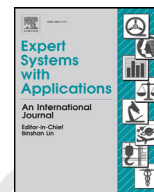




ELSEVIER

Contents lists available at ScienceDirect

Expert Systems With Applications

journal homepage: www.elsevier.com/locate/eswa

Social networks and genetic algorithms to choose committees with independent members

Eduardo Zamudio*, Luis S. Berdún, Analía A. Amandi

ISISTAN Research Institute (UNICEN/CONICET), Campus Universitario, Paraje Arroyo Seco, Tandil, Buenos Aires B7001BBO, Argentina

ARTICLE INFO

Keywords:

Committee
Group selection
Independence
Social network
Genetic algorithm

ABSTRACT

Choosing committees with independent members in social networks can be regarded as a group selection problem where independence, as the main selection criterion, can be measured by the social distance between group members. Although there are many solutions for the group selection problem in social networks, such as target set selection or community detection, none of them have proposed an approach to select committee members based on independence as group performance measure. In this work, we propose a novel approach for independent node group selection in social networks. This approach defines an independence group function and a genetic algorithm in order to optimize it. We present a case study where we build a real social network with on-line available data extracted from a Research and Development (R&D) public agency, and then we compare selected groups with existing committees of the same agency. Results show that the proposed approach can generate committees that improve group independence compared with existing committees.

© 2015 Elsevier Ltd. All rights reserved.

1. Introduction

Organizations need representative individuals to make decisions about particular concerns. These representative individuals are appointed in committees, and we expect from these members to make decisions based on the benefit of the whole community they are representing, avoiding bias that could arise from closeness between them. In this context, the best committees are those which show the greatest independence between his members. How to choose these members based on objective criteria could be a difficult task, either because of the definition of the criteria or because of the analysis of the community from where members are chosen. Therefore, a committee in which some of its members are closely related is an unbalanced committee.

Fig. 1 shows a graphical example of difference between balanced and unbalanced committees that allow us to appreciate the distribution of selected nodes within a graph. A balanced distribution is essential to improve desirable features, such as independence. For instance, a committee to discuss about budget allocation must avoid biased decisions by ensuring that committee members are not closely related.

As mentioned before, Fig. 1 shows a simple example of individuals and their dispersal. Fig. 2 shows a graphical representation of the community used to evaluate this approach. This graph allows us to understand the problem complexity and critical importance of choosing the best committee members to maximize independence.

Initially, the committee member selection problem can be solved by a mathematical combination, but the computational cost associated to this approach could be very high. For instance, given a community with n members, the maximal number of groups is given by $2^n - 1$, and complexity is $O(2^n)$. In addition, if committees are r size groups, the number of possible solutions is given by applying binomial coefficient ${}_nC_r$ and complexity is $O(n!)$.

If there is no polynomial function to solve the problem, an alternative could be to adopt a non deterministic approach to approximate optimal solutions. For instance, a stochastic approach could produce random solutions, and then apply an independence function to rank these solutions. This approach is subjective because of the probability in selecting random committee members, and because of the joint probability of the committee.

However, the problem can be addressed by implementing some optimization strategy to approximate optimal solutions, such as genetic algorithms. A genetic algorithm could be implemented to search for the greatest independence between committee members, but not necessarily to guarantee the best solution. In other words, could be enough to approximate an optimal solution. For committee selection problem, the best solutions will be determined by the maximal independence between his members.

* Corresponding author. Tel.: +54 249 4439882/3764206786.

E-mail addresses: eduardo.zamudio@isistan.unicen.edu.ar, eduardozamudio@gmail.com (E. Zamudio), lberdun@exa.unicen.edu.ar (L.S. Berdún), amandi@exa.unicen.edu.ar (A.A. Amandi).

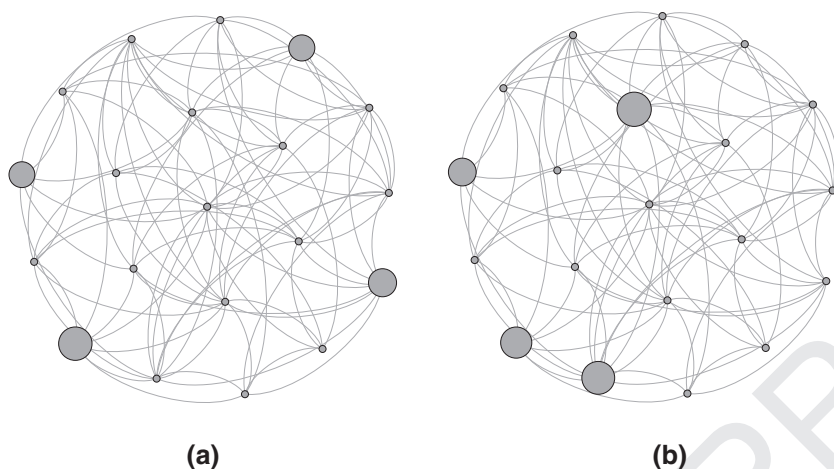


Fig. 1. Difference between balanced (a) and unbalanced (b) committees, where selected members are the largest 4.

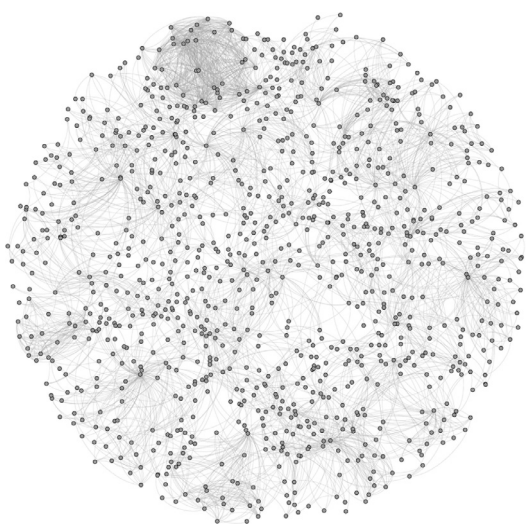


Fig. 2. Graphical representation of the community used to evaluate the approach.

to evaluate the proposed approach. Then, we compare results with current committees of the same public agency.

This document is organized as follows. Section 2 describes the construction process of the social network. Section 3 describes the implementation of the genetic algorithm and the function to evaluate group independence. Section 4 describes a case study and the configurations of the genetic algorithm, along with a discussion of the experiment results. Section 5 presents a discussion of the current literature. Finally, Section 6 presents conclusions as well as future work.

2. Social network construction

In order to choose committee members, we propose to build a social network to calculate distances between candidates, and then apply a genetic algorithm to get potential committees with the greatest distances between their members.

A social network is a set of individuals (actors) and relations (ties) between them; the social network analysis is used to study structures created by these relations and individuals.

We are particularly interested in the construction of a social network for its ability to represent analysis criteria based on ties. To clarify this concept, we built a network of researchers related through co-authorship and workplace. In this network, actors are the researchers, and ties are the criteria for calculating distance between each pair of researchers.

