



# Relevance of Negative Links in Graph Partitioning: A Case Study Using Votes From the European Parliament

Israel Mendonça, Rosa Figueiredo, Vincent Labatut, Philippe Michelon

## ► To cite this version:

Israel Mendonça, Rosa Figueiredo, Vincent Labatut, Philippe Michelon. Relevance of Negative Links in Graph Partitioning: A Case Study Using Votes From the European Parliament. 2nd European Network Intelligence Conference (ENIC), Sep 2015, Karlskrona, Sweden. pp.122-129, 10.1109/ENIC.2015.25 . hal-01176090v2

**HAL Id: hal-01176090**

**<https://hal.science/hal-01176090v2>**

Submitted on 21 Jul 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - ShareAlike 4.0 International License

# Relevance of Negative Links in Graph Partitioning: A Case Study Using Votes From the European Parliament

Israel Mendonça, Rosa Figueiredo, Vincent Labatut, Philippe Michelon

Université d'Avignon, LIA EA 4128, France

Emails: {mendonci, rosa.figueiredo, vincent.labatut, philippe.michelon}@univ-avignon.fr

**Abstract**—In this paper, we want to study the informative value of negative links in signed complex networks. For this purpose, we extract and analyze a collection of signed networks representing voting sessions of the European Parliament (EP). We first process some data collected by the VoteWatch Europe Website for the whole 7<sup>th</sup> term (2009-2014), by considering voting similarities between Members of the EP to define weighted signed links. We then apply a selection of community detection algorithms, designed to process only positive links, to these data. We also apply Parallel Iterative Local Search (Parallel ILS), an algorithm recently proposed to identify balanced partitions in signed networks. Our results show that, contrary to the conclusions of a previous study focusing on other data, the partitions detected by ignoring or considering the negative links are indeed remarkably different for these networks. The relevance of negative links for graph partitioning therefore is an open question which should be further explored.

**Index Terms**—signed graphs, structural balance, graph partition, European Parliament.

## I. INTRODUCTION

In *signed* graphs, each link is labeled with a sign + or −, which indicates the nature of the relationship between the considered adjacent nodes. This type of graphs was primarily introduced in Psychology, with the objective of describing the relationship between people belonging to distinct social groups [1]. More generally, a signed graph can be used to model any system containing two types of antithetical relationships, such as like/dislike, for/against, etc. This work and its extensions by Cartwright and Harary in the 1950's [2], [3], [4] are the basis for the concept of *Structural balance*. A signed graph is considered *structurally balanced* if it can be partitioned into two [2] or more [5] mutually hostile subgroups each having internal solidarity. Here, the words *hostile* and *solidary* mean: *connected by negative* and *positive* links, respectively.

However, it is very rare for a real-world network to have a perfectly balanced structure: the question is then to quantify how balanced it is. For this purpose, one must first define a measure of balance, and then apply a method to evaluate the network balance according to this measure. For instance, one could consider counting the numbers of positive links located inside the groups, and of negative links located between them. Such a measure is expressed relatively to a graph partition, so processing the graph balance amounts to identifying the partition corresponding to the highest balance measure. In

other words, calculating the graph balance can be formulated as an optimization problem.

By using different variants of the balance measure and/or by introducing some additional constraints, one can express various versions of the notion of balance. Each one potentially leads to a different optimization problem to be solved. However, besides the very classic measures such as the one mentioned previously, only a few recent works explored this aspects from an Operations Research perspective [6], [7], [8], [9], [10]. A deep investigation of efficient approaches and mathematical formulations to problems related with signed graph balance is therefore still missing.

Independently from the Operations Research domain, the study and partition of signed graph has also recently been the object of several works in the domain of Complex Network Analysis, and more particularly community detection. The community detection problem originally concerns unsigned graphs. It consists in partitioning such a graph, in a way such that most links are located inside the groups (aka. communities) and only few remain between them. By definition, an unsigned graph focuses on a single type of relationships, say the positive ones. A signed graph representing the same system can therefore be considered as more informative, since it additionally contains the links of the other type (in our example, the negative ones). For this reason, a few authors tried to adapt existing community detection methods, in order to take advantage of this additional information [11], [12], [13], [6].

Other authors tried to study how informative these additional links really are [14]. Indeed, retrieving a signed network is a task potentially more costly than for an unsigned network, be it in terms of time, money, or methodological complications. For example, in the context of a ground survey, it is much easier to get people to name their friends than their foes. So, the question to know whether this extra cost is worth it is extremely relevant. In their work, Esmailian *et al.* [14] suggested that if one detects the communities based only on positive links (by ignoring negative ones), most negative links are already placed between the communities, and that the few ones located inside do not significantly affect the communities. The latter point is tested by checking that no additional division of the community allows increasing the overall balance. Consequently, using algorithms that do not

take negative links into consideration, such as InfoMap [15], it is possible to obtain a reasonably well partitioned network. However, we see two limitations to this work. First, in order to assess the significance of the negative links located inside the communities, Esmailian *et al.* considered each community separately, instead of the graph as a whole. Second, these observations were made only for two datasets, both representing Social Networking Services (Slashdot and Epinions), so they do not necessarily apply to all networks, or even to all types of networks.

In this paper, we want to explore further the informative value of negative links in the context of graph partitioning. To this purpose, we present a method to extract signed networks from voting data describing the activity of the *Members of the European Parliament* (MEPs). Based on this new data, we apply state-of-the-art tools in order to partition the graph, on the one hand in terms of community structure, and on the other hand according to the notion of structural balance. We then compare the obtained partitions and show the presence of significant differences between them.

