# (Pro-)Social Learning and Strategic Disclosure[*]

Roland Bénabou[†]     Nikhil Vellodi[‡]

*American Economic Journal: Micro,* forthcoming

**Abstract**

We study a sequential experimentation model with endogenous feedback. Agents choose between a safe and risky action, the latter generating stochastic rewards. When making this choice, each agent is selfishly motivated (myopic). However, agents can disclose their experiences to a public record, and when doing so are pro-socially motivated (forward-looking). Disclosure is both *polarized* (only extreme signals are disclosed) and *positively biased* (no feedback is bad news). The extent of disclosure is non-monotone in prior uncertainty. Subsidizing disclosure costs can paradoxically lead to less disclosure, but more experimentation.

**Keywords**: social learning, experimentation, dynamic disclosure, consumer reviews, time-inconsistent preferences, motivated beliefs. **JEL Classification**: D82, L11, L12.

## 1 INTRODUCTION

In many settings, agents face a choice between safe and risky actions, with different individuals facing these choices in sequence. Agents might benefit from the information generated by those who preceded them. For instance, consider consumers choosing whether or not to dine at a restaurant with unknown quality, or to watch a new movie. Those who do so can then leave feedback, helping later-arriving consumers make more informed choices. Similar

[†]Princeton University, NBER, CEPR, IZA, BREAD, and briq. rbenabou@princeton.edu
[‡]Paris School of Economics. n.vellodi@gmail.com

settings include the adoption of new products and technologies, employment choices, and sequential voting.

A well-known dynamic externality emerges in such settings, namely that agents do not internalize the benefit to future consumers of taking the risky action, leaving feedback and thereby generating socially valuable information. To remedy this inefficient under-exploration, a planner would direct agents to take the risky action even when it is unprofitable to them, provided the informational gain to future agents more than compensates. This question has been studied extensively in economics and computer science under the label of "incentivized exploration (IE)" (Kremer et al., 2014; Che and Hörner, 2018), itself part of the broader literatures on social learning and sequential experimentation (Banerjee, 1992; Bikhchandani et al., 1992; Smith and Sørensen, 2000; Smith et al., 2021).[1] Papers in the IE literature largely take a *normative* approach to the problem. Namely, they assume the presence of a benevolent designer who can control the provision of incentives either via dynamic information provision or direct recommendations. Furthermore, these works largely assume that, once generated, individual signals are perfectly observed by the planner or designer in charge of public information provision.

Such studies are thus silent on a particularly salient issue within the online feedback setting — *why* and *when* do people leave feedback in the first place? For instance, consumers might be driven by a desire to help future consumers make informed choices, or to reward or punish a seller for a positive or negative experience. In practice, while feedback is often highly valued by consumers, the vast majority fail to provide it,[2] and those that do provide feedback display well-known biases such as *positive selection* (Nosko and Tadelis, 2015; Hui et al., 2024) — undisclosed experiences are on average negative — and *polarization* (Schoenmueller et al., 2020) — extreme reviews are more prevalent than average reviews.

We take a step towards addressing these questions by providing a *positive* theory of IE. Namely, we propose a simple three-period model of sequential experimentation, in which

---

[1]See Slivkins (2022) and Bikhchandani et al. (2022) for recent surveys on IE and social learning more generally.

[2]Recent surveys report that only around 10% of consumers regularly leave reviews. See https://tinyurl.com/mrrsf9v5, https://tinyurl.com/ux3zyuem.

we endow our agents with *pro-social* motives to leave feedback.[3] In our model, three agents arrive sequentially and have the choice between a safe (S) and risky action (R). Action S generates a deterministic reward, whereas R yields a random reward correlated to an underlying hidden state, for instance, the unknown quality of a product. If an agent plays R, they can disclose their signal. Our main analysis models disclosure via hard evidence with noisy transmission (Dye, 1985a); an agent is able to leave feedback with probability $\alpha \in [0,1]$, but can either truthfully report their signal or not report it at all. Two crucial assumptions determine an agent's payoff. First, when choosing their *action*, we assume that they are fully *self-interested*, maximizing only their personal reward. Second, when making their subsequent *disclosure* choice, we assume agents are fully *pro-social*. Formally, they have lexicographic preferences over making optimal consumption decisions for themselves, and then transmitting useful information to help others do the same. We discuss these modeling choices further below.

