

# Controlling Fake Reviews

Yuta Yasui

September 29, 2020

## Abstract

In this paper, I theoretically analyze fake reviews on a platform market, where a seller makes fake reviews through incentivized transactions and the sales depends on its rating based on a review history. The platform can control the incentive of fake reviews by changing parameters of the rating system such as its filtering policy and weights on the past reviews. At the equilibrium, the number of fake reviews is increasing in the current quality and decreasing in the current reputation. Since the fake reviews have a positive relationship with the underlying quality, rational consumers can find a rating more informative with fake reviews than one without fake reviews. On the other hand, credulous consumers suffer from a bias caused by boosted reputations. Stringent filtering policy can decrease the expected amount of fake reviews and the bias of the credulous consumers, but it can decrease the informativeness for the rational consumers at the same time. In terms of the weights on the review history, the rational consumers benefit from higher weights on the past reviews than the optimal weights without fake reviews.

[This is a preliminary draft. The latest version is found here.]

*“Fully ref\*\*ded after RE\*\*W*

*If you are interested dm and comment”*

— a post on Facebook

## 1 Introduction

Online platform markets are growing all over the world, so both businesses and their customers increasingly rely on reviews on the platforms.<sup>1</sup> At the same time, incentives for sellers to make fake reviews is also growing. Washington Post (Dwoskin and Timberg, 2018) reports that, based on fake review detection algorithms, 50.7% of reviews for Bluetooth headphones, 58.2% of Bluetooth speakers, 55.6% of weight-loss pills, 67.0% of testosterone booster on Amazon are suspicious. How do the sellers make fake reviews? The sellers can post information of their products with refund offers, which are typically finalized via PayPal after purchases and positive reviews on Amazon.

---

<sup>1</sup>Hollenbeck (2018); Hollenbeck et al. (2019) show that ratings work as a substitute of other form of advertisement or brand names, and this pattern is getting stronger over time in the hotel industry. Reimers and Waldfogel (2020) exhibit that the existence of star ratings has 15 times as the impact on consumer surplus as the professional reviews on New York Times. For the institutional details and data analysis on platforms and ratings, see also Belleflamme and Peitz (2018)

Such a review is considered as a review with a verified purchase and reflected to the star rating (until it is detected by Amazon).<sup>2</sup> He et al. (2020) connects such refund offers on Facebook with the product listings on Amazon and show that the positive correlation between the refund offers on Facebook and the product’s performance on Amazon such as ratings, sales ranking, and the number of reviews. The regulators have been concerned about the fake reviews and the attitude toward the fake reviews are getting stringent. For instance, in 2019, Federal Trade Commission (FTC) filed the first case against paid fake reviews, by CureEncapsulations on Amazon. Online platforms have restricted fake reviews in their own ways, but the regulators puts increasingly high pressure to online platforms to have more strict attitude against the fake reviews.<sup>3</sup>

However, the impact of the fake reviews on consumers on platform is not very clear. First, consumers might not be fooled by fake reviews if they know that there are fake reviews. In a standard work of Holmström (1999), the market can correctly anticipate the long-lived player’s behavior and debias the signal. Furthermore, customers might be able to elicit some extra information from fake reviews. If only high quality sellers make fake reviews to boost their initial reputation, the boosted rating can be an even better signal of good quality. Such a behavior might be possible if low quality will be revealed through the word-of-mouth and only high quality seller can reap benefits from the future sales as suggested by Nelson (1970;1974) in the context of advertising.<sup>4</sup>

In this paper, examine a theoretical model where the sales is determined by the seller’s reputation level and the seller chooses the amount of positive fake reviews at each moment. Consumers form a reputation based on the potentially boosted rating displayed on the platform. The platform can control how strictly it filters fake reviews and how much the rating reflects the information of the past feedback (ie, how fast the rating evolves). A key assumption in this paper is that it gets harder for the seller to make fake reviews when its reputation gets higher. In the main part, this is explained by the higher reimbursement to incentivized reviewers due to higher price.<sup>5</sup> This brings more fake reviews from the seller with low reputation. Furthermore, it is shown that such dependence will be greater with more stringent censoring policy (Proposition 3). In the dynamic model, fake reviews today have several impacts in the future: (i) increase of the future sales; (ii) increase of the cost of the fake reviews in the future; (iii) decrease in the number of the future fake reviews. This also feeds into the reaction of fake reviews to the seller’s type since the high type seller benefits more from the high reputation. As a result, the high quality type generates more fake reviews in the equilibrium. Because of this, the consumers sometimes benefits from lenient policy on fake reviews.

---

<sup>2</sup>Offers of such fake reviews from fake reviewers have been found on eBay.

<sup>3</sup>For instance, in 2019, the Competition and Markets Authority (CMA) in U.K. launched work programme “has written to Facebook and eBay this week urging them to conduct an urgent review of their sites to prevent fake and misleading online reviews from being bought and sold”. In responses, both Facebook and eBay have immediately deleted posts identified by CMA, and updated their policy to explicitly prohibit offers of fake reviews. In 2020, May, CMA has launched new investigation into online websites on how they currently detect fake or misleading reviews.

<sup>4</sup>Ananthakrishnan et al. (2020) analyze the display of fake reviews from a different perspective and show that the consumers form more trust on the platform if it shows the fake reviews with flags indicating them as fake reviews, rather than deleting them from the platform.

<sup>5</sup>We can see the interaction between fake reviews and reputation more commonly. For instance, fake reviews might be crowded out if the seller receives many organic feedback due to large demand caused by high reputation. Then, the effective fake review would be costly for such a seller.

The cause of the complementarity between quality and reputation is new in the literature. In this paper, the complementarity comes from a cost saving effect in the future rather than the increase of the revenue. This arises because the promotional behavior is explicitly modeled and plays a key role in contrast the previous research on the signaling promotion such as Kihlstrom and Riordan (1984) and Milgrom and Roberts (1986a).

The opposite dependence of fake reviews on the reputation and on the underlying true quality also gives some implication to the empirical literature on the signaling promotion: The reputation-based index such as the customer rating can be a bad proxy for the underlying quality. Researchers can estimate opposite results if they use the customer rating as a proxy for the quality. Furthermore, even if the true quality is measured somehow, it is important to control reputation level when we estimate the relationship between the promotional behavior and the underlying quality. Fig. 1 exemplify the possibility of the omitted variable issue: The level of promotion and the true quality can be negatively correlated without conditioning on the reputation level, even though the quality and the promotion has positive relationship, *ceteris paribus*.

The negative relationship between the fake reviews and the reputation level also increases the transition speed of the rating. When the rating goes down (up), it more quickly goes up (down) than the rating system without fake reviews. This distort the informativeness of the rating system. The transition speed of the rating is interpreted as a relative weight on the new information in the rating. Higher the weight on new information (and lower the weight on the old information), faster the transition of the rating. Thus, the equilibrium effect which makes the transition faster is distorting the weight on the new information upward (and the weight on the old information downward). Therefore, given the existence of the fake reviews, the platform should adjust it back. The platform should set a lower weight on new information (and higher weight on old information) compared to the rating system without fake reviews.

The above discussion is based on the assumption of the rational consumers, who know the seller's strategy. However, the regulator's concern is not necessarily such a sophisticated consumers, but more naive consumers, who is vulnerable to fake reviews.<sup>6</sup> In this paper, we also incorporate such consumers, and show how much they are biased due to the fake reviews by the sellers. Even though the relationship between the bias and the censorship policy is not monotone in general, the stringent censorship reduces the naive consumer's bias under a reasonable range of parameters.

Thus, the regulator might face a trade-off between the precision of the information for the rational consumers and the bias which the credulous consumers suffer from. This paper provides a framework to analyze such a trade-off.

The remaining of the paper is organized as follows. In Section 2, we review the related literature. In Section 3, we analyze a model with rational buyers. In section 4, we introduce credulous consumers. Section 5 concludes. Most of the proofs are deferred to the Appendix.

---

<sup>6</sup>For instance, Federal Trade Commission (FTC)'s mission is "[p]rotecting consumers and competition by preventing anticompetitive, deceptive, and unfair business practices through ...". (<https://www.ftc.gov/about-ftc>)

## 2 Literature Review

This paper mainly contributes two lines of the literature: rating design and signaling through promotion. The literature of the rating design can be divided to two strands: (i) how to reveal the known quality levels or an estimated quality index (i.e. whether to reveal the full information or add noise/coarsen the information), and (ii) how to make the index of unknown quality based on the multiple sources of information on the performance of the player.

The first strand is often framed in the context of certification such as Lizzeri (1999), Ostrovsky and Schwarz (2009), Boleslavsky and Cotton (2015), Harbaugh and Rasmusen (2018), Hopenhayn and Saeedi (2019), Hui et al. (2018). Some models are made tractable by the representation with posterior distribution in the line of Bayesian persuasion proliferated by Rayo and Segal (2010) and Kamenica and Gentzkow (2011). Saeedi and Shourideh (2020) extend the framework where the quality is endogenously chosen by the seller rather than the exogenous variable.

This paper fits in another strand analyzing how to aggregate the player's actions into a single index. In a one-shot model, Ball (2019) analyzes the optimal way to aggregate the various sources of potentially manipulated signals. In a dynamic setting based on Holmström (1999)'s signal jamming/career concern model, Hörner and Lambert (2018) show that the effort level of the long-lived player is maximized by a rating linear in the past observations. Vellodi (2020) analyzes the impact of the rating on the entry/exit behavior of the firm, and derive an optimal rating preventing high-quality seller exiting from the market due to reputation trap of failing to accumulate good reputation because of initial bad luck. Bonatti and Cisternas (2019) examine a long-lived consumer's Ratchet effect. The consumers try to hide its willingness to pay to avoid the personalized pricing by short-lived monopolist, so that the consumption does not perfectly reflect their willingness to pay. Like as Hörner and Lambert (2018) and Bonatti and Cisternas (2019), this paper examines a relationship between a signal-jamming structure and a linear rating system. In contrast to Hörner and Lambert (2018), the equilibrium strategy is dependent on the hidden quality and the reputation so that the seller's strategy changes the informativeness of the rating on equilibrium path as in Bonatti and Cisternas (2019).<sup>7</sup> In contrast to Bonatti and Cisternas (2019), where the effect of the manipulation is totally endogenously determined via the belief of the short-lived player, in this paper, the platform controls the effectiveness of the manipulation so that we can analyze the impact of the censorship by the platform. Besides, this paper depart from the literature by analyzing the impact of the manipulation on naive/credulous consumers, which is often the concern of the regulators.

This paper also contributes the literature on the promotion and signaling. Nelson (1970, 1974) argues that even if the promotion does not have any intrinsic information, the fact of burning money itself can be a signal of a good quality because such a signal pays off only for high quality type

---

<sup>7</sup>Another contrast to Hörner and Lambert (2018) is that they start from a general information structure so that they can represent any reputation by changing the information structure. Then, they can focus on the resulted process of reputation level in a similar way that researchers focus on the resulted outcome by the revelation principle in the context of the mechanism design. On the other hand, this paper and Bonatti and Cisternas (2019) use more specific information structure, so that we should examine how the consumers interpret the resulted rating.

through the repeated purchase in the future. This idea is formalized later by Kihlstrom and Riordan (1984), Milgrom and Roberts (1986a) and many others as separating equilibria in signaling models. Via a one-shot signal jamming framework rather than a signaling model, Mayzlin (2006) showed a negative relationship between a promotion through fake reviews and quality, and Dellarocas (2006) generalizes conditions for the positive/negative correlation in one-shot signal jamming model. Bar-isaac and Deb (2014) examine effects of vertically/horizontally heterogeneous preference, and Grunewald and Kräkel (2017) examine the effect of competition between firms. Most of the research on the signaling role of the promotion are based on models with one-shot promotion. An exception is Horstmann and MacDonald (1994), where the experience of the product is an imperfect signal of the quality and the signaling via the advertising is done only after establishing some reputation so that it is hard for the low-quality seller to mimic the high-quality seller’s behavior.<sup>8</sup> In this paper, I examine a dynamic signal-jamming model, where a reputational concern is a driving force of the positive correlation between the quality and the promotion. It also generate a non-degenerate dynamics consistent with an observation by Luca and Zervas (2016) that strategic manipulation increases after a drop of reputation.

