

# Extensive Experimental Validation of a Personalized Approach for Coping with Unfair Ratings in Reputation Systems

Jie Zhang

Nanyang Technological University, School of Computer Engineering, zhangj@ntu.edu.sg

Received 11 September 2010; received in revised form 25 May 2011; accepted 1 July 2011

## Abstract

The unfair rating problem exists when a buying agent models the trustworthiness of selling agents by also relying on ratings of the sellers from other buyers in electronic marketplaces, that is in a reputation system. In this article, we first analyze the capabilities of existing approaches for coping with unfair ratings in different challenging scenarios, including ones where the majority of buyers are dishonest, buyers lack personal experience with sellers, sellers may vary their behavior, and buyers may provide a large number of ratings. We then present a personalized modeling approach (PMA) that has all these capabilities. Our approach allows a buyer to model both the private reputation and public reputation of other buyers to determine whether these buyers' ratings are fair. More importantly, in this work, we focus on experimental comparison of our approach with two key models in a simulated dynamic e-marketplace environment. We specifically examine the above mentioned scenarios to confirm our analysis and to demonstrate the capabilities of our approach. Our study thus provides the extensive experimental support for the personalized approach that can be effectively employed by reputation systems to cope with unfair ratings.

**Keywords:** Trust and reputation systems, Unfair ratings, Electronic marketplaces, Probabilistic reasoning approaches, Multi-agent systems

## 1 Introduction

In electronic marketplaces populated by self-interested agents, buying agents would benefit by modeling the trustworthiness of selling agents, in order to make effective decisions about which agents to trust. How to effectively represent the trustworthiness of sellers then becomes a challenge that must be addressed in order to ensure that users feel secure when engaging in commerce online. One method for representing sellers' trustworthiness is to ask other buying agents in the system (called advisors) to provide ratings of sellers, that is a reputation system [1]. For example, the REGRET model of Sabater et al. [26] offers a multi-dimensional view of trust that includes a social dimension, where the ratings of a selling agent provided by other members (advisors) in a buying agent's group are also considered for evaluating the trustworthiness of the selling agent. However, the problem of unfair ratings may then arise. Advisors may provide unfairly high ratings to promote sellers. This is referred to as "ballot stuffing" [2]. Advisors may also provide unfairly low ratings, in order to cooperate with other sellers to drive a seller out of the marketplace. This is referred to as "bad-mouthing".

To cope with the problem of unfair ratings, researchers have been developing proactive approaches that create incentives for buyers to provide fair ratings [4], [12], [41]. These approaches have to be deployed in the marketplaces since the very beginning. They cannot deal with the already existing unfair ratings. Note that the approach proposed in [41] can in fact deal with the existing unfair rating because it employs the reactive approaches as its basis for creating incentives. However, the discussion of this approach is not the focus of this paper, and the proactive approaches in general cannot deal with the already existing unfair ratings. In this paper, we focus on another set of approaches that are reactive. These approaches can be used to filter out the existing unfair ratings. Some reactive approaches [2], [6], [18], [19] apply the clustering (grouping) of ratings or buyers to deal with the unfair rating problem. For example, the approach of Dellarocas [2] uses a divisive clustering algorithm to separate ratings for a seller into two clusters: the one containing lower ratings, and the one containing higher ratings. The ratings in the higher ratings cluster are considered as unfairly high ratings, and therefore are discarded. However, this approach cannot effectively handle unfairly low ratings. The approach of Liu et al. [18], [19] uses clustering to find the other buyers that are similar to a buyer based on the buyer's own ratings and other buyers' ratings for their commonly rated sellers. This approach becomes ineffective when the buyer has limited experience with sellers. Despite the rich literature for coping with the problem of unfair ratings, in this paper we focus more on a set of probabilistic reasoning approaches, for the purpose of demonstrating the benefits of our particular probabilistic approach.

A variety of reactive approaches have been proposed to use probabilistic reasoning for addressing the problem of unfair ratings. For example, the beta reputation system (BRS) [34] estimates the trustworthiness of sellers using a probabilistic model based on the beta probability density function. It filters out the ratings that are not in the majority amongst other ones. However, the BRS system is only effective when the significant majority of ratings are fair. TRAVOS, developed by Teacy et al. [31] proposes that possibly unreliable ratings of sellers should be discounted when the buying agent tries to reason about the trustworthiness of the sellers. However, this model does not work well when sellers vary their behavior widely. In this article, we begin by surveying some of these existing probabilistic approaches to the unfair rating problem and analyzing their capabilities.

We present a personalized modeling approach (PMA) for coping with unfair ratings in reputation systems. This approach allows a buyer to model the trustworthiness of advisors by combining the buyer's personal experience with the advisors and the public knowledge about them held by the system. More specifically, our approach allows buyers to first represent private reputation values of advisors, based on what is known about the advisors' ratings for sellers with which the buyers have already had some experience. Next, buyers construct a public model of trustworthiness of advisors based on common, centrally held knowledge of sellers and the ratings provided by advisors, including the ratings of sellers totally unknown to the buyer. We use the terms private and public reputation to reflect the advisor's trustworthiness as assessed by the buyer using an accumulation of either private or public knowledge. Then both private and public models can be combined, in order to obtain a value for the trustworthiness of each advisor. Our approach also offers more flexibility for the buyer to weigh the values of both the private and public reputation of advisors.

We focus on extensive experimental comparison of our approach with two key models, the BRS system [34] and the TRAVOS model [31] in a simulated dynamic e-marketplace environment involving possibly deceptive buying and selling agents. Because BRS, TRAVOS and our approach work for only binary ratings, the environment allows only binary ratings to represent simple and objective results of transactions between sellers and buyers (advisors). Advisors in the environment do not make profit from providing advice or pay cost to generate advice. We specifically examine different scenarios, including ones where the majority of buyers are dishonest, buyers lack personal experience with sellers, sellers may vary their behavior, and buyers may provide a large number of ratings. Experimental results show that our personalized approach in general performs better than the TRAVOS model and the BRS system. Our approach shows its advantages especially when buyers do not have much experience with sellers and sellers vary their behavior widely. Our personalized model can therefore be seen as a valuable approach to use when introducing the sharing of seller ratings among buyers in order to model the trustworthiness of sellers in electronic marketplaces.

The rest of the paper is organized as follows. Section 2 presents a survey of some existing probabilistic approaches for coping with unfair ratings and summarizes their capabilities. Section 3 describes the formalization of our personalized approach. Section 4 provides the framework used for simulating an e-marketplace and for conducting experiments. Section 5 presents comparative results. Finally, Section 6 concludes the paper and proposes future work.

## 2 Related Work

In this section, we provide a summary of some existing probabilistic approaches for coping with the unfair rating problem. Advantages and disadvantages of these approaches are also pointed out. We then provide deep analysis of their capabilities.

### 2.1 Beta Reputation System

The beta reputation system (BRS) proposed by Jøsang and Ismail [11] estimates reputation of selling agents using a probabilistic model. This model is based on the beta probability density function, which can be used to represent probability distributions of binary events. The beta distributions are a family of statistical distribution functions that are characterized by two parameters  $\alpha$  and  $\beta$ . The beta probability density function is defined as follows:

$$beta(p|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1}(1-p)^{\beta-1} \quad (1)$$

where  $\Gamma$  is the gamma function,  $p \in [0,1]$  is a probability variable, and  $\alpha, \beta > 0$ . This function shows the relative likelihood of the values for the parameter  $p$ , given the fixed parameters  $\alpha$  and  $\beta$ .

This model is able to estimate the reputation of a seller by propagating ratings provided by multiple advisors. Ratings are binary in this model (1 or 0, to represent that the advisor considers the seller to be satisfactory or dissatisfactory in a transaction). Individual ratings received are combined by simply accumulating the number of ratings ( $m$ ) supporting the conclusion that the seller has good reputation and the number of ratings ( $n$ ) supporting the conclusion that the seller has bad reputation. To ensure  $\alpha, \beta > 0$ , the values for  $\alpha$  and  $\beta$  are then set as follows:

$$\alpha = m + 1, \quad \beta = n + 1 \quad (2)$$

The prior distribution of the parameter  $p$  is assumed to be the uniform beta probability density function with  $\alpha = 1$  and  $\beta = 1$ . The posteriori distribution of  $p$  is the beta probability density function after observing  $\alpha - 1$  ratings of 1 and  $\beta - 1$  ratings of 0. An example of the beta probability density function when  $m = 7$  and  $n = 1$  is shown in Figure 1. This curve expresses the relative likelihood of the probability  $p$  that the seller will have good reputation in the future. When  $m > n$ , it is more likely that the probability value  $p > 0.5$ . For example, from the curve in Figure 1, we can see that it is more likely that  $p = 0.6$  than that  $p = 0.2$ .

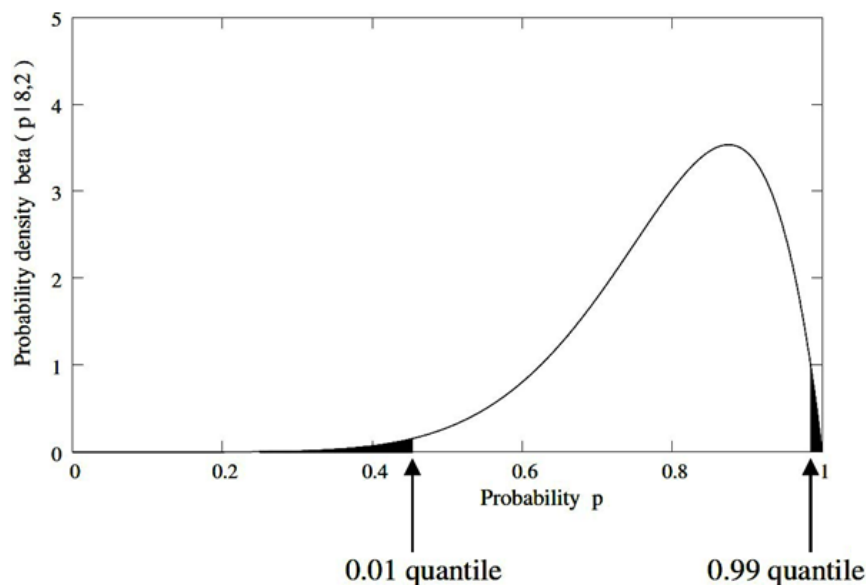


Figure 1: PDF when  $m = 7$  and  $n = 1$  [34]

The reputation of the seller  $s$  can then be represented by the probability expectation value of the beta distribution, which is the most likely frequency value, used to predict whether the seller will act honestly in the future. The formalization of this is given as follows:

$$Tr(s) = E(p) = \frac{\alpha}{\alpha + \beta} \quad (3)$$

According to this calculation, the reputation of the seller  $s$  in Figure 1 is 0.8.

