

A peer-reviewed version of this preprint was published in PeerJ on 14 October 2015.

[View the peer-reviewed version](https://peerj.com/articles/cs-27) (peerj.com/articles/cs-27), which is the preferred citable publication unless you specifically need to cite this preprint.

Palleschi V, Pagani L, Pagnotta S, Amato G, Tofanelli S. 2015.
Application of Graph Theory to the elaboration of personal genomic data for genealogical research. PeerJ Computer Science 1:e27
<https://doi.org/10.7717/peerj-cs.27>

Application of Graph Theory to the elaboration of personal genomic data for genealogical research

Vincenzo Palleschi, Luca Pagani, Stefano Pagnotta, Giuseppe Amato, Sergio Tofanelli

In this communication a representation of the links between DNA-relatives based on Graph Theory is applied to the analysis of personal genomic data to obtain genealogical information. The method is tested on real data and discussed its applicability to the field of genealogical research. We envisage the proposed approach as a valid tool for a streamlined application to the publicly available data generated by many online personal genomic companies. By this way, anonymized matrices of pairwise genome sharing counts will enable to improve the retrieval of genetic relationship between customers who provided explicit consent to the treatment of their data .

1 Application of Graph Theory to the Elaboration of Personal Genomic
2 Data
3 for Genealogical Research

4 Vincenzo Palleschi^{1,2}, Luca Pagani^{3,4*}, Stefano Pagnotta¹, Giuseppe Amato⁵, Sergio Tofanelli⁶

5 ¹Institute of Chemistry of Organometallic Compounds
6 Research Area of National Research Council
7 Via G. Moruzzi, 1 – 56124 Pisa (ITALY)

8
9 ²Department of Civilizations and Forms of Knowledge
10 University of Pisa
11 Via G. Galvani, 1 – 56126 Pisa (ITALY)

12
13 ³Division of Biological Anthropology, University of Cambridge, UK

14
15 ⁴Department of Biological, Geological and Environmental Sciences,
16 University of Bologna
17 Via Selmi 3, 40126 Bologna (ITALY)

18
19 ⁵Institute of Sciences and Technology of Information
20 Research Area of National Research Council
21 Via G. Moruzzi, 1 – 56124 Pisa (ITALY)

22
23 ⁶Department of Biology
24 University of Pisa
25 Via L. Ghini, 13 – 56126 Pisa (ITALY)

26
27 **ABSTRACT**

28 *In this communication a representation of the links between DNA-relatives based on Graph*
29 *Theory is applied to the analysis of personal genomic data to obtain genealogical information.*
30 *The method is tested on real data and discussed its applicability to the field of genealogical*
31 *research. We envisage the proposed approach as a valid tool for a streamlined application to*
32 *the publicly available data generated by many online personal genomic companies. By this way,*
33 *anonymized matrices of pairwise genome sharing counts will enable to improve the retrieval of*
34 *genetic relationship between customers who provided explicit consent to the treatment of their*
35 *data.*

36 **Keywords:** DNA Analysis, Personal Genomics, Genealogy, Genetic Genealogy, Statistical
37 methods, Graph Theory, Ancestry reconstruction.

39 *Corresponding Author: Luca Pagani – lp.lucapagani@gmail.com

40 **1. Introduction**

41 In recent years, a number of companies started offering commercial services based on DNA
42 analysis for genealogical research^[1-3]. The informatic tools available to interpret such results,
43 usually provided by the same companies or by external services^[4], are mainly focused on
44 general population studies (Paternal and Maternal lineages based on Y chromosome and
45 mitochondrial haplogroups, Ancestry Composition/Admixture, etc.). On the other hand, very
46 few tools are provided to investigate the links of one's DNA profile with the relatives made
47 recognizable through personal genomic data. Notably, these pre-compiled tools are often the
48 only way to access the data provided by the DNA testing companies for a panel of hundreds or
49 thousands of individuals. Therefore, the starting point of any downstream analysis based on
50 this kind of data can only rely on the semi-processed input provided by the aforementioned
51 tools. The introduction by the genetic service providers of a wrapped application tool would
52 facilitate users' interpretations and unearth hidden genealogical information. Such tool should
53 enable to implement the mass of data each single DNA test makes available in an easy-to-grasp
54 graphical form. This would be particularly useful to detect the provenience of distant autosomic
55 DNA-relatives from either the paternal or the maternal lineage. In fact this task is often made
56 difficult by the links that might exist between the two parental genealogies due to the custom
57 in closed communities to marry between relatives, especially in the past.