The dependence of fake reviews on the reputation also gives some implication to the empirical literature on the signaling promotion. The literature have had a weak support on the correlation of the quality and the promotion. For instance, Kwoka (1984) observes that optometrists with more advertisement provide less thorough eye examination, and Horstmann and Moorthy (2003) observe that the advertising is hump-shaped in quality among restaurants in New York. Recently, Sahni and Nair (2019) implement a quasi-experiment to isolate the intrinsic information and the signaling effect of burning money and show that the consumer positively respond to the money burning. They pointed out that it is difficult to show the relationship between the quality and the promotion level since it is hard to obtain a reliable measure of quality. This paper emphasizes this point. The reputation-based index such as the customer rating can be a bad proxy for the underlying quality. The reputation level and the underlying quality level have the opposite impact on the promotion level in the equilibrium. Furthermore, even if the true quality is measured somehow, it is important to control reputation level. As shown in fig. 1, the level of promotion and the true quality can be negatively correlated without conditioning on the reputation level, even though the quality and the promotion has positive relationship, *ceteris paribus*.

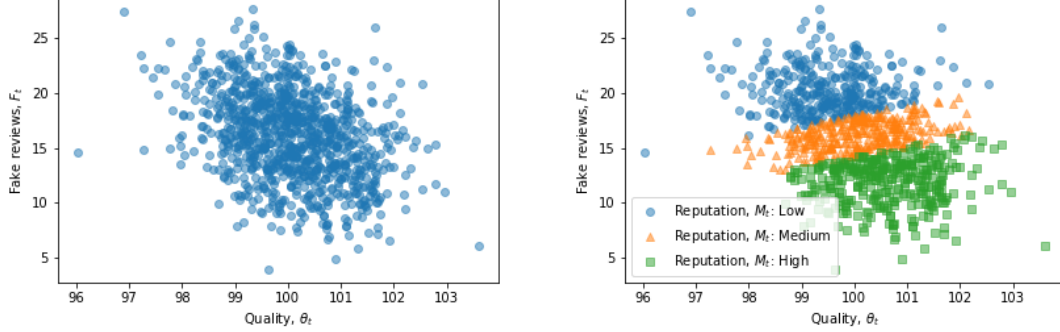
### 3 Rating Design for Rational Consumers

In this paper, we examine both models with rational consumers and with naive consumers. In this section, we first introduce a baseline model with a mass of rational consumers. The consumers

---

<sup>8</sup>Aside from the context of the rating system or the signaling promotion, Grugov and Troya-Martinez (2019) examine the biasing behavior of the seller in a model a. la. Holmström (1999) incorporating a detection rule and credulous consumers, and show that the biasing behavior increases as the authority requires stricter rule and the share of credulous consumer increases.

Figure 1: A simulated distribution of quality levels and the amount of fake reviews



The left panel shows that the amount of the fake reviews is negatively correlated with the quality level, unconditional on the level of reputation. On the other hand, the right panel shows that the amount of the fake reviews is increasing in the quality level, conditional on the reputation level.

rationally expect that the long-lived seller makes fake reviews following a linear strategy. However, they cannot induce the exact seller's action at time  $t$  since the quality is still hidden even though the strategy and the current reputation are known to the consumers.

Then, in the next section, we introduce a market with naive consumers, who does not expect any fake reviews, while the seller still makes fake reviews so the reputation is biased upward. In each model, we examine the impact of the platform's filtering/censoring policy on reviews, the weights on new and old reviews, and the precision of the genuine reviews.

### 3.1 Model

The model is in continuous time and infinite horizon,  $t \in [0, \infty)$ . At each instance  $t$ , a long-lived seller sells  $q$  units of its product, whose quality is denoted as  $\theta_t$ , and makes  $F_t$  units of fake reviews. A sufficiently large mass,  $n$ , of consumers forms a demand function such that the price  $p_t = E[\theta_t | Y_t] \equiv M_t$  clears the market, where  $Y_t$  is a rating of the product at time  $t$ .<sup>9</sup> The price being the reputation about the hidden quality is the standard assumption in the literature of the reputation. The quality  $\theta_t$  governs the willingness to pay to the product, so the price will be high when the expected quality of the product is high. A more specific underlying model, which can incorporate naive consumers is suggested in the Appendix.

The quality  $\theta_t$  and the rating  $Y_t$  change over time. The quality  $\theta_t$  follows an exogenous mean reverting process:

$$d\theta_t = (-\kappa\theta_t + \mu) dt + \sigma_\theta dZ_t^\theta \quad (1)$$

while the rating  $Y_t$  is characterized by the following differential equation:

$$dY_t = -\phi Y_t + d\xi_t \quad (2)$$

<sup>9</sup>Saeedi (2019) showed that the reputation is the measure determinant of the price on eBay market.

where  $d\xi_t$  is defined as

$$d\xi_t = aF_t dt + bq\theta_t dt + \sqrt{bq}\sigma_\xi dZ_t^\xi \quad (3)$$

where  $(Z_t^\theta, Z_t^\xi)$  is a standard Brownian motion,  $a$  is the effectiveness of the fake review,  $b$  is the feedback rate from customers,  $\mu$  is the mean of  $\theta_t$  in the stationary distribution,  $\sigma_\theta$  and  $\sigma_\xi$  govern the standard deviations of the disturbance. The exogenous mean-reverting process of  $\theta_t$  is understood as a result of the competition over quality among sellers. Relative quality of the own product might go down due to the rise other sellers with even higher quality. Own relative quality might go up when the competitor increase its own price. The transition of the rating  $Y_t$  is interpreted in discrete time analogue that the future rating  $Y_{t+dt}$  is a weighted sum of the new reviews  $d\xi_t$  and the previous reviews  $Y_t$  with weights of 1 and  $1 - \phi dt$ , respectively. The new reviews consist of two components: “organic” reviews and remaining fake reviews after filtering. The second and third term of eq. (3) corresponds to organic reviews. Higher quality tends to generate high reviews, and the information gets precise when there are many feedback from many transactions  $q$  or high response rate  $b$ . The disturbance  $\sigma_\xi dZ_t^\xi$  is caused by the heterogeneity of criteria among the customers.<sup>10</sup> The first term is the effect of the fake reviews. The seller tries to boost the average review through fake reviews, but some of them are detected by the platform and the remaining enters as  $aF_t dt$ . Thus, small  $a$  implies stringent censorship. As in Hörner and Lambert (2018); Vellodi (2020); Bonatti and Cisternas (2019), the rating  $Y_t$  is not exactly capturing 5-star rating on Amazon, Yelp, or some other online platforms. The mean of  $Y_t$  is determined by the mean of  $\theta_t$  and other parameters. By this specification of the rating, we can rely on the normality to simplify the analysis.

The seller’s instantaneous payoff is defined as:

$$\pi_t = (1 - \tau) p_t (q + F_t) - p_t \cdot F_t - \frac{c}{2} F_t^2$$

where  $\tau$  is transaction fees imposed by the platform. The first term is the total revenue from all the transactions including the own fake reviews, and the second term is the reimbursement cost to the fake reviewers. The last term is expressing that making more fake reviews is harder. The seller might find it difficult to search incentivized reviewers through communities such as Facebook, or some fake review services ask the seller a higher price for fake reviews because making more fake reviews comes with higher risk of being detected by the platform. The cost of the production is abstracted out from the model.<sup>11</sup> The long-lived seller maximizes its discounted present value by choosing  $(F_t)_{t \geq 0}$ .

The instantaneous profit becomes easier to compare with the previous research by rewriting it

---

<sup>10</sup>In this paper, the mechanism behind the customer feedback is abstracted and assumed that the fixed portion of consumers keep reviewing. For detailed analysis on the customer feedback, see Chevalier et al. (2018) and the literature cited in it. They analyze the relationship with managerial responses to reviews.

<sup>11</sup>Whether the high quality seller or low quality seller face high costs of production is arguable by itself. If high quality come from the seller’s high productivity, the high quality seller can produce with lower costs. If the low quality is by the seller’s choice rather than the difference in the production technology among sellers, the low quality product would be associated with low production cost. The different specifications on the production costs can cause different pattern in fake reviews, but those extensions are deferred to the future research.

as

$$\pi_t = (1 - \tau) M_t \cdot q - \tau M_t \cdot F_t - \frac{c}{2} F_t^2. \quad (4)$$

Without the second term in eq. (4), the model becomes effectively a special case of Hörner and Lambert (2018), which is based on Holmström (1999)'s signal jamming model and uses a general information structure as a rating. However, due to the existence of this term, the marginal cost of the manipulation depends on the current reputation level, so the equilibrium manipulation level depends on the current rating in contrast to Hörner and Lambert (2018) where the equilibrium action turns out to be state independent. In stead of relying on the time- and state-invariant action, we apply the idea of Bonatti and Cisternas (2019) to focus on a linear strategy and Gaussian stationary distribution of  $(\theta_t, Y_t)$ . Then, the HJB equation gives a simple the quadratic value function, solved by guess-and-verify method. It is verified that as  $\tau$  goes to zero, the equilibrium strategy becomes invariant to  $\theta_t, Y_t$ , (and  $t$ ).

The interaction between the current reputation and the current action is considered as a driving force of non-degenerate Markov equilibrium strategy. In this paper, this interaction between the reputation and the manipulation is derived from the reimbursement to the fake reviewers, however, such an interaction can be more commonly observed in the context of fake reviews. For instance, if the reputation is high, then large demand can crowd out fake reviews, so that the effective fake reviews can be more costly given the high reputation. In the appendix, an alternative model with such an interpretation is discussed. A model with the changing quantity which is isomorphic to the main model is discussed in the Appendix C.

**Definition of the Equilibrium** As mentioned above, we focus on a linear Markov strategy equilibria, where a linear Markov strategy maximizes the seller's discounted present value among any admissible strategies.

A linear strategy (in  $\theta_t$  and  $Y_t$ ) is defined as:

$$F_t = \hat{\alpha}\theta_t + \hat{\beta}Y_t + \hat{\gamma}$$

Note that  $\theta_t$  does not directly appear in the instantaneous payoff function, but it appears in the transition of the payoff relevant state variable,  $Y_t$ . Thus, the seller is potentially sensitive to the level of  $\theta_t$ . Now the equilibrium is defined as follows:

**Definition 1.** A linear Markov strategy  $F = (F_t)_{t \geq 0}$  s.t.  $F_t = \hat{\alpha}\theta_t + \hat{\beta}Y_t + \hat{\gamma}$  is a stationary Gaussian linear Markov equilibrium if

1.  $F = \arg \max_{(\tilde{F}_t)_{t \geq 0}} E_0 [\int_0^\infty e^{-tr} \pi_t]$  where  $(\tilde{F}_t)_{t \geq 0}$  is admissible,
2.  $M_t = E[\theta_t | Y_t]$ , and
3.  $(\theta_t, Y_t)_{t \geq 0}$  induced by  $F$  is stationary Gaussian



We don't know that  $(\theta_t, Y_t)_{t \geq 0}$  is stationary or Gaussian *ex ante* because  $Y_t$  is endogenously determined by  $F_t$ , however, the condition for  $(\theta_t, Y_t)_{t \geq 0}$  to be stationary Gaussian is simply characterized by an inequality, similarly to Bonatti and Cisternas (2019). By eqs. (2) and (3), and the definition of the linear strategy,

$$\begin{aligned} dY_t &= -\phi Y_t dt + a F_t dt + b q \theta_t dt + \sqrt{b q} \sigma_\xi dZ_t^\xi \\ &= -\left(\phi - a\hat{\beta}\right) Y_t dt + (a\hat{\alpha} + b q) \theta_t dt + a \delta \mu dt + \sqrt{b q} \sigma_\xi dZ_t^\xi \end{aligned} \quad (5)$$

Thus, an inequality  $\phi - a\hat{\beta} > 0$  must hold for the  $(\theta_t, Y_t)_{t \geq 0}$  to have the stationary distribution. (Otherwise, the process of  $Y_t$  diverges.) When  $(\theta_t, Y_t)$  is stationary Gaussian, by the projection theorem on the Gaussian distribution,

$$M_t \equiv E[\theta_t | Y_t] = E[\theta_t] + \frac{Cov(\theta_t, Y_t)}{Var(Y_t)} [Y_t - E[Y_t]] \quad (6)$$

. Furthermore, if it's stationary, all the expectations in eq.(6) are constants. By letting  $\lambda \equiv \frac{Cov(\theta_t, Y_t)}{Var(Y_t)}$  and  $\nu \equiv E[Y_t]$  (and  $\mu = E[\theta_t]$  by construction), eq.(6) is written as  $M_t = \mu + \lambda[Y_t - \nu]$ . In the following part of this section, we use  $M_t$  instead of  $Y_t$  as a state variable for the sake of expositional simplicity. Then, the linear strategy is redefined as

$$F_t = \alpha \theta_t + \beta M_t + \delta \mu$$

The stationary condition is summarized as follows:

**Lemma 1.** (*Stationality and the characterization of the long-run moments*) Suppose  $F_t = \alpha \theta_t + \beta M_t + \delta \mu$  where  $M_t \equiv E[\theta_t | Y_t]$  for all  $t \geq 0$ . Then, a process  $(\theta_t, Y_t)_{t \geq 0}$  is stationary Gaussian if and only if

- i.  $M_t = \mu + \lambda[Y_t - \nu]$  for all  $t$
- ii.  $a\lambda\beta - \phi < 0$ , and
- iii.  $(\theta_0, Y_0)' \sim \mathcal{N}([\mu, \nu]', \Gamma)$  is independent of  $(Z_t^\theta, Z_t^\xi)_{t \geq 0}$  where  $\Gamma$  is the variance-covariance matrix in the stationary distribution.