To handle unfair ratings provided by advisors, Whitby et al. [34] extend BRS to filter out those ratings that are not in the majority amongst other ones. More specifically, feedback provided by each advisor is represented by a beta distribution. If the cumulated reputation of the seller falls between the lower and upper boundaries of the feedback, this feedback will be considered as fair. Figure 1 shows a demonstration of this process when the lower and upper boundaries are 0.01 and 0.99 respectively. When the cumulated reputation of the seller is within the black area ( $Tr(s) > 0.98$  or  $Tr(s) < 0.45$  in this case), the advisor's ratings will be considered as unfairly high or unfairly low ratings.

However, this approach is only effective when a significant majority of the ratings are fair. Suppose there are 4 advisors,  $a_1, a_2, a_3$  and  $a_4$ . Each advisor has provided one rating for a dishonest seller. The rating provided by advisor  $a_1$  is 0, which is fair. The other three advisors' ratings are all 1, which is unfair. In this case, the cumulated reputation of the seller is calculated as  $\frac{3+1}{4+2} = 0.67$  (Equation 3). By setting the lower and upper boundaries to 0.01 and 0.99 as suggested by the authors of [34], the cumulated reputation of the seller falls between the lower and upper boundaries of the ratings of advisors  $a_2, a_3$  and  $a_4$ . The unfair ratings of advisors  $a_2, a_3$  and  $a_4$  will then be incorrectly considered as fair ratings.

## 2.2 TRAVOS Model

Teacy et al. [31] propose the TRAVOS model, which is a trust and reputation model for agent-based virtual organizations. This approach is also based on the beta probability density function. It copes with inaccurate reputation advice by accomplishing two tasks. The first task is to estimate the accuracy of the current reputation advice (ratings of 1 or 0) provided by the advisor about the seller, based on the buyer's personal experience with the advisor's previous advice. More specifically, the TRAVOS model divides the interval of [0,1] into  $N_{bin}$  number of equal bins. It then finds out all the previous advice provided by the advisor that is similar to the advice being currently given by the advisor. The two pieces of advice are similar if they are within the same bin. The accuracy of the current advice will be the expected value of the beta probability density function representing the amount of the successful and unsuccessful interactions between the buyer and the seller when the buyer follows the previous advice.

Let us consider an example of estimating the trustworthiness of an advisor. Suppose the interval of [0,1] is divided into two bins, [0,0.5] and [0.5,1]. The current advice provided by the advisor about a seller consists of 7 ratings of 1 and 1 rating of 0. This indicates that the trustworthiness of the seller is 0.8 (using the calculations in the previous section). The current advice is then within the bin of [0.5,1]. Thus, the previous advice of the advisor that is also between 0.5 and 1 will be considered similar to the current advice. Suppose that by following this similar advice, the buyer has had 3 successful interactions and 0 unsuccessful interactions with the seller. The trustworthiness of the advisor is then calculated as  $\frac{3+1}{3+2} = 0.8$ .

The second task is to adjust reputation advice according to its accuracy. The aim of this task is to reduce the effect of inaccurate advice. This task is necessary because it can deal with the situation where an advisor unfairly rates a seller a large number of times. Experimental results show that TRAVOS has better performance in estimating sellers' trustworthiness than the BRS system [31]. However, this model also has some weaknesses. It assumes that selling agents act consistently and thus does not discount old ratings provided by advisors for sellers. This assumption might not be true. A seller may change its behavior from being trustworthy to being untrustworthy. Suppose an advisor has done business with the seller before and their interaction is successful. The fair advice provided by the advisor then indicates that the seller is trustworthy. However, this advice will be incorrectly considered as unfair when a buyer takes this advice and does business with the seller after the seller changes its behavior. The second problem is that this model relies only on the buyer's personal experience with the advisor's advice. This will be problematic when the buyer does not have much experience with selling agents, for example if the buyer is new to the system. In this case, it is difficult for the buyer to determine whether the advisor is trustworthy.

## 2.3 Reinforcement Learning Model

Tran and Cohen [32] have buying agents use reinforcement learning to determine with which selling agents to do business, in order to maximize the buyers' expected profit. They also have selling agents use the same learning method to maximize the sellers' profit by adjusting product prices and altering product quality offered to different buyers. To avoid doing business with possibly dishonest sellers, buyers in the market determine the trustworthiness of the sellers using an incremental updating approach motivated by that proposed in [35], after the true value of

delivered products is evaluated and compared to the buying agent's expected value for the products. This approach updates the trustworthiness of sellers based on their previous trust values after examination of goods. The formulae proposed adhere to the principle that trust is difficult to build up but easy to lose. Selling agents can be classified as untrustworthy if their reputation values fall below a certain threshold and buyers try to select only the selling agents with the highest expected value for the goods from the set of selling agents not yet labeled as untrustworthy. This approach of modeling trustworthiness of sellers relies only on buyers' personal experience with sellers. However, a (new) buyer may not have much personal experience with some sellers.

Regan et al. [24] extend this work of Tran and Cohen to allow buyers to share models of sellers. Advisors are then modeled in order to determine whose advice will be considered, using a similar approach of modeling trustworthiness of sellers based on whether buyers are satisfied with the advisors' advice. This modeling only depends on buyers' personal experience with advisors' advice. Reputation advice from selected (trustworthy) advisors, however, is treated equally in this model when estimating an aggregated trust value for each seller.

## 2.4 Bayesian Network Model

Wang and Vassileva [33] propose a Bayesian network-based trust model in a peer-to-peer file sharing system. In this system, file providers' capabilities are evaluated according to different aspects, including download speed, file quality, and file type. A naive Bayesian network is constructed to represent conditional dependencies between the trustworthiness of file providers and the aspects. Each user holds a naive Bayesian network for each file provider. If a user has no personal experience with a file provider, he may ask other users (advisors) for recommendations. A recommendation provided by an advisor will be considered by the user according to the trust value he has of the advisor. The trust value is updated by a reinforcement learning formula. More specifically, it will be increased/decreased after each comparison between the naive Bayesian networks held by the user and the advisor for the file provider. The Bayesian network-based trust model takes into account preference similarity between users and advisors. However, this approach assumes that the aspects of file providers' capabilities are conditionally independent. This assumption may be unrealistic. For instance, users may prefer high quality video and picture files, but do not care much about the quality of text files.

## 2.5 Weighted Majority Algorithm

Yu and Singh [36] propose to use Dempster-Shafer theory as the basis for computing the trustworthiness of an agent. More specifically, they define belief, disbelief and uncertainty parameters  $b, d, u \in [0,1]$  respectively, for the proposition that the agent is trustworthy. These parameters sum up to 1. An orthogonal sum function is also defined to combine beliefs of any two other agents (advisors) about the trustworthiness of the agent that is currently being evaluated. This function yields the same aggregated value regardless of the order in which the beliefs of multiple advisors are combined.

To handle possibly unfair reporting from advisors, Yu and Singh propose an algorithm that uses a version of the weighted majority algorithm (WMA) [37]. In their algorithm, weights are assigned to the advisors. These weights are initialized to be 1 and can be considered as the trustworthiness of the corresponding advisors. The algorithm predicts the trustworthiness of sellers based on the weighted sum of the ratings provided by those advisors. The weight of an advisor's ratings is determined by the trustworthiness of the advisor.

Yu and Singh propose to tune the weights of the advisors after an unsuccessful prediction so that the weights assigned to the advisors are decreased. They assume that the ratings of dishonest advisors may conflict with the observations of the buyers receiving these ratings. By decreasing the weights of these advisors over time, unfair ratings are filtered. Their approach determines the weights of the advisors based only on the buyers' personal experience with the advisors' ratings. If the buyers do not have much personal experience with the advisors' ratings, the weights of the advisors will not be decreased. These weights remain high and the advisors' ratings will then be heavily considered by the buyers. Another problem is that once the weights of the advisors are decreased, the advisors will not be able to gain trust back from the buyers by providing fair ratings to the buyers.

## 2.6 Capabilities

In this section, we compare the different approaches summarized above. Different from the work of [23] that proposes a framework to compare general features of trust and reputation systems, we specifically analyze the capabilities the different approaches have, by first listing the following four capabilities that an effective approach should have.

- *Majority*: An effective approach should be able to cope with unfair ratings even when the majority of the ratings of a seller is unfair. Some approaches assume that unfair ratings can be recognized by their statistical properties, and therefore may suffer in this situation. For example, the performance of BRS largely decreases when the majority of ratings are unfair, which will be demonstrated in Sections 5.1 and 5.2.1;

- *Flooding*: An approach should also be able to deal with the situation where advisors may provide a large number of ratings within a short period of time. The approach of BRS is affected by this situation and the reason for this will be further explained in Section 5.2.4. The Bayesian network-based model is also affected because one advisor may be able to quickly build up its reputation by providing a large number of fair ratings within a short period. One possible way to cope with this is to consider only a limited number of ratings from each advisor within the same period of time, as suggested by Zacharia et al. [38]. In the WMA approach, fair ratings do not increase advisors' trustworthiness, and therefore WMA is not affected by this situation;
- *Lack of Experience*: An approach should still be effective even when buyers do not have much experience with sellers. The approaches, such as TRAVOS, Bayesian, and WMA, suffer from this type of situation. BRS is able to deal with this situation because it can rely on the all the ratings provided for sellers;
- *Varying*: An approach should be able to deal with changes of selling agents' behavior. Because of changes of selling agents' behavior, buying agents may provide different ratings for the same seller. Even though two ratings provided within different periods of time are different, it does not necessarily mean that one of them must be unfair. TRAVOS assumes that selling agents act consistently and it suffers from this problem. Different ways are proposed to deal with this situation. BRS [34] uses a forgetting factor  $\lambda$  ( $0 \leq \lambda \leq 1$ ) to dampen ratings according to the time when they are provided. Older ratings are dampened more heavily than more recent ones.

Table 1: Capabilities of approaches

Approaches	Majority	Flooding	Lack of Experience	Varying
BRS			√	√
TRAVOS	v	√		
Bayesian Network	√			
WMA	√	√		
Reinforcement Learning	√			

Table 1 lists capabilities of the approaches summarized in the previous sections. In this table, the mark “v” indicates that an approach has the capability. For example, the BRS approach is capable of dealing with changes of sellers' behavior and is still effective when buyers do not have much experience. As will be discussed in the next sections, our personalized model has all these capabilities. These capabilities of our approach will be further demonstrated through experiments.

### 3 Our Personalized Approach

In this section, we describe our personalized approach for modeling the trustworthiness of advisors. The early version of this approach was introduced in [39], [40]. The approach is used as part of a centralized reputation system. We assume that all buyers can play the role of advisors to other buyers. We assume as well that advisors provide ratings only when a transaction occurs and these are stored with the central server. This may be kept in check by the centralized system where all buyers agree to have their interactions with sellers known, for instance. We also assume a marketplace where sellers are offering similar kinds of goods.