The contributions of this paper are essentially two-fold. First, we constitute a new dataset of signed networks and make it publicly available to the community, with the scripts used to obtain it. We treat the voting patterns using several parameters, leading to a collection of signed networks describing the behavior of MEPs according to various modes (time, topic...). Second, based on these data, we experimentally show that negative links *can* be essential when partitioning networks. We see our work as complementary to that of Esmailian *et al.*, first because the use of a method taking negative links into account as a reference allows us to avoid the issue regarding the assessment of intra-community negative links; and second because we treat a different type of signed real-world networks, in which the links represent vote similarity instead of self-declared social relationships.

The rest of this paper is organized as follows. Section II presents a review of the literature regarding the graph partition task. Section III describes the method we used to extract signed networks from the raw data constituted of the sequences of MEPs votes. Section IV summarizes the algorithms we selected to partition our signed networks. In Section V, we present and discuss our experimental results regarding network extraction and network partition. Finally, we conclude by highlighting the main points of the article, and identifying some possible perspectives.

## II. RELATED WORKS

As mentioned before, the concepts of *signed graph* and *structural balance* were introduced by Heider [1]. Later, Cartwright *et al.* [2] formalized Heider's theory, stating that a balanced social group could be partitioned into two mutually hostile subgroups, each having internal solidarity. Observing that a social group may contain more than two hostile subgroups, Davis [5] proposed the notion of *clusterable* signed graph.

The clustering problem consists in finding the most balanced partition of a signed graph. Evaluating this balance according to the structural balance (SB) measure amounts to solving an optimization problem called *Correlation Clustering* (CC) [16]. This problem was addressed first by Doreian & Mrvar [17], who proposed an approximate solution and used it to analyze the structural balance of real-world social networks. In [11], Yang *et al.* called the CC problem *Community Mining*, and proposed an agent-based heuristic called FEC to find an approximate solution. Elsner & Schudy performed a comparison of several strategies for solving the CC problem in [18], and applied them to document clustering and natural language processing issues. In this context, these authors identified the best strategy as a greedy algorithm able to quickly achieve good objective values with tight bounds. The solution of the CC problem and of some of its variants has already been used as a criterion to measure the balance of signed social networks [19], [17], [20], [8], [6], and as a tool to identify relations contributing to their imbalance [21]. In [6], the authors provide an efficient solution of the CC problem, by the use of a ILS metaheuristic. The proposed algorithm outperforms other methods from the literature on 3 huge signed social networks. In this work, we will use this tool to evaluate the imbalance of the MEPs networks.

In the complex networks field, works dedicated to signed networks focus only on the clustering problem, as defined by Davis [5]. Various methods were proposed for this purpose: evolutionary approaches [22], [23], [24], [25], agent-based systems [11], matrix transformation [26], extensions of the Modularity measure [27], [9], [28], [29], [30], simulated annealing [31], spectral approaches [32], [33], [34], particle swarm optimization [35], [36], and others. Some authors performed the same task on bipartite networks [37], while others relaxed the clustering problem in order to identify overlapping communities [38]. Although the methods listed here were applied to networks representing very different systems, authors did not investigate the possibility that some alternative versions of the clustering problems were more appropriate to certain data.

Few works tried to compare the CC and community detection approaches. As mentioned in the introduction, Esmailian *et al.* showed that, in certain cases, partitions estimated in signed networks by community detection methods, i.e. based only on the positive links, can be highly balanced [14]. However, this work was conducted only on two networks of self-declared social interaction networks (Epinions and Slashdot), and using a single community detection method (InfoMap [15]). Moreover, they did not compare their results to partitions detected by algorithms designed to solve the CC problem. We investigate if this statement also holds for other real-world networks and community detection methods, and how these compare to results obtained with CC methods.

## III. NETWORK EXTRACTION

In this section, we describe the source we used to retrieve our raw data, and the process we applied to extract signed

networks from these data.

#### A. European Parliament Votes

In order to be able to conduct our experiments, we were looking for data allowing to extract some form of signed network of interactions. Moreover, in future works, we want to study how the network and the structural balance evolve, so the data had to be longitudinal, with stable nodes (nodes should not change too much through time). The best data we could find relatively to these criteria are those describing the activity of the *European Parliament*<sup>1</sup>. More precisely, we focused on the votes of the Members of the European Parliament (MEPs).

The Website *VoteWatch Europe*<sup>2</sup> is a non-partisan international non-governmental organization, completely independent from national and local governments, from the European Union, as well as from political parties, institutions, agencies, businesses and all other bodies. Their goal is to provide easy access to the votes and other activities of the European parliament (among other European institutions). Votewatch compiles data provided by the EP to give a full overview of the MEPs activity. In particular, they describe the vote cast by each MEP for each document considered at the EP. Each MEP is also described through his name, country and political group. Other fields are available too, which we have not used yet, such as how loyal the MEP is to his political group. To summarize, the behavior of a MEP is represented by the series of votes he cast over a certain time period (e.g. a year, a term).

For a given document, a MEP can express his vote in one of the three following ways: FOR (the MEP wants the document to be accepted), AGAINST (he wants the document to be rejected) and ABSTAIN (he wants to express his neutrality). Besides these expressed votes, it is also possible for the MEP not to vote at all, leading to the following possibilities: ABSENT (the MEP was not present during the vote), DID NOT VOTE (he was there, but did not cast his vote), and DOCUMENTED ABSENCE (he was not there but justified his absence).

For each document, we also have access to the category it belongs to, called *Policy*. It corresponds roughly to the main theme treated in the considered document. All the policies treated during the 7th term of the EP are listed in Table I, with the numbers of documents they concern.

VoteWatch gives us access to raw data, which could be described as individual data, in the sense they describe the state and behavior of the MEPs when considered independently from each others. However, a network is by nature relational, i.e. it represents the relationships between some objects of interest. Thus, we need to process the VoteWatch data in order to retrieve the networks we want.

