# Conformity versus Manipulation in Reputation Systems*

Seyed Rasoul Etesami, Sadegh Bolouki, Angelia Nedić, Tamer Başar
Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801
Email: (etesami1,bolouki,angelia,basar1)@illinois.edu

*Abstract*— In this paper, we consider a reputation system, where a number of individuals express their opinions, modeled by discrete scalars in the interval [0,1], about an object and the object's score (reputation) is determined as the arithmetic mean of all expressed opinions. An individual's expressed opinion may or may not be consistent with her actual opinion, a continuous scalar in [0,1], for a variety of reasons. In this paper, we address in a unified, game-theoretic framework the influence of two opposing social behaviors, namely conformity and manipulation, on the outcome of a reputation system. For the purposes of this paper, conformity as a social behavior refers to the tendency of an individual to express an opinion that matches the public opinion, whereas manipulation refers to the tendency of an individual to express an opinion so as to manipulate the public opinion toward her actual opinion.

*Index Terms*— Conformity; manipulation; reputation system; Nash equilibrium; actual opinion; virtual opinion.

## I. Introduction

In recent years, there has been a wide range of studies to understand and analyze the underlying factors affecting opinion formation in social networks. In fact, the importance of such studies is more pronounced when they concern the outcome of critical elections such as political elections in human societies, or evaluating the worthiness of various candidates based on the average opinion of the electorate. These have led to introduction and study of different models for so-called reputation systems where the goal is to evaluate the reputation scores for a set of objects based on a collection of individuals' opinions. As some of the earlier models on this topic, one can consider the DeGroot model [1], Friedkin model [2], and the Hegselmann-Krause model [3], among several others [4], [5].

In this context, building on the earlier results, we introduce here a new model to study and explain some of the behavioral patterns which arise in many real reputation systems. Our work is in part motivated by the fact that in many real world reputation systems individuals have the tendency to announce their opinions in such a manner to be as close to the average public opinion as possible; a type of social behavior referred to as *conformity*. On the other hand, there still exist numerous examples which show that individuals express their opinions somehow to change the average score of the system toward their own actual belief; a type of social behavior herein referred to as *manipulation*. Although

there are many results in the mode of conformity in the opinion dynamics literature [6]–[8], there are not many comparable results for the role of manipulation in reputation systems, with some exceptions being [9], [10]. To provide a real evidence for such manipulative behavior in reputation systems, we have listed in the following three quotes from different reviewers who rated between 1 and 10 the movie "Interstellar" on IMDB official website (1 for the worst and 10 for the best):

- "...I give 1 star to bring balance to the current rating, in reality this movie is of course not that bad."
- "My honest rating would be 6 for that movie but I rated it 1 to balance the 'emotional' ratings."
- "...It's like 5-7ish depending on how you want to weight various factors. I'm just doing this because all the hype, and all the 'stellar' (pun intended) reviews are blowing my mind and need to be balanced out..."

What can be understood from these quotes is the fact that sometimes or quite often the social entities express their opinions in such a way to manipulate the outcome of the system rather than reporting their actual opinions or attempting to conform with the average public opinion. On the other hand, trust, local experiences and feedback are important factors which gradually influence individuals to report their opinion more toward the average opinion of the others. In other words, the entities have up to some degree conformity in their evaluation process in the reputation systems.

In this paper, we consider a *reputation system* that provides the trend of an object's score over time, such as *tracking polls*. For instance, see Fig. 1 as an example of a tracking poll about president Obama's job approval. Our aim is to model the coexistence of conformity and manipulation in such a system and analyze the evolution of the outcome of the system over time. For this purpose, we propose a game-theoretic framework in which an individual's actual opinion is modeled by a continuous scalar in the interval $[0, 1]$, while her expressed opinion at any time instant can be either 0 or 1. The choice of such discrete variables for the expressed opinions is again inspired by the vast majority of opinion polls, that constitute the most popular type of reputation systems. We first investigate a model in which each individual is only trying to manipulate the average score of the system toward her own actual opinion by expressing an opinion which may be inconsistent with her actual opinion. We interpret the expressed opinions as players' actions.
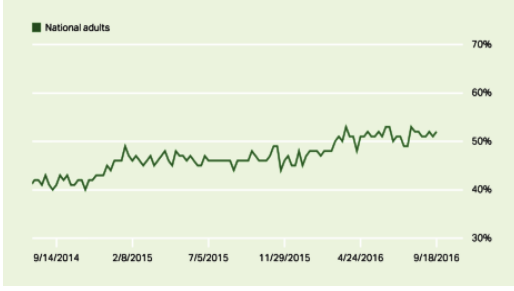
Fig. 1. *Gallup Daily*'s tracking poll result on President Obama's job approval in the past two years with one week tracking window [11].

We will see that when everyone is playing based on such manipulative behavior, the average weight of the system will always converge although the players' actions will converge to a unique pure-strategy Nash equilibrium only when a certain condition holds. This game is then generalized so as to also account for conformity, and the game dynamics and the outcome of the system are once again investigated accordingly.