The third condition is required so that the game starts from stationary distribution. Now, HJB

equation is simply written by using Ito's lemma:

$$\begin{aligned}
rV(\theta, M) = & \sup_{F \in \mathbb{R}} (1 - \tau) M \cdot q - \tau M \cdot F - \frac{c}{2} F^2 \\
& - \kappa (\theta - \mu) V_\theta \\
& + \{a\lambda F + bq\lambda\theta - \phi [M - \bar{\theta} + \lambda\bar{Y}]\} V_M \\
& + \frac{\sigma_\theta^2}{2} V_{\theta\theta} \\
& + \frac{bq\lambda^2\sigma_\xi^2}{2} V_{MM}
\end{aligned} \tag{7}$$

By guessing the quadratic form of the value function,  $V = v_0 + v_1\theta + v_2M + v_3\theta^2 + v_4M^2 + v_5\theta M$ , and the linear strategy, we can verify the existence and uniqueness of the value function and the linear strategy via the matching coefficient.

### 3.2 Equilibrium Characterization

The equilibrium strategy is characterized by guessing the quadratic value function and the linear strategy and matching coefficients  $\alpha, \beta, \delta, (v_k)_{k=0}^5$  of the first order conditions, envelop conditions, and the stationarity condition characterized in Lemma 1. In the proof, the characterizing conditions are summarized into one equation  $h(L) = 0$  with an aggregator  $L \equiv a\lambda\beta$ , and then all the equilibrium coefficients are derived as a function of  $L$ . The aggregator  $L$  is interpreted as an equilibrium effect on the speed of the rating transition, or the equilibrium effect on the relative weight on new information. When  $L$  is positive, the rating transition effectively speeds up since the low rating is soon boosted back to the average rating by fake reviews.

By analyzing the existence and uniqueness of the aggregator  $L$  and examining the corresponding equilibrium coefficients, we obtain the following theorem:

**Theorem 1** (Existence and uniqueness). *There always exists a stationary linear Markov equilibrium. For any equilibrium,  $\alpha > 0$ ,  $\beta \in (-\frac{\tau}{c}, 0)$ ,  $\lambda > 0$  and  $L > 0$  hold. Furthermore, if  $h'(L) < 0$  holds, then such an equilibrium is unique and the equilibrium strategy  $\alpha, \beta, \delta$  is differentiable in parameters.*

*$h'(L) < 0$  holds for any  $L > 0$  if  $6\kappa\phi + 4r^2 + 2\kappa r + 17r\phi + 19\phi^2 > \kappa^2$ .*

Note that  $6\kappa\phi + 4r^2 + 2\kappa r + 17r\phi + 19\phi^2 > \kappa^2$  is a loose and reasonable condition.  $\phi$  is a transition speed of the rating and  $\kappa$  is the transition speed of the quality. The required inequality is reasonable as long as the rating system is meant to help estimating the current quality rather than the previous quality. For instance, even if the true quality does not drift much (ie,  $\kappa \simeq 0$ ), the rating should drift to reflect new information about the unknown quality (ie,  $\phi > 0$ ).

**Intuition of the Equilibrium Strategy** In Theorem 1, it is shown that high-quality type makes more fake reviews ( $\alpha > 0$ ), conditional on its reputation level. and high reputation type makes less fake reviews ( $\beta < 0$ ) conditional on the quality type. Given the logic of Nelson (1970; 1074),  $\alpha > 0$

(and  $\beta < 0$ ) might look intuitive but this model add slightly different reasons than the previous research.

I start from the negative  $\beta$ . From the first order condition, the optimal strategy is expressed as

$$F_t = -\frac{\tau}{c}M + a\lambda \underbrace{\{v_2 + 2M_tv_4 + \theta v_5\}}_{=V_M}$$

Then,  $\beta = -\frac{\tau}{c} + \frac{2a\lambda}{c}v_4$ , furthermore, the envelop condition gives an expression for  $v_4$  so that it is rewritten as  $\beta = -\frac{\tau}{c} - \frac{\tau}{c} \frac{a\beta\lambda}{(-a\beta\lambda+r+2\phi)}$ . The first term comes from the interaction of the reputation level and the fake reviews in the cost term,  $\tau M_t F_t$ . If the reputation is high, then the marginal cost of fake review is high. Therefore, the seller will make less fake reviews given the higher reputation. The second term is corresponding to the fake review's the marginal benefit in the future. Given the equilibrium strategy,  $v_4 = -\frac{\beta\tau}{2(-a\beta\lambda+r+2\phi)}$  is positive, meaning that the marginal benefit in the future is increasing in the reputation. This is because the future self will reduce the amount of the fake reviews after observing the boosted reputation due to today's fake reviews. Furthermore, this effect is increasing in  $M_t$ , because the future reputation  $M_{t+dt}$  tends to be high given high  $M_t$ , so the interaction term

$$\tau M_{t+dt} F_{t+dt} = \alpha \tau M_{t+dt} \theta_{t+dt} + \tau \beta M_{t+dt}^2 + \delta \mu \tau M_{t+dt} \quad (8)$$

reduces quadratically given negative  $\beta$ . It turns out that the first term dominates the second term so that  $\beta$  remains negative.

The intuition of the positive  $\alpha$  comes the complementarity between the quality  $\theta$  and the reputation  $M$  in the seller's value function. Given high quality  $\theta_t$  today, the reputation in the future tends to be higher than the case with low quality today, given the same level of the reputation  $M_t$  today. Furthermore, as shown in the above paragraph, the future benefit from the reputation boost is higher given higher reputation in the future. Thus, the high quality results in high incentive of fake reviews. Mathematically, the equilibrium coefficient  $\alpha$  is characterized as

$$\alpha = a\lambda v_5 = \frac{a\lambda}{\kappa + r + \phi} \{2(a\alpha + bq)\lambda v_4 - \alpha\tau\} \quad (9)$$

. The first equality captures that the sign of  $\alpha$  is coming from the complementarity of the  $\theta$  and  $M$  in the value function. In the last expression,  $(a\alpha + bq)\lambda$  represents that the high  $\theta_t$  results in high  $M_{t+dt}$ . It is multiplied with positive  $v_4$  representing increasing marginal value with respect to  $M_{t+dt}$ . This is the driving force of the positive  $\alpha$ . The remaining term of eq. (9),  $-\alpha\tau$ , states that such an incentive is attenuated because the quality in the near future  $\theta_{t+dt}$  tends to be high given high  $\theta_t$ , so today's fake review increases the cost in the future through via the first term of eq. (8).

In summary, the driving force of  $\beta < 0$  is the incentive to reduce  $\tau M_t F_t$  today given high  $M_t$ .  $\alpha$  is positive because of the complementarity of  $\theta_t$  and  $M_t$  through the cost saving. Readers might wonder why the increasing revenue (like as Nelson; 1970, 1974)) does not arise in the intuition. If

$\theta_t$  is high, the boosted revenue would stay high for a long time, but in this model, such a product would eventually achieve high reputation through the organic feedback even without fake reviews. Therefore, the *marginal future revenue*  $\frac{dp_s}{dF_t}$  ( $s \geq t$ ) is independent of  $\theta_t$ . It is worth noting that the same intuition applies even in a variant of the model with a fixed price  $p$  and time-varying quantity  $q_t$  discussed in the appendix.

### 3.2.1 Properties of the equilibrium

Before examining normative properties of the equilibrium, we check some positive properties of the equilibrium.

First, the expected amount of the fake reviews is increasing in  $a$ . This is simply because the marginal benefit of the fake review in the future would increase if the platform loosen the censorship policy. The model does not guarantee the positive amount of the fake review in general, but it is also shown that the expected amount of the fake review is positive under some parameters.

**Proposition 1.**  *$E[F_t]$  is increasing in  $L$ , furthermore, increasing in  $a$ . Furthermore,  $E[F_t] \geq 0$  holds for sufficiently large  $a$ .*

Thus, the model can represent reasonable situation under some parameters where the fake review has non-trivial effect (ie,  $a$  is significantly high). There still remains a small probability that  $F_t$  becomes negative due to the normal distribution, but the model can approximate a reasonable distribution of the fake reviews under which the negative revenue is rarely observed as shown in fig. 1.

The precision of “organic” feedback from normal customers also monotonically changes the expected amount of the fake reviews. When the organic feedback from customers varies a lot, it is hard for the seller to manipulate the reputation since a boosted rating is attributed to a large variation in the feedback.

**Proposition 2.**  *$E[F_t]$  is decreasing in  $\left(\frac{\sigma_\xi}{\sigma_\theta}\right)$ .*

Even though a stringent policy decrease the expected amount of the fake reviews as shown in Proposition 1, it does not imply that the seller’s strategy gets closer to no-fake strategy of  $\{\alpha, \beta, \delta\} = \{0, 0, 0\}$ . Moreover, the stringent policy might have unintentional effects of increasing the absolute value of the equilibrium coefficients.

**Proposition 3.**  *$|\alpha|$  is increasing in  $\frac{\tau}{c}$  and decreasing in  $\frac{\sigma_\xi}{\sigma_\theta}$ .  $|\beta|$  is decreasing in  $a$  and increasing in  $\left(\frac{\sigma_\xi}{\sigma_\theta}\right)$ .*

Under a stringent policy (small  $a$ ), the marginal benefit of the fake review decreases since the fake reviews are reflected less in the rating, but at the same time, the dependence of the marginal benefit on the current reputation also decreases. Mathematically, the second term of  $\beta = -\frac{\tau}{c} + \frac{\tau}{c} \frac{-a\beta\lambda}{(-a\beta\lambda + r + 2\phi)}$  decreases while the marginal cost still depends on the current reputation regardless of the censoring policy. Therefore,  $|\beta|$  gets larger due to the less countervailing effect.

In the proof of the proposition, the intensity of dynamic consideration is also captured by an aggregator  $L = -a\lambda\beta$ , which is the equilibrium effect on the reputation transition speed.  $L$  becomes smaller when the dynamic incentive becomes smaller, so the  $\alpha$  which only comes from the future marginal benefit becomes smaller, and  $|\beta|$ , to which the future marginal benefit only works as a counteracting effect, becomes greater because the present cost reduction incentive prevails.  $L$  is shown to be increasing in  $\frac{a\tau}{c}$  and decreasing in  $\frac{\sigma_\xi}{\sigma_\theta}$ .

**Lemma 2.**  *$L$  at the equilibrium is increasing in  $\frac{a\tau}{c}$  and decreasing in  $\frac{\sigma_\xi}{\sigma_\theta}$ . Furthermore,  $L \rightarrow 0$  as  $\frac{a\tau}{c} \rightarrow 0$  and  $L \rightarrow \infty$  as  $\frac{a\tau}{c} \rightarrow \infty$ .*

This concludes Proposition 3.  $\alpha$  is not necessarily increasing in  $a$  because  $\alpha$  is a function in  $a$  and  $L$ , so the change in  $a$  impact directly, and indirectly via  $L$ , and the net impact is not clear.  $|\beta|$  is not necessarily decreasing in  $\frac{\tau}{c}$  by the analogous reason even though a limit of  $\tau \rightarrow 0$  is known.