As mentioned above, relations between actors define what can be analyzed in the network. The aim of this analysis is to calculate distances between a set of actors. In order to do this, we built a consolidated graph. This graph contains every kind of relation proposed as analysis criterion. Fig. 3 shows a unified graph from two kinds of relations (*coauthor* and *same workplace*) of five researchers (A, B, C, D, and E) where relations are binary (relation is present or not), undirected (direction is meaningless), and irreflexive (a researcher does not publish with himself or does not work with himself).

Our proposal is to establish the greatest independence between committee members based on their distances. Thus, we need to calculate distances between committee members, for which we use the *shortest path* and *geodesic distance* (length of the shortest path) (Freeman, 1977) over the unified graph.

The graph must be connected to apply this metrics, which means that every actor must be reachable from every other actor in the network. This can be determined through a reachability matrix, which can be obtained through matrix multiplication (Wasserman & Faust, 1994).

Distances between each pair of actors is represented by a proximity matrix, obtained by applying power to the matrix

If we consider committee candidates as individuals connected with each other through ties, it is possible to determine which of these ties could be relevant to analyze independence. These individuals and their relationship represent the basic elements of a social network; therefore, we can apply social network analysis to select committee members with the greatest independence. However, a social network approach requires a social network, and data to represent its elements, such as actors, ties, kind of network, and analysis object.

The current social network analysis techniques aim to identify the value or number of relations, roles or prominence of nodes, and to discover hidden groups or cohesive groups. The aim of this work is to present an alternative to the committee selection problem by choosing a set with maximal independence between members. To do this, we build a social network and then we define an independence group performance function and a genetic algorithm, to obtain n member committees with the greatest independence between members.

The main contributions of this work are summarized as follows. (1) We propose an approach for the committee selection problem with independent members as a group selection problem in social networks. (2) We define a novel group independence performance function to assess group fitness in social networks. Then, such a measure was optimized by means of a genetic algorithm. (3) We build a social network from a Research and Development (R&D) public agency with on-line available data. (4) We use such a social network

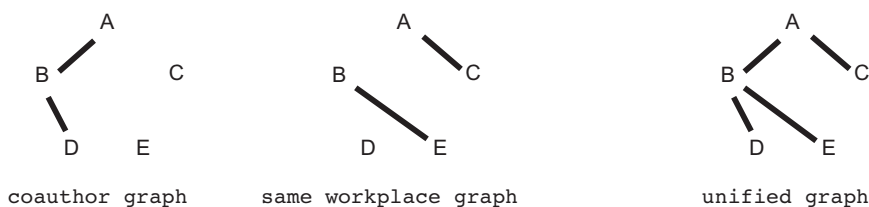


Fig. 3. Unified graph representing two kinds of relationships (coauthor and same workplace).

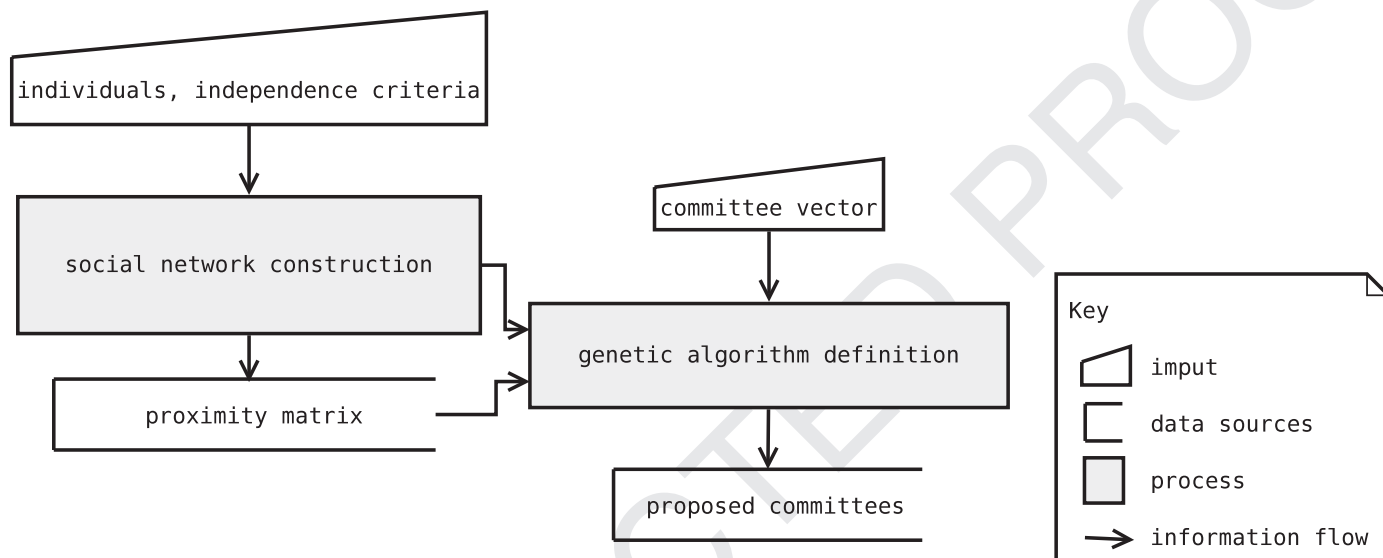


Fig. 4. Flow chart of the proposed approach showing inputs, datasources production, and processes related to the social network construction and the genetic algorithm definition.

116 representation of the unified graph. Proximity matrix contains input
 117 data for the algorithm whose aim is to choose a committee (a group
 118 of actors) with the greatest independence between its members. In
 119 this case, we work with a genetic algorithm which defines a function
 120 to optimize this distance to the largest one.

121 Fig. 4 shows the proposed approach in a flow chart, in which in-
 122 dividuals and independence criteria are the inputs. Then, we gener-
 123 ate the unified graph to determine relations between actors, and thus
 124 to build the social network. Next, we build the proximity matrix by
 125 calculating geodesic distances; then, the proximity matrix and the
 126 network data are put together into the genetic algorithm to produce
 127 optimized solutions.

128 **3. Genetic algorithm definition**

129 A genetic algorithm is a type of evolutionary algorithm that can be
 130 considered as a function optimization method (Smith & Eiben, 2008).
 131 Even though there is no definitive genetic algorithm, it is possible to
 132 adapt one using representations and operators considered suitable to
 133 the modeled problem. As an analogy of the biological model, chro-
 134 mosomes are the elements used in genetics algorithms to represent
 135 configurations, which contain genetic information represented by lo-
 136 cation and value of their genes. These chromosomes stand for solu-
 137 tions to the modeled problem.

138 In order to choose a subset of actors from a social network, we
 139 have defined an ad-hoc function to calculate distances between com-
 140 mittee members. Consequently, we have defined a genetic algorithm
 141 to approximate solutions to an optimum by maximizing this function.

142 The development of a genetic algorithm requires defining repre-
 143 sentation, fitness function, parent selection and survivor selection
 144 mechanisms as well as mating and mutation operators. Next, we
 145 present selected configurations to the modeled problem.

