# Rational Trust Modeling

Mehrdad Nojoumian

Department of Computer & Electrical Engineering and Computer Science
Florida Atlantic University, Boca Raton, Florida, USA
mnojoumian@fau.edu

**Abstract.** Trust models are widely used in various computer science disciplines. The primary purpose of a trust model is to continuously measure the trustworthiness of a set of entities based on their behaviors. In this article, the novel notion of *rational trust modeling* is introduced by bridging trust management and game theory. Note that trust models/reputation systems have been used in game theory (e.g., repeated games) for a long time, however, game theory has not been utilized in the process of trust model construction; this is the novelty of our approach. In our proposed setting, the designer of a trust model assumes that the players who intend to utilize the model are rational/selfish, i.e., they decide to become trustworthy or untrustworthy based on the utility that they can gain. In other words, the players are incentivized (or penalized) by the model itself to act properly. The problem of trust management can be then approached by game theoretical analyses and solution concepts such as Nash equilibrium. Although rationality might be built-in in some existing trust models, we intend to formalize the notion of rational trust modeling from the designer's perspective. This approach will result in two fascinating outcomes. First of all, the designer of a trust model can incentivize trustworthiness in the first place by incorporating proper parameters into the trust function, which can be later utilized among selfish players in strategic trust-based interactions (e.g., e-commerce scenarios). Furthermore, using a rational trust model, we can prevent many well-known attacks on trust models. These two prominent properties also help us to predict the behavior of the players in subsequent steps by game theoretical analyses.

**Keywords:** trust management, reputation system, game theory, and rationality.

## 1   Introduction

The primary purpose of a trust model is to continuously measure the trustworthiness of a set of entities (e.g., servers, sellers, agents, nodes, robots, players, etc) based on their behaviors. Indeed, scientists across various disciplines have conducted research on trust over decades and produced fascinating discoveries, however, there is not only a huge gap among findings in these research communities but also these discoveries have not been properly formalized to have a better understanding of the notion of trust, and consequently, practical computational models of trust. We therefore intend to look at the problem of trust modeling from an interdisciplinary perspective that is more realistic and closer to human comprehension of trust.

From a social science perspective, *trust* is the willingness of a person to become vulnerable to the actions of another person irrespective of the ability to control those

actions [1]. However, in the computer science community, *trust* is defined as a personal expectation that a player has with respect to the future behavior of another party, i.e., a personal quantity measured to help the players in their future dyadic encounters. On the other hand, *reputation* is the perception that a player has with respect to another player's intention, i.e., a social quantity computed based on the actions of a given player and observations made by other parties in an electronic community that consists of interacting parties such as people or businesses [2].

From another perspective [3], *trust* is made up of underlying beliefs and it is a function based on the values of these beliefs. Similarly, *reputation* is a social notion of trust. In our lives, we each maintain a set of reputation values for people we know. Furthermore, when we decide to establish an interaction with a new person, we may ask other people to provide recommendations regarding the new party. Based on the information we gather, we form an opinion about the reputation of the new person. This decentralized method of reputation measurement is called *referral chain*. Trust can be also created based on both local and/or social evidence. In the former case, trust is built through direct observations of a player whereas, in the latter case, it is built through information from other parties. It is worth mentioning that a player can gain or lose her reputation not only because of her cooperation/defection in a specific setting but also based on the ability to produce accurate referrals.

Generally speaking, the goal of reputation systems is to collect, distribute and aggregate feedback about participants' past behavior. These systems address the development of reputation by recording the behavior of the parties, e.g., in e-commerce, the model of reputation is constructed from a buying agent's positive or negative past experiences with the goal of predicting how satisfied a buying agent will be in future interactions with a selling agent. The ultimate goal is to help the players decide whom to trust and to detect dishonest or compromised parties in a system [4]. There exist many fascinating applications of trust models and reputation systems in various engineering and computer science disciplines.

In fact, trust models are widely used in scientific and engineering disciplines such as electronic commerce [5–10], computer security and rational cryptography [11–16], vehicular ad-hoc networks [17, 18], social and semantic web [19, 20], multiagent systems [21–23], robotics and autonomous systems [24, 25], game theory and economics [26, 27]. To the best of our knowledge, there is no literature on *rational trust modeling*, that is, using game theory during the construction of a trust model. Note that game theoretic models have been used for management and analyses of trust-based systems [28, 29].

## 1.1 Our Motivation and Contribution

As our motivation, we intend to provide a new mechanism for trust modeling by which:

1. The trust model incentivizes trustworthiness in the first place, i.e., *self-enforcing*.
2. The model is naturally resistant to attacks on trust models, i.e., *resistant*.

We therefore introduce the novel notion of *rational trust modeling* by bridging trust management and game theory. We would like to emphasize that trust models have been used in game theory for a long time, for instance, in repeated games to incentivize the

players to be cooperative and not to deviate from the game's protocol. However, game theory has not been utilized in the process of trust model construction; in fact, this is the novelty of our proposed approach.

In our setting, the designer of a trust model assumes that the players who intend to utilize the model are rational/selfish meaning that they cooperate to become trustworthy or defect otherwise based on the utility (to be defined by the trust model) that they can gain, which is a reasonable and standard assumption. In other words, the players are incentivized (or penalized) by the model itself to act properly. The problem of trust modeling can be then approached by strategic games among the players using utility functions and solution concepts such as Nash equilibrium.

Although rationality might be built-in in some existing trust models, we formalize the notion of rational trust modeling from the model designer's perspective. This approach results in two fascinating outcomes. First of all, the designer of a trust model can incentivize trustworthiness in the first place by incorporating proper parameters into the trust function, which can be later utilized among selfish players in strategic trust-based interactions (e.g., e-commerce scenarios between sellers and buyers). Furthermore, using a rational trust model, we can prevent many well-known attacks on trust models, as we describe later. These two prominent properties also help us to predict behavior of the players in subsequent steps by game theoretical analyses.

