

How good is the Shapley value-based approach to the influence maximization problem?

Kamil Adamczewski^{1,3} and Szymon Matejczyk^{2,3} and Tomasz P. Michalak⁴

1 INTRODUCTION

This paper studies the process of information diffusion over a network, where new nodes as time passes become “informed” (or equivalently, influenced or infected), and a related problem of influence maximization that concerns finding the set of initial nodes that will lead to the most widespread effect of diffusion in the network [4, 6]. Influence maximization is an NP-hard problem [4], and therefore, literature has focused on approximating the optimal solution with greedy algorithms and various heuristics [4, 2, 3].

In this spirit, recent research [5, 1] proposed the Shapley value, a concept borrowed from game theory, as a measure of node centrality. The key idea is to define a cooperative game over a network in which players are nodes, coalitions are groups of nodes, and payoffs of coalitions depend on how many nodes these coalitions can infect as a group. Next, a game-theoretic solution concept—most often the Shapley value—quantifies a role or importance of a player (e.g. in the diffusion process) by its average marginal contribution in the game. In this paper, we also briefly test the Banzhaf index (BI) as an alternative solution concept.

Nevertheless, unlike other centrality measures, the game-theoretic centrality for influence maximization has not been yet rigorously evaluated. The only experimental work is that of Ramasuri and Narahari [5] who proposed the SPIN algorithm built upon the Shapley value. This evaluation is, however, very preliminary, as it relies on the approximated Shapley value and is not scalable even for the graphs larger than 15K nodes Chen et al. [2].

In this paper, we firstly correct the above shortcomings building upon the work of Aadithya et al. [1]. We further show that the use of the Shapley value is not restricted to those games and propose the use of the Shapley value in the LDAG model from Chen et al. [2]. Finally, we verify the usefulness of the algorithms in [1] in information diffusion application and compare it against current state-of-the-art algorithms for estimation of top k nodes.

2 DISCOUNT SHAPLEY VALUE CENTRALITY

The importance of a node as described in the previous section can be modeled by a game where the marginal contribution of a node is the average probability that a node contributes its neighbors. A basic characteristic function can be defined as

$$v_1(C) = \begin{cases} 0 & \text{if } C = \emptyset \\ \text{size}(\text{surrounding}(C)) & \text{otherwise} \end{cases} \quad (1)$$

where $\text{surrounding}(C)$ is the set of nodes which are the neighbours of the nodes in the coalition, formally $v \in \text{surrounding}(C)$ if $\{v \in V : \exists u \in C \text{ such that } (u, v) \in E \text{ and } v \notin C\}$.

It can be easily shown that the Shapley value of this game can be computed in polynomial time by reducing this game to the Game 1 in Aadithya et al. [1] and it is equal to $\sum_{v_k \in \{v \cup N(v)\}} \frac{1}{1 + \text{deg}(v_k)} - 1$.

In the influence maximization problem we desire to find seeds that maximize the probability that a node in a network infects its neighbors given a changing and unknown set of already infected nodes. The surrounding model addresses this issue partially and calculates the expected node contribution. However, it neglects two key aspects of information diffusion: 1) a node cannot infect itself, 2) diffusion is a process over time, not limited to one-step infection at a single point in time (thus, expected value should also be conditioned on nodes that are infected but are not part of the initial seed). These two issues are addressed in the Discounted Shapley value model.

In resolving the first issue we take advantage of the fact that in the the above equation the Shapley value is the sum of probabilities that the node contributes each of its neighbors and itself. As we only want to consider the influence of a node on others, Algorithm 1 ignores the probability that a given node contributes itself (lines 2-5).

We subsequently address the second issue, that is we attempt to account for active nodes that have been infected in the time steps $t > 1$. From among the uninfected nodes we pick the node with the highest Shapley value, add it to the set of the top nodes A and “remove” it and its neighbors from the network. We do it because adjacent nodes are likely to share a substantial number of neighbours, meaning that when both nodes are chosen as top nodes, they will trigger influence wave in the same region of the graph, leaving more distant nodes unaffected. Subsequently for the uninfected nodes, we update their SV by subtracting the probability that they influence “removed” nodes.