#### B. Extraction Process

Our extraction process is two-stepped. As mentioned before, in the data received from VoteWatch, the behavior of each MEP is represented by a series of votes, corresponding to

TABLE I  
LIST OF ALL POLICIES RELATIVE TO THE DOCUMENTS VOTED AT THE EUROPEAN PARLIAMENT, WITH THE CORRESPONDING NUMBERS OF DOCUMENTS, FOR THE 7TH TERM

| Policy   | Number of documents |
|--|---------------------|
| Agriculture                                    | 53                  |
| Budget   | 179                 |
| Budgetary control                              | 113                 |
| Civil liberties, justice & home affairs        | 99                  |
| Constitutional and inter-institutional affairs | 40                  |
| Culture & education                            | 19                  |
| Development                                    | 29                  |
| Economic & monetary affairs                    | 128                 |
| Employment & social affairs                    | 44                  |
| Environment & public health                    | 100                 |
| Fisheries                                      | 53                  |
| Foreign & security policy                      | 191                 |
| Gender equality                                | 28                  |
| Industry, research & energy                    | 51                  |
| Internal market & consumer protection          | 39                  |
| Internal regulations of the EP                 | 7                   |
| International trade                            | 106                 |
| Legal affairs                                  | 67                  |
| Petitions                                      | 5                   |
| Regional development                           | 35                  |
| Transport & tourism                            | 40                  |

all the documents reviewed by the EP during one term. In this article, we focused on the 7th term (from june 2009 to june 2014). We first filter these data depending on temporal and topical criteria. In other terms, if required, it is possible to focus only on the documents related to a specific policy and/or a specific period of the term, for instance a given year. The second step consists in comparing individually all MEPs in terms of similarity of their voting behaviors. The result of this process is what we call the *agreement* matrix  $M$ . Each numerical value  $m_{uv}$  contained in the matrix represents the average agreement between two MEPs  $u$  and  $v$ , i.e. how similarly they vote over all considered documents.

The filtering step is straightforward, however the agreement processing constitutes a major methodological point: depending on how it is conducted, it can strongly affect the resulting network. For a pair of MEPs  $u$  and  $v$  and a given document  $d_i$ , we define the *document-wise agreement score*  $m_{uv}(d_i)$  by comparing the votes of both considered MEPs. It ranges from  $-1$  if the MEPs fully disagree, i.e. one voted FOR and the other AGAINST, to  $+1$  if they fully agree, i.e. they both voted FOR or AGAINST.

However, as we mentioned previously, a vote can take other values than just FOR and AGAINST, and those must also be treated. Let us consider first the non-expressed votes: ABSENT, DID NOT VOTE and DOCUMENTED ABSENCE. The EU distinguishes these different forms of absence not for political, but rather for administrative reasons, so we decided to consider them all simply as absences. The common approach when treating this type of vote data [39], [40] is to ignore all documents  $d_i$  for which at least one of the considered MEPs was absent. However, certain MEPs are absent very often, which mean they would share a very small number of documents with others. This approach could therefore artificially produce

<sup>1</sup><http://www.europarl.europa.eu/>

<sup>2</sup><http://www.votewatch.eu/>



extremely strong agreement or disagreement scores. To avoid this, we assign a neutral score of 0 when at least one person is absent for a given document.

Handling the abstentions is a bit trickier, because such a behavior can mean different things. A MEP can choose not to vote because he is personally FOR or AGAINST, but undergoes some pressure (from his political group, his constituents, etc.) to vote the other way: in this case, voting ABSTAIN is a way of expressing this conflicting situation. But abstaining could also simply represent neutrality, meaning the MEP is neither FOR nor AGAINST the considered document. There is no consensus in the literature, and different approaches were proposed to account for ABSTAIN-FOR, ABSTAIN-AGAINST and ABSTAIN-ABSTAIN situations [29], [39], [40]. Here, we present two variants corresponding to different interpretations. In the first, described in Table II, an abstention is considered as half an agreement with FOR or AGAINST, leading to a score of +0.5. In the second, described in Table III, two abstaining MEPs are considered to fully agree (+1). But, when only one of them abstains, we consider there is not enough information to determine whether they agree or disagree, and we therefore use a 0 score. Note absences were left out of the tables for clarity.

TABLE II  
VOTE WEIGHTS REPRESENTING ABSTENTION AS HALF AN AGREEMENT

|         | FOR  | ABSTAIN | AGAINST |
|---------|------|---------|---------|
| FOR     | +1   | +0.5    | -1      |
| ABSTAIN | +0.5 | +0.5    | +0.5    |
| AGAINST | -1   | +0.5    | +1      |

TABLE III  
VOTE WEIGHTS REPRESENTING ABSTENTION AS AN ABSENCE OF OPINION

|         | FOR | ABSTAIN | AGAINST |
|---------|-----|---------|---------|
| FOR     | +1  | 0       | -1      |
| ABSTAIN | 0   | +1      | 0       |
| AGAINST | -1  | 0       | +1      |

The document-wise agreement score is completely defined by selecting one of the proposed tables. The average agreement is then obtained by averaging this score over all considered documents. More formally, let us consider two users  $u$  and  $v$  and note  $d_1, \dots, d_\ell$  the documents remaining after the filtering step, and for which  $u$  and  $v$  both cast their votes. The *average agreement*  $m_{uv}$  between these two MEPs is:

$$m_{uv} = \frac{1}{\ell} \sum_{i=1}^{\ell} m_{uv}(d_i) \quad (1)$$

#### IV. PARTITION METHODS

In this section, we present the methods used to partition the signed network extracted from the VoteWatch data. We first introduce the community detection approaches we selected for our experiment. Then we formally define the Correlation Clustering problem and describe the algorithm we used in this article to estimate its solutions.

