

Protecting Online Rating Systems from Unfair Ratings

Jianshu Weng, Chunyan Miao, and Angela Goh

School of Computer Engineering, NTU
{weng0004, ASCYMiao, ASESgoh}@ntu.edu.sg

Abstract. Online rating systems have been widely adopted by online trading communities to ban “bad” service providers and prompt them to provide “good” services. However, the performance of the online rating systems is easily compromised by various unfair ratings, e.g. balloting, badmouthing, and complementary unfair ratings. How to mitigate the influence of the unfair ratings remains an important issue in online rating systems. In this paper, we propose a novel entropy-based method to measure the rating quality as well as to screen the unfair ratings. Experimental results show that the proposed method is both effective and practical in alleviating the influence of different types of unfair ratings.

1 Introduction

With the development of Internet, a large number of people carry out transactions through online trading communities such as eBay. Nevertheless, people still regard online trading as a risk since it is hard to determine whether to trust various online sellers before the transactions [1]. Therefore, reputation mechanisms have been introduced into online trading communities to establish trust between sellers and buyers [2]. One way to establish such mechanism is the online rating systems [2]. The main idea of online rating systems is to allow each buyer give a rating for the seller after each transaction. The existing ratings of a particular seller will then be used by the potential buyers to derive the seller’s reputation score, which serves as an indicator whether the seller will provide “good” service or not in future transactions.

Online rating systems have already been adopted by many online trading communities, e.g. eBay [1], and have been credited to their success. Despite the wide adoption of online rating systems, there are still some open issues, especially the issue of *unfair ratings*. That is, buyers might give ratings which are different from their real experiences, e.g. although seller provides “good” service in one transaction, buyer gives rating as “bad”, and vice versa. The performance of online rating systems would easily be compromised by unfair ratings [2].

Finding effective ways to guard against unfair ratings has attracted many research efforts in recent years, e.g. [3,4,5]. Most of the existing methods depend on assumptions that sellers’ behaviors (as well as buyers’ ratings) follow a particular distribution, which hinders their general application to other settings,

e.g. Beta distribution [4,5]. In this paper, we propose an entropy-based method to tackle the issue of unfair ratings in online rating systems. An entropy-based metric is designed to measure the rating quality based on which unfair ratings can be further screened. Unlike existing methods, the proposed entropy-based method is distribution-free. It does not make any assumption regarding the distribution of the ratings. In our current research, the proposed method is explored in context of *Bayesian* rating system. Nevertheless, the proposed method is not limited to *Bayesian* rating system. It can be easily extended to other types of rating systems due to its distribution-independent nature.

The remainder of this paper is organized as follows. A review of related work is given in Section 2. Section 3 gives a brief review of the *Bayesian* rating system. Section 4 presents the proposed method to screen unfair ratings in the context of *Bayesian* rating system. The effectiveness of the proposed method has been shown through experimental results in Section 5. Finally, Section 6 concludes the paper with an overview of future work.

2 Related Work

Online rating systems have played an important role in many online trading communities, e.g. eBay [1]. The presence of unfair ratings is a threat to online rating systems [2]. Some methods have been proposed to address this issue.

Whitby et al. [4] propose to screen ratings reported by others (i.e. testimonies) in Bayesian rating system by determining whether a testimony is outside the $q\%$ quantile and $(1 - q)\%$ quantile of the majority opinion. If it is, the testimony is considered as an unfair rating and will be excluded. Then the majority opinion will be calculated again with the remaining testimonies. This process is carried out iteratively until no other testimony can be excluded. One major limitation of this method is that it does not scale well with the increase in the number of testimonies due to its iterative nature. Moreover, this work depends on an assumption that the ratings follow a Beta distribution. However, it is not easy to justify this assumption especially in the scenarios where few observations are available in the ratings (either local rating or testimonies).

Buchegger and Boudec [5] propose a method to address the issue of unfair ratings in the context of mobile ad-hoc network. This method has two main limitations. First, this method does not consider the majority opinions when screening testimonies. Instead, testimonies that are different from a node's own experience (i.e. local rating) are rejected. This may not be true in general, since one single node's experience might not reflect the target node's behavior. Secondly, it is also based on an assumption that nodes' behavior follows a Beta distribution.

