

This the author's accepted version of the following article Calderoni, F., Catanese, S., De Meo, P., Ficara, A., Fiumara, G. (2020). Robust link prediction in criminal networks: A case study of the Sicilian Mafia. *Expert Systems with Applications*, 161, 113666. <https://doi.org/10.1016/j.eswa.2020.113666>

Robust Link Prediction in Criminal Networks: a case study of the Sicilian Mafia

Francesco Calderoni^a, Salvatore Catanese^b, Pasquale De Meo^{c,*}, Annamaria Ficara^b, Giacomo Fiumara^b

^aTranscrime, Università Cattolica del Sacro Cuore, Largo Gemelli 1, 20123, Milan, Italy

^bMIFT Department, University of Messina, V.le F. Stagno D'Alcontres, 31, 98166, Messina, Italy

^cDepartment of Modern and Ancient Civilizations, Polo Universitario Annunziata, 98122, Messina, Italy

Abstract

Link prediction exercises may prove particularly challenging with noisy and incomplete networks, such as criminal networks. Also, the link prediction effectiveness may vary across different relations within a social group. We address these issues by assessing the performance of different link prediction algorithms on a mafia organization. The analysis relies on an original dataset manually extracted from the judicial documents of operation “Montagna”, conducted by the Italian law enforcement agencies against individuals affiliated with the Sicilian Mafia. To run our analysis, we extracted two networks: one including meetings and one recording telephone calls among suspects, respectively. We conducted two experiments on these networks. First, we applied several link prediction algorithms and observed that link prediction algorithms leveraging the full graph topology (such as the Katz score) provide very accurate results even on very sparse networks. Second, we carried out extensive simulations to investigate how the noisy and incomplete nature of criminal networks may affect the accuracy of link prediction algorithms. The experimental findings suggest the soundness of link predictions is relatively high provided that only a limited amount of knowledge about connections is hidden or missing, and the unobserved edges follow some kind of generative law. The different results on the meeting and telephone call networks indicate that the specific features of a network should be taken into careful consideration.

Keywords: Criminal Networks, Social Network Analysis, Network Science, Link Prediction in Uncertain Graphs.

1. Introduction

Methods from Social Network Analysis (in short, SNA) (Sparrow, 1991; Klerks, 2001; Xu & Chen, 2005; Van der Hulst, 2009; Agreste et al., 2016; Berlusconi et al., 2016) greatly contributed to intelligence and criminal investigations: for example, SNA allows to identify, within a criminal network, the most central members in terms of connections or information flow (Calderoni & Superchi, 2019), the presence of different communities (Catanese et al., 2014; Calderoni et al., 2017), and the most efficient strategies for dismantling the network (Agreste et al., 2016).

A crucial application of SNA methods to intelligence is the so-called *link prediction problem* (Liben-Nowell & Kleinberg, 2003; Pandey et al., 2019): given a graph G which describes interactions between pairs of criminals, we wish to predict which edges are more likely to appear in G in the near future. Algorithms to solve the link prediction problem may hugely impact police activities: in fact, if we would be able to accurately predict the formation of new links, we would be able to discover pairs of criminals who are likely to collaborate and, thus, we could early detect and prevent crimes.

*Corresponding author

Email addresses: francesco.calderoni@unicatt.it (Francesco Calderoni), salvocatanese@gmail.com (Salvatore Catanese), pdemeo@unime.it (Pasquale De Meo), aficara@unime.it (Annamaria Ficara), gfiumara@unime.it (Giacomo Fiumara)

13 Many link prediction methods have been designed and implemented in a broad range of domains (see the excellent
14 reviews by Liben-Nowell & Kleinberg (2003) and Pandey et al. (2019)) and, more recently, Berlusconi et al. (2016)
15 applied link prediction algorithms on a dataset derived from an Italian criminal case against a Mafia group.

16 Almost all of the existing approaches to link prediction focus on maximizing the accuracy and they overlook
17 fundamental aspects such as the *robustness* of predictions, namely the extent to which the incompleteness of informa-
18 tion about relations may affect the quality of predictions. By construction, in fact, datasets associated with criminal
19 networks are noisy and incomplete: on one hand, investigations often encounter individuals unrelated to the crimi-
20 nal organization (e.g. friends, relatives, and other frequent contacts) and, on the other hand, some members of the
21 organization actively attempt to avoid detection, e.g. by refraining from the use of telephone, using intermediaries,
22 and coding messages. As a consequence, imprecise and incomplete information is a critical impediment to under-
23 stand network boundaries and topology and, ultimately, it constitutes a main challenge for law enforcement agencies
24 (hereafter LEAs) which plan to get reliable results from the application of link prediction algorithms.

25 In this paper we tackle the problem of estimating the robustness of link prediction algorithm in criminal networks.
26 To do so, we analysed judicial sources on operation “Montagna”, a long investigation on a large criminal organization
27 belonging to Cosa Nostra (i.e., the Sicilian Mafia) (Paoli, 2004, 2008) active in the north of Sicily. We extracted two
28 graphs (Ficara et al., 2020): the former (called *Meeting Graph* G_M) maps meetings between person under investiga-
29 tions and the latter (called *Phone Call Graph* G_P) is built on the monitoring of phone communications (also known
30 as *wiretapping*). Our dataset is unique and we believe it might represent a valuable resource for better understanding
31 complex criminal phenomena from a quantitative standpoint.

32 We applied many classical methods of link prediction (Liben-Nowell & Kleinberg, 2003) on both G_M and G_P such
33 as the Common Neighbors (CN), the Jaccard Coefficient (JC), the Adamic Adar (AA) coefficient, the Preferential
34 Attachment (PA), the Katz score, and more recent methods such as the *Node2Vec* (Grover & Leskovec, 2016) graph
35 embedding algorithm and the Personalized PageRank (PPR) similarity score (Avrachenkov et al., 2019). We used
36 the *Area Under the Receiving Operating Curve* (AUROC) (Fawcett, 2006) to assess the accuracy of a prediction
37 algorithm. We recall that the AUROC ranges from 0 to 1 and the larger the AUROC, the more accurate a link
38 prediction algorithm.

39 Subsequently, we assumed that both G_M and G_P graphs are not completely known and introduced some generative
40 models to study the incomplete information about connections among criminals that we define as *uncertainty*. Specif-
41 ically, we considered G_M (resp., G_P) as a sample of a true graph G'_M (resp., G'_P) and we introduced a scoring function
42 (called *likelihood*) to decide whether an edge non observed in G_M (resp., G_P) actually exists in G'_M (resp., G'_P). A
43 core assumption of our method is that we know all nodes of the “real” networks, but we have an incomplete knowl-
44 edge of the edges of the observed networks, namely of the interconnections among criminals. Our assumptions are
45 backed on previous research results showing that LEAs may rarely miss important individuals in a well-built criminal
46 investigation (Campana & Varese, 2012; Berlusconi, 2013) as well as on the length and relevance of the investigation
47 pursued in “Montagna” operation. Conversely, LEAs often need to identify the relevant ties among thousands of
48 communications and meetings and this process may be biased by lack of resources or by criminals’ strategies to pre-
49 vent detection. We resort to simulation to create graphs G'_M and G'_P and, in our experimental analysis, we considered
50 multiple likelihood functions. Our simulation method allowed also to specify the fraction p of non observed edges in
51 G'_M (resp., G'_P) which are actually placed in G'_M (resp., G'_P).

52 The main findings of our analysis are as follows:

- 53 1. The Katz and PPR scores prove as the most successful method to predict missing edges and they achieve an
54 AUROC larger than 0.95, thus signalling a high degree of accuracy. It is worth observing that the highest
55 AUROC is achieved when only short paths are taken into account, which is consistent with the short-range
56 structure of criminal networks.
- 57 2. Graph topology significantly affects the accuracy of a link prediction algorithm: specifically, algorithms which
58 are very accurate on G_M performs badly on G_P and vice versa. In detail, if a graph is poorly connected (i.e.,
59 it displays a low edge density and a small clustering coefficient), then *local methods* (i.e., link prediction algo-
60 rithms which rely only on the local knowledge of graph topology) are to be preferred to *global* ones (i.e., link
61 prediction algorithms leveraging the knowledge of the full graph topology). Vice versa, global methods such
62 as the Katz score achieve their best AUROC on graphs with greater levels of connectivity (i.e. in graph with
63 higher edge density and clustering coefficients).

64 3. The knowledge LEAs have acquired about the Meeting Network G_M is quite complete whereas the Phone
65 Call Network G_P is more susceptible to uncertainty, encoded in our study through the parameter p . We can
66 generalize such a result: if the amount of uncertainty is relatively small and non-observed edges derive from a
67 specified generative model we can hope for robust edge prediction. In the light of our studies, we recommend
68 LEAs not only to build a detailed map of connections between criminals but also to investigate how such a map
69 evolves over time: in this way we would be able to design sophisticated likelihood functions which fits fairly
70 well experimental observations and help LEAs to early detect and prevent crimes.

71 The rest of this paper is organized as follows: in Section 2 we present the related literature, whereas in Section
72 3 we outline our research questions. In Section 4 we introduce some basic definitions from graph theory. In Section
73 5 we introduce the operation “Montagna” and provide some details about the Meeting and Phone Call Graphs. In
74 Sections 6 we compare the accuracy of some link prediction methods on our graphs. We then present our experiments
75 aimed at modelling missing edges in Section 7. We draw our conclusions and highlight some possible future work in
76 Section 8.