The paper is organized as follows: We formulate in Section II a reputation system purely based on manipulation as a repeated game with a finite number of players. In particular, we show the convergence of the best response dynamics to a unique pure-strategy Nash equilibrium with an upper bound on its convergence time when the limiting average weight of the system does not coincide with the actual opinions of the players. In Section III, we turn our attention to characterizing the stationary distribution of a reputation system when the players express their virtual opinions via a mixture of manipulative and conformist behaviors. We conclude the paper by identifying some future direction of research in Section IV.

**Notations**: For a positive integer $n$, we let $[n] := \{1, 2, \ldots, n\}$. For a vector $v \in \mathbb{R}^n$, we let $v_i$ be the $i$th entry of $v$, and $\|v\|$ be its Euclidean norm. We denote the transpose of a vector $v$ by $v'$. We let $\mathbf{1}$ and $\mathbf{0}$ be vectors, of appropriate sizes, of all ones and zeros, respectively. We use $|S|$ to denote the cardinality of a finite set $S$. For a stochastic $n \times n$ matrix $Q$, we denote its second largest eigenvalue by $\lambda_2(Q)$ and the ergodicity coefficient of $Q$ by $\tau(Q)$, i.e., $\tau(Q) = 1 - \min_{i \neq j}(\sum_{\ell=1}^{n} \min(Q_{i\ell}, Q_{j\ell}))$.

## II. ABSOLUTE MANIPULATION SYSTEM WITH BINARY ACTION SPACE

In this section, we formulate a reputation system as a repeated game in which each player's objective is to manipulate the average weight of the system toward her own actual belief. We consider a set of $n$ players (agents). For a given object, we assume that each player $i \in [n]$ has an actual opinion, denoted by $y_i \in [0, 1]$, which is private and only known to that player. Moreover, we assume that each player $i$ has a virtual opinion $x_i \in \{0, 1\}$, which is the opinion expressed by him at each round of the game and can be seen by others. One can view $x_i$ as the action taken by player $i$ in order to force the average weight of the system toward her private belief $y_i$.

We now define a sequence of inductive games $\mathcal{G}_1, \mathcal{G}_2, \ldots$, each of which has the same set of players, that is $[n]$, and the same action sets, that is $\{0, 1\}^n$. The cost functions associated with the players are defined as follows: Denoting the actions of player $i$ in games $\mathcal{G}_1, \ldots, \mathcal{G}_t$ by $x_i(1), \ldots, x_i(t)$, we let $w_t$ be the average weight of the system up to time $t$, i.e., $w_t = \frac{1}{nt} \sum_{\ell=0}^{t} \sum_{j=1}^{n} x_j(\ell)$ (by convention we set $w_0 = 0$). Moreover, we define the cost associated with player $i$ in game $\mathcal{G}_{t+1}$ by

$$C_i^{(t+1)}(x_i, w_t) = \left( x_i + y_i - w_t - \frac{1}{2} \right)^2. \tag{1}$$

In other words, each game $\mathcal{G}_t$ is played assuming that the total average of players' actions in all the previous games $\mathcal{G}_1, \ldots, \mathcal{G}_{t-1}$ is available.[1] We say that the sequence of games $\mathcal{G}_1, \mathcal{G}_2, \ldots$ converges if $\lim_{t \to \infty} C_i^{(t+1)}(x_i, w_t)$ exists, i.e., the sequence of plays on the game sequence results in a well-defined game in the limit.[2]

We point out that the reason for formulating the cost function as (1) is the following: in order for a player $i$ to have the most influence on the average score of the system toward her actual opinion, she will look at the difference between her actual opinion and the average weight of the system up to time $t$ and depending on whether this difference is less than or greater than $\frac{1}{2}$, she will choose action 0 or 1, respectively, at the next time step. It is worth noting that the average of the system at some stage depends not only on the players' actions at that stage, but also on the entire history of the players' actions.

*Definition 1:* We say that $(x_1^*, \ldots, x_n^*) \in \{0, 1\}^n$ constitutes a *pure-strategy Nash equilibrium* for the sequence of games $\mathcal{G}_1, \mathcal{G}_2, \ldots$ (or simply a pure-strategy NE), if i) this sequence of games converges, and ii) $(x_1^*, \ldots, x_n^*)$ forms a pure-strategy Nash equilibrium for the limit game. ∎

*Remark 1:* Using the cost function (1), we can write the best response of player $i$ in game $\mathcal{G}_t$ by $x_i(t+1) = \mathrm{sgn}(y_i - w_t)$, where $\mathrm{sgn}(a)$ equals to 1 if $a > 0$, and 0, otherwise. ∎

In the following theorem, we show that if at each game $\mathcal{G}_t, t = 1, 2, \ldots$, every player plays its best response with respect to the cost function (1), the average weight of the reputation system will eventually converge to a unique limit point, even though the actions of the players may or may not converge to a pure-strategy Nash equilibrium of the limit game. We define $k^*$ as the unique integer in $[n]$ for which at least one of the following inequality pairs hold:

$$y_{k^*} < \frac{n - k^*}{n} < y_{k^*+1}, \tag{2}$$

$$\frac{n - k^*}{n} \leq y_{k^*} < \frac{n - k^* + 1}{n}. \tag{3}$$

We note that both (2) and (3) cannot hold simultaneously.