58 Here we describe and annotate an artificial intelligence tool that helps exploiting the
59 information provided to customers by genealogical genetic services. The original approach of
60 this work is the use of cross-information about the links between the living DNA-relatives of the
61 test user (TU) for obtaining hints about the possible connections with other individuals, in the
62 absence of a-priori genetic or genealogic evidence.

63

64 **2. Data**

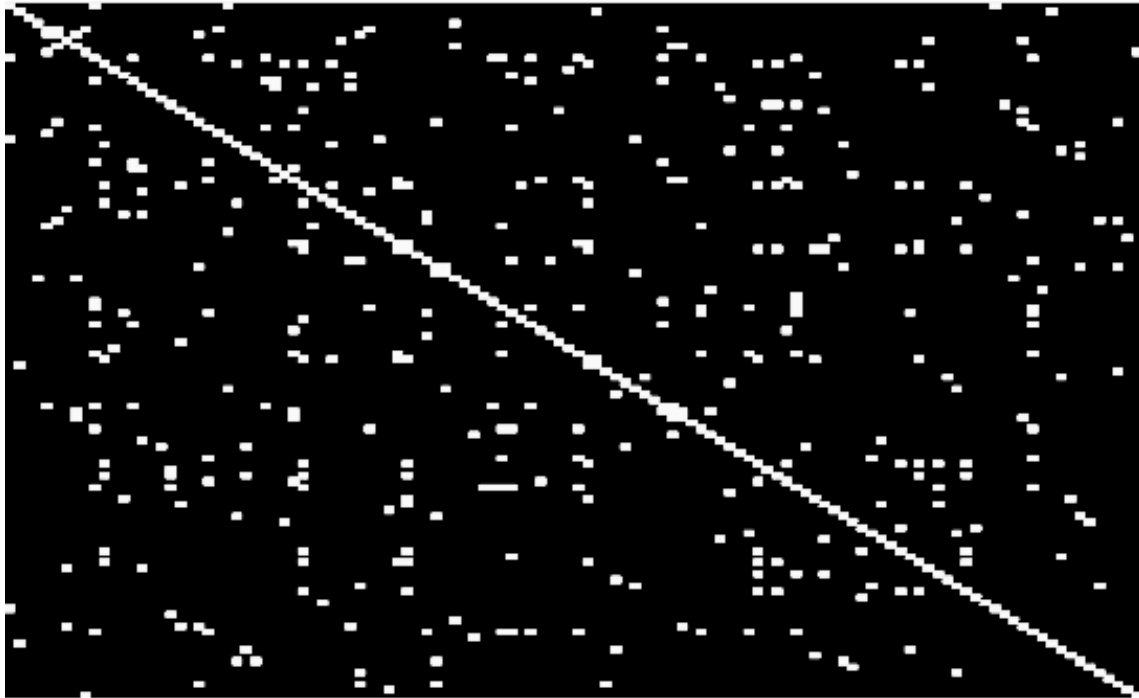
65 The data, consensually provided and anonymously treated, are derived from the results
66 obtained by a test-user (TU) from the personal genomic service 23andMe^[2]. Such results
67 typically consist of summary statistics on about one million single nucleotide polymorphisms
68 (SNPs)^[5].

69 A total of 120 anonymized individuals (progressively numbered with an ID from 1 to 120) were
70 considered in the analyses. All of them are ‘DNA-relatives’ of the TU according to the 23andMe
71 criteria and accepted the invitation to share their DNA information (excluding data related to
72 health conditions). The raw data is available as an Excel online matrix. Since this is a secondary
73 analysis of pre-existing data and the samples are treated in an anonymised version we did not
74 apply for an ethical clearance.