This simple combination of ingredients delivers a rich theory of selective disclosure, which both accords with well-documented phenomena and provides new testable predictions. In particular, our first main result (Theorem 1) demonstrates that equilibrium disclosure is both positively selected and polarized. The intuition is simple. When player 1 (P1) plays R, their disclosure choice is governed equally by the subsequent payoffs of P2 and P3. In contrast, P2 is guided purely by their own payoff when taking their action. For instance, an informed P2 might fail to experiment (play R), even though the loss to themselves is smaller than the gain to P3 that the information generated by doing so would provide. P1 would ideally like to avoid such instances, and thus strategically conceals their own experience, inducing P2 to experiment against their interests for the sake of P3.

Simply put, an early adopter would rather not take responsibility for causing the untimely demise of a new product, if there is a reasonable chance the product is in fact worth a second chance, and in this case they keep quiet. On the other hand, when experiences are sufficiently negative, P1 is convinced that no further experimentation should occur and

---

[3]We discuss the limitations of our three-period specification in Section 6.3.

thus terminates it by posting their feedback, while for (even marginally) positive experiences, there is no downside to disclosure. Thus, strategic non-disclosure is used exclusively by P1 to foster efficient experimentation by P2. Of course, models that assume that leaving feedback is costly and done only when sufficiently informative also generate polarized feedback, but struggle to also deliver positive selection from a single behavioral foundation.[4]

Beyond these, our model delivers further predictions. For instance, we fully characterize how equilibrium non-disclosure, and thus experimentation, varies with the prior belief regarding risky payoffs (Theorem 2). We view this exercise as capturing, in a reduced-form manner, how disclosure varies with how old or well-established the product market in question is. We show that the extent of experimentation is hump-shaped in the prior. Moreover, equilibrium experimentation disappears as prior uncertainty vanishes.

We also show that the extent of experimentation is also hump-shaped in $\alpha$, the feedback opportunity parameter. This insight has important implications for real-world interventions; practitioners argue that the lack of feedback in online markets leads to biased inference, and that making feedback less costly (e.g., by providing explicitly monetary incentives) would lead to more information and thus experimentation (Marinescu et al., 2021). If we take the natural interpretation that $\alpha$ corresponds to the fraction of agents for whom feedback is costless and $1 - \alpha$ the fraction for whom it is prohibitively costly, our result suggests that making feedback less costly could paradoxically lead to *less* disclosure, and more generally that the effectiveness of such interventions in stimulating feedback rates will vary by products and markets.

The joint assumption of selfish consumption and pro-social disclosure is appealing on three separate fronts. First, endowing agents with benevolent preferences in this manner allows our theory to be viewed as a minimal departure from the normative analyses in the IE literature. That is, our agents are effectively mini-planners when disclosing, facing the same trade-off between *exploration* (long-run information gains) and *exploitation* (short-run consumption gains) as in previous work, but they also face additional constraints imposed

---

[4]For instance, Hui et al. (2024) allow feedback to be positively biased for unmodelled reasons, suggesting reasons such as fear of retaliation or a simple aversion to providing negative criticism.

on them in equilibrium, such as ex-post optimality of disclosure rules. Our results thus demonstrate how such constraints shape the degree to which disclosure can be used to incentivize exploration.

Second, from a positive perspective, recent surveys suggest that the welfare of other consumers is a key driver when leaving feedback.[5] At the same time, empirical evidence suggests that incentives to provide feedback are divorced from actual consumption choices in online settings (Cabral and Li, 2015). We present a first attempt at formalizing these arguments, with a view to understanding both their theoretical foundations and their ability to organize empirical findings.