[1] In the case of tracking polls as reputation system, such an assumption corresponds to having an unbounded tracking window.

[2] Note that in the definition of cost function (1), the actions of the players are not coupled through their instantaneous costs, but through past history of actions through $w_t$.

*Theorem 1:* If all the players play their best responses at each game $\mathcal{G}_t$, $t = 1, 2, \ldots$, then the average weight $w_t$ always converges to a unique point $w^* \in [0, 1]$ as time grows. Consequently, the sequence of games $\mathcal{G}_t$, $t = 1, 2, \ldots$ converges. Furthermore, $w^*$ can be expressed as

$$w^* = \begin{cases} \frac{n - k^*}{n} & \text{if (2) holds,} \\ y_{k^*} & \text{if (3) holds.} \end{cases}$$

Lastly, in case (2) holds, (i) the best response strategies of the players converge to a pure-strategy Nash equilibrium, (ii) $|w_t - w^*| = O(\frac{1}{t})$.

*Proof:* For any $w \in [0, 1]$ let us define $k(w)$ to be the number of players whose actual opinions are no larger than $w$, i.e., $k(w) = |\{i : y_i \leq w\}|$. Since each player is playing its best response based on (1), using Remark (1) at the next time step exactly $k(w_t)$ of the players will announce $x_i(t+1) = 0$ as their virtual opinion and the remaining will announce $x_i(t+1) = 1$. Hence, at the next time step we will have

$$w_{t+1} = \frac{k(w_t)}{n} \times \frac{t w_t + 0}{t + 1} + \frac{n - k(w_t)}{n} \times \frac{t w_t + 1}{t + 1}$$
$$= \frac{t}{t + 1} w_t + \frac{n - k(w_t)}{n(t + 1)}. \tag{4}$$

Next we show that $\lim_{t \to \infty} w_t$ always exists. To see this let us define $f(w) : [0, 1] \to [-1, 1]$ by $f(w) = w - \frac{n - k(w)}{n}$. First we note that $f(0) \leq 0$ and $f(1) = 1$. Moreover, if $w < w'$, then we have $k(w) \leq k(w')$, and hence $f(w') - f(w) = (w' - w) + \frac{k(w') - k(w)}{n} > 0$, which shows that $f(w)$ is a strictly increasing function. Let us define $w^* = \sup\{w : f(w) \leq 0\}$. We will next show that $\lim_{t \to \infty} w_t = w^*$. We consider two cases:

- **Case 1**: $w_t < w^*$. In this case, by definition of $w^*$ we must have $f(w_t) < 0$. Therefore, using (4) we can write

$$w_{t+1} - w_t = -\frac{w_t}{t + 1} + \frac{n - k(w_t)}{n(t + 1)} = -\frac{1}{t + 1} f(w_t) > 0.$$

- **Case 2**: $w_t > w^*$. Again, by definition of $w^*$ we must have $f(w_t) > 0$. Hence, we have

$$w_{t+1} - w_t = -\frac{1}{t + 1} f(w_t) < 0.$$

Therefore, depending on whether $w_t$ is smaller or larger than $w^*$, at the next time step $w_{t+1}$ strictly increases or strictly decreases. Since the sequence $\{w_t\}$ lies inside the compact set $[0, 1]$, it must have at least one accumulation point.

First let us assume that $w^*$ is an accumulation point. We show that all the accumulation points of $\{w_t\}$ are indeed the same and equal to $w^*$. Otherwise, let us assume that, to the contrary, there exists another accumulation point $u^* \neq w^*$. Let us first assume that $u^* < w^*$ and set $\delta := \frac{w^* - u^*}{2} > 0$ (the proof for the case where $u^* > w^*$ follows exactly similar lines of argument). Since $\lim_{t \to \infty}(w_{t+1} - w_t) = \lim_{t \to \infty} -\frac{1}{t + 1} f(w_t) = 0$, there exists $N \in \mathbb{N}$ such that $|w_{t+1} - w_t| < \delta$, $\forall t \geq N$. Next we claim that $w_t \leq u^* + \delta$, $\forall t > N$. Otherwise, let us assume that for some

$t_1 > N$ we have $w_{t_1} > u^* + \delta = w^* - \delta$. Now if $w_{t_1} < w^*$, then using case 1, at the next time step $w_{t_1+1}$ will increase, and hence, increase its distance with $u^*$ even further, i.e., $w_{t_1+1} > w_{t_1} > u^* + \delta$. Moreover, if $w_{t_1} > w^*$, we have $w_{t_1+1} > w_{t_1} - \delta > w^* - \delta = u^* + \delta$. These together show that if for some $t_1 > N$ we have $w_{t_1} > u^* + \delta$, then for any time $t > t_1$ we must have $w_t > u^* + \delta$, which is in contradiction with $u^*$ being an accumulation point.