Our personalized approach allows a buyer to estimate the reputation (referred to as private reputation) of an advisor based on their ratings for commonly rated sellers. We call this type of reputation private reputation because it is based on the buyer's own experience with the advisor's advice, and is not shared with the public. The private reputation value of the advisor may vary for different buyers. When the buyer has limited private knowledge of the advisor, the public reputation of the advisor will also be considered. We call this type of reputation public reputation because it is based on the public's opinions about the advisor's advice, and it is shared by all of the public. The public reputation value of the advisor is the same for every buyer; it is estimated based on all ratings for the sellers ever rated by the advisor. Finally, the trustworthiness of the advisor will be modeled by combining the weighted private and public reputations. These weights are determined based on the estimated reliability of the private reputation.

#### 3.1 Private Reputation of Advisor

Our personalized approach allows a buying agent  $b$  to evaluate the private reputation of an advisor  $a$  by comparing their ratings for commonly rated sellers  $\{s_1, s_2, \dots, s_m\}$ . For one of the commonly rated sellers  $s_i$  ( $1 \leq i \leq m$ ), advisor  $a$  has the rating vector  $R_{a, s_i}$ , and buyer  $b$  has the rating vector  $R_{b, s_i}$ . A rating for  $s_i$  from  $b$  and  $a$  is binary, in which 1 means that  $s_i$  is trustworthy and 0 means that  $s_i$  is untrustworthy. In the current work, we assume that ratings from advisors are objective. Dealing with subjective ratings is left for future work [22]. For the remainder of this paper, we

assume ratings for sellers are binary. Possible ways of extending our approach to accept ratings in different ranges are left for future work.

The ratings in  $R_{a,s_i}$  and  $R_{b,s_i}$  are ordered according to the time when they are provided. The ratings are then partitioned into different elemental time windows. The length of an elemental time window may be fixed (e.g. three days) or adapted by the frequency of the ratings to the seller  $s_i$ , similar to the way proposed in [2]. A window should also be sufficiently small so that there is no need to worry about the changes of sellers' behavior within the time window. We define a pair of ratings  $(r_{a,s_i}, r_{b,s_i})$ , such that  $r_{a,s_i}$  is one of the ratings of  $R_{a,s_i}$ ,  $r_{b,s_i}$  is one of the ratings of  $R_{b,s_i}$ , and  $r_{a,s_i}$  corresponds to  $r_{b,s_i}$ . The two ratings,  $r_{a,s_i}$  and  $r_{b,s_i}$ , are correspondent only if the rating  $r_{b,s_i}$  is the most recent rating in its time window, and the rating  $r_{a,s_i}$  is the closest and prior to the rating  $r_{b,s_i}$ . We consider ratings provided by buyer  $b$  after those by advisor  $a$ , in order to incorporate into buyer  $b$ 's ratings anything learned from advisor  $a$ , before taking an action. According to the solution proposed by Zacharia et al. [38], by keeping only the most recent ratings, we can avoid the issue of advisors "flooding" the system. No matter how many ratings are provided by one advisor in a time window, we only keep the most recent one.

We define the rating pair  $(r_{a,s_i}, r_{b,s_i})$  as a positive rating pair if  $r_{a,s_i}$  is the same value as  $r_{b,s_i}$ . Otherwise, the pair is a negative rating pair. We assume that  $r_{b,s_i}$  is provided within the time window  $T_b$  and  $r_{a,s_i}$  is within the time window  $T_a$ . We also assume that each time window is identified by an integer value, where 1 is the most recent time window with a rating, 2 is the time window just prior, and so on until the oldest time window. So,  $T_a$  is always greater than or equal to  $T_b$  because  $r_{a,s_i}$  is prior to the rating  $r_{b,s_i}$ . As also pointed out by Jøsang and Ismail [11], old ratings may not always be relevant for sellers' actual trustworthiness because sellers may change their behavior over time. Older ratings should be given less weight than more recent ones. In our case, if  $r_{a,s_i}$  and  $r_{b,s_i}$  are within the same time window, it is more relevant to compare them and the rating pair will be given more weight; otherwise, the rating pair will be given less weight.

We then examine rating pairs for  $s_i$ . We define  $N_{s_i}$  as the sum of the weights of all rating pairs for  $s_i$ . The sum of weights  $N_{all}$  of all rating pairs for sellers rated by both the buyer and the advisor will then be calculated as follows:

$$N_{all} = \sum_{i=1}^m N_{s_i} \quad (4)$$

We also define  $N_p$  as the sum of the weights of all positive rating pairs for all commonly rated sellers.

If the two ratings in a rating pair are within the same time window, the weight of the rating pair is 1. In a simple case where each of all rating pairs has two ratings that are within the same time window, we only need to count the number of rating pairs for  $s_i$  to calculate  $N_{s_i}$  and the total number of rating pairs for all commonly rated sellers to calculate  $N_{all}$ .  $N_p$  is the number of all positive rating ratings for all commonly rated sellers in this case.

For the more general case where a rating pair  $(r_{a,s_i}, r_{b,s_i})$  may have two ratings that are within different time windows, we calculate the weight of the rating pair, as follows:

$$z = \lambda^{T_a - T_b} \quad (5)$$

where  $\lambda$  is a forgetting factor (a concept used by BRS [8]) and  $0 \leq \lambda \leq 1$ . Note that when  $\lambda = 1$  there is no forgetting (i.e. older ratings supplied by advisors will be accepted and compared to the buyer's rating in the closest time window). Note as well that when  $\lambda > 0$ , the higher the value of  $\lambda$ , the greater the weight placed on the ratings provided by the advisor. When  $\lambda = 0$ , we are in the simple case described above; ratings that are not in the same window will not be considered.

The private reputation of the advisor  $a$  is estimated as the probability that advisor  $a$  will provide fair ratings to the buyer  $b$ . Because there is only incomplete information about the advisor, the best way of estimating the probability is to use the expected value of the probability. The expected value of a continuous random variable is dependent on a probability density function, which is used to model the probability that a variable will have a certain value. Because of its flexibility and the fact that it is the conjugate prior for distributions of binary events [25], the beta family of probability density functions is commonly used to represent probability distributions of binary events (see, e.g. the generalized trust models BRS [11] and TRAVOS [31]). Therefore, the private reputation of advisor  $a$  can be calculated as follows:

$$\alpha = N_p + 1, \quad \beta = N_{all} - N_p + 1$$

$$R_{pri}(a) = E(Pr(a)) = \frac{\alpha}{\alpha + \beta} \quad (6)$$

where  $Pr(a)$  is the probability that advisor  $a$  will provide fair ratings to buyer  $b$ , and  $E(Pr(a))$  is the expected value of the probability, which is the most likely probability value that the advisor will be honest in the future. An advisor's

rating is considered to be a fair rating if it is the same as the buyer's rating. As explained, the advisor's rating is examined either in the same or the closest time window, and is submitted prior to the buyer's experience. The buyer's experience is used to judge the fairness of the rating. The buyer may decide not to trust the advisor if they have a different view of sellers.

### 3.2 Public Reputation of Advisor

When there are not enough rating pairs, the buyer  $b$  will also consider advisor  $a$ 's public reputation. This is determined by Equations 8 and 9 for calculating the weight of private reputation, which will be explained later in this section. When the weight is less than 1, there are not enough rating pairs and public reputation will also be considered. The public reputation of advisor  $a$  is estimated based on her ratings and other ratings for the sellers rated by advisor  $a$ . Each time advisor  $a$  provides a rating  $r_{a,s}$  for any seller  $s$ , the rating will be judged centrally as a consistent or inconsistent rating. We define a rating for a seller as a consistent rating if it is consistent with the majority of the ratings of the seller up to the moment when the rating is provided. Determining consistency with the majority of ratings can be achieved in a variety of ways, for instance averaging all the ratings and seeing if that is close to the advisor's rating, which is the method used in our experiments in Section 5. The development of more comprehensive methods is left for future work. We consider only the ratings within a time window prior to the moment when the rating  $r_{a,s}$  is provided, and we only consider the most recent rating from each advisor. In so doing, as sellers change their behavior and become more or less trustworthy to each advisor, the majority of ratings will be able to change.

Suppose that the advisor  $a$  provides  $N'_{all}$  ratings in total. If there are  $N_c$  consistent ratings, the number of inconsistent ratings provided by advisor  $a$  will be  $N'_{all} - N_c$ . In a similar way as estimating the private reputation, the public reputation of the advisor  $a$  is estimated as the probability that advisor  $a$  will provide consistent ratings. It can be calculated as follows:

$$\alpha' = N_c + 1, \quad \beta' = N'_{all} - N_c + 1$$

$$R_{pub}(A) = \frac{\alpha'}{\alpha' + \beta'} \quad (7)$$

which also indicates that the greater the percentage of consistent ratings advisor  $a$  provides, the more reputable she will be considered.

### 3.3 Modeling Trustworthiness of Advisor

To estimate the trustworthiness of advisor  $a$ , we combine the private reputation and public reputation values together. The private reputation and public reputation values are assigned different weights. The weights are determined by the reliability of the estimated private reputation value.

We first determine the minimum number of rating pairs needed for buyer  $b$  to be confident about the private reputation value he has of advisor  $a$ . The Chernoff Bound theorem [21] provides a bound for the probability that the estimation error of private reputation exceeds a threshold, given the number of rating pairs. Accordingly, the minimum number of pairs can be determined by an acceptable level of error and a confidence measurement as follows:

$$N_{min} = -\frac{1}{2\varepsilon^2} \ln \frac{1-r}{2} \quad (8)$$

Where  $\varepsilon \in (0,1)$  is the maximal level of error that will be accepted by  $b$  and  $\gamma \in (0,1)$  is the level of confidence buyer  $b$  would like to attain. If the total weight of all rating pairs  $N_{all}$  is larger than or equal to  $N_{min}$ , buyer  $b$  will be confident about the private reputation value estimated based on his ratings and the advisor  $a$ 's ratings for all commonly rated sellers. Otherwise, there are not enough rating pairs, the buyer will not be confident about the private reputation value, and it will then also consider public reputation. The reliability of the private reputation value can be measured as follows:

$$w = \begin{cases} \frac{N_{all}}{N_{min}} & \text{if } N_{all} < N_{min}; \\ 1 & \text{otherwise} \end{cases} \quad (9)$$

The trust value of advisor  $a$  will be calculated by combining the weighted private reputation and public reputation values as follows:

$$Tr(a) = wR_{pri}(a) + (1-w)R_{pub}(a) \quad (10)$$



The buyer will consider the public reputation value less when the private reputation value is more reliable. Note that when  $w = 1$ , the buyer relies only on private reputation. This can be used as well if the majority rating is suspect. The buyer can rely on its own private knowledge and allow for a difference of opinion. Once a buyer has had personal experience, it will know better whether the majority opinion is acceptable. For future work, a learning approach may need to be developed to optimally learn the weight. This approach will adjust the weight of different parts of the personalized approach for each buyer, by learning the two parameters  $\varepsilon$  and  $\gamma$  (see Equations 8 and 9). Once the optimal values for these two parameters are learned based on the amount of the private knowledge the buyer has about advisors and the buyer's estimation whether the majority of advisors is lying, the weight can then be effectively determined.