147 **3.1. Representation**

148 The problem requires defining a representation of the chromo-
 149 some. In this work, we do permutations of a vector of integers (chro-
 150 mosome), where each element references to only one node (gene). In
 151 this vector, every node in the network under study is included. Thus,
 152 a chromosome has as many genes as a community has individuals.
 153 Also, the participation in the committee is given by a vector with the
 154 same size as the vector of nodes, the vector of committee members,
 155 in which every location is binary valued. Therefore, if value = 1, then
 156 the node with same position in the vector of nodes must be included
 157 in the committee, and if value = 0, then the node is excluded from the
 158 committee. With this representation, a member appears only once in
 159 a given committee. It is important to note that in the modeled prob-
 160 lem the order of members is not relevant. Fig. 5 shows a graphical
 representation of these vectors.

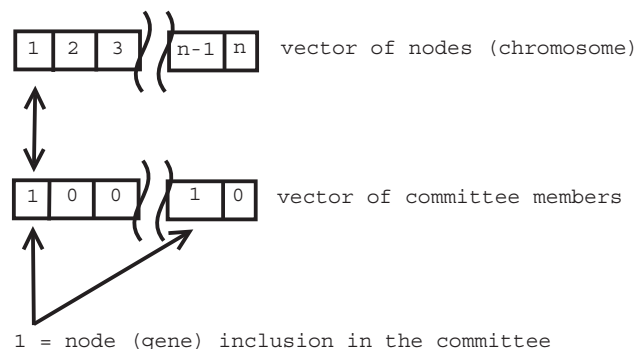


Fig. 5. Representation of the genetic algorithm through a vector of nodes which contains every node in the network under study, and a vector of committee members which indicates the elements of the vector of nodes to be included in the committee.

161 3.2. Fitness function

162 The aim of the fitness function is to calculate the solution value.
 163 In this work, we developed an ad-hoc fitness function to maximize
 164 distances, represented by the cumulative sum of distances between
 165 each pair of committee members. In order to get relative values, we
 166 consider the size of the committee and the network diameter. To im-
 167 prove results, we set a parameter to maximize minimum distances of
 168 the committees, defined as follows:

$$f = \frac{\left[\left(\sum_{i,j=0}^k d(i, j) \right) / k \right] + m}{2 * D}$$

169 Where d is the distance function between two members, $\forall i, j | i \neq j$
 170 and $i, j \in S$, S represents the whole nodes set, k is the number of com-
 171 mittee members, m is the minimum distance between each pair of
 172 members in the committee, and D is the network diameter. As previ-
 173 ously established, it is necessary for the network to be connected.

174 3.3. Parent selection

175 The genetic information is obtained from the parents, which are
 176 chromosomes (solutions) of the previous generation. To this end,
 177 we need to define a strategy of parents selection by adopting one
 178 of the mechanisms suitable to the modeled problem. In this work,
 179 the mechanisms selected include *Stochastic Universal Sampling* (SUS)
 180 since we need to choose several parents from a community; and *Tour-*
 181 *namment*, since in both cases global fitness is unknown.

182 3.4. Crossover

183 Genetic information of new generations is determined by their
 184 parents. This process called genetic recombination is produced
 185 through crossover mechanisms. For instance, having two chromo-
 186 somes representing distinct solutions, crossover implies that the new
 187 generation inherit genetic information from both parents.

188 To keep a valid permutation we have chosen recombination op-
 189 erators *Partially Mapped Crossover* (PMX) and *Order Crossover* (OX).
 190 Since the former is an algorithm designed for adjacency problems it
 191 is suitable to the modeled problem, and even though the latter is de-
 192 signed for order problems, the order in the second parent could be
 193 beneficial in new chromosome production.

194 3.5. Mutation

195 The other mechanism used in this work for genetic recombination
 196 is mutation, which implies to alter the genes within a chromosome.
 197 In permutations, mutation alters location of the values in the solution
 198 vector of the new generation.

199 We have selected *Swap Mutation* and *Insert Mutation*, since both
 200 operators are accepted to keep a valid permutation.

201 3.6. Survivor selection

202 Once a new generation is produced, the survivors must be selected
 203 in order to keep the number of solutions in every generation.

204 We have selected *Steady-state* and *Generational* mechanisms to
 205 keep solutions with the best fitness in the succeeding generations.

206 4. Case study

207 To evaluate the proposed approach, we decided to build a social
 208 network based on public information about researchers published by
 209 the National Council of Technical and Scientific Research (CONICET).

This organization establishes committees for specific areas with dif- 210
 ferent responsibilities. For instance, in the Informatics and Commu- 211
 nications area there are 3 committees to evaluate *Admissions, Reports,* 212
and Fellowship awards. 213

The prospective committee members are chosen from a set of ex- 214
 perts in the field that could be internal or external to the organization. 215

We calculated fitness for distinct configurations of the genetic al- 216
 gorithm to propose committees based on the greatest distances. With 217
 the same criteria, we calculated fitness for existing committees. 218

219 4.1. Dataset

The dataset used here to produce the social network based on re- 220
 searchers (actors) information was built by applying *web crawling,* 221
 which consists in gathering information from web pages. In this case, 222
 we used basic information to characterize actors and their informa- 223
 tion about contributions and workplaces in order to discover ties be- 224
 tween those actors. This process required disambiguation of actors 225
 and ties, since most of the information presented for every researcher 226
 is produced by themselves, particularly contribution data. 227

In addition, not every actor in the network is considered as can- 228
 didate. For the Informatics and Communication area there is a list of 229
 qualified specialists that fulfill some requirements (i.e., to have a hier- 230
 archical degree), which means that only a limited set of actors qualify 231
 as committee members. 232

Thus, the social network in the case study is composed by 1293 233
 nodes and 4322 ties, which produces 74 components (subgroups of 234
 actors disconnected from the rest of the network). From those com- 235
 ponents, the bigger one has 1058 ($\approx 82\%$) actors (75 of them are qual- 236
 ified specialist), and 3878 ($\approx 90\%$) ties. 237

238 4.2. Configuration

Having established the social network, we set up the genetic algo- 239
 rithm to evaluate groups of actors with the largest distances between 240
 them, which we assume as an independence criterion. This configura- 241
 tion has the following parameters: 242

- Community size: The number of solutions in every moment was 243
 given by P/n , where P is the set of all researchers, and n the size of 244
 the committees. 245
- Crossover probability: A generational parameter, selected from 246
 range [0.6; 0.9]. 247
- Mutation probability: A mutation operator parameter, selected 248
 from range [0.01; 0.15]. 249
- Stop condition: A generational parameter, set in 25 generations. 250
- Configurations: Sixteen different configurations emerged from 251
 the combination of the selected mechanisms in this approach (se- 252
 lection, mutation, and crossover). In addition, we use *Steady-state* 253
 and *Generational* as selection mechanisms. Table 1 shows these 254
 configurations. 255
- Runs: 40 runs produced by 5 runs per configuration. Average val- 256
 ues and standard deviation (σ) were calculated. 257

258 4.3. Results

Here we show a fitness evaluation and social network centrality 259
 metric values for current committees of the Informatics and Com- 260
 munications area, and then we show results of the genetic algorithm 261
 runs. 262

263 4.3.1. Fitness of current committees

The current committees of the Informatics and Communications 264
 area had 6 members in 2014. In order to evaluate committee fitness 265
 we initially decided to apply the fitness function to committee mem- 266
 bers. This approach was modified since some members of the current 267

Table 1
16 proposed configurations for the genetic algorithm describing operators and selection mechanisms.