Garg et al. [6] developed a reputation system in context of structured P2P network. After one peer interacts with the target peer, it rates the target peer and sends the rating to all the M score managers who are responsible for calculating and answering other peers' query of the target node's reputation score. The M score managers then aggregate ratings from all peers who have report

testimonies on the target peer and calculate the target peer’s reputation score. When calculating the reputation score, each testimony is given a weight based on the credibility of the peer who reports the rating. The credibility is determined based on the difference between this peer’s testimony and score managers’ aggregated rating (i.e. majority opinion). The limitation of this work is that it also assumes that peers’ behaviors in the reputation system follow a normal distribution.

In contrast, the proposed method in this paper applies an entropy-based metric to screen the testimonies. It does not make any assumption regarding the distribution of ratings. The proposed method also takes the majority opinion into account to make the screening more accurate. And more desirable, it scales linearly with the increase in the number of available testimonies.

3 Bayesian Rating System

In our current research, we explore an entropy-based method for filtering unfair ratings in the context of Bayesian rating systems. Before presenting the proposed method, this section reviews Bayesian rating systems. We reiterate that the proposed method is not limited to Bayesian rating system. It can be easily extended to other types of rating systems due to its distribution-independent nature.

There are primarily two components in *Bayesian* rating system [7]: one for collecting seller’s behaviors in the past transactions, another for predicting seller’s behaviors in the future transactions.

3.1 Collecting Seller’s Past Behaviors

With Bayesian rating systems, buyers give feedbacks of the seller’s behavior after a transaction is cleared. Buyer assigns a positive rating of “1” to the seller if he thinks that the seller provided a “good” service; otherwise it assigns a negative rating of “0”. Buyer B ’s rating for seller S in transaction T can be presented in vector notion as:

$$r_{BS}^T = \begin{bmatrix} p \\ n \end{bmatrix}, \quad \text{where } \begin{bmatrix} p \\ n \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (1)$$

Instead of maintaining ratings of all the past transactions, a buyer usually only maintains a summary of ratings for the past transactions within a window of size W . This is reasonable since a seller’s behavior is usually changing from one transaction to another. Moreover, buyers generally choose to “care” more about a seller’s recent behavior and “forget” its past behavior. By introducing a forgetting factor (λ), which controls the rate that the seller’s old behaviors are forgotten, the summary of ratings within the window can be represented in vector notion as:

$$r_{BS} = \begin{bmatrix} pf \\ nf \end{bmatrix} = \sum_{T=T_c-W+1}^{T_c} \lambda^{T_c-T} \begin{bmatrix} p \\ n \end{bmatrix}, \quad (T_c - T) \leq W \quad (2)$$

where W is the window size, T_c is the latest transaction, and T is the transaction after which rating was collected, and r_{BS} is termed as the B 's local rating of S .

3.2 Predicting Seller's Behaviors

In *Bayesian* rating system, it is assumed that seller's behavior (and buyer's rating as well) follows a Beta distribution. The probability density function (PDF) of Beta distribution is given by:

$$\text{beta}(Pr|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} Pr^{\alpha-1} (1 - Pr)^{\beta-1}; \quad 0 < Pr < 1, \alpha \geq 1, \beta \geq 1 \quad (3)$$

where $B(\alpha, \beta)$ is the beta function. This PDF expresses the probability (Pr) that a seller will provide "good" services in future transactions. Then buyers predict the probability that the seller will provide "good" service in the next transaction as the expectation value of the Beta distribution, which is given by:

$$E(Pr) = \frac{\alpha}{\alpha + \beta}, \quad \alpha = pf + 1, \quad \beta = nf + 1 \quad (4)$$

When a buyer has not transacted with a particular seller before, $pf = nf = 0$, α and β are set to be 1 correspondingly, which causes $E(Pr) = 0.5$. It is interpreted that the seller has equal probabilities of providing "good" or "bad" service. Then with the update of pf and nf after each transaction, the buyer also updates its prediction of the seller's behavior.

3.3 Problems Caused by Unfair Ratings

When making prediction of a seller's behavior, a buyer will request and aggregate ratings from other buyers who have transacted with the same seller before [7,4]. The initiating buyer B , who is predicting the seller's behavior, will send out a requesting message first. Upon receiving the requesting message from the initiating buyer, the answering buyers will simply reply as testimonies their local ratings of the target seller (if any)¹. The initiating buyer then updates α and β by aggregating all the returned testimonies with its local rating (if any). That is:

$$\alpha = 1 + pf_B + \sum_{X \in C} pf_X, \quad \beta = 1 + nf_B + \sum_{X \in C} nf_X \quad (5)$$

where pf_B and nf_B denote the number of positive and negative ratings in B 's local rating respectively, C denotes the set of the answering buyers, and pf_X, nf_X refer to the numbers of positive and negative ratings in the testimony returned by a particular buyer X in the set C . Then the initiating buyer updates its prediction of the seller's behavior using Eq. (4).