Algorithm 1: Discounted Shapley Value

```

1 for  $i$  to  $n$  do
2   foreach  $u \in \text{neighbor}(v_i)$  do
3      $\text{shapley}[i]_+ = \frac{1}{1 + \text{deg}(u)}$ ;
4   end
5 end
6  $A \leftarrow \emptyset$ ;  $\text{infected} \leftarrow \emptyset$ ;
7 for  $1$  to  $k$  do
8   if not all nodes are infected then
9      $\text{topnode} \leftarrow \text{argmax}_{i \notin \text{infected}} \{\text{shapley}[i]\}$ ;
10     $A \leftarrow A \cup \{\text{topnode}\}$ ;
11     $\text{infected} \leftarrow \text{infected} \cup \{\text{topnode}\}$ ;
12    foreach  $u \in \text{neighbor}(\text{topnode})$  do
13       $\text{infected} \leftarrow \text{infected} \cup \{u\}$ ;
14      foreach  $i \in \text{neighbor}(u)$  do
15         $\text{shapley}[i]_- = \frac{1}{1 + \text{deg}(u)}$ ;
16      end
17    end
18  end
19  else
20    Choose node  $\notin A$  with highest initial Shapley value and add to  $A$ ;
21  end
22 end
23 return  $A$  containing top  $k$  nodes ;

```

¹ University of Oxford, Seoul National University, email: kamil.m.adamczewski@gmail.com

² Institute of Computer Science, Polish Academy of Sciences, email: s.matejczyk@phd.ipipan.waw.pl

³ Both first authors contributed equally to this work.

⁴ University of Oxford, University of Warsaw, email: tomasz.michalak@cs.ox.ac.uk

3 SHAPLEY VALUE AND BANZHAF INDEX CENTRALITY IN LOCAL DAGs

We also propose an algorithm which incorporates game theoretical solution concepts of the Shapley value and the Banzhaf index into the greedy approach by [2] to find the most influential nodes called local DAG (LDAG).

The motivation behind the LDAG method comes from two observations. 1) Computing seed spread function is $\#P$ -hard under both models main models. 2) Finding set that maximizes spread function is NP-hard even if we can compute the spread function in polynomial time. Thus, Chen et al. [2] propose to reduce the network and form a LDAG, a directed acyclic graph which encapsulates the approximate influence exerted on a given node in the network. LDAG is chosen for each node in order to capture as much “influence” from the entire network as possible. Although, Chen et al. [2] prove that finding such a graph is NP-hard, they observe that a greedy algorithm performs very well in practice.

Since we assume we can reduce the entire network to a set of LDAGs, we find the initial seed set by analyzing the most influential nodes in all the LDAGs. While, in order to achieve this, Chen et al. [2] use a greedy approach, we propose an alternative approach that uses the Shapley value and the Banzhaf index as measures of node centrality in the LDAGs. Since the computation of both solution concepts is usually challenging, we use Monte Carlo simulations where we approximate the solution by sampling permutations. Furthermore, we take advantage of the LDAG structure which is comparatively small compared to the size of the network. We also reduce the number of input nodes for the computation of the solution concepts (this reduced the number of necessary Monte Carlo simulations) using properties of power indices.

We particularly take advantage of the additive nature of these two solution concepts. As a result, this game-theoretic approach, as opposed to the greedy approach in Chen et al. [2] is particularly suitable for distributed systems, because the resulting power indices can be computed independently on LDAGs and easily merged.

We compute the Shapley value and the Banzhaf index assuming that the characteristic function is the approximated influence spread in LDAG. In the Banzhaf index case a node v is influenced independently by its predecessors and ancestors in a given LDAG. This makes possible to run Monte Carlo simulations for these sets independently and when calculating $BI(v)$ in $ldag(u)$ (DAG directed at u) we can forget about nodes that are neither v 's ancestors nor predecessors. Using this fact we can reduce the number of MC iterations even further.

4 EXPERIMENTS

The experiments consist of two parts, 1) finding k most influential nodes according to each algorithm (k is 2-30% of the network size), 2) testing the performance of the seed set by means of Monte Carlo simulations. We conduct the experiments on two diffusion models: Independent Cascade and Linear Threshold Kempe et al. [4].