Algorithm 1 (Figure 2) is a pseudo code summary of the personalized approach for modeling the trustworthiness of an advisor.

---

```

//Buyer estimates private reputation of advisor
{ $s_1, s_2, \dots, s_m$ }: sellers commonly rated by buyer  $b$  and advisor  $a$ ;
Set  $N_{all} = 0$ : sum of weights of all rating pairs for  $b$  and  $a$ ;
Set  $N_p = 0$ : sum of weights of all positive rating pairs for  $b$  and  $a$ ;
foreach  $s_i$  in  $\{s_1, s_2, \dots, s_m\}$  do
    //comparing ratings for commonly rated sellers
     $R_{b, s_i}$ : buyer  $b$ 's ratings for seller  $s_i$ ;
     $R_{a, s_i}$ : buyer  $a$ 's ratings for seller  $s_i$ ;
    foreach rating  $r_{b, s_i}$  in  $R_{b, s_i}$  do
        if a rating  $r_{a, s_i}$  of advisor  $a$  in  $R_{a, s_i}$  corresponds to  $r_{b, s_i}$  then
            //checking time windows
             $N_{all} = N_{all} + z$ ; //z is calculated using Equation 5
            if  $r_{a, s_i} = r_{b, s_i}$  then
                 $N_p = N_p + z$ ;
    _
_

Private reputation is then calculated using Equation 6;
Calculate weight  $w$  using Equations 8 and 9;
Set public reputation = 0;

if weight  $w < 1$  then
    // private knowledge is limited, buyer also estimates public reputation
    Set  $N'_{all} = 0$ : number of all ratings provided by advisor  $a$ ;
    Set  $N_c = 0$ : number of ratings by advisor  $a$  consistent with majority;;
    { $s_1, s_2, \dots, s_m$ }: sellers ever rated by advisor  $a$ ;
    foreach  $s_j$  in  $\{s_1, s_2, \dots, s_n\}$  do
         $R_{a, s_j}$ : advisor  $a$ 's ratings for seller  $s_j$ ;
        foreach rating  $r_{a, s_j}$  in  $R_{a, s_j}$  do
            //Comparing  $r_{a, s_j}$  with other ratings of seller  $s_j$ 
            if  $r_{a, s_i}$  is consistent then
                 $N_c = N_c + 1$ ;
        _
    _

Public reputation is then calculated using Equation 7;

Trustworthiness = weighted combination of private and public reputation
    
```

---

Figure 2: Algorithm 1 - buyer  $b$  modeling trustworthiness of an advisor  $a$

## 4 Experimental Framework

In this section, we introduce a framework for conducting experiments to compare our personalized approach with the other approaches for handling unfair ratings. The marketplace environment used for experiments is populated with self-interested buying and selling agents. The buyers and sellers are brought together by a procurement (reverse) auction, where the auctioneer is a buyer and bidders are sellers. There is a central server that runs the auction. In the marketplace, a buyer  $b$  that wants to purchase a product  $p$  sends a request to the central server. Sellers interested in selling the product to the buyer will register to participate in the auction. The buyer will first limit the sellers it will consider for the auction, by modeling their trustworthiness. To directly compare the performance of the approaches for coping with unfair ratings, we use an algorithm for the buyer to model the trustworthiness of the sellers, only making use of ratings from advisors. Note that in this framework, we keep the roles of trust and auctions separated for the purpose of simplicity. More specifically, trust is used to limit the sellers that buyers will consider for auctions, and auctions are used to determine the price of sellers' products. A deeper study on how to effectively use

trust to enhance auction-based seller selection is left for future work. Detailed discussion on this direction can be found in Section 6.

#### 4.1 Modeling the Trustworthiness of Seller

Assume that a buyer  $b$  considers ratings provided by advisors that are trustworthy. It sends a request to the central server to ask for all the ratings provided by the trustworthy advisors  $\{a_1, a_2, \dots, a_m\}$  ( $m \geq 1$ ) for the seller  $s$ . Suppose that the advisor  $a_i$  ( $1 \leq i \leq m$ ) provided  $N_{pos}^{a_i}$  positive ratings and  $N_{neg}^{a_i}$  negative ratings. In this work, we consider only ratings that are binary because the approaches we compare are all using the beta density function for representing binary ratings. These ratings will be discounted based on the trustworthiness of the advisor, so that the ratings from less trustworthy advisors will carry less weight than ratings from more trustworthy ones.

Jøsang [9] provides a mapping from beliefs defined by the Dempster-Shafer theory to the beta function as follows:

$$\begin{cases} f = \frac{N_{pos}^{a_i}}{N_{pos}^{a_i} + N_{neg}^{a_i} + 2} \\ d = \frac{N_{neg}^{a_i}}{N_{pos}^{a_i} + N_{neg}^{a_i} + 2} \\ u = \frac{2}{N_{pos}^{a_i} + N_{neg}^{a_i} + 2} \end{cases} \quad (11)$$

where  $f$ ,  $d$  and  $u$  represent belief, disbelief and uncertainty parameters, respectively. For our setting of trust modeling,  $f$  represents the probability that the proposition that the seller is trustworthy is true, and  $d$  represents the probability that the proposition is false. Note that  $f + d + u = 1$  and  $f, d, u \in [0,1]$ . As also pointed out in [9] and [37], beliefs and disbeliefs can be directly discounted by the trustworthiness of the advisor  $a_i$  as follows:

$$\begin{cases} f' = Tr(a_i)f \\ d' = Tr(a_i)d \end{cases} \quad (12)$$

where  $Tr(a_i)$  is the trustworthiness of  $a_i$ . From Equations [11] and [12], we then can derive a discounting function for the amount of ratings provided by  $a_i$  as follows:

$$\begin{cases} D_{pos}^{a_i} = \frac{2Tr(a_i)N_{pos}^{a_i}}{(1 - Tr(a_i))(N_{pos}^{a_i} + N_{neg}^{a_i}) + 2} \\ D_{neg}^{a_i} = \frac{2Tr(a_i)N_{neg}^{a_i}}{(1 - Tr(a_i))(N_{pos}^{a_i} + N_{neg}^{a_i}) + 2} \end{cases} \quad (13)$$

The trustworthiness of seller  $s$  can be calculated as follows:

$$Tr(s) = \frac{[\sum_{i=1}^m D_{pos}^{a_i}] + 1}{[\sum_{i=1}^m (D_{pos}^{a_i} + D_{neg}^{a_i})] + 2} \quad (14)$$

Algorithm 2 (Figure 3) is a pseudo code summary of the method for modeling the trustworthiness of a seller. A seller is considered trustworthy if its trust value is greater than a threshold  $\theta$ . It will be considered untrustworthy if the trust value is less than  $\delta$ . The buyer in our framework will allow only a limited number of the most trustworthy sellers to join the auction. This can be achieved by using the trust thresholds. If there are no trustworthy sellers, the sellers with trust values between  $\theta$  and  $\delta$  may also be allowed to join the auction. Note that these thresholds are subjectively determined by individual buyers. For example, a risk adverse buyer may choose higher values for both  $\theta$  and  $\delta$  so that sellers are required to have very high trust values to be considered as trustworthy, but sellers with slightly low trust values may be considered as untrustworthy. On another hand, a risk taking buyer may choose lower values for the two threshold values. In our experiments, we simply set uniform values for  $\theta$  and  $\delta$ . This simple but fair setting may still satisfy the need of evaluating and comparing different approaches for coping with unfair ratings. For future work, it would be interesting to look into the method that can adaptively update the threshold values according to the different situations.

The buyer will then convey to the central server which sellers it is willing to consider, and the pool of possible sellers is thus reduced. Sellers  $\{s_1, s_2, \dots, s_n\}$  ( $n \geq 1$ ) allowed to join the auction submit their bids by setting the prices and values for the non-price features of the product  $p$ . The buyer will select the winner of the auction as the seller whose product (described in its bid) gives the buyer the largest profit, based on the buyer's valuation of the product  $V_b$ , formalized as follows:

$$S_{win} = \arg \max_{j=1}^n (V_b - P_{s_j}) \quad (15)$$

where  $P_{s_j}$  is the price of product offered by seller  $s_j$ .

Once the buyer has selected the winning seller, it pays that seller the amount indicated in the bid. The winning seller is supposed to deliver the product to the buyer. However, it may decide not to deliver the product. The buyer will report the result of conducting business with the seller to the central server, registering a rating for the seller. It is precisely these ratings of the seller that can then be shared with other buyers.

---

```

Set Trustworthiness of seller  $s = 0$ ;
//Trustworthiness of seller is modeled based on advisors' ratings for the seller
 $\{a_1, a_2, \dots, a_k\}$ : advisors that have provided ratings for seller  $s$ ;
Set  $N_{pos}^a = 0$ : amount of all discounted positive ratings of advisors;
Set  $N_{neg}^a = 0$ : amount of all discounted negative ratings of advisors;
foreach advisor  $a_i$  in  $\{a_1, a_2, \dots, a_k\}$  do
|   Set  $N_{pos}^{a_i} = 0$ : amount of all discounted positive ratings of  $a_i$ ;
|   Set  $N_{neg}^{a_i} = 0$ : amount of all discounted negative ratings of  $a_i$ ;
|   Count  $N_{pos}^{a_i}, N_{neg}^{a_i}$ : number of  $a_i$ 's positive/negative ratings;
|   Set  $D_{pos}^{a_i}$  based on  $D_{pos}^{a_i}$  using Equation 13;
|   Set  $D_{neg}^{a_i}$  based on  $D_{neg}^{a_i}$  using Equation 13;
|_   $N_{pos}^a = N_{pos}^a + D_{pos}^{a_i}; \quad N_{neg}^a = N_{neg}^a + D_{neg}^{a_i}$ 

```

---

*Trustworthiness* is then calculated using Equation 14;

---

Figure 3: Algorithm 2 - buyer  $b$  modeling trustworthiness of a seller  $s$

In this work, we compare our personalized approach with the other two approaches: BRS, TRAVOS. These two approaches are also based on the beta density function. They are also useful for demonstrating the importance of the capabilities, which some of the approaches have and others do not. We implement the TRAVOS model and the personalized approach for modeling the trustworthiness of advisors. Note that TRAVOS does not discount older ratings of sellers. We also implement the BRS approach to filter out unfair ratings for each seller. The aggregation of fair ratings is slightly different from Equation 13 by assuming  $Tr(a_i)$  is always 1 because trustworthiness of advisors is not modeled by BRS.

## 4.2 Simulation Setting

The marketplace operates for a period of 60 days. It involves 90 buyers. These buyers are grouped into three groups. They have different numbers of requests. Each group of buyers has a different number (20, 40 and 60) of requests. In our experiments, we assume that there is only one product in each request and each buyer has a maximum of one request each day. For the purpose of simplicity, we also assume that the products requested by buyers have the same valuation for buyers. After they finish business with sellers, buyers rate sellers. Some dishonest buyers from each group will provide unfair ratings. We allow 2 buyers from each group to leave the marketplace at the end of each day. Accordingly, we also allow 6 buyers to join the marketplace at the end of each day. Some of them may also provide unfair ratings, to keep the percentage of dishonest buyers in each group the same in each day. There are also 6 sellers in total in the marketplace. Each 2 sellers acts dishonestly in different percentages (0%, 25% and 50%) of their business with buyers.