## 1.2 Our Approach in Nutshell

Suppose there exist two sample trust functions: The first function $f_1(\mathscr{T}_i^{p-1}, \alpha_i)$ receives the previous trust value $\mathscr{T}_i^{p-1}$ and the current action $\alpha_i$ of a seller $S_i$ (i.e., cooperation or defection) as two inputs to compute the updated trust value $\mathscr{T}_i^p$ for the next round. However, the second function $f_2(\mathscr{T}_i^{p-1}, \alpha_i, \ell_i)$ has an extra input value known as the seller's lifetime denoted by $\ell_i$. Using the second trust function, a seller with a longer lifetime will be rewarded (or penalized) more (or less) than a seller with a shorter lifetime assuming that the other two inputs (i.e., current trust value and the action) are the same. In this scenario, "reward" means gaining a higher trust value and becoming more trustworthy, and "penalty" means otherwise. In other words, if two sellers $S_i$ and $S_j$ both cooperate $\alpha_i = \alpha_j = \mathscr{C}$ and their current trust values are equal $\mathscr{T}_i^{p-1} = \mathscr{T}_j^{p-1}$ but their lifetime parameters are different, for instance, $\ell_i > \ell_j$, the seller with a higher lifetime parameter, gains a higher trust value for the next round, i.e., $\mathscr{T}_i^p > \mathscr{T}_j^p$. This may help $S_i$ to sell more items and accumulate more revenue because buyers always prefer to buy from trustworthy sellers, i.e., sellers with a higher trust value.

Now consider a situation in which the sellers can sell defective versions of an item with more revenue or non-defective versions of the same item with less revenue. If we utilize the first sample trust function $f_1$, it might be tempting for a seller to sell defective items because he can gain more utility. Furthermore, the seller can return to the community with a new identity (a.k.a, re-entry attack) after selling defective items and accumulating a large revenue. However, if we use the second sample trust function $f_2$, it's no longer in a seller's best interest to sell defective items because if he returns to the community with a new identity, his lifetime parameter becomes zero and he loses all the credits that he has accumulated overtime. As a result, he loses his future

customers and a huge potential revenue, i.e., buyers may prefer a seller with a longer lifetime over a seller who is a newcomer. The second trust function not only incentivizes trustworthiness but also prevents the re-entry attack.

Note that this is just an example of rational trust modeling for the sake of clarification. The second sample function here utilizes an extra parameter $\ell_i$ in order to incentivize trustworthiness and prevent the re-entry attack. In fact, different parameters can be incorporated into trust functions based on the context (whether it's a scenario in e-commerce or cybersecurity and so on), and consequently, different attacks can be prevented, as discussed in Section 4.

## 2 Preliminaries: Game Theory

In this section, preliminary materials regarding game-theoretic concepts are provided for further technical discussions.

A *game* consists of a set of *players*, a set of *actions* and *strategies* (i.e., the method of selecting actions), and finally, a *utility function* that is used by each player to compute how much benefit he can gain by choosing a certain action. In *cooperative games*, players collaborate and split the aggregated utility among themselves, that is, cooperation is incentivized by agreement. However, in *non-cooperative games*, players can not form agreements to coordinate their behavior, that is, any cooperation must be self-enforcing.

The *prisoner's dilemma*, as illustrated in Figure 1, is an example of non-cooperative games. In this setting, two possible actions are considered: $\mathscr{C}$: *keep quiet* (or cooperation) and $\mathscr{D}$: *confess* (or defection). In the payoff (utility) matrix, $+1, 0, -1$, and $-2$ denote freedom, jail for one year, jail for two years, and jail for three years respectively. The outcome of this game will be $(\mathscr{D}, \mathscr{D})$ because of the *Nash equilibrium* concept, while the ideal outcome is $(\mathscr{C}, \mathscr{C})$. To better understand the notion of Nash equilibrium, and consequently, why the game has such an outcome, consider the following two possible scenarios:

1. If player $P_1$ selects $\mathscr{C}$ (the first row), then player $P_2$ will select $\mathscr{D}$ (the second column) since $+1 > 0$.

2. If player $P_1$ selects $\mathscr{D}$ (the second row), then player $P_2$ will select $\mathscr{D}$ (the second column) since $-1 > -2$.

In other words, regardless of whether player $P_1$ cooperates or defects, player $P_2$ will always defect. Since the payoff matrix is symmetric, player $P_1$ will also defect regardless of whether $P_2$ cooperates or defects. In fact, since the players are in two different locations and are not able to coordinate their behavior, the final outcome will be $(\mathscr{D}, \mathscr{D})$.

We now briefly review some well-known game-theoretic concepts and definitions [30] for our further analysis and discussions.

**Definition 1.** *Let* $\mathscr{A} \stackrel{\text{def}}{=} \mathscr{A}_1 \times \cdots \times \mathscr{A}_n$ *be an action profile for n players, where* $\mathscr{A}_i$ *denote the set of possible actions of player* $P_i$. *A game* $\Gamma = (\mathscr{A}_i, u_i)$ *for* $1 \leq i \leq n$, *consists of* $\mathscr{A}_i$ *and a utility function* $u_i : \mathscr{A} \mapsto \mathbb{R}$ *for each player* $P_i$. *We refer to a vector of actions* $a = (a_1, \ldots, a_n) \in \mathscr{A}$ *as an* outcome *of the game.*

|     | **P₂** | |
| --- | --- | --- |
|     | **C:** Quiet | **D:** Confess |
| **C:** Quiet | 0 , 0 | -2 , +1 |
| **D:** Confess | +1 , -2 | -1 , -1 |

(with $P_1$ labeling the rows)

**P₂**

| 0 , 0 | -2 , +1 |
| --- | --- |
| +1 , -2 | -1 , -1 |

$P_1$: "what if I **cooperate**"