##### A. Community Detection

In the literature, the problem of community detection is usually defined in an informal way. It consists in finding a partition of the node set of a graph, such that many links lie inside the parts (communities) and few lie in-between them. An other way of putting it is that we are looking for groups of densely interconnected nodes, relatively to the rest of the network [41]. It is difficult to find a formal definition of this problem, or rather, to find a *unique* formal definition: many authors present and solve their own variant. Because of this, the algorithms presented in the literature do not necessarily solve the exact same problem, although it is still named community detection. To account for this variance, we selected several methods for our experiments, which we briefly present here. All of them are able to process weighted networks.

**InfoMap** [15]. The community structure is represented through a two-level nomenclature based on Huffman coding: one level to distinguish communities in the network and the other to distinguish nodes in a community. The problem of finding the best community structure is expressed as minimizing the quantity of information needed to represent some random walk in the network using this nomenclature. With a partition containing few intercommunity links, the walker will probably stay longer inside communities; therefore only the second level will be needed to describe its path, leading to a compact representation. The authors optimize their criterion using simulated annealing.

**EdgeBetweenness** [42]. This divisive hierarchical algorithm adopts a top-down approach to recursively split communities into smaller and smaller node groups. The split is performed by iteratively removing the most central link of the network. This centrality is expressed in terms of edge-betweenness, i.e. number of shortest paths running through the considered link. The idea behind this method is that links connecting different communities tend to be present in the many shortest paths connecting one community to the other. Once the network has been split in two separate components, each one is split again applying the same process, and so on. The resulting components correspond to communities in the original network.

**WalkTrap** [43]. To the contrary of EdgeBetweenness, this is an agglomerative hierarchical algorithm, which means it uses a bottom-up approach to merge communities into larger and larger groups, starting from singletons. To select the communities to merge, WalkTrap uses a random walk-based distance. Indeed, random walkers tend to get trapped into communities, because most locally available links lead to nodes from the same communities, while only a few links all to escape this community (by definition). If two nodes  $u$  and  $v$  are in the same community, the probability to reach, through a random walk, a third node located in the same community should not be very different for  $u$  and  $v$ . The distance is constructed by summing these differences over all nodes, with a correction for degree.

**FastGreedy** [44]. Like WalkTrap, this algorithm adopts an agglomerative hierarchical approach. But this time, the

merges are not decided using a distance measure, but rather by locally optimizing the well-known objective function called *Modularity* [45]. Briefly, this measure compares the proportion of intra-community links present in the network of interest, to the expectation of the same quantity for a randomly generated network of similar size and degree distribution. The process stops when it is not possible to improve the modularity anymore, or when there is no more communities to merge.

### B. Correlation Clustering

Before formally describing the CC problem, we need to introduce some notations and definitions first. Let  $G = (V, E, s, w)$  be a weighted undirected signed graph. The sets  $V$  and  $E$  correspond to the nodes and links constituting the graph. The functions  $s : E \rightarrow \{+, -\}$  and  $w : E \rightarrow [0; +1]$  assign a sign and a positive weight to each link in  $E$ , respectively.

A link  $e \in E$  is called *negative* if  $s(e) = -$  and *positive* if  $s(e) = +$ . Let  $E^- \subset E$  and  $E^+ \subset E$  denote the sets of negative and positive links in  $G$ , respectively. Notice that, according to the above definitions,  $E = E^- \cup E^+$ . We define the negative and positive subgraphs of  $G$  as  $G^- = (V, E^-)$  and  $G^+ = (V, E^+)$ , respectively. The *complementary negative graph* is  $\overline{G^-} = (V, \overline{E^-})$ , where  $\overline{E^-} = \mathcal{P}_2(V) \setminus E^-$ ,  $\mathcal{P}_2(V)$  being the set of all unordered pairs from  $V$ .

Let us consider a partition  $P$  of  $V$  such that  $P = \{V_1, \dots, V_k\}$ . A link is said to be *cut* if it connects nodes from two different parts. We note  $E[V_i : V_j] \subset E$  the set of links connecting two nodes from  $V_i$  and  $V_j$  (cut links), and  $E[V_i] \subset E$  the set of links connecting two nodes from  $V_i$  (so,  $E[V_i] = E[V_i : V_i]$ ) (uncut links).

As mentioned before, negative links located inside parts (uncut negative links) and positive links located between parts (cut positive links) are considered to lower the graph balance. For  $V_i$ , the total weight of *uncut negative links* is:

$$\Omega^-(V_i) = \sum_{e \in E^- \cap E[V_i]} w_e \quad (2)$$

And for two parts  $V_i$  and  $V_j$ , the total weight of *cut positive links*  $\Omega^+$  is:

$$\Omega^+(V_i, V_j) = \sum_{e \in E^+ \cap E[V_i : V_j]} w_e \quad (3)$$

The *Imbalance*  $I(P)$  of a partition  $P$  can be defined as the sum of uncut negative and cut positive links over the whole graph:

$$I(P) = \sum_{1 \leq i \leq k} \Omega^-(V_i) + \sum_{1 \leq i < j \leq k} \Omega^+(V_i, V_j) \quad (4)$$

Finally, the *Correlation Clustering problem* is the problem of finding a partition  $P$  of  $V$  such that the *imbalance*  $I(P)$  is minimized.

In this work we will solve the CC problem using the *Parallel ILS algorithm* presented in [6], which was designed to solve the CC problem in large real-world networks. ILS is itself a metaheuristic approach allowing to obtain good quality solutions by applying iteratively greedy search methods

[46]. Starting from an initial solution estimated through a greedy method, the general principle is two-stepped: first, some perturbations are introduced to modify the current best solution; second, some local searches are performed to find better solutions within the neighborhood. This iterative process is stopped when some condition is met (minimal quality, time limit, etc.). This specific implementation is parallelized, in order to improve speed.