Thus we have shown that $w_t \leq u^* + \delta$, $\forall t > N$. In particular, we must have $f(w_t) \leq f(u^* + \delta) < 0$ (recall that $f(\cdot)$ is a strictly increasing function with $f(w) < 0$, for any $w < w^*$ ). Now, for any integer $M > N$, we can write

$$w_M = w_N + \sum_{t=N}^{M-1}(w_{t+1} - w_t) = w_N + \sum_{t=N}^{M-1} -\frac{1}{t + 1} f(w_t)$$
$$\geq w_N + \sum_{t=N}^{M-1} -\frac{1}{t + 1} f(u^* + \delta)$$
$$= w_N - f(u^* + \delta) \sum_{t=N}^{M-1} \frac{1}{t + 1} \tag{5}$$

But for large $M$ the right hand side of the above relation goes to infinity while the left hand side is always in $[0, 1]$. This contradiction shows that there cannot be another accumulation point $u^* < w^*$, which shows that $w^*$ is a unique accumulation point of $\{w_t\}$

Next, we argue that $w^*$ is indeed an accumulation point of $\{w_t\}$. Otherwise, since the set of accumulation points of $\{w_t\}$ is a closed set, there must exist $r > 0$ small enough and $K_1 \in \mathbb{N}$ sufficiently large such that $(w^* - r, w^* + r) \cap \{w_t, : t \geq K_1\} = \emptyset$. For such $r$, and since $\lim_{t \to \infty}(w_{t+1} - w_t) = 0$, there must exist $K_2 \in \mathbb{N}$ such that for $t \geq K_2$ we have $|w_{t+1} - w_t| < r$. Taking $K = \max(K_1, K_2)$, we have

$$\forall t \geq K, \quad (w^* - r, w^* + r) \cap \{w_t\} = \emptyset, \ |w_{t+1} - w_t| < r.$$

In other words, $\{w_t, t \geq K\}$ must be totally either smaller than $w^* - r$, or larger than $w^* + r$. But using a similar argument as in (5) we reach a contradiction since starting from any time, the sequence $\{w_t\}$ cannot always stay away by a constant distance $r > 0$ smaller or larger than $w^*$. This shows that, indeed $w^*$ is an accumulation point, and hence, i.e., $\lim_{t \to \infty} w_t = w^*$.

Next we proceed to characterize $w^*$ based on the actual opinions $y_i, i \in [n]$. Let us sort the actual opinions in a non decreasing order $y_1 \leq y_2 \leq \ldots \leq y_n$ and define $A$ to be the set of indices of all the players whose actual opinions are exactly equal to $w^*$ i.e., $A = \{i : y_i = w^*\}$. Now for any $0 < \epsilon < \min\{|w^* - y_i| : i \in [n] \setminus A\}$ and using the fact that $w^* = \sup\{w : f(w) \leq 0\}$, we can write

$$f(w^* + \epsilon) = w^* + \epsilon - \frac{n - k(w^* + \epsilon)}{n} > 0,$$
$$\Rightarrow w^* + \epsilon > \frac{n - k(w^* + \epsilon)}{n} = \frac{n - k(w^*)}{n}. \tag{6}$$

Moreover, since $k(w^* - \epsilon) = k(w^*) - |A|$, we get

$$f(w^* - \epsilon) = w^* - \epsilon - \frac{n - k(w^* + \epsilon)}{n} < 0,$$

$$\Rightarrow w^* - \epsilon < \frac{n - k(w^* - \epsilon)}{n} = \frac{n - k(w^*) + |A|}{n}. \quad (7)$$

Combining (6) and (7), we have

$$-\epsilon < w^* - \frac{n - k(w^*)}{n} < \epsilon + \frac{|A|}{n}. \quad (8)$$

By letting $\epsilon \to 0$ in (8), we have two possibilities:

- $|A| = 0$. In this case $w^* = \frac{n - k(w^*)}{n}$. Moreover, since we have $k(w^*) = |\{i : y_i \leq w^*\}| = |\{i : y_i \leq \frac{n - k(w^*)}{n}\}|$, there exists a unique integer $k^* := k(w^*)$ such that $y_{k^*} < \frac{n - k^*}{n} < y_{k^*+1}$, in which case $w^* = \frac{n - k^*}{n}$.