Proposition 3 implies less signaling (smaller  $\alpha$ ) and more distortion in the effective transition speed of the rating (greater  $|\beta|$ ) when the aggregator on the strategic effect  $L$  is small. This suggests less information from the rating system when the strategic effect  $L$  is small. In the following section, we formally examine this effect.

Some limits of the equilibrium strategy are worth noting before jumping into normative analysis. Since the negative  $\beta$  is coming from the interaction term in the cost of the fake reviews, whose coefficient is  $\tau$ ,  $\beta$  goes zero as  $\tau$  goes to zero. At the same time,  $\alpha$  also goes to zero because the complementarity of  $\theta$  and  $M$  is caused by the future cost saving via negative  $\beta$ . In this limit, the fake reviews becomes constant a. la. Holmström (1999). This is summarized in the following proposition.

**Proposition 4.**  $|\alpha|, |\beta| \rightarrow 0$  as  $\tau \rightarrow 0$ .

### 3.3 Optimal Rating System for Rational Consumers

In this paper, we focus on the informativeness of the rating system as a normative criteria. There are two reason to do this. First, from the viewpoint of consumer protection: as the rating system gets more informative about the quality of the product, the price is likely to be close to the underlying quality. Thus, it becomes less likely that consumer faces a huge regret from the purchase of the product. Second, from the viewpoint of the platform, the informativeness of the rating is crucial to attract consumers in the long-run. If consumers find it not informative, they can move to other platform. Then, so does sellers in the platform given less consumers in the market. So the informativeness of the rating would be the first priority when the platform controls it.

Since the rational customers can form an unbiased estimate from any current rating,  $M_t = E[\theta_t|Y_t]$ , the informativeness of the rating is defined by the variance of the customer's estimate on the quality. Due to the normality assumption, This is rewritten as  $Var(\theta_t|Y_t) = Var(\theta_t)(1 - \rho^2)$  where  $\rho^2$  is the correlation between  $\theta_t$  and  $Y_t$ . Therefore, we use  $\rho^2$  as the criteria on the informativeness of the rating.

Given an equilibrium strategy, the stochastic differential equation eqs. (1) and (5) gives us  $\rho^2$  as a function of the parameters and the equilibrium strategy. Therefore, change of a parameter directly affects  $\rho^2$  and indirectly affects through change of the equilibrium strategy. Fortunately, by representing equilibrium coefficients  $\alpha, \beta$  as a function of the equilibrium aggregator  $L = a\beta\lambda$ , all the direct and indirect effects of the censorship ( $a$ ) are expressed as an effect through  $L$ . Comparative statics about other parameters such as  $\phi$  and  $\sigma_\xi/\sigma_\theta$  can be also examined by the indirect effect through  $L$  and the direct effect.

**Lemma 3.** *At the equilibrium,  $\rho^2$  is expressed as a function:*

$$\rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta, r, bq) = \frac{(\phi + L)}{(\kappa + \phi + L)} \frac{(A(L; \phi, \kappa, r, bq) + 1)^2}{((A(L; \phi, \kappa, r, bq) + 1)^2 + \kappa(\sigma_\xi/\sigma_\theta)^2(\kappa + \phi L))}$$

on which  $a, c, \tau$  affect only through  $L$ .

$A(L; \phi, \kappa, r, bq)$  summarizes all the direct and the indirect effect of  $a$  on the informativeness as a function of  $L$ .

### 3.3.1 Filtering/Censoring Reviews

First, we analyze the impact of filtering/censoring policy,  $a$ . Does the fake reviews damage the informativeness of the rating system compared to the case without fake reviews? Does the filtering/censoring the reviews (ie, decrease of  $a$ ) increase the informativeness?

As a benchmark, we derive the informativeness *without* fake reviews. By construction, we can do this by letting  $\alpha = \beta = \delta = 0$ .<sup>12</sup> The same informativeness is also replicated by setting  $L = 0$  in  $\rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta)$ , so that it becomes easier to compare with the informativeness at the equilibrium.

**Lemma 4.**  $\rho^2(0; \phi, \kappa, \sigma_\xi, \sigma_\theta)$  coincides with  $\rho^2$  under the no-fake strategy.

Note that  $L = 0$  does not necessarily mean  $\alpha = \beta = \delta = 0$ . For instance,  $L$  goes to 0 as  $a$  goes to 0, but at the same time,  $\beta$  converges to some negative value. The lemma says that even under such a situation, the informativeness is the same as the one without fake reviews. By combining with Lemma 2 about the relationship between  $L$  and parameters, this gives us the following proposition:

**Proposition 5.** *The informativeness of the rating system in the equilibrium converges to that of the “no-fake” strategy as  $\frac{a\tau}{c} \rightarrow 0$ .*

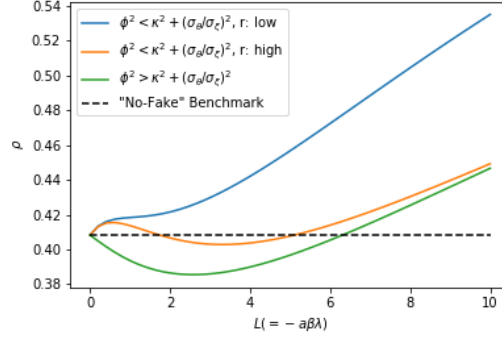
Thus, even though the equilibrium strategy at the limit of  $\frac{a\tau}{c}$  is not necessarily the no-fake strategy, the informativeness converges to that of no-fake strategy.

Now, by analyzing the behavior of  $\rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta)$  with respect to  $L$ , we can conclude that the informativeness can be even higher under some parameters where the positive amount of the fake reviews are expected. In other words, the stringent censorship can decrease the informativeness of the rating system.

---

<sup>12</sup>Actually,  $\delta$  does not enter in the formula for the informativeness, so  $\delta = 0$  does not matter in terms of the informativeness.

Figure 2: Change of the informativeness in the aggregator  $L$



The graph indicates that the informativeness is (i) increasing in  $L$  if  $\phi$  and  $r$  are relatively low, (ii) increasing in  $L$  around zero, then decreasing, and then increasing if  $\phi$  is relatively low but  $r$  is relatively high, and (iii) decreasing in  $L$  around zero and then increasing in  $L$  if  $\phi$  is relatively high. It also indicates the rating becomes more informative than the no-fake benchmark as  $L$  gets large.

**Proposition 6.** *The equilibrium strategy is more informative than no-fake strategy under a set of parameters such that*

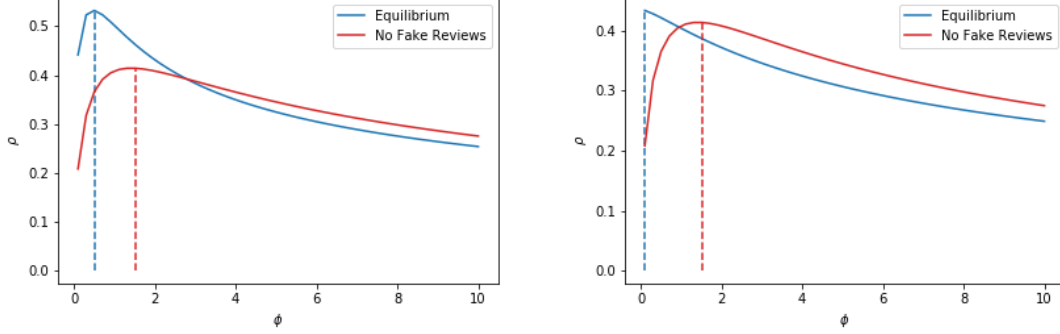
1.  $\frac{a\tau}{c}$  is sufficiently large, or
2.  $\frac{a\tau}{c}$  is sufficiently small and  $\phi^2 < \kappa^2 + \frac{\sigma_\theta^2}{\sigma_\xi^2}$

Fig. 2 exhibits the behavior of  $\rho^2$  with respect to  $L$ . The first part of the proposition comes from the fact that  $\rho^2$  converges to 1 as  $L$  goes to infinite. Thus, it surpasses  $\rho^2$  of the no-fake benchmark at some point. Since  $L$  is increasing  $\frac{a\tau}{c}$  from zero to infinite, the equilibrium surpasses  $\rho^2$  of the no-fake benchmark at some point as  $\frac{a\tau}{c}$  increases. The second part is derived from the behavior of  $\rho^2$  around  $L = 0$ . The derivative of  $\rho^2$  w.r.t.  $L$  is determined by the relative size of  $\phi^2$  and  $(\kappa^2 + \sigma_\theta^2/\sigma_\xi^2)$ : If  $\phi^2 < \kappa^2 + \frac{\sigma_\theta^2}{\sigma_\xi^2}$ , then  $\rho^2$  is decreasing in  $L$ , so is decreasing in  $\frac{a\tau}{c}$ , *vice versa*.<sup>13</sup>

Intuition of the proposition consists of two parts: (i) As mentioned in Subsection 3.2.1, the fake review's sensitivity to  $\theta_t$  decreases as the strategic effect  $L$  decreases. Thus, the strict censorship policy, which reduces the equilibrium aggregator  $L$ , decreases signaling effect of the fake reviews. (ii) On the other hand,  $L > 0$  increases the effective transition speed of reputation to  $\phi + L$ . It can be good or bad in terms of the informativeness, depending on the original transition speed,  $\phi$ . More specifically, the threshold of  $\sqrt{\kappa^2 + \sigma_\theta^2/\sigma_\xi^2} \equiv \phi^0$  is the informativeness-maximizing  $\phi$ , given no fake reviews. Therefore, if  $\phi$  is smaller than  $\phi^0$ , the faster transition benefits the informativeness, and vice versa. It turns out that the first effect dominates for large  $L$  and the second effect dominates for  $L$  close to zero.

<sup>13</sup>Note that  $E[F_t]$  is increasing in  $L$  and positive for large  $L$  (by Proposition 1). Thus, the high informativeness is not due to negative fake reviews, but associated with the positive amount of fake reviews.

Figure 3: Change of the informativeness in  $\phi$



The left panel shows change of the informativeness in  $\phi$  when  $r$  is relatively low, while the right panel shows that of a relatively high  $r$ . The informativeness is maximized at a lower  $\phi$  under the equilibrium than the maximizer under the no-fake benchmark.

### 3.3.2 Weights on New/Previous Reviews

Next, we analyze the optimal weights on the new and old reviews. Again, the informativeness without fake reviews is expressed by  $\rho^2(0; \phi, \kappa, \sigma_\xi, \sigma_\theta)$ . Therefore, the optimal weight at this benchmark is simply characterized by  $\frac{\partial \rho^2}{\partial \phi}(0; \phi, \kappa, \sigma_\xi, \sigma_\theta) = 0$ . Let  $\phi^0$  be the solution of this equation. On the other hand, at the equilibrium,  $\phi$  changes the equilibrium aggregator  $L$  as well. Thus, the optimal weight at the equilibrium is characterized by  $\frac{d\rho^2}{d\phi} = \frac{\partial \rho^2}{\partial \phi}(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) + \frac{\partial \rho^2}{\partial L}(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) \frac{dL}{d\phi} = 0$ . Let the solution of this equation be  $\phi^*$ . Now, we have the following proposition.<sup>14</sup>

**Proposition 7.**  $\frac{d\rho^2}{d\phi} < 0$  at  $\phi = \phi^0$ . Furthermore, if  $r$  is sufficiently small, then  $\rho^2(L(\phi^*); \phi^*, \kappa, \sigma_\xi, \sigma_\theta) > \rho^2(0; \phi^0, \kappa, \sigma_\xi, \sigma_\theta)$ .

The first part of the proposition states that the platform should reduce the speed of transition  $\phi$ , given the existence of the fake reviews. Intuitively, this is explained as follows: At the equilibrium, the transition of the rating score  $Y_t$  is  $\phi + L$  where  $L$  is non-negative. Therefore, to cancel the strategic impact on the transition speed, the platform should decrease  $\phi$  compared to the no-fake benchmark,  $\phi^0$ . Again, the transition speed is interpreted as a relative weight of the new information. At the equilibrium, the number of fake reviews is decreasing in the current rating, so the fake reviews are canceling the past performance. In other words, the new information is effectively weighted more than the platform intends. Thus, the platform can increase the informativeness by adjusting it downward.