Considering that the networks extracted from the VoteWatch data are very dense, we had to perform some minor modifications on the original Parallel ILS algorithm, so that the processing time was acceptable. First, the search space used in the local search was reduced by adding a probably (0.7) of visiting a neighbor solution. In other terms, in average we limit the search to only a part of the neighborhood. Second, the perturbation level had to be reduced to 15, half of the maximum run number in the original work.

## V. RESULTS AND DISCUSSION

In this section, we first describe the networks extracted from the VoteWatch data, and how they are affected by the parameters controlling this extraction process. Then, we discuss the results obtained with the partition methods presented in section IV.

In order to process the VoteWatch data, we developed a tool called *NetVotes*, which takes the form of a collection of R scripts. It implements the method described in section III-B, and additionally calculates some metrics describing the studied networks and their partitions. It is generic enough to treat any type of data of the same form. To perform the community detection, we used the *igraph* R package, which contains all the algorithms we selected. For the CC problem, we used the author's version of Parallel ILS, which we modified as explained in section IV-B. All our source code, as well as the data it outputs, are publicly available on GitHub<sup>3</sup> and FigShare<sup>4</sup>, respectively.

### A. Networks Extraction

As described in section III-B, our extraction method takes three parameters: the table used to process the agreement scores, the policy and the time period. We proposed 2 different tables, there are 21 policies and we also considered all documents independently from their policies, and we considered each year separately as well as the whole 5-year long 7<sup>th</sup> term (2009-2014). This amounts to a total of 264 different modalities. However, in certain cases, the filtering step led to less than 2 documents, so we were not able to extract networks for all combinations of policies and time periods.

We first study how the choice of the table used to process the agreement scores affects the extracted network. Figure 1 shows the average agreement distribution for Table II (top plot) and Table III (bottom), using all documents for the whole term (i.e., not applying any filter). Both distributions are very similar, with a clear separation between the negative

<sup>3</sup><https://github.com/CompNet/NetVotes/>

<sup>4</sup>[http://figshare.com/articles/NetVotes\\_Data/1456268](http://figshare.com/articles/NetVotes_Data/1456268)

and positive values. The agreement side is bimodal, with larger frequencies around 0 and 0.6–0.7. The right peak can be explained by the fact the majority of MEPs tend to vote similarly most of the time. The other peak, located at zero, is due to the frequent absence of a certain number of MEPs. When a MEP is absent for a given document, his agreement score with all other MEPs is zero. It is not as contrasted on the disagreement side, with a much flatter distribution. Moreover, there is no strong disagreement since the smaller values are around  $-0.5$  (by comparison, the agreement values can get close to 0.9). This means only a small proportion of MEPs systematically disagree with the rest of the EP.

The same observations can be made when considering the different policies independently, as well as when considering each year separately. There are some variations in terms of position and amplitude of the right peak, but this is mainly due to large differences in the number of documents discussed for each policy. The nature of the table used to process the agreement scores does not seem to have any clear effect on the average agreement distribution. We additionally tried to use a variant of Table II, replacing the  $+0.5$  (half-agreement) by  $-0.5$  (half-disagreement) for the situations involving ABSTAIN vs. FOR or AGAINST. The results were extremely similar, confirming our observation. Consequently, in the rest of the article, we present only the results obtained with Table III.

### B. Partition Comparison

We now want to study how the selected community detection methods behave on the CC problem, when compared to Parallel ILS, an algorithm specifically designed to treat this problem. However, as mentioned in Section IV-A, these methods can only take positive links into account, so they cannot be applied directly to our signed networks, unlike Parallel ILS. To solve this issue, we proposed to consider two subgraphs of the original signed networks: the signed graph and the complementary negative graph, noted  $G^+$  and  $G^-$  in Section IV-B, respectively. The former is a version of the original graph retaining only its positive links. The latter contains all possible links but the ones labeled negative in the original graph. In both cases, the result is a graph with only one type of unlabeled links, representing a part of the information originally conveyed by the original graph. This is very consistent with our objective, since we want to study if the information loss translates in terms of detected partitions.

We applied all the selected community detection algorithms to both types of graphs, for all the modalities described in the previous subsection. For space matters, it is not possible to display and comment all of them, so we decided to focus on two policies of interest, over all years and over the whole term. We picked *Foreign & Security Affairs* because it is the most frequent, with 191 documents, and *Agriculture* because it is also well represented (51 document) and topically very distant from the former.

The obtained results are shown in Figure 2. The top plot is dedicated Agriculture and the bottom one to Foreign &