Configuration	Crossover		Mutation		Parent selection		Survivor selection	
	PMX	OX	Swap	Insert	SUS	Tournament	Steady-state	Generational
1	X		X		X		X	
2		X	X		X		X	
3	X			X	X		X	
4		X		X	X		X	
5	X		X			X	X	
6		X	X			X	X	
7	X			X		X	X	
8		X		X		X	X	
9	X		X		X			X
10		X	X		X			X
11	X			X	X			X
12		X		X	X			X
13	X		X			X		X
14		X	X			X		X
15	X			X		X		X
16		X		X		X		X

268 committees were not present in the dataset. This situation occurs be-
 269 cause of the low number of specialists in the area belonging to CON-
 270 ICET (actually there are 87 specialists in the Informatics and Com-
 271 munications area), which means that committees usually incorporate
 272 external researchers from other areas. Therefore, we have identified
 273 the current committees members present in the largest component
 274 of the proposed social network. In the *Admissions* committee, only
 275 3/6 members are present in the social network; in the *Reports* com-
 276 mittee, only 4/6 members are present in the social network; and in
 277 the *Fellowship awards* committee, only 5/6 members are present in
 278 the social network. Since names of the committee members are not
 279 relevant in this study, we enumerated members from 1 to 6 for each
 280 committee.

281 The *Admissions* committee of the Informatics and Communication
 282 area has fitness = 0.65152 for members A1–A3, since A4 is present in
 283 another component and A5 and A6 are not classified as specialists.
 284 The other 2 committees are in similar situation. The *Reports* commit-
 285 tee has fitness = 0.36364 for members R1–R4, since the other mem-
 286 bers of the committee do not belong to CONICET (R5) or are not clas-
 287 sified as specialists in the area (R6). And the *Fellowship awards* has
 288 fitness = 0.38636 for members F1–F5, since F6 does not belong to
 289 CONICET. Table 2 shows current committee members with centrality
 290 metric values for those members present in the largest component of
 291 the social network.

4.3.2. Social network metrics

292 The social network metrics for current committees shown in
 293 Table 2 can be compared with metrics of the whole component,
 294 which average degree = 7.316, network diameter = 11, and average
 295 path length = 5.76. This indicates that almost every member (except
 296 for F2) of current committees has degree over the average component
 297 degree, but far away from the highest degree (80) in the component.
 298 Some committee members (A3 and F2) show very low betweenness,
 299 but their closeness is more balanced between each other.

4.3.3. Genetic algorithm runs

301 In order to compare the fitness of current committees with the fit-
 302 ness of the members proposed by the genetic algorithm, we decided
 303 to modify the genetic algorithm to generate committees of 3, 4, and 5
 304 members.

305 For the *Admissions* committee, we set up the genetic algorithm in
 306 order to produce committees with 3 members. Table 3 shows results
 307 where maximal average fitness ≈ 0.72727 and minimal $\sigma = 0$ for con-
 308 figurations 9 and 11. Maximal fitness ≈ 0.72727 was reached by con-
 309 figurations 9, 11, 12, and 13, from which we infer that a local optimum
 310 is reached in these cases.

Table 2
Current *Admissions*, *Reports*, and *Fellowship awards* committees with each
 degree, betweenness and closeness (last two metrics expressed in relative
 values).

Committee	Node	Degree	Betweenness	Closeness
<i>Admissions</i> (fitness = 0.65152)	A1	49	0.05293	0.22404
	A2	21	0.02593	0.17283
	A3	5	0.00001	0.15606
	³ A4	–	–	–
	² A5	–	–	–
	² A6	–	–	–
<i>Reports</i> (fitness = 0.36364)	R1	35	0.11858	0.20596
	R2	51	0.11909	0.25101
	R3	37	0.03512	0.19495
	R4	34	0.14864	0.20989
	¹ R5	–	–	–
	² R6	–	–	–
<i>Fellowship awards</i> (fitness = 0.38636)	F1	22	0.01246	0.15696
	F2	6	0.00001	0.16307
	F3	19	0.00595	0.19317
	F4	42	0.07272	0.22751
	F5	46	0.06920	0.23731
	¹ F6	–	–	–

¹ Does not belong to CONICET.

² Not marked as specialist.

³ Present in another component.

Table 3
Fitness of proposed configurations with average fitness, standard deviation, and max-
 imal fitness for 3-member committees in 5 runs (best values in **bold**).

Configuration	Average fitness (runs = 5)	σ	Maximal fitness (with the shortest time in seconds)
1	0.58788	0.01134	0.59091
2	0.57879	0.01134	0.59091
3	0.60303	0.02607	0.65152
4	0.60909	0.02938	0.66667
5	0.61818	0.02938	0.65152
6	0.62424	0.03090	0.66667
7	0.62121	0.02710	0.65152
8	0.62121	0.03711	0.66667
9	0.72727	0.00000	0.72727
10	0.64545	0.00742	0.65152
11	0.72727	0.00000	0.72727
12	0.67879	0.02607	0.72727
13	0.70606	0.02642	0.72727
14	0.63939	0.02607	0.66667
15	0.67879	0.00606	0.68182
16	0.61515	0.02642	0.65152

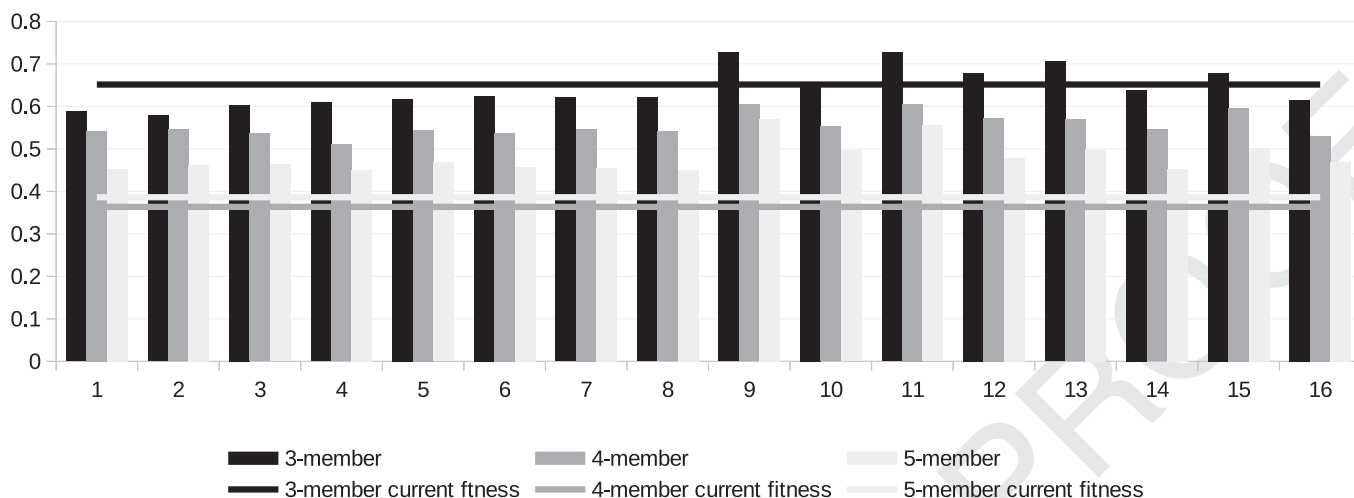


Fig. 6. Average fitnesses of 3-member, 4-member, and 5-member committees for the 16 configurations.