¹ A good testimony propagation algorithm is expected to scale with the size of the community [2]. Although it is also a very important issue, it is not the focus of this paper. Instead, we assume in this paper that the initiating buyer can always receive the testimonies they need.

However, if the initiating buyer aggregates all returned testimonies blindly, the answering buyers can easily diverge the initiating buyer's predictions by reporting unfair ratings that are different from their real experiences. For example, buyer B is now evaluating whether to buy from a potential seller S . B requests testimonies from other buyers who have interacted with S before. S colludes with some buyers, who report unfairly higher ratings than the real quality of services that S delivered. Those unfair positive ratings will increase the value of α and decrease the value of β in Eq. (5), which immediately leads to an increase in B 's predicted probability that S will provide "good" service. As this simple example shows, the performance of online rating systems would easily be compromised by the presence of *unfair ratings* [4,8]. Making the rating systems robust to avoid or mitigate the influence of unfair ratings is a fundamental issue in building online rating systems [4,8].

4 Entropy-Based Ratings Screening

Motivated by the problems caused by unfair ratings in current online rating system, in this paper, we propose an entropy-based metric to measure the quality of the ratings (both local rating and testimonies), based on which to screen ratings and to mitigate the influence of unfair ratings. The basic idea of the proposed method is that: if, compared with the quality of the already-aggregated testimonies (i.e. majority opinion), there is a significant quality improvement or downgrade in the testimony from a particular buyer, the testimony is away from the majority opinion. Thus it can be considered as a possible unfair rating.

Entropy, a measure of *uncertainty* contained in information [9], is employed as the basis of the rating quality metric. The entropy of a variable V can be calculated as: $H(V) = -\sum Pr(v)\log(Pr(v))$, where v is a possible value of variable V , and $Pr(v)$ is the corresponding probability of V taking the value v .

Since rating in *Bayesian* rating system is basically binary, it can be seen as a discrete variable taking two possible values. Consequently, uncertainty (of seller's behavior in future transactions) observed in buyer B 's rating can be measured as: $H(r_B) = -Pr_p\log(Pr_p) - Pr_n\log(Pr_n)$. Here Pr_p and Pr_n denote the probabilities of positive ratings and negative ratings observed in the window of past W transactions, which are given by:

$$Pr_p = \frac{\alpha}{\alpha + \beta}, \quad Pr_n = \frac{\beta}{\alpha + \beta}. \quad (6)$$

Here α and β share the same meanings as in Eq. (4).

The maximum uncertainty $H_{max}(r_B)$ occurs when there are identical probabilities of positive and negative rating in the past W transactions [9]. In this case $H_{max}(r_B) = 1$. Minimum uncertainty $H_{min}(r_B)$ appears when only positive (or negative) ratings are observed in all the past W transactions.

Now, we can measure the **Quality** of the rating as:

$$Q(r_B) = 1 - \frac{H(r_B) - H_{min}(r_B)}{H_{max}(r_B) - H_{min}(r_B)} = \frac{H_{max}(r_B) - H(r_B)}{H_{max}(r_B) - H_{min}(r_B)} \quad (7)$$

Quality of testimonies from other buyers can be measured likewise. Then buyer B aggregates testimony from X if:

$$|Q(r_X) - Q(r)| \leq \varepsilon$$

where $Q(r_X)$ is the quality of the testimony reported by X , $Q(r)$ is the quality of buyer B 's current aggregated rating ($Q(r) = Q(r_B)$ initially). ε is a screening threshold (usually $\varepsilon \in [0, 1]$), which controls the sensitivity to the presence of unfair ratings. With a larger ε , the screening is less sensitive to unfair testimonies, whereas with a smaller ε , the screening is more sensitive. Both cases may lead to divergent prediction. A balanced selection of ε is necessary to make the screening work effectively. Experimental results show that $\varepsilon \in [0.35, 0.45]$ generally shows a good balance (See Section 5.4).