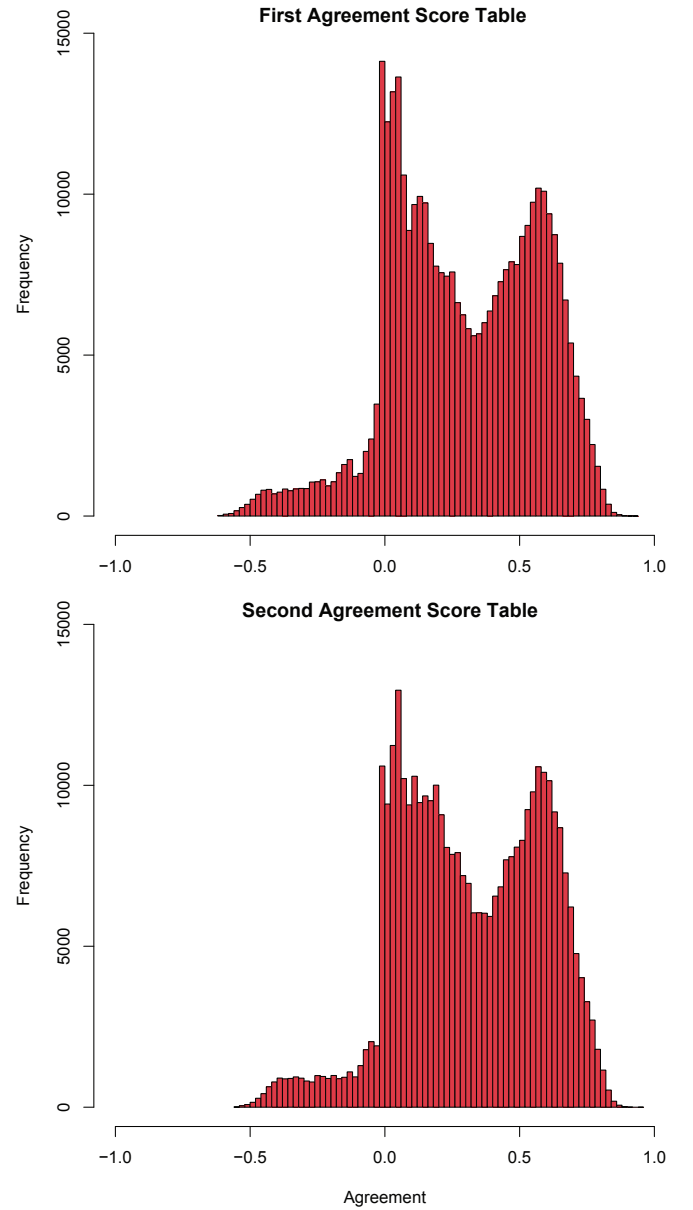


Fig. 1. Agreement distribution for the whole term and all policies when using Table II (top) and Table III (bottom)

Security Affairs. Each group of bars represents the results obtained by one algorithm for each year taken independently, and for the whole term (see the legend). The bar heights are proportional to the imbalance of the estimated partitions, as described in equation (4), only they are expressed in terms of percents relatively to  $|E|$ . The numbers on top of the bars indicate how many parts (communities) the corresponding partitions contain. Note the displayed results are representative of the other policies.

Let us compare the algorithms performances. EdgeBetweenness, FastGreedy and WalkTrap are far from finding optimal results when processing the positive subgraphs: they obtain scores ranging from 20% to more than 60% imbalance, and

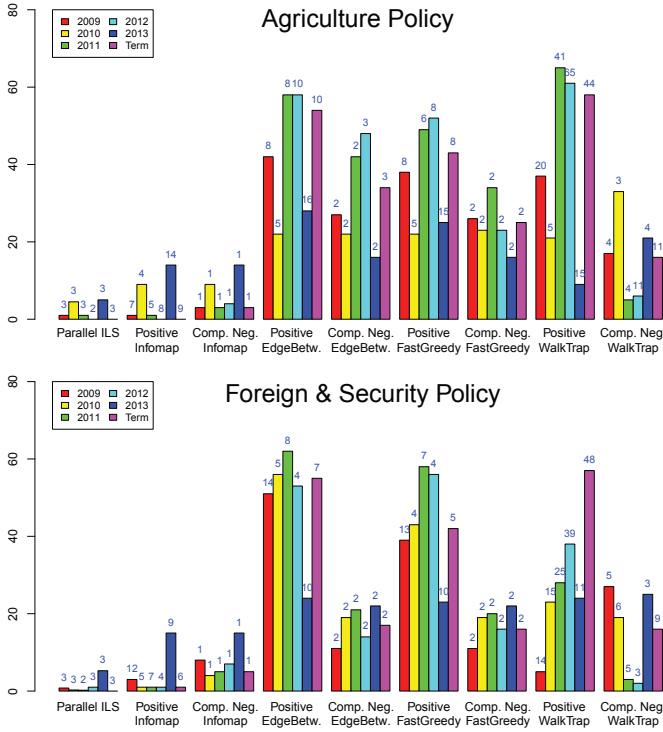


Fig. 2. Imbalance of the partitions (bars) and numbers of detected clusters (blue values), obtained through Parallel ILS (left bar group) and community detection methods (other bar groups), for each year and the whole term (see legend), processed for the *Agriculture* (top plot) and *Foreign & Security* (bottom plot) policies

generally find a high number of clusters. The multitude of clusters is certainly the cause for these large imbalances. Note this observation is not inconsistent with being efficient at detecting communities, since this task implies taking link density into account. The behavior of the same algorithms is very different when applied to the complementary negative subgraph. The number of detected clusters is much smaller (generally around 2–5), and the imbalance is smaller, but still around 20%. The reason for that is certainly that the graphs being much denser, it becomes harder to distinguish dense subgroups, i.e. communities.

The InfoMap algorithm is much more successful at detecting balanced partitions, and reaches much smaller imbalance than the other community detection algorithms (always less than 20%, often less than 5%). However, on the negative complementary graphs, InfoMap simply puts all the nodes in the same cluster, so these results cannot be considered as relevant. On the positive graphs, the imbalance is very low (with the exception of the year 2013), close to 1%, and the algorithm finds 4–14 clusters. The results obtained with Parallel ILS are even better, in terms of imbalance, since they consistently get close to 0%. Moreover, the number of clusters is relatively low (2–3), which corresponds to what we were expecting *a priori*. Indeed, the EP is known to be split in two major political sides (EPP and S&P), with some punctual alliances of smaller parties, leading to the formation of third

or fourth groups. It is worth noticing that the imbalance is more marked for both algorithms for the year 2013, for both considered policy. This might be due to this year being the last in the 7<sup>th</sup> term, and therefore coinciding with the negotiation of the 8<sup>th</sup> term budgets and changes in the policies orientation. For instance, the CAP (Common Agricultural Policy) was made greener<sup>5</sup>. Such changes lead to stronger discussions in the EP, and may challenge the balance of certain political groups.

