



**HAL**  
open science

## Personalised rating

Umberto Grandi, James Stewart, Paolo Turrini

► **To cite this version:**

Umberto Grandi, James Stewart, Paolo Turrini. Personalised rating. *Autonomous Agents and Multi-Agent Systems*, 2020, 34 (article 55), pp.1-42. 10.1007/s10458-020-09479-2 . hal-03066902

**HAL Id: hal-03066902**

**<https://hal.science/hal-03066902>**

Submitted on 15 Dec 2020

**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.

## Personalised Rating

**Umberto Grandi · James Stewart ·  
Paolo Turrini**

Received: date / Accepted: date

**Abstract** We introduce Personalised Rating, a network-based rating system where individuals, connected in a social network, decide whether or not to consume a service (e.g., a restaurant) based on the evaluations provided by their peers. We compare Personalised Rating with the more widely used Objective Rating where, instead, customers receive an aggregate evaluation of what everybody else has declared so far. We focus on the manipulability of such systems, allowing a malicious service provider (e.g., the restaurant owner) to transfer monetary incentive to the individuals in order to manipulate their rating and increase the overall profit. We study manipulation under various constraints, such as the proportion of individuals who evaluate the service and, in particular, how much the attacker knows of the underlying customers' network, showing the conditions under which the system is bribery-proof, i.e., no manipulation strategy yields a strictly positive expected gain to the service provider. We also look at manipulation strategies that are feasible in theory but might, in general, be infeasible in practice, deriving a number of algorithmic properties of manipulation under Personalised Rating. In particular we show that establishing the existence of a rewarding manipulation strategy for

---

This work revises and extends papers presented at IJCAI-2016 (Grandi and Turrini, 2016) and at AAAI-2018 (Grandi et al., 2018). We are grateful for the feedback received by several anonymous reviewers as well as the audiences of COMSOC-2018 in Troy, U.S.A., the Workshop on Theoretical Aspects of e-Democracy in Toulouse, France, in 2017, and the 6th World Congress of the Game Theory Society in Maastricht, The Netherlands, in 2016.

---

Institut de Recherche en Informatique de Toulouse (IRIT)  
University of Toulouse, France  
E-mail: umberto.grandi@irit.fr

Department of Computer Science  
University of Oxford, United Kingdom  
E-mail: james.stewart@cs.ox.ac.uk

Department of Computer Science  
University of Warwick, United Kingdom  
E-mail: p.turrini@warwick.ac.uk

the attacker—and, notably, an optimal one—is **NP**-complete, even with full knowledge of the underlying network structure.

## 1 Introduction

We use online reviews all the time: for food, movies and even doctors. But can we trust them? Online rating systems such as TripAdvisor, Amazon or Netflix, where a small proportion of users writes reviews which are read by a potentially large number of others, are clearly manipulable: each service provider is able to offer a compensation—monetary or otherwise—in exchange for a positive review, having an impact on the whole set of potential customers. These systems are based on what we call *Objective Rating*, or O-RATING: individual evaluations are aggregated into a single figure, which is seen by, and thus influences, every potential user.

The New York Times (Streitfeld, 2016) has recently argued how the most commonly used rating systems come with “persistent controversies over how many of the reviews on the internet were bought by the subject rather than written as finely reasoned opinions from a neutral party, and whether that distorts all results”. In particular, relating to what is demonstrated by de Lange et al. (2016), their emphasis on “the average user rating as a cue for objective quality appears to be based on an *illusion of validity*” and shows “a substantial disconnect between the objective quality of information that online user ratings actually convey and the extent to which consumers trust them as indicators of objective quality” (de Lange et al., 2016). As the number of decisions taken based on rating systems grows by the day, it is fair to say that the problem of deterring their manipulation is one of the big challenges faced by companies and governments today.

Distributed Artificial Intelligence has seen a sharp rise in the study of recommendation systems, i.e., platforms that construct (often learning-based) protocols to match users and provide accurate suggestions, as an effort to improve the trustworthiness of online rating (for a survey see, e.g., Bobadilla et al., 2013). Users’ evaluations can be carefully screened, and obvious biases (e.g., ethnicity-based discrimination) can be detected, but the AI of recommendation systems still needs to be combined with an analysis of the *incentives* needed to avert manipulation. In particular, as also noted by Tennenholtz (2008) and Alon et al. (2015), we need to establish the theoretical guarantees that are required for such systems to be resistant to manipulation. Recently, researchers in algorithmic mechanism design have begun to focus their attention on recommendation systems (see, e.g., the emerging subject of economic recommendation systems, Bahar et al., 2016) as an effort to close this gap and to start studying rating formation from a strategic point of view. Our paper should be seen as part of this effort.

What we study here is a rating system in which each individual sees *only* the evaluation given by the set of trusted peers, his or her friends, and only this aggregated opinion influences his or her decision. This is what we call

*Personalised Rating*, or P-RATING, which can be seen as a generalisation of O-RATING in which influence has a complex network structure. In particular, while in the case of O-RATING, the restaurant owner *knows* exactly how influence flows among the customers, this might not be the case with P-RATING. It is therefore important to assess the capacity of the external attacker to modify the network as a function of what he or she knows about its structure, which is the point of view we take in this paper.

**Our contribution** We analyse the effect of bribing strategies in the case of O-RATING and P-RATING under various constraints depending on the presence of customers who do not express any opinion and the knowledge of the network that the service provider has. In particular we investigate the cases of whether the exact network is known to the attacker, the network is known but not the customers' exact positions, and the network is completely unknown. We show under which conditions the system is *bribery-proof*, i.e., there is no bribe yielding a strictly positive gain for the service provider, and we provide algorithms for computing (all) optimal bribing strategies when they do exist. Intuitively, being able to know and *bribe* influential customers is crucial for guaranteeing that a bribing strategy yields a positive reward. However, while with large populations of non-voters *random* bribes can still be profitable, the effect of P-RATING is largely different from that of O-RATING and, as we show, the expected profit with the former can be severely limited and drops below zero in all networks, under certain (mild) conditions on the cost of bribes.

We also study manipulation from an algorithmic point of view. In doing so, we follow the standard approach of computational voting theory (for a survey see, e.g., Conitzer and Walsh, 2016), which has shown that even if strategy-proof voting rules cannot be designed, as a consequence of the Gibbard-Satterthwaite Theorem (Gibbard, 1973; Satterthwaite, 1975), some voting rules exist that are safe from manipulation *in practice*, as computing manipulation strategies is **NP**-hard. We show that even if a personalised rating system is manipulable in theory, the problem of manipulation, in general, might be infeasible in practice. In particular, we establish that even when the attacker has full knowledge of the network the problem of determining the existence of a manipulation strategy that guarantees a given reward—and, in particular, an optimal one—is **NP**-complete. We do so by giving a polynomial-time reduction from the problem of finding an independent set of a given size  $k$  in a 3-regular graph.

**Related research lines** Our approach relates to several research lines in the field of multi-agent systems:

*Network-based voting and mechanism design* We study social networks in which individuals' local decisions can be manipulated to modify the resulting global properties. A similar approach is taken by Apt and Markakis (2014) and Simon and Apt (2015), which study the changes on a social network needed to make users adopt a certain product. Further contributions include rational secret sharing and multi-party computation (Abraham et al., 2006), the strategic manipulation of peer reviews (Kurokawa et al., 2015),

and the growing literature on voting in social networks (see, e.g., the survey by Grandi (2017)). A highly relevant line of research is the work of Lev and Tennenholtz (2017), which studies theoretical guarantees for group recommendations, as well as papers that have looked at social network-based recommendations, such as the work of Andersen et al. (2008).

*Lobbying and Bribery* Our framework features an external agent trying to influence individual decisions to reach his or her private objectives. Lobbying in decision-making is an important problem in the area of social choice which is evident from the seminal contribution of Helpman and Persson (1998) to more recent studies in multi-issue voting (Christian et al., 2007; Grandi et al., 2019). Lobbying and bribery are also established concepts in computational social choice, with their computational complexity being analysed extensively in multiple recent papers (Faliszewski et al., 2009; Baumeister et al., 2011; Bredereck et al., 2014, 2016). Although our focus is manipulation by incentives, there are a number of relevant approaches that have close connections to our work, notably the work of Conitzer et al. (2010), Waggoner et al. (2012), Todo and Conitzer (2013) and Brill et al. (2016), which study the effect of adding fake profiles to a social network, a closely related problem to bribery. The importance of this line of research is to demonstrate the role of the graph structure in resisting manipulation, with applications to opinion spreading Alon et al. (2015) and community detection Todo and Conitzer (2013).

*Reputation-based systems* We study the aggregation of possibly insincere individual evaluations by agents that can influence one another through trust relations. In this sense our framework can be seen as a study of reputation in Multi Agent Systems, which has been an important concern of MAS in past decades (Conte and Paolucci, 2002; Sabater and Sierra, 2005; Garcin et al., 2009). In particular, our framework treats reputation as a manipulable piece of information, not just a static aggregation of individual opinions, much as is the case in the work of Conte et al. (2008) and Pinyol and Sabater-Mir (2013).

*Agent-Mediated Electronic Commerce* Ever since the seminal contributions by Sierra (2004), Feigenbaum et al. (2009), and Dash et al. (2003), the concept of agent-mediated market has taken a central role in distributed AI. It goes without saying that the increased level of automation in these markets requires increasingly robust mechanisms to avoid manipulation. The stream of literature in agent-mediated electronic commerce is often tied to game-theoretic modelling and our work aligns to it by proposing a personalised rating model that is trustworthy by design.

**Paper structure** Section 2 presents the formal setup of the framework, introducing the notions of objective and personalised ratings, together with some basic observations concerning their manipulability. Section 3 is devoted to objective rating, and shows the effect of bribes both in the presence of non-voters and when all customers vote. Section 4 moves onto the study of personalised rating. In particular, it compares the effect of bribes to the results on objective

rating, emphasising the importance of a suitably defined notion of influence weights between customers. It then shows how knowledge of the network plays a role in determining the optimal bribing strategies and the cases in which one does not exist. Our results are constructive: we give algorithms for determining effective bribing strategies and the computational complexity of doing so. Section 5 extensively discusses the consequences of our definitions and results, paying particular attention to the underlying assumptions on customers' behaviour. Finally, Section 6 concludes the paper and presents a number of future research directions.

## 2 Objective and Personalised Rating

In this section we provide the basic formal definitions of objective and personalised rating, and we give a formal definition of a bribing strategy and its effect on ratings. We also define the revenue resulting from the execution of a given bribing strategy, and we present some basic examples and first results on the manipulability of the two rating systems that we have defined.

### 2.1 Restaurants and customers

Our framework features an external object  $r$ , the *restaurant*, which is evaluated by a finite non-empty set of individuals  $C = \{c_1, \dots, c_n\}$ , called *customers*. Customers are connected in a network  $E \subseteq C \times C$ , called the *customers network*. We assume  $E$  to be an undirected graph, i.e., the relation  $E$  is symmetric. We also assume  $E$  to be reflexive, i.e.,  $(c, c) \in E$  for all  $c \in C$ . Given  $c \in C$  we call  $N(c) = \{x \in C \mid (x, c) \in E\}$  the *neighbourhood* of  $c$ . Note that this always includes  $c$  itself.

Customers concurrently submit an *evaluation* of the restaurant, drawn from a set of values  $Val \subseteq [0, 1]$ , together with a distinguished element  $\{*\}$ , which represents no opinion. Examples of sets of values are the set  $[0, 1]$  itself, or a discrete assignment of 1 to 5 stars, as is common in online rating systems. We make the assumption that  $\{0, 1\} \subseteq Val$  and that  $Val$  is closed under the operation  $\min\{1, x + y\}$  for all  $x, y \in Val$ . These assumption allow us to map the most common rating methods, e.g., 1 to 5 stars or any continuous bounded set of values, onto the  $[0, 1]$  interval and analyse them within our framework.

We represent the customer evaluations as a function  $eval : C \rightarrow Val \cup \{*\}$  and define  $V \subseteq C$  as the subset of customers that express an evaluation of the restaurant, i.e.,  $V = \{c \in C \mid eval(c) \neq *\}$ . We refer to this set as the set of *voters* and we assume it to be always non-empty, i.e., there is at least one customer that expresses an opinion.

## 2.2 Two rating systems

In online rating systems such as Tripadvisor<sup>®</sup> every interested customer can see—and is therefore influenced by—the (average of) what the other customers have written. We call this method O-RATING, which stands for *objective rating*. Given an evaluation function  $eval$  for a restaurant, the associated O-RATING is defined as

$$\text{O-RATING}(eval) = \text{avg}_{c \in V} eval(c),$$

where  $\text{avg}$  is the average function across real-valued  $eval(c)$ , disregarding  $*$  and thus restricted to the set of voters  $V$ .

O-RATING flattens individual evaluations into a unique objective aggregate, the rating that a certain restaurant is given. What we propose is a refinement of O-RATING, which takes the network of influence into account. In this system customers are *only* interested in the evaluation of other customers they can trust, e.g., their friends. We call our method P-RATING, which stands for *personalised rating*. It is defined for a customer-evaluation pair  $(c, eval)$  as follows.

$$\text{P-RATING}(c, eval) = \text{avg}_{k \in N(c) \cap V} eval(k)$$

So, the  $\text{P-RATING}(c, eval)$  give the opinion of a customer  $c$  on the restaurant, taking the average of the opinions of the (voting) customers that  $c$  is connected to. We omit  $eval$  whenever clear from the context, writing simply  $\text{P-RATING}(c)$ .

Observe that in case a customer is not connected to a voter, then P-RATING is not defined. To facilitate the analysis we make the technical assumption that *each customer is connected to at least one voter*. Also observe that when  $E = C \times C$ , i.e., in the case that the network is complete and each individual is influenced by every other individual, then for all  $c \in C$  and  $eval$  we have that  $\text{P-RATING}(c, eval) = \text{O-RATING}(eval)$ .

Finally, note that while O-RATING and P-RATING are defined in terms of average, different aggregators are possible and will be discussed in Section 5.