- $|A| > 0$. In this case $\frac{n - k(w^*)}{n} \leq w^* < \frac{n - k(w^*) + |A|}{n}$. Now since all the $|A|$ values $y_{k(w^*) - |A| + 1} = \ldots = y_{k(w^*)}$ are equal to $w^*$, and they are in the semi-closed interval $[\frac{n - k(w^*)}{n}, \frac{n - k(w^*) + |A|}{n}) = \bigcup_{j=k(w^*)}^{k(w^*) + |A| - 1} [\frac{n - j}{n}, \frac{n - j + 1}{n})$, we have that there exists a unique integer index $k^* \in \{k(w^*) - |A| + 1, \ldots, k(w^*)\}$ such that $\frac{n - k^*}{n} \leq y_{k^*} \leq \frac{n - k^* + 1}{n}$, in which case $w^* = y_{k^*}$.

Finally, we show that a pure-strategy Nash equilibrium in the limit game exists only when $|A| = 0$, or equivalently, $w^* \neq y_i, \forall i \in [n]$. This is because if $w^* \neq y_i, \forall i \in [n]$, for a sufficiently large constant $L$ we have $|w_t - w^*| < \min\{|w^* - y_i|, i \in [n]\}$. In particular, for $t > L$, we have $k(w_t) = k(w^*)$, i.e., the actions of the players from time $L$ onward do not change, and hence

$$x_i^* = \begin{cases} 0 & \text{if } i \leq k(w^*), \\ 1 & \text{if } i > k(w^*), \end{cases}$$

constitutes a pure-strategy Nash equilibrium (in the case where $|A| > 0$, the actions of those players whose indices belong to $A$ will oscillate between 0 and 1, and will never converge to a pure-strategy Nash equilibrium, despite the fact that $w_t$ always converges). Now for any $t > L$ we can write

$$\begin{aligned} w_{t+1} - w^* &= \frac{t}{t+1}(w_t - w^*) + \frac{1}{t+1}\left(\frac{n - k(w_t)}{n} - w^*\right) \\ &= \frac{t}{t+1}(w_t - w^*) + \frac{1}{t+1}\left(\frac{n - k(w^*)}{n} - w^*\right) \\ &= \frac{t}{t+1}(w_t - w^*). \end{aligned}$$

Therefore, repeating the above process inductively for $t > L$, and since $L$ is a constant, we get

$$\begin{aligned} |w_{t+1} - w^*| &= \left(\prod_{\ell=T}^{t} \frac{\ell}{\ell+1}\right)|w_L - w^*| = \frac{L}{t+1}|w_L - w^*| \\ &\leq \frac{L}{t+1} = O\left(\frac{1}{t}\right). \end{aligned}$$

∎

In fact the condition of existence of a pure-strategy Nash equilibrium in the sequence of absolute manipulation games is not a very restrictive condition. In the following theorem we compute the probability of existence of an equilibrium when the actual opinions are distributed based on some continuous distribution in the interval $[0, 1]$. In particular, we will see that for the case of uniform distribution of actual opinions almost half of the times an equilibrium exists.

*Theorem 2:* If the actual opinions of the agents are i.i.d, continuous random variables with a cumulative function $F(y)$ on the interval $[0, 1]$, then the probability of existence of a pure-strategy Nash equilibrium in the absolute manipulation system equals to $\sum_{k=0}^{n} \binom{n}{k} F^k(\frac{n-k}{n})(1 - F(\frac{n-k}{n}))^{n-k}$.

*Proof:* From Theorem 1 we know that the absolute manipulation system admits a pure-strategy Nash equilibrium if and only if there exists an integer $k$ such that $y_1 \leq \ldots \leq y_k < \frac{n-k}{n} < y_{k+1} \leq \ldots \leq y_n$. Since $y_i \sim F(y)$, without any loss of generality in probability computation, we can assume that all the $y_i, i \in [n]$ are not equal to $\{\frac{n-k}{n}, k = 1, \ldots, n-1\}$ (the probability that a continuous random variable takes finitely many specific numbers equals to 0). Using Theorem 1 and for any arbitrary but fixed $k = 1, 2, \ldots, n-1$, the probability that $\frac{n-k}{n}$ is an equilibrium equals to choosing exactly $k$ of $y_i$ and putting them on the left side of $\frac{n-k}{n}$ and the remaining $n - k$ actual opinions on the right side of $\frac{n-k}{n}$. Since the former happens with probability $F^k(\frac{n-k}{n})$ and the latter happens with probability $(1 - F(\frac{n-k}{n}))^{n-k}$, the probability that $\frac{n-k}{n}$ is an equilibrium is exactly $\binom{n}{k} F^k(\frac{n-k}{n})(1 - F(\frac{n-k}{n}))^{n-k}$. Finally, by uniqueness of the equilibrium, the events that both $\frac{n-k}{n}$ and $\frac{n-k'}{n}$ for some $k \neq k'$ are equilibrium points are disjoint. Therefore, we can write

$$\begin{aligned} \mathbb{P}\{\text{Pure NE exists}\} &= \sum_{k=1}^{n-1} \mathbb{P}(\frac{n-k}{n} \text{ is a pure NE}) \\ &= \sum_{k=0}^{n} \binom{n}{k} F^k(\frac{n-k}{n})(1 - F(\frac{n-k}{n}))^{n-k}. \end{aligned}$$

∎

*Corollary 1:* For the uniform distribution, the probability of existence of a pure-strategy Nash equilibrium approaches $\frac{1}{2}$ for large $n$.

*Proof:* For the uniform distribution we have $F(y) = y$. Thus the probability of existence of a pure-strategy equilibrium equals to $\sum_{k=0}^{n} \binom{n}{k} \left(\frac{k}{n}\right)^{n-k} \left(1 - \frac{k}{n}\right)^k$. Using Stirling approximation it can be shown that this expression is an increasing function of $n > 2$ and approaches $\frac{1}{2}$ as $n$ gets sufficiently large. ∎

In Fig. 2 we have illustrated the probability of existence of a pure-strategy Nash equilibrium in the absolute manipulation system for uniform distribution $F(y) = y$ (red curve) and normalized exponential distribution with cumulative function $F(y) = \frac{e}{e-1}(1 - e^{-y})$ (blue curve). As it can be seen, in both cases the probability of existence of an equilibrium is more than $0.45$ and monotonically increases when $n$ gets larger. In particular, for the uniform distribution this probability approaches $0.5$.

In fact, increasing the number of players to infinity can be viewed as a continuum of opinions in the interval $[0, 1]$ such
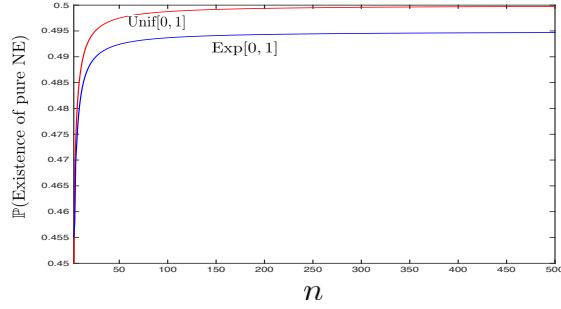
Fig. 2. Probability of existence of a pure-strategy Nash equilibrium in the limit game for uniform (red) and normalized exponential (blue) distribution.

that the proportion of players whose opinion belong to the infinitesimal interval $[s, s + \delta)$ is given by $g(s)$, where $g(\cdot)$ is a probability density function such that $\int_0^1 g(s)ds = 1$. Again, one can consider a variant of the absolute manipulation system with a continuum of players. With some abuse of notation, if we denote the average weight of the system at time $t$ by $w_t$, then at the next time step the average weight of the system can be written as

$$w_{t+1} = \frac{t}{t+1}w_t + \frac{1}{t+1}\int_{w_t}^1 g(s)ds. \tag{9}$$

Now assuming that $g(\cdot) : [0, 1] \to \mathbb{R}^{>0}$ is a continuous function, using a similar argument as in the proof of Theorem 1 one can show that the recursion of (9) has a unique limit point which is given by the solution of $w^* = \int_{w^*}^1 g(s)ds$ (note that $h(w)$ defined by $h(w) := w - \int_w^1 g(s)ds$ is continuous and strictly increasing. Moreover, $h(0) = -1$ and $h(1) = 1$, which by the Mean Value Theorem implies that $h(w) = 0$ at some unique point $w^* \in [0, 1]$).

It is worth noting that one of the main reasons why the dynamics in (4) or (9) converge is because they take into account the history of players' actions in their cost functions, and hence, in computing their next actions. To see why this is crucial, let us consider a scenario in which the dynamics of (9) are only based on the average of the current actions. Then, we have

$$v_{t+1} = \int_{v_t}^1 g(s)ds, \tag{10}$$

where $v_t$ denotes the average instant weight of the system at time $t$. Although the dynamics of (10) admit a unique fixed point satisfying $w^* = \int_{w^*}^1 g(s)ds$ (by a similar argument as above), it is not clear a priori whether the iterations in (10) will converge to such a unique fixed point or not. In fact, the convergence of $\{v_t\}$ heavily depends on the initial instant weight of the system, $v_0$, and the density function $g(s)$. As an example, defining $G(v) := \int_v^1 g(s)ds$, if there exists $v_0$ such that $v_0 < G(G(v_0)) < G(v_0)$, then starting from the initial weight $v_0$, the iteration of the system (10) will converge to the point $w^*$. To resolve this issue and in order to define the players' cost functions based on the instant average of the actions (and not the average of the entire history), in the next section we introduce a variant of the absolute manipulation game by introducing noisy actions for the players.

## III. CONFORMITY V.S. MANIPULATION WITH INSTANT AVERAGE OVER NOISY ACTION SPACE

In this section we consider an effect of so called noise in the actions of the players in a reputation system. This kind of noise can be thought of as capturing an uncertainty players face about the history of the game, trembles in strategy choices, unfamiliarity of the new players with respect to the reputation system, or deliberate experimentation, see, e.g., [12]–[14].

In this model, as opposed to the scenario considered in Section II, we assume that only the *instant average* of the actions is available to the players. In the case of tracking polls, such an assumption can be interpreted as having a small tracking window. The general idea here is that at each time instant $t = 1, 2, \ldots$ each player plays manipulative or conformist with respective probabilities $p$ and $1 - p$ for some $p \in (0, 1)$ (note that $p \neq 0, 1$, and it is uniform across players). In other words, denoting the instant average weight of the system at time $t$ by $v_t$, i.e., $v_t = \frac{\sum_{i=1}^n x_i(t)}{n}$, player $i \in [n]$ chooses her next action according to:

$$x_i(t+1) = \begin{cases} \underset{x \in \{0,1\}}{\text{argmin}} \ (x + y_i - v_t - \frac{1}{2})^2 & \text{w.p.} \ \ p, \\ \underset{x \in \{0,1\}}{\text{argmin}} \ (x - v_t)^2 & \text{w.p.} \ \ 1 - p. \end{cases} \tag{11}$$

In other words, at each time instant $t$, player $i$ has two evaluations for her cost: $(x + y_i - v_t - \frac{1}{2})^2$ and $(x - v_t)^2$, in which with probability $p$ he plays her best response with respect to the first cost function and with probability $1 - p$ he will take her best action with respect to her second cost function.

*Theorem 3:* The noisy dynamics defined by (11) will converge to a unique stationary distribution with geometric convergence rate of $O\left((1 - p^n(1-p)^n)^{\frac{t}{2}}\right)$.

*Proof:* To analyze the dynamics of (11), we formulate them as an evolution of a Markov chain over the entire action space. Let us consider a Markov chain over the product action space $\mathcal{X} := \{0, 1\}^n$. For every two states $x_1, x_2 \in \mathcal{X}$, the probability of moving from $x_1$ to $x_2$ is fully determined using (11) by $x_1$, $x_2$, and $\{y_i\}_{i \in [n]}$ (note that $v_t = \frac{x(t)'\mathbf{1}}{n}$, where $x(t) = (x_1(t), \ldots, x_n(t))'$ denotes the state of the Markov chain at time $t$ and $\mathbf{1}$ is the vector of all 1s). Therefore, we can rewrite the dynamics of (11) as evolution of a Markov chain as

$$x(t+1) = Px(t), \quad x(t) \in \mathcal{X}, \tag{12}$$

where $P$ is a stochastic matrix of dimension $2^n \times 2^n$.

Next we argue that indeed $P$ is an irreducible and aperiodic matrix, which in turn implies that (12) converges to a unique stationary distribution. To show this, let us define $\mathcal{X}_1$ to be the set of all the states whose instant averages are less than or equal to $\frac{1}{2}$, i.e., $\mathcal{X}_1 = \{x \in \mathcal{X} : \frac{x'\mathbf{1}}{n} \le \frac{1}{2}\}$. We establish the following facts:

- Fact 1) $\mathbb{P}(x(t+1) = x | x(t) = \mathbf{0}) > 0, \forall x \in \mathcal{X}$. In other words, from the state $\mathbf{0} := (0, 0, \ldots, 0)$ with positive

probability we can go to any other state $x \in \mathcal{X}$. This is because the average instant weight at state $\mathbf{0}$ is zero, thus manipulative play for a player means choosing action 1, and conformist play means choosing action 0. Therefore, we have

$$\mathbb{P}(x(t+1) = x | x(t) = \mathbf{0}) = p^{x'\mathbf{1}}(1-p)^{n-x'\mathbf{1}} > 0.$$

- Fact 2) $\mathbb{P}(x(t+1) = \mathbf{0} | x(t) = x) > 0, \forall x \in \mathcal{X}_1$. This is because the average instant weight at state $x \in \mathcal{X}_1$ is at most $\frac{1}{2}$, thus conformist play for a player means choosing 0, hence,

$$\mathbb{P}(x(t+1) = \mathbf{0} | x(t) = x) \geq (1-p)^n > 0.$$

- Fact 3) $\mathbb{P}(x(t+1) = \mathbf{1} | x(t) = x) > 0, \forall x \in \mathcal{X} \setminus \mathcal{X}_1$. This is because the average instant weight at state $x \in \mathcal{X} \setminus \mathcal{X}_1$ is more than $\frac{1}{2}$, thus conformist play for a player means choosing 1. Thus we get

$$\mathbb{P}(x(t+1) = \mathbf{1} | x(t) = x) \geq (1-p)^n > 0.$$

- Fact 4) $\mathbb{P}(x(t+1) = \mathbf{0} | x(t) = \mathbf{1}) > 0$. This is because the average instant weight at state $\mathbf{1}$ is equal to 1. Since $y_i \leq 1, \forall i \in [n]$, manipulative play for a player means choosing 0. Therefore,

$$\mathbb{P}(x(t+1) = \mathbf{0} | x(t) = \mathbf{1}) = p^n > 0.$$

From Facts 1-4 one can easily check that from every state the chain can go to any other state with positive probability. This shows that indeed the Markov chain (12) is irreducible. Moreover, there is a self-loop at state $\mathbf{0}$ (since by Fact 1 we have $\mathbb{P}(x(t+1) = \mathbf{0} | x(t) = \mathbf{0}) = (1-p)^n > 0$), which shows that the chain is aperiodic as well. These together show that the Markov chain (12) will converge to a unique stationary distribution.