We also set different parameters in the experiments. We set the lower and upper boundaries for BRS to be 0.1 and 0.99 respectively, as recommended in [34]. The number of bins  $N_{bin}$  used by the TRAVOS model is chosen to produce the best results in our experiments. The weight of private reputation used by the personalized approach is also selected to produce the best performance. Note that the personalized approach may not be able to produce the best performance when buyers do not have much experience with the environment. After they obtain enough knowledge about the environment, the best performance becomes achievable. A concrete method for this goal is left for future work (see Section 3.3). In our evaluation, we fairly compare the three approaches by showing their best performance, otherwise, these approaches will be unfairly compared. We set the threshold  $\theta$  to be 0.7 and  $\delta$  to be 0.3. Therefore, a seller is considered trustworthy if its trust value is greater than 0.7 and untrustworthy if it is below 0.3. In our experiments, a buyer is considered to be honest if its trust value is greater than 0.5; otherwise, it is dishonest.

## 4.3 Performance Measurement

We measure the performance of an approach for coping with unfair ratings in two ways. One is its ability to detect dishonest advisors. An effective approach should be able to correctly detect dishonest advisors. This performance can be measured by the false positive rate (FPR) and false negative rate (FNR). A false positive represents that an honest advisor is incorrectly detected as a dishonest advisor. A false negative represents that an advisor is

misclassified as honest but actually is dishonest. The lower values of FPR and FNR imply better performance. We also use Matthew's correlation coefficient (MCC) [20] to measure the approaches' performance in detecting dishonest advisors. MCC is a convenient measure because it gives a single metric for the quality of binary classifications, and is computed as follows:

$$MCC = \frac{(t_p t_n - f_p f_n)}{\sqrt{(t_p + f_p)(t_p + f_n)(t_n + f_p)(t_n + f_n)}} \quad (16)$$

where  $f_p$  = false position,  $t_p$  = true positives,  $f_n$  = false negatives,  $t_n$  = true negatives. An MCC value is between -1 and +1. A coefficient of +1 represents a perfect detection, 0 an average random detection and -1 the worst possible detection.

We also measure the performance of an approach based on how much buyers can benefit if the approach is employed. We use two metrics to represent this benefit, the profit of buyers and the ratio of buyers' successful business with sellers. Eventually, the higher the ratio of successful business the buyers can have with sellers, the larger the profit they will be able to gain.

## 5 Experimental Results and Analysis

In this section, we present experimental results comparing the three approaches, BRS, TRAVOS and the personalized approach. We first provide the comparison of their overall performance. We then analyze how these approaches perform in different scenarios.

### 5.1 Overall Performance Comparison

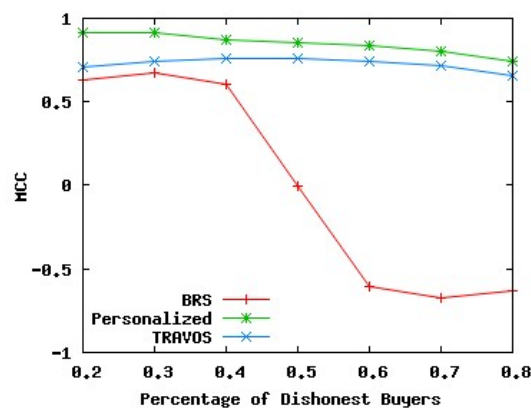


Figure 4: Detecting dishonest buyers

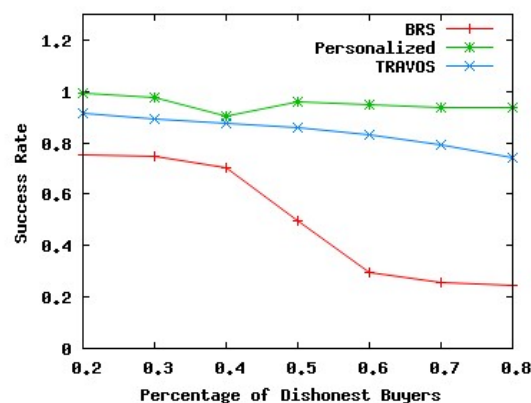


Figure 5: Ratio of successful business

In this experiment, we vary the percentage of dishonest buyers (from 20% to 80%) in the marketplace environment. We then measure the average MCC values for TRAVOS, BRS and the personalized approach for the period of 60 days. Results are shown in Figure 4. From this figure, we can see that the personalized approach produces the highest MCC values for different percentages of dishonest buyers. TRAVOS performs better than BRS. The performance of these approaches will generally decrease when more buyers are dishonest. Note that the

performance of BRS is close to random classification when 50% of buyers are dishonest and becomes much worse when the majority of buyers are dishonest. This result confirms our argument in Sections 2.1 and 2.6.

We measure the ratio of buyers' successful business with sellers. We call a transaction between a buyer and a seller successful business if the seller is honest and delivers what it promised. We measure the success ratio of buyers after 60 days. We then average the success ratio over the total number of buyers in the marketplace (90 in our experiments). In this experiment, we also measure the average total profit of buyers after 60 days.

The profit of a buyer is based on the buyer's valuation for the good and the price of the good. If a buyer does business with an honest seller, the profit of the buyer from this transaction will be calculated as the difference between the value of the product and the price of the product set by the seller. If the buyer does business with a dishonest seller, the profit of the buyer will be reduced by the price of the product.

The results are shown in Figures 5 and 6. These two figures are very similar and also confirm the results shown in Figure 4. Note that the performance of the personalized approach decreases when 40% of the buyers are dishonest. This is because the public reputation component of the personalized approach does not perform well when a large number of buyers are dishonest. When 40% of buyers are dishonest, the personalized approach still considers the public reputation part. Its performance is then affected by the public part. When more than 50% of buyers are dishonest, the personalized approach will rely only on the private component.

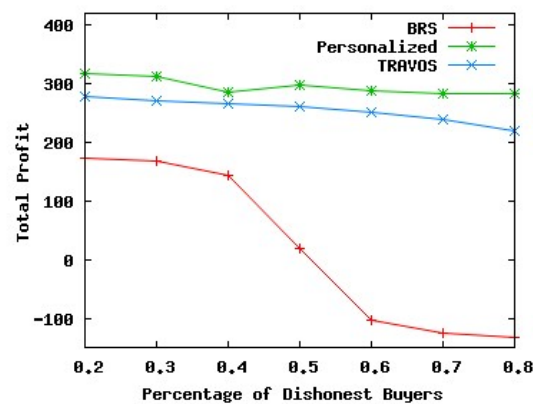


Figure 6: Total profit of buyer

In summary, the personalized approach performs the best. The TRAVOS model performs better than BRS, which is similar to the results in [31]. BRS performs much worse when the majority of buyers are dishonest, which will be further analyzed in depth in the next section. We will also analyze how the three approaches perform in different scenarios.

## 5.2 Analysis of Different Scenarios

In order to further compare the three approaches and analyze their capabilities, we simulate different scenarios where the majority of buyers are dishonest, buyers do not have much experience with sellers in the marketplace, sellers may vary their behavior widely, and buyers may provide a large number of ratings in a short period of time. Note that in this section we will only present the performance of the approaches in detecting dishonest buyers because this performance is correlated with the results of total profit and success ratio of buyers, as presented in the previous section.

### 5.2.1 Dishonest Majority

BRS assumes that a significant majority of the buyers are honest. This is why the performance of BRS decreases dramatically when half of the buyers are liars as shown in Figures 4, 5 and 6.

In order to better see the reasons behind this performance decrease, we show the error of BRS in detecting dishonest buyers when 50% of buyers are dishonest in a period of 120 days, in Figure 7. From this figure, we can see that the ratio of false negatives approaches 0. However, the ratio of false positives continuously increases and approaches 1. This means that BRS tends to label every buyer as dishonest.

Figure 8 explains the statistical foundation of BRS's behavior when 50% of buyers are dishonest. For a honest seller, dishonest buyers provide unfairly low ratings and their Beta distributions reside near 0, according to Equation 3 when  $\beta$  increases. However, for the same seller, honest buyers provide high ratings that make their Beta distributions reside near 1. Overall, the expected value of the aggregated Beta distribution becomes 0.5 and it does not stay

within the margins defined by the lower and upper boundaries of the buyers' Beta distributions. Hence, both the dishonest and honest buyers are regarded as dishonest.

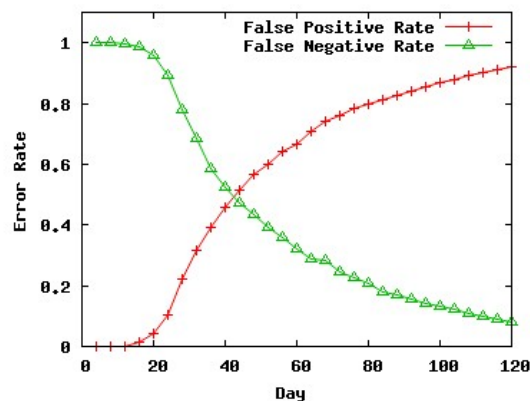


Figure 7: Error rate of BRS

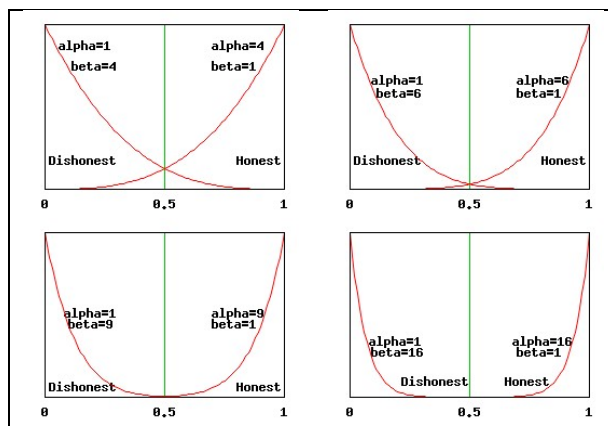


Figure 8: BRS for 50% of dishonest buyers

### 5.2.2 Lack of Personal Experience

The TRAVOS model relies only on buyers' personal knowledge with advisors' advice, whereas BRS and the personalized approach also considers public knowledge of advisors' advice. The public knowledge is useful especially when buyers do not have much experience with sellers, and in consequence do not have much personal knowledge with advisors' advice. In this experiment, we demonstrate the performance of these three approaches in detecting dishonest buyers when 30% of buyers are dishonest. We plot the MCC values of their performance over 60 days, as shown in Figure 9. We can see that both BRS and the personalized approach perform much better than the TRAVOS model in the beginning 10 days. This confirms our argument that buyers should rely on public knowledge about advisors when they do not have much experience with sellers. We also can see from Figure 9 that the performance of BRS will decrease after 30 days and become worse than that of TRAVOS. The reason for this will be further analyzed and explained in Section 5.2.4.