Compared with current committee fitness ≈ 0.65152 , maximal average fitness shows a fitness improvement of ≈ 8 points.

For the *Reports* committee, we set up the genetic algorithm in order to produce committees with 4 members. Results show maximal average fitness ≈ 0.60606 and minimal $\sigma = 0$ for configuration 11. Maximal fitness ≈ 0.60606 was reached by configurations 9, 11, and 15, from which we infer that a local optimum is reached in these cases.

Compared with the current committee fitness ≈ 0.36364 , maximal average fitness shows a fitness improvement of ≈ 24 points.

For the *Fellowship awards* committee, we set up the genetic algorithm in order to produce committees with 5 members. Results show maximal average fitness ≈ 0.57091 for configuration 9, minimal $\sigma \approx 0.00530$ for configuration 4, and maximal fitness ≈ 0.59091 for configurations 9 and 11.

Compared with current committee fitness ≈ 0.38636 , maximal average fitness shows a fitness improvement of ≈ 20 points.

As shown in Fig. 6, *Generational* (configurations 9–16) selection mechanism produced better results than *Steady-state* (configurations 1–8), but Fig. 7 shows that *Generational* required more time than other configurations. For instance, in 5-member committees, the minimal time for *Steady-state* = 4.73 s. (seconds) and for *Generational* = 67.049 s. This situation is similar for 3-member and

4-member committees. In order to reach the time required by *Generational* configurations, we extended *Steady-state* stop condition to 25,000 generations, resulting always in lower fitnesses than those obtained with *Generational* mechanism configurations.

For 3-member and 5-member committees, configuration 9 presents the fullest average fitness, and for all committees, configurations 9 and 11 show the highest maximal fitness values, from which we infer that in searching for optimal values in similar studies we should prefer the *Generational* selection mechanism and the PMX operator. In addition, in this case the mutation operator does not produce relevant differences. However, in bigger or more complex networks, computational cost improvement may be a requirement, in which cases we should prefer *Steady-state* selection mechanism instead of *Generational* selection mechanism. In addition, Fig. 8 shows that 3-member configurations 9 and 11 reached $\sigma = 0$, and that 3-member and 4-member configuration 9 reached $\sigma = 0$, from which we infer the stability of these configurations, at least for 3-member and 4-member committees.

Fig. 9 shows the social network built for the case study, in which current committee members are closer than the best fitness committee members obtained in experimentation. This representation shows a balance improvement of distances between the best fitness committee members compared to current committee members.

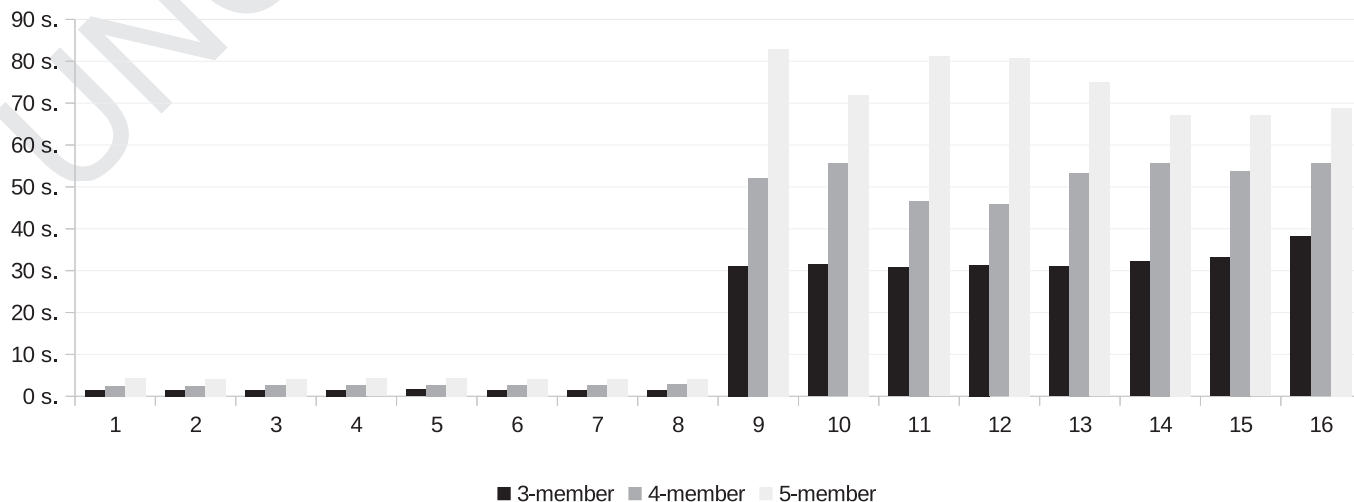


Fig. 7. Shortest times of 3-member, 4-member, and 5-member committees for the 16 configurations.

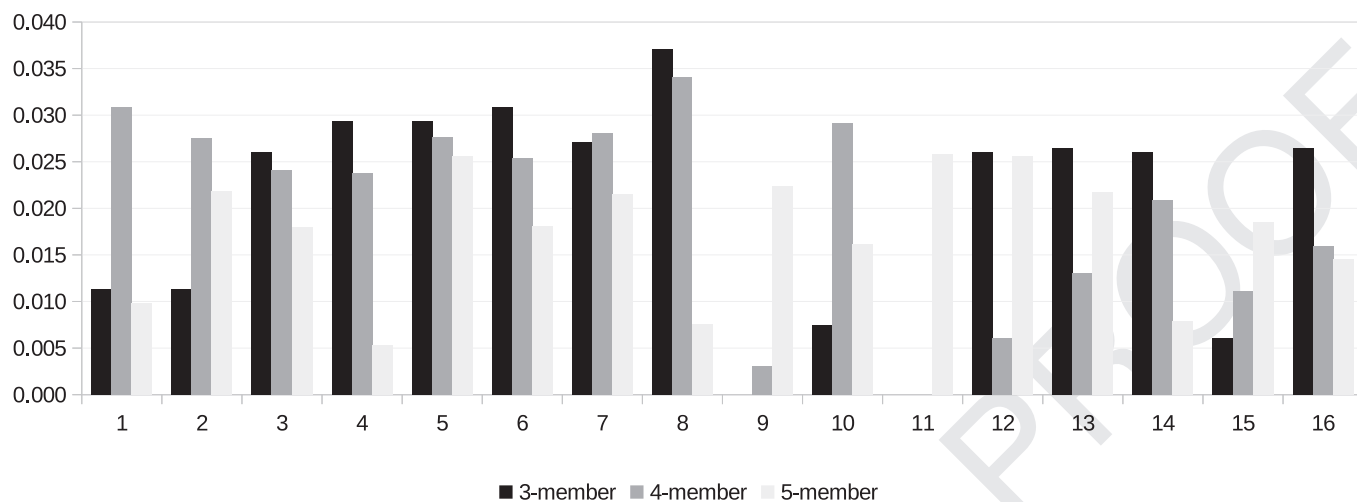


Fig. 8. Standard deviations of 3-member, 4-member, and 5-member committees for the 16 configurations.