The second part of the proposition is even more striking. If the seller is sufficiently concerned about the future, the platform can achieve higher informativeness than no-fake review by adjusting the update speed of the rating. The implication looks similar to Proposition 5, but is slightly different from that. The right panel of fig. 3 exemplify that the equilibrium informativeness is

<sup>14</sup> $\phi$  corresponding to disaggregated information,  $\phi^d$ , is an alternative benchmark as in Bonatti and Cisternas (2019). In this model, we obtain a mixed result for the comparison of  $\phi^*$  and  $\phi^d$ . See the appendix for more details.



greater than the informativeness without fake reviews under some parameters (say  $\phi = 0.9$ ) as shown in Proposition 5, but it can be still lower than the maximum of the informativeness without fake reviews (maximized around  $\phi = 1.6$ ). The second part of Proposition 6 states that, even when we compare maximum of informativeness of the rating with and without fake reviews, one with fake reviews will be higher if the seller cares about the future enough as shown in the left panel of fig. 3..

### 3.3.3 Precision of the genuine reviews

Lastly, we examine the impact of the precision of the organic feedback,  $\frac{\sigma_\xi}{\sigma_\theta}$ . As discussed in Subsection 3.2.1, increasing  $\frac{\sigma_\xi}{\sigma_\theta}$  and decreasing  $a$  have the similar effects on the equilibrium strategy. However, they depart in terms of the impacts on the informativeness. This is because  $a$  affects the informativeness only through the equilibrium aggregator  $L$ , but  $\frac{\sigma_\xi}{\sigma_\theta}$  affects the informativeness directly as well. Intuitively, if the reviews consists of the less precise feedback (ie, higher  $\frac{\sigma_\xi}{\sigma_\theta}$ ), the rating score is less informative about the quality by definition. The indirect effect consists of two parts like as the comparative statics over  $a$ : (i) Higher  $\frac{\sigma_\xi}{\sigma_\theta}$  implies smaller strategic effect  $L$ , which then implies less signaling effect. (ii)  $L > 0$  effectively increases the rating transition to  $\phi + L$ . The following proposition shows that the direct effect and the first indirect effect dominate the second indirect effect for any parameters.

**Proposition 8.** *The informativeness at the equilibrium is decreasing in  $\frac{\sigma_\xi}{\sigma_\theta}$ .*

Thus, the precise organic feedback increases the informativeness even though it comes with more fake reviews.

## 4 Rating Design for Naive Consumers

The model with rational consumers is a standard starting point for any economic models, but in the context of customer reviews, it is natural to consider the impact on naive consumers who don't expect any fake reviews. The regulator often tries to protect customers from the fake reviews with an assumption that the fake reviews can fool or at least confuse consumers. In this section, we assume that some consumers do not expect any fake reviews on the platform. Such customers are modeled by assuming that the reputation (and the price) is characterized as  $\widetilde{M}_t = \mu + \widetilde{\lambda}[Y_t - \widetilde{v}]$  where  $\widetilde{\lambda}$  and  $\widetilde{v}$  are characterized by the stochastic differential equation eqs. (1) and (5) where  $\alpha = \beta = \delta = a = 0$ . On the other hand, the long-lived seller face the same problem as in the previous chapter except for the definition  $p_t$ .<sup>15</sup>

<sup>15</sup>**Note to be added:** Similarity to Milgrom and Roberts (1986b) RAND “Relying on the Information of Interested Parties”]

#### 4.1 Model / Equilibrium Characterization

In this section, the price is assumed to be a convex combination of a rational reputation  $M$  and a naive reputation  $\widetilde{M}$ .

$$\begin{aligned} p &= \eta M + (1 - \eta) \widetilde{M} \\ &= \eta \{\mu + \lambda [Y_t - \nu]\} + (1 - \eta) \{\mu + \lambda^{naive} [Y_t - \nu^{naive}]\} \\ &= \mu - (\eta\lambda\nu + (1 - \eta)\lambda^{naive}\nu^{naive}) + (\eta\lambda + (1 - \eta)\lambda^{naive}) Y_t \end{aligned}$$

One interpretation is that each consumer can be partially rational, so their expectation about the quality of the product is somewhere in between the totally sophisticated expectation and the totally naive expectation. Thus, the rationality of each consumer is captured by  $\eta$ .

Another interpretation is that  $\eta$  is a ratio of the rational consumers among all the consumers, so the market price is set somewhere in between the rational expectation and the naive expectation. When the ratio of rational consumer goes up, it converges to the rational expectation. The linear specification captures such a relationship in a simple manner, and can be founded by a specific utility function of the buyer. Suppose that there are  $n$  consumers in the market and  $\eta \cdot n$  of them are rational and the other  $(1 - \eta) \cdot n$  are naive. Consumer  $i \in [0, n]$  feels  $u_{t,i} = \theta_t + \epsilon_{t,i} - p_t$  if the consumer buy the product, and 0 otherwise, where  $\epsilon_{t,i}$  is identically and independently distributed. The rational consumers and the naive consumers differ only in terms of how they form the expectation based on the same observation of the rating  $Y_t$ . Conditional on  $Y_t$ , a rational consumer purchases the product if and only if  $M_t + \epsilon_i - p \geq 0$ , while a naive consumer purchase it if and only if  $\widetilde{M}_t + \epsilon_i - p \geq 0$ . Therefore, the total demand function is expressed as

$$\eta \cdot n \cdot (1 - F(p - M)) + (1 - \eta) \cdot n \cdot (1 - F(p - \widetilde{M}))$$

where  $F(p)$  is a c.d.f. of the random variable  $\epsilon_i$ . By letting  $n = 2q$  and assuming that  $\epsilon_i$  is distributed uniformly and symmetrically around zero. We obtain  $p = \eta M + (1 - \eta) \widetilde{M}$  to clear the market.

In this section, we consider a linear strategy  $F_t = \hat{\alpha}\theta_t + \hat{\beta}Y_t + \hat{\gamma}$  and the HJB equation with state variables  $\theta$  and  $Y$ , since  $Y$  keeps track of both  $M$  and  $\widetilde{M}$  in a simple manner:

$$\begin{aligned}
rV(\theta, Y) = & \sup_{F \in \mathbb{R}} (1 - \tau) p \cdot q - \tau p \cdot F - \frac{c}{2} F^2 \\
& - \kappa(\theta - \mu) V_\theta \\
& + \{-\phi Y_t + a F_t dt + b q \theta_t\} V_Y \\
& + \frac{\sigma_\theta^2}{2} V_{\theta\theta} \\
& + \frac{b^2 q^2 \sigma_\xi^2}{2} V_{YY}
\end{aligned} \tag{10}$$

The following theorem states that, even with credulous consumers, we have the existence and the uniqueness given the same condition as the baseline model.

**Theorem 2.** *For any  $\eta \in [0, 1]$ , there always exists a stationary linear Markov equilibrium. For any equilibrium,  $\alpha > 0$ ,  $\beta \in (-\frac{\tau}{c}, 0)$ ,  $\lambda > 0$  and  $L > 0$  hold. Furthermore, if  $h'(L) < 0$  holds, then such an equilibrium is unique and the equilibrium strategy  $\alpha, \beta, \delta$  is differentiable in parameters.*

*$h'(L) < 0$  holds for any  $L > 0$  if  $6\kappa\phi + 4r^2 + 2\kappa r + 17r\phi + 19\phi^2 > \kappa^2$ .*

In addition, somehow surprisingly, the existence of the naive consumers reduce the sellers strategic behavior.

**Proposition 9.** *The equilibrium with naive consumer ( $\eta \in [0, 1)$ ) generates smaller  $|\alpha|$ , larger  $|\beta|$ , and smaller  $E[F_t]$  compared to the equilibrium without naive consumers ( $\eta = 1$ ).*

This is because the rational consumers are more sensitive to the change of ratings compared to the rational consumers. Rational consumers know that the rating is boosted, but they also know that the rating is boosted more by high quality type. Therefore, rational consumers attribute the boosted rating to high quality type, and set high price for such a boosted rating. Therefore, with naive consumers, the price is less responsive to the boost of the ratings since the naive consumers do not know such a strategic correlation between the quality and the rating, so the *marginal* benefit of the fake review for the seller is smaller. Thus, with naive consumers, the seller makes less fake reviews in expectation.

Readers might wonder why the seller does not become more exploitative with naive consumers. It is simply because the fake review strategy against rational consumers generates more fake reviews by the different reasons than exploiting consumers. If only a small number of naive consumers exist and observe the ratings, the naive consumers would form even more biased estimate since the seller makes more fake reviews trying to send a good signal to rational consumers.

## 4.2 Optimal Rating System for Naive Consumers

**Criteria: Bias in the Reputation** In this part we evaluate the impact of fake reviews on naive consumers. To do so, we introduce a bias in naive consumer's expectation caused by the boosted

rating:

$$\begin{aligned}
Bias &\equiv E \left[ \widetilde{M}_t - \theta_t \right] \\
&= E \left[ \mu - \theta_t + \widetilde{\lambda} [Y_t - \widetilde{\nu}] \right] \\
&= \widetilde{\lambda} [\nu - \widetilde{\nu}]
\end{aligned}$$

where  $\widetilde{\lambda}$  is the sensitivity of the reputation to the rating, and  $\nu$  and  $\widetilde{\nu}$  are the actual mean of the rating and the estimate of mean of the rating by the naive consumers. The above decomposition of the bias is intuitive: the positive bias is due to the boosted reputation. Since the naive consumers do not expect any fake reviews, they interpret a high rating ( $Y_t > \widetilde{\nu}$ ) as a result of the high quality even though it is actually the average level of the rating at the equilibrium ( $Y_t = \nu > \widetilde{\nu}$ ).

Therefore, as long as the seller makes the positive amount of fake reviews (in expectation) to boost the rating, the naive consumers are positively biased. This intuition is verified in the following lemma.

**Lemma 5.** *Bias  $\geq 0$  if and only if  $E[F_t] \geq 0$ .*

#### 4.2.1 Filtering/Censoring Reviews

In the following part, for the sake of the tractability, I focus on the case of  $\eta = 0$ , where the only naive consumers exist in the market. Numerical exercises in cases of  $\eta \in (0, 1)$  can be found in the Appendix.

First, we examine the impact of a filtering policy, which regulators are arguably concerned about the most. The following proposition provides a theoretical background of a stringent policy protecting the naive customers. Note that even though the statement seems pretty intuitive, it is not trivial since the model predicts non-monotone relationship between the censorship and the bias in general. Fortunately, in a realistic range of parameters where naive consumers suffer from the positive bias in their reputation, the stringent censorship will reduce such bias.

**Proposition 10.** *Suppose  $Bias \geq 0$ . Then,  $Bias$  is increasing in  $a$ .*

Combining with Lemma 5, the condition for the stringent policy to work for the naive consumers is translated as the condition of a measure observable by the platform.

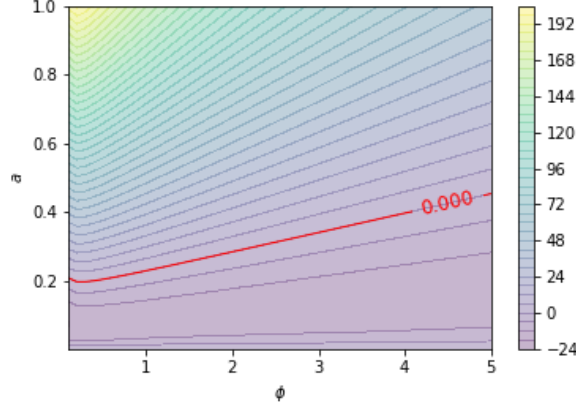
**Corollary 1.** *The stringent censorship reduces the naive consumers' bias whenever the expected amount of the fake reviews is positive.*

Thus, as long as the positive fake reviews are observed, the stringent policy is beneficial for naive consumers even though it can reduce the informativeness of the rating for the rational consumers.

#### 4.2.2 Weights on New/Previous Reviews

As shown in fig. 4, the bias tends to be hump shaped in  $\phi$ . This is intuitive because the fake reviews would be effective only when the rating is believed to be informative by the consumers so that the

Figure 4: Impact of censorship intensity and the weights of reviews on naive consumer's bias.



consumers react to the rating. Since the informativeness is hump shaped in  $\phi$ , so is the bias caused by the fake reviews. This emphasizes that the trade off between the bias and the informativeness can be an inherent feature of the fake reviews.