The proposed screening method can be outlined as Algorithm 1:

Algorithm 1. Entropy-based rating screening algorithm

B denotes the buyer initiating the testimony aggregation

C denotes the set of buyers whose testimonies are requested

X denotes a particular buyer in the set C

```

1: measure Quality of buyer  $B$ 's local rating  $Q(r_B)$  using Eq. (7)
2:  $Q(r) = Q(r_B)$ 
3: for all  $X$  in  $C$  do
4:   measure the Quality of the testimony  $Q(r_X)$  reported by  $X$  using Eq. (7)
5:   if  $|Q(r_X) - Q(r)| \leq \varepsilon$  then
6:     aggregate  $X$ 's testimony by updating  $\alpha$  and  $\beta$  using Eq. (5) accordingly, and then update
       the quality of the aggregated rating  $Q(r)$ 
7:   else
8:     discard  $X$ 's testimony
9:   end if
10: end for

```

5 Experimental Results

5.1 Setup

We simulate a trading community, in which there is one seller² and 100 buyers. There are a total of 1000 transactions in each round of simulation. The seller's behavior is mainly controlled by its loyalty, which denotes its willingness to provide "good" services. The seller may change its loyalty from one transaction to another due to many reasons, e.g. the fluctuation of the profit by providing services. In each round of simulation, seller's initial loyalty is set to be 0.9. We simulate three styles of changes of the seller's loyalty in the course of each simulation: increases and decreases from the one in previous transaction, and remains same as the one in previous transaction. The ratios of different styles of changes are chosen to be 1/3 respectively. The window size of the past transactions is set

² The goal of the experiment is to investigate whether *Bayesian* rating system can predict seller's behaviors truly even with presence of unfair ratings. *Bayesian* rating system is a distributed rating system [2], in which each seller (and buyer) is treated equally. One seller is sufficient to meet our goal.

to 50 (i.e. $W=50$ in Eq. (2)). And the forgetting factor is set to 0.9 (i.e. $\lambda = 0.9$ in Eq. (2)). P_{unfair} is set to 70%.

Before each transaction, buyer will predict the probability that the seller will provide “good” service based on other buyers’ testimonies and its local rating (if any). As the unfair ratings usually lead to divergent prediction of the seller’s behaviors, we can measure the effectiveness of the proposed method by measuring how close it is between the predicted probability and the seller’s loyalty for each transaction. We can measure the “closeness” as the *Mean Squared Error* (MSE) between the predicted probabilities and seller’s loyalties averaged over all the 1000 transactions. Ideally $MSE = 0$, which means the predicted probabilities are always equal to the seller’s loyalties in all transactions.

5.2 Types of Unfair Ratings

There are three types of unfair ratings studied in the experiments [8,10]:

- **Ballot-stuffing.** A buyer, with a probability P_{unfair} , reports that seller provides “good” service regardless of its real experience.
- **Badmouthing.** A buyer, with a probability P_{unfair} , reports that seller provides “bad” service regardless of its real experience.
- **Complementary.** A more general type of unfair rating is the Complementary unfair rating. That is, a buyer, with a probability P_{unfair} , reports a rating opposite to the real experience.

Before proceeding, we demonstrate the influence of unfair ratings first. Fig. 1(a) shows the seller’s loyalties and predicted probabilities by one buyer³ over 1000 transactions. It can be seen that *Bayesian* rating system predicts the seller’s behaviors quite close to the seller’s loyalties without the presence of unfair ratings. Fig. 1(b) shows seller’s loyalties and predicted probabilities over the 1000 transactions with presence of *badmouthing* unfair rating. It can be observed that the predicted probabilities now deviate from seller’s loyalties.

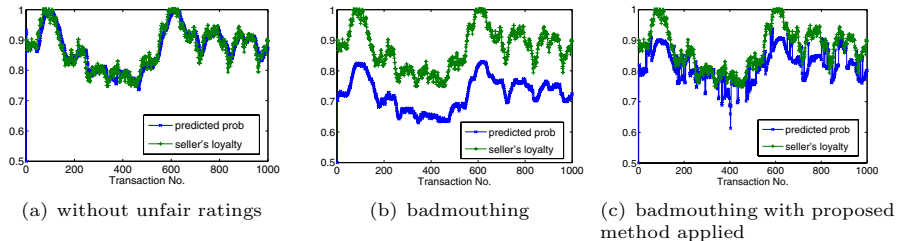


Fig. 1. Influence of unfair ratings and effectiveness of ratings screening

³ Since *Bayesian* rating system is a distributed rating system, each buyer maintains a local view of the seller’s behaviors, one buyer’s prediction is enough for studying the effectiveness of the proposed method.