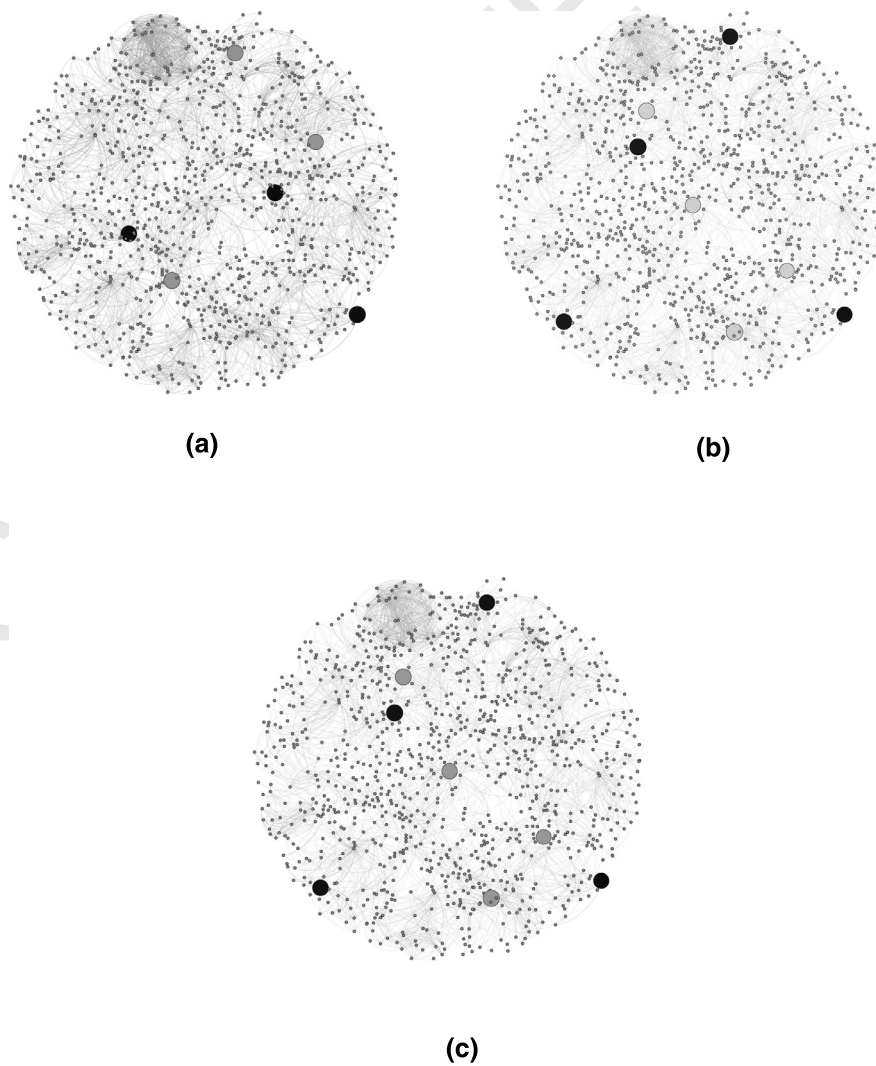


Fig. 9. Current committees members (big gray nodes) versus the best fitness committees members (big black nodes) for 3-member (a), 4-member (b), and 5-member (c) committees.

4.4. Discussion

The proposed social network is intentionally simple about tie complexity and node complexity. Here, ties are binary edges, and nodes do not have attributes considered in the committee setup. On real scenarios, other criteria could be taken into account, such as node prominence, related topic, or skill, in searching to fulfill certain requirements.

To test this approach, we used a new dataset based on public on-line available data. For simplicity, we built the social network starting from a set of specialists (those in the Informatics and Communication area), and then we created nodes and ties based on co-authorship and workplace information. To analyze other kind of specialists, social network should be built from all actors in the community or a new social network should be built starting from a new set of specialists from the area for which the committee is needed.

5. Related work

Previous work have contributed in the field of creating people committees applied to different areas, such as audit (Abbott & Parker, 2000), board directors (Shivdasani & Yermack, 1999; Westphal & Zajac, 1995), or public agencies (Loewenberg, Patterson, & Jewell, 1985). Some approaches have been focused on the diversity of the members (Aksela & Laaksonen, 2006; Hadjitodorov, Kuncheva, & Todorova, 2006; Kuncheva, 2005; Kuncheva & Whitaker, 2003; Shin & Sohn, 2005; Zouari, Heutte, & Lecourtier, 2005), while other approaches have been based different voting techniques (Bock, Day, & McMorris, 1998; Fishburn, 1981; Gehrlein, 1985). However, to the best of our knowledge there are no precedents in committee selections with independent members by using social networks.

Choosing committees with independent members in social networks can be regarded a group selection problem. Generally, this problem includes node group selection, structural consideration such as cohesion or centrality measures, and some optimization strategy since most of them are classified as NP problems.

Two well-known group selection problems in social networks are the target set selection problem and the community detection problem, however these problems present some differences with committee selection problem. The target set selection problem aims to select nodes that maximize influence in order to spread something in a network, such as information. Here, the focus is on the network, since the problem is determined by which set of nodes increase the influence. The community detection problem aims to discover node sets based on node relations or structural properties. Here, the focus is on the set and its internal structural properties, since the problem is determined by which nodes belong to a group or community.

However, committee member selection problem focuses on the group and the network, since the group considers relations between committee members and the group independence considers the whole network.

Current literature about target set selection problem shares some elements with this work. Wang, Deng, Zhou, and Jiang (2014) develop a set-based coding genetic algorithm (SGA) that converges in probability to the problem optimal solution. Here, the authors code chromosomes as sets, and choose operators based on the chromosome representation. However, SGA mainly differs with this work in the use of diffusion dynamics to measure performance. Cao, Wu, Wang, and Hu (2011) propose a transformation of the target selection problem into an optimal resource allocation problem. Here, the authors make use of the modular structural property of social networks, and propose a dynamic programming algorithm to solve the problem, which was proved to be NP-hard.