Some readers might want an integrated criteria of the bias and the informativeness. The mean squared error (MSE) is a natural candidate. It does not give us a clear-cut prediction, but a simulation about MSE is provided in the Appendix.

## 5 Summary

In this paper, effects of fake reviews on rational and credulous consumers are analyzed. The key assumption is that the high reputation causes high cost of the fake reviews. This is rationalized by the high reimbursement to reviewers or high demand and many authentic feedback crowding out the fake reviews.

At the equilibrium, the amount of the fake reviews is increasing in quality and decreasing in reputation level, which imply difficulties of empirical analysis on the signaling promotion. The stringent censorship reduces the expected amount of fake reviews while it can decrease the signaling effect and increase the effective transition speed of the rating.

This leads to a normative result that the rating under a less strict filtering policy can be more informative compared to the rating under a strict policy or the rating with no fake reviews. In terms of the weights on new and old information in the rating, the platform should reduce the weight on the new information to maximize the informativeness of the rating compared to one without fake reviews. Since the fake reviews attenuate the old information and increase the relative weight of the new information effectively, the platform should adjust it back.

The existence of credulous consumers decreases the expected amount of fake reviews since they are less responsive to the rating without taking the positive relationship between the fake review and the quality. On the other hand, they are vulnerable to fake reviews and pay more than the true quality in expectation. The model predict that, as long as the positive amount of the fake reviews

are observed, the regulator or the platform can reduce such a biased behavior by enhancing the censorship.

The results emphasizes that the regulator or the platform would face a trade-off between the precision of the informativeness and the bias caused by the fake reviews. As long as the rating is considered informative, the incentive to make fake reviews arise.

## References

- Ananthakrishnan, Uttara M., Beibei Li, and Michael D. Smith**, “A Tangled Web: Should Online Review Portals Display Fraudulent Reviews?,” *Information Systems Research*, jun 2020, *Article in*, 1–22.
- Ball, Ian**, “Scoring Strategic Agents,” 2019, (November), 1–57.
- Bar-isaac, Heski and Joyee Deb**, “What is a Good Reputation? Career Concerns with Heterogeneous Audiences,” *International Journal of Industrial Organization*, 2014, *34*, 44–50.
- Belleflamme, Paul and Martin Peitz**, “Inside the Engine Room of Digital Platform: Reviews, Ratings, and Recommendations,” in “Economic Analysis of the Digital Revolution” 2018, pp. 75–114.
- Boleslavsky, Raphael and Christopher Cotton**, “Grading standards and education quality,” *American Economic Journal: Microeconomics*, 2015, *7* (2), 248–279.
- Bonatti, Alessandro and Gonzalo Cisternas**, “Consumer Scores and Price Discrimination,” *Review of Economic Studies*, 2019.
- Chevalier, Judith A., Yaniv Dover, and Dina Mayzlin**, “Channels of Impact: User Reviews When Quality Is Dynamic and Managers Respond,” *Marketing Science*, 2018, *37* (5), 688–709.
- Dellarocas, Chrysanthos**, “Strategic Manipulation of Internet Opinion Forums: Implications for Consumers and Firms,” *Management Science*, 2006, *52* (10), 1577–1593.
- Dwoskin, Elizabeth and Craig Timberg**, “How merchants use Facebook to flood Amazon with fake reviews,” 2018.
- Grugov, Mikhail and Marta Troya-Martinez**, “Vague Lies and Lax Standards of Proof: On the Law and Economics,” *Journal of Economics & Management Strategy*, 2019, *28*, 298–315.
- Grunewald, Andreas and Matthias Kräkel**, “Advertising as signal jamming,” *International Journal of Industrial Organization*, 2017, *55*, 91–113.
- Harbaugh, Rick and Eric Rasmusen**, “Coarse Grades: Informing the Public by Withholding Information,” *American Economic Journal: Microeconomics*, 2018, *10* (1), 210–235.

- He, Sherry, Brett Hollenbeck, and Davide Proserpio**, “The Market for Fake Reviews,” *SSRN Electronic Journal*, 2020, pp. 1–38.
- Hollenbeck, Brett**, “Online Reputation Mechanisms and the Decreasing Value of Chain Affiliation,” *Journal of Marketing Research*, oct 2018, *55* (5), 636–654.
- , **Sridhar Moorthy, and Davide Proserpio**, “Advertising Strategy in the Presence of Reviews: An Empirical Analysis,” *Marketing Science*, 2019, *38* (5), 793–811.
- Holmström, Bengt**, “Managerial Incentive Problems: A Dynamic Perspective,” *Review of Economic Studies*, 1999, *66* (1), 169–182.
- Hopenhayn, Hugo and Maryam Saeedi**, “Optimal Ratings and Market Outcomes,” *NBER Working Paper Series*, 2019, pp. 1–39.
- Hörner, Johannes and Nicolas Lambert**, “Motivational Ratings,” *Review of Economic Studies*, 2018, *forthcomin*.
- Horstmann, Ignatius J. and Glen M. MacDonald**, “When is Advertisign a Signal of Product Quality?,” *Journal of Economics and Management Strategy*, 1994, *3* (3), 561–584.
- **and Sridhar Moorthy**, “Advertising Spending and Quality for Services: The Role of Capacity,” *Quantitative Marketing and Economics*, 2003, *1* (3), 337–365.
- Hui, Xiang, Maryam Saeedi, Giancarlo Spagnolo, and Steve Tadelis**, “Certification, Reputation and Entry: An Empirical Analysis,” 2018, pp. 1–59.
- Kamenica, Emir and Matthew Gentzkow**, “Bayesian Persuasion,” *American Economic Review*, 2011, *101* (6), 2590–2615.
- Kihlstrom, Richard E. and Michael H. Riordan**, “Advertising as a Signal,” *Journal of Political Economy*, 1984, *92* (3), 427–450.
- Lizzeri, Alessandro**, “Information Revelation and Certification Intermediaries,” *RAND Journal of Economics*, 1999, *30* (2), 214–231.
- Luca, Michael and Georgios Zervas**, “Fake It Till You Make It: Reputation, Competition, and Yelp Review Fraud,” *Management Science*, 2016, *62* (12), 3412–3427.
- Mayzlin, Dina**, “Promotional Chat on the Internet,” *Marketing Science*, 2006, *25* (2), 155–163.
- Milgrom, Paul and John Roberts**, “Price and Advertising Signals of Product Quality,” *Journal of Political Economy*, 1986, *94* (4), 796–821.
- **and —**, “Relying on the Information of Interested Parties,” *RAND Journal of Economics*, 1986, *17* (1), 18–32.

- Ostrovsky, Michael and Michael Schwarz**, “Information disclosure and unraveling in matching markets,” *American Economic Journal: Microeconomics*, 2009, 2 (2), 34–63.
- Rayo, Luis and Ilya Segal**, “Optimal Information Disclosure,” *Journal of Political Economy*, 2010, 118 (5), 949–987.
- Reimers, Imke C. and Joel Waldfogel**, “Digitization and Pre-Purchase Information: The Causal and Welfare Impacts of Reviews and Crowd Ratings,” *National Bureau of Economic Research*, 2020.
- Saeedi, Maryam**, “Reputation and Adverse Selection: Theory and Evidence from eBay,” *RAND Journal of Economics*, 2019, 50 (4), 822–853.
- **and Ali Shourideh**, “Optimal Rating Design,” 2020, pp. 1–47.
- Sahni, Navdeep S and Harikesh S Nair**, “Does Advertising Serve as a Signal? Evidence from a Field Experiment in Mobile Search,” *The Review of Economic Studies*, 2019, (October 2019), 1529–1564.
- Vellodi, Nikhil**, “Ratings Design and Barriers to Entry,” *SSRN Electronic Journal*, 2020, pp. 1–63.

## A Proofs

*Proof of Theorem 1.* By  $M_t = \mu + \lambda[Y_t - \nu] \Leftrightarrow \lambda Y_t = M_t - \mu + \lambda\nu$ , and the linear strategy  $F_t = \alpha\theta_t + \beta M_t + \delta\mu$ , the increment of  $M_t$  is written as

$$\begin{aligned}
dM_t &= d(\lambda Y_t) \\
&= (-\phi + a\lambda\beta) M_t dt \\
&\quad + (a\lambda\alpha + bq\lambda) \theta_t dt \\
&\quad + (\phi\mu - \phi\lambda\nu + a\lambda\delta\mu) dt \\
&\quad + bq\lambda\sigma_\xi dZ_t^\xi
\end{aligned}$$

Now, we look for a quadratic value function

$$V = v_0 + v_1\theta + v_2M + v_3\theta^2 + v_4M^2 + v_5\theta M \quad (11)$$



satisfying the HJB equation:

$$\begin{aligned}
rV(\theta, M) = & \sup_{F \in \mathbb{R}} (1 - \tau) M \cdot q - \tau M \cdot F - \frac{c}{2} F^2 \\
& - \kappa(\theta - \mu) V_\theta \\
& + \{a\lambda F + bq\lambda\theta - \phi[M - \bar{\theta} + \lambda\bar{Y}]\} V_M \\
& + \frac{\sigma_\theta^2}{2} V_{\theta\theta} \\
& + \frac{bq\lambda^2\sigma_\xi^2}{2} V_{MM}
\end{aligned}$$

By the first-order condition,

$$\begin{aligned}
0 &= -\tau M - cF + a\lambda V_M \\
\Leftrightarrow F &= -\frac{\tau}{c} M + \frac{a\lambda}{c} V_M \\
&= \frac{a\lambda}{c} v_5 \theta + \left(2\frac{a\lambda}{c} v_4 - \frac{\tau}{c}\right) M + \frac{a\lambda}{c} v_2
\end{aligned}$$

By matching coefficients with  $F = \alpha\theta + \beta M + \delta\mu$ ,

$$\begin{aligned}
\alpha &= \frac{a\lambda}{c} v_5 \\
\beta &= 2\frac{a\lambda}{c} v_4 - \frac{\tau}{c} \\
\delta\mu &= \frac{a\lambda}{c} v_2
\end{aligned}$$

By solving them for  $v_k$ 's,

$$\frac{c}{a\lambda} \alpha = v_5 \tag{12}$$

$$\frac{c}{2a\lambda} \left(\beta + \frac{\tau}{c}\right) = v_4 \tag{13}$$

$$\frac{\delta\mu c}{a\lambda} = v_2 \tag{14}$$

By the Envelop condition w.r.t.  $M$ ,<sup>16</sup>

$$\begin{aligned}
rV_M &= (1 - \tau) q - \tau F \\
&- \kappa(\theta - \mu) V_{\theta M} \\
&- \phi V_M \\
&+ \{a\lambda F + bq\lambda\theta - \phi[M - \mu + \lambda\nu]\} V_{MM}
\end{aligned}$$

---

<sup>16</sup>The envelop condition w.r.t.  $\theta$  gives conditions characterizing  $v_1$  and  $v_3$ , and one characterizing  $v_5$ , which coincides with the condition from the envelop condition w.r.t.  $M$ .