Finally, Facts 1-4 show that after at most two steps, every node can reach state $\mathbf{0}$ with a positive probability which is lower bounded by $p^n(1-p)^n$. In other words, $p^n(1-p)^n$ is a lower bound for all the entries of the column corresponding to the state $\mathbf{0}$ in $P^2$. Since we always have $|\lambda_2(P^2)| \leq \tau(P^2)$ (see, e.g., [15]), we can write

$$|\lambda_2(P^2)| \leq \tau(P^2) = 1 - \min_{i \neq j}\left(\sum_{\ell=1}^{2^n} \min(P_{i\ell}^2, P_{j\ell}^2)\right)$$
$$\leq 1 - p^n(1-p)^n.$$

This, in view of the fundamental theorem of Markov chains, shows that

$$\|x(t) - x^*\| \leq (\lambda_2(P^2))^{\frac{t}{2}} \|x(0) - x^*\|$$
$$\leq \left(1 - p^n(1-p)^n\right)^{\frac{t}{2}} \|x(0) - x^*\|,$$

where $x^*$ is the unique stationary distribution of the chain. ∎

## IV. Conclusion

In this paper we have studied the role of manipulation versus conformity in reputation systems from a game-theoretic perspective. We have first analyzed an absolute manipulation game in which every player is only trying to change the average weight of the system toward her own actual opinion. We have observed that although the average weight of the system always converges, depending on the actual opinions the limit system may or may not have a pure-strategy NE. Furthermore, we have observed that taking the history of the players' actions in our model is crucial for convergence of best-response dynamics. We then removed the dependency of the cost functions on the entire history by introducing a noisy action game so that the best responses of the players now depend on only the last stage and not the entire history of the action. In particular, we have characterized the stationary distribution point of the system using a Markov model.

As a future research direction one can generalize the results in this paper into higher dimensions, where the opinion of the players can take vector form, e.g., when there are multiple correlated topics of interest being scored. Moreover, studying the mean field behavior of the above models is another interesting problem. Finally, in Section III, we had assumed homogeneous probability parameter $p$ across all players. Extending our results to an in-homogeneous case where there are finitely many different types would be an interesting research direction.

## References

[1] M. H. DeGroot, "Reaching a consensus," *Journal of the American Statistical Association*, vol. 69, no. 345, pp. 118–121, 1974.

[2] N. E. Friedkin and E. C. Johnsen, "Social influence networks and opinion change," *Advances in Group Processes*, vol. 16, no. 1, pp. 1–29, 1999.

[3] R. Hegselmann and U. Krause, "Opinion dynamics and bounded confidence models, analysis, and simulation," *Artificial Societies and Social Simulation*, vol. 5, pp. 1–33, 2002.

[4] J. B. Stiff and P. A. Mongeau, *Persuasive Communication*. Guilford Press, 2003.

[5] N. E. Friedkin and E. C. Johnsen, "Social influence network theory," *Cambridge Univ. Press, New York*, 2011.

[6] J. Ghaderi and R. Srikant, "Opinion dynamics in social networks: A local interaction game with stubborn agents," in *American Control Conference (ACC), 2013*. IEEE, 2013, pp. 1982–1987.

[7] S. R. Etesami and T. Başar, "Game-theoretic analysis of the Hegselmann-Krause model for opinion dynamics in finite dimensions," *IEEE Transactions on Automatic Control*, vol. 60, no. 7, pp. 1886–1897, 2015.

[8] Z. Xu, J. Liu, and T. Başar, "On a modified DeGroot-Friedkin model of opinion dynamics," in *American Control Conference (ACC), 2015*. IEEE, 2015, pp. 1047–1052.

[9] M. Förster, A. Mauleon, and V. J. Vannetelbosch, "Trust and manipulation in social networks," 2014.

[10] A. Gionis, E. Terzi, and P. Tsaparas, "Opinion maximization in social networks," *SIAM*, 2013.

[11] *http://www.gallup.com/poll/125729/obama-job-approval-weekly.aspx*.

[12] G. Ellison, "Learning, local interaction, and coordination," *Econometrica: Journal of the Econometric Society*, pp. 1047–1071, 1993.

[13] M. Kandori, G. J. Mailath, and R. Rob, "Learning, mutation, and long run equilibria in games," *Econometrica: Journal of the Econometric Society*, pp. 29–56, 1993.

[14] T. Başar and G. J. Olsder, *Dynamic Noncooperative Game Theory*. Siam, vol 23, 1999.

[15] E. Seneta, *Non-negative Matrices and Markov Chains*. Springer Science & Business Media, 2006.