Third, the informational externality described above derives fundamentally from the structure of intertemporal preferences, namely that agents are "present-biased" when making their consumption choices. This gives rise to an alternative, psychological interpretation of the model. Instead of a sequence of agents, consider a single decision maker with the following dynamically inconsistent preferences. When taking actions that affect current payoffs, they are myopic (completely present biased), whereas when deciding what available information to store in memory to inform future choices, they are patient. This corresponds to a limiting case of quasi-hyperbolic, or $\beta\delta$ (Laibson, 1997) preferences where $\beta$ is arbitrarily small. By modeling the disclosure objective as altruistic, our framework permits this application to an individual who selectively encodes their experiences in order to become less "conservative" — that is, more open to trying and learning from new experiences. Our work thus demonstrates a close conceptual connection between IE and motivated reasoning (Bénabou and Tirole, 2002, 2004; Carrillo and Mariotti, 2000).

Our model of feedback imposes two important constraints: feedback must be both ex-post optimal and truthful. To explore how the first constraint shapes our results, we analyze communication under commitment, and find that it is again polarized and positively selected However, in contrast to the disclosure benchmark (Theorem 2), this pattern does not vanish with prior uncertainty.

---

[5]For example, see https://tinyurl.com/mrybw969.

5

Turning to the second constraint, our focus on disclosure of hard evidence as the channel through which feedback takes place is primarily motivated by the case of online reviews. While these share elements of both disclosure and cheap talk, a number of features make them closer to the former. First, the fact that a large share of consumers choose to not leave a review calls for a disclosure model. Second, platforms indicate which are "verified purchase" or "verified traveler" reviews and, conversely, take strong measures against fake reviews: using artificial intelligence to detect and remove them, and taking legal action against the intermediaries who sell such fake reviews.[6] Vendors can also appeal to the platform to remove false criticism. Third, reviewers often post specific facts, photos and videos, book commentaries, etc. to support their evaluations; relatedly, a large experimental literature documents significant and widespread lying aversion (e.g., Abeler et al. (2014, 2019)). Finally, reviews are themselves evaluated by other customers, who can tag a review as helpful or on the contrary report it as fake. Amazon also materially incentivizes and then highlights informative reviews with Amazon Vines, a program that selects "customers who consistently write [the most] insightful reviews"; they can then request for free products from thousands of brands, on which they then write reviews that are distinguished by a special badge. Admittedly, even all these measures still leave room for some fake reviews and biased reviews by real consumers (He et al., 2022). Therefore, we also analyze the case of cheap talk, comparing and contrasting its implications with those of disclosure

In concurrent and independent work, Smirnov and Starkov (2024) analyze a very similar model, focusing on the persuasion benchmark and cheap talk. Like us, they mainly study the three-period case, but also obtain some partial results for an infinite horizon (see Section 6 for further discussion). Analyzing disclosure allows us to uncover tight comparative-statics implications on the nature and degree of equilibrium communication; our results regarding non-monotone disclosure (Theorem 2 and Corollary 3) have no analog under either of the other two alternative forms of information transmission.

The paper proceeds as follows. After introducing the model (Section 2), we fully char-

---

[6]See https://tinyurl.com/2md6zhnp.

acterize equilibrium disclosure in Section 3. In Sections 4 and 5 we derive key comparative static results regarding both prior uncertainty and disclosure opportunities, respectively. In Section 6 we discuss our modeling choices, and in particular contrast verifiable disclosure with persuasion and cheap talk, with details in Online Appendices B and C. We conclude with thoughts on future research in Section 7. Unless otherwise mentioned, proofs are gathered in the Appendix.

## 2 Model

**Players and signals** – At each date $t = 1, 2, 3$, a short-lived agent arrives and takes a binary decision $a_t \in \{0, 1\}$, corresponding to safe and risky actions respectively. The safe action generates a payoff 0. The risky action incurs a cost $c \in (0, 1)$ and generates a payoff distributed according to $F_\theta$, where the distribution $F_\theta$ depends on a state $\theta \in \{H, L\}$, admits a density $f_\theta$ and is supported on the compact real interval $X = [\underline{x}, \bar{x}]$. We shall often refer to $a_t = 1$ as "consuming", and to the realized payoff $x \in X$ as a signal.