Table 1. Comparison of MSEs

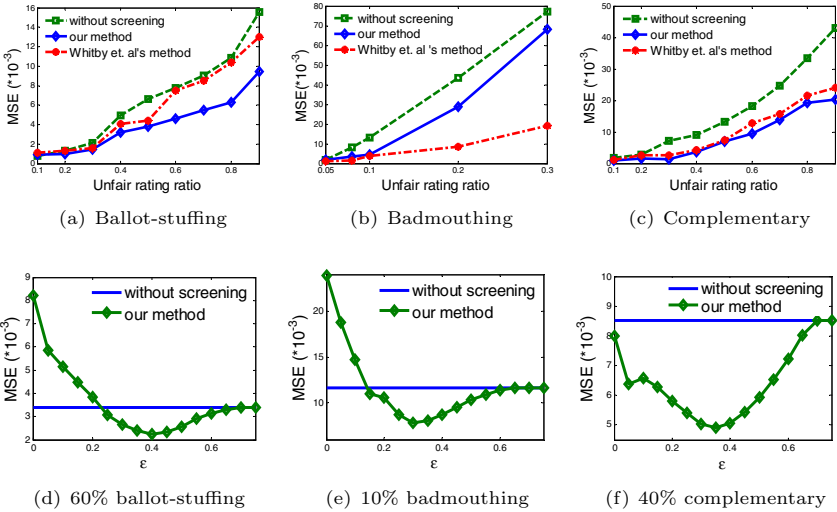
Types of unfair rating	proposed method not applied	proposed method applied
ballot-stuffing	$1.7883 * 10^{-3}$	$1.7816 * 10^{-3}$
badmouthing	$22.4487 * 10^{-3}$	$4.6240 * 10^{-3}$
complementary	$4.9041 * 10^{-3}$	$1.3683 * 10^{-3}$

5.3 Effectiveness of the Proposed Method

Fig. 1(c) shows the seller's loyalties and the predicted probabilities in the presence of *badmouthing* with the proposed screening method applied. Compared with Fig. 1(b), the predicted probabilities follow seller's loyalties more closely. MSEs in scenarios with and without the proposed screening method applied are listed in Table 1.

Compared with corresponding scenarios without proposed method applied, the proposed method manages to reduce the MSEs by 0.4%, 79.4%, and 72.1% in scenarios with the presence of ballot-stuffing, badmouthing, complementary unfair ratings respectively. The relatively lower effectiveness in the presence of ballot-stuffing is due to the reason that unfair ratings (i.e. 1 in this case) are already quite close to the seller's loyalties (around 0.9 in the simulation). The proposed method's effectiveness in mitigating the influence of the unfair ratings is thus justified.

Experiments are also conducted to study the proposed screening method's effectiveness in the presence of various ratios of unfair ratings. MSEs in different scenarios are plotted in Fig. 2(a)-2(c). With the increase of unfair rating ratio, the predicted probabilities deviate from seller's loyalties more significantly in scenarios both *with* and *without* the proposed method applied. However, with the proposed method applied, improvement over the scenarios without the pro-

**Fig. 2.** Change of MSE w.r.t. different unfair rating ratios and ϵ

posed method applied is still observable. It can be observed from Fig. 2(a)-2(c) that the proposed method mitigates the influence of the unfair ratings most effectively in scenarios with less than 60% ballot-stuffing, 10% badmouthing, and 40% complementary respectively.

We also compare the performance of the proposed screening method against the one proposed by Whitby et al. in [4]. We implement their method with q value instantiated as 0.01 since $q = 0.01$ is a good balance as reported in [4]. MSEs by applying their method are also plotted in Fig. 2(a)-2(c). It can be observed that our method outperforms Whitby et al.'s method in the presence of both ballot-stuffing and complementary unfair ratings. However, our method is not as effective as Whitby et al.'s method in the presence of badmouthing unfair ratings. This is because in our experiments, the seller's loyalties are around 0.9, which means seller would provide "good" service 90% of all the transactions. However high ratio of badmouthing might make the Pr_p and Pr_n in Eq. (6) swap their values, which makes the quality of the unfair ratings same as the honest ones. For example, the majority opinion reports that the seller provides "good" service in 8 out of 10 transactions, while an unfair rating reports the seller provides "bad" service in 8 out of 10 transactions. Qualities of both the majority opinion and the unfair rating are $-0.8\log(0.8) - 0.2\log(0.2)$. In this case, our method becomes ineffective in screening unfair ratings. However, the proposed method is much faster than Whitby et al.'s method. It takes about 0.0217 second to screen 100 testimonies for one transaction by average. In contrast, Whitby et al.'s method takes about 1.4577 seconds. Moreover, their method does not scale well with the increase in the number of testimonies due to its iterative nature, whereas the proposed method scales linearly with the increase in the number of available testimonies.