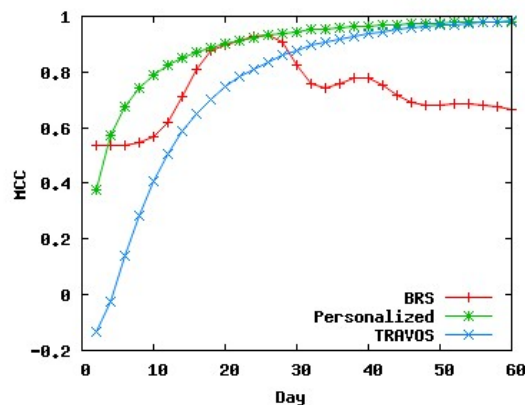


Figure 9: Detecting dishonest buyers

We also carry out an experiment to demonstrate the advantage of our personalized approach that combines private and public reputation components (the combined version) versus that uses only the private reputation (the private version), in the scenarios where buyers have and do not have much experience with sellers respectively. We measure the false negative rates of the combined version and the private version of our approach. Two conclusions can be drawn from the results shown in Figure 10. One is that the combined version has lower false negative rate and thus is more effective than the private version no matter how much experience buyers have with sellers. We can also see that the difference between the combined and the private versions is larger when buyers do not have much experience with sellers. This shows the greater advantage of our personalized approach in the lack of experience scenario. The second conclusion can be further clearly seen from the following experiment when comparing our personalized approach with TRAVOS, because TRAVOS is similar to the private version of our approach in the sense that they both rely only on buyers' personal knowledge with advisors' advice for modeling the trustworthiness of advisors.

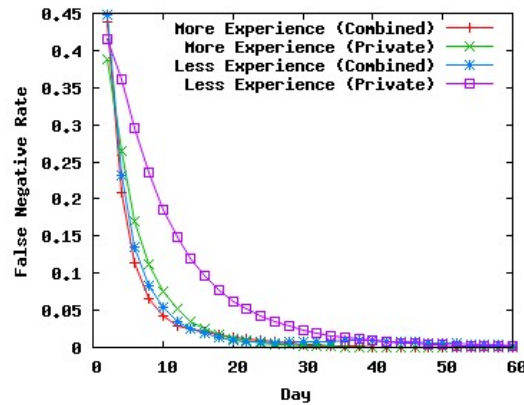


Figure 10: Combined vs. private for the personalized approach

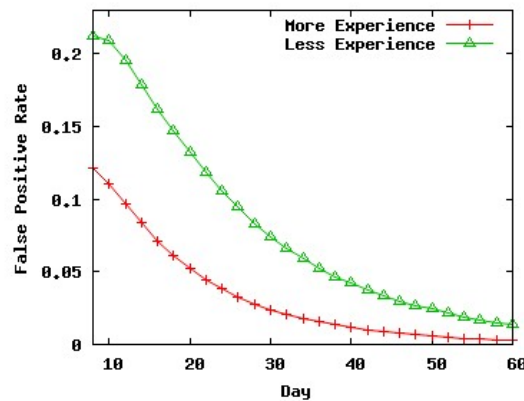


Figure 11: FPR: Personalized vs. TRAVOS

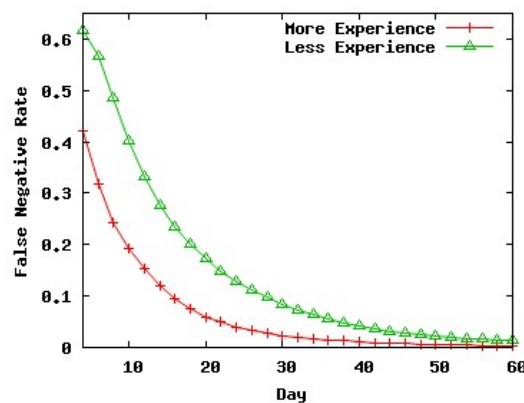


Figure 12: FNR: Personalized vs. TRAVOS

In the third experiment, we directly compare the performance of the personalized approach with that of TRAVOS in the scenario where buyers do not have much experience with sellers. In the experimental setting, 30% of buyers are dishonest. Half of all buyers have more requests for products and another half have fewer requests. Buyers having more requests will have more experience with sellers. We measure how much the personalized approach outperforms TRAVOS in detecting dishonest buyers.

Results are shown in Figures 11, 12 and 13. In both cases when buyers have more or less experience with sellers, the personalized approach outperforms TRAVOS. From the figures, we can see that the difference in FPR, FNR and MCC is larger when buyers do not have much experience with sellers. The performance difference will decrease day after day because buyers will have more and more experience with sellers. This suggests that an approach of modeling the trustworthiness of advisors for coping with unfair ratings should rely on public knowledge of advisors' advice as well as when buyers do not have much experience with sellers.

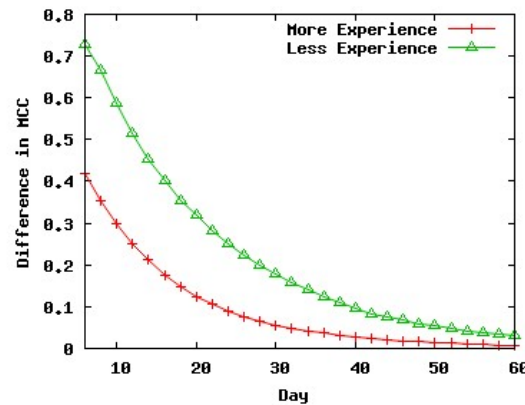


Figure 13: MCC: Personalized vs. TRAVOS

### 5.2.3 Seller Varying Behavior

The personalized approach introduces the concept of a time window when evaluating the trustworthiness of advisors. For example, it only compares a buyer's and an advisor's ratings if these two ratings are within the same time window when computing the private reputation of the advisor, by setting  $\lambda$  in Equation 5 to be 0. This is to deal with the problem when sellers vary their behavior widely. However, as we point out in Section 2.2, the TRAVOS model is not able to deal with this problem. In this section, we present experimental results to confirm this argument.

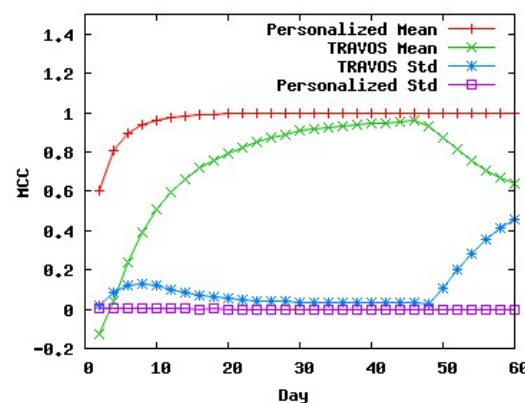


Figure 14: Personalized vs. TRAVOS

We first carry out an experiment to compare the personalized approach with the TRAVOS model in the situation where sellers may change their behavior. In this experiment, the sellers that vary their behavior will be dishonest in 25% or 50% of the period of 60 days. We also have three types of sellers. The first type of sellers act dishonestly in a uniform manner. The second type of sellers is honest first and then becomes dishonest. The third type of sellers acts dishonestly first and then honestly later on. We run simulations separately 500 times for each type of seller and average the results. We then calculate the mean and standard deviation of the two approaches' performance in detecting dishonest buyers. Note that  $\lambda$  in Equation 5 is set to 0, because the seller behavior is varying so much.

From the results shown in Figure 14, we can see that the mean performance of the personalized approach consistently increases after each day. The standard deviation of its performance stays nearly at 0, which implies that the performance of the personalized approach is not affected by sellers' varying behavior. However, the



mean performance of the TRAVOS model decreases heavily after 45 days and the standard deviation of its performance is considerably large for the first 15 days and the last 15 days. Therefore, TRAVOS does not perform well when sellers change their behavior widely.

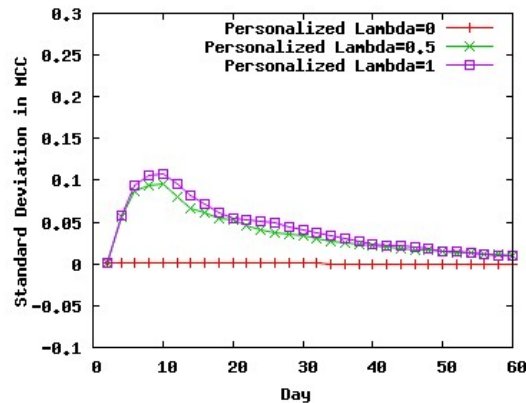


Figure 15: Different  $\lambda$  values for the personalized approach

We then explore how different values of  $\lambda$  may affect our personalized approach's capability in coping with sellers' varying behavior. Results in Figure 15 show that when  $\lambda$  value is larger, the standard deviation of our personalized approach's performance is higher, which implies that the approach is affected more by sellers' varying behavior. Note that because the standard deviation of our personalized approach is generally small (less than 0.12), the effect of different  $\lambda$  values on the mean performance of the personalized approach is also small and not shown in the figure.

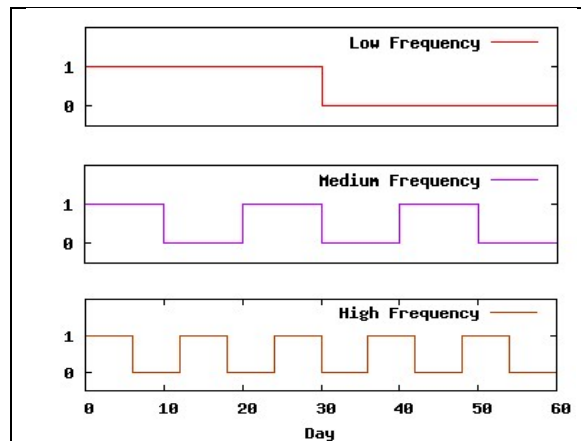


Figure 16: Seller varying behavior

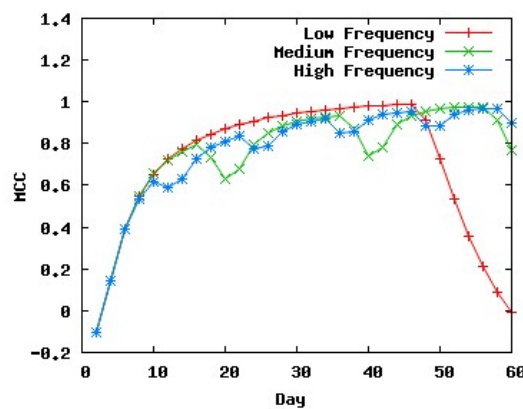


Figure 17: Performance of TRAVOS

We also carry out another experiment to analyze in depth how the TRAVOS model will be affected by different types of seller varying behavior. In this experiment, we have sellers vary their behavior in different frequencies. All sellers in this experiment will act honestly first and then dishonestly later on. These different types of sellers vary their behavior for 1, 3 and 5 times respectively within the period of 60 days, as shown in Figure 16. This figure shows an example how a seller that is dishonest in 50% of the period of 60 days will vary its behavior. A seller's honesty of 1 on the vertical axis means that the seller acts honestly in the corresponding day and 0 represents dishonest behavior.

The performance of TRAVOS for different frequencies of seller changing behavior is presented in Figure 17. When sellers change their behavior very frequently, the performance of TRAVOS will also change more often. The change of its performance is less than that when sellers vary behavior less frequently. When the sellers change their behavior only once from being honest to being dishonest, the performance decreases to a great extent to nearly a random classification.

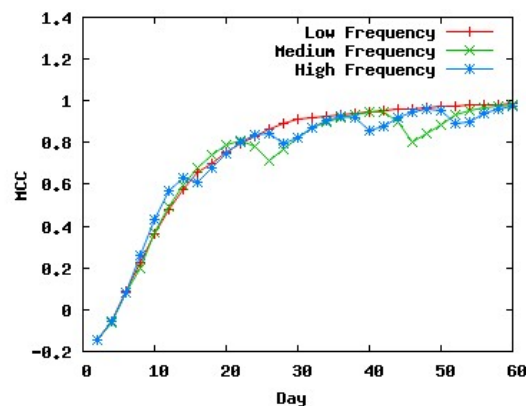


Figure 18: Performance of TRAVOS

We also show the results of the performance of TRAVOS when all sellers act dishonestly first and then honestly later on. Similarly, sellers vary their behavior in different frequencies. The results are shown in Figure 18. Comparing this figure with Figure 17, we can see that the performance of TRAVOS is affected less than that in the situation where sellers act honestly first and then dishonestly. Especially when sellers vary their behavior at a low frequency, the performance of TRAVOS does not have much change compared to that in Figure 17. In the simulation framework, sellers acting dishonestly at the beginning will have very low trust values and be prevented from doing business with buyers. The changes of their behavior will no longer affect the performance of detecting dishonest buyers. This also implies that a more effective varying behavior for a seller is to be honest first to build up its trustworthiness, and to then act dishonestly to exploit the marketplace (a behavior explored by such trust researchers as Tran and Cohen [32], and Sen and Banerjee [28]).

### 5.2.4 Buyers' Flooding

Buyers' flooding is the situation where buyers (advisors) may provide a large number of ratings for a seller in a short period of time. To deal with this situation, for example, the personalized approach uses the concept of a time window and considers only a limited number of ratings from one buyer for the seller within the same time window. As discussed in Section 2.6, the BRS approach will be heavily affected by buyers' flooding. In the case where buyers provide a large number of unfair ratings, BRS will suffer from the dishonest majority problem as demonstrated in previous sections. In this section, we carry out experiments to show that BRS is affected even when buyers provide a large number of fair ratings within a short period of time.

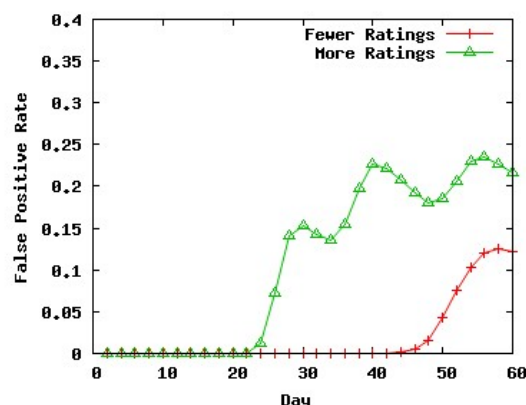


Figure 19: False positive rate of BRS

In this experiment, we involve two types of buyers. The first type of buyers has much more requests and therefore will provide a lot of ratings to sellers. The second type of buyers provide fewer ratings. In both cases, 20% of buyers are dishonest. We run simulations for the two cases separately and measure the false positive rate of BRS in detecting dishonest buyers. Results are shown in Figure 19. We can see that after 20 or 40 days, BRS will start incorrectly classifying honest buyers as dishonest. The false positive rate is higher when buyers provide more ratings. Therefore, BRS is even affected by the situation where buyers may provide a large number of fair ratings.

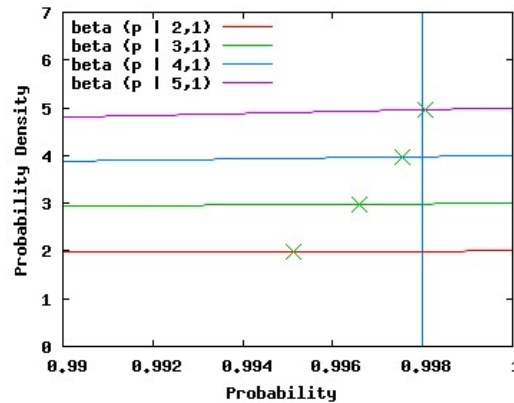


Figure 20: BRS unable to cope with flooding

We further analyze the statistical foundation of this phenomenon, as shown in Figure 20. The vertical line on the figure represents the expected value (trustworthiness) of a seller when there are 500 positive ratings and 0 negative ratings provided by buyers for the seller. This figure also shows the beta distributions for buyers that provide 1, 2, 3, and 4 positive ratings respectively, and 0 negative ratings for the seller. The “x” symbols on the distributions represent the cut-off points of upper bounds of these distributions. We can see from the figure that the seller's expected value only falls within the upper bounds of the distribution with 4 positive ratings. Therefore, the honest buyers that have only provided 1, 2 or 3 positive ratings will be incorrectly classified as dishonest buyers. This therefore increases the false positive rate of BRS.

### 5.3 Summary of Results

We have carried out experiments to compare the overall performance of the three representative approaches, TRAVOS, BRS and the personalized approach. We measure their accuracy in detecting dishonest buyers, the ratio of buyers' successful business with sellers when these approaches are employed, and the total profit of buyers. Results show that the personalized approach performs the best, TRAVOS performs better than BRS, and BRS performs much worse when the majority of buyers are dishonest.

We also analyze how these three approaches perform in different scenarios. Results show that the personalized approach performs much better than TRAVOS especially when buyers do not have much experience with sellers. In this case, BRS also performs better than TRAVOS when the majority of buyers are honest. TRAVOS suffers from the situation where sellers may vary their behavior, and is heavily affected especially when sellers first build up their trust by being honest and then act dishonestly. BRS is shown to be ineffective when buyers provide a large number of ratings for a seller.

## 6 Conclusions and Future Work

In this paper, we presented a personalized approach for effectively handling unfair ratings in centralized reputation systems. It allows a buying agent to estimate the private reputation of an advisor agent based on their ratings for commonly rated selling agents. When the buying agent is not confident with the private reputation value, it can also use the public reputation of the advisor. The public reputation of the advisor is evaluated based on all ratings for the selling agents rated by the advisor agent. Compared with other trust and reputation modeling approaches summarized in Section 2, our personalized approach for modeling the trustworthiness of advisors has all of the desirable features that we outlined in Section 2.6. It is able to cope with unfair ratings even when the majority of the ratings of a seller is unfair. It is able to deal with the situation where advisors may provide a large number of ratings within a short period of time. It is effective even when buyers do not have much experience with sellers and is also able to deal with changes of agents' behavior over time. These capabilities of our approach are further demonstrated through experiments.

We then focus on experimental comparison with the representative approaches, including BRS and TRAVOS. Instead of using the ART Testbed [5] that is proposed to provide unified performance benchmarks for comparing trust and reputation modeling approaches, we propose a framework that simulates a dynamic electronic marketplace

environment involving possibly deceptive buying and selling agents. The current ART Testbed specification is in an artwork appraisal domain where appraisers want to buy artwork about which they may have limited knowledge. They may then seek information about artwork from other appraisers (opinion providers). Opinion providers may choose to lie about the true value of the artwork. The appraisers will model the trustworthiness of opinion providers based on their own knowledge about the opinion providers or reputation opinion of other appraisers (reputation providers). These reputation providers may choose to lie about opinion providers' true trust values. An approach for coping with untruthful reputation opinions from opinion providers may then be integrated and evaluated by the ART Testbed. However, integrating TRAVOS, BRS and the personalized approach into the testbed is challenging. These approaches are developed for a rather simpler e-marketplace environment. They allow only binary ratings to represent simple and objective results of transactions between sellers and buyers (advisors). Advisors modeled by these approaches do not make profit from providing advice or pay cost to generate advice. Overly simplifying the ART Testbed may lose its advantages, and adapting these approaches to the complicated testbed may change their original design. Furthermore, the winning approach IAM [30] for the 2006 ART Testbed competition does not even consider reputation opinions from other appraisers. This decision raises the concern about the importance of an approach for coping with untruthful reputation opinions in this testbed, and whether the results of comparing the approaches based on this testbed will be significant.

The approaches of BRS, TRAVOS and our personalized approach are compared for the first time in terms of their capabilities for detecting dishonest buyers. Total profit of buyers is also the most direct and important measure used in the comparison between these approaches. We further specifically examine different scenarios, including ones where the majority of buyers are dishonest, buyers lack personal experience with sellers, sellers may vary their behavior, and buyers may provide a lot of ratings. Such an empirical study is useful for highlighting the importance of the capabilities of our personalized approach.

In the current evaluation framework, the decision of a winning seller is based only on the bids submitted by sellers once the sellers are allowed to join the auction because they are considered as trustworthy by the buyer. Adding information about sellers' trustworthiness into winner selection for the auctions may introduce some challenges. For example, when sellers submit bids, they may do so trying to incorporate reasoning about the trustworthiness of their competitors (as well as their view of their own trustworthiness). The work in [16] provides a comprehensive study on how to effectively use the information of sellers' trustworthiness to enhance the auctions. The work of Hazard and Singh [7] also provides some hints about how to discount buyer utility based on sellers' trustworthiness. Their study results may be applied in the future improvement on our evaluation framework.

For future work, in our evaluations, it would be worthwhile to explore the case where some dishonest buyers lie only for some sellers while being honest for other sellers. It would also be worthwhile to consider other types of dishonest buyers from the literature, such as the Exaggerated Positive and Exaggerated Negative types defined in [37], [29]. The performance of detecting these types of dishonest buyers would then be evaluated and compared for those approaches.

We may want to investigate more advanced dishonest buyers that are strategic [14], [15]. For example, some dishonest buyers may have mixed lying types. Inspired by the evaluation in [34], a marketplace may involve some buyers that have an adaptive lying strategy where buyers may learn from the marketplace and build some strategies to adapt their lying types or lying frequency. A similar idea can be found in the work of Sen and Banerjee [28] and the work of Feng et al. [3], where strategic agents may exploit the marketplace. We are interested in demonstrating how the existing approaches perform in this kind of marketplace environment. We are also interested in seeing how well they handle marketplaces where strategic agents collude with each other [13].