The state $\theta$ is initially unknown, with all agents sharing the common prior $p = \mathbb{P}(\theta = H)$. Let $p^x$ denote the posterior belief formed by combining the belief $p$ with the signal $x \in X$. That is,

$$p^x \equiv \frac{pf_H(x)}{f_p(x)} \equiv \frac{pf_H(x)}{pf_H(x) + (1-p)F_L(x)} \quad \text{for } x \in X. \tag{1}$$

Note that for all $p \in (0, 1)$, $p^x = p$ if and only if $f_H(x) = f_L(x)$. Let $\hat{x}$ denote the "neutral" signal that satisfies this equality, and more generally, let $x(p, q)$ solve $p^{x(p,q)} = q$, i.e. it is the signal required to achieve posterior $q$ starting from prior $p$.[7] We will sometimes use a natural transformation from signal space $X$ into belief space $[0, 1]$. Namely, we denote by $G$ the distribution (with density $g$) over posterior beliefs induced by the signal distribution: for each $p, q \in [0, 1]$, let $G_p(q) \equiv F_p(x(p, q))$, where $F_p \equiv pF_H + (1-p)F_L$.

**Disclosure** – Conditional on receiving outcome $x$, the agent may then have the opportunity to provide feedback regarding their experience, via direct communication. We

---

[7]We will impose assumptions on $f_\theta$ that ensure that both $\hat{x}$ and $x(p, q)$ are guaranteed to exist and be unique for all $p, q \in (0, 1)$.

assume hard evidence and verifiable disclosure (Dye, 1985a; Jung and Kwon, 1988), wherein a player: (i) with probability $\alpha \in [0, 1]$, is able to freely disclose their signal $x$, and chooses whether or not to do so; (ii) with probability $1 - \alpha$, has no such opportunity, for instance due to a prohibitively high disclosure cost.[8]

**Payoffs** – Each agent values the payoffs to both themselves and future agents, but very differently. We assume a form of lexicographic preferences, in which players care infinitely more about their own consumption than that of any other consumer.[9] Formally, given a belief $p_t$, agent $t$ chooses $a_t$ to maximize their expected consumption payoff $a_t \mathbb{E}(x_t - c)$, so that by Assumption $(\mathbf{1.c})$, $a_t(p_t) = \mathbb{I}_{p_t \geq c}$.[10]

On the other hand, once their consumption choice has been made, agents value the welfare of future consumers equally when making their disclosure choice. Thus, if P2 consumes and obtains the signal $x$, their value from inducing a belief $r$ upon P3 through their disclosure choice, while themselves holding belief $q$, is $V_2(r \mid q) = u(r \mid q) \equiv \mathbb{I}_{r \geqslant c}(q - c)$, i.e. the utility, as judged by P2, that P3 will derive from their own consumption decision.

We will restrict our attention throughout the paper to equilibria in which P2 fully reveals. This is natural for several reasons. First, truthful revelation is weakly dominant for P2, as P2 and P3 have fully aligned preferences. Second, we show in Online Appendix A that truthful revelation by P2 is strictly dominant in the presence of (possibly arbitrarily small) shocks to players' payoffs, and is thus uniquely selected by an argument of robustness to such perturbations.

Turning now to P1, they value the consumption outcomes of both P2 and P3 equally, hence their continuation value $V_1(r \mid q)$ is

$$V_1(r \mid q) = \begin{cases} u(r \mid q) + \alpha \Lambda(r \mid q) + (1 - \alpha)u(r \mid q) & \text{if} \quad r \geqslant c \\ 0 & \text{if} \quad r < c, \end{cases} \tag{2}$$

---

[8]In a slight variant of the model, the arrival of agents is random and unobservable to others, occurring with probability $\alpha$ in each period.

[9]We discuss this assumption in Section 6.3.

[10]The weak inequality implies that each agent is assumed to consume when indifferent.

where

$$\Lambda(r \mid q) \equiv \mathbb{E}\left[u(r^z \mid q^z)\right] = \int_{x(r,c)}^{\bar{x}} (q^z - c) \, f_q(z) \, dz$$

denotes the expected consumption value of P3 from P1's perspective, given that P1 holds private belief $q$ and that P2 both holds belief $r$ and consumes.

The disclosure rule is a function $d : X \to \{0, 1\}$, where $d(x) = 1$ denotes disclosure by P1 of signal $x$ at prior $p$ and $d(x) = 0$ denotes non-disclosure. We will typically use $p$ to denote P1's prior belief, $q$ to denote P1's posterior belief, and $r$ to denote the public prior held by P2, which is ultimately determined by P1's disclosure rule $d$.

**Equilibrium** – In order to describe incentive compatible disclosure rules, we must develop our analysis of belief formation under non-disclosure. If the signal $x$ is disclosed, it is simply combined with the current belief according to Bayes' rule (1). If it is not disclosed, then the update rule must account for all other signals at which non-disclosure also occurs, as well as the possibility that disclosure was not feasible. For a disclosure rule $d$, let $D(d) = \{x \in X \mid d(x) = 1\}$ and $N(d) \equiv X \backslash D(d)$.[11] We have:

$$p^{\emptyset} \equiv \frac{\mathbb{P}(d = \emptyset \mid \theta = H)}{\mathbb{P}(d = \emptyset)} = \frac{(1 - \alpha)p + \alpha \int_{N(d)} p^x f_p(x) \, dx}{(1 - \alpha) + \alpha \int_{N(d)} f_p(x) \, dx}. \tag{3}$$

The relevant incentive compatibility (IC) constraint for the disclosure choice by P1 is then: for all $x \in X$, $d(x) = 1$ if and only if

$$V_1(p^x \mid p^x) \geqslant V_1(p^{\emptyset} \mid p^x). \tag{4}$$

An equilibrium is simply a disclosure rule $d$ for P1 such that: 1) given the non-disclosure belief $p^{\emptyset}$, $d$ is incentive compatible, and 2) given $d$, $p^{\emptyset}$ is correctly computed:

**Definition 1.** An equilibrium is a disclosure rule $d$ such that (3) and (4) are satisfied.

---

[11]For any $p < c$, $P_1$ abstains from consuming ($a_1 = 0$) and thus has no signal to report, making $D(d)$ irrelevant. In what follows we will therefore focus on values $p \geq c$.

Let us define the *experimentation region* of a disclosure rule $d$ as

$$X_E(d) = \{x \in X \mid x \in N(d) \text{ and } p^\emptyset \geq c > p^x\}.$$

A signal $x \in X_E(d)$ if under $d$, P1 chooses not to disclose it, and by so doing induces P2 to consume when they wouldn't if P1 had disclosed. An equilibrium $d$ is an *experimentation equilibrium* (EE) if $X_E(d)$ has strictly positive measure. Let $\mathcal{E}$ denote the set of all such equilibria.[12] An equilibrium $d$ is a *maximal experimentation equilibrium* (MEE) if $d \in \mathcal{E}$ and $d' \in \mathcal{E}$ implies $X_E(d') \subset X_E(d)$. Thus, the MEE contains the largest experimentation region out of all EE. As we show below in Lemma 2, the MEE is the welfare-optimal equilibrium and thus forms a natural benchmark, on which we will later on perform comparative statics.

## 3 Positively Biased and Polarized Disclosure

We now analyze equilibrium disclosure rules. First, we introduce further natural assumptions on the signal structure (Smith et al., 2021):

**Assumption 1.** *(1.a) $F_H, F_L$ satisfy the monotone likelihood ratio property (MLRP).*

*(1.b)* $\inf_x \left( \frac{f_L}{f_H} \right)(x) = 0$, $\sup_x \left( \frac{f_L}{f_H} \right)(x) = \infty$.

*(1.c)* $\mathbb{E}(x \mid \theta = H) = 1$ *and* $\mathbb{E}(x \mid \theta = L) = 0$.

Assumption (**1.a**) states that higher signals are more likely in the high state, and that no perfectly revealing signal exists in either state. Assumption (**1.b**) is the "unbounded beliefs" assumption of Smith and Sørensen (2000), stating that there always exists a signal strong enough to almost completely overturn any prior belief. Assumption (**1.c**) is a normalization ensuring that beliefs and expected payoffs coincide, i.e. $\mathbb{E}(x \mid p) = p$, and is made simply for algebraic convenience.[13] As in Smith et al. (2021), we further assume that the distribution

---

[12]Focusing on experimentation equilibria rules out pathological equilibria that turn on the indifference P1 has over disclosure of extreme signals. For instance, any disclosure profile $d$ such that $N(d) \subset [\underline{x}, x(p, c))$ and $p^\emptyset < c$ is an equilibrium; for $x < x(p, c)$, both disclosure and non-disclosure lead to P2 not consuming, while for $x \geq x(p, c)$, truth-telling is strictly optimal, as shown below in the proof of Theorem 1.

[13]Assumptions (**1.a**) and (**1.c**) jointly imply that $\underline{x} < 0, 1 < \bar{x}$.

10

of the log-likelihood ratio of signals is log-concave. This ensures an intuitive feature of belief updating known as "posterior monotonicity" (PM) holds under Bayesian updating.

**Assumption 2.** *Let $\phi_\theta(l)$ denote the state-contingent densities for the transformed variable $l = \log(x/(1-x))$. Then $\phi_\theta(\cdot)$ is log-concave for $\theta \in \{0,1\}$.*

Our first main result below — Theorem 1 — characterizes the structure of equilibrium disclosure. In order to interpret its content, we introduce two key concepts, *polarity bias* and *positively selected disclosure*:

**Definition 2.** A disclosure rule $d$ is:

1. *Polarized* if there exist $\underline{\varepsilon}, \overline{\varepsilon} > 0$ such that $d(p, x) = 1$ for all $x \in [0, \underline{\varepsilon}) \cup [1 - \overline{\varepsilon}, 1]$ and $N(d)$ has strictly positive measure.

2. *Positively selected at $p$ if $p^\emptyset < p$.*

A polarized disclosure rule is one where extreme signals are disclosed. A positively selected rule is one where the posterior belief formed upon observing no feedback is strictly lower than the prior, so that "no news is bad news". Theorem 1 shows that *any* experimentation equilibrium exhibits both of these features:

**Theorem 1.** *In any EE, player 1 adopts the disclosure strategy:*

$$
d(x) = \begin{cases} 1 & if \quad p^x \geqslant c \\ 0 & if \quad p^x \in [\underline{q}, c) \\ 1 & if \quad p^x < \underline{q}, \end{cases}
$$

*for some $\underline{q} \in (0, c)$.*

In equilibrium, $P1$ discloses only those signals that lie on either side of the interval $[x(p, \underline{q}), x(p, c)]$, thus exhibiting both polarity and positive selection (since $x(p, c) \leq x(p, p) = \hat{x}$). P1 thus thinks along the following lines. If disclosing their experience does

11

not affect P2's demand, then P1 is happy to do so. This is the case when P1's experience is "good enough", so that leaving feedback does no harm and improves public information. However, if disclosing leads P2 to not consume (and thus P3 subsequently), P1 discloses only if they are sufficiently convinced that the product's quality is low; in this case, P1 would rather terminate future consumption. Otherwise, they keep their opinion to themselves, as they would rather give the product a "second chance" by having P2 consume and generate further information.[14] Put simply, P1 is always happy to leave a good review, but thinks twice about leaving a bad review, and only does so if their experience was sufficiently bad.

## 3.1  EXPERIMENTATION VERSUS ACCURACY

To provide further intuition for the proof of Theorem 1, we identify the key tradeoff facing P1 when disclosing, namely fostering experimentation versus improving accuracy. Since disclosure is verifiable, if P1 wants to distort the actions of P2 they must do so by not disclosing their experience, causing a rift between their posterior belief and P2's prior. The benefit of doing so is that P2 will experiment when they would not have done otherwise. The cost is that this rift in beliefs will propagate through to P3, as P3 will combine P2's disclosed signal with P2's (incorrect) prior belief. Consequently, from P1's perspective, there is a chance that P3 will make consumption errors, i.e. consume when they shouldn't or not consume when they should.

To understand the role of such consumption errors more formally, we characterize the properties of the value function $V_1(r \mid q)$ as both P1's posterior belief $q$ and P2's prior belief $r$ vary, rather than studying disclosure rules directly. From equation (2), it is clear that the function $\Lambda$ is crucial in determining P1's preferences for strategic disclosure. The following lemma provides a complete characterization of $\Lambda$.

**Lemma 1.** *1. $r \mapsto \Lambda(r \mid q)$ is strictly increasing on $[0, q)$ and strictly decreasing on $(q, 1]$.*

*2. $q \mapsto \Lambda(r \mid q)$ is strictly increasing (and in particular affine) for all $r \geqslant c$.*

---

[14]Note that the restriction to EE's ensures that $p^\emptyset \geq c$ so that P2 consumes conditional on non-disclosure by P1. There may exist pathological, non-experimentation equilibria wherein $q$ is sufficiently low that $p^\emptyset < c$ and P2 does not experiment.

*3.* $\Lambda(c \mid c) > 0$.

Importantly, the map $r \mapsto \Lambda(r \mid q)$ is single-peaked at $q$. Thus, $\Lambda$ encodes the loss (from P1's perspective) from inducing an incorrect belief, due to consumption errors by P3. The further is $r$ from $q$, the greater is the likelihood that P3 makes consumption errors. For instance, when $r > q$, P3 might consume when they shouldn't (in the event that P2's signal $x$ results in $q^x < c \le r^x$), while conversely if $r < q$, P3 might not consume when they should. Only if $r = q$ do neither of these errors occur. $\Lambda(c \mid c) > 0$ quantifies the *option value* from P2's consumption; note that $u(c \mid c) = 0$, so while the immediate return from P2 consuming at belief $c$ is 0, the gain to P3 from such consumption is strictly positive, as there is a chance P2 receives a positive outcome, acquiring useful information and thus providing an expected gain to P3.

The question remains whether there exist situations in which P1 resolves this trade-off in favor of fostering experimentation. Consider posterior beliefs $q$ just below $c$. Ideally, P1 would like P2 to consume, but hold the correct belief to minimize consumption errors as discussed above. Formally, since $\Lambda(c \mid c) > 0$, $\Lambda(c \mid q) > 0$ for $q$ just below $c$ by part 1) of the lemma. However, since inducing a belief below $c$ leads to non-consumption, P1 understands that P2 must necessarily hold an incorrect belief for consumption to occur. Theorem 1 says that when P1's posterior is sufficiently close to $c$, they would rather suffer the loss in accuracy than terminate consumption.

This reasoning also reveals why multiple equilibria may exist. Since $r \mapsto \Lambda(r \mid q)$ is decreasing for $r > q$, the lowest posterior $q < c$ at which P1 is indifferent between disclosing and not is increasing in the non-disclosure belief $p^\emptyset$. Intuitively, a higher $p^\emptyset$ implies a greater chance of consumption errors by P3, which dampens P1's incentive to foster experimentation through non-disclosure, causing $q$ to be higher and thus sustaining the higher $p^\emptyset$ in equilibrium.

We can show, however, that the MEE is ex-ante welfare maximizing across all equilibria,

experimentation or otherwise; that is, it maximizes

$$\mathcal{W}(d;p) \equiv \int_{D(d)} V_1(q \mid q)g_p(q)\,dq + \int_{N(d)} V_1(p^{\emptyset} \mid q)g_p(q)\,dq. \tag{5}$$

For intuition, note that by fostering maximal experimentation, the MEE also exhibits another key feature; it induces the non-disclosure belief closest to $c$ across all experimentation equilibria. Lemma 1 part 1 tells us that this property is desirable to P1, as it minimizes consumption errors made by P3. It is this combination of minimizing errors and maximizing experimentation that renders the MEE welfare optimal.[15] Henceforth, we will therefore focus on the MEE.

**Lemma 2.** *The MEE maximizes $\mathcal{W}(d;p)$, as given by (5), across all equilibrium disclosure rules $d$ and prior beliefs $p \in [c,1)$.*

## 4  U-shaped Disclosure with respect to prior

Theorem 1 tells us that non-disclosure occurs on an interior interval of signals, if at all. But how does this interval depend on primitives, such as P1's prior belief $p$, and the disclosure parameter $\alpha$?