In average, InfoMap identifies partitions 3 times more imbalanced than Parallel ILS and also tends to partition the graph in more clusters. Table IV compares the InfoMap and Parallel ILS partitions in terms of *Normalized Mutual Information*, which is the standard measure to compare partitions in the domain of unsupervised classification [47]. This measure ranges from  $-1$  (completely different) to  $+1$  (completely identical), whereas 0 represents statistical independence. The values obtained for both considered policies, and for all the time periods, are extremely close to zero. This means the partitions detected by the two algorithms have little in common, even though their number of clusters and/or imbalance level are sometimes similar.

TABLE IV  
COMPARISON OF THE INFOMAP AND PARALLEL ILS PARTITIONS IN TERMS OF NMI

| Policy          | 2009 | 2010 | 2011 | 2012 | 2013 | Term |
|-----------------|------|------|------|------|------|------|
| Agriculture     | 0.01 | 0.04 | 0.01 | 0.02 | 0.01 | 0.02 |
| Foreign Affairs | 0.01 | 0.01 | 0.01 | 0.01 | 0.03 | 0.02 |

We can conclude by stating that, on these data, our results do not confirm the findings of Esmailian *et al.* [14] regarding the low informative value of negative links. Taking negative links into account leads to a lower imbalance and a different partition, containing larger clusters. Moreover, among our selection of community detection algorithms, InfoMap is the only one to exhibit a behavior comparable to that of Parallel ILS. This means the notion of community implemented in this algorithm, which relies on an information compression-based approach, can be considered as compatible enough with the concept of structural balance. However, this is not the case for the other considered methods, based on link centrality, node distance and modularity. Discussing collectively the different methods proposed to solve the community detection problem might not be relevant, since the notions of community they rely upon are different (despite a common name).

## VI. CONCLUSION

In this article, we have investigated some of the aspects inherent to the partition of signed networks, using data from the European Parliament (EP). We first extracted a collection of networks using the voting patterns of the Members of the EP. Then, we applied a selection of community detection methods to these networks, as well as Parallel ILS, an algorithm

<sup>5</sup>[http://www.europarl.europa.eu/pdfs/news/expert/infopress/20131118IPR25538/20131118IPR25538\\_en.pdf](http://www.europarl.europa.eu/pdfs/news/expert/infopress/20131118IPR25538/20131118IPR25538_en.pdf)



specifically designed to treat signed graphs. Among the former, the best results in terms of structural balance are obtained, by far, by InfoMap. However, in average, Parallel ILS detected partitions three times more balanced. This seems to be due to the fact community detection methods ignore negative links and focus instead on link density. Independently from the balance aspect, the number of clusters detected by ILS is lower, which is more consistent with the studied system.

These results are in opposition with the finding of Esmailian *et al.* [14], however they do not invalidate them. Indeed, in both cases, the experiments were performed on a very limited number of networks. The process should be conducted on a large number of different datasets in order to draw more reliable conclusions. In our future work, we plan to constitute a collection of real-world signed networks in order to perform this task. We also want to continue studying the MEPs voting data in further details, focusing on the interpretation of the identified balanced clusters.