By inserting the derivatives of eq.(11) and equating the coefficients of  $\theta$ ,  $M$ , and constants on LHS and RHS,

$$\begin{aligned}(r + \phi) v_5 &= -\tau\alpha - \kappa v_5 + \{a\lambda\alpha + bq\lambda\} 2v_4 \\ 2(r + \phi) v_4 &= -\tau\beta + \{a\lambda\beta - \phi\} 2v_4 \\ (r + \phi) v_2 &= (1 - \tau)q - \tau\delta\bar{\theta} + \kappa\mu v_5 + \{a\lambda\delta\mu + \phi\mu - \phi\lambda\nu\} 2v_4\end{aligned}$$

Then, inserting eq.(12) to eq (14),

$$(r + \phi + \kappa) \frac{c}{a\lambda} \alpha = -\tau\alpha + \{a\lambda\alpha + bq\lambda\} 2 \frac{c}{2a\lambda} \left(\beta + \frac{\tau}{c}\right) \quad (15)$$

$$2(r + \phi) \frac{c}{2a\lambda} \left(\beta + \frac{\tau}{c}\right) = -\tau\beta + \{a\lambda\beta - \phi\} 2 \frac{c}{2a\lambda} \left(\beta + \frac{\tau}{c}\right) \quad (16)$$

$$(r + \phi) \frac{\delta\mu c}{a\lambda} = (1 - \tau)q - \tau\delta\mu + \kappa\mu \frac{c}{a\lambda} \alpha + \{a\lambda\delta\mu + \phi\mu - \phi\lambda\nu\} 2 \frac{c}{2a\lambda} \left(\beta + \frac{\tau}{c}\right) \quad (17)$$

By combining with the consistency of  $\lambda$ :  $\lambda = \frac{(a\alpha+bq)\sigma_\theta^2(\phi-a\beta\lambda)}{(\phi-a\beta\lambda+\kappa)\kappa bq\sigma_\xi^2+\sigma_\theta^2(a\alpha+bq)^2}$ , we can characterize  $\alpha$ ,  $\beta$ ,  $\delta$ ,  $\lambda$ . In the following, I do so by using an aggregator  $L = -a\beta\lambda$  so that the stationarity condition is easier to verify. First, by replacing  $\lambda$  to  $-\frac{L}{a\beta}$  in the above four equations,

$$0 = -\frac{bq(\beta c + \tau)}{a} + \alpha\tau - \alpha(\beta c + \tau) - \frac{\alpha\beta c\kappa}{L} - \frac{\alpha\beta c\phi}{L} - \frac{\alpha\beta cr}{L} \quad (18)$$

$$0 = \beta\tau - \beta(\beta c + \tau) - \frac{2\beta\phi(\beta c + \tau)}{L} - \frac{\beta r(\beta c + \tau)}{L} \quad (19)$$

$$0 = \frac{\nu\phi(\beta c + \tau)}{a} - \delta\mu(\beta c + \tau) + \frac{\alpha\beta c\kappa\mu}{L} - \frac{\beta c\delta\mu\phi}{L} + \frac{\beta\mu\phi(\beta c + \tau)}{L} - \frac{\beta c\delta\mu r}{L} + \delta\mu\tau + q\tau - q \quad (20)$$

$$-\frac{L}{a\beta} = \frac{\sigma_\theta^2(L + \phi)(a\alpha + bq)}{\sigma_\theta^2(a\alpha + bq)^2 + \kappa bq\sigma_\xi^2(\kappa + L + \phi)} \quad (21)$$

By solving (19) for  $\beta$ , we get  $\beta = -\frac{\tau}{c} \left(\frac{r+2\phi}{r+2\phi+L}\right) \equiv B(L)$ . By inserting this into (18) and solving it for  $\alpha$ , we get  $\alpha = \frac{bq}{a} \frac{L^2}{(r+2\phi)(r+\phi+\kappa+L)} \equiv A(L)$ . By plugging  $\beta = B(L)$  and  $\alpha = A(L)$  into (21), we obtain an equation characterizing  $L$ :

$$-\frac{L}{aB(L)} = \frac{\sigma_\theta^2(L + \phi)(aA(L) + bq)}{\sigma_\theta^2(aA(L) + bq)^2 + \kappa bq\sigma_\xi^2(\kappa + L + \phi)}$$

Rearranging it , we get

$$\begin{aligned}1 &= \frac{\sigma_\theta^2(L + \phi)(aA(L) + bq)}{\sigma_\theta^2(aA(L) + bq)^2 + \kappa bq\sigma_\xi^2(\kappa + L + \phi)} \frac{-aB(L)}{L} \\ &\equiv h(L)\end{aligned}$$

To evaluate  $h(L)$ , the sign of  $L$  is useful to characterize.

**Lemma 6.**  $\beta < 0$  and  $L > 0$  under the linear stationary Gaussian equilibrium.

*Proof.* By the stationarity, we must have  $\phi + L > 0$ . Then,

$$\begin{aligned}\beta &= -\frac{\tau}{c} \left( \frac{r + 2\phi}{r + 2\phi + L} \right) \\ &= -\frac{\tau}{c} \left( \frac{r + 2\phi}{r + \phi + \phi + L} \right) \\ &< 0\end{aligned}$$

Then,  $\alpha = \frac{bq}{a} \frac{L^2}{(r+2\phi)(r+\phi+\kappa+L)} > 0$  and  $\lambda = \frac{(a\alpha+bq)\sigma_\theta^2(\phi+L)}{(\phi+L+\kappa)b^2q^2\kappa\sigma_\xi^2+\sigma_\theta^2(a\alpha+bq)^2} > 0$ . Now, we can conclude  $-a\beta\lambda \equiv L > 0$ .  $\square$

Now, it is shown that  $\lim_{L \rightarrow 0} h(L) = \infty$  and  $\lim_{L \rightarrow \infty} h(L) = 0$ . Then, combined with the continuity of  $h(L)$ , there exist some  $L$  such that  $h(L) = 1$ . The uniqueness is proved by checking whether  $h'(L) < 0$  holds. It is shown that

$$h'(L) = -h_1(L) \{h_2(L) + L^4(-\kappa^2 + 6\kappa\phi + 4r^2 + 2\kappa r + 17r\phi + 19\phi^2)\}$$

where  $h_1(L), h_2(L) > 0$  for all  $L > 0$ . Thus,  $6\kappa\phi + 4r^2 + 2\kappa r + 17r\phi + 19\phi^2 > \kappa^2$  is sufficient for  $h'(L) < 0$ .  $\square$

*Proof of Lemma 2.* By plugging  $\alpha(L)$  and  $\beta(L)$  in to  $h$ , it can be written as  $h(L) = \frac{\frac{a\tau}{c} \frac{h_3}{L(L+r+2\phi)(h_4+(\sigma_\xi/\sigma_\theta)^2 h_5)}}{\square}$

where  $h_3 = (L + \phi)(r + 2\phi)^2(\kappa + L + r + \phi)(L^2 + L(r + 2\phi) + (r + 2\phi)(\kappa + r + \phi))$ ,  $h_4 = bq(L^2 + L(r + 2\phi) + (r + 2\phi)(\kappa + r + \phi))^2$ ,  $h_5 = \kappa(r + 2\phi)^2(\kappa + L + \phi)(\kappa + L + r + \phi)^2$ . Note that  $h_3, h_4, h_5$  are positive and independent of  $a$  and  $\sigma_\xi/\sigma_\theta$ . Thus,  $h$  is increasing in  $\frac{a\tau}{c}$  and decreasing in  $\sigma_\xi/\sigma_\theta$ . Since  $h'(L) < 0$  is shown in the proof of Theorem 1, the implicit function theorem tells that  $L$  is increasing in  $a$  and decreasing in  $\sigma_\xi/\sigma_\theta$ . Furthermore,  $h(L) \rightarrow \infty$  if  $L$  is bounded above and  $\frac{a\tau}{c} \rightarrow \infty$ . Thus, to satisfy the equilibrium condition:  $1 = h(L)$ ,  $L$  goes infinite as  $\frac{a\tau}{c}$  goes infinite. Similarly,  $h(L) \rightarrow 0$  if  $L$  is bounded away from zero and  $\frac{a\tau}{c} \rightarrow 0$ . Thus,  $L$  goes infinite as  $\frac{a\tau}{c}$  goes infinite to satisfy the equilibrium condition.

*Proof of Proposition 1 and 2.* Since  $E[M_t] = E[E[\theta_t|Y_t]] = \mu$ , we have  $E[F_t] = E[\alpha\theta_t + \beta M_t + \delta\mu] = (\alpha + \beta + \delta)\mu$ . By expressing  $\alpha, \beta, \delta$  as a function of the equilibrium aggregator  $L$ , it is written as  $E[F_t] = \frac{cLq(1-\tau)(L+r+2\phi)-\mu\tau^2(r^2+3r\phi+2\phi^2)}{c\tau(L^2+L(r+2\phi)+r^2+3r\phi+2\phi^2)}$  and the partial derivative with respect to  $L$  is  $\frac{\partial E[F_t]}{\partial L} = \frac{(r^2+3r\phi+2\phi^2)(2L+r+2\phi)(cq(1-\tau)+\mu\tau^2)}{c\tau(L^2+L(r+2\phi)+r^2+3r\phi+2\phi^2)^2} > 0$ .

Since  $a$ ,  $\sigma_\xi$ , and  $\sigma_\theta$  affects  $E[F_t]$  only through the aggregator  $L$ , we can show the effects of  $a$  and  $\frac{\sigma_\xi}{\sigma_\theta}$  by analyzing the sign of  $\frac{dL}{da}$  and  $\frac{dL}{d(\sigma_\xi/\sigma_\theta)}$ . By Lemma 2, we can conclude  $E[F_t]$  increasing in  $a$  and decreasing in  $\frac{\sigma_\xi}{\sigma_\theta}$ .

Since  $E[F_t] > 0$  for sufficiently large  $L$  and  $L \rightarrow \infty$  as  $a \rightarrow \infty$ ,  $E[F_t] > 0$  holds for sufficiently large  $a$ .  $\square$

*Proof of Proposition 3.* The equilibrium condition gives  $\alpha = \frac{bq}{a} \frac{L^2}{(r+2\phi)(r+\phi+\kappa+L)}$  and  $\beta = -\frac{\tau}{c} \left( \frac{r+2\phi}{r+2\phi+L} \right)$ . Furthermore, it is shown that  $\frac{\partial \alpha}{\partial L} > 0$  and  $\frac{\partial \beta}{\partial L} > 0$ . Then, Lemma 2 concludes the proposition.  $\square$

*Proof of Lemma 3 and 4.* An arbitrary strategy  $\alpha, \beta, \delta$  satisfying  $\phi - a\beta\lambda$  (not necessarily the equilibrium strategy) generates a stationary distribution. Using the variance-covariance matrix of the stationary distribution, the informativeness is written as

$$\rho^2 = \frac{(\phi - a\beta\lambda)(a\alpha + bq)^2}{(\kappa + \phi - a\beta\lambda) \left( (a\alpha + bq)^2 + \kappa bq (\sigma_\xi/\sigma_\theta)^2 (\kappa + \phi - a\beta\lambda) \right)}$$

Thus, the informativeness without fake reviews is

$$\rho^2 = \frac{\phi(bq)^2}{(\kappa + \phi) \left( (bq)^2 + \kappa bq (\sigma_\xi/\sigma_\theta)^2 (\kappa + \phi) \right)}$$

.On the other hand, at the equilibrium,  $-a\beta\lambda$  can be replaced to  $L$ , and  $a\alpha$  is written as a function in  $L$ :  $a\alpha = bq \frac{L^2}{(r+2\phi)(r+\phi+\kappa+L)}$  such that  $a\alpha = 0$  when  $L = 0$ . Note that  $a$  does not appear in the RHS, so the direct and indirect effects of  $a$  on  $a \cdot \alpha$  are all captured by  $L$ . Now the equilibrium informativeness is written as:

$$\rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) = \frac{(\phi + L)(a\alpha + bq)^2}{(\kappa + \phi + L) \left( (a\alpha + bq)^2 + \kappa bq (\sigma_\xi/\sigma_\theta)^2 (\kappa + \phi + L) \right)}.$$

Note that  $\rho^2(0; \phi, \kappa, \sigma_\xi, \sigma_\theta) = \frac{\phi(bq)^2}{(\kappa + \phi) \left( (bq)^2 + \kappa bq (\sigma_\xi/\sigma_\theta)^2 (\kappa + \phi) \right)}$  coincides with the informativeness without fake reviews. This concludes Lemma 4.  $\square$

*Proof of Proposition 5.* The first part is proved by the limit as  $L \rightarrow \infty$ :

$$\begin{aligned} & \lim_{L \rightarrow \infty} \rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) \\ &= \lim_{L \rightarrow \infty} \frac{(\phi + L)}{(\kappa + \phi + L)} \frac{(a\alpha + bq)^2}{\left((a\alpha + bq)^2 + \kappa bq (\sigma_\xi/\sigma_\theta)^2 (\kappa + \phi + L)\right)} \\ &= 1 \end{aligned}$$

The second part comes from the derivative of  $\rho^2$  with respect to  $L$  around zero.  $\square$

*Proof of Proposition 6.* The optimal  $\phi$  without fake reviews is characterized by  $\frac{\partial}{\partial \phi} \rho^2(0; \phi, \kappa, \sigma_\xi, \sigma_\theta) = 0$ , which yields  $\phi^0 = \sqrt{bq(\sigma_\theta/\sigma_\xi)^2 + \kappa^2}$  as the optimal level. On the other hand, the effect of  $\phi$  at the equilibrium is

$$\begin{aligned} \frac{d\rho^2}{d\phi} &= \frac{\partial}{\partial \phi} \rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) + \frac{\partial}{\partial L} \rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) \frac{dL}{d\phi} \\ &= \frac{\partial}{\partial \phi} \rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) - \frac{\partial}{\partial L} \rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) \frac{\partial h}{\partial \phi} / \frac{\partial h}{\partial L} \end{aligned}$$

By evaluating this at  $\phi = \phi^0$ , we obtain  $\frac{d\rho^2}{d\phi}|_{\phi=\phi^0} < 0$ .