$P_1$: "what if I **defect**"

| 0 , 0 | -2 , +1 |
| --- | --- |
| +1 , -2 | -1 , -1 |

**P₂**

**Fig. 1.** Nash Equilibrium in Prisoner's Dilemma

**Definition 2.** *The* utility function $u_i$ *illustrates the preferences of player $P_i$ over different outcomes. We say $P_i$ prefers* outcome a to a' iff $u_i(a) > u_i(a')$, *and he* weakly prefers *outcome a to a' if $u_i(a) \geq u_i(a')$.*

To allow the players to follow randomized strategies (strategy defines how to select actions), we define $\sigma_i$ as a probability distribution over $\mathscr{A}_i$ for a player $P_i$. This means that he samples $a_i \in \mathscr{A}_i$ according to $\sigma_i$. A strategy is said to be a *pure-strategy* if each $\sigma_i$ assigns probability 1 to a certain action, otherwise, it is said to be a *mixed-strategy*. Let $\sigma = (\sigma_1, \ldots, \sigma_n)$ be the vector of players' strategies, and let $(\sigma'_i, \sigma_{-i}) \stackrel{\text{def}}{=} (\sigma_1, \ldots, \sigma_{i-1}, \sigma'_i, \sigma_{i+1}, \ldots, \sigma_n)$, where $P_i$ replaces $\sigma_i$ by $\sigma'_i$ and all the other players' strategies remain unchanged. Therefore, $u_i(\sigma)$ denote the expected utility of $P_i$ under the strategy vector $\sigma$. A player's goal is to maximize $u_i(\sigma)$. In the following definitions, one can substitute an action $a_i \in \mathscr{A}_i$ with its probability distribution $\sigma_i \in S_i$ or vice versa.

**Definition 3.** *A vector of strategies $\sigma$ is* Nash equilibrium *if, for all i and any $\sigma'_i \neq \sigma_i$, it holds that $u_i(\sigma'_i, \sigma_{-i}) \leq u_i(\sigma)$. This means no one gains any advantage by deviating from the protocol as long as the others follow the protocol.*

**Definition 4.** *Let $S_{-i} \stackrel{\text{def}}{=} S_1 \times \cdots \times S_{i-1} \times S_{i+1} \times \cdots \times S_n$. A strategy $\sigma_i \in S_i$ (or an action) is* weakly dominated *by a strategy $\sigma'_i \in S_i$ (or another action) with respect to $S_{-i}$ if:*

1. *For all $\sigma_{-i} \in S_{-i}$, it holds that $u_i(\sigma_i, \sigma_{-i}) \leq u_i(\sigma'_i, \sigma_{-i})$.*
2. *There exists a $\sigma_{-i} \in S_{-i}$ s.t. $u_i(\sigma_i, \sigma_{-i}) < u_i(\sigma'_i, \sigma_{-i})$.*

*This means that $P_i$ can never improve its utility by playing $\sigma_i$, and he can sometimes improve it by not playing $\sigma_i$. A strategy $\sigma_i \in S_i$ is* strictly dominated *if player $P_i$ can always improve its utility by not playing $\sigma_i$.*
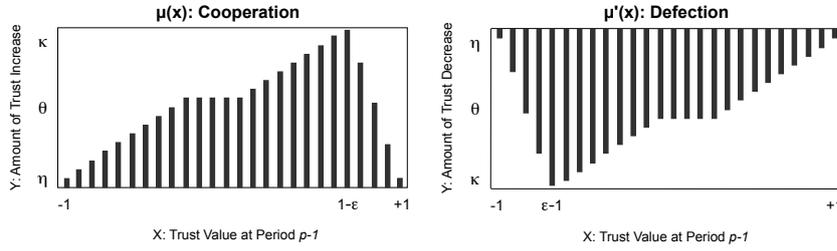
## 3 Rational Trust Modeling

We stress that our goal here is not to design specific trust models or construct certain utility functions. Our main objective is to illustrate the high-level idea of *rational trust modeling* through examples/analyses without loss of generality.

### 3.1 Trust Modeling: Construction and Evaluation

To construct a quantifiable model of trust, a mathematical function or model for trust measurement in a community of *n* players must be designed. First of all, a basic trust function is defined as follows:

**Definition 5.** *Let $\mathcal{T}_i^p$ denote trust value of player $P_i$ in period $p$ where $-1 \leq \mathcal{T}_i^p \leq +1$ and $\mathcal{T}_i^0 = 0$ for newcomers. A trust function is a mapping from $\mathbb{R} \times \mathbb{N}$ to $\mathbb{R}$ which is defined as follows: $(\mathcal{T}_i^{p-1}, \alpha_i) \mapsto \mathcal{T}_i^p$, where $\mathcal{T}_i^{p-1}$ denote the trust value of player $P_i$ in period $p-1$ and $\alpha_i \in \{0,1\}$ denote whether $P_i$ has cooperated, i.e., $\alpha_i = 1$, or defected, i.e., $\alpha_i = 0$, in period p.*

As an example, we can refer to the following mathematical model [9, 31]. In this model, $\mathcal{T}_i^p = \mathcal{T}_i^{p-1} + \mu(x)$ or $\mathcal{T}_i^p = \mathcal{T}_i^{p-1} - \mu'(x)$ for $\alpha_i = 1$ or $\alpha_i = 0$ respectively, shown in Figure 2. Parameters $\eta$, $\theta$ and $\kappa$ are used to reward or penalize players based on their actions (for instance, as defined in [31], $\eta = 0.01$, $\theta = 0.05$ and $\kappa = 0.09$). Note that in $[1-\varepsilon, +1]$ and $[-1, \varepsilon - 1]$, $\mu(x)$ and $\mu'(x)$ both converge to zero, as required by Definition 5, i.e., $-1 \leq \mathcal{T}_i^p \leq +1$.



**Fig. 2.** Trust Adjustment by $\mu(x)$ and $\mu'(x)$

After designing a mathematical function, it must be assessed and validated from different perspectives for further improvement. We provide high-level descriptions of some validation procedures to be considered for evaluation of a trust model, that is, *behavioral*, *adversarial* and *operational* methodologies.

1. *Behavioral*: how the model performs among a sufficient number of players by running a number of standard tests, i.e., executing a sequence of "cooperation" and "defection" (or no-participation) for each player. For instance, how fast the model can detect defective behavior by creating a reasonable trust margin between cooperative and non-cooperative parties.

2. *Adversarial*: how vulnerable the trust model is to different attacks or any kinds of corruption by a player or a coalition of malicious parties. Seven well-known attacks on trust models are listed below. The first five attacks are known as *single-agent attacks* and the last two are known as *multi-agent* or *coalition attacks* [32].

   (a) **Sybil:** forging identities or creating multiple false accounts by one player.

   (b) **Lag:** cooperating for some time to gain a high trust value and then cheat.

   (c) **Re-Entry:** corrupted players return to the scheme using new identities.

   (d) **Imbalance:** cooperating on cheap transactions; defecting on expensive ones.

   (e) **Multi-Tactic:** any combination of attacks mentioned above.

   (f) **Ballot-Stuffing:** fake transactions among colluders to gain a high trust value.

   (g) **Bad-Mouthing:** submitting negative reviews to non-coalition members.

3. *Operational*: how well the future states of trust can be predicted with a relatively accurate approximation based on possible action(s) of the players (prediction can help us to prevent some well-known attacks), and how well the model can incentivize cooperation in the first place.

In the next section, we clarify what considerations should be taken into account by the designer in order to construct a proper trust model that resists against various attacks and also encourages trustworthiness in the first place.

### 3.2   Rational Trust Modeling Illustration: Seller's Dilemma

We now illustrate a dilemma between two sellers by considering two different trust functions. In this setting, each seller has defective and non-defective versions of an item for sale. We consider the following two possible actions:

1. *Cooperation*: selling the non-defective version of the item for \$3 to different buyers.
2. *Defection*: selling the defective version of the item for \$2 to different buyers.

    Assuming that the buyers are not aware of the existence of the defective version of the item, they may prefer to buy from the seller who offers the lowest price. This is a pretty natural and standard assumption. As a result, the seller who offers the lowest price has the highest chance to sell the item, and consequently, he can gain more utility.

    An appropriate payoff function can be designed for this seller's dilemma based on the probability of being selected by a buyer since there is a correlation between the offered price and this probability, Figure 3. In other words, if they both offer the same price (\$2 or \$3), they have an equal chance of being selected by a buyer, otherwise, the seller who offers a lower price (\$2) will be selected by the probability of 1.
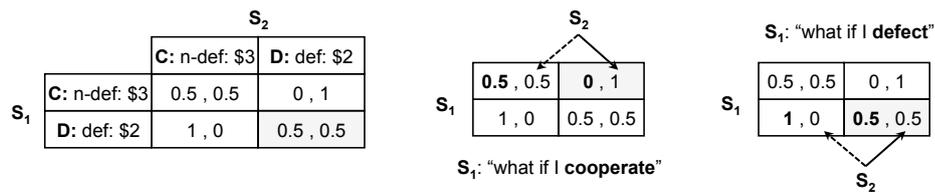


**Fig. 3.** Seller's Dilemma

    Similar to the prisoner's dilemma, defection is Nash equilibrium meaning that it is in the best interest of each seller to maximize his utility by selling the defective version of the item. For instance, suppose $S_1$ cooperates by selling the non-defective item for

$3, $S_2$ will then offer the defective item for $2 to have the highest chance to sell the item, and consequently, he can gain more utility. On the other hand, suppose $S_1$ defects by selling the defective item for $2, $S_2$ will then offer the defective item for $2 to compete with $S_1$. That is, regardless of whether seller $S_1$ cooperates or defects, seller $S_2$ will always defect, and since the payoff matrix is symmetric, defection (selling the defective item) is always Nash equilibrium.

Without loss of generality, we now show how a proper trust model can deal with this dilemma; note that this is just an example for the sake of clarification. We first consider two different trust functions, as we described earlier:

1. The first function $f_1 : (\mathscr{T}_i^{p-1}, \alpha_i) \mapsto \mathscr{T}_i^p$, where $\mathscr{T}_i^{p-1}$ denote the trust value of seller $S_i$ in period $p-1$ and $\alpha_i \in \{0,1\}$ denote whether seller $S_i$ has cooperated or defected in the current period $p$.

2. The second function $f_2 : (\mathscr{T}_i^{p-1}, \alpha_i, \ell_i) \mapsto \mathscr{T}_i^p$, where $\mathscr{T}_i^{p-1}$ and $\alpha_i \in \{0,1\}$ denote the same notions as of the previous function and $\ell_i \geq 0$ denote the lifetime of seller $S_i$ as a new input in the trust function. This parameter defines how long a seller with a reasonable number of transactions has been in the market.

For the sake of simplicity, we didn't consider two different parameters for the lifetime and the number of transactions, however, two separate parameters can be simply incorporated into our trust function and we can still achieve the same game theoretical result. The main reason is because we want to make sure the sellers who have been in the market for a long time but have been inactive or have had a limited number of transactions cannot obtain a high trust value, which is a reasonable assumption.

Now considering the seller's dilemma that we illustrated in Figure 3, the first function $f_1$ is significantly vulnerable to *re-entry* attack. That is, a seller $S_i$ may defect on a sequence of transactions in the middle of his lifetime to gain substantial revenues (utility). He can then return to the market with a new identity as a newcomer.

However, the lifetime $\ell_i$ is part of the second trust function $f_2$ meaning that a seller $S_i$ with a longer lifetime is more reliable/trustworthy from the buyers' perspective. As a result, he has a higher chance to be selected by the buyers, and consequently, he can gain more utility. This is a very realistic assumption in the e-marketplace. Therefore, it's not in the best interest of a seller to sacrifice his lifetime indicator (and correspondingly his trustworthiness) for a short term utility through defection and then re-entry attack.

It is not hard to show that, by using function $f_2$ rather than function $f_1$, "defection" is no longer Nash equilibrium in the seller's dilemma, as we illustrate in Section 3.3. When we assume the sellers are rational/selfish and they decide based on their utility functions, we can then design a proper trust function similar to $f_2$ to incentivize cooperation in the first place. Furthermore, we can deal with a wide range of attacks, as we mentioned earlier. Finally, at any point, the behavior of a seller can be predicted by estimation of his payoff through trust and utility functions.

### 3.3 Rational Trust Modeling: Design and Analysis

In our setting, the utility function $u_i : \mathscr{A} \times \mathscr{T}_i \mapsto \mathbb{R}$, which depends on the seller's action and his trust value. This function computes the utility that each $S_i$ gains or loses by

selecting a certain action. If we consider the 2nd trust function $f_2 : (\mathscr{T}_i^{p-1}, \alpha_i, \ell_i) \mapsto \mathscr{T}_i^p$, the trust value then depends on the seller's lifetime $\ell_i$ as well. As a result, the lifetime of the seller directly affects the utility that the seller can gain or lose. Now consider the following simple utility function:

$$u_i = \Omega \times g(\mathscr{T}_i^p) \quad \text{where } 0 \leq g(\mathscr{T}_i^p) \leq 1, \Omega \text{ is a constant} \tag{1}$$

As stated earlier, we first define the following parameters, where $-1 \leq \mathscr{T}_i^p \leq +1$ and $\alpha_i \in \{0, 1\}$ denote whether $S_i$ has cooperated or defected in the previous period:

$$\tau_i = \mathscr{T}_i^p - \mathscr{T}_i^{p-1} \quad \text{where} \quad \frac{|\tau_i|}{\tau_i} = \begin{cases} +1 & \text{if } \alpha_i = 1 \\ -1 & \text{if } \alpha_i = 0 \end{cases} \tag{2}$$

In the following equations, the first function $f_1$ does not depend on the seller's lifetime $\ell_i$, however, the second function $f_2$ has an extra factor that is defined by lifetime $\ell_i$ and constants $\rho$. We can assume that $\rho \ell_i$ is in the same range as of $\mu$ depending on the player's lifetime; that is why $\ell_i$ is multiplied by multiplicative factor $\rho$. Also, it's always positive meaning that no matter if a player cooperates or defects, he will always be rewarded by $\rho \ell_i$. We stress that parameter $\ell_i$ in function $f_2$ is just an examples of how a rational trust function can be designed. The designer can simply consider various parameters (that denote different concepts) as additive or multiplicative factors based on the context in which the trust model is supposed to be utilized. We discuss this issue later in Section 4 in detail.

$$f_1 \ : \ \mathscr{T}_i^p = \mathscr{T}_i^{p-1} + \frac{|\tau_i|}{\tau_i} \mu \tag{3}$$

$$f_2 \ : \ \mathscr{T}_i^p = \mathscr{T}_i^{p-1} + \frac{|\tau_i|}{\tau_i} \mu + \rho \ell_i \tag{4}$$

$-1 \leq \mathscr{T}_i^p \leq +1$, E.g.: $0 \leq \mu < 0.1$ is a unified function in $f_1$ and $f_2$

The first function $f_1$ rewards or penalizes the sellers based on their actions and independent of their lifetimes. This makes function $f_1$ vulnerable to different attacks such as the re-entry attack because a malicious seller can always come back to the scheme with a new identity, and then, starts re-building his reputation for another malicious activity. It is possible to make the sign-up procedure costly but it is out of the scope of this paper.

On the other hand, the second trust function $f_2$ has an extra term that is defined by the seller's lifetime $\ell_i$. This term will be adjusted by $\rho$ as an additional reward or punishment factor in the trust function. In other words, the seller's current lifetime $\ell_i$ in addition to a constant $\pm \beta$ (in the case of cooperation/defection) determine the extra reward/punishment factor. As a result, it is not in the best interest of a seller to reset his lifetime indicator $\ell_i$ to zero because of a short-term utility. This lifetime indicator can increase the seller's trustworthiness, and consequently, his long-term utility overtime.

Let assume our sample utility function is further extended as follows, where $\Omega$ is a constant, for instance, $\Omega$ can be \$100:

$$u_i = \Omega \left( \frac{\mathscr{T}_i^p + 1}{2} \right) \quad \text{where } 0 \leq \frac{\mathscr{T}_i^p + 1}{2} \leq 1, -1 \leq \mathscr{T}_i^p \leq +1 \tag{5}$$

The utility function simply indicates a seller with a higher trust value (which depends on his lifetime indicator as well) can gain more utility because he has a higher chance to be selected by the buyers. In other words, Eqn. (5) maps the current trust value $\mathscr{T}_i^p$ to a value between zero and one, which can be also interpreted as the probability of being selected by the buyers. For the sake of simplicity, suppose $\mathscr{T}_i^{p-1}$ is canceled out in both $f_1$ and $f_2$ as a common factor. The overall utility $U_i^{f_1}$ is shown below when $f_1$ is used. Note that $u_i$ computes the utility of a seller in the case of cooperation or defection whereas $U_i$ also takes into account the *external utility* or *future loss* that a seller may gain or lose. For instance, more savings through selling the defective version of an item instead of its non-defective version.