## REFERENCES

- [1] F. Heider. Attitudes and cognitive organization. *Journal of Psychology*, 21:107–112, 1946.
- [2] D. Cartwright and F. Harary. Structure balance: A generalization of heider's theory. *Psychological Review*, 63(5):277–293, 1956.
- [3] J. A. Davis. Structural balance, mechanical solidarity, and interpersonal relations. *American Journal of Sociology*, 68:444–462, 1963.
- [4] F. Harary. On the notion of balance of a signed graph. *Michigan Math. Journal*, 2(2):143–146, 1953.
- [5] J. A. Davis. Clustering and structural balance in graphs. *Human Relations*, 20:181–187, 1967.
- [6] M. Levorato, L. Drummond, Y. Frota, and R. Figueiredo. An ils algorithm to evaluate structural balance in signed social networks. In *Symposium on Applied Computing*, pages 1117–1122, 2015.
- [7] R. Figueiredo and Y. Frota. The maximum balanced subgraph of a signed graph: Applications and solution approaches. *European Journal of Operational Research*, 236:473–487, 2014.
- [8] R. Figueiredo and G. Moura. Mixed integer programming formulations for clustering problems related to structural balance. *Social Networks*, 35:639–651, 2013.
- [9] J. Bruggeman, V. A. Traag, and J. Uitermark. Detecting communities through network data. *American Sociology Review*, 77(6):1050–1063, 2012.
- [10] M. Brusco and D. Steinly. Integer programs for one- and two-model blockmodeling based on prespecified image matrices for structural and regular equivalence. *Journal of Mathematical Psychology*, 53:577–585, 2009.
- [11] B. Yang, W. K. Cheung, and J. Liu. Community mining from signed social networks. *IEEE Transactions on Knowledge and Data Engineering*, 19(10):1333–1348, 2007.
- [12] J. Leskovec, D. Huttenlocher, and J. Kleinberg. Predicting positive and negative links in online social networks. *Social Networks*, 31:1–11, 2009.
- [13] P. Doreian and A. Mrvar. Partitioning signed social networks. In *ACM International Conference on World Wide Web*, pages 641–650, 2010.
- [14] P. Esmailian, S. E. Abtahi, and M. Jalili. Mesoscopic analysis of online social networks - the role of negative ties. *Phys. Rev. E*, 90:042817, 2014.
- [15] M. Rowvall and C. T. Bergstrom. Maps of random walks on complex networks reveal community structure. *PNAS*, 105(4):1118–1123, 2008.
- [16] N. Bansal, A. Blum, and S. Chawla. Correlation clustering. In *43rd IEEE FOCS*, pages 238–250, 2002.
- [17] P. Doreian and A. Mrvar. A partitioning approach to structural balance. *Social Networks*, 18:149–168, 1996.
- [18] M. Elsner and W. Schudy. Bounding and comparing methods for correlation clustering beyond ilp. In *Workshop on Integer Linear Programming for Natural Language Processing*, pages 19–27, 2009.
- [19] P. Doreian. A multiple indicator approach to blockmodeling signed networks. *Social Networks*, 30:247–258, 2008.
- [20] P. Doreian and A. Mrvar. Partitioning signed social networks. *Social Networks*, 31:1–11, 2009.
- [21] P. Abell and M. Ludwig. Structural balance: a dynamic perspective. *Journal of Mathematical Sociology*, 33:129–155, 2009.
- [22] Y. Li, J. Liu, and C. Liu. A comparative analysis of evolutionary and memetic algorithms for community detection from signed social networks. *Soft Computing*, 18(2):329–348, 2014.
- [23] C. Liu, J. Liu, and Z. Jiang. A multiobjective evolutionary algorithm based on similarity for community detection from signed social networks. *IEEE Transactions on Cybernetics*, 42(12):2274–2287, 2014.
- [24] B. Yang. Self-organizing network evolving model for mining network community structure. *Lecture Notes in Computer Science*, 4093:404–415, 2006.
- [25] Y. Zeng and J. Liu. Community detection from signed social networks using a multi-objective evolutionary algorithm. *Proceedings in Adaptation, Learning and Optimization*, 1:259–270, 2015.
- [26] B. Yang and D.-Y. Liu. A heuristic clustering algorithm for mining communities in signed networks. *Journal of Computer Science and Technology*, 22(2):320–328, 2007.
- [27] A. Amelio and C. Pizzuti. Community mining in signed networks: A multiobjective approach. In *IEEE-ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 95–99, 2013.
- [28] S. Gomez, P. Jensen, and A. Arenas. Analysis of community structure in networks of correlated data. *Phys. Rev. E*, 80:016114, 2009.
- [29] K. T. Macon, P. J. Mucha, and M. A. Porter. Community structure in the united nations general assembly. *Physica A*, 391(1-2):343–361, 2012.
- [30] V. A. Traag and J. Bruggeman. Community detection in networks with positive and negative links. *Phys. Rev. E*, 80:036115, 2008.
- [31] P. Bogdanov, N. D. Larusso, and A. Singh. Towards community discovery in signed collaborative interaction networks. In *IEEE ICDM SIASP*, pages 288–295, 2010.
- [32] P. Anchuri and M. Magdon-Ismael. Communities and balance in signed networks: A spectral approach. In *International Conference on Advances in Social Networks Analysis and Mining*, pages 259–270, 2012.
- [33] J. Kunegis, S. Schmidt, A. Lommatzsch, J. Lerner, E. W. de Luca, and S. Albayrak. Spectral analysis of signed graphs for clustering, prediction and visualization. In *SDM*, pages 559–570, 2010.
- [34] L. Wu, X. Ying, X. Wu, A. Lu, and Z.-H. Zhou. Spectral analysis of k-balanced signed graphs. In *15th Pacific-Asia Conference on Knowledge Discovery and Data Mining*, volume 6635, pages 11–12, 2011.
- [35] Q. Cai, M. Gong, B. Shen, L. Ma, and L. Jiao. Discrete particle swarm optimization for identifying community structures in signed social networks. *Neural Networks*, 58:4–13, 2014.
- [36] M. Gong, Q. Cai, X. Chen, and L. Ma. Complex network clustering by multiobjective discrete particle swarm optimization based on decomposition. *IEEE Transactions on Evolutionary Computation*, 18(1):82–97, 2013.
- [37] A. Mrvar and P. Doreian. Partitioning signed two-mode networks. *Journal of Mathematical Sociology*, 33(3):196–221, 2009.
- [38] X. L. Chen, Y. Wang, B. Yuan, and B. Z. Tang. Overlapping community detection in networks with positive and negative links. *Journal of Statistical Mechanics*, 03:P03021, 2014.
- [39] M. A. Porter, P. J. Mucha, M. E. J. Newman, and C. M. Warmbrand. A network analysis of committees in the u.s. house of representatives. *PNAS*, 102(20):7057–7062, May 2005.
- [40] C. dal Maso, G. Pompa, M. Puliga, G. Riotta, and A. Chessa. Voting behavior, coalitions and government strength through a complex network analysis. *PLoS ONE*, 9(12):e116046, 2014.
- [41] S. Fortunato. Community detection in graphs. *Physics Reports*, 486(3-5):75–174, February 2010.
- [42] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Phys. Rev. E*, 69:026113, 2003.
- [43] P. Pons and M. Latapy. Computing communities in large networks using random walks. *Lecture Notes in Computer Science*, 3733:284–293, 2005.
- [44] A. Clauset, M. E. J. Newman, and C. Moore. Finding community structure in very large networks. *Phys. Rev. E*, 70:066111, 2004.
- [45] M. E. J. Newman. Modularity and community structure in networks. *PNAS*, 103(23):8577–8582, 2006.
- [46] H. R. Loureno, O. Martin, and T. Sttze. Iterated local search: Framework and applications. In *Handbook of Metaheuristics*, volume 146 of *International Series in Operations Research & Management Science*, pages 363–397. Kluwer Academic Publishers, 2010.
- [47] A. L. N. Fred and A. K. Jain. Robust data clustering. In *Computer Vision and Pattern Recognition*, volume 2, pages 128–133, 2003.