75 As reference parameter we considered the total amount of autosomal DNA in common
76 between pairs of individuals, calculated as the total length of shared SNP haplotype blocks in
77 mega base-pairs (Mbp) units. This amount, once converted into proportion of shared genome,
78 provides a rough estimate of the number of generations separating any two individuals, under a
79 simple model of “infinite number of ancestors” (Supplementary Table 1). Information either on
80 the relevant chromosomes where the match occurs, or on the number of segments in common
81 was not used. This choice is justified by the fact that only a minimum percentage of the
82 individuals considered shows DNA matches on more than one chromosome. Furthermore, the
83 information about the specific segment of the chromosome where such match occurs is not
84 easily obtainable from the data made available to the users by 23andMe.

85 Using the Genome-Wide Comparison option in the 23andMe ‘Family Traits’ feature, the input
86 data were prepared in the form of a symmetric square matrix C , whose $C(i,j)$ elements
87 correspond to the total length of shared SNP haplotype blocks between the individual i and the
88 individual j , expressed in Mbp units. Most elements of the matrix are equal to zero,
89 corresponding to the fact that the majority of the individuals does not result genetically related.
90 The sparsity of the matrix $C(i,j)$ is visually shown in Figure 1, where the white points indicate a

91 mutual match of any magnitude between two individuals, and the black correspond to no
92 genetic relation at all.



93
94 **Figure 1** – Visual representation of the correlation between the individuals considered in this work.

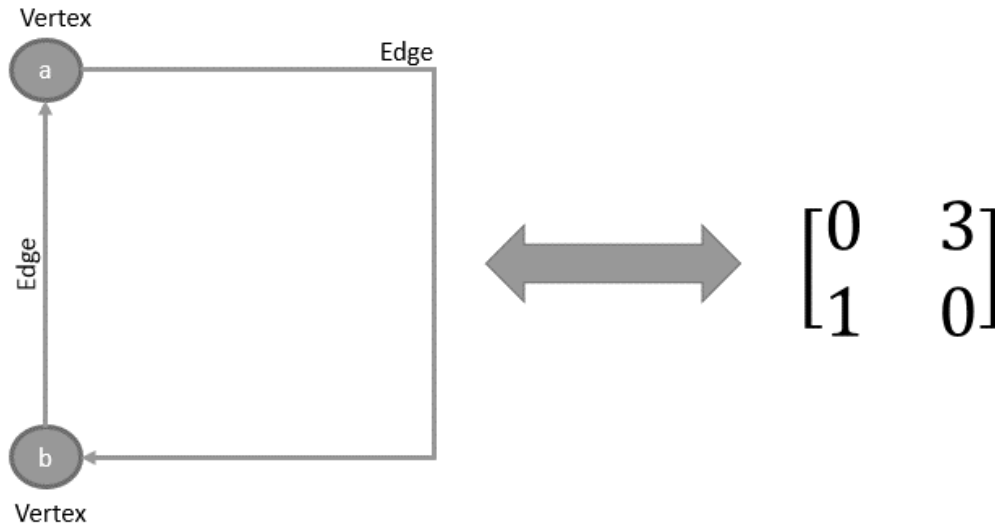
95 **3. Classification**

96 The matrix depicted in Figure 1 can be alternatively interpreted as a correlation matrix, a
97 covariance matrix, a similarity matrix^[6] or it can be transformed in a distance matrix^[7].
98 Accordingly, the way to elaborate and manipulate the associated information varies depending
99 on the interpretation tasks. Given that the statistical analysis is aimed at simplifying data
100 outputs, a loss of information with respect to the original data has to be expected. The
101 effectiveness of the analysis thus depends on the amount of ‘interesting’ information
102 unearthed out of the bulk of ‘redundant’ information. It follows that different methods can be
103 more or less effective according to what is considered, from time to time, interesting or
104 redundant.

105 To this extent, a number of potential confounders must be considered when dealing with the
106 available genetic similarity matrix. First of all the genetic information on which the analysis is
107 based is intrinsically fuzzy, because of the uncertainty in the data obtained by the service
108 provider (a few 'no-called' SNPs should be routinely expected). Additionally, the presence of
109 identical by state (IBS) other than identical by descent (IBD)^[8] SNPs could potentially bias the
110 genealogical interpretation, especially the one associated with distant relationships (Most
111 Recent Common Ancestors distant more than 6/7 generations). Finally, as opposed to
112 uniparental markers, the diploid autosomic data combine information inherited from the
113 paternal and maternal genealogy that should be kept separated when tracing one's ancestry.
114 Therefore, the analysis must be performed using statistical techniques robust enough to sustain
115 these unavoidable uncertainties.

116 4. Graph theory approach

117 The ideal framework for studying the complex network of links between the DNA-
118 relatives of a TU is the Graph Theory^[9,10]. This approach, widely used in Mathematics,
119 Engineering, and Computer Science, allows the analysis and graphical representation of the
120 links between different entities in a network. In synthesis, the Graph Theory represents the
121 elements in a network as *vertices* (or nodes) connected by *edges*. Edges are often associated
122 with a value representing a *weight*. In our case, the weight of an edge connecting two vertexes
123 is related to the genetic distance between them. A couple of vertexes *a* and *b* can be
124 connected, in principle, by more than one edge. Graphs can be generally oriented, so that the
125 edge from *a* to *b* is different from that linking *b* to *a*. In this way, the distance between the
126 vertexes *a* and *b* can be different from the distance between *b* and *a* (a typical example is
127 driving a car between two points in a city, where the traffic regulations might impose different
128 routes for the direct and return trip, see figure 2).



129

130 **Figure 2** – Graphic representation of a Graph with two vertexes and two edges (oriented Graph).

131 At the right, the corresponding adjacency matrix.

132 The relation between the vertexes is usually represented in matrix form (adjacency

133 matrix^[11]) where the elements out of the diagonal are the weights of the corresponding edges.

134 If the adjacency matrix is symmetric (the distance between two nodes is the same in both the

135 directions) the resulting graph is called *unoriented*.

136 In our scenario, the correlation matrix $C(i,j)$ between the DNA-relatives of the TU is

137 interpreted as a symmetric adjacency matrix. Therefore, we will use unoriented graphs,

138 implemented using the Matlab® code provided in Supplementary Materials.

139

140 5. Results

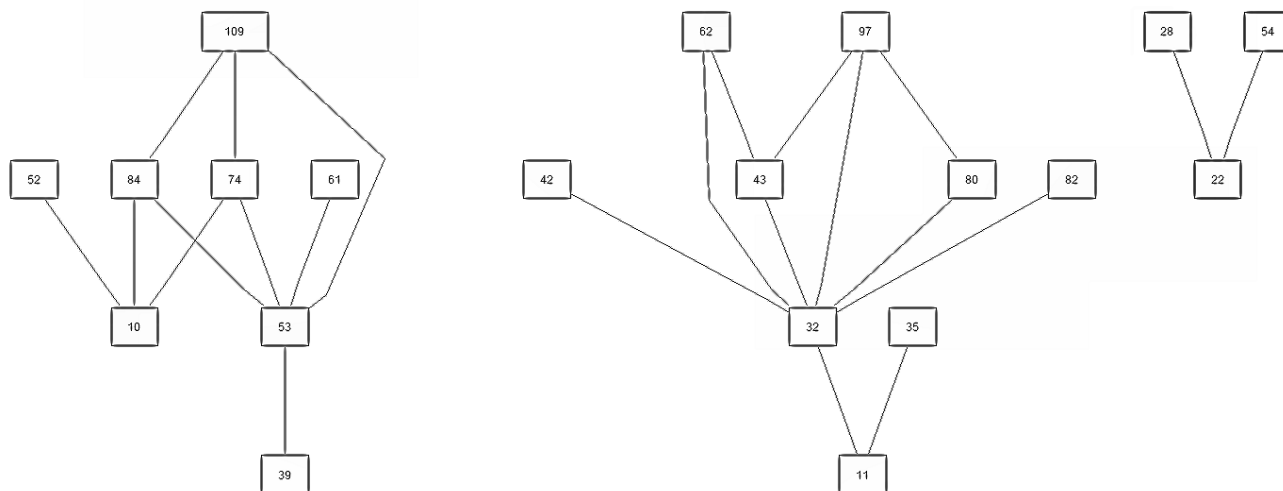
141 In the dataset analysed here, the adjacency matrix is described by an unweighted Graph
142 with 120 vertexes (individuals) and 196 edges (DNA links between them). The graphical
143 representation of the Graph described by this matrix is shown in Supplementary Figure 1.

144 The main network connects 100 vertexes (83% of the total) by 190 edges (97% of the total) and
145 sets aside only a few individuals, singularly (10 individuals) or in small groups of two or three
146 persons. A strict interpretation of Supplementary Figure 1 would thus bring to the conclusion
147 that all the individuals belonging to the main group should be considered as somehow related,
148 directly or indirectly, to all the other members of the group. To reduce this connectivity and to
149 assign the various individuals to the TU paternal and maternal ancestries, a further treatment
150 of the input data is thus necessary.

151 5.1 Pruning

152 The strength of the DNA cross-links between the individuals can be used to reduce (prune) the
153 connections highlighted in Supplementary Figure 1. Since all the 120 individuals included in this
154 study are, by design, related with the TU, no information can be derived from those that are
155 connected only to the TU. They are represented, in graphical form, as isolated vertexes with no
156 edges associated. Therefore, these individuals can be safely removed from the adjacency matrix
157 without any loss of information. Moreover, as already discussed in Section 3, spurious
158 connections could be introduced by fuzziness of the genetic data and the occurrence of IBS
159 SNPs. These connections can be excluded via the application of an upper threshold on the
160 genetic distances between the individuals. The threshold amount of shared genome for a link to
161 be considered 'real' (i.e., corresponding to IBD SNPs) can be easily converted into expected
162 number of generations, using Supplementary Table 1.

163 Figure 3 shows the Graph corresponding to the adjacency matrix $C(i,j)$ where only the edge
164 weights greater or equal to 24 Mbp (roughly a 8 generations distance between the
165 vertexes/individuals, see Supplementary Table 1) are considered.



166

167 **Figure 3** – Graphic representation of the Graph described by the adjacency matrix $C(i,j)$ considering only
 168 the edges corresponding to DNA-matches greater or equal to 24 Mbp. Isolated individuals and groups of
 169 two are not reported in the figure
 170

171 Figure 3 corresponds to the idea of unconnected graph that we associate with the separation of
 172 the different ancestral lines of the TU. Surprisingly enough, when the results of the Graph
 173 Theory are compared with the pre-existing genealogical information on some of the matching
 174 individuals, it turns out that the two large groups correspond to relatives of the TU related to
 175 the maternal grandfather (at the center of the figure) and maternal grandmother (at the left).
 176 Another small group of three individuals, at the right in figure 3, shows up, containing an
 177 individual associated to the maternal grandmother's lineage of the TU (n.22). The two
 178 individuals that can be identified with reasonable certainty as belonging to the paternal
 179 grandfather's (n. 118) and grandmother's (n. 96) lineage of the TU, remains unconnected.
 180 These results are summarized in Table 1. The individuals underlined and marked in bold are the
 181 ones for whom a genealogical evidence exists, and therefore can be assigned with certainty to a
 182 given lineage. The ones underlined and marked in *italic*, on the other hand, cannot be assigned
 183 with similar certainty, although there are strong independent clues suggesting that they would
 184 actually belong to that lineage.

185

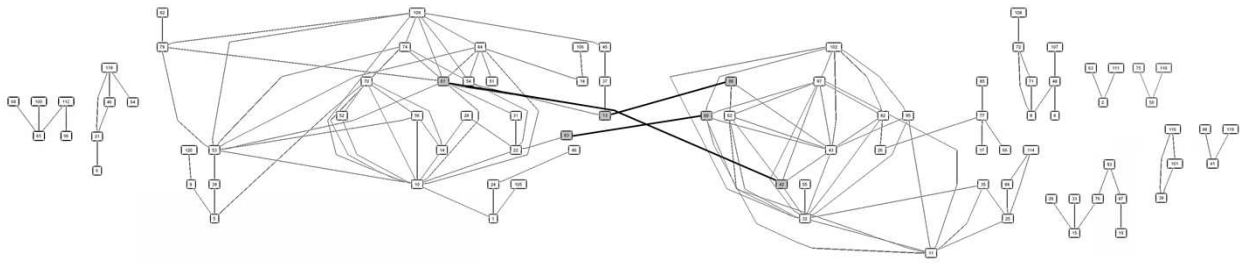
186 **Table 1** – Classification of the individuals according to their lineage (24 Mbp threshold). Individuals
 187 underlined and marked in bold are the ones for whom a genealogical evidence exists, and
 188 therefore can be assigned with certainty to a given lineage. The ones underlined and marked in
 189 italic, on the other hand, cannot be assigned with similar certainty, although there are strong
 190 independent clues suggesting that they would actually belong to that lineage.