Coping with unfair ratings from advisors in e-marketplaces by a modeling of their trustworthiness has some similarity with the challenge of addressing shilling attacks in recommender systems [27]. The research of [17] suggests that the general algorithms used by attackers (i.e. the kind of attacks) may be useful to model and that the areas being attacked (e.g. low use items) may influence the possible damage that can be inflicted. For future work, it would be useful to simulate these attacks and the attacks summarized in [8], [10], to compare the robustness of the approaches against the attacks.

Introducing innovation to the design of trust modeling systems used in agent-oriented e-marketplaces is a crucial concern, as part of the ongoing effort to promote electronic commerce to businesses and organizations. This paper has demonstrated some key shortcomings of existing trust modeling systems and has discussed the advantages introduced by our particular personalized approach. As a result, specific directions are now available for users who are selecting trust modeling algorithms to run in e-marketplaces.

## References

- [1] R. Alnemr, S. Koenig, T. Eymann, and C. Meinel, Enabling usage control through reputation objects: A discussion on e-commerce and the internet of services environments, *Journal of Theoretical and Applied Electronic Commerce Research*, vol. 5, no. 2, pp. 59-76, 2010.

- [2] C. Dellarocas, Immunizing online reputation reporting systems against unfair ratings and discriminatory behavior, in Proceedings of the Second ACM Conference on Electronic Commerce (EC), Minneapolis, 2000, pp. 150-157.
- [3] Q. Feng, Y. Sun, L. Liu, Y. Yang, and Y. Dai, Voting systems with trust mechanisms in cyberspace: Vulnerabilities and defenses, IEEE Transactions on Knowledge and Data Engineering, vol. 22, no. 12, pp. 1766-1780, 2010.
- [4] A. Fernandes, E. Kotsovinos, S. Östring, and B. Dragovic, Pinocchio: Incentives for honest participation in distributed trust management, in Proceedings of the Second International Conference on Trust Management (iTrust), Oxford, U.K., 2004, pp. 63-77.
- [5] K. K. Fullam, T. B. Klos, G. Muller, J. Sabater, A. Schlosser, Z. Topol, K. S. Barber, J. S. Rosenschein, L. Vercouter, and M. Voss, A specification of the agent reputation and trust (art) testbed: Experimentation and competition for trust in agent societies, in Proceedings of the Fourth International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS), Utrecht, Netherlands, 2005, pp. 512-518.
- [6] N. Gal-Oz, E. Gudes, and D. Hendler, A robust and knot-aware trust-based reputation model, in Proceedings of the IFIP WG 11.11 International Conference on Trust Management (IFIPTM), Trondheim, Norway, 2008, pp. 167-182.
- [7] C. J. Hazard and M. P. Singh, Intertemporal discount factors as a measure of trustworthiness in electronic commerce, IEEE Transactions on Knowledge and Data Engineering, vol. 23, no. 5, pp. 699-712, 2011.
- [8] K. Hoffman, D. Zage, and C. Nita-Rotaru, A survey of attack and defense techniques for reputation systems, ACM Computing Surveys, vol. 42, no. 1, pp. 1-31, 2009.
- [9] A. Jøsang, A logic for uncertain probabilities, International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, vol. 9, no. 3, pp. 279-311, 2001.
- [10] A. Jøsang and J. Golbeck, Challenges for robust trust and reputation systems, in Proceedings of 5th International Workshop on Security and Trust Management (STM), Saint Malo, France, 2009.
- [11] A. Jøsang and R. Ismail, The beta reputation system, in Proceedings of the Fifteenth Bled Electronic Commerce Conference, Bled, Slovenia, 2002, pp. 324-337.
- [12] R. Jurca and B. Faltings, Eliciting truthful feedback for binary reputation mechanisms, in Proceedings of the 2004 IEEE/WIC/ACM International Conference on Web Intelligence, Washington, DC, USA, 2004, pp. 214-220.
- [13] R. Kerr, Coalition detection and identification, in Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems (AAMAS), Toronto, Canada, 2010, pp. 1657-1658.
- [14] R. Kerr and R. Cohen, Smart cheaters do prosper: Defeating trust and reputation systems, in Proceedings of the 8th International Conference on Autonomous Agents and Multiagent Systems (AAMAS), Budapest, Hungary, 2009, pp. 993-1000.
- [15] R. Kerr and R. Cohen, An experimental testbed for evaluation of trust and reputation systems, in Proceedings of the Third IFIP WG 11.11 International Conference on Trust Management (IFIPTM'09), West Lafayette, IN, 2009, pp. 252-266.
- [16] S. König, S. Hudert, T. Eymann, and M. Paolucci, Towards reputation enhanced electronic negotiations for service oriented computing, in Proceedings of the 2008 10th IEEE Conference on E-Commerce Technology and the Fifth IEEE Conference on Enterprise Computing, E-Commerce and E-Services, Washington, DC, 2008, pp. 273-291.
- [17] S. K. Lam and J. Riedl, Shilling recommender systems for fun and profit, in Proceedings of the Thirteenth International Conference on World Wide Web, New York, USA, 2004, pp. 393-402.
- [18] S. Liu, C. Miao, Y.-L. Theng, and A. C. Kot, A clustering approach to filtering unfair testimonies for reputation systems, in Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems (AAMAS), Toronto, Canada, 2010, pp. 1577-1578.
- [19] S. Liu, J. Zhang, C. Miao, Y.-L. Theng, and A. C. Kot, iCLUB: An integrated clustering-based approach to improve the robustness of reputation systems, in Proceedings of the 10th International Conference on Autonomous Agents and Multiagent Systems (AAMAS), Taipei, Taiwan, 2011, pp. 1151-1152.
- [20] B. W. Matthews, Comparison of the predicted and observed secondary structure of t4 phage lysozyme, Biochim. Biophys. Acta, vol. 405, no. 2, pp.442-451, 1975.
- [21] L. Mui, M. Mohtashemi, and A. Halberstadt, A computational model of trust and reputation, in Proceedings of the Thirty Fifth Hawaii International Conference on System Science (HICSS), Big Island, Hawaii, pp. 2431-2439, 2002.
- [22] Z. Noorian, S. Marsh, M. Fleming, Multi-layered cognitive filtering by behavioural modeling, in Proceedings of the 10th International Conference on Autonomous Agents and Multiagent Systems (AAMAS), Taipei, Taiwan, 2011, pp. 871-878.
- [23] Z. Noorian and M. Ulieru, The state of the art in trust and reputation systems: A framework for comparison, Journal of Theoretical and Applied Electronic Commerce Research (JTAER), vol. 5, no. 2, pp. 97-117, 2010.
- [24] K. Regan, T. Tran, and R. Cohen, Sharing models of sellers amongst buying agents in electronic marketplaces, in Proceedings of the 10th International Conference on User Modeling (UM'2005) Workshop on Decentralized Agent-Based and Social Approaches to User Modeling, Edinburgh, UK, 2005, pp. 75-79.
- [25] S. Russell and P. Norvig, Artificial Intelligence: A Modern Approach. Second Edition, Prentice Hall, Englewood Cliffs, New Jersey, 2002.
- [26] J. Sabater and C. Sierra, Regret: A reputation model for gregarious societies, in Proceedings of the Fifth International Conference on Autonomous Agents Workshop on Deception, Fraud and Trust in Agent Societies, Montreal, Canada, 2001, pp. 61-69.

- [27] J. J. Sandvig, B. Mobasher, and R. Burke, A survey of collaborative recommendation and the robustness of model-based algorithms, *IEEE Data Engineering Bulletin*, vol. 31, no. 1, pp. 3-13, 2008.
- [28] S. Sen and D. Banerjee, Monopolizing markets by exploiting trust, in *Proceedings of the Fifth International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, Hakodate, Japan, 2006, pp. 1249-1256.
- [29] M. Sensoy and P. Yolum, Experimental evaluation of deceptive information filtering in context-aware service selection, in *Proceedings of the 11th International Workshop on Trust in Agent Societies, Lecture Notes in Artificial Intelligence*, Estoril, Portugal, 2008, pp. 326-347.
- [30] W. T. L. Teacy, T. D. Huynh, R. K. Dash, N. R. Jennings, J. Patel, and M. Luck, The art of IAM: The winning strategy for the 2006 competition, in *The Workshop on Trust in Agent Societies at The Sixth International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, Honolulu, Hawaii, 2007, pp. 102-111.
- [31] W. T. L. Teacy, J. Patel, N. R. Jennings, and M. Luck, Coping with inaccurate reputation sources: Experimental analysis of a probabilistic trust model, in *Proceedings of Fourth International Autonomous Agents and Multiagent Systems (AAMAS)*, Netherlands, 2005, pp. 997-1004.
- [32] T. Tran and R. Cohen, Improving user satisfaction in agent-based electronic marketplaces by reputation modeling and adjustable product quality, in *Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, New York, 2004, pp. 828-835.
- [33] Y. Wang and J. Vassileva, Bayesian network-based trust model, in *Proceedings of the Sixth International Workshop on Trust, Privacy, Deception and Fraud in Agent Systems*, Halifax, NS, Canada, 2003, pp. 372-380.
- [34] A. Whitby, A. Jøsang, and J. Indulska, Filtering out unfair ratings in bayesian reputation systems, in *Proceedings of the 7th Int. Workshop on Trust in Agent Societies (at AAMAS'04)*, Rome, 2004, pp. 106-117.
- [35] B. Yu and M. P. Singh, A social mechanism of reputation management in electronic communities, in *Proceedings of the 4th International Workshop on Cooperative Information Agents*, Boston, USA, 2000, pp. 154-165.
- [36] B. Yu and M. P. Singh, An evidential model of distributed reputation management, in *Proceedings of International Autonomous Agents and Multi Agent Systems (AAMAS)*, Bologna, Italy, 2002, pp. 294-301.
- [37] B. Yu and M. P. Singh, Detecting deception in reputation management, in *Proceedings of International Autonomous Agents and Multi Agent Systems (AAMAS)*, Melbourne, Australia, 2003, pp. 73-80.
- [38] G. Zacharia, A. Moukas, and P. Maes, Collaborative reputation mechanisms in electronic marketplaces, in *Proceedings of the 32nd Hawaii International Conference on System Sciences (HICSS-32)*, pp. 8026-8032, Maui, Hawaii, 1999, pp. 8026-8032.
- [39] J. Zhang and R. Cohen, Trusting advice from other buyers in e-marketplaces: The problem of unfair ratings, in *Proceedings of the Eighth International Conference on Electronic Commerce (ICEC'06)*, Fredericton, New Brunswick, Canada, 2006, pp. 225-234.
- [40] J. Zhang and R. Cohen, Evaluating the trustworthiness of advice about selling agents in e-marketplaces: A personalized approach, *Electronic Commerce Research and Applications*, vol. 7, no. 3, pp. 330-340, 2008.
- [41] J. Zhang, R. Cohen, and K. Larson, Leveraging a social network of trust for promoting honesty in e-marketplaces, in *Proceedings of the IFIP WG 11.11 International Conference on Trust Management (IFIPTM)*, Morioka, Iwate, Japan, 2010, pp. 216-231.