77 2. Related Work

78 Social Network Analysis (SNA) is increasingly used by law enforcement agencies (LEAs) to analyze criminal
79 networks as well as to investigate on the relations among criminals based on calls, meetings and other events derived
80 from investigations (Sparrow, 1991; Xu & Chen, 2005; Van der Hulst, 2009; Strang, 2014).

81 Given the social embeddedness of organized crime and, in particular, of Mafia-like organizations, the analysis
82 of the social structure of Sicilian Mafia syndicates generated a great scientific interest (Kleemans & Bunt, 1999;
83 Kleemans & De Poot, 2008). For instance, Morselli (2003) studied the connections within a New York-based family
84 (the Gambino family). The study focused on the career of one of its members, Saul Gravano. One of the main
85 findings is that Gravano’s ability of building and extending over time his personal network of contacts was a key
86 factor to climbing the Gambino’s family organization. Calderoni (2012) showed that high status Mafia members were
87 able to indirectly manage illicit drug traffics leaving in more central and visible position middle-level criminals.

88 SNA is not only a tool to describe the structure and functioning of a criminal organizations but it has been largely
89 employed in the construction of crime prevention systems (Chen et al., 2004). For instance, Xu & Chen (2005) jointly
90 applied SNA with hierarchical clustering algorithms. The proposed approach worked in two stages: first, a criminal
91 network was partitioned into subgroups by means of a clustering algorithm. Then, block modelling techniques have
92 been used to extract interaction patterns between these subgroups. Agreste et al. (2016) applied percolation theory to
93 efficiently dismantle Mafia syndicates. Calderoni & Superchi (2019) showed that the node’s betweenness centrality
94 in a meeting network is evidence of Mafia leadership, suggesting that this variable could be exploited by LEAs in
95 selecting the most suitable targets for additional investigations and disruption. Duxbury & Haynie (2019) used an
96 agent-based model to evaluate how criminal networks recover from disruption and identified which disruption strate-
97 gies are most effective at damaging various criminal networks. Grassi et al. (2019) explored different betweenness
98 centrality including the classic betweenness by Freeman (1979) and three inspired by the dual projection approach
99 recently suggested by Everett & Borgatti (2013), which proved to be more successful than classic approaches in iden-
100 tifying the criminal leaders. Bouchard (2020) used a network approach to specify and model collaboration among
101 people involved in organized crime. His approach provides methodological guidelines for clarifying boundaries and
102 helps solve four puzzles: social boundaries, boundaries of group membership, ethnic boundaries and recruitment.
103 Overall, while the studies above provided insight into the social organization and possible countermeasures against
104 criminal groups, the application of SNA to them nearly inevitably faces problems of *noisy or incomplete information*.

105 Information on a criminal network is often likely to be missing or hidden, due to the covert and stealthy nature of
106 criminal actions (Krebs, 2002; Xu & Chen, 2005). Consequently, the derived networks are incomplete, incorrect, and
107 inconsistent, either due to deliberate deception on the part of criminals, or to limited resources or unintentional errors
108 by LEAs (Calderoni, 2010; Campana & Varese, 2012; Catanese et al., 2014; Ferrara et al., 2014; Agreste et al., 2016).
109 These limitations may bias the analysis and they cause problems of uncertain information, potentially jeopardizing
110 the effectiveness of the investigations (Strang, 2014).

111 In the analysis of criminal networks, missing data can refer to missing *nodes* and/or missing *edges* (Calderoni,
112 2010). The problem of missing nodes has already received attention (Kim & Leskovec, 2011; Hric et al., 2016) but,

113 from our perspective, missing nodes in criminal networks are not particularly relevant, in that it is quite improbable
 114 that LEAs may disregard central criminals during such prolonged investigations. On the other hand, while it is possible
 115 to predict some missing edges among already detected criminals, it is impossible to detect missing criminals relying
 116 only on pre-trial detention orders. Missing edges refer to the lack of information on the relations between two known
 117 criminals. LEAs, in fact, may miss a lot of criminal activities such as meetings or phone calls, and therefore relevant
 118 plans of the criminal organization (Campana & Varese, 2012; Ferrara et al., 2014; Catanese et al., 2014; Agreste et al.,
 119 2016). For instance, criminals may use different telephone lines, according to the nature of the conversation and the
 120 interlocutor, and investigators may be able to identify only some of them. The frequent change of mobile phones and
 121 SIM cards and the use of particular lines to communicate with high-ranking affiliates may also prevent LEAs from
 122 identifying all conversations among suspects.

123 Several recent contributions have addressed the problem of missing links with particular attention to criminal
 124 networks. Fan et al. (2017) proposed a combined link prediction index considering both the nodes' types effects and
 125 nodes' structural similarities, and demonstrated that it is remarkably superior to all the 25 existing similarity-based
 126 methods both in predicting missing links and identifying spurious links in a real military network data. This method is
 127 also suitable to many other social organizations, such as criminal networks. Marciani et al. (2017) proposed three new
 128 similarity social network metrics, specifically tailored for criminal link detection and prediction and evaluated them
 129 through a flexible data stream processing application observing that the new metrics could reach up to 83% accuracy
 130 in detection and 82% accuracy in prediction and be competitive with the state of the art metrics. Bahulkar et al. (2018)
 131 described a framework which predicts the missing links in the social network data and then algorithms are applied
 132 to the augmented data to detect the communities of a transnational criminal organization (TCO). Parisi et al. (2018)
 133 proposed an entropy-based method to predict a given percentage of missing links from an observed network structure,
 134 by identifying them with the most probable non-observed ones. Diviák (2019) tried and followed a systematic and
 135 general solution to deal with the problem of missing data. He proposed three potentially synergistic and combinable
 136 techniques for data collection for each stage of data collection – biographies for data extraction, graph databases
 137 for data storage, and checklists for data reporting. Lim et al. (2019) explored the application of deep reinforcement
 138 learning (DRL) in developing a criminal network hidden links prediction model from the reconstruction of a corrupted
 139 criminal network dataset. De Moor et al. (2020) also considered the problem of missing data in criminal networks.
 140 They compared the statistics on a reduced or incomplete network with those from a known network, integrating
 141 police data on known offenders with DNA data on unknown offenders and showing how networks with both known
 142 and unknown offenders are bigger but also have a different structure to networks with only known offenders.

143 Despite the growing scholarly attention to missing links in criminal networks and the important consequences of
 144 missing relations, there is no previous research work aiming at modelling uncertainty in criminal networks and how
 145 such an uncertainty affects the analysis of a criminal network. Our work aims at filling this gap and, to the best of our
 146 knowledge, it represents the first step toward an objective assessment of the robustness of a link prediction algorithm
 147 in a criminal network.

148 3. Research Questions

149 We assume that *multiple interactions* can be observed among the members of a Mafia group: typical interactions
 150 are phone calls, meetings, interactions on Social Media platforms, financial transactions and so on. For each interac-
 151 tion, we represent a Mafia group as a graph in which a node identifies an individual who belongs (or is close) to the
 152 group and an edge indicates that an interaction occurred between two individuals represented by the nodes tied by that
 153 edge.

154 If we assume that K type of interactions exist, we can understand a Mafia group as a collection of graphs $\mathcal{G} =$
 155 $\{G_A, G_B, \dots, G_K\}$. We call each graph $G_D \in \mathcal{G}$ as *interaction graph*. The topology of each interaction graph G_D has a
 156 remarkable value for investigation purposes: for instance, we could use the topology of G_D to: (i) infer collaboration
 157 between criminals, (ii) to identify individuals who intercept most of the information flow in the criminal organization
 158 and (iii) to design police operations to dismantle the underlying criminal organization.

159 In this paper we will concentrate on task (i), and, more specifically, we seek an answer to the following question:

160 **Q₁** *Given an interaction graph $G_D \in \mathcal{G}$, can we design algorithms to accurately predict edges between*
 161 *pairs of criminals?*

162 Our research question is strictly linked to a popular problem in Network Science, namely the *link prediction prob-*
 163 *lem in networks* (Liben-Nowell & Kleinberg, 2003). The link prediction problem has been extensively studied in a
 164 number of domains, such as e-commerce (Chen et al., 2005), homeland security (Hasan et al., 2006) and bioinfor-
 165 matics (Menon & Elkan, 2011). In the criminology the problem of predicting links is well known but there are few
 166 studies approaching it (Rhodes & Jones, 2009; Berlusconi et al., 2016).

167 Given a link prediction algorithm \mathcal{A} , one could ask if the accuracy of \mathcal{A} depends on the topological features of the
 168 interaction graph, say G_D , on which we decide to run \mathcal{A} . More formally, we are interested in answering the following
 169 question:

170 **Q₂** *Given a Mafia network $\mathcal{G} = \{G_A, G_B, \dots, G_K\}$ and a link prediction algorithm \mathcal{A} , how does the*
 171 *topological features of a dimension graph $G_D \in \mathcal{G}$ impact on the accuracy of \mathcal{A} ?*

172 Finally, a fundamental issue of criminal networks is that they are incomplete and noisy and, thus, a key scientific
 173 challenge is about the *robustness* of the results that the algorithm \mathcal{A} produces. In other words, we are interested in
 174 estimating how uncertainties in the topology of G_D impact of the accuracy of \mathcal{A} and such a reasoning leads us to
 175 formulate the following question:

176 **Q₃** *How sensitive is the link prediction algorithm \mathcal{A} to uncertainty in G_D ? Can police forces get reliable*
 177 *outcomes when they apply the algorithm \mathcal{A} on a specified dimension graph G_D ?*

178 In the next subsections we illustrate the outcomes of our study to answer **Q₁-Q₃**; our study builds upon a case
 179 study draw from “Montagna”, a law enforcement operation tackling Mafia gangs in the North of Sicily.

180 4. Fundamentals of Graph Theory

181 In this section we introduce some basic definitions from graph theory which will be largely used throughout the
 182 paper.

183 **Definition 1 (Graphs).** *A graph $G = \langle N, E \rangle$ is a pair in which N is the set of nodes and $E \subseteq N \times N$ is the set of*
 184 *edges. A graph is undirected if $\langle i, j \rangle \in E$ implies that $\langle j, i \rangle \in E$ for each pair of nodes i and j , directed otherwise. The*
 185 *non-edge set $T \subseteq N \times N$ is the complement of E , i.e., $T = \{\langle i, j \rangle : i \in N, j \in N \wedge \langle i, j \rangle \notin E\}$.*

186 Given a graph $G = \langle N, E \rangle$, we say that a graph $G' = \langle N', E' \rangle$ is a *subgraph* of G if $N' \subseteq N$ and $E' \subseteq E$.

187 In this paper we consider only undirected graphs. A graph G is associated with an adjacency matrix \mathbf{A} , whose
 188 entries are defined as follows: $\mathbf{A}_{ij} = 1$ if and only if $\langle i, j \rangle \in E$, 0 otherwise. The *order* of a graph G is the number
 189 $n = |N|$ of its nodes and the *size* of G is defined as the number $m = |E|$ of its edges. A graph with n nodes may contain
 190 at most $\binom{n}{2}$ edges and it is said *complete*. The ratio $\delta = \frac{m}{\binom{n}{2}}$ is known as *graph density*: if δ is $O(n^{-1})$ we say that the
 191 graph is *sparse*, *dense* otherwise. We define the *neighbour-set* $N(i)$ of i as the set of nodes connected to i and the
 192 *degree* d_i of i as $d_i = |N(i)|$. The average degree \bar{d} of a graph G is defined as $\bar{d} = \frac{1}{n} \sum_{i=1}^n d_i$. We also define the *local*
 193 *clustering coefficient* (Watts & Strogatz, 1998) of G as follows:

194 **Definition 2 (Clustering Coefficient).** *Let G be a graph with non-edge set T . Let $i \in N$ be a node in G with*
 195 *neighbour-set $N(i)$. Let us define the set $S(i) = \{\langle j, k \rangle : j \in N(i), k \in N(i), \langle j, k \rangle \in E\}$, i.e., $S(i)$ contains pairs of*
 196 *nodes j and k which are both connected to i and which are connected through an edge.*

197 *The local clustering coefficient $lc(i)$ of i is defined as follows:*

$$lc(i) = \frac{|S(i)|}{\binom{d_i}{2}} \quad (1)$$

198 *The average clustering coefficient ac is defined as:*

$$ac = \frac{1}{|N|} \sum_{i \in N} lc(i) \quad (2)$$

| <i>Parameter</i> | G_M | G_P |
|-------------------------------------|-------|-------|
| Number of Nodes $ N $ | 101 | 100 |
| Number of Edges $ E $ | 256 | 124 |
| Average Degree \bar{d} | 5.07 | 2.48 |
| Density δ | 0.051 | 0.025 |
| Average Clustering Coefficient ac | 0.656 | 0.105 |

Table 1: Statistics of G_M and G_P graphs.

The clustering coefficient $lc(i)$ quantifies to which extent the neighbours of the node i tend to form a tie. It is of crucial importance in the study of criminal networks because, due to previous criminological studies (Berlusconi, 2013; Agreste et al., 2016), we expect that criminals form dense clusters, thus implying large values of clustering coefficients.

For our research purposes we are also interested in more complex structures, such as *walks* and *paths*, which are specified below:

Definition 3 (Walks and Paths). Let G be a graph. A walk of length $r - 1$ in G is a sequence i_0, i_1, \dots, i_r of nodes such that $\langle i_x, i_{x+1} \rangle \in E$ for each $x \in \{0, 1, \dots, r\}$. A walk with no-repeated nodes is called path.

5. The Montagna Police Operation

Our case study concerns the anti-mafia operation called “Montagna” concluded in 2007 by the Public Prosecutor’s Office of Messina (Sicily) and conducted by the R.O.S. (*Reparto Operativo Speciale*, or Special Operations Group, a specialized anti-mafia police unit of the Italian Carabinieri). The investigation, one of the most important of the period, focused on the Cosa Nostra groups known as “*Mistretta*” family (hereafter, clan A) and the “*clan Batanesi*” (hereafter, clan B).

From 2003 to 2007, these families had infiltrated several economic activities including the public works in the area, through a cartel of entrepreneurs close to Cosa Nostra. The groups engaged in extortion racketeering and provided illegal protection to achieve illegal profits from the public construction works, with dynamics similar to those described by Gambetta (1993) and Gambetta & Reuter (1995). Furthermore, the investigation showed that the “*Mistretta*” family had taken on the role of mediator between the Cosa Nostra families of Palermo and Catania and the other criminal organizations around Messina. Indeed both the *Mistretta* and *Batanesi* families had close connections with other Cosa Nostra families located in the province of Messina, namely the “*Barcellona*” family (hereafter clan C), and the “*Caltagirone*” family (hereafter clan D). The charges were upheld by several trials and the majority of the individuals have been sentenced to long prison terms.

The main data source is the pre-trial detention order by the Court of Messina’s preliminary investigation judge issued on March 14, 2007 towards the end of the investigation. The order concerned a total of 52 suspects, all charged with the crime of participation in a Mafia clan (Article 416 bis of the Italian Criminal Code) as well as other crimes (e.g., theft, extortion, damaging followed by arson). According to the Italian Criminal code, the affiliation to a Mafia clan carries a penalty of between ten and fifteen years of imprisonment. The Court ordered the pre-trial detention for 38 individuals and provided detailed motivations for the decision in a document of more than two hundred pages with an important amount of information about the suspects’ crimes, activities, meetings, and calls.

Most of the information from judicial documents were about clan A and clan B. From the analysis of legal documents we built two graphs: *a) the Meeting Graph G_M* , in which nodes are uniquely associated with suspected criminals and edges specify meetings among individuals *b) the Phone Call Graph G_P* , in which nodes are uniquely associated with suspected criminals and edges records phone calls between pairs of individuals (Ficara et al., 2020).

The G_M graph had 101 nodes and 256 edges while G_P had 100 nodes and it contained only 124 edges. There were 47 individuals who jointly belonged to G_M and G_P ; some statistics about G_M and G_P are displayed in Table 1.

Nodes in G_M and G_P can take active roles in the criminal organization: for instance, some nodes correspond to individuals who can be classified as “boss” (i.e., leaders of the criminal organization) while others are classified as “picciotti” (i.e., soldiers of the organization). Of course, some individuals can have regular contacts with members of

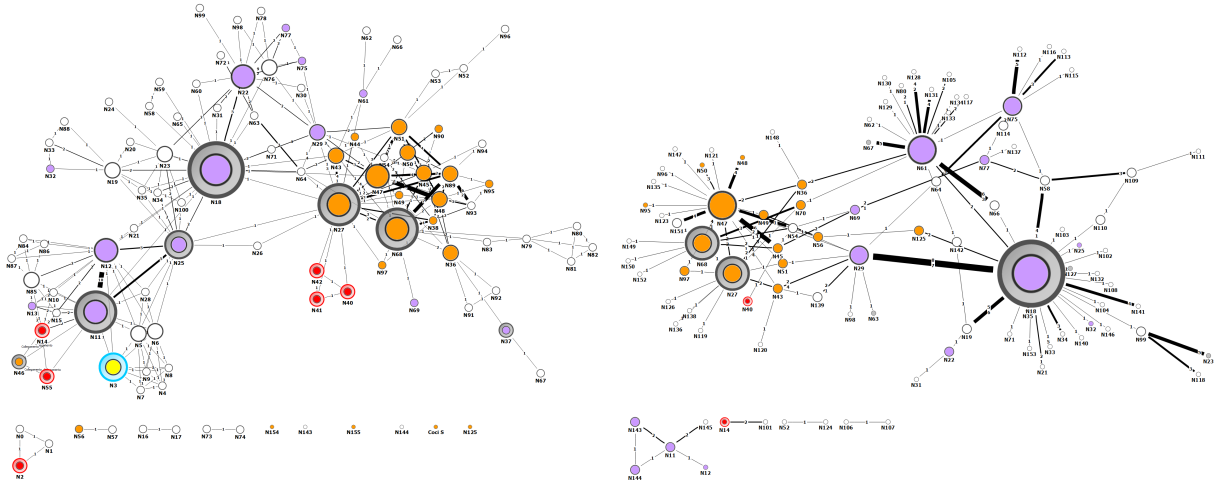


Figure 1: **Left Panel.** A graphical representation of the G_M graph. In G_M , nodes which represent the members of the “Mistretta” family and the “Batanesi” family are highlighted in violet and orange, respectively. Circled nodes correspond to the subjects investigated for having promoted, organized and directed the Mafia association (leaders). The red and yellow circled nodes refer to bosses of Mafia families of other districts. The white knots represent the other subjects considered to be: i) close to the association and ii) not classifiable in any of the previous categories, but nevertheless useful for the purposes of the Mafia-type association and the realization of its plans. **Right Panel.** A graphical representation of the G_M graph. The color of nodes has the same meaning as in the G_M . In both G_M and G_P , the width of the edges is proportional to the number of meetings (or phone calls) and the size of the nodes to their degree.

238 a criminal organization (due, for instance, to kinship relations) but they are not involved in any criminal activity. To
 239 guarantee anonymity, we used the symbol Nx (being x an integer) to identify an individual in both G_P and G_M .

240 In Figure 1 we graphically report the G_M graph (left panel) and the G_P graph (right panel): here the width of
 241 a node is proportional to its degree, while the width of an edge is proportional to the total number of meetings (or
 242 telephone calls) recorded between the nodes that edge connects. Members of the “Mistretta” and “Batanesi” families
 243 are colored in orange and purple, respectively.

244 6. Link Prediction in Montagna

245 In this section we consider the problem of predicting links in the G_M and G_P graphs.

246 The link prediction problem (Liben-Nowell & Kleinberg, 2003) is defined as follows:

247 **Definition 4.** Let $G = \langle N, E \rangle$ be an undirected graph and let $G' = \langle N, E' \rangle$ be a subgraph of G which contains all
 248 nodes in G and a subset $E' \subseteq E$ of its edges. The link prediction problem consists of printing a list of non-edges in G'
 249 which are edges in G .

250 We will call the set E' as the *training set* and the set $E - E'$ as the *test set*.

251 In practice, algorithms to solve the link prediction problem build a matrix Ω in which the entry $\Omega_{ij} = \sigma_{ij}$ specifies
 252 the *degree of similarity* between the nodes i and j ; all pair of non edges $\langle i, j \rangle$ in G are thus ranked in decreasing order
 253 of similarity and non-edges with the largest similarity scores are the most likely to exist (Liben-Nowell & Kleinberg,
 254 2003).

255 We can define many similarity scores to compute the similarity degree of two nodes in G . In what follows we
 256 first illustrate some of these similarity metrics and, thus, we analyse their accuracy in predicting edges in G_M and G_P .
 257 Methods to compute node similarity can be classified into *local* and *global* methods

258 6.1. Local Methods to calculate node similarity

259 A first class of methods to calculate node similarity in graphs is known as *local methods* (Liben-Nowell & Klein-
 260 berg, 2003; Leicht et al., 2006). because they only require the knowledge of the neighbours of two nodes i and j .
 261 Some of the most popular local methods are as follows:

262 1. *Jaccard Coefficient* (JC) (Jaccard, 1912; Liben-Nowell & Kleinberg, 2003):

$$JC(i, j) = \frac{|N(i) \cap N(j)|}{|N(i) \cup N(j)|} \quad (3)$$

263 2. *Common Neighbors* (CN) (Newman, 2001; Liben-Nowell & Kleinberg, 2003):

$$CN(i, j) = |N(i) \cap N(j)| \quad (4)$$

264 3. *Preferential Attachment* (PA) (Newman, 2001; Liben-Nowell & Kleinberg, 2003):

$$PA(i, j) = d_i \times d_j \quad (5)$$

265 4. *Adamic-Adar coefficient* (AA) (Adamic & Adar, 2003; Liben-Nowell & Kleinberg, 2003):

$$AA(i, j) = \sum_{x \in N(i) \cap N(j)} \frac{1}{\log |N(x)|} \quad (6)$$

266 6.2. Global Methods to calculate node similarity

267 Observe that both the G_M and G_P graphs are highly sparse and, thus, we expect that the task of predicting edges
268 is hard if we would rely only on local information.

269 However, both G_M and G_P display a very high clustering coefficient (see Table 1), which is much higher than that
270 we observe in other type of real-life social networks of roughly equal size. A large clustering coefficient implies that
271 if two nodes i and j share at least one neighbour, then there is a high chance that i and j will be linked by an edge
272 too. Therefore, methods to calculate node similarity which leverage higher order structures (e.g., as walks or paths)
273 or, more in general, the *full knowledge* of the graph topology, might be more accurate than local methods in predicting
274 edges.

275 We will call these methods as *global methods* and one of the most popular global methods is the so called *Katz*
276 *score* (Katz, 1953).

277 The Katz score $\kappa(i, j)$ associated with a pair of nodes i and j considers the whole ensemble of walks connecting
278 i and j and it assumes that each walk provides a contribution to determine the degree of similarity between i and j .
279 A core assumption in the calculation of $\kappa(i, j)$ is that long walks are to be penalized with respect to short ones, which
280 implies that two nodes are highly similar if they are connected by many short walks in G . To formally encode such a
281 principle, we introduce a discount factor α and we denote $w_k(i, j)$ as the number of walks of length $k = 0, 1, \dots$, from
282 i to j . The Katz coefficient score is then computed as follows:

$$\kappa(i, j) = w_0(i, j) + \alpha w_1(i, j) + \alpha^2 w_2(i, j) + \dots + \alpha^k w_k(i, j) + \dots = \sum_{k=0}^{\infty} \alpha^k w_k(i, j) \quad (7)$$

283 Observe that $w_0(i, j) = 1$ if and only if nodes i and j coincide, 0 otherwise. If we let \mathbf{A} be the adjacency matrix of
284 G and suppose that α is less than $\frac{1}{\lambda_1}$, λ_1 being the largest eigenvalue of \mathbf{A} ¹, then the Katz score between any pair of
285 nodes in G can be seen as a matrix \mathcal{K}_α , which can be computed as follows:

$$\mathcal{K}_\alpha = (\mathbf{I} - \alpha \mathbf{A})^{-1} - \mathbf{I} \quad (8)$$

286 Here \mathbf{I} is the identity matrix.

287 In our analysis we consider also *Node2Vec* (Grover & Leskovec, 2016), a recent but promising approach for
288 embedding graphs onto vectors. More specifically, given a graph $G = \langle N, E \rangle$, Node2Vec seeks at finding out a
289 function $f : N \rightarrow \mathbb{R}^k$ where k is a fixed constant and \mathbb{R}^k is the set of k -th dimensional arrays of real numbers. The
290 main requirement we impose on f is that if two nodes i and j are “close” in G , then their representations $f(i)$ and

¹The parameter λ_1 is also known as the *spectral radius* of \mathbf{A} .

291 $f(j)$ should be close in \mathbb{R}^k too. To detect pairs of close nodes, Node2Vec simulates a random walk on G which can
 292 be thought as an interpolation of two popular procedures to explore a graph, namely the Breadth First Search (BFS)
 293 and the Depth First Search (DFS). More specifically, such a random walk is regulated by two parameters, namely
 294 the *return parameter* p (which specifies the likelihood the random walk will immediately revisiting a node) and the
 295 *in-out parameter* q : if $q > 1$, the random walk acts as a BFS because it tends to visit nodes which are close to the
 296 currently visited node; vice versa, if $q < 1$, the walk tends to move to nodes that are farther away from the current,
 297 thus simulating a DFS.

298 After applying the Node2Vec algorithm, each node i is associated with a vector \mathbf{v}_i and the similarity of two nodes i
 299 and j is defined as the cosine similarity of vectors \mathbf{v}_i and \mathbf{v}_j .

300 A further method to compute node similarity is the *Personalized PageRank similarity score* (PPR) (Avrachenkov
 301 et al., 2019), which, in matrix form, is defined as follow:

$$PPR_\alpha = (\mathbf{I} - \alpha\mathbf{P})^{-1} \quad (9)$$

302 The matrix \mathbf{P} is a row-stochastic matrix defined as $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$: here \mathbf{D} is a diagonal matrix storing the degrees of
 303 nodes in G and, \mathbf{A} is the adjacency matrix of G . Therefore, the sum of the elements within each row of \mathbf{P} is 1 and we
 304 can interpret \mathbf{P} as the transition probability of a random walk over G in which the random walker, at any step, chooses
 305 uniformly at random one of its neighbors.

306 6.3. Simulation setup

307 As a first step of our analysis, we compare local and global methods. Specifically, let $\sigma_{\text{met}}(i, j)$ be the similarity
 308 score between a pair of nodes i and j calculated by means the method *met*, where *met* is one of the methods previously
 309 introduced. Similarity scores generated by each method were normalized to range from 0 to 1. We claim that i and
 310 j are connected if and only if $\sigma_{\text{met}}(i, j)$ is bigger than a threshold θ and we negate the existence of that edge if
 311 $\sigma_{\text{met}}(i, j) < \theta$. In this way, we were able to map continuous similarity scores onto discrete labels (i.e., 0 and 1 to
 312 claim/negate the existence of an edge).

313 We used two metrics to evaluate the level of association between a particular measure of similarity and the ex-
 314 istence of an edge: (i) the *True Positive Rate (TPR)* and (ii) the *True Negative Rate (TNR)*. The TPR measures the
 315 proportion between the number of edges that a similarity measure claims exist and the real number of edges. The TNR
 316 is the proportion between the number of node pairs that according to a particular similarity measure are not connected
 317 and the actual number of pairs of nodes not connected. Space limitations preclude us from reporting the TPR and TNR
 318 of local methods for a broad range of values of θ . However, it is instructive to comment the configuration $\theta = 0.5$:
 319 here we observe that the TPR of all local methods was around 0 and their TNR were close to 1. Such a result implies
 320 that local methods almost always negate the existence of an edge and, thus, due to the sparsity of G_M and G_P , their
 321 guesses are almost always exact. Of course, local methods fail to identify edges actually existing. The Katz and PPR
 322 scores, instead, work much better than local methods. Because of space limitations, in Figure 2 we plot only the TPR
 323 and TNR for Katz score as function of α ; similar results hold true for PPR as function of α .

324 The main conclusions we drew from our analysis are as follows: (i) An increase of α yields a decrease in TPR. (ii)
 325 The TNR achieved by the Katz score in G_M and G_P is generally very large (bigger than 0.99) even if slightly smaller
 326 that achieved by local methods. Specifically, Figure 2 (Left Panel) indicates the presence of a turning point $\bar{\alpha}$ (with
 327 $\bar{\alpha} \simeq 0.52$ in case of the G_M graph and $\bar{\alpha} \simeq 0.44$ in case of the G_P graph) beyond which the TPR quickly drops. The Katz
 328 score thus perfectly addresses issues we highlighted above and, with a suitable choice of α , all highly-scored pair of
 329 nodes are actually tied by an edge. Such a result agrees fairly well with our model about information flow in criminal
 330 networks: criminal often do not communicate directly each other but they prefer to make use of intermediaries to
 331 convey messages, both in face-to-face meetings and in case of phone calls; however, the chain of intermediaries is
 332 generally very short for security reasons.

333 6.4. Accuracy of Link Prediction Methods in Montagna

334 As a further step of our analysis, we analyse the accuracy of the methods to calculate node similarity.

335 We applied 10-fold cross validation to quantify the predictive accuracy of each previous predictor. Cross-validation
 336 is a procedure used to assess the accuracy of a Machine Learning algorithm which, in the latest years, gained an
 337 astonishing popularity (Hastie et al., 2009). The main reasons explaining the popularity of k -fold cross validation

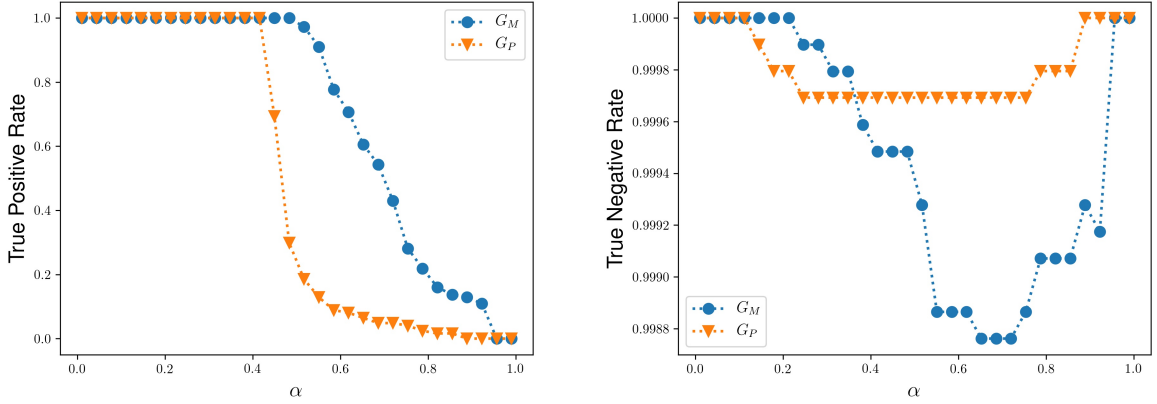


Figure 2: **Left Panel.** True Positive Rate (TPR) associated with the Katz score as function of α on G_M and G_P graphs. **Right Panel.** True Negative Rate (TNR) associated with the Katz score as function of α on G_M and G_P graphs.

338 are its simplicity as well as its ability of producing less optimistic accuracy assessment than methods based on the
 339 (random) division of a dataset into a training and a test part. In short, in the k -fold cross validation we randomly
 340 shuffle a dataset \mathcal{D} and divide it into k groups, say, G_1, \dots, G_k ; common choices for k are $k = 5$ and $k = 10$. For each
 341 group G_i , we take G_i as test dataset and we use the remaining G_1, G_2, \dots, G_k groups as training set: in other words
 342 we use all groups G_j (but G_i) to fit our model; once our model has been fitted, we evaluate its accuracy on G_i . Such a
 343 procedure is repeated for each group G_i and, consequently, any sample in the original dataset is used $k - 1$ times for
 344 training purposes and one time for testing purposes. At the end of evaluation procedure we obtain k values of accuracy
 345 (one for each group used as test set); we thus take the average of the accuracy scores on each group G_i as the accuracy
 346 of the algorithm to evaluate.

347 The prediction accuracy is evaluated by a standard metric, the Area Under the Receiving Operating Curve (AU-
 348 ROC) ². We repeated the calculation of AUROC $n = 50$ times, thus generating a sample of the true AUROC scores.
 349 We then calculated the empirical mean m and the empirical standard deviation s of the sample above; if we denote as
 350 μ the true AUROC, it is well-known that the random variable $t = \frac{\sqrt{n}(m-\mu)}{s}$ follows a t -student distribution with $n - 1$
 351 degrees of freedom (Ross, 2017). We then calculated the value of A for which $P(-A \leq t \leq A) = 0.95$ and we take the
 352 interval $\left(m - A \frac{s}{\sqrt{n}}, m + A \frac{s}{\sqrt{n}}\right)$ as the 95% confidence interval associated with the true AUROC score.

353 In Table 2 we report the confidence intervals (CI) associated with AUROC for AA, CN, PA, JC, Node2Vec, Katz
 354 and PPR methods on G_M and G_P graphs. We report the AUROC 95% confidence intervals for the Katz and PPR
 355 methods. Moreover, we considered some specified values of the parameter α (namely $\alpha = 0.1, 0.3, 0.5, 0.7$ and 0.9)
 356 and we investigated how the α parameter affected the AUROC.

357 In case of G_M graph, the AUROC is generally very high for all methods under investigation and the worst-
 358 performing method is PA. The Katz score and the PPR score generally outperform but their AUROC tends to slightly
 359 decrease as α increases: for instance, if $\alpha > 0.7$ the AUROC achieved by Katz score ranges from 0.893 to 0.918 while
 360 the AUROC measured for PPR ranges between 0.922 and 0.94. The JC, CN and AA methods achieve an AUROC
 361 which is slightly smaller than that of the Katz and the PPR score. In contrast, PA displays the worst performance and
 362 its AUROC is 17.51% less than that of AA and 17.25% smaller than that of JC.

363 On the G_P network, instead, the PA method achieves the highest AUROC and the performances of all other
 364 methods significantly deteriorate. For instance, the AA method achieves an AUROC ranging from 0.602 to 0.643,
 365 with a loss of more than 30% with respect to the G_M graph.

366 We are therefore able to answer research questions \mathbf{Q}_1 and \mathbf{Q}_2 . As for \mathbf{Q}_1 , we observe a few methods introduced
 367 in this paper are very and they achieve an AUROC, which, in some cases, is higher than 0.9.

²The AUROC is understood as the probability that a randomly chosen edge in the test set gets a higher score than a randomly chosen non-edge.

| Method | | CI | |
|---------------|----------|----------------|----------------|
| Name | α | G_M | G_P |
| JC | | (0.917, 0.938) | (0.57, 0.623) |
| CN | | (0.938, 0.953) | (0.595, 0.634) |
| PA | | (0.754, 0.789) | (0.872, 0.913) |
| AA | | (0.939, 0.957) | (0.602, 0.643) |
| Node2Vec | | (0.899, 0.919) | (0.577, 0.646) |
| \mathcal{K} | 0.1 | (0.946, 0.959) | (0.706, 0.753) |
| | 0.3 | (0.95, 0.965) | (0.711, 0.772) |
| | 0.5 | (0.939, 0.955) | (0.701, 0.754) |
| | 0.7 | (0.927, 0.946) | (0.73, 0.778) |
| | 0.9 | (0.927, 0.946) | (0.696, 0.748) |
| PPR | 0.1 | (0.939, 0.956) | (0.66, 0.724) |
| | 0.3 | (0.954, 0.968) | (0.696, 0.752) |
| | 0.5 | (0.942, 0.958) | (0.698, 0.747) |
| | 0.7 | (0.939, 0.955) | (0.687, 0.743) |
| | 0.9 | (0.922, 0.94) | (0.688, 0.75) |

Table 2: AUROC Confidence Intervals for the AA, CN, PA, JC, Node2Vec, Katz Coefficient and PPR methods computed on the G_M and G_P graphs.

As for \mathbf{Q}_2 , we report that graph topology actually plays an important role on the process of predicting edges: specifically, methods which are very accurate on G_M performs badly on G_P (and vice versa). In detail if a graph is poorly connected (with a low edge density and a small clustering coefficient) local methods are to be preferred to global ones. Vice versa, global methods as the Katz score achieve their best accuracy on graphs which display a better level of connectivity (i.e. in graph with larger edge density and clustering coefficients).

7. Robustness of the Link Prediction algorithms

In this section we aim at answering our research question \mathbf{Q}_3 . We recall that both G_M and G_P are built upon the evidence collected by police forces and, therefore, they are an incomplete sample of true graphs G'_M and G'_P . An important discrepancy between G_M and G'_M (resp., G_P and G'_P) might significantly alter the conclusions we can draw from the analysis of G_M (resp. G_P) and, in particular, it might severely alter our ability of predicting edges between criminals.

We run our analysis in parallel for the two networks G_M and G_P and rely on the methods achieving the highest prediction accuracy in the analysis of the previous section: for G_M , we concentrate on the Katz score with different levels of parameter α ; for G_P , we focus on PA.³ Let us consider the G_M graph and our aim is to quantify the difference between $\mathcal{K}_\alpha(G_M)$ and $\mathcal{K}_\alpha(G'_M)$. At an aggregate level, we introduce the parameter $\rho_\alpha(G_M, G'_M)$ to quantify such a difference:

$$\rho_\alpha(G_M, G'_M) = \frac{\|\mathcal{K}_\alpha(G_M) - \mathcal{K}_\alpha(G'_M)\|_2}{\|\mathcal{K}_\alpha(G_M)\|_2} \quad (10)$$

Equation 10 can be applied to graphs G_P and G'_P and method PA, which yielded the highest prediction accuracy in the telephone call network. However, the equation is unapplicable in practice because we do not know the true graphs G'_M and G'_P . We can overcome this issue by assuming that missing edges – i.e., those edges in G'_M (resp., G'_P) but not observed in G_M (resp., G_P) – have been generated using a suitable *probabilistic model*.

Our probabilistic model assumes that non-observed edges in G'_M (resp. G'_P) are non-edges in G_M (resp., G_P); each non-edge in G_M (resp., G_P) is associated with a parameter ℓ , called *likelihood*, such that the higher the likelihood, the

³For G_M , we have also run our analysis in the case of PPR score with similar results. Due to space limitations we report only results in case of the Katz score.

| p | CN | JC | Random |
|------|-------|-------|--------|
| 1.0 | 0.134 | 0.107 | 0.121 |
| 5.0 | 0.157 | 0.174 | 0.179 |
| 10.0 | 0.232 | 0.235 | 0.201 |
| 15.0 | 0.291 | 0.296 | 0.233 |

Table 3: ρ as function of p , percentage of added edges, in the G_p graph.

390 more likely a non-edge in G_M (resp., G_P) will correspond to an edge in G'_M (resp., G'_P). If the likelihood ℓ is specified,
391 we can select non-edges from G_M (resp., G_P) on the basis of their likelihood and we can incrementally insert them
392 into G_M (resp., G_P) until a pre-defined stop condition is satisfied. At the end of this procedure we obtain G'_M (resp.,
393 G'_P).

394 We considered multiple strategies to model the likelihood ℓ , namely: (i) Common Neighbors (see Equation 4),
395 (ii) Jaccard's Coefficient (see Equation 3), and (iii) a Random model, a baseline where ℓ is distributed as a uniform
396 random variable in the interval $[0, 1]$.

397 The model above resemble *network-growth models* (Newman, 2010) which describe the creation/evolution of a
398 network (for example, a mechanism similar to the preferential attachment is at the base of the generation of Barabasi-
399 Albert networks). However, in network-growth models we assume that new nodes arrive and join the network and the
400 last node can decide which other nodes to connect to. In contrast, in our model, there are no new nodes that can be
401 added to the network: this is equivalent to the simplifying hypothesis that the network is *perfectly observable* about
402 what concerns the subjects in it (that is, the investigation has not excluded any criminal subject) and the possible lack
403 of information only concerns the relations observed by the investigators.

404 Our experimental protocol consists of the following steps:

405
406 **Step 1** Let T be the list of non-edges in G_M (resp., G_P) sorted by decreasing likelihood scores. We took 50% of
407 the top elements in T , i.e., we chose half of the non-edges that have the highest likelihood. This step is
408 necessary to create a group of potential non-edges that is sufficiently large but, at the same time, which is
409 reliable enough because non-edges with low values of likelihood are filtered out. We call $C \subseteq T$ the set of the
410 non-edges generated at the end of Step 1.

411 **Step 2** We randomly choose a sample of $R(p) \subseteq C$ with size equal to p from C . In our experiments, we set $p =$
412 $\{1\%, 5\%, 10\%, 15\%\}$. Of course, the larger p , the higher the number of missing edges.

413 **Step 3** We add elements in $R(p)$ to G_M (resp., G_P), thus creating a new graph $G'_M(p) = \langle N, E \cup R(p) \rangle$ (resp., $G'_P(p) =$
414 $\langle N, E \cup R(p) \rangle$).

415 **Step 4** We calculate the relative variation ρ using Equation 10, where the graph G'_M (resp., G'_P) is replaced by $G'_M(p)$
416 (resp, $G'_P(p)$).

417
418 Steps 2-4 have been repeated 30 times to avoid statistical fluctuations. The results are shown in Figure 3 for G_M ,
419 and in Table 3 and in Figure 4 for G_P .

420 As for the meeting network, the Random strategy clearly induces the highest values of ρ for any value of α and p .
421 This is a largely predictable result: if the probability of the existence of a non-edge follows one of the other strategies
422 (i.e., JC and CN), then the network structure is somehow able to predict the existence of missing edges. On the other
423 hand, if edges were randomly placed, the network structure would not offer any insight to predict the existence of
424 missing edges and, thus, the parameter ρ significantly grows: for instance, it suffices to set $p = 5\%$ and $\alpha = 0.2$ to
425 obtain $\rho \simeq 1$.

426 The growth of α implies a growth of ρ for the Random strategy. For CN, ρ is relatively stable, only slightly
427 increasing for higher values of α . A limit case happens when we decide to adopt JC as the likelihood function: in this
428 case, the result is anti-intuitive because when α increases, a reduction of ρ occurs (with peaks up to 18%). In practice,
429 if $\alpha \rightarrow 1$, the contribution of relatively long walks is not-negligible, and, thus, long walks are capable of contrasting

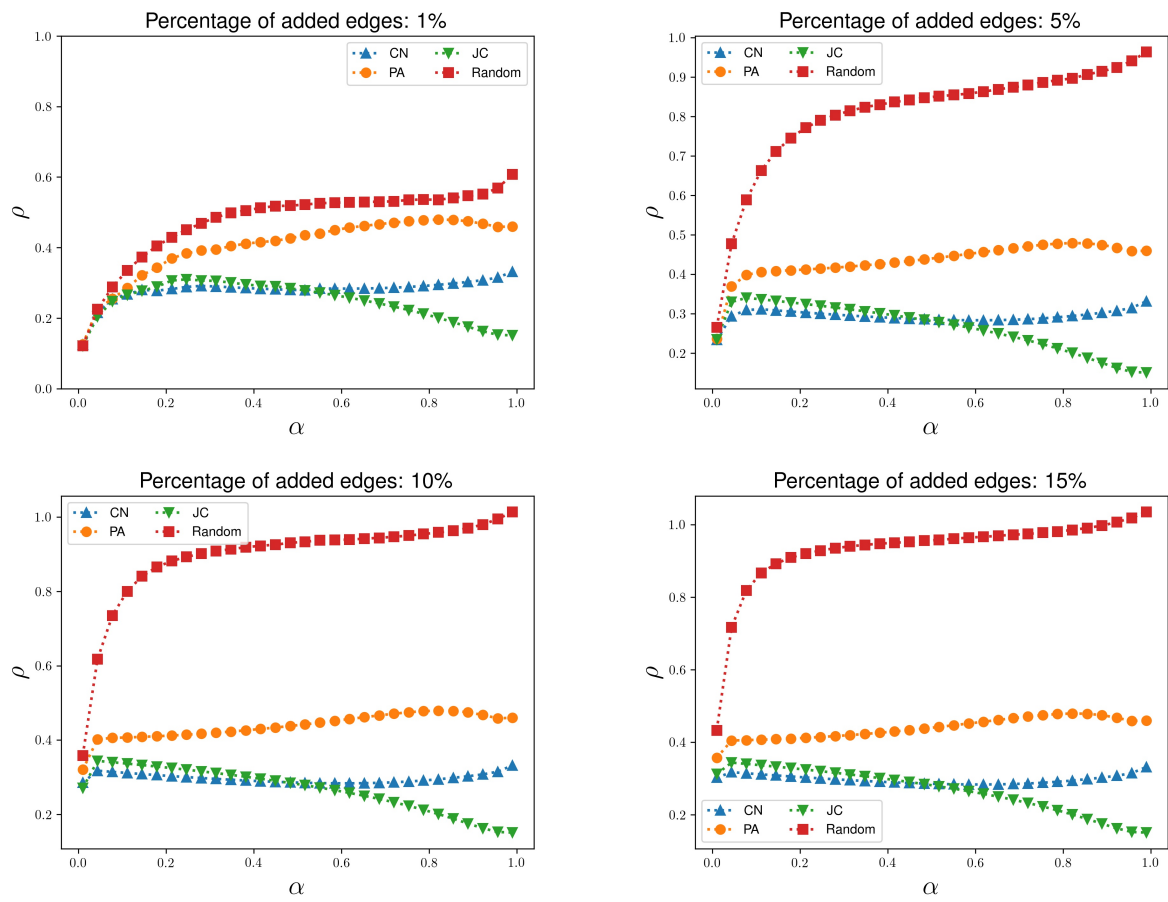


Figure 3: Variation of ρ as function of α in the G_M graph for $p = 1\%$ (Top Left), 5% (Top Right), 10% (Bottom Left), 15% (Bottom Right).

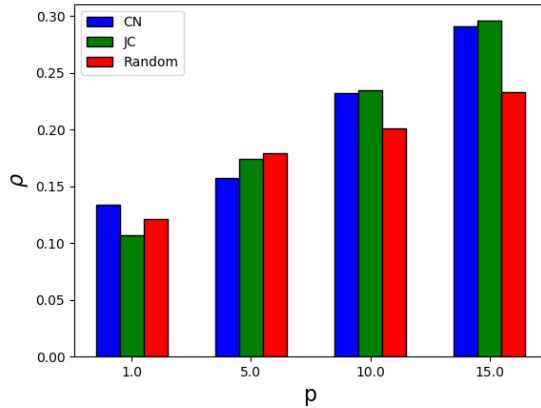


Figure 4: ρ as function of p , percentage of added edges, in the G_p graph.

430 high level of uncertainty associated with larger values of p . For a fixed α , the parameter p plays a key role on the
 431 value of ρ and, obviously, the higher p , the higher ρ .

432 We obtained totally different results in the G_p graph. The Random and JC likelihood functions are the only
 433 strategies that generate the highest value of ρ and there is a crossover point $\bar{\alpha}$ in the JC likelihood beyond which we
 434 observe a variation of ρ greater than that detected in the random generative model. The variations of ρ in the CN
 435 generative model become almost imperceptible if α gets larger than 0.1, and, thus, CN seems an unhappy choice to
 436 analyse telephone conversation flows. The trend of ρ is relatively little affected by p if $p < 15\%$; however, if $p > 15\%$
 437 the value of ρ in all the generative models analysed undergoes significant changes.

438 We are now in the position of answer **Q₃**: link prediction algorithms are sensitive not only to uncertainty in G_M
 439 and G_p (captured by the parameter p) but also on the type of graph they operate on. If the amount of uncertainty
 440 is relatively small and non-observed edges derive from a specified generative model, we can hope for robust edge
 441 prediction. This is confirmed by the results in G_M , where both CN and JC generate lower values of ρ than the Random
 442 strategy for different combination of p and α . Conversely, there is no clear indication from G_p , suggesting that the
 443 network growth does not follow a specific strategy. We can thus conclude that the robustness of link prediction is not
 444 only dependent on the amount of uncertainty (i.e. p) but also on type of network and underlying relations. In the light
 445 of our studies, we recommend law enforcement agencies not only to build a detailed map of connections between
 446 criminals but also to investigate how such a map evolves over time: in this way we would be able to determine which
 447 of the likelihood functions described in this section better fit experimental observations and, if required, we could
 448 design more sophisticated likelihood functions to help law enforcement agencies to detect and prevent crimes.

449 8. Conclusions and Future Works

450 We presented a study of two criminal networks extracted from the outcome of an anti-mafia law enforcement oper-
 451 ation called “Montagna” against individuals charged for participating in a mafia association. This study is interesting
 452 *per se* as the pre-trial detention order from which the networks under study have been extracted, namely the network
 453 of meetings and phone calls, concerns the birth and growth of a branch of Sicilian Cosa Nostra in the North-Eastern
 454 part of Sicily, a territory historically under the control of the Palermo and Catania families.

455 We first applied some of the most widely used similarity criteria to both networks in order to perform link predic-
 456 tion. The most accurate results were obtained by applying the Katz score, and our experimental finding confirm that
 457 the predictions heavily rely upon short-range interactions. This is consistent with the structure of Mafia families and
 458 the average clustering coefficient of the two networks under study.

459 Next, we investigated on the robustness of link prediction algorithms in presence of network uncertainties. To this
 460 end, we carried out an experiment in which the observed networks were regarded as the starting point of a growing

461 mechanism during which some missing edges were added. The impact of non-observed edges was measured in terms
462 of the difference of the Katz scores.

463 Our experiment shows that the Meeting Graph G_M slightly differs from graphs in which up to 15% of new edges
464 were added. On the contrary, the Phone Call Graph G_P exhibits strong differences between the observed and the new
465 graphs. This may induce to ask whether the information relative to this network may actually be incomplete and may
466 need some more edges which were neglected during the investigations.

467 Even more interestingly, the experiments we carried out clearly show that metrics of accuracy, the most widely
468 used measures able to assess the quality of a link prediction method, should be integrated with a new measure, the
469 stability, which takes into account the extent with which the insertion of unobserved edges modifies the network.

470 As for future work, we plan to apply network embedding methods to better analyse criminal networks. For our
471 purposes, we need to resort to approaches capable of handling heterogeneous networks (Li & Tang, 2019), i.e. graphs
472 in which nodes and edges carry specific information (e.g., nodes may be labelled with the role of an individual in a
473 criminal group while edges specify the type of interaction between two criminals). We also plan to consider temporal
474 information in the link prediction task, as described by Soares & Prudencio (2013).

475 References

- 476 Adamic, L. A., & Adar, E. (2003). Friends and neighbors on the web. *Social Networks*, 25(3), 211–230. [https://doi.org/10.1016/S0378-8733\(03\)00009-1](https://doi.org/10.1016/S0378-8733(03)00009-1).
- 477 Agreste, S., Catanese, S., De Meo, P., Ferrara, E., & Fiumara, G. (2016). Network structure and resilience of mafia syndicates. *Information Sciences*, 351, 30–47. <https://doi.org/10.1016/j.ins.2016.02.027>.
- 478 Avrachenkov, K., Chebotarev, P., & Rubanov, D. (2019). Similarities on graphs: Kernels versus proximity measures. *European Journal of Combinatorics*, 80, 47–56. <https://doi.org/10.1016/j.ejc.2018.02.002>.
- 479 Bahulkar, A., Baycik, N. O., Sharkey, T., Shen, Y., Szymanski, B., & Wallace, W. (2018). Integrative analytics for detecting and disrupting transnational interdependent criminal smuggling, money, and money-laundering networks. In *2018 IEEE International Symposium on Technologies for Homeland Security (HST)* (pp. 1–6). <https://doi.org/10.1109/THS.2018.8574121>.
- 480 Berlusconi, G. (2013). Do all the pieces matter? Assessing the reliability of law enforcement data sources for the network analysis of wire taps. *Global Crime*, 14(1), 61–81. <https://doi.org/10.1080/17440572.2012.746940>.
- 481 Berlusconi, G., Calderoni, F., Parolini, N., Verani, M., & Piccardi, C. (2016). Link prediction in criminal networks: A tool for criminal intelligence analysis. *PLOS ONE*, 11(4), 1–21. <https://doi.org/10.1371/journal.pone.0154244>.
- 482 Bouchard, M. (2020). Collaboration and boundaries in organized crime: A network perspective. *Crime and Justice*, 49. <https://doi.org/10.1086/708435>.
- 483 Calderoni, F. (2010). Inside criminal networks. *European Journal on Criminal Policy and Research*, 16(1), 69–70. <https://doi.org/10.1007/s10610-010-9118-7>.
- 484 Calderoni, F. (2012). The structure of drug trafficking mafias: the ‘ndrangheta and cocaine. *Crime, Law and Social Change*, 58(3), 321–349. <https://doi.org/10.1007/s10611-012-9387-9>.
- 485 Calderoni, F., Piccardi, C., & Brunetto, D. (2017). Communities in criminal networks: A case study. *Social Networks*, 48, 116–125. <https://doi.org/10.1016/j.socnet.2016.08.003>.
- 486 Calderoni, F., & Superchi, E. (2019). The nature of organized crime leadership: criminal leaders in meeting and wiretap networks. *Crime, Law and Social Change*, (pp. 1–26). <https://doi.org/10.1007/s10611-019-09829-6>.
- 487 Campana, P., & Varese, F. (2012). Listening to the wire: criteria and techniques for the quantitative analysis of phone intercepts. *Trends in Organized Crime*, 15(1), 13–30. <https://doi.org/10.1007/s12117-011-9131-3>.
- 488 Catanese, S., De Meo, P., Ferrara, E., & Fiumara, G. (2014). Detecting criminal organizations in mobile phone networks. *Expert Systems with Applications*, 41(13), 5733–5750. <https://doi.org/10.1016/j.eswa.2014.03.024>.
- 489 Chen, H., Chung, W., Xu, J., Wang, G., Qin, Y., & Chau, M. (2004). Crime data mining: A general framework and some examples. *IEEE Computer*, 37, 50–56. <https://doi.org/10.1109/MC.2004.1297301>.
- 490 Chen, H., Li, X., & Huang, Z. (2005). Link prediction approach to collaborative filtering. In *Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries (JCDL '05)* (pp. 141–142). Denver, Colorado, USA: IEEE. <https://doi.org/10.1145/1065385.1065415>.
- 491 De Moor, S., Vandeviver, C., & Vander Beken, T. (2020). Assessing the missing data problem in criminal network analysis using forensic DNA data. *Social Networks*, 61, 99–106. <https://doi.org/10.1016/j.socnet.2019.09.003>.
- 492 Diviák, T. (2019). Key aspects of covert networks data collection: Problems, challenges, and opportunities. *Social Networks*. <https://doi.org/10.1016/j.socnet.2019.10.002>.
- 493 Duxbury, S. W., & Haynie, D. L. (2019). Criminal network security: An agent-based approach to evaluating network resilience*. *Criminology*, 57(2), 314–342. <https://doi.org/10.1111/1745-9125.12203>.
- 494 Everrett, M., & Borgatti, S. (2013). The dual-projection approach for two-mode networks. *Social Networks*, 35(2), 204–210. <https://doi.org/10.1016/j.socnet.2012.05.004>.
- 495 Fan, C., Liu, Z., Lu, X., Xiu, B., & Chen, Q. (2017). An efficient link prediction index for complex military organization. *Physica A: Statistical Mechanics and its Applications*, 469, 572–587. <https://doi.org/10.1016/j.physa.2016.11.097>.
- 496 Fawcett, T. (2006). An introduction to ROC analysis. *Pattern recognition letters*, 27(8), 861–874. <https://doi.org/10.1016/j.patrec.2005.10.010>.
- 497 Ferrara, E., De Meo, P., Catanese, S., & Fiumara, G. (2014). Visualizing criminal networks reconstructed from mobile phone records. *CEUR Workshop Proceedings*, 1210.

- 520 Ficara, A., Cavallaro, L., De Meo, P., Fiumara, G., Catanese, S., Bagdasar, O., & Liotta, A. (2020). Social network analysis of sicilian mafia
521 interconnections. In H. Cherifi, S. Gaito, J. F. Mendes, E. Moro, & L. M. Rocha (Eds.), *Complex Networks and Their Applications VIII* (pp.
522 440–450). Cham: Springer International Publishing. https://doi.org/10.1007/978-3-030-36683-4_36.
- 523 Freeman, L. C. (1979). Centrality in social networks' conceptual clarification. *Social Networks*, 1(3), 215–239.
- 524 Gambetta, D. (1993). *The sicilian mafia*. Harvard University Press Cambridge, MA.
- 525 Gambetta, D., & Reuter, P. (1995). Conspiracy among the many: the mafia in legitimate industries. In *The Economic Dimensions of Crime* (pp.
526 99–120). Springer. https://doi.org/10.1007/978-1-349-62853-7_5.
- 527 Grassi, R., Calderoni, F., Bianchi, M., & Torriero, A. (2019). Betweenness to assess leaders in criminal networks: New evidence using the dual
528 projection approach. *Social Networks*, 56, 23–32. <https://doi.org/10.1016/j.socnet.2018.08.001>.
- 529 Grover, A., & Leskovec, J. (2016). node2vec: Scalable feature learning for networks. In *Proc. of the 22nd ACM SIGKDD International Conference*
530 *on Knowledge Discovery and Data Mining* (pp. 855–864). San Francisco, CA, USA: ACM. <https://doi.org/10.1145/2939672.2939754>.
- 531 Hasan, M., Chaoji, V., Salem, S., & Zaki, M. (2006). Link prediction using supervised learning. In *Proc. of the SDM 2006 Workshop on link*
532 *analysis, counter-terrorism and security*.
- 533 Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer Science &
534 Business Media.
- 535 Hric, D., Peixoto, T., & Fortunato, S. (2016). Network structure, metadata, and the prediction of missing nodes and annotations. *Physical Review*
536 *X*, 6(3), 031038. <https://doi.org/10.1103/PhysRevX.6.031038>.
- 537 Jaccard, P. (1912). The distribution of the flora in the alpine zone. *New Phytologist*, 11(2), 37–50. <https://doi.org/10.1111/j.1469-8137.1912.tb05611.x>.
- 538 Katz, L. (1953). A new status index derived from sociometric analysis. *Psychometrika*, 18(1), 39–43. <https://doi.org/10.1007/BF02289026>.
- 539 Kim, M., & Leskovec, J. (2011). The network completion problem: Inferring missing nodes and edges in networks. In *Proc. of the SIAM*
540 *International Conference on Data Mining (SDM 2011)* (pp. 47–58). Mesa, Arizona, USA: SIAM. <https://doi.org/10.1137/1.9781611972818.5>.
- 541 Kleemans, E., & Bunt, H. (1999). The social embeddedness of organized crime. *Transnational Organized Crime*, 5, 19–36.
- 542 Kleemans, E., & De Poot, C. (2008). Criminal careers in organized crime and social opportunity structure. *European Journal of Criminology*, 5,
543 69–98. <https://doi.org/10.1177/1477370807084225>.
- 544 Klerks, P. (2001). The network paradigm applied to criminal organisations: Theoretical nitpicking or a relevant doctrine for investigators? recent
545 developments in the Netherlands. *Connections*, 24(3), 53–65.
- 546 Krebs, V. (2002). Mapping networks of terrorist cells. *Connections*, 24(3), 43–52.
- 547 Leicht, E., Holme, P., & Newman, M. (2006). Vertex similarity in networks. *Physical review. E, Statistical, nonlinear, and soft matter physics*,
548 73(2), 026120. <https://doi.org/10.1103/PhysRevE.73.026120>.
- 549 Li, C., & Tang, Y. (2019). Efficient heterogeneous proximity preserving network embedding model. *Expert Systems with Applications*, 134,
550 201–208. <https://doi.org/10.1016/j.eswa.2019.05.044>.
- 551 Liben-Nowell, D., & Kleinberg, J. (2003). The link prediction problem for social networks. In *Proceedings of the Twelfth International Con-*
552 *ference on Information and Knowledge Management CIKM '03* (p. 556–559). New York, NY, USA: Association for Computing Machinery.
553 <https://doi.org/10.1145/956863.956972>.
- 554 Lim, M., Abdullah, A., Zaman, N., & Supramaniam, M. (2019). Hidden link prediction in criminal networks using the deep reinforcement learning
555 technique. *Computers*, 8, 8. <https://doi.org/10.3390/computers8010008>.
- 556 Marciani, G., Porretta, M., Nardelli, M., & Italiano, G. F. (2017). A data streaming approach to link mining in criminal networks. In *2017 5th Inter-*
557 *national Conference on Future Internet of Things and Cloud Workshops (FiCloudW)* (pp. 138–143). <https://doi.org/10.1109/FiCloudW.2017.88>.
- 558 Menon, A. K., & Elkan, C. (2011). Link prediction via matrix factorization. In D. Gunopulos, T. Hofmann, D. Malerba, & M. Vazirgiannis (Eds.),
559 *Machine Learning and Knowledge Discovery in Databases* (pp. 437–452). Berlin, Heidelberg: Springer Berlin Heidelberg.
- 560 Morselli, C. (2003). Career opportunities and network-based privileges in the Cosa Nostra. *Crime, Law and Social Change*, 39(4), 383–418.
561 <https://doi.org/10.1023/A:1024020609694>.
- 562 Newman, M. (2001). Clustering and preferential attachment in growing networks. *Phys. Rev. E*, 64(2), 025102.
563 <https://doi.org/10.1103/PhysRevE.64.025102>.
- 564 Newman, M. (2010). *Networks: an Introduction*. Oxford University Press.
- 565 Pandey, B., Bhanodia, P. K., Khamparia, A., & Pandey, D. (2019). A comprehensive survey of edge prediction in social networks: Techniques,
566 parameters and challenges. *Expert Systems with Applications*, 124, 164–181. <https://doi.org/10.1016/j.eswa.2019.01.040>.
- 567 Paoli, L. (2004). Italian organised crime: Mafia associations and criminal enterprises. *Global Crime Today: The Changing Face of Organised*
568 *Crime*, 6(1), 19–32. <https://doi.org/10.1080/1744057042000297954>.
- 569 Paoli, L. (2008). *Mafia brotherhoods: Organized crime, Italian style*. Oxford University Press.
- 570 Parisi, F., Caldarelli, G., & Squartini, T. (2018). Entropy-based approach to missing-links prediction. *Applied Network Science*, 3(1), 17.
571 <https://doi.org/10.1007/s41109-018-0073-4>.
- 572 Rhodes, C., & Jones, P. (2009). Inferring missing links in partially observed social networks. *JORS*, 60, 1373–1383.
573 <https://doi.org/10.1057/jors.2008.110>.
- 574 Ross, S. (2017). *Introductory statistics*. Academic Press.
- 575 Soares, P., & Prudencio, R. (2013). Proximity measures for link prediction based on temporal events. *Expert Systems with Applications*, 40(16),
576 6652–6660. <https://doi.org/10.1016/j.eswa.2013.06.016>.
- 577 Sparrow, M. (1991). The application of network analysis to criminal intelligence: An assessment of the prospects. *Social Networks*, 13(3), 251–274.
578 [https://doi.org/10.1016/0378-8733\(91\)90008-H](https://doi.org/10.1016/0378-8733(91)90008-H).
- 579 Strang, S. (2014). Network analysis in criminal intelligence. In *Networks and network analysis for defence and security* (pp. 1–26). Springer.
580 https://doi.org/10.1007/978-3-319-04147-6_1.
- 581 Van der Hulst, R. (2009). Introduction to social network analysis (SNA) as an investigative tool. *Trends in Organized Crime*, 12(2), 101–121.
582 <https://doi.org/10.1007/s12117-008-9057-6>.
- 583 Watts, D., & Strogatz, S. (1998). Collective dynamics of 'small-world' networks. *Nature*, 393(6684), 440–442. <https://doi.org/10.1038/30918>.
- 584

585 Xu, J., & Chen, H.-c. (2005). Criminal network analysis and visualization. *Commun. ACM*, 48, 100–107. <https://doi.org/10.1145/1064830.1064834>.