Similar to the target set selection problem is the key player problem (KPP) (Ballester, Calvó-Armengol, & Zenou, 2006; Borgatti, 2006;

Everett & Borgatti, 2010). KPP identifies key player sets with two different approaches, KPP-Neg and KPP-Pos. KPP-Neg searches for key players sets that if removed will disrupt the network. KPP-Pos searches for key players sets optimally connected to all other nodes in a network. The main difference with this work is on the structural property, since KPP-Pos uses set cohesion and KPP-Neg uses closeness centrality. Also, the authors suggest some evolutionary strategies for function optimization.

An early effort on maximizing the impact in social networks is presented in Liberman and Wolf (1997) that proposes a strategy to increase impact of information flow on scientific communities. This work has historical value, but it shows that similar problems in social networks have had different names over time.

Current literature about community detection problem shows a growing interest in topics such as social circles, topic models, or complex networks. However, there still are community detection approaches mainly based on structural properties. Bhattacharyya and Bickel (2014) use graph distances to detect communities in graphs by using a block model approach. The authors use geodesic distances which have underlying problems, such as the impossibility to measure geodesic distances in unconnected graphs. The authors solve this constraint by replacing distances of disconnected pair of nodes with the largest geodesic distance in the graph.

About the use of genetic algorithms as an optimization strategy for community detection, Freeman (1993) presents a review of the group selection problem and recognizes the computation constraint of uncovering groups based on proximity matrix representation. He also recognizes the need for a search strategy, therefore he proposes a simple genetic algorithm. The main differences with our work are in the chromosome representation and in the fitness function, which uses the proximity matrix information and a binary node classification.

As a precedent on using a structural approach to select people groups, Burt (1978) proposes a process that uses sociometric measures for sampling firm representatives of interlocking directorates to overcome profit constraints of an industry.

We found other areas that use distance as social network structural property for group selection. For instance, in the recommendation area, Hwang, Wei, and Liao (2010) suggest articles based on a co-authorship network and different schemes to measure the closeness of author sets. Here, the social network graph representation includes directed and valued ties which affect closeness measure implementation. In the social network analysis homophily area, Preciado, Snijders, Burk, Stattin, and Kerr (2012) take geographical proximity as distance in order to analyze likelihood of friendship existence and dynamics within social networks. A related approach is presented by Morgan and Carley (2011, 2014) which uses social distance as part of an impact factor set to candidate selection for hiring processes.

As another group selection approach, Wi, Mun, Oh, and Jung (2009a, 2009b) use social network structural properties along with genetic algorithms. The authors propose a quantitative method for the team member selection problem based on knowledge and collaboration of candidates. This problem aims to select teams based on abilities of candidates to fulfill project requirements and to predict team performance. Network structural properties are used to measure familiarity between candidates which is translated in what they call knowledge competence. Also, they use structural properties to select project managers from teams.

A previous work that uses geodesic paths as structural property for group selection (Kolaczyk, Chua, & Barthélemy, 2009) proposes a metric called co-betweenness, which extends betweenness centrality to sets of nodes in order to measure the information flow of the set. Co-betweenness considers the geodesic paths that pass through all nodes in the set.

487 Out of the social network scope, some works in artificial intelli-
 488 gence use a committee based concept to select other kinds of groups,
 489 such as classification (Aksela, 2003; Argamon-Engelson & Dagan,
 490 1999; Li, Zou, Hu, Wu, & Yu, 2013; Wang & Wang, 2006; Zheng, 1998)
 491 or clustering (Hadjitodorov et al., 2006; Tao, Ma, & Qiao, 2013).

492 6. Conclusions

493 A novel social network approach to the committee member selec-
 494 tion problem has been proposed. This approach consists in a mech-
 495 anism that models the problem as a social network group selection
 496 problem.

497 In this group selection problem for committee member selection,
 498 independence is the main selection criterion, for which a novel group
 499 independence function is defined. This group independence func-
 500 tion uses geodesic distances to measure social distances between all
 501 node pairs in the social network. Also, a genetic algorithm is defined
 502 to generate committee candidates. Then, the group independence
 503 function is maximized to choose candidate groups with the best
 504 fitness.

505 A case study is presented where the proposed approach is applied
 506 to a real social network. The social network was built with on-line
 507 available data extracted from a public R&D funding agency. Further,
 508 results were compared with current committees of the same agency.
 509 Results show that the proposed approach can generate committees
 510 that improve group independence compared to the current commit-
 511 tee performances.

512 Assisting committee selection processes may be the greatest com-
 513 petitive advantage offered by the proposed approach, since we have
 514 proved that the best performance groups can be selected within
 515 seconds for a real scenario. Also, alternative group selections can
 516 be preferred by experts in charge for committee appointments.
 517 Moreover, this work is built upon a simple infrastructure because
 518 there are many genetic algorithm implementations, and social net-
 519 work manipulation software, that allow the implementation and
 520 the execution of the approach in standard hardware and software
 521 configurations. As practical usage, this approach can be implemented
 522 in recommendation processes to propose alternative group selec-
 523 tions, or even group member replacements in order to improve group
 524 performances. Also, this approach can be used in opinion polls where
 525 there is a need to select less related respondents, such as focus
 526 groups.

527 Although this approach is presented as a simple alternative to the
 528 committee selection problem, there still are some limitations. These
 529 limitations include an underlying problem, which implies that the
 530 geodesic distances must be calculated between every node pair in
 531 the network. Another limitation of the geodesic distance as under-
 532 lying measure is that distance between nodes from different com-
 533 ponents cannot be determined. Also, despite the proposed genetic
 534 algorithm returns the best performance solutions, it is still an ap-
 535 proximation strategy to the global optimum. Finally, the proposed
 536 approach is intentionally designed for simple social networks with
 537 unidirectional and unvalued ties, therefore its application in other sce-
 538 narios, such as complex networks, may require some modifications.

539 Future works aim to test the proposed approach in other domains
 540 that require committee member selection. Despite this approach uses
 541 a simple network representation, more complex committee member
 542 selection processes may include criteria other than the group inde-
 543 pendence, therefore future works may include multiple criteria in
 544 group selection for the committee member selection problem. Fur-
 545 ther, other optimization strategies could be evaluated, particularly for
 546 scalability scenarios. Moreover, a complex social network representa-
 547 tion will allow to include other kinds of network properties, such as
 548 directed ties or node attributes.