Paternal GF	Paternal GM	Maternal GF	Maternal GM	Unclassified
		62	52	28
		<u>97</u>	<i><u>109</u></i>	54
		<u>42</u>	84	<i><u>22</u></i>
		43	<i><u>74</u></i>	
		80	61	
		82	<i><u>10</u></i>	
		<u>32</u>	<i><u>53</u></i>	
		<u>35</u>	39	
		<u>11</u>		

196 The adoption of a conservative threshold (24 Mbp / approx. 8 generations distance / 3rd – 4th
 197 cousin range) to define a link between the individuals produced the classification reported in
 198 Table I, which is robust and reliable. However, only 17 individuals over a total of 120 (110 with
 199 at least one DNA match besides the TU) are attributed to the corresponding ancestral lineage.

200 Reducing the level of the threshold to 12 Mbp (approx. 9 generations distance) increases the
 201 number of individuals that can be associated to the different groups (Supplementary Figure 2).
 202 Individual 22 is now correctly associated to the maternal grandmother's group, along with the
 203 other members of his/her subgroup. Most importantly the graph now shows an additional
 204 group of three individuals (21, 46 and 118) that can be associated to the TU paternal
 205 grandfather's lineage, on the basis of independent genealogical information existing for
 206 individual 118.

207 Further lowering the threshold to 6 Mbp (approx. 10 generations distance, i.e. a 4th – 5th cousin
 208 range, which is usually considered the lower limit for having a significant DNA match between
 209 two individuals) allows to recover important information, graphically represented in figure 4.



210

211 **Figure 4** – The same as in figure 3, considering only the edges corresponding to DNA-matches greater or
 212 equal to 6 Mbp. The individuals connecting the two main groups and their links are evidenced. Isolated
 213 individuals and groups of two are not reported in the figure

214 From the analysis of figure 4 it is evident that after lowering the threshold to 6 Mbp, a
 215 connection appears between the two main groups. The key elements which are linked to both
 216 the groups (corresponding to the maternal grandparents of the TU) are individual 83 (initially
 217 classified in the maternal GM group) which connects with individual 80 in the maternal GF
 218 group), individual 61 of the maternal GM group which connects with individual 42 in the
 219 maternal GF group, and individual 86 of the maternal GF group which connects with individual
 220 13 of the maternal GM group.

221 Lowering the threshold also increased the number of individuals associated to the paternal
 222 grandfather of the TU, which at this level formed a group of five persons (118, 46, 21, 6 and 64)
 223 connected by the same sub-graph, and recovered a new group of five individuals (96, 112, 65,
 224 100 and 68) that can be associated to the TU paternal grandmother's lineage on the basis of
 225 independent genealogical information existing for individual 96.

226 The main information that can be derived by the comparison of the Graphs obtained using
 227 different thresholds on the edge weight is a classification of the individuals according to the
 228 different ancestral lineages, with increasing 'levels of confidence'. In that respect,
 229 Supplementary Figure 1 would give a minimum level of information, providing classification at
 230 the confidence level of the minimum match in the $C(i,j)$ matrix, which in our case is 2 Mbp,
 231 subsequently refined at higher thresholds of genomic sharing in Figure 3, Supplementary Figure
 232 2 and Figure 4.

233 The most important results of this paper are shown in Table II, where the classification of the
 234 DNA-relatives of the TU is reported according to his maternal and paternal ancestral lineages,
 235 with the corresponding confidence level, or 'strength', in brackets. The individuals connecting
 236 the groups corresponding to the two maternal grandparents are assigned to both the groups
 237 and marked in gray.

238 **Table 2** – Classification of the individuals according to their ancestral lineage. The corresponding level of
 239 confidence of the classification is reported in brackets. The individuals connecting the two groups of the
 240 maternal grandparents are marked in gray.