*Example 1* Figure 1 shows three customers connected in a network. Two of them,  $c_1$  and  $c_2$ , have provided their own score of 0.9 and 0.3 respectively, while  $c_3$  hasn't. Given these evaluations we have

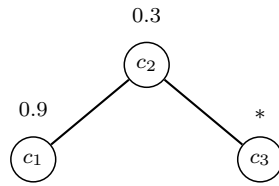
$$\text{O-RATING}(eval) = \text{avg}_{c \in V} eval(c) = \text{avg}\{0.9, 0.3\} = 0.6$$

while

$$\begin{aligned} \text{P-RATING}(c_1, eval) &= \text{P-RATING}(c_2, eval) = \\ &\text{avg}_{k \in N(c_1) \cap V} eval(k) = \text{avg}\{0.9, 0.3\} = 0.6 \end{aligned}$$

and

$$\text{P-RATING}(c_3, eval) = \text{avg}_{k \in N(c_3) \cap V} eval(k) = 0.3$$



**Fig. 1** A network of three customers and their evaluations. The symbols above each node represent the individual evaluations, the straight lines are the customers' connections.

Notice how, despite a fairly positive objective rating of 0.6, customer  $c_3$  is influenced by  $c_2$  only and therefore forms a more negative opinion of the restaurant.

The definition of Personalised Rating comes with a number of simplifications and constraints that affect its structural properties, perhaps the most basic of which being the fact that each individual attaches equivalent importance to the opinions of itself and each of its neighbours. There are more complex alternative ways in which an individual might aggregate the opinions of those they are influenced by, however many of the interesting properties and limitations of network-based rating systems can be observed even in this setting, as defined above. Alternative methods of aggregation, for example, weighted variants of the P-rating, and ideas for future work will be discussed in detail in Sections 5 and 6.

### 2.3 Utilities and strategies

We interpret a customer evaluation as a measure of his or her *propensity* to go to the restaurant. We therefore assume that the utility that a restaurant gets is proportional to its rating. To simplify the analysis we assume a factor 1 proportionality, that is to say the restaurant's utility is equivalent to the sum of all customers' (personalised) ratings. We devote Section 5 to discuss the impact and the relaxation of this and other simplifying assumptions.

#### 2.3.1 The case of O-RATING

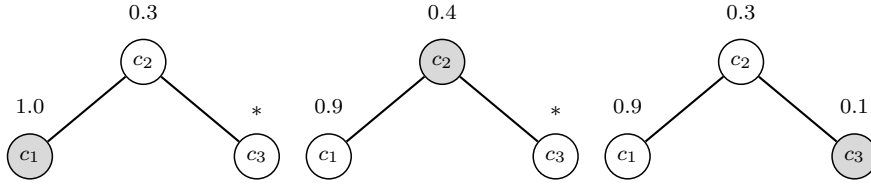
For O-RATING, we assume that the restaurant's initial utility is defined as

$$u_O^0 = |C| \text{O-RATING}(eval).$$

Intuitively, the initial utility amounts to the number of customers that actually go to the restaurant, weighted by their (average) predisposition. For example, the initial utility of the restaurant with the customers' evaluations as in Figure 1 is  $u_O^0 = 3 \times 0.6 = 1.8$ .

At the initial stage of the game, the restaurant owner receives  $u_O^0$ , and can then decide to invest a part of it to influence a subset of customers and improve upon the initial utility. We assume utility to be fully transferable and,





**Fig. 2** The effect of a single bribe by 0.1 on each of the customers in Figure 1. The left ( $\sigma_A$ ), middle ( $\sigma_B$ ) and right ( $\sigma_C$ ) strategies, although costing the same (0.1) to the restaurant induce different returns. In particular  $u_O^{\sigma_A} = u_O^{\sigma_B} = 3 \times 0.65 - 0.1 = 1.85$  while  $u_O^{\sigma_C} = 3 \times 0.43 - 0.1 = 1.2$ . Observe how the first two give a positive change with respect to the original O-RATING while the third one does not. Hence,  $\sigma_A$  and  $\sigma_B$  are profitable strategies while  $\sigma_C$  is not.

to facilitate the analysis, that such transfers translate directly into changes in customers' evaluations.

**Definition 1** A strategy is a function  $\sigma : C \rightarrow Val$  such that  $\sum_{c \in C} \sigma(c) \leq u_O^0$ .

Definition 1 imposes that strategies are *weakly budget balanced*, i.e., restaurants can only pay with resources they have. We denote by  $\sigma^0$ , the strategy that assigns 0 to all customers and we call a *bribing strategy* any strategy that is different from  $\sigma^0$ . After the execution of a bribing strategy, the evaluation is updated according to the following definition.

**Definition 2** The evaluation  $eval^\sigma(c)$  after the execution of  $\sigma$  is  $eval^\sigma(c) = \min\{1, eval(c) + \sigma(c)\}$ , where  $* + \sigma(c) = \sigma(c)$ , if  $\sigma(c) \neq 0$ , and  $* + \sigma(c) = *$ , if  $\sigma(c) = 0$ .

In this definition we are making the assumption that the effect of bribing a non-voter to vote is equivalent to that of bribing a voter that has an evaluation of 0 as, intuitively, the individual has no associated predisposition to go to the restaurant. Section 5 will discuss the extent of such assumption further.

A strategy is called *efficient* if  $\sigma(c) + eval(c) \leq 1$  for all  $c \in C$ . Let  $B(\sigma) = \{c \in C \mid \sigma(c) \neq 0\}$  be the set of bribed customers, and let  $V^\sigma$  be the set of voters after the execution of  $\sigma$ . Executing  $\sigma$  induces the following change in utility:

$$u_O^\sigma = |C| \text{O-RATING}(eval^\sigma) - \sum_{c \in C} \sigma(c).$$

Intuitively,  $u_O^\sigma$  is obtained by adding to the initial utility of the restaurant the utility obtained from paying each customer, minus the amount of money spent. Figure 2 shows an example of strategy executions and their effect on O-RATING.

We define the revenue of a strategy  $\sigma$  as the marginal utility obtained by executing it.

**Definition 3** Let  $\sigma$  be a strategy. The *revenue* of  $\sigma$  is defined as  $\mathbf{r}_O(\sigma) = u_O^\sigma - u_O^0$ . We say that  $\sigma$  is *profitable* if  $\mathbf{r}_O(\sigma) > 0$ .

Considering the strategies in Figure 2, we have a revenue of 0.05 for both  $\sigma_A$  and  $\sigma_B$  while the revenue is of  $-0.6$  for  $\sigma_C$ , which is not a profitable bribing strategy. Finally, we recall the standard notion of dominance and optimality.

**Definition 4** A strategy  $\sigma$  *weakly dominates* a strategy  $\sigma'$  if  $u_O^\sigma \geq u_O^{\sigma'}$ . A strategy  $\sigma$  is *weakly dominant* or *optimal* if it weakly dominates every other strategy  $\sigma'$ .

### 2.3.2 The case of P-RATING

The previous definitions is lifted to the case of P-RATING in the expected way. In this case, we view the rating of the restaurant, which is now personalised for each customer, as the propensity of the customer to visit the restaurant. The initial utility received can therefore be defined as follows:

$$u_P^0 = \sum_{c \in C} \text{P-RATING}(c, \text{eval}).$$

In the situation of Figure 1 we have that  $u_P^0 = \sum_{c \in C} \text{P-RATING}(c, \text{eval}) = 2 \times 0.6 + 0.3 = 1.5$ . Assuming the same effect of a bribe on the customers' evaluations, we can define the utility generated by a bribing strategy as:

$$u_P^\sigma = \sum_{c \in C} \text{P-RATING}(c, \text{eval}^\sigma) - \sum_{c \in C} \sigma(c).$$

Finally, let the revenue of a bribing strategy  $\sigma$  be  $\mathbf{r}_P(\sigma) = u_P^\sigma - u_P^0$ , with profitable strategies being those generating a positive revenue. If it is clear from the context, we use  $\text{P-RATING}^\sigma(c)$  for  $\text{P-RATING}(\text{eval}^\sigma, c)$ , i.e., the P-RATING obtained after the execution of bribing strategy  $\sigma$ .

*Example 2* Let us consider the strategies displayed in Figure 2 and their effect on P-RATING. Let us start with  $\sigma_A$ . We have that  $\text{P-RATING}^{\sigma_A}(c_1) = \text{P-RATING}^{\sigma_A}(c_2) = 0.65$  while  $\text{P-RATING}^{\sigma_A}(c_3) = 0.3$ . Note already how bribing  $c_1$  does not affect  $c_3$ , whose P-RATING stays the same. We thus have that  $u_P^{\sigma_A} = 1.5 = u_P^0$ , which means  $\mathbf{r}_P(\sigma_A) = 0$  and thus  $\sigma_A$  is not a profitable strategy. Observe also that any strategy that behaves like  $\sigma_A$ , bribing only customer  $c_1$ , gives either a zero or a negative revenue. Conversely,  $\sigma_B$  is different. In this case, we also have that  $\text{P-RATING}^{\sigma_B}(c_1) = \text{P-RATING}^{\sigma_B}(c_2) = 0.65$ . However,  $\text{P-RATING}^{\sigma_B}(c_3) = 0.4$ , which means that  $u_P^{\sigma_B} = 1.6$  and thus  $\mathbf{r}_P(\sigma_B) = 0.1$ . So, the strategy  $\sigma_B$  yields a positive revenue to the restaurant. Finally,  $\sigma_C$  is such that  $\text{P-RATING}^{\sigma_C}(c_1) = 0.6$ ,  $\text{P-RATING}^{\sigma_C}(c_2) = 0.4\bar{3}$  and  $\text{P-RATING}^{\sigma_C}(c_3) = 0.2$ , which means  $u_P^{\sigma_C} = 1.1\bar{3}$  and thus  $\mathbf{r}_P(\sigma_C) = -0.3\bar{6}$ . So, the strategy  $\sigma_C$  yields a negative revenue to the restaurant.

In later sections, in order to determine the optimal bribing strategies we need to establish how the customers vote, how they are connected, and what the restaurant owner knows. We assume that the restaurant knows *eval*, leaving the interesting case of unknown *eval* to future work. We focus instead on the following cases: the restaurant knows the network, the restaurant knows the shape of the network but not the individuals' positions, and the network is unknown. We analyse the effect of bribing strategies on P-RATING in each case. Notice how for the case of O-RATING all collapse to the first. We also look at the special case in which every customer is a voter. Given some set of assumptions, we say that O-RATING (or P-RATING) is *bribery-proof* under those assumptions if  $\sigma^0$  is optimal.

## 2.4 Charting the boundaries of manipulation

Bribing strategies have a different effect on the revenue obtained under the two rating systems that we have defined. Whilst the utility of O-RATING is a sum of the *global* average of voters' evaluations, the utility of P-RATING is a sum of *local* averages of voters' evaluations based on the evaluations of their peers. Therefore, a strategy bribing one voter affects everybody in the case of O-RATING, but it can be shown to have a limited effect in the case of P-RATING. In this section we provide initial results on evaluating the revenue that can be obtained through bribing strategies in the two rating systems. We begin with the following proposition.

**Proposition 1** *Let  $\sigma$  be an efficient strategy s.t.  $|B(\sigma)| = 1$ , and let  $\bar{c}$  be such that  $\sigma(\bar{c}) \neq 0$ . Then  $\mathbf{r}_P(\sigma) < |N(\bar{c})|$ .*

*Proof* By calculation, we have that:

$$\begin{aligned}
 \mathbf{r}_P(\sigma) &= u_P^\sigma - u_P^0 \\
 &= \sum_{c \in C} \text{P-RATING}^\sigma(c) - \sigma(\bar{c}) - \sum_{c \in C} \text{P-RATING}(c) \\
 &= \sum_{c' \in N(\bar{c})} \text{P-RATING}^\sigma(c') - \sigma(\bar{c}) - \sum_{c' \in N(\bar{c})} \text{P-RATING}(c') \\
 &\leq 1 \times |N(\bar{c})| - \sigma(\bar{c}) - \sum_{c' \in N(\bar{c})} \text{P-RATING}(c') \\
 &< |N(\bar{c})|,
 \end{aligned}$$

where the last step is obtained by observing that  $\text{P-RATING}(c) \in [0, 1]$  and  $\sigma(\bar{c}) > 0$ .

The previous result shows that increasing the number of individuals that are not connected to an agent that is being bribed, even if these are non-voters, does not increase the upper bound on the revenue of the bribing strategy. This is not true when we use O-RATING.

**Proposition 2** *Let  $\sigma$  be an efficient strategy bribing voters only. The revenue  $\mathbf{r}_O(\sigma)$  of  $\sigma$  is monotonically increasing in the number of non-voters.*

*Proof* It follows from our definitions that:

$$\begin{aligned} \mathbf{r}_O(\sigma) &= \\ u_O^\sigma - u_O^0 &= \\ |C| \frac{\sum_{c \in V} \text{eval}(c) + \sigma(c)}{|V|} - \sum_{c \in C} \sigma(c) - |C| \frac{\sum_{c \in V} \text{eval}(c)}{|V|} &= \\ \left( \frac{|C|}{|V|} - 1 \right) \sum_{c \in V} \sigma(c). \end{aligned}$$

The above figure is monotonically increasing in the number of non-voters, which can be obtained by increasing  $C$  keeping  $V$  fixed.

A similar, albeit more complex, point can be made regarding strategies that also bribe non-voters, as is explained in the following section. All in all, under realistic assumptions on the network structure, such as a the presence of a very large proportion of non-voters and with customers having a small number of connections, we can already show that bribing under O-RATING is increasingly rewarding, while under P-RATING this is no longer the case.

### 3 Manipulation of Objective Rating

In this section we consider bribing strategies under O-RATING, first focussing on the case where everyone expresses an opinion, then moving on to the more general case. In the first case we show that O-RATING is bribery-proof, and for the second we give an efficient algorithm to compute optimal bribing strategies.

#### 3.1 All vote

Let us now consider the case in which  $V = C$ . First, we show the following result.

**Proposition 3** *If  $V = C$ , then no strategy is profitable.*

*Proof (Proof)* Let  $\sigma$  be an arbitrary bribing strategy, and  $C = V$ . We therefore have that

$$\begin{aligned}
\mathbf{r}_O(\sigma) &= |C| \text{O-RATING}(eval^\sigma) - \sum_{c \in C} \sigma(c) - |C| \text{O-RATING}(eval) \\
&= |C| \frac{\sum_{c \in C} \min\{1, eval(c) + \sigma(c)\}}{|C|} - \sum_{c \in C} \sigma(c) - |C| \frac{\sum_{c \in C} eval(c)}{|C|} \\
&\leq |C| \frac{\sum_{c \in C} eval(c) + \sigma(c)}{|C|} - \sum_{c \in C} \sigma(c) - |C| \frac{\sum_{c \in C} eval(c)}{|C|} \\
&= |C| \frac{\sum_{c \in C} eval(c)}{|C|} + \sum_{c \in C} \sigma(c) - \sum_{c \in C} \sigma(c) - |C| \frac{\sum_{c \in C} eval(c)}{|C|} \\
&= 0.
\end{aligned}$$

Each step is a straightforward consequence of our definitions. Observe that any efficient bribing strategy has revenue of exactly zero. From this it follows that  $\sigma^0$  is optimal and therefore O-RATING is bribery-proof in the case that all customers vote.