The second part is proved by two inequalities:  $\rho^2(0; \phi^0, \kappa, \sigma_\xi, \sigma_\theta) < \rho^2(L(\phi^0); \phi^0, \kappa, \sigma_\xi, \sigma_\theta) \leq \rho^2(L(\phi^*); \phi^*, \kappa, \sigma_\xi, \sigma_\theta)$ . The first inequality is proved as follows. For any  $L > 0$ ,

$$\begin{aligned} & \rho^2(L; \phi^0, \kappa, \sigma_\xi, \sigma_\theta) - \rho^2(0; \phi^0, \kappa, \sigma_\xi, \sigma_\theta) \\ &= r \cdot g_1 + g_2 \end{aligned}$$

where  $g_1$  is polynomial in  $r$  and  $L$  and  $g_2 > 0$  is polynomial in  $L$  and does not depend on  $r$ . Since  $L \rightarrow C$  for some  $C > 0$  as  $r \rightarrow 0$ ,  $r \cdot g_1 + g_2$  converges to a positive number. Thus, for sufficiently small  $r$ , the first inequality holds. The second inequality holds by definition.  $\square$

*Proof of Proposition 7.* Similarly to Proposition 6, the total effect of  $\sigma_\xi/\sigma_\theta$  is written as  $\frac{d\rho^2}{d(\sigma_\xi/\sigma_\theta)} = \frac{\partial}{\partial(\sigma_\xi/\sigma_\theta)} \rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) - \frac{\partial}{\partial L} \rho^2(L; \phi, \kappa, \sigma_\xi, \sigma_\theta) \frac{\partial h}{\partial(\sigma_\xi/\sigma_\theta)} / \frac{\partial h}{\partial L}$ . It is shown that  $\frac{d\rho^2}{d(\sigma_\xi/\sigma_\theta)} < 0$ .  $\square$

*Proof of Theorem 2.* Now, we look for a quadratic value function

$$V = v_0 + v_1\theta + v_2Y + v_3\theta^2 + v_4Y^2 + v_5\theta Y \quad (22)$$

satisfying the HJB equation:

$$\begin{aligned}
rV(\theta, Y) = & \sup_{F \in \mathbb{R}} (1 - \tau) p \cdot q - \tau p \cdot F - \frac{c}{2} F^2 \\
& - \kappa(\theta - \mu) V_\theta \\
& + (aF + bq\theta - \phi Y) V_Y \\
& + \frac{\sigma_\theta^2}{2} V_{\theta\theta} \\
& + \frac{bq\sigma_\xi^2}{2} V_{YY} \\
\text{s.t. } & p = \mu - \left( \eta\lambda + (1 - \eta) \tilde{\lambda} \right) Y + \left( \eta\lambda\nu + (1 - \eta) \tilde{\lambda}\tilde{\nu} \right)
\end{aligned}$$

The first order condition and gives

$$v_5 = \frac{\alpha c}{a} \quad (23)$$

$$v_4 = \frac{\beta c + \hat{\lambda}\tau}{2a} \quad (24)$$

$$v_2 = \frac{c\delta\mu + \mu\tau - \widehat{\lambda\nu}\tau}{a} \quad (25)$$

where  $\hat{\lambda} = \left( \eta\lambda + (1 - \eta) \tilde{\lambda} \right)$  and  $\widehat{\lambda\nu} = \left( \eta\lambda\nu + (1 - \eta) \tilde{\lambda}\tilde{\nu} \right)$ , and the envelop condition gives

$$0 = \hat{\lambda}\alpha\tau - 2a\alpha v_4 - 2bqv_4 + rv_5 + \kappa v_5 + v_5\phi \quad (26)$$

$$0 = -2a\beta v_4 + \beta\hat{\lambda}\tau + 2rv_4 + 4v_4\phi \quad (27)$$

$$0 = -2a\delta\mu v_4 + \delta\mu\hat{\lambda}\tau + \hat{\lambda}q\tau - \hat{\lambda}q + rv_2 - \kappa\mu v_5 + v_2\phi \quad (28)$$

By inserting eq.(24) into (27) and solving it for  $\hat{\lambda}$  and by letting  $L = a\beta$ , we obtain

$$\hat{\lambda} = \frac{cL(L + r + 2\phi)}{a\tau(r + 2\phi)} \equiv \lambda(L)$$

On the other hand, the stochastic differential equation for  $(\theta, Y)$  gives

$$\lambda = \frac{bq\sigma_\theta^2(L + \phi)(A(L) + 1)}{\sigma_\theta^2(bqA(L) + bq)^2 + \kappa bq\sigma_\xi^2(\kappa + L + \phi)} \equiv \lambda(L)$$

$$\tilde{\lambda} = \frac{bq\sigma_\theta^2\phi}{\sigma_\theta^2(bq)^2 + \kappa bq\sigma_\xi^2(\kappa + \phi)} = \lambda(0)$$

Then, by rearranging

$$\begin{aligned}\hat{\lambda} &= \left( \eta\lambda + (1 - \eta)\tilde{\lambda} \right) \\ \Rightarrow 1 &= \frac{\eta\lambda(0) + (1 - \eta)\lambda(L)}{\hat{\lambda}(L)} \equiv h(L; \eta)\end{aligned}$$

Note that  $\lim_{L \rightarrow 0} h(L; \eta) = \infty$  and  $\lim_{L \rightarrow \infty} h(L; \eta) = 0$ . Then,  $h_L(L; \eta) < 0$  holds for any  $\eta \in [0, 1]$  as long as  $h_L(L; 1) < 0$ .

[Proof of Proposition 8] Since  $\lambda(0) \leq \lambda(L)$  for any  $L \geq 0$ , we have  $h(L; \eta) \leq h(L; 1)$  for any  $\eta \in [0, 1]$ . Thus, the equilibrium  $L$  will be smaller given  $\eta < 1$  than the equilibrium  $L$  given  $\eta = 1$ .

The expected amount of the fake reviews is

$$E[F_t] = \alpha\mu + \beta\nu + \delta\mu$$

By plugging the equilibrium conditions and taking derivative with respect to  $L$ , we can show  $\frac{\partial}{\partial L} E[F_t] \geq 0$ .

[Proof of Proposition 9] At the equilibrium,  $\frac{\partial bias}{\partial L} \geq 0$  always holds and  $\frac{\partial bias}{\partial a} \geq 0$  holds if  $bias \geq 0$ .  $\square$

## B An interpretation of the pricing rule

this pricing rule as a result of competition among heterogeneous consumers, to which we can easily introduce a mixture of rational and naive consumers in the next section. Suppose that consumer  $i \in [0, n]$  feels  $u_{t,i} = \theta_t + \epsilon_{t,i} - p_t$  if the consumer buy the product, and 0 otherwise, where  $\epsilon_{t,i}$  is identically and independently distributed. Then, given the rating shown on the platform,  $Y_t$ , the consumer will choose to purchase the product if and only if  $E[\theta_t|Y_t] + \epsilon_{t,i} - p_t \geq 0$ . Therefore, the demand function is expressed as  $n \cdot (1 - F(p_t - M_t))$  where  $F(\cdot)$  is a c.d.f. of the random variable  $\epsilon_{t,i}$ . By letting  $n = 2q$  and assuming that  $\epsilon_{t,i}$  is distributed symmetrically around zero. We obtain  $p_t = M_t$  as the market clearing price.

## C An Alternative Model with Changing $q$

The same results with the base line model can be generated with a slightly different specification of the model with the quantity level dependent on the reputation level.

Now, suppose that the seller sells  $q_t$  units of the product at a fixed price of  $p$ , and makes  $F_t$  units of fake reviews. The quality of the product is denoted as  $\theta_t$ . A sufficiently large mass of consumers forms a belief on the quality  $E[\theta_t|Y_t] \equiv M_t$  and the demand function based on that. Since the price is fixed, high reputation results in large quantity:  $q_t = M_t$ .

The quality  $\theta_t$  evolves in the same way as the main model. The new information as

$$aF_t dt + bq_t \left( \theta_t dt + \sigma_\xi dZ_t^\xi \right) \quad (29)$$

The difference from the main model is that the quantity varies over time and the coefficient of  $dZ_t^\xi$  is now defined as  $bq_t\sigma_\xi$  instead of  $\sqrt{bq_t}\sigma_\xi$ . In this specification, we can analyze the effect of the organic reviews crowding out the fake reviews, but not the effect of the large transaction generating intrinsically more precise information by the large sample.

The seller's instantaneous payoff is defined as:

$$\pi_t = (1 - \tau) p (q_t + F_t) - p \cdot F_t - \frac{c}{2} \left( \frac{F_t}{q_t} \right)^2$$

where  $\tau$  is transaction fees imposed by the platform. The specification of the quadratic cost is now different from the base line model: the seller needs to pay a large cost if the seller tries to increase the share of the fake reviews among the all the reviews. The revenue and the reimbursement cost is still the same as the baseline model.

$$\begin{aligned} \pi_t &= (1 - \tau) p q_t - \tau p \cdot F_t - \frac{c}{2} \left( \frac{F_t}{q_t} \right)^2 \\ &= (1 - \tau) p M_t - \tau p \cdot M_t \frac{F_t}{M_t} - \frac{c}{2} \left( \frac{F_t}{M_t} \right)^2 \end{aligned}$$

By changing the choice variable of the seller from  $F_t$  to  $\frac{F_t}{M_t}$ , which is the combination of the original variable and a constant at time  $t$ , we can write the instantaneous profit isomorphic to one in the baseline model. To simplify the analysis, we assume that the platform use an average information at time  $t$  to update the ratings:

$$d\xi = \frac{a}{b} \frac{F_t}{M_t} dt + \theta_t dt + \sigma_\xi dZ_t^\xi \quad (30)$$

The model is then isomorphic to the baseline model, so generates the same results as those from the baseline model.

## D Simulation Results

### D.1 Mixture of the Rational and Naive Consumers

It is natural to question how the bias and the correlation changes as the consumer changes from totally naive to totally rational. We also examine mean squared error since it integrates both bias and the variance of the estimate in a natural way.



### D.1.1 Mean Squared Error

Since the price is modeled as the market's estimate of the quality, we consider the mean squared errors of the price

$$\begin{aligned}
MSE_p &= E \left[ (p_t - \theta_t)^2 \right] \\
&= E \left[ \left( \eta \{ \mu + \lambda [Y_t - \nu] \} + (1 - \eta) \{ \mu + \tilde{\lambda} [Y_t - \tilde{\nu}] \} - \theta_t \right)^2 \right] \\
&= Var(Y_t) \left\{ \left( \eta\lambda + (1 - \eta)\tilde{\lambda} \right)^2 - 2 \left( \eta\lambda + (1 - \eta)\tilde{\lambda} \right) \lambda \right\} + (1 - \eta)^2 Bias^2 + Var(\theta_t)
\end{aligned}$$

Note that, when  $\eta = 1$ , it is reduced to

$$\begin{aligned}
MSE_p &= -\lambda^2 Var(Y_t) + Var(\theta_t) \\
&= Var(\theta_t) \left\{ 1 - \frac{Cov(Y_t, \theta_t)^2}{Var(Y_t)^2} \frac{Var(Y_t)}{Var(\theta_t)} \right\} \\
&= Var(\theta_t) \{1 - \rho^2\}
\end{aligned}$$

For different levels of  $\eta$ , we calculate the correlation of  $Y$  and  $\theta$  as a criteria for the rational consumers, the bias as a criteria for naive consumers, and the mean squared error as a criteria for the whole market. See fig.5 for the simulation results.

Figure 5:

