304, 2004.
- 12 Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.
- 13 Ernesto Estrada and Juan A Rodriguez-Velazquez. Subgraph centrality in complex networks. *Physical Review E*, 71(5):056103, 2005.
- 14 Mohammad Reza Faghani and Uyen Trang Nguyen. A study of XSS worm propagation and detection mechanisms in online social networks. *IEEE transactions on information forensics and security*, 8(11):1815–1826, 2013.
- 15 Linton C Freeman. A set of measures of centrality based on betweenness. *Sociometry*, pages 35–41, 1977.
- 16 Linton C Freeman. Centrality in social networks conceptual clarification. *Social networks*, 1(3):215–239, 1978.
- 17 Manuj Garg. Axiomatic foundations of centrality in networks. Available at SSRN 1372441, 2009.

- 18 Georg Gottlob, Gianluigi Greco, Nicola Leone, and Francesco Scarcello. Hypertree Decompositions: Questions and Answers. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 57–74, 2016.
- 19 Hawoong Jeong, Sean P Mason, A-L Barabási, and Zoltan N Oltvai. Lethality and centrality in protein networks. *Nature*, 411(6833):41, 2001.
- 20 Leo Katz. A new status index derived from sociometric analysis. *Psychometrika*, 18(1):39–43, 1953.
- 21 Harold J Leavitt. Some effects of certain communication patterns on group performance. *The Journal of Abnormal and Social Psychology*, 46(1):38, 1951.
- 22 Johannes Lorey and Felix Naumann. Detecting SPARQL query templates for data prefetching. In *Extended Semantic Web Conference*, pages 124–139. Springer, 2013.
- 23 Silviu Maniu, Pierre Senellart, and Suraj Jog. An Experimental Study of the Treewidth of Real-World Graph Data. In *22nd International Conference on Database Theory, ICDT 2019, March 26-28, 2019, Lisbon, Portugal*, pages 12:1–12:18, 2019.
- 24 Gonzalo Navarro. *Compact data structures: A practical approach*. Cambridge University Press, 2016.
- 25 Mark Newman. *Networks: an introduction*. Oxford university press, 2010.
- 26 Neil Robertson and Paul D Seymour. Graph minors. III. Planar tree-width. *Journal of Combinatorial Theory, Series B*, 36(1):49–64, 1984.
- 27 Gert Sabidussi. The centrality index of a graph. *Psychometrika*, 31(4):581–603, 1966.
- 28 Oskar Skibski, Talal Rahwan, Tomasz P Michalak, and Makoto Yokoo. Attachment centrality: An axiomatic approach to connectivity in networks. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, pages 168–176. International Foundation for Autonomous Agents and Multiagent Systems, 2016.
- 29 Oskar Skibski and Jadwiga Sosnowska. Axioms for distance-based centralities. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- 30 Leslie G Valiant. The complexity of computing the permanent. *Theoretical computer science*, 8(2):189–201, 1979.
- 31 Tomasz Waś and Oskar Skibski. An axiomatization of the eigenvector and Katz centralities. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.