$$U_i^{f_1} = \Omega \times \begin{cases} \frac{+\mu+1}{2} & \text{using } f_1 \text{ when } \alpha_i = 1 \\[2ex] \frac{-\mu+1}{2} + \beta & \text{using } f_1 \text{ when } \alpha_i = 0 \text{ plus } \beta, \beta \text{ is the } \textit{external utility} \\ & \text{that the seller obtains by selling the defective item} \end{cases}$$

As shown in $U_i^{f_1}$, function $f_1$ rewards/penalizes sellers in each period by factor $\pm\frac{\mu}{2}$. Accordingly, we can assume *external utility* $\beta$ that the seller obtains by selling the defective item is slightly more than (as much as $\sigma$) the utility that the seller may lose because of defection; otherwise, the seller wouldn't defect, that is, $\beta = \frac{\mu}{2} + |-\frac{\mu}{2}| + \sigma = \mu + \sigma$ (note that the seller not only loses a potential reward $\frac{\mu}{2}$ but also he is penalized by factor $-\frac{\mu}{2}$ when he defects.) In other words, external utility $\beta$ not only compensates for loss $\frac{\mu}{2} + |-\frac{\mu}{2}|$ but also provides additional gain $\sigma$.

As a result, $\frac{-\mu+1}{2} + \beta = \frac{-\mu+1}{2} + (\mu + \sigma) = \frac{\mu+1}{2} + \sigma$. Therefore, $\mathscr{D}$efection is always Nash Equilibrium when $f_1$ is used, as shown in Table 1. We can assume the seller cheats on $\delta$ rounds until he is labeled as an untrustworthy seller. At this point, he leaves and returns with a new identity with the same initial trust value of newcomers, i.e., re-entry attack. Our analysis remains the same even if cheating is repeated for $\delta$ rounds.

| $S_1$ \ $S_2$ | $\mathscr{C}$ooperation | $\mathscr{D}$efection |
|---|---|---|
| $\mathscr{C}$ooperation | $\frac{\mu+1}{2}, \frac{\mu+1}{2}$ | $\frac{\mu+1}{2}, \frac{\mu+1}{2} + \sigma$ |
| $\mathscr{D}$efection | $\frac{\mu+1}{2} + \sigma, +\frac{\mu+1}{2}$ | $\frac{\mu+1}{2} + \sigma, \frac{\mu+1}{2} + \sigma$ |

**Table 1.** Seller's Dilemma: $\mathscr{D}$efection is always Nash Equilibrium using $f_1$.

Similarly, function $f_2$ rewards/penalizes sellers through $U_i^{f_2}$ in each period by factor $\pm\frac{\mu}{2}$. Furthermore, this function also has a positive reward (or forgiveness) factor $\frac{\rho\ell_i}{2}$ for cooperative (or non-cooperative) sellers, which is defined by their lifetime factors. Likewise, we can assume *external utility* $\beta$ that the seller obtains by selling the defective item is slightly more than the utility that the seller may lose by defection ($\beta = \mu + \sigma$).

The overall utility $U_i^{f2}$ will be as follows when the $f_2$ is used:

$$U_i^{f2} = \Omega \times \begin{cases} \frac{(+\mu + \rho \ell_i) + 1}{2} & \text{using } f_1 \text{ when } \alpha_i = 1 \\[2em] \frac{(-\mu + \rho \ell_i) + 1}{2} + \beta - \gamma & \text{using } f_1 \text{ when } \alpha_i = 0 \text{ plus } \beta \text{ as before,} \\ & \text{where } \gamma \text{ is the } \textit{future loss} \text{ due to the impact of } \ell_i \end{cases}$$

Without loss of generality, suppose the seller defects, leaves and then comes back with a new identity. As a result the lifetime index $\ell_i$ becomes zero. Let assume this index is increased by the following arithmetic progression to reach to where it was: $0$, $\frac{1}{5}\ell_i$, $\frac{2}{5}\ell_i$, $\frac{3}{5}\ell_i$, $\frac{4}{5}\ell_i$, $\ell_i$. In reality, it takes a while for a seller to accumulate this credit based on our definition, i.e., *years of existence* and *number of transactions*. Therefore,

$$\gamma \approx \frac{\rho}{2}\left((\ell_i - 0) + (\ell_i - \frac{1}{5}\ell_i) + (\ell_i - \frac{2}{5}\ell_i) + (\ell_i - \frac{3}{5}\ell_i) + (\ell_i - \frac{4}{5}\ell_i) + (\ell_i - \ell_i)\right)$$
$$= \frac{\rho}{2}\left(\ell_i + \frac{4}{5}\ell_i + \frac{3}{5}\ell_i + \frac{2}{5}\ell_i + \frac{1}{5}\ell_i + 0\right) = \frac{3}{2}\rho \ell_i$$

E.g., $(\ell_i - \frac{1}{5}\ell_i)$ denote the *lifetime* could be $\ell_i$, or even more, but it's now $\frac{1}{5}\ell_i$ meaning that the seller is losing $\frac{4}{5}\ell_i$, and so on. We now simplify the $U_i^{f2}$ when $\alpha_i = 0$ as follows:

$$U_i^{f2} \; : \; \frac{(-\mu + \rho \ell_i) + 1}{2} + \beta - \gamma$$
$$= \frac{(-\mu + \rho \ell_i) + 1}{2} + \mu + \sigma - \frac{3}{2}\rho \ell_i = \overbrace{\frac{(+\mu + \rho \ell_i) + 1}{2}}^{\Psi} + \sigma - \frac{3}{2}\rho \ell_i$$

This is a simple but interesting result that shows, as long as $\frac{3}{2}\rho \ell_i > \sigma$, $\mathscr{C}$ooperation is always Nash Equilibrium when $f_2$ is used, Table 2. In other words, as long as future loss $\gamma$ is greater than the short-term gain through defection, it's not in the best interest of the seller to cheat and commit to the re-entry attack, that is, the seller may gain a small *short-term* utility by cheating, however, he loses a larger *long-term* utility because it takes a while to reach to $\ell_i$ from 0. The analysis will be the same if the seller cheats on $\delta$ rounds before committing to the re-entry attack as long as the future loss is greater than the short-term gain. In fact, the role of parameter $\ell_i$ is to make the future loss costly.

| $S_1$ \ $S_2$ | $\mathscr{C}$ooperation | $\mathscr{D}$efection |
|---|---|---|
| $\mathscr{C}$ooperation | $\Psi$ , $\Psi$ | $\Psi$ , $\Psi + \sigma - \frac{3}{2}\rho \ell_i$ |
| $\mathscr{D}$efection | $\Psi + \sigma - \frac{3}{2}\rho \ell_i$ , $\Psi$ | $\Psi + \sigma - \frac{3}{2}\rho \ell_i$ , $\Psi + \sigma - \frac{3}{2}\rho \ell_i$ |

**Table 2.** Seller's Dilemma: $\mathscr{C}$ooperation is always Nash Equilibrium using $f_2$ when $\frac{3}{2}\rho \ell_i > \sigma$.

## 4   Technical Analysis and Discussion

As stated earlier, we would like to emphasize that our intention here was not to design specific trust models, construct utility functions (which is hard in many cases), target a certain set of attacks, or focus on particular assumptions/games/dilemmas. Our main objective was to illustrate the high-level idea of *rational trust modeling* by some examples and analyses without loss of generality. The presented models, functions, dilemma scenarios, attack strategies, assumptions and parameters can be modified as long as the model designers utilize the technical approach and strategy of rational trust modeling.

As we illustrated, by designing a proper trust function and using a game-theoretical analysis, not only *trustworthiness can be incentivized* but also well-known *attacks on trust functions can be prevented*, such as re-entry attack in our example. Furthermore, *behavior of the players can be predicted* by estimating the utility that each player may gain. In this section, we further discuss on these issues while focusing on other types of attacks against trust models. As shown in Table 3, all single-agent attacks can be simply prevented if the designer of the model incorporates one or more extra parameters (in addition to the previous trust value $\mathscr{T}_i^{p-1}$ and the current action $\alpha_i$) into the function.

| Attacks | Parameter | Description |
|---------|-----------|-------------|
| **Sybil** | Total number of Past Transactions | Prevent the players to create multiple false accounts |
| **Lag** | High Expectancy | Prevent the players to cheat after gaining a high trust value |
| **Re-Entry** | Lifetime of the Player | Prevent the players to return with a new identity |
| **Imbalance** | Transaction Cost | Prevent the players to cheat on expensive transactions |
| **Multi Tactic** | Combination of Parameters | Prevent the players to defect in various circumstances |

**Table 3.** Sample Parameters: To deal with single-agent attacks during rational trust modeling.

For instance, to deal with the Sybil attack, we can consider a parameter that only reflects the total number of past transactions. In that case, it's not in the best interest of a player to create multiple accounts and divides his total number of transactions among different identities. For imbalance attack, we can consider a parameter for transaction cost, i.e., if the player defects on an expensive transaction, his trust value declines with much faster ratio. For other attacks and their corresponding parameters, see Table 3.

It is worth mentioning that when the trust value reaches to the saturated region, e.g., very close to $+1$, a player may not have any interest to accumulate more trust credits. However, in this situation, the *high expectancy* parameter (as shown in Table 3) can be simply utilized in the trust function to warn the fully trusted players that they can sustain this credibility as long as they remain reliable, and if they commit to defections, they will be negatively and significantly (more than others) affected due to high expectancy.

Similarly, we can consider more complicated parameters to incentivize the players not to collude, and consequently, deal with multi-agent/coalition attacks. It is also worth mentioning that consideration should be given to the context in which the trust model is supposed to be used. Some of these parameters are context-oriented and the designer of the model should take this fact into account when designing a rational trust function.

## 5   Concluding Remarks

In this paper, the novel notion of *rational trust modeling* was introduced by bridging trust management and game theory. In our proposed setting, the designer of a trust model assumes that the players who intend to utilize the model are rational/selfish, i.e., they decide to become trustworthy or untrustworthy based on the utility that they can gain. In other words, the players are incentivized (or penalized) by the model itself to act properly. The problem of trust management can be then approached by strategic games among the players using utility functions and solution concepts such as NE.

Our approach resulted in two fascinating outcomes. First of all, the designer of a trust model can incentivize trustworthiness in the first place by incorporating proper parameters into the trust function. Furthermore, using a rational trust model, we can prevent many well-known attacks on trust models. These prominent properties also help us to predict the behavior of the players in subsequent steps by game theoretical analyses. As our final remark, we would like to emphasize that our rational trust modeling approach can be extended to any mathematical modeling where some sorts of utility and/or rationality are involved.

## References

1. Roger C Mayer, James H Davis, and F David Schoorman. An integrative model of organizational trust. *Academy of Management Review*, 20(3):709–734, 1995.

2. Lik Mui, Mojdeh Mohtashemi, and Ari Halberstadt. Notions of reputation in multi-agents systems: a review. In *1st ACM Int Joint Conf on Autonomous Agents and Multiagent Systems, AAMAS'02*, pages 280–287, 2002.

3. C. Castelfranchi and R. Falcone. Principles of trust for mas: Cognitive anatomy, social importance, and quantification. In *3rd International Conference on Multi Agent Systems*, pages 72–79. IEEE, 1998.

4. P. Resnick, K. Kuwabara, R. Zeckhauser, and E. Friedman. Reputation systems: facilitating trust in internet interactions. *Comm of the ACM*, 43(12):45–48, 2000.

5. Jie Zhang, Robin Cohen, and Kate Larson. Combining trust modeling and mechanism design for promoting honesty in e-marketplaces. *Computational Intelligence*, 28(4):549–578, 2012.