Paternal GF	Paternal GM	Maternal GF	Maternal GM	Unclassified
<u>118</u> (12)	<u>96</u> (6)	<u>97</u> (24)	52 (24)	116
46 (12)	112 (6)	62 (24)	<u>109</u> (24)	101
21 (12)	65 (6)	42 (24)	84 (24)	38
6 (6)	100 (6)	43 (24)	<u>74</u> (24)	-----
64 (6)	68 (6)	80 (24)	61 (24)	29
		82 (24)	<u>10</u> (24)	33
		32 (24)	<u>53</u> (24)	15
		35 (24)	39 (24)	76
		<u>11</u> (24)	<u>70</u> (12)	93
		102 (12)	83 (12)	87
		95 (12)	54 (12)	19
		25 (6)	51 (12)	-----
		89 (6)	92 (12)	63
		114 (6)	79 (12)	111
		83 (6)	<u>14</u> (12)	2
		61 (6)	<u>22</u> (12)	-----
		85 (6)	9 (12)	75
		77 (6)	5 (12)	110
		17 (6)	28 (12)	50
		66 (6)	1 (6)	-----
		26 (6)	24 (6)	88
		55 (6)	56 (6)	119
		13 (6)	42 (6)	41
		86 (6)	18 (6)	-----
			106 (6)	108
			31 (6)	72
			13 (6)	71
			86 (6)	8
			45 (6)	48
			37 (6)	4
			120 (6)	107
			105 (6)	

241

			40 (6)	
			80 (6)	

242

243 The Graph Theory method here proposed is capable of reliably classifying 62 individuals at
244 strength 6 (Mbp) over a total of 110 DNA-relatives of the TU (56%). Six other unclassified
245 groups with more than two members can also be determined. Some of them could be
246 connected to the main groups if additional information from new DNA relatives of the TU will
247 become available in the future.

248

249 **Conclusion**

250 The statistical method presented in this work can be usefully exploited for extracting
251 genealogical information from genetic/genomic data. The input data are usually 'fuzzy' and,
252 therefore, the methods used for their analysis should be robust enough for providing useful
253 information. The approach proposed, based on the Graph representation of the adjacency
254 matrix built from the mutual matches between the DNA-relatives of the test user, after the
255 setting of a suitable threshold fulfils this requirement. The method, for which the code is
256 provided at the bottom of this paper, could be easily implementable by the genetic service
257 providers for an easy visualization of the DNA-links existing between the customer and the
258 other users of the service, at different levels of confidence.

259

260 **References**

- 261 [1] <https://genographic.nationalgeographic.com/>
- 262 [2] <https://www.23andme.com/>
- 263 [3] <https://www.familytreedna.com/>
- 264 [4] <http://www.gedmatch.com/>
- 265 [5] Nachman, M. W. (2001) Single nucleotide polymorphisms and recombination rate in
266 humans, *Trends in genetics*, 17 (9), 481–485.
- 267 [6] Muni S. Srivastava, *Methods of Multivariate Statistics* (2002) Wiley.
- 268 [7] Smouse, P. E. and Long, J. C. (1992) Matrix correlation analysis in anthropology and
269 genetics, *American Journal of Physical Anthropology*, 35, 187–213.
- 270 [8] Stevens, E.L., Heckenberg, G., Roberson, E.D.O., Baugher, J.D., Downey, T.J., Pevsner, J.
271 (2011) Inference of relationships in population data using identity-by-descent and identity-by-
272 state, *PLoS Genetics*, 7 (9), art. no. e1002287.
- 273 [9] Bondy, A., Murty, U.S.R., *Graph Theory* (2008) Springer-Verlag London.
- 274 [10] Pavlopoulos, G. A., Secrier, M., Moschopoulos, C. N., Soldatos, T. G., Kossida, S., Aerts, J.,
275 Schneider, R., Bagos, P. G. (2011) Using graph theory to analyze biological networks, *BioData*
276 *Mining*, 4:10.
- 277 [11] Gehlenborg, N., Wong, B. (2012) Points of View: Networks, *Nature Methods*, 9, 115.