As far as the quality of seed set is concerned, the greedy LDAG performs consistently best across all the data sets, seed sizes and on both models (Chen et al. [2] only test it on the LT model).

The performance of CELF++ and Shapley value LDAG approach are similar on IC model, where CELF++ performs slightly better for smaller seed size and the roles reverse for larger seed size. SV LDAG performs better on the LT model which makes sense since LDAG is designed for LT model. DSV and [1] perform similarly in the IC model and DSV is slightly better in the LT model. In the larger net-

works with thousands of nodes, the performance of the Shapley value LDAG and DSV is only preceded by the greedy LDAG. The three algorithms perform substantially better than the Degree Discount algorithm.

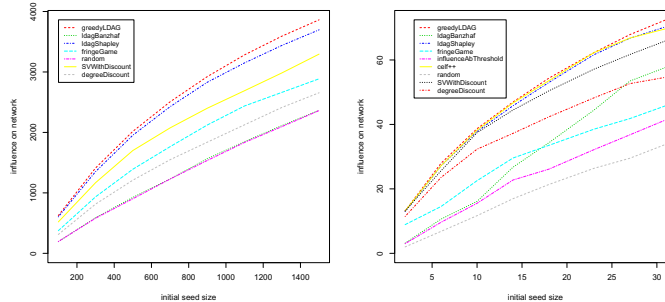


Figure 1. Comparison of various methods seed set quality (expected spread) as a function of its size (k) for two different real life data.

5 CONCLUSION

Our result show that the Shapley value is a competitive centrality measure for information diffusion. Specifically, we presented two algorithms that use the Shapley value in the two approaches recently proposed in the literature to determine the most influential nodes in a network: a greedy approach which relies on repeated computation of the information spread, and the heuristics which uses the Shapley value (which is exact and computable in polynomial time) as a centrality measure. We also verified the performance of the Shapley value-based centrality proposed in the work of Aadithya et al. [1]. The experimental result show that the greedy LDAG approach comes up with the highest quality seed set. Yet our proposed heuristic based on the Shapley value performs almost as good as the greedy algorithm in terms of solution quality, and it can be easily adopted to Map-Reduce scheme. Finally, the Discount Shapley value centrality heuristic performs better than the models from Aadithya et al. [1] and the current state-of-the-art Discount Degree heuristic; thus, narrowing the gap between the centrality based heuristics and greedy approximations.

ACKNOWLEDGEMENTS

Kamil Adamczewski & Tomasz Michalak were supported by the Polish National Science Centre grant DEC-2013/09/D/ST6/03920 (University of Warsaw). Szymon Matejczyk was supported by the European Union from resources of the European Social Fund. Project PO KL “Information technologies: Research and their interdisciplinary applications”. Tomasz Michalak was supported by the European Research Council under Advanced Grant 291528 (RACE) at the University of Oxford.

REFERENCES

- [1] K. V. Aadithya, B. Ravindran, T. P. Michalak, and N. R. Jennings. Efficient computation of the shapley value for centrality in networks. In *Internet and Network Economics*, pages 1–13. Springer, 2010.
- [2] Wei Chen, Yifei Yuan, and Li Zhang. Scalable influence maximization in social networks under the linear threshold model. In *Data Mining (ICDM), 2010 IEEE*, pages 88–97. IEEE, 2010.
- [3] Amit Goyal, Wei Lu, and Laks VS Lakshmanan. Celf++: optimizing the greedy algorithm for influence maximization in social networks. In *Proceedings of the 20th WWW conference*, pages 47–48. ACM, 2011.
- [4] D. Kempe, J. Kleinberg, and E. Tardos. Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 137–146. ACM, 2003.
- [5] N Ramasuri and Y Narahari. Determining the top-k nodes in social networks using the shapley value. In *Proceedings of the 7th international conference on AAMAS-Vol. 3*, pages 1509–1512, 2008.
- [6] M. Richardson and P. Domingos. Mining knowledge-sharing sites for viral marketing. In *8th ACM SIGKDD*, pages 61–70, 2002.