### 3.2 Non-voters

We now consider the case of  $V \subset C$ , i.e., when there is at least one customer who is not a voter. The following example shows that O-RATING, in this case, is not bribery-proof, and that the order in which customers are bribed matters, suggesting that finding an optimal bribing strategy might be a non-trivial task.

*Example 3* Let  $C = \{A, B, C\}$ , and let  $eval(A) = 0.5$ ,  $eval(B) = 0.5$ , and  $eval(C) = *$ . The initial resources are  $u^0 = \text{O-RATING} \times 3 = 1.5$ . Now let  $\sigma_1(A) = 0.5$  and  $\sigma_1(B) = \sigma_1(C) = 0$ , and let  $\sigma_2(C) = 0.5$  and  $\sigma_2(A) = \sigma_2(B) = 0$ . Now,  $u_O^{\sigma_1} = 0.75 \times 3 - 0.5 = 1.75$  and  $u_O^{\sigma_2} = 0.5 \times 3 - 0.5 = 1$ , but  $u_O^{\sigma_1 \circ \sigma_2} = 0.6 \times 3 - 1 = 1$ .

This example (in particular  $\sigma_1$ ) also shows that O-RATING, in this case, is not bribery-proof.

We now characterise the set of optimal bribing strategies. We begin by showing that bribing a non-voter is always dominated. First, let  $\sigma$  be a strategy such that  $\sigma(\bar{c}) \neq 0$  for some  $\bar{c} \in C \setminus V$  and recall that  $V^\sigma$  is the set of voters after execution of  $\sigma$ . Let us define the  $\bar{c}$ -greedy restriction of  $\sigma$  to be any strategy  $\sigma^{-\bar{c}}$  such that:

- $V^{\sigma^{-\bar{c}}} = V^\sigma \setminus \{\bar{c}\}$ , i.e., the greedy restriction eliminates  $\bar{c}$  from the set of voters.
- For each  $c \in V^\sigma \setminus \{\bar{c}\}$ ,  $\max\{1, eval(c) + \sigma(c)\} = \max\{1, eval(c) + \sigma^{-\bar{c}}(c)\}$ , i.e., the greedy restriction does not waste further resources.
- If there exists  $c \in V^\sigma \setminus \{\bar{c}\}$  such that  $eval(c) + \sigma^{-\bar{c}}(c) < 1$  then  $\sum_{c \in C} \sigma^{-\bar{c}}(c) = \sum_{c \in C} \sigma(c)$ , i.e., the  $\sigma^{-\bar{c}}$  redistributes  $\sigma(\bar{c})$  among the remaining voters.

We now show that each strategy bribing a non-voter is strictly dominated by any of its greedy restrictions.

**Proposition 4** *Let  $V \neq C$ , and  $\bar{c} \in C \setminus V$ . Then each  $\sigma$  with  $\sigma(\bar{c}) \neq 0$  is strictly dominated by  $\sigma^{-\bar{c}}$ .*

*Proof* Let  $\sigma$  be a strategy with  $\sigma(\bar{c}) \neq 0$  for some non-voter  $\bar{c}$ , and let  $\sigma^{-\bar{c}}$  be one of its greedy restriction defined above.

$$\begin{aligned} u_O^{\sigma^{-\bar{c}}} - u_O^\sigma &= |C| \left( \text{O-RATING}^{\sigma^{-\bar{c}}} - \text{O-RATING}^\sigma \right) + \sum_{c \in C} \sigma(c) - \sum_{c \in C} \sigma^{-\bar{c}}(c) \\ &= |C| \left( \frac{\sum_{c \in V^\sigma \setminus \{\bar{c}\}} \text{eval}^{\sigma^{-\bar{c}}}(c)}{|V^\sigma \setminus \{\bar{c}\}|} - \frac{\sum_{c \in V^\sigma} \text{eval}^\sigma(c)}{|V^\sigma|} \right) + \\ &\quad + \left( \sum_{c \in C} \sigma(c) - \sum_{c \in C} \sigma^{-\bar{c}}(c) \right) \end{aligned}$$

Observe first that  $\sigma^{-\bar{c}}$  is a redistribution, hence  $\sum_c \sigma(c) - \sum_c \sigma^{-\bar{c}}(c) \geq 0$ , i.e., the second addendum in the above equation is positive. Consider now the case where there exists  $c \in V^\sigma \setminus \{\bar{c}\}$  such that  $\text{eval}(c) + \sigma^{-\bar{c}}(c) < 1$ . Then by the definition of  $\sigma^{-\bar{c}}$  we have that  $\sum_{c \in V^\sigma} \text{eval}^\sigma(c) = \sum_{c \in V^{\sigma^{-\bar{c}}}} \text{eval}^{\sigma^{-\bar{c}}}(c)$ , i.e., the greedy restriction preserves the overall evaluation. By straightforward calculation this gives that  $u_O^{\sigma^{-\bar{c}}} - u_O^\sigma > 0$ . If no such  $c$  exists, and therefore  $\text{O-RATING}^{\sigma^{-\bar{c}}} = 1$  we have that either  $\text{O-RATING}^\sigma < 1$  or, by the efficiency requirement and the fact that  $\sigma(\bar{c}) \neq 0$ , we have that  $\sum_{c \in C} \sigma(c) > \sum_{c \in C} \sigma^{-\bar{c}}(c)$ . In either case we have that  $u_O^{\sigma^{-\bar{c}}} - u_O^\sigma > 0$ .

Let an *O-greedy strategy* be any efficient strategy that redistributes all the initial resources  $u_O^0$  among voters. Making use of the previous result, we are able to characterise the set of all optimal strategies for O-RATING.

**Proposition 5** *Let  $V \neq C$  and  $\sigma$  a strategy. The following three statements are equivalent:*

- (i)  $\sigma$  is optimal for O-RATING,
- (ii)  $\sigma$  is an O-greedy strategy,
- (iii)  $\sigma$  yields a payoff of  $\left( \frac{|C|}{|V|} - 1 \right) \times \min\{u_O^0, \sum_{c \in V} (1 - \text{eval}(c))\}$ .

*Proof* By Proposition 4 we know that strategies bribing non-voters are dominated. Therefore we can restrict ourselves to strategies bribing only voters. Clearly, inefficient strategies are dominated. By Proposition 2 the revenue of an efficient strategy bribing only voters is  $\left( \frac{|C|}{|V|} - 1 \right) \times \sum_{c \in C} \sigma(c)$ . Observe that to maximise this quantity, we need to maximise  $\sum_{c \in C} \sigma(c)$ . This means bribes need to reach a total of  $\min\{u_O^0, \sum_{c \in V} (1 - \text{eval}(c))\}$ , which corresponds to the definition of an O-greedy strategy.

While there may be cases in which the number of optimal strategies under O-RATING is exponential, all such strategies yields the same revenue, and Proposition 5 gives us a polynomial-time algorithm to find one of them: starting from an evaluation vector  $eval$ , distribute all available resources  $u_O^0$  to the voters, without exceeding the maximal evaluation of 1. By either exhausting the available budget or distributing it all, we are guaranteed the maximum gain by Proposition 5.

#### 4 Manipulation of Personalised Rating

In this section we analyse the effect of bribing strategies under P-RATING. Unlike the case of O-RATING, where a complete underlying graph is implicitly assumed, P-RATING is defined upon an underlying network that may or may not be known to the external attacker.

We begin by looking at the case where the network is known and, as we have done in Section 3, we first study the case when everyone votes and then move on to allowing non-voters. Subsequently, we lift the assumption that the network is known to the attacker, first looking at the case in which the network shape is known but the customers' positions are unknown, and then remove any knowledge of the underlying network altogether. For each of these cases we look at the potential for manipulation and its feasibility in practice, i.e., the computational complexity of computing a profitable bribing strategy.

Before we proceed, we introduce a useful graph-theoretic measure of influence and a closed form for calculating the utility of a bribing strategy.

**Definition 5** The *influence weight* of a customer  $c \in C$  in a network  $E$  with respect to a set of designated voters  $V$  is defined as follows:

$$w_c^V = \sum_{k \in N(c)} \frac{1}{|N(k) \cap V|}$$

Recall, we assumed that every customer can see a voter, thus  $w_c^V$  is well-defined for every  $c$ . If  $V = C$ , i.e., when everybody voted, we let  $w_c = w_c^C$ . In this case, we obtain  $w_c = \sum_{k \in N(c)} \frac{1}{deg(k)}$ , where  $deg(c) = |N(c)|$  is the *degree* of  $c$  in  $E$ . When  $V$  is defined by a bribing strategy  $\sigma$ , we write  $w_c^\sigma = w_c^{V^\sigma}$ .

Intuitively, each individual's individual rating influences the rating of each of its connections, with a factor that is inversely proportional to the number of second-level connections that have expressed an evaluation. We formalise this statement in the following lemma.

**Lemma 1** The utility obtained by playing bribing strategy  $\sigma$  with P-RATING on network  $E$  is  $u_P^\sigma = \sum_{c \in V^\sigma} w_c^\sigma \times eval^\sigma(c) - \sum_{c \in C} \sigma(c)$ .

*Proof* By calculation, we have the following

$$\begin{aligned}
u_P^\sigma + \sum_{c \in C} \sigma(c) &= \sum_{c \in C} \text{P-RATING}^\sigma(c) = \sum_{c \in C} \text{avg}_{k \in N(c) \cap V^\sigma} \text{eval}^\sigma(k) \\
&= \sum_{c \in C} \left[ \frac{1}{|N(c) \cap V^\sigma|} \sum_{k \in N(c) \cap V^\sigma} \text{eval}^\sigma(k) \right] \\
&= \sum_{k \in V^\sigma} \left[ \text{eval}^\sigma(k) \times \sum_{k' \in N(k)} \frac{1}{|N(k') \cap V^\sigma|} \right] \\
&= \sum_{c \in V^\sigma} w_c^\sigma \times \text{eval}^\sigma(c)
\end{aligned}$$

#### 4.1 All vote, known network

We begin by studying the simplest case in which every customer has expressed an opinion, the restaurant owner knows the evaluation  $eval$ , and the restaurant owner knows the position of each customer on the network, and therefore the influence between customers. Recall that  $B(\sigma)$  is the set of customers bribed by  $\sigma$ . We say that two strategies  $\sigma_1$  and  $\sigma_2$  are *disjoint* if  $B(\sigma_1) \cap B(\sigma_2) = \emptyset$ . Given two disjoint strategies  $\sigma_1$  and  $\sigma_2$ , we denote  $\sigma_1 \circ \sigma_2$  their concatenation, i.e., the strategy  $(\sigma_1 \circ \sigma_2)$  such that  $(\sigma_1 \circ \sigma_2)(c) = \sigma_1(c)$ , whenever  $\sigma_1(c) \neq 0$ , and  $(\sigma_1 \circ \sigma_2)(c) = \sigma_2(c)$ , otherwise. The following corollary is a straightforward consequence of Lemma 1:

**Corollary 1** *Let  $V = C$  and let  $\sigma_1$  and  $\sigma_2$  be two disjoint strategies, then  $\mathbf{r}_P(\sigma_1 \circ \sigma_2) = \mathbf{r}_P(\sigma_1) + \mathbf{r}_P(\sigma_2)$ .*

We are now able to give a precise characterisation of the revenue obtained by any efficient strategy  $\sigma$ .

**Proposition 6** *Let  $V = C$ , let  $E$  be a known network, and let  $\sigma$  be an efficient strategy. Then  $\mathbf{r}_P(\sigma) = \sum_{c \in C} (w_c - 1)\sigma(c)$ .*

*Proof* By calculation, where Step (2) uses Lemma 1, and Step (4) uses the fact that  $\sigma$  is efficient:

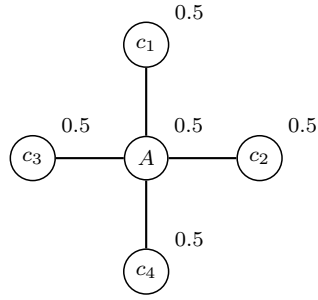
$$\mathbf{r}_P(\sigma) = u_P^\sigma - u_P^0 \quad (1)$$

$$= \left[ \sum_{c \in C} w_c \text{eval}_c^\sigma - \sum_{c \in C} \sigma(c) - \sum_{c \in C} w_c \text{eval}(c) \right] \quad (2)$$

$$= \sum_{c \in C} [w_c [\min\{1, \text{eval}(c) + \sigma(c)\} - \text{eval}(c)]] - \sum_{c \in C} \sigma(c) \quad (3)$$

$$= \sum_{c \in C} (w_c - 1)\sigma(c). \quad (4)$$





**Fig. 3** A four-arm star.

Proposition 6 tells us that the factors  $w_c$  are crucial in determining the revenue of a given bribing strategy. Bribing a customer  $c$  is profitable whenever  $w_c > 1$  (provided its evaluation was not 1 already), whilst bribing a customer  $c$  with  $w_c \leq 1$  is at most as profitable as doing nothing, as can be seen in the example below. Most importantly, it shows that P-RATING is *not* bribery-proof when the restaurant knows both the network and the customers' evaluations.

*Example 4* Let  $E$  be the four-arm star depicted in Figure 3, with  $A$  being the distinguished individual in the centre and with each customer giving an evaluation of 0.5. We have that  $w_A = 2.2$  and  $w_c = 0.7$  for all  $c$  different from  $A$ . Consider now two bribing strategies:  $\sigma^A$  which bribes  $A$  with 0.5, and  $\sigma^B$  which bribes a single individual  $B \neq A$  with the same amount. By Lemma 1, the utility we obtain is  $\mathbf{r}_P(\sigma^A) = 0.6$ , whilst  $\mathbf{r}_P(\sigma^B) = -0.15$ . This shows that a bribe is profitable only if the influence weight of the bribed customer is bigger than 1.

Given a network  $E$  and an evaluation vector  $eval$ , let Algorithm 1 define the *P-greedy bribing strategy*. Note that we make use of the notion of influence weight of Definition 5.

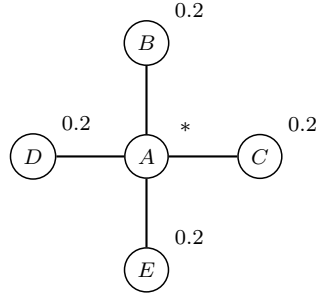
As a direct consequence of Proposition 6 we obtain the following corollary.

**Corollary 2** *The P-greedy bribing strategy defined in Algorithm 1 is optimal.*

As is the case for O-RATING, Corollary 2 has repercussions on the computational complexity of bribery: it shows that computing an optimal strategy can be done in polynomial-time. Notice how the most costly operation lies in the computation of the influence weights  $w_c$ , which can be performed only once, assuming the network is static. Similar problems, such as recognising whether bribing a certain individual is profitable, or estimating whether individuals on a network can be bribed above a certain threshold, are also efficiently computable.

**Input:** Evaluation function  $eval$  and network  $E$   
**Output:** A bribing strategy  $\sigma_P^G : C \rightarrow Val$   
 $Budget = u_P^0$   
 $\sigma_P^G(c) = 0$  for all  $c \in C$   
 Compute  $w_c$  for all  $c \in C$   
 Sort  $c \in C$  in descending order  $c_0, \dots, c_m$  based on  $w_c$   
**for**  $i=0, \dots, m$  **do**  
   **if**  $Budget \neq 0$  **then**  
     **if**  $w_{c_i} > 1$  **then**  
        $\sigma_P^G(c_i) = \min\{1 - eval(c_i), Budget\}$   
        $Budget = Budget - \sigma_P^G(c_i)$   
     **end**  
   **end**  
**return**  $\sigma_P^G$   
**end**

**Algorithm 1:** The  $P$ -greedy bribing strategy  $\sigma_P^G$



**Fig. 4** A star with a non-voter in the centre.

#### 4.2 Non-voters, known network

While we showed in Section 3 that bribing a non-voter is never optimal under objective rating, the following example shows that under P-RATING there exist networks for which the optimal bribe is to a non-voter.

*Example 5* Consider four individuals  $\{B, C, D, E\}$  connected only to a non-voter  $A$  in the middle, with  $eval(j) = 0.2$  for all  $j \neq A$ , as shown in Figure 4. We have  $u_P^0 = 1$ . Let  $\sigma_1(A) = 1$  and 0 otherwise. The utility of  $\sigma_1$  is:

$$\text{P-RATING}^{\sigma_1}(A) + 4\text{P-RATING}^{\sigma_1}(j) - 1 = 1.76$$

where  $j = B, C, D, E$ . All other strategies can be shown to be dominated by  $\sigma_1$ . Take for instance a strategy  $\sigma_2$  such that  $\sigma_2(B) = 0.8$ ,  $\sigma_2(C) = 0.2$  and 0 otherwise. The utility of  $\sigma_2$  is  $u_P^{\sigma_2} = 1.25$ .

The computation of bribing strategies involving non-voters is highly non-trivial. We begin by investigating restricted classes of strategies, in particular the ones only bribing voters and the ones bribing a single customer only.

#### 4.2.1 Voter-only strategies

Let *voter-only strategies* be bribing strategies  $\sigma$  such that  $\sigma(c) = 0$  for all  $c \notin V$ . In this case, a similar proof to Proposition 6 gives the following result.

**Proposition 7** *Let  $V \neq C$ ,  $E$  be a known network, and let  $\sigma$  be an efficient bribing strategy such that  $B(\sigma) \subseteq V$ . Then,  $\mathbf{r}_P(\sigma) = \sum_{c \in V} (w_c^V - 1)\sigma(c)$ .*

The difference with the case of  $V = C$  is that  $w_c^V$  can be arbitrarily large in the presence of non-voters, such as in Example 5.

#### 4.2.2 Single Bribes

We now investigate how much revenue can be gained by bribing a single customer with an efficient strategy. To ease the notation, for each  $y \in N(x)$  set  $\nu_y = |N(x) \cap V|$ . We begin by reformulating Proposition 7 for the case of bribing one single voter:

**Proposition 8** *Let  $V \subseteq C$ , let  $x \in V$ , and let  $\sigma$  be an efficient bribing strategy such that  $\sigma(x) = b > 0$  and  $\sigma(y) = 0$ , for all  $y \in C \setminus \{x\}$ . Then,*

$$\mathbf{r}_P(\sigma) = b \left[ \left( \sum_{y \in N(x)} \frac{1}{\nu_y} \right) - 1 \right].$$

We then move to computing the revenue gained by (efficiently) by bribing a solo non-voter:

**Proposition 9** *Let  $V \subseteq C$ , let  $x \in C \setminus V$ , and let  $\sigma$  be an efficient bribing strategy such that  $\sigma(x) = b > 0$  and  $\sigma(y) = 0$ , for all  $y \in C \setminus \{x\}$ . Then,*

$$\mathbf{r}_P(\sigma) = b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} \text{eval}(k) \right).$$

*Proof* By definition of revenue of a bribing strategy, we have that:

$$r_P(\sigma) = u_P^\sigma - u_P^0 = \left( \sum_{c \in C} \text{P-RATING}(c, \text{eval}^\sigma) \right) - b - \left( \sum_{c \in C} \text{P-RATING}(c, \text{eval}) \right).$$

Given that  $\sigma$  bribes exactly one customer  $x$ , only the P-RATING of customers of  $N(x)$  changes; consequently, we have that:

$$\begin{aligned} r_P(\sigma) &= \sum_{y \in N(x)} \left( \text{P-RATING}(y, \text{eval}^\sigma) - \text{P-RATING}(y, \text{eval}) \right) - b \\ &= \sum_{y \in N(x)} \left( \frac{1}{|N(y) \cap V^\sigma|} \left( \sum_{k \in N(y) \cap V^\sigma} \text{eval}^\sigma(k) \right) + \right. \\ &\quad \left. - \frac{1}{|N(y) \cap V|} \left( \sum_{k \in N(y) \cap V} \text{eval}(k) \right) \right) - b, \end{aligned}$$

where the last equality follows by simply applying the definition of P-RATING. For any  $y \in N(x)$ , let us denote  $|N(y) \cap V^\sigma|$  by  $\nu_y^\sigma$ ; so,  $\nu_y^\sigma = \nu_y + 1$ . Hence, we have the following

$$\begin{aligned} r_P(\sigma) &= \sum_{y \in N(x)} \left( \left[ \frac{1}{\nu_y^\sigma} \sum_{k \in N(y) \cap V^\sigma} eval^\sigma(k) \right] - \left[ \frac{1}{\nu_y} \sum_{k \in N(y) \cap V} eval(k) \right] \right) - b \\ &= \sum_{y \in N(x)} \left( \left[ \frac{1}{\nu_y^\sigma} \sum_{k \in N(y) \cap V} eval(k) \right] + \frac{b}{\nu_y^\sigma} - \left[ \frac{1}{\nu_y} \sum_{k \in N(y) \cap V} eval(k) \right] \right) - b \\ &= b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} eval(k) \right) \end{aligned}$$

and the result follows.

Proposition 8 shows that for some fixed amount, the extent to which a voter is profitable to bribe can be expressed as a function of only the network structure (by this we refer to the topology of the network and the positions of non-voters on the network). Contrary to this, Proposition 9 shows that in order to express the revenue obtained by bribing a non-voter it is also necessary to know the evaluation of the neighbouring customers.

#### 4.2.3 Independence of Bribing Order

We now explore whether the order of bribing customers impacts the resulting revenue, obtaining the following result:

**Proposition 10** *Let  $V \subseteq C$ , let  $x, x' \in C$  be distinct customers, and let  $\sigma$  be an efficient strategy such that  $\sigma(x) = b > 0$ ,  $\sigma(x') = b' > 0$ , and  $\sigma(y) = 0$ , for all  $y \in C \setminus \{x, x'\}$ . No matter whether we bribe  $x$  before  $x'$  or  $x'$  before  $x$ , the resulting cumulative revenue will be the same.*

The proof of Proposition 10, reported in Appendix A, shows that, despite the fact that bribing non-voters transforms the set of voters, we can ignore the order of bribes when evaluating the effect of a strategy. A useful reformulation of this fact is the following statement: given two bribing strategies  $\sigma_1$  and  $\sigma_2$  and two distinct customers  $x, x' \in C$ , if  $\sigma_1$  only bribes customer  $x$  and  $\sigma_2$  only bribes  $x'$ , then  $\mathbf{r}_P(\sigma_1 \circ \sigma_2) = \mathbf{r}_P(\sigma_2 \circ \sigma_1)$ .

#### 4.2.4 A (Non-Optimal) Greedy Algorithm

By Proposition 10, bribing strategies involving multiple customers can be decomposed by looking at strategies bribing a single one, which can then be executed without attention to the order. This fact suggests that a greedy algorithm could be proposed to find optimal strategies. First, select the customer who will yield the highest revenue when bribed the maximal amount allowed by the initial budget and their own evaluation, using the formulas obtained in

Propositions 8 and 9. Note that this could be either a voter or a non-voter. Then, repeat the process until the initial budget is exhausted, or until all individuals on the network who do not have maximal evaluation yield a negative revenue when bribed. Call this procedure the *non-voters greedy algorithm*.

This simple idea, which was shown to yield an optimal bribing strategy in case every customer votes (Corollary 2), does not yield an optimal strategy in the presence of non-voters, as the following proposition shows.

**Proposition 11** *The non-voters greedy algorithm is not optimal.*

The proof of Proposition 11 is reported in Appendix B. What is left open is the question of whether we can find an optimal bribing strategy in polynomial-time, or decide whether there exists a successful manipulation strategy in polynomial-time, or, instead, whether there might be a complexity-theoretic barrier to doing so. As the the following section shows, the latter is true.

#### 4.2.5 Finding optimal manipulation strategies for P-RATING

We now investigate, from a complexity theoretic point of view, the problem of computing a bribing strategy yielding at least some given revenue, when not every customer votes and the restaurant has full knowledge of each customer's position. This, notice, will allow us to determine the existence of both a successful manipulation strategy and an optimal strategy. Firstly, we reformulate the above optimization problem as a decision problem.

BRIBE-NVKL

**Instance:** Network  $(C, E)$ , evaluation  $eval_0$ ,  $\rho \in \mathbb{Q}$

**Yes-Instance:** An instance of BRIBE-NVKL s.t. there exists a strategy  $\sigma$  with  $\mathbf{r}(\sigma) \geq \rho$

Any instance of the above problem should adhere to the usual restrictions of the framework. These are, most importantly, that the initial evaluation is such that every customer  $c \in C$  is adjacent to at least one customer  $c' \in C$  such that  $eval(c') \neq *$  (recall that every customer is adjacent to itself). Also, any strategy  $\sigma$  is such that  $\sum_{c \in C} \sigma(c)$  is at most the initial utility resulting from  $eval_0$ .

In what follows, we show that BRIBE-NVKL is **NP**-complete, by giving a reduction from the known **NP**-complete problem of finding an independent set on 3-regular graphs, aka ISREG(3) (Garey and Johnson, 1990).

Recall that a graph  $G$  is 3-regular if the degree of every vertex is 3, and an independent set of  $G$  is a subset  $X$  of its vertices such that there is no edge of  $G$  joining any pair of vertices in  $X$ . We can now give the following definition:

ISREG(3)

**Instance:** A 3-regular graph  $G$ ,  $k \in \mathbb{N}$

**Yes-Instance:** An instance of ISREG(3) such that  $G$  has an independent set of size at least  $k$

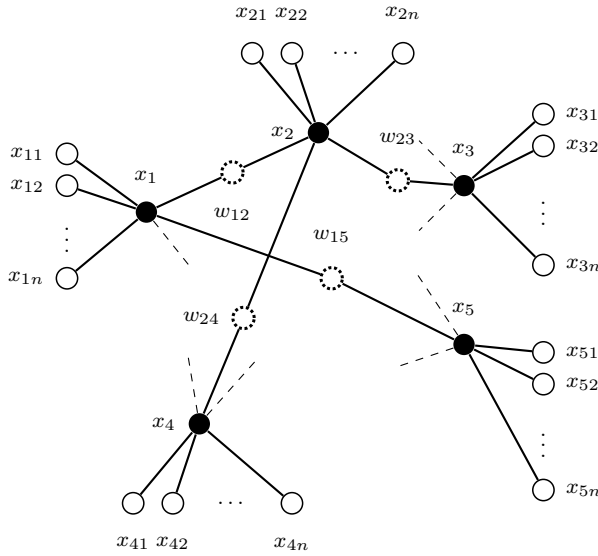
We now prove the main result of this section (part of the proof is reported in Appendix C):

**Theorem 1** BRIBE-NVKL is NP-complete.

*Proof* We begin by giving a reduction from an arbitrary instance of ISREG(3) to an instance of BRIBE-NVKL. That is, given a 3-regular graph  $G$  and  $k \in \mathbb{N}$ , we construct a network  $(C, E)$ , an initial evaluation  $eval_0$ , and  $\rho \in \mathbb{Q}$  such that  $G$  has an independent set of size at least  $k \iff$  there exists a strategy on  $((C, E), eval_0)$  that yields a revenue of at least  $\rho$ . Given a 3-regular graph  $G$ , we define a network of customers as follows:

**Customers** The set  $C$  of customers is composed of *old*, *pendant*, and *edge* customers. For all vertices  $v \in G$ , we create an *old customer*  $v \in C$ , as well as a set of *pendant customers*  $v_1, \dots, v_n \in C$ , where  $n$  is the number of vertices of  $G$ . For each edge  $(u, v)$  of  $G$ , we introduce an *edge customer*  $w_{u,v} \in C$ .

**Network** The network  $E$  relating customers is defined as follows. For each old customer  $v$ , there is an edge  $(v, v_i) \in E$  for  $i = 1, 2, \dots, n$ , connecting it to the related pendant customers. For every edge  $(u, v)$  of  $G$ , we add  $(u, w_{u,v})$  and  $(w_{u,v}, v)$  to  $E$ , relating the two old customers with the corresponding edge customer.



**Fig. 5** The network of customers associated to a portion of a 3-regular graph, focusing on vertices  $\{x_1, x_2, x_3, x_4, x_5\}$  and edges  $\{(x_1, x_2), (x_2, x_3), (x_2, x_4), (x_1, x_5)\}$ . Full nodes are old customers, empty nodes are pendant customers, and dashed nodes are edge customers.

For any such network as constructed above, we can define an initial evaluation  $eval_0$  as follows, where  $0 < \epsilon < 1$  is some value that will be set later in the proof.

- If  $c \in C$  is an old customer then  $eval_0(c) = *$  (non-voter).
- If  $c \in C$  is an edge or pendant customer then  $eval_0(c) = \epsilon$ .

An example of the construction of the customer network and evaluation from a graph  $G$  can be seen in Figure 5.

By the construction of the network, we have that for all  $c \in C$ , the  $P$ -rating( $c, eval_0$ ) =  $\epsilon$  (recall that we assumed  $c \in N(c)$  for all customers). Every customer of the newly constructed customer network contributes  $\epsilon$  to the initial utility of the network and therefore  $u_P^0 = \epsilon(n + n^2 + \frac{3n}{2})$ . Note that the additive factor  $\frac{3n}{2}$  is a consequence of the assumption of 3-regularity of the graph. We now choose  $\epsilon$  so that  $u_P^0 = k$ ; that is, so that

$$\epsilon = \frac{k}{n + n^2 + \frac{3n}{2}}.$$

By assumption the restaurant owner can only make bribes totalling at most  $k$ . Furthermore, note that the initial evaluation is a valid one in that every customer of the network is adjacent to at least one voter. Finally, let

$$\rho = k(1 - \epsilon) \left( \frac{1}{n + 4} + \frac{n + 3}{2} \right) - k.$$

( $\implies$ ) Suppose that our instance  $(G, k)$  of ISREG(3) is a yes-instance; that is, there is a set  $I$  of at least  $k$  vertices such that no two vertices of  $I$  are adjacent in  $G$ . Consider the bribing strategy for  $(C, E)$  (as constructed above) where  $\sigma(c) = 1$ , for every old customer corresponding to some vertex of  $I$ , and  $\sigma(c') = 0$  for all other  $c' \in C$ .

Let us now compute the revenue obtained by  $\sigma$ . Recall that the revenue is equal to the increase in P-RATING of the bribed customers and their neighbourhoods (old, pendant, and edge customers), minus the cost of the bribe. The cumulative increase in rating of bribed old customers is:

$$k \left( \frac{1 + (n + 3)\epsilon}{n + 4} - \epsilon \right) = k \frac{1 - \epsilon}{n + 4}.$$

The cumulative increase in rating of pendant customers is:

$$nk \left( \frac{1 + \epsilon}{2} - \epsilon \right) = nk \frac{1 - \epsilon}{2}.$$

Finally, the increase in rating due to edge customers is:

$$3k \left( \frac{1 + \epsilon}{2} - \epsilon \right) = 3k \frac{1 - \epsilon}{2}.$$

Recall that bribed old customers correspond to an independent set in  $G$ . Summing up, the revenue of strategy  $\sigma$  is:

$$k(1 - \epsilon) \left( \frac{1}{n+4} + \frac{n+3}{2} \right) - k = \rho.$$

Therefore  $((C, E), eval_0, \rho)$  is a yes-instance of NVKL.

( $\Leftarrow$ ) To ease readability we leave this more articulate part to Appendix C.

Finally, it is easy to see that BRIBE-NVKL is in **NP**. Given a customer network  $(C, E)$ , an evaluation  $eval$ , and  $\rho \in \mathbb{Q}$ , we can clearly decide whether a given strategy  $\sigma$  yields a revenue of at least  $\rho$  in polynomial-time (we simply evaluate the strategy).

In summary, we have been able to prove the **NP**-completeness of BRIBE-NVKL by giving a reduction from ISREG(3). This is an important finding, that significantly strengthens the value of personalised rating systems and their resistance to bribery, as we have demonstrated that we cannot compute an optimal bribing strategy, nor any strategy guaranteeing at least a given reward, in a reasonable amount of time; that is, of course, unless  $\mathbf{P} = \mathbf{NP}$ .

#### 4.3 All vote, unknown network

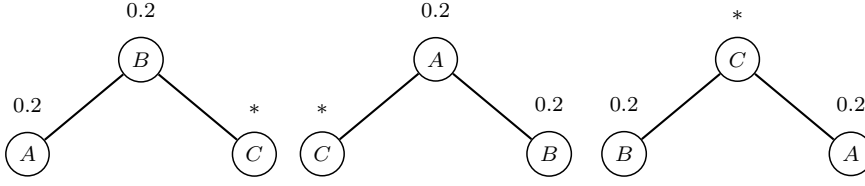
We now move onto studying the more complex case of an unknown network, starting from a situation in which every customer expresses an opinion. Surprisingly, we are able to show that no bribing strategy is profitable (in expectation), and hence P-RATING is bribery-proof in this case.

We begin by assuming that the restaurant knows the structure of the network, but not the position of each customer. Formally, the restaurant knows  $E$ , but considers any permutation of the customers in  $C$  over  $E$  as being possible. Let us define the expected revenue of a strategy  $\sigma$  over a given network  $E$  as the average over all possible permutations of customers:  $\mathbb{E}[\mathbf{r}_P(\sigma)] = \sum \frac{1}{n!} [u_\rho^\sigma - u_\rho^0]$ , where we abuse notation by writing  $u_\rho^\sigma$  as  $u_\rho^\sigma$  under permutation  $\rho$  over the network  $E$ . What we are able to show is that, in expectation, no strategy is more profitable than  $\sigma^0$ .

**Proposition 12** *Let  $V = C$ , let the network structure of  $E$  be known but not the relative positions of customers on  $E$ . Then  $\mathbb{E}[\mathbf{r}_P(\sigma)] = 0$  for all strategies  $\sigma$ .*

*Proof* Let  $|C| = n$ . We first show that the result for any strategy  $\sigma$  that bribes a single customer  $\bar{c}$  is zero, and the general statement follows from the linearity of  $\mathbb{E}[\mathbf{r}(\sigma)]$  and the additivity of the revenue. Let therefore  $\sigma$  be a single bribe to a customer. By Proposition 6 we can compute the expected revenue of  $\sigma$  for each permutation  $\rho$  of customers  $C$  on the network, assuming that each





**Fig. 6** Customers permutations in Example 6.

permutation is equally likely:

$$\mathbb{E}[\sigma] = \sum_{\rho} \frac{1}{n!} (u_{\rho}^{\sigma} - u_{\rho}^0) = \sum_{\rho} \frac{1}{n!} (w_{\rho(\bar{c})} - 1) \sigma(\bar{c}) \quad (5)$$

$$= \sum_{c \in C} \frac{(n-1)!}{n!} (w_c - 1) \sigma(c) = \frac{(n-1)!}{n!} \sum_{c \in C} (w_c - 1) = 0 \quad (6)$$

Where  $\bar{c}$  in equation (5) is the only customer receiving a bribe under permutation  $\rho$ , and the equation (6) follows from the observation that, under the assumption that every customer voted,  $\sum_c w_c = |C|$  and hence  $\sum_c (w_c - 1) = 0$ .

Hence, if we assume a uniform probability over all permutations of customers on the network, a straightforward consequence of Proposition 12 concludes that it is not profitable (in expectation) to bribe any customer.

**Corollary 3** *If  $V = C$  and the network is unknown, then no bribing strategy for P-RATING is profitable in expected return.*

#### 4.4 Non voters, unknown network

Unlike the case of  $V = C$ , in this case it is possible to define bribing strategies for P-RATING that are profitable (in expected return).

*Example 6* Let  $C = \{A, B, C\}$ , and the initial evaluation  $eval(A) = eval(B) = 0.2$  and  $eval(C) = *$ . Assume that the structure of the network is known, but the position of the individuals is not. Let the three possible network positions (without counting the symmetries) be depicted in Figure 6. Let  $\sigma(B) = 0.2$  and  $\sigma(A) = \sigma(C) = 0$ . In the first case

$$\begin{aligned} \mathbf{r}_P^1(\sigma) &= \text{P-RATING}(A) + \dots + \text{P-RATING}(C) - 0.2 - u_P^0 \\ &= 0.3 + 0.3 + 0.4 - 0.2 - 0.6 \\ &= 0.2, \end{aligned}$$

in the second case  $\mathbf{r}_P^2(\sigma) = 0$ , whilst in the third case

$$\mathbf{r}_P^3(\sigma) = 0.4 + 0.3 + 0.2 - 0.2 - 0.6 = 0.1.$$

Therefore, P-RATING is not bribery-proof (in expectation) in the presence of non-voters when the network is unknown. Interesting computational problems therefore open up in this setting. Finding optimal bribing strategies in the presence of non-voters is already computationally hard when the network is known, as shown by Theorem 1. We therefore focus on the problem of finding profitable (but not necessarily optimal) strategies under the assumption of incomplete knowledge of the network.

We start from the assumption that the restaurant knows the network topology of the network but does not know the location of each customer on the network. We also assume that the restaurant owner knows the overall utility, and we ask the question of whether it is possible to compute in polynomial-time an assignment of evaluations to customers such that the resulting utility of this assignment is equal to the one the restaurant owner knows about.

More formally, we assume that the restaurant knows the topology of the network of influence, that is,  $E \subseteq N \times N$ , where here  $N$  is the set of (unnamed) customers. So, we can think of our network as a graph where there are no customer names on the vertices. The restaurant knows the evaluation function  $eval : C \rightarrow Val \cup \{*\}$ , however it does not know where each customer name sits within the network. That is, the restaurant owner is missing a bijection between  $N$  and  $C$ . We define the following problem:

UTILITY PLACEMENT (UPLACE)

**Instance:** An unlabelled network  $(N, E)$ , evaluation  $eval$ , target utility  $b \in \mathbb{Q}$

**Yes-Instance:** An instance of UPLACE s.t. there is a bijection  $\rho$  from  $N$  to  $C$  such that (a) every customer is adjacent to at least one voting customer on the network  $E$  labelled by  $\rho$ , and (b) the utility resulting from the assignment is at least  $b$ .

We now show that UPLACE is **NP**-complete by giving a simple reduction from the vertex dominating set problem, which is a known **NP**-hard problem defined as follows (Garey and Johnson, 1990):

VERTEX DOMINATING SET (VDS)

**Instance:** A graph  $G$  with  $n$  vertices,  $k \in \mathbb{N}$  such that  $k \leq n$

**Yes-Instance:** An instance of VDS s.t. there is a subset  $S$  of  $k$  vertices of  $G$  such that every vertex not in  $S$  is adjacent to at least one vertex of  $S$ .

**Proposition 13** UPLACE is **NP**-complete.

*Proof* Given any instance  $(G, k)$  of VDS we construct a customer network  $C = G$ . We define  $eval = (1, \dots, 1) \in \mathbb{R}^k$  and let  $b = k$ . Let  $(C, eval, b)$  be an instance of UPLACE. If  $(G, k)$  has a vertex dominating set of size  $k$  then clearly the assignment of 1 to each customer corresponding to a vertex of the dominating set would give us a utility of  $b \geq k$  in  $C$ . Conversely, if  $G$  has no vertex dominating set of size  $k$  then there is no assignment of evaluations

1's to customers such that condition (a) in the definition of UPLACE holds. That is,  $(C, eval, b)$  is not a yes-instance of UPLACE. Since VDS is **NP**-hard it follows that UPLACE is **NP**-hard. Given an instance of UPLACE we can clearly check in non-deterministic polynomial time whether the utility of  $C$  with respect to  $eval$  is  $\geq b$ . We can therefore deduce that UPLACE is in **NP** thus concluding the proof.

In conclusion, despite the fact that P-RATING is theoretically manipulable with full knowledge of the network and even more so in presence of non-voters, the lack of knowledge of the underlying network structure contributes to its non-manipulability in practice, even when the network structure is known but the attacker is faced with the problem of deducing the individuals' locations.

#### 4.5 Stocktaking

We have analysed the manipulability of P-RATING with full knowledge of the underlying network, with and without non-voters. We have seen that in this case there are situations in which the rating system is manipulable. However, while we were able to give a polynomial-time algorithm for computing optimal bribing strategies when everyone voted (Corollary 2), the presence of non-voters makes this problem **NP**-complete (Theorem 1). Finding profitable, albeit suboptimal, strategies is easy, and we provided closed formulas for the computation of the generated revenue (Propositions 8 and 9).

Finally, we investigated the case of incomplete knowledge of the network, showing that P-RATING is bribery-proof if everyone votes (Proposition 12), but that reconstructing the evaluation of customers when their location on the network is unknown is an intractable computational problem (Proposition 13).

### 5 Discussing the Underlying Assumptions

Our model is built upon a number of simplifying assumptions, some of which require non-trivial extensions, while some others can be generalised without affecting the overall applicability of the framework. Among the former we include: the effect of time on agents' decisions, both of the customers, as their view of the restaurant and of the other customers may change after observing the quality of the service provided, and the restaurant itself, as the knowledge of the network is bound to change with new information coming in; the presence of multiple restaurants, therefore giving rise to a game, albeit not necessarily fully competitive; the graph dynamics, as agents' might reconsider their friendships as well. There are however simplifying assumptions that can be lifted without altering the main message: i.e., that uncertainty of the underlying network structure makes P-RATING resistant to manipulation. The following is a non-exhaustive list of the main ones.

*Directed weighted graphs* We assumed undirected connections where the relative importance of customers in their neighbours' decisions is the same. This can be generalised to networks  $E \subseteq C \times C$  that are not necessarily symmetric. This means that each customer  $c$  will have a neighbour  $N^{in}(c)$  of incoming connections that influence them. At the same time,  $c$  will influence a set  $N^{out}(c)$  of outgoing connections, possibly different than  $N^{in}(c)$ . On top of that, for each  $c$ , members  $k$  of  $N^{in}(c)$  also have a relative importance, which we denote as  $I_{kc}$ , meaning how important  $k$ 's evaluation is for  $c$ .  $I_{kc}$  comes with some natural assumptions, in particular that  $\forall k, c : I_{kc} \in [0, 1]$  and  $\sum_{k \in N^{in}(c)} I_{kc} = 1$ . Observe that this more general model allows us to raise arbitrarily the relative weight of each loop  $I_{cc}$ , lifting our assumption that a customer is as influenced by his or her neighbours' opinions as by his or her own one.

Despite the fact that a customers' importance is normalised, the calculation of P-RATING needs to take into account the fact that some of the incoming connections to a customer may not be from voters:

$$\text{P-RATING}(c, eval) = \sum_{k \in N^{in}(c) \cap V} \frac{I_{kc} eval(k)}{\sum_{k \in N^{in}(c) \cap V} I_{kc}}$$

Considering the degrees of importance that we just introduced, the customer influence weights from Definition 5 can be redefined as follows:

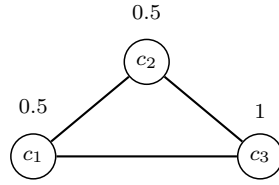
$$w^\sigma(c) = \sum_{k \in N^{out}(c)} \frac{I_{ck}}{\sum_{k' \in N^{in}(k) \cap V^\sigma} I_{k'k}}$$

in order to obtain, once again, by similar calculation as in the proof of Proposition 1, that  $u_p^\sigma = \sum_{c \in V^\sigma} w_c^\sigma \times eval^\sigma(c) - \sum_{c \in C} \sigma(c)$ . In particular, whenever  $V = C$ , we obtain that  $w_c^\sigma = \sum_{k \in N^{out}(c)} I_{ck}$ , i.e., the influence weight of a customer is its weighted out-degree.

In conclusion, our results and techniques can be applied to the more general case of weighted directed networks, adapting the calculations as described above. In particular, we remark that the hardness of bribery established in Section 4 carries through to the weighted setting.

*Aggregation rules* Customers' personalised ratings aggregate the average of their voting connections. The choice of the average as aggregation rule is a simplifying one and the previous paragraph has shown that it can be naturally lifted to weighted graphs, and thus to weighted averages. We also note how average-based aggregators have a long tradition in social network analysis and a number of impactful contributions have used it to formulate plausible computational models. Notable examples include the polarization model of Axelrod (1997), the belief mixtures of Deffuant et al. (2000) and the opinion aggregation dynamics model in Hegselmann and Krause (2002).

There are a number of different approaches one can take towards aggregating evaluations from trusted peers. Garcin et al. (2009) suggested the use of median voter (the median evaluation among the given ones) to aggregate



**Fig. 7** Personalised ratings with median aggregator is *not* bribery proof

reputation feedback, emphasising the fact that the median voter is strategy-proof. Although we acknowledge their point, we can show that the use of the median in the calculation of the P-RATING is not bribery proof either. Figure 7 gives an example of a fully connected network where the median vote is 0.5 for all customers. However any bribe of 0.5 to any of the customers voting 0.5 already brings the utility from 1.5 to 3. Notice how this shows that, unlike in the case of the average, the median is not even bribery proof with O-RATING when all customers vote.

Similar points can be made for variants of personalised rating built on the mode, i.e., the most frequent appearing evaluation and, as shown, weighted average, concluding the list of four most common aggregators for reputation feedback (Garcin et al., 2009).

*Linear responses* Our framework assumes linearity of behaviour in various ways, notably: (i) customers' ratings correspond to their propensity to go to the restaurant, and (ii) bribe  $\sigma(c)$  affects evaluation  $eval(c)$  linearly. These assumptions can be generalised by multiplicative factors, such as an average price  $R$  paid at the restaurant, a propensity to use the service, and a customer price.

Let us begin with the customer price, encoding the response to bribes. Let  $\mu_c \in [0, \infty)$  encode the relative effect of a bribed sum on a customer. This is such that, for any strategy  $\sigma$ ,

$$eval^\sigma(c) = \min(1, eval(c) + \mu_c \sigma(c)).$$

So, the old evaluation is updated with the bribe multiplied by the customer response, capped at 1. Note that  $\mu_c = 0$  encodes customers who do not accept bribes. Our methods also work in this case. To showcase this fact, we can generalise Proposition 8 to include customers' prices by observing that

$$\mathbf{r}_P(\sigma) = \sum_{c \in C} (w_c \mu_c - 1) \sigma(c)$$

which means that a bribe is profitable if and only if it is given to customers with  $w_c \mu_c > 1$ .

Regarding the propensity to use the service, we can associate to each customer a propensity value  $\pi_c \in [-1, 1]$ , which encodes their reaction to their

own P-RATING which is then reflected on the utility perceived by the restaurant. With this extension the utility of the restaurant is calculated as follows:

$$u_P = \sum_{c \in C} (\text{P-RATING}(c) + \delta(c))$$

where

$$\delta(c) = \begin{cases} \pi_c \text{P-RATING}(c) & \text{if } \pi_c \leq 0 \\ \pi_c (1 - \text{P-RATING}(c)) & \text{if } \pi_c > 0. \end{cases}$$

In other words, each customer will go to the restaurant based on their P-RATING, discounted by their own reaction to it—increasing the original P-RATING if  $\pi_c$  is positive, and decreasing it otherwise. Profitable bribing strategies can be found by similar calculations to those in previous sections, taking the discounting factor into account.

*Various other constraints* Another assumption that can be removed without loss of generality is the reflexivity of connections, i.e., the fact that a customer is necessarily influenced by its own personal opinion. Notice this is equivalent to setting  $I_{cc} = 0$  in the previously discussed case of weighted networks. Likewise, the same reasoning applies for customers that are already voters, and thus we might want to have them place higher importance on their own personal opinion.

Finally, the constraint of a non-voter being connected to at least one voter might seem demanding. To tackle this with more generality, one could define a  $k$ -neighbourhood notion of influence—with P-RATING calculations lifted accordingly—or simply provide these customers with the O-RATING only. This would also provide a solution to the cold-start problem affecting any personalised approach to ratings, driving customers to submit reviews until a sufficient number of voters filled most neighbourhoods. Once again, this constraint is a technical assumption that heavily facilitates calculations, but could in general be dispensed with.

## 6 Conclusions and Future Work

We introduced Personalised Rating, a network-based rating system which generalises the commonly used Objective Rating, and analysed its resistance to external attacks under various conditions.

Our results show that the introduction of an underlying network structure to constrain the available information has clear advantages in terms of the trustworthiness of the resulting rating. In particular, we have shown that while P-RATING and O-RATING are not bribery-proof in general, the former has a clear upper bound on the profit that can be made by external influence (Propositions 1 and 2), even if the external attacker has access to the underlying network. The reliability of P-RATING increases with the incomplete knowledge of the attacker. We have shown that as long as the restaurant

owner does not know the customers' positions in a network — and, therefore, even if they know the full structure of the network — P-RATING is not manipulable when all customers vote, (Proposition 12) which, coupled with our findings in the presence of non-voters, shows a clear advantage to the use of network-based rating systems. We have also looked at situations in which manipulation is possible in theory, but may, in general, be computationally infeasible in practice. Specifically, unlike the case of O-RATING, the problem of manipulation under P-RATING with full knowledge of the network is intractable, as we know from showing that BRIBE-NVKL is **NP**-complete. This, we find, is a major strengthening for the practical applicability of the personalised rating framework.

It has to be emphasised that our results are confined to worst-case complexity analysis and it is therefore necessary to analyse alternative methods for manipulation. These include studying the parameterized complexity analysis of various subproblems (Faliszewski and Niedermeier, 2014). Alternatively, we can investigate the problem of approximate bribing strategies. It could be possible that something can still be salvaged of the *P*-greedy approach, or, that we find it is still computationally hard to even compute a bribing strategy that yields a revenue that is within some factor of a given amount. We saw through our Proposition 11 that whilst not yielding the optimal amount of revenue, we can still compute a profitable return. We can approach this question from a slightly less formal but nevertheless important angle and seek to obtain a number of experimental results concerning the performance of greedy algorithms. The results on Personalised Rating demonstrate the impact of incomplete knowledge in deterring manipulation. This suggests the importance of devising ways to *induce* uncertainty for the attacker, with respect to the network structure and the individuals' evaluations. The study of probabilistic aggregation methods, along the lines of Barberà (1979) and Fishburn (1984), is therefore an important future research direction.

## Acknowledgments

Umberto Grandi acknowledges the support of the ANR JCJC project SCONE (ANR 18-CE23-0009-01).

## References

- Abraham I, Dolev D, Gonen R, Halpern J (2006) Distributed computing meets game theory: robust mechanisms for rational secret sharing and multiparty computation. In: Proceedings of the 25th ACM symposium on Principles of Distributed Computing (PODC)
- Alon N, Feldman M, Lev O, Tennenholtz M (2015) How robust is the wisdom of the crowds? In: Proceedings of the 24th International Conference on Artificial Intelligence (IJCAI)

- Andersen R, Borgs C, Chayes J, Feige U, Flaxman A, Kalai A, Mirrokni V, Tennenholtz M (2008) Trust-based recommendation systems: An axiomatic approach. In: Proceedings of the 17th International Conference on World Wide Web (WWW)
- Apt KR, Markakis E (2014) Social networks with competing products. *Fundamenta Informaticae* 129(3):225–250
- Axelrod R (1997) The dissemination of culture: A model with local convergence and global polarization. *Journal of Conflict Resolution* 41(2):203–226
- Bahar G, Smorodinsky R, Tennenholtz M (2016) Economic recommendation systems (one page abstract). In: Proceedings of the 2016 ACM Conference on Economics and Computation (EC)
- Barberà S (1979) Majority and positional voting in a probabilistic framework. *Review of Economic Studies* 46:379–389
- Baumeister D, Erdélyi G, Rothe J (2011) How hard is it to bribe the judges? A study of the complexity of bribery in judgment aggregation. In: Proceedings of the Second International Conference on Algorithmic Decision Theory (ADT)
- Bobadilla J, Hernando A, Ortega F, Gutiérrez A (2013) Recommender systems survey. *Knowledge-Based Systems* 46:109 – 132
- Bredereck R, Chen J, Hartung S, Kratsch S, Niedermeier R, Suchý O, Woeginger GJ (2014) A multivariate complexity analysis of lobbying in multiple referenda. *Journal of Artificial Intelligence Research* 50:409–446
- Bredereck R, Faliszewski P, Niedermeier R, Talmon N (2016) Large-scale election campaigns: Combinatorial shift bribery. *Journal of Artificial Intelligence Research (JAIR)* 55:603–652
- Brill M, Conitzer V, Freeman R, Shah N (2016) False-name-proof recommendations in social networks. In: Proceedings of the 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)
- Christian R, Fellows M, Rosamond F, Slinko A (2007) On complexity of lobbying in multiple referenda. *Review of Economic Design* 11(3):217–224
- Conitzer V, Walsh T (2016) Barriers to manipulation in voting. In: Brandt F, Conitzer V, Endriss U, Lang J, Procaccia AD (eds) *Handbook of Computational Social Choice*, Cambridge University Press, chap 6, pp 127–145
- Conitzer V, Immorlica N, Letchford J, Munagala K, Wagman L (2010) False-Name-Proofness in Social Networks. In: 6th International Workshop on Internet and Network Economics (WINE)
- Conte R, Paolucci M (2002) *Reputation in Artificial Societies: Social Beliefs for Social Order*. Kluwer
- Conte R, Paolucci M, Sabater-Mir J (2008) Reputation for innovating social networks. *Advances in Complex Systems* 11(2):303–320
- Dash RK, Jennings NR, Parkes DC (2003) Computational-mechanism design: A call to arms. *IEEE Intelligent Systems* 18(6):40–47
- Deffuant G, Neau D, Amblard F, Weisbuch G (2000) Mixing beliefs among interacting agents. *Adv Complex Syst* 3(1–4):87–98
- Faliszewski P, Niedermeier R (2014) Parameterization in computational social choice. In: Kao MY (ed) *Encyclopedia of Algorithms*, Springer Berlin



Heidelberg

- Faliszewski P, Hemaspaandra E, Hemaspaandra LA (2009) How hard is bribery in elections? *Journal of Artificial Intelligence Research (JAIR)* 35:485–532
- Feigenbaum J, Parkes DC, Pennock DM (2009) Computational challenges in e-commerce. *Communications of the ACM* 52:70–74
- Fishburn P (1984) Probabilistic social choice based on simple voting comparisons. *Review of Economic Studies* 51(4):683–692
- Garcin F, Faltings B, Jurca R, Joswig N (2009) Rating aggregation in collaborative filtering systems. In: *Proceedings of the Third ACM Conference on Recommender Systems (RecSys)*
- Garey MR, Johnson DS (1990) *Computers and Intractability; A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., New York, NY, USA
- Gibbard A (1973) Manipulation of voting schemes: A general result. *Econometrica* 41(4):587–601
- Grandi U (2017) Social choice and social networks. In: Endriss U (ed) *Trends in Computational Social Choice*, AI Access, chap 9, pp 169–184
- Grandi U, Turrini P (2016) A network-based rating system and its resistance to bribery. In: *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI)*
- Grandi U, Stewart J, Turrini P (2018) The complexity of bribery in network-based rating systems. In: *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI)*
- Grandi U, Grossi D, Turrini P (2019) Negotiable votes. *Journal of Artificial Intelligence Research (JAIR)* 64:895–929
- Hegselmann R, Krause U (2002) Opinion dynamics and bounded confidence: models, analysis and simulation. *Journal of Artificial Societies and Social Simulation* 5(3)
- Helpman E, Persson T (1998) Lobbying and legislative bargaining. Working Paper 6589, National Bureau of Economic Research
- Kurokawa D, Lev O, Morgenstern J, Procaccia AD (2015) Impartial peer review. In: *Proceedings of the 24th International Joint Conference on Artificial Intelligence (IJCAI)*
- de Lange B, Fernbach PM, Lichtenstein DR (2016) Navigating by the stars: Investigating the actual and perceived validity of online user ratings. *Journal of Consumer Research* 42(6):817–833
- Lev O, Tennenholtz M (2017) Group recommendations: Axioms, impossibilities, and random walks. In: *Proceedings of the 16th Conference on Theoretical Aspects of Rationality and Knowledge (TARK)*
- Pinyol I, Sabater-Mir J (2013) Computational trust and reputation models for open multi-agent systems: a review. *Artificial Intelligence Review* 40(1):1–25
- Sabater J, Sierra C (2005) Review on computational trust and reputation models. *Artificial Intelligence Review* 24(1):33–60
- Satterthwaite MA (1975) Strategy-proofness and Arrow’s conditions: Existence and correspondence theorems for voting procedures and social welfare functions. *Journal of Economic Theory* 10(2):187 – 217

- Sierra C (2004) Agent-mediated electronic commerce. *Autonomous Agents and Multi-Agent Systems* 9(3):285–301
- Simon S, Apt KR (2015) Social network games. *Journal of Logic and Computation* 25(1):207–242
- Streitfeld D (2016) Online reviews: researchers give them a low rating. *The New York Times* URL <https://www.nytimes.com/2016/06/09/technology/online-reviews-researchers-give-them-a-low-rating.html>
- Tennenholtz M (2008) Game-theoretic recommendations: some progress in an uphill battle. In: *Proceedings of the 7th International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)*
- Todo T, Conitzer V (2013) False-name-proof matching. In: *International conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*
- Waggoner B, Xia L, Conitzer V (2012) Evaluating resistance to false-name manipulations in elections. In: *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence (AAAI)*

## Appendix A - Proof of Proposition 10

**Proposition** *Let  $V \subseteq C$ , let  $x, x' \in C$  be distinct customers, and let  $\sigma$  be an efficient strategy such that  $\sigma(x) = b > 0$ ,  $\sigma(x') = b' > 0$ , and  $\sigma(y) = 0$ , for all  $y \in C \setminus \{x, x'\}$ . No matter whether we bribe  $x$  before  $x'$  or  $x'$  before  $x$ , the resulting cumulative revenue will be the same.*

*Proof* We show that, if  $\sigma$  is a strategy bribing only two customers  $x$  and  $x'$ , the order of bribes does not influence the resulting revenue. We reason by cases. If  $x, x' \in V$ , i.e., both bribed customers are voters, the independence of the order is a straightforward consequence of Proposition 7.

Let us then consider the case of bribing two non-voters. Let therefore  $x, x' \in C \setminus V$ , and let  $\sigma$  be an efficient strategy such that  $\sigma(x) = b > 0$ ,  $\sigma(x') = b' > 0$ , and  $\sigma(y) = 0$ , for all  $y \in C \setminus \{x, x'\}$ . Suppose that we bribe  $x$  first. By Proposition 9, the revenue  $r_1$  acquired can be computed as follows

$$r_1 = b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} eval(k) \right),$$

where we recall that for each  $y \in N(x) \cup N(x')$ ,  $\nu_y$  is defined as  $|N(y) \cap V|$ . Denote the evaluation obtained after bribing  $x$  by  $eval'$  and the resulting set of voters by  $V'$ ; in particular,  $eval'$  differs from  $eval$  only on  $x$  and  $V' = V \cup \{x\}$ . For each  $y \in N(x')$ , define  $\nu'_y$  as  $|N(y) \cap V'|$ . So, again by Proposition 9, the additional revenue  $r_2$  acquired by bribing  $x'$  after  $x$  is

$$r_2 = b' \sum_{y \in N(x')} \left( \frac{1}{\nu'_y + 1} - \frac{1}{|N(x')|} \right) - \sum_{y \in N(x')} \left( \frac{1}{\nu'_y(\nu'_y + 1)} \sum_{k \in N(y) \cap V'} eval'(k) \right).$$

Note that any  $\nu'_y$  depends upon whether  $y \in N(x') \cap N(x)$  or whether  $y \in N(x') \setminus N(x)$ : in the former case,  $\nu'_y = \nu_y + 1$ ; and in the latter case,  $\nu'_y = \nu_y$ . Consequently

$$\begin{aligned}
r_2 &= b' \sum_{y \in N(x') \cap N(x)} \left( \frac{1}{\nu_y + 2} - \frac{1}{|N(x')|} \right) + b' \sum_{y \in N(x') \setminus N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x')|} \right) \\
&\quad - \sum_{y \in N(x') \cap N(x)} \left( \frac{1}{(\nu_y + 1)(\nu_y + 2)} \sum_{k \in N(y) \cap V'} \text{eval}'(k) \right) \\
&\quad - \sum_{y \in N(x') \setminus N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V'} \text{eval}'(k) \right) \\
&= b' \sum_{y \in N(x') \cap N(x)} \left( \frac{1}{\nu_y + 2} - \frac{1}{|N(x')|} \right) + b' \sum_{y \in N(x') \setminus N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x')|} \right) \\
&\quad - \sum_{y \in N(x') \cap N(x)} \left( \frac{1}{(\nu_y + 1)(\nu_y + 2)} \left( b + \sum_{k \in N(y) \cap V} \text{eval}(k) \right) \right) \\
&\quad - \sum_{y \in N(x') \setminus N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} \text{eval}(k) \right).
\end{aligned}$$

Now we have that

$$\begin{aligned}
r_1 + r_2 &= b \left[ \sum_{y \in N(x) \cap N(x')} \left( \frac{1}{\nu_y + 2} - \frac{1}{|N(x)|} \right) + \sum_{y \in N(x) \setminus N(x')} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) \right] \\
&\quad + b' \left[ \sum_{y \in N(x') \cap N(x)} \left( \frac{1}{\nu_y + 2} - \frac{1}{|N(x')|} \right) + \sum_{y \in N(x') \setminus N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x')|} \right) \right] \\
&\quad - \sum_{y \in N(x) \cap N(x')} \left( \frac{1}{\nu_y(\nu_y + 2)} \sum_{k \in N(y) \cap V} \text{eval}(k) \right) \\
&\quad - \sum_{y \in N(x) \setminus N(x')} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} \text{eval}(k) \right) \\
&\quad - \sum_{y \in N(x') \setminus N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} \text{eval}(k) \right).
\end{aligned}$$

We can now conclude the first part of the proof, since by the symmetry of the expression for  $r_1 + r_2$ , were we to first bribe  $x'$  and then  $x$ , we would get exactly the same revenue.

Now we consider the case where we compare bribing a non-voter  $x$  and then a voter  $x'$  with bribing the voter  $x'$  and then the non-voter  $x$ . Let therefore  $x \in C \setminus V$ , let  $x' \in V$ , and let  $\sigma$  be an efficient strategy such that  $\sigma(x) = b > 0$ ,  $\sigma(x') = b' > 0$ , and  $\sigma(y) = 0$ , for all  $y \in C \setminus \{x, x'\}$ . Suppose that we

first bribe the voter  $x'$  with  $b'$ , to acquire revenue  $r_1$ , before bribing the non-voter  $x$  with  $b$ , to acquire revenue  $r_2$ . Denote the evaluation obtained from  $eval$  by bribing  $x'$  by  $eval'$ ; so,  $eval'$  is identical to  $eval$  except on  $x'$  where  $eval'(x') = eval(x') + b'$ . Note also that the set of voters has not changed. By Proposition 8 and Proposition 9, we obtain

$$\begin{aligned}
r_1 + r_2 &= b' \left[ \left( \sum_{y \in N(x')} \frac{1}{\nu_y} \right) - 1 \right] + b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) \\
&\quad - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} eval'(k) \right) \\
&= b' \left[ \sum_{y \in N(x') \cap N(x)} \left( \frac{1}{\nu_y + 1} \right) + \sum_{y \in N(x') \setminus N(x)} \left( \frac{1}{\nu_y} \right) - 1 \right] \\
&\quad + b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) \\
&\quad - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} eval(k) \right).
\end{aligned}$$

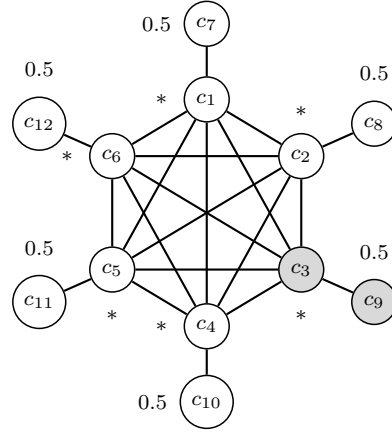
Now suppose that we first bribe the non-voter  $x$  with  $b$ , to acquire revenue  $s_1$ , before bribing the voter  $x'$  with  $b'$ , to acquire revenue  $s_2$ . By Proposition 9, we have that

$$s_1 = b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} eval(k) \right).$$

Denote the evaluation obtained from  $eval$  by bribing  $x$  by  $eval'$ , denote the set of voters after bribing  $x$  by  $V'$ , and for every  $y \in N(x')$ , set  $\nu'_y$  as  $|N(y) \cap V'|$ ; so,  $eval'$  is identical to  $eval$  except on  $x$ , where  $eval'(x) = b$ , and  $V' = V \cup \{x\}$ . By Proposition 8, we have that

$$\begin{aligned}
s_1 + s_2 &= b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} eval(k) \right) \\
&\quad + b' \left[ \left( \sum_{y \in N(x')} \frac{1}{\nu'_y} \right) - 1 \right] \\
&= b \sum_{y \in N(x)} \left( \frac{1}{\nu_y + 1} - \frac{1}{|N(x)|} \right) - \sum_{y \in N(x)} \left( \frac{1}{\nu_y(\nu_y + 1)} \sum_{k \in N(y) \cap V} eval(k) \right) \\
&\quad + b' \left[ \left( \sum_{y \in N(x') \cap N(x)} \frac{1}{\nu_y + 1} \right) + \left( \sum_{y \in N(x') \setminus N(x)} \frac{1}{\nu_y} \right) - 1 \right]
\end{aligned}$$

Hence, we have that  $r_1 + r_2 = s_1 + s_2$  and the result follows.



**Fig. 8** A network and an evaluation for which the non-voter greedy algorithm does not yield an optimal bribing strategy. The numbers or the symbol \* above or left of each node, indicate the initial evaluation of the corresponding customer.

## Appendix B - Proof of Proposition 11

**Proposition** *The non-voters greedy algorithm is not optimal.*

*Proof* We present the following counter-example to the optimality of the non-voter greedy algorithm. Consider a 6-clique  $\mathcal{X}$  of non-voters, each connected to an associated voter with evaluation  $\frac{1}{2}$  as is depicted in Figure 8.

The initial utility of the network is as follows:

$$u_P^0 = \sum_{c \in C} \text{P-RATING}(c, \text{eval}) = \sum_{c \in \mathcal{X}} \frac{1}{2} + \sum_{c \in C \setminus \mathcal{X}} \frac{1}{2} = 6.$$

Suppose that we bribe some clique customer  $x \in \mathcal{X}$  its maximal amount of 1. By Proposition 9 the revenue gained is a monotonic formula in the amount bribed, and in this case it is equal to:

$$= \sum_{y \in N(x)} \left( \frac{1}{2} - \frac{1}{7} \right) - \sum_{y \in N(x)} \left( \frac{1}{2} \sum_{k \in N(y) \cap V} \frac{1}{2} \right) = \frac{3}{4}.$$

Alternatively, suppose we bribe some non-clique customer  $x \in C \setminus \mathcal{X}$  its maximal amount  $\frac{1}{2}$ . The revenue gained from doing so, by Proposition 8, is:

$$\frac{1}{2} \left[ \left( \sum_{y \in N(x)} \frac{1}{\nu_y} \right) - 1 \right] = \frac{1}{2} \left[ \left( \sum_{y \in N(x)} 1 \right) - 1 \right] = \frac{1}{2}.$$

Since  $\frac{3}{4} > \frac{1}{2}$ , a greedy algorithm would first bribe a clique-customer the amount 1.

Consider now any unbribed non-voter and the pair it makes with the unique voter it is adjacent to, as is depicted in Figure 8 by the vertices coloured in gray. As long as the non-voter remains unbribed, we can bribe the voter by  $\frac{1}{2}$  and gain an increase in revenue. Consequently, the non-voter greedy algorithm will bribe at least one customer of each pair of non-voter in the clique and corresponding voter. Since there are five remaining such pairs, the least amount that the greedy algorithm bribes, from the first step on, is  $\frac{1}{2}$  per pair. To see this, observe that all cutomers evaluations are above 0.5, hence any profitable bribe to a non-voter must also be above 0.5. Therefore, the revenue produced by the strategy computed by the greedy algorithm is bounded by:

$$r_P(\sigma) \leq 12 - 6 - 1 - \frac{5}{2} = \frac{5}{2}.$$

This is due to the fact that the maximum utility of the network after executing any strategy  $\sigma$  is  $12 = |C|$ , the initial utility of the network is 6, the amount 1 is spent on the first bribe, and at least  $\frac{1}{2}$  is spent on bribing the remaining five voter/non-voter pairs, as established above. Consider now the strategy  $\sigma'$  that bribes all non-clique customers fully:

$$\sigma'(c_6)=\sigma'(c_7)=\sigma'(c_8)=\sigma'(c_9)=\sigma'(c_{10})=\sigma'(c_{11})=\sigma'(c_{12})=\frac{1}{2}$$

and  $\sigma'(x) = 0$  for all other customers  $x \in \mathcal{X}$ . The revenue gained by playing this strategy is

$$r_P(\sigma') = \sum_{c \in C} \text{P-RATING}(c, eval) - u_P^0 - 3 = 12 - 6 - 3 = 3.$$

### Appendix C - Proof of Theorem 1, left-to-right direction

We begin by showing two technical lemmas:

**Lemma 2** *Let  $\sigma$  be an optimal bribing strategy. Let  $X$  be the set of customers for which  $eval(x) < eval^\sigma(x) < 1$ . Let  $v_y = |N(y) \cap V|$  for any  $y \in C$ . For all  $x, y \in X$  we have that:*

$$\sum_{z \in N(x)} \frac{1}{v_z} = \sum_{z \in N(y)} \frac{1}{v_z}.$$

*Proof* We consider the effect of transferring  $\delta > 0$  of the bribe of  $x$  to the bribe of  $y$ , where  $\delta$  is as small as we like and where  $eval(x) < eval^\sigma(x) - \delta < eval^\sigma(x) + \delta < 1$  and  $eval(y) < eval^\sigma(y) - \delta < eval^\sigma(y) + \delta < 1$ . By Proposition 8, we have that the utility change is as follows.

$$\delta \left( \sum_{z \in N(y)} \frac{1}{v_z} \right) - \delta \left( \sum_{z \in N(x)} \frac{1}{v_z} \right) = \delta \left( \sum_{z \in N(y)} \frac{1}{v_z} - \sum_{z \in N(x)} \frac{1}{v_z} \right)$$

Since we know that  $\sigma$  is an optimal bribing strategy, this change cannot yield a positive change in network utility. Thus, we can bound this quantity as follows.

$$\delta \left( \sum_{z \in N(y)} \frac{1}{v_z} - \sum_{z \in N(x)} \frac{1}{v_z} \right) \leq 0$$

Now, we consider the effect of transferring  $\delta$  amount of the bribe of  $y$  to the bribe of  $x$ . The change in network utility is given by the following.

$$\delta \left( \sum_{z \in N(x)} \frac{1}{v_z} \right) - \delta \left( \sum_{z \in N(y)} \frac{1}{v_z} \right) = \delta \left( \sum_{z \in N(x)} \frac{1}{v_z} - \sum_{z \in N(y)} \frac{1}{v_z} \right)$$

With similar reasoning to the above, we have that

$$\delta \left( \sum_{z \in N(x)} \frac{1}{v_z} - \sum_{z \in N(y)} \frac{1}{v_z} \right) \leq 0.$$

Combining these two inequalities yields

$$\sum_{z \in N(x)} \frac{1}{v_z} \leq \sum_{z \in N(y)} \frac{1}{v_z}, \text{ and } \sum_{z \in N(y)} \frac{1}{v_z} \leq \sum_{z \in N(x)} \frac{1}{v_z},$$

which finally imply that

$$\sum_{z \in N(y)} \frac{1}{v_z} = \sum_{z \in N(x)} \frac{1}{v_z}.$$

In words, Lemma 2 shows that given an optimal bribing strategy, we can move bribes amongst non-fully bribed voters arbitrarily without affecting the revenue acquired so long as we do not totally remove all the bribe from a customer that was not originally a non-voter, and we do not turn a non-voter into a voting one.