5.4 Effectiveness with Different Screening Threshold ϵ

In order to study the influence of ϵ on the proposed method's effectiveness, we choose different screening thresholds ϵ in scenarios with 60% ballot-stuffing, 10% badmouthing, and 40% complementary. MSEs between the predicted probabilities and the seller's loyalties with different ϵ are plotted in Fig. 2(d)-2(f). With a larger threshold (e.g. $\epsilon = 0.6, 0.65$), the proposed method is less sensitive to the presence of the unfair ratings. Larger MSEs are thus observed. The extreme of this case is that all testimonies are not discarded (e.g. $\epsilon \geq 0.7$), which has the same effect as without proposed method applied. On the other hand, with a smaller threshold (e.g. $\epsilon = 0.25, 0.3$), the proposed method is more sensitive to the presence of unfair ratings, more testimonies (even some honest ones) are discarded, thus the predicted probabilities depend more on the buyer's local rating and may not reflect the seller's loyalties truly. In this case, it may even make the MSEs larger than the scenarios without the proposed method applied. It can be observed from Fig. 2(d)-2(f) that $\epsilon \in [0.35, 0.45]$ generally shows a good balance.

6 Conclusions and Future Work

To the best of our knowledge, the proposed method in this paper is the first one to tackle the issue of unfair ratings from a perspective of entropy. It is

distribution-independent, and it scales linearly with the increase in the number of testimonies. Our experimental results showed that it manages to mitigate the influence of different types of unfair ratings. However, as there is no unified platform and benchmark available, a comprehensive comparison of the proposed method between other existing methods is not practical for the time being. With the planned release of “Trust Competition Testbed⁴” in July 2005, a more detailed comparison is planned as future work.

The proposed method is not effective enough in some scenarios, e.g. high loyalty with large ratios of badmouthing unfair ratings as shown by the experimental results. We plan to improve the proposed method’s performance in those scenarios in our future work. The rationale of the proposed method is that sellers provide indiscriminate services to all buyers. However, there are also cases that sellers provide “good” service to everyone except a few specific buyers that they do not “like”. In those cases, even the majority opinion might not reflect the seller’s real behavior, and the proposed method would become ineffective. Effectiveness of the proposed method in those cases is to be further investigated.

References

1. Resnick, P., Zeckhauser, R.: Trust among strangers in internet transactions: Empirical analysis of ebay’s reputation system. In Baye, M.R., ed.: *The Economics of the Internet and E-commerce*. Volume 11 of *Advances in Applied Microeconomics*. Elsevier (2002) 127–157
2. Jøsang, A., Ismail, R., Boyd, C.: A survey of trust and reputation systems for online service provision. *Decision Support System To appear* (2005)
3. Ekström, M.A., Björnsson, H.C.: A rating system for AEC e-bidding that accounts for rater credibility. In: *Proceedings of 10th Joint W055 - W065 International Symposium on Construction Innovation and Global Competitiveness (CIB W055, W065)*. (2000)
4. Whitby, A., Jøsang, A., Indulska, J.: Filtering out unfair ratings in bayesian reputation systems. In: *Proceedings of the Workshop on Trust in Agent Societies, at the 3rd International Joint Conference on Autonomous Agents and Multi Agent Systems (AAMAS2004)*. (2004)
5. Buchegger, S., Boudec, J.Y.L.: A robust reputation system for mobile ad-hoc networks. Technical Report IC/2003/50, EPFL-IC-LCA (2003)
6. Garg, A., Battiti, R., Costanzi, G.: Dynamic self-management of autonomic systems: The reputation, quality and credibility (RQC) scheme. (2004)
7. Jøsang, A., Ismail, R.: The beta reputation system. In: *Proceedings of the 15th Bled Conference on Electronic Commerce*. (2002)
8. Dellarocas, C.: Immunizing online reputation reporting systems against unfair ratings and discriminatory behavior. In: *Proceedings of the 2nd ACM conference on Electronic commerce*, ACM Press (2000) 150–157
9. Cover, T.M., Thomas, J.A.: *Elements of Information Theory*. Wisley (1991)
10. Yu, B., Singh, M.P., Sycara, K.: Developing trust in large-scale peer-to-peer systems. In: *Proceedings of 1st IEEE Symposium on Multi-Agent Security and Survivability*. (2004)

⁴ <http://www.lips.utexas.edu/kfullam/competition/>