References

- Abbott, L. J., & Parker, S. (2000). Auditor selection and audit committee characteristics. *AUDITING: A Journal of Practice & Theory*, 19(2), 47–66. doi:10.2308/aud.2000.19.2.47.
- Aksela, M. (2003). Comparison of classifier selection methods for improving committee performance. In *Multiple Classifier Systems* (pp. 84–93). Springer. http://link.springer.com/chapter/10.1007/3-540-44938-8_9.
- Aksela, M., & Laaksonen, J. (2006). Using diversity of errors for selecting members of a committee classifier. *Pattern Recognition*, 39(4), 608–623. doi:10.1016/j.patcog.2005.08.017.
- Argamon-Engelson, S., & Dagan, I. (1999). Committee-based sample selection for probabilistic classifiers. *Journal of Artificial Intelligence Research*, 11, 335–360. doi:10.1613/jair.612.
- Ballester, C., Calvó-Armengol, A., & Zenou, Y. (2006). Who's who in networks. wanted: the key player. *Econometrica*, 74(5), 1403–1417. doi:10.1111/j.1468-0262.2006.00709.x.
- Bhattacharyya, S., & Bickel, P. J. (2014). *Community detection in networks using graph distance* arXiv:1401.3915 [cs, stat], Netherlands.
- Bock, H.-H., Day, W. H., & McMorris, F. (1998). Consensus rules for committee elections. *Mathematical Social Sciences*, 35(3), 219–232. doi:10.1016/S0165-4896(97)00033-4.
- Borgatti, S. P. (2006). Identifying sets of key players in a social network. *Computational & Mathematical Organization Theory*, 12(1), 21–34. doi:10.1007/s10588-006-7084-x.
- Burt, R. S. (1978). A structural theory of interlocking corporate directorates. *Social Networks*, 1(4), 415–435. doi:10.1016/0378-8733(78)90006-0.
- Cao, T., Wu, X., Wang, S., & Hu, X. (2011). Maximizing influence spread in modular social networks by optimal resource allocation. *Expert Systems With Applications*, 38(10), 13128–13135. doi:10.1016/j.eswa.2011.04.119.
- Everett, M. G., & Borgatti, S. P. (2010). Induced, endogenous and exogenous centrality. *Social Networks*, 32(4), 339–344. doi:10.1016/j.socnet.2010.06.004.
- Fishburn, P. C. (1981). Majority committees. *Journal of Economic Theory*, 25(2), 255–268. doi:10.1016/0022-0531(81)90005-3.
- Freeman, L. C. (1977). A set of measures of centrality based on betweenness. *Sociometry*, 40(1), 35–41. doi:10.2307/3033543.
- Freeman, L. C. (1993). Finding groups with a simple genetic algorithm. *Journal of Mathematical Sociology*, 17(4), 227–241. doi:10.1080/0022250X.1993.9990109.
- Gehrlein, W. V. (1985). The Condorcet criterion and committee selection. *Mathematical Social Sciences*, 10(3), 199–209. doi:10.1016/0165-4896(85)90043-5.
- Hadjitodorov, S. T., Kuncheva, L. I., & Todorova, L. P. (2006). Moderate diversity for better cluster ensembles. *Information Fusion*, 7(3), 264–275. doi:10.1016/j.inffus.2005.01.008.
- Hwang, S.-Y., Wei, C.-P., & Liao, Y.-F. (2010). Coauthorship networks and academic literature recommendation. *Electronic Commerce Research and Applications*, 9(4), 323–334. doi:10.1016/j.elerap.2010.01.001.
- Kolaczyk, E. D., Chua, D. B., & Barthélemy, M. (2009). Group betweenness and co-betweenness: inter-related notions of coalition centrality. *Social Networks*, 31(3), 190–203. doi:10.1016/j.socnet.2009.02.003.
- Kuncheva, L. I. (2005). Using diversity measures for generating error-correcting output codes in classifier ensembles. *Pattern Recognition Letters*, 26(1), 83–90. doi:10.1016/j.patrec.2004.08.019.
- Kuncheva, L. I., & Whitaker, C. J. (2003). Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy. *Machine Learning*, 51(2), 181–207. doi:10.1023/A:1022859003006.
- Li, L., Zou, B., Hu, Q., Wu, X., & Yu, D. (2013). Dynamic classifier ensemble using classification confidence. *Neurocomputing*, 99, 581–591. doi:10.1016/j.neucom.2012.07.026.
- Liberman, S., & Wolf, K. B. (1997). The flow of knowledge: scientific contacts in formal meetings. *Social Networks*, 19(3), 271–283. doi:10.1016/S0378-8733(96)00303-6.
- Loewenberg, G., Patterson, S. C., & Jewell, M. E. (1985). *Handbook of legislative research* (1st ed.). Cambridge, MA: Harvard University Press.
- Morgan, G. P., & Carley, K. M. (2011). Exploring the impact of a stochastic hiring function in dynamic organizations. In *Proceedings of BRIMS* (pp. 106–113).
- Morgan, G. P., & Carley, K. M. (2014). Comparing hiring strategies in a committee with similarity biases. *Computational and Mathematical Organization Theory*, 20(1), 1–19. doi:10.1007/s10588-012-9130-1.
- Preciado, P., Snijders, T. A. B., Burk, W. J., Stattin, H., & Kerr, M. (2012). Does proximity matter? Distance dependence of adolescent friendships. *Social Networks*, 34(1), 18–31. doi:10.1016/j.socnet.2011.01.002.
- Shin, H., & Sohn, S. (2005). Selected tree classifier combination based on both accuracy and error diversity. *Pattern Recognition*, 38(2), 191–197. doi:10.1016/j.patcog.2004.06.008.
- Shivdasani, A., & Yermack, D. (1999). CEO involvement in the selection of new board members: an empirical analysis. *The Journal of Finance*, 54(5), 1829–1853. doi:10.1111/0022-1082.00168.
- Smith, J. E., & Eiben, A. E. (2008). *Introduction to evolutionary computing*. Springer.
- Tao, H., Ma, X.-p., & Qiao, M.-y. (2013). Subspace selective ensemble algorithm based on feature clustering. *Journal of Computers*, 8(2). doi:10.4304/jcp.8.2.509-516.
- Wang, C., Deng, L., Zhou, G., & Jiang, M. (2014). A global optimization algorithm for target set selection problems. *Information Sciences*, 267, 101–118. doi:10.1016/j.ins.2013.09.033.
- Wang, X., & Wang, H. (2006). Classification by evolutionary ensembles. *Pattern Recognition*, 39(4), 595–607. doi:10.1016/j.patcog.2005.09.016.
- Wasserman, S., & Faust, K. (1994). *Social network analysis: Methods and applications*. Cambridge University Press.

- 633 Westphal, J. D., & Zajac, E. J. (1995). Who shall govern? CEO/board power, demographic
634 similarity, and new director selection. *Administrative Science Quarterly*, 40(1), 60.
635 doi:10.2307/2393700.
- 636 Wi, H., Mun, J., Oh, S., & Jung, M. (2009a). Modeling and analysis of project team forma-
637 tion factors in a project-oriented virtual organization (ProVO). *Expert Systems with*
638 *Applications*, 36(3, Part 2), 5775–5783. doi:10.1016/j.eswa.2008.06.116.
- 639 Wi, H., Oh, S., Mun, J., & Jung, M. (2009b). A team formation model based on
640 knowledge and collaboration. *Expert Systems with Applications*, 36(5), 9121–9134.
641 doi:10.1016/j.eswa.2008.12.031.
- Zheng, Z. (1998). Naive Bayesian classifier committees. In *Naive Bayesian classifier* 642
643 *committees*. Springer. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.53.3003&rep=rep1&type=pdf>.
- Zouari, H., Heutte, L., & Lecourtier, Y. (2005). Controlling the diversity in classifier en- 645
646 sembles through a measure of agreement. *Pattern Recognition*, 38(11), 2195–2199.
647 doi:10.1016/j.patcog.2005.02.012.

UNCORRECTED PROOF