6. Xin Liu, Anwitaman Datta, Hui Fang, and Jie Zhang. Detecting imprudence of 'reliable' sellers in online auction sites. In *11th IEEE International Conference on Trust, Security and Privacy in Computing and Communications, TrustCom'12*, pages 246–253, 2012.

7. Lizi Zhang, Siwei Jiang, Jie Zhang, and Wee Keong Ng. Robustness of trust models and combinations for handling unfair ratings. In *6th IFIP Int Conference on Trust Management*, volume 374, pages 36–51. Springer, 2012.

8. Joshua Gorner, Jie Zhang, and Robin Cohen. Improving the use of advisor networks for multi-agent trust modelling. In *9th IEEE Annual Conference on Privacy, Security and Trust, PST'11*, pages 71–78, 2011.

9. Mehrdad Nojoumian and Timothy C. Lethbridge. A new approach for the trust calculation in social networks. In *E-business and Telecom Networks: 3rd International Conference on E-Business*, volume 9 of *CCIS*, pages 64–77. Springer, 2008.

10. Audun Jøsang, Roslan Ismail, and Colin Boyd. A survey of trust and reputation systems for online service provision. *Decision Support Systems*, 43(2):618–644, 2007.

11. Mehrdad Nojoumian. *Novel Secret Sharing and Commitment Schemes for Cryptographic Applications*. PhD thesis, Department of Computer Science, UWaterloo, Canada, 2012.

12. Mehrdad Nojoumian and Douglas R. Stinson. Socio-rational secret sharing as a new direction in rational cryptography. In *3rd International Conference on Decision and Game Theory for Security (GameSec)*, volume 7638 of *LNCS*, pages 18–37. Springer, 2012.

13. Mehrdad Nojoumian and Douglas R. Stinson. Social secret sharing in coud computing using a new trust function. In *10th IEEE Annual International Conference on Privacy, Security and Trust (PST)*, pages 161–167, 2012.

14. Carol J. Fung, Jie Zhang, and Raouf Boutaba. Effective acquaintance management based on bayesian learning for distributed intrusion detection networks. *IEEE Trans on Network & Service Management*, 9(3):320–332, 2012.

15. Mehrdad Nojoumian, Douglas R. Stinson, and Morgan Grainger. Unconditionally secure social secret sharing scheme. *IET Information Security (IFS), Special Issue on Multi-Agent and Distributed Information Security*, 4(4):202–211, 2010.

16. Mehrdad Nojoumian and Douglas R. Stinson. Brief announcement: Secret sharing based on the social behaviors of players. In *29th ACM Symposium on Principles of Distributed Computing (PODC)*, pages 239–240, 2010.

17. Q. Li, A. Malip, K. Martin, S. Ng, and J. Zhang. A reputation-based announcement scheme for VANETs. *IEEE Trans on Vehicular Technology*, 61:4095 – 4108, 2012.

18. J. Zhang. A survey on trust management for VANETs. In *25th IEEE International Conf. on Advanced Information Networking and Applications, AINA'11*, pages 105–112, 2011.

19. J. Gorner, J. Zhang, and R. Cohen. Improving trust modelling through the limit of advisor network size and use of referrals. *E-Commerce Research & Applications*, 2012.

20. L. Zhang, H. Fang, W.K. Ng, and J. Zhang. Intrank: Interaction ranking-based trustworthy friend recommendation. In *10th IEEE International Conference on Trust, Security and Privacy in Computing and Communications, TrustCom'11*, pages 266–273, 2011.

21. Yonghong Wang and Munindar P. Singh. Evidence-based trust: A mathematical model geared for multiagent systems. *ACM Trans on Auto & Adaptive Sys.*, 5(4):1–28, 2010.

22. Yonghong Wang and Munindar P. Singh. Formal trust model for multiagent systems. In *20th International Joint Conference on Artificial Intelligence, IJCAI'07*, pages 1551–1556, 2007.

23. Yonghong Wang and Munindar P. Singh. Trust representation and aggregation in a distributed agent system. In *21st National Conf. on AI, AAAI'06*, pages 1425–1430, 2006.

24. Matthew Aitken, Nisar Ahmed, Dale Lawrence, Brian Argrow, and Eric Frew. Assurances and machine self-confidence for enhanced trust in autonomous systems. In *RSS 2016 Workshop on Social Trust in Autonomous Systems*, 2016.

25. Scott R Winter, Stephen Rice, Rian Mehta, Ismael Cremer, Katie M Reid, Timothy G Rosser, and Julie C Moore. Indian and american consumer perceptions of cockpit configuration policy. *Journal of Air Transport Management*, 42:226–231, 2015.

26. G.J. Mailath and L. Samuelson. *Repeated games and reputations: long-run relationships*. Oxford University Press, USA, 2006.

27. L. Mui. *Computational models of trust and reputation: Agents, evolutionary games, and social networks*. PhD thesis, Massachusetts Institute of Technology, 2002.

28. Renjian Feng, Shenyun Che, Xiao Wang, and Jiangwen Wan. An incentive mechanism based on game theory for trust management. *Security and Communication Networks*, 7(12):2318–2325, 2014.

29. M Harish, GS Mahalakshmi, and TV Geetha. Game theoretic model for p2p trust management. In *International Conference on Computational Intelligence and Multimedia Applications*, volume 1, pages 564–566. IEEE, 2007.

30. Martin J. Osborne and A. Rubinstein. *A course in game theory*. MIT press, 1994.

31. Mehrdad Nojoumian. Trust, influence and reputation management based on human reasoning. In *4th AAAI Workshop on Incentives and Trust in E-Communities*, pages 21–24, 2015.

32. Reid Kerr. *Addressing the Issues of Coalitions & Collusion in Multiagent Systems*. PhD thesis, UWaterloo, 2013.