**Lemma 3** *Let  $(C, E)$  be some network with initial evaluation  $eval_0$  and let  $\sigma$  be a bribing strategy. Let  $c \in C$  be such that  $eval_0(c) \neq *$  and  $eval^\sigma(c) = \delta > 0$ , but where for every customer  $c'' \in \bigcup \{N(c') : c' \in N(c)\}$ , we have that  $\delta < eval^\sigma(c'')$ . If  $\sigma_{-c}$  is the bribing strategy obtained from  $\sigma$  by removing the bribe from  $c$ , we have that  $\mathbf{r}(\sigma_{-c}) \geq \mathbf{r}(\sigma)$ .*

*Proof* By definition we have that

$$u_P^0 = \sum_{z \in C} \text{P-RATING}(z, eval).$$

For all  $z \in N^2(x)$ , we know that  $eval(z) > \delta$ . In order to examine the change in utility after setting  $eval(x) = *$ , we should consider the P-RATINGS of the neighbours of  $x$ . Let  $y$  be a neighbour of  $x$ . The average P-RATING of its voting neighbours is clearly higher when  $eval(x) = *$ , since the previous evaluation of  $x$  ( $\delta$ ) is lower than the evaluation of all of its other neighbours. Therefore, all P-RATINGS of all neighbours of  $x$  *increase*; that is, the utility also increases.

We are now ready to detail the second part of the proof ( $\Leftarrow$ ) of the following statement:

**Theorem** BRIBE-NVKL is NP-complete.

*Proof* ( $\Leftarrow$ ) Suppose that  $((C, E), eval_0, \rho)$  is a yes-instance of NVKL and that  $\sigma$  is a bribing strategy that yields a revenue of at least  $\rho$ . We will assume that  $\sigma$  is also optimal, i.e., that there is no other strategy  $\sigma'$  yielding a higher revenue. We make this assumption so that we are able to apply Lemma 2, which gives an important property of optimal bribing strategies that we will repeatedly exploit throughout the remainder of this reduction.

We will now show that  $\sigma$  can be transformed into a revenue-equivalent strategy such that (a) only old customers are bribed, (b) all bribed old customers are bribed fully, and (c) exactly  $k$  old customers are bribed. Recall the terminology from the reduction detailed in the proof of the right-to-left direction of the theorem.

*Revenue equivalent strategy - new customers* Let us call *new customers*, the set of edge and pendant customers. We begin by showing that  $\sigma$  can be modified into an optimal strategy that does not bribe any new customer.

Recall that the initial evaluation of any new customer is  $\epsilon > 0$ . If a new customer is bribed then the bribes to new customers can be enumerated in descending order as  $1 - \epsilon, 1 - \epsilon, \dots, 1 - \epsilon, \epsilon_1, \epsilon_2, \dots, \epsilon_s$ , for some  $s \geq 0$  and where  $0 < \epsilon_i < 1 - \epsilon$  for each  $i = 1, 2, \dots, s$ . Note that it is possible there is no bribe of  $1 - \epsilon$ ; that is, that no new customer is fully bribed. We may assume that  $s \leq 1$ , since if  $s \geq 2$  then by Lemma 2, we could reduce the bribes  $\epsilon_2, \epsilon_3, \dots, \epsilon_s$  without making any equal to zero, so as to increase the bribe  $\epsilon_1$  to  $1 - \epsilon$  and secure another fully bribed new customer.

We begin by proving that *there exists an old customer  $c'$  who has not been bribed and where at most one of its adjacent new pendant customers has been bribed*. Our initial choice of  $\epsilon$  was such that

$$\epsilon = \frac{k}{n + n^2 + \frac{3n}{2}}.$$

Furthermore, the amount invested  $k$ , is certainly less than  $n$ . Therefore we have that

$$\epsilon \leq \frac{n}{n + n^2 + \frac{3n}{2}} = \frac{2}{2n + 5}, \text{ and } 1 - \epsilon > \frac{2n + 3}{2n + 5}.$$

Let  $p$  be the number of old customers that have been fully bribed and let  $t$  be the number of new customers that have been bribed. By summing the total amount bribed, we have that  $p + 1 + t(1 - \epsilon) < k$ . Thus

$$t < \frac{k - p - 1}{1 - \epsilon} < \frac{(2n + 5)(n - p - 1)}{2n + 3},$$

and therefore  $t < 2(n - (p + 1))$ . This says that the number of bribed new customers is strictly less than twice the number of old customers that remain



unbribed. If every unbribed old customer was adjacent to two or more bribed new pendant customers then we would have that  $t \geq 2(n - (p + 1))$  which is clearly not possible. Therefore we can conclude that there exists an unbribed old customer  $c'$  such that at most one of its adjacent new pendant vertices has been bribed.

We now consider two cases. **Case 1a:** we suppose that there exists a fully bribed new customer, and derive a contradiction with the optimality of  $\sigma$ . Consider the bribe of  $1 - \epsilon$  to  $c$ , and consider the increase in P-RATING generated by this single bribe. If  $c$  is a new pendant customer then this contribution is certainly less than 2 as  $|N(c)| = 2$ , and if  $c$  is a new edge customer then this contribution is less than 3 as  $|N(c)| = 3$ . Therefore, in all cases, the bribe of  $1 - \epsilon$  to  $c$  contributes less than 3 units to the overall utility accrued from  $\sigma$ .

Let  $c'$  be an old customer that is not bribed and that is adjacent to at most one new pendant customer that has been bribed (such customer exists, as shown above). Consider moving the  $1 - \epsilon$  bribe from  $c$  to  $c'$ ; so, we obtain a new (efficient) strategy  $\sigma'$ . Let us examine the increase in P-RATING generated by this new  $1 - \epsilon$  bribe.

At least  $n - 1$  of the new pendant customers adjacent to  $c'$  have not been bribed and so the associated cumulative increase in rating is given by  $(n - 1)\frac{1}{2} - (n - 1)\epsilon$ . Given that  $\epsilon \leq \frac{2}{2n+5}$  then the cumulative increase in utility is

$$(n - 1) \left( \frac{1}{2} - \epsilon \right) > \frac{n - 1}{2} - 1.$$

Bribing  $c'$  might reduce the  $P$ -ratings of  $c'$  and its adjacent new edge customers. However, this reduction is certainly less than 4 units. Therefore we may conclude that the movement of  $1 - \epsilon$  of bribe from  $c$  to  $c'$  increases the overall utility by an amount greater than  $(\frac{n-1}{2} - 1) - 7$  units. This amount is strictly positive for  $n$  sufficiently large ( $n \geq 14$ ). Therefore the strategy  $\sigma'$  that we have constructed yields a revenue greater than that of  $\sigma$ , in contradiction with its optimality.

**Case 2a:** we suppose that some new customer  $c$  has been bribed some amount  $\delta$  such that  $0 < \delta < 1 - \epsilon$ , and by a detailed case study (omitted) we derive a contradiction with the optimality of  $\sigma$ .

We can now conclude that *no new customer have been bribed* in the revenue-equivalent optimal strategy  $\sigma$ .

*Revenue equivalent strategy - old customers* We now turn our attention to old customers. The bribes on old customers can be enumerated in descending order as  $1, 1, \dots, 1, \delta_1, \delta_2, \dots, \delta_m$ , for some  $m \geq 0$  and where  $0 < \delta_i < 1$ , for each  $i = 1, 2, \dots, m$ , with possibly no bribes of 1. Without loss of generality, we may assume that  $\sum_{i=1}^m \delta_i \leq 1$ ; otherwise, we would have that  $m \geq 2$  and we could reduce the bribes  $\delta_2, \delta_3, \dots, \delta_m$ , without making any equal to zero, so as to increase the bribe  $\delta_1$  to 1 and secure another fully bribed customer.

Observe that by Lemma 3, we can assume that at most one old customer has not been fully bribed. We now reason by cases.

**Case 1b:** suppose now that there is in fact *one bribed old customer that has not been fully bribed*. Let us call this old customer  $c$  and further suppose that it has been bribed  $\delta$  where  $0 < \delta < 1$ . We will again show that this yields yet another contradiction with the optimality of  $\sigma$ . We have the capacity to increase this bribe to 1 at a cost of  $1 - \delta$  (which we can do, given the remaining resource). The P-RATING of all the customers within  $N(c)$  will increase with the cumulative increase (only due to new pendant neighbours) being

$$n \frac{1 + \epsilon}{2} - n \frac{\delta + \epsilon}{2} = n \frac{1 - \delta}{2}.$$

Hence we obtain an increase in revenue for  $n$  sufficiently large ( $n \geq 3$ ). This contradicts the optimality of  $\sigma$ . Henceforth, we assume that, without loss of generality, any optimal bribing strategy  $\sigma$  on  $(C, E)$ , with initial evaluation  $eval_0$ , is necessarily such that only old customers are bribed and bribed old customers are fully bribed.

**Case 2b:** suppose now that the bribing strategy  $\sigma$  *bribes less than  $k$  old customers*; so, there is an old customer  $c$  that has not been bribed. Let us amend  $\sigma$  to obtain a new bribing strategy  $\sigma'$  by bribing  $c$  so that  $\sigma'(c) = 1$ . This costs us 1 unit of resource. There is no customer of  $C$  such that its P-RATING decreases, and the cumulative increase in P-RATING of the  $n$  new pendant customers adjacent to  $c$  is

$$n \left( \frac{1 + \epsilon}{2} - \epsilon \right) = n \left( \frac{1 - \epsilon}{2} \right) > \frac{n(2n + 3)}{2(2n + 5)} > \frac{n}{4}$$

which is strictly greater than 1 (the amount invested) for  $n$  sufficiently large ( $n \geq 5$ ). This contradicts the optimality of  $\sigma$ . Furthermore, it is clear that more than  $k$  old customers could not have been bribed since the initial utility of the network totals only  $k$  and each old customer is bribed by 1.

*Finding an independent set of size  $k$*  We have shown above that the optimal bribing strategy  $\sigma$  on  $(C, E)$  is such that only old customers are bribed, all bribed old customers are fully bribed, and exactly  $k$  old customers are bribed. Consider now the revenue accruing from our optimal bribing strategy  $\sigma$ . Irrespective of which  $k$  old customers are fully-bribed, the increase in P-RATING due to these old customers is equal to:

$$\frac{(1 + (n + 3)\epsilon)}{n + 4} - \epsilon = \frac{1 - \epsilon}{n + 4},$$

and the P-RATING of the pendant customers adjacent to each of these bribed old customers increases by:

$$\frac{1 + \epsilon}{2} - \epsilon = \frac{1 - \epsilon}{2}.$$

All that remains is to compute the revenue accruing due to the new edge customers adjacent to each of these bribed old customers (as the P-RATING of

any other old or new customer does not change). However, this depends upon how many bribed old customers each new edge customer is adjacent to. Let  $m_i$  denote the number of new edge customers adjacent to  $i$  bribed old customers, for  $i = 1, 2$ . If a new edge customer  $c$  is adjacent to 1 bribed old customer then its increase in P-RATING is

$$\frac{(1 + \epsilon)}{2} - \epsilon = \frac{(1 - \epsilon)}{2}$$

and if it is adjacent to two bribed old customers then its increase in P-RATING is

$$\frac{(2 + \epsilon)}{3} - \epsilon = \frac{2(1 - \epsilon)}{3}.$$

So, the total increase in revenue is

$$m_1 \frac{(1 - \epsilon)}{2} + m_2 \frac{2(1 - \epsilon)}{3}.$$

We also know that by counting the edges joining bribed old customers and their adjacent new edge customers, we obtain that  $3k = 2m_2 + m_1$ . Hence, the total increase in P-RATING due to new edge customers is:

$$m_1 \frac{(1 - \epsilon)}{2} + m_2 \frac{2(1 - \epsilon)}{3} = (3k - 2m_2) \frac{(1 - \epsilon)}{2} + m_2 \frac{2(1 - \epsilon)}{3} = \frac{3k(1 - \epsilon)}{2} - m_2 \frac{(1 - \epsilon)}{3}.$$

So, the revenue due to the bribing strategy  $\sigma$  is:

$$\begin{aligned} & \frac{k(1 - \epsilon)}{n + 4} + \frac{nk(1 - \epsilon)}{2} + \frac{3k(1 - \epsilon)}{2} - m_2 \frac{(1 - \epsilon)}{3} - k \\ &= (1 - \epsilon) \left[ \frac{k}{n + 4} + \frac{k(n + 3)}{2} - \frac{m_2}{3} \right] - k. \end{aligned}$$

Clearly this revenue is largest when  $m_2$  is 0, and if  $m_2 > 0$  then the revenue is less than this maximal value. Also, when  $m_2$  is 0 this revenue is exactly equal to  $\rho$ . Hence, as we started with a yes-instance of NVKL, we must have that  $m_2 = 0$ , i.e., that no edge customer is adjacent to two bribed old customers. Thus, the  $k$  vertices of  $G$  corresponding to the  $k$  bribed old customers in  $C$  form an independent set, and  $(G, k)$  is a yes-instance of ISREG(3).