In this section, we show how the equilibrium disclosure threshold $q$ varies with the prior belief $p$. Practically speaking, one can view this exercise as comparing products that differ in how well-established they are. For instance, if $p$ is close to 1 then the product is well-established, while for $p$ close to $c$, the product is close to exit. For $p$ in the interior of this region, products can be viewed as novel. The following result, illustrated in Figure 1, summarizes these findings and constitutes our second main result:
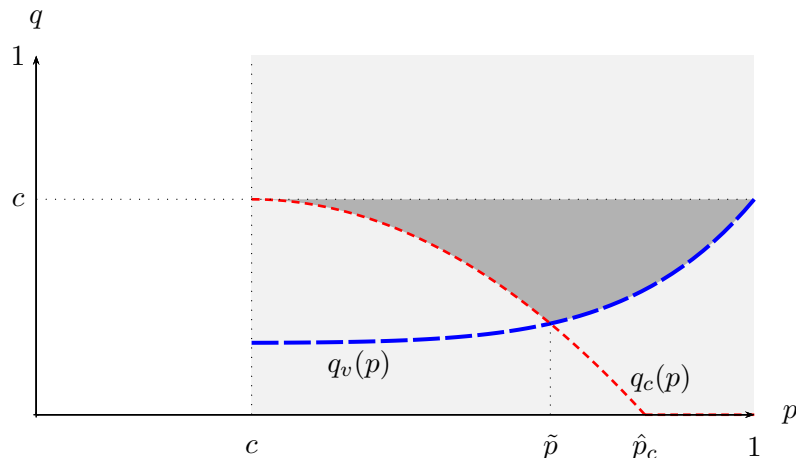
**Theorem 2.** *There exist a unique belief $\tilde{p} \in (c,1)$ and functions $q_c(p), q_v(p) : [c,1] \to [0,c]$, respectively weakly decreasing and strictly increasing, such that:*

*1. If $p \in [c,\tilde{p}]$ then setting $\underline{q} = q_c(p)$ constitutes the MEE.*

---

[15]The proof of Lemma 2 uses some notation introduced for Theorem 2, and as such it can be found after the proof of Theorem 2 in the Appendix.

Figure 1: Equilibrium Non-Disclosure as prior $p$ varies.



Non-disclosure in the MEE. x-axis: prior $p$; y-axis: posterior $q$. Dark-shaded region: non-disclosed posterior beliefs; light-shaded region: disclosed posterior beliefs. Long-dashed line: incentive constraint, $q_v(p)$; short-dashed line: belief constraint, $q_c(p)$. Signals inducing posteriors in the interval $[\underline{q}(p), c)$ are not disclosed, where $\underline{q}(p)$ is the greater of $q_v(p)$ and $q_c(p)$.

2. *If $p \in [\tilde{p}, 1]$, then setting $\underline{q} = q_v(p)$ constitutes the MEE.*

The functions $q_c, q_v$ represent two important constraints on strategic non-disclosure. The first, $q_c(p)$, is a "belief constraint" given by

$$q_c(p) = \inf \{q \in [0,1] \mid \phi(p,q) \geqslant c\}. \tag{6}$$

where

$$(q \leqslant c): \quad \phi(p,q) = \frac{(1-\alpha)p + \alpha \int_{x(p,q)}^{x(p,c)} p^z f_p(z)\, dz}{(1-\alpha) + \alpha \int_{x(p,q)}^{x(p,c)} f_p(z)\, dz}, \quad (q > c): \quad \phi(p,q) = p \tag{7}$$

is the belief P2 forms if they do not see a signal and P1's non-disclosure set is the interval $[x(p,q), x(p,c))$. In words, $q_c(p)$ is the lowest belief such that, if all signals that lead to posteriors in $[q_c(p), c)$ are concealed, the belief $p^\emptyset$ following non-disclosure remains above $c$.

This constraint arises due to a classic form of unraveling: in Theorem 1, we showed that non-disclosure by P1 is positively selected, and so in order to induce experimentation by

P2, P1 cannot conceal signals that are too negative, or else the resultant posterior $p^\emptyset$ would drop below $c$. This constraint tightens as $p$ gets closer to $c$ ($q_c$ is decreasing); the closer is $p$ to $c$, the less room there is for negatively selected non-disclosure to keep $p^\emptyset$ above $c$ and thus induce experimentation. Intuitively, when P1's prior is $c$, even the slightest downgrade in P2's belief will stop experimentation.

The second, $q_v(p)$, is an "incentive constraint" defined by

$$q_v(p) \equiv \inf \left\{ q \in [0, c] \mid W_1\left(\phi(p, q), q\right) = 0 \right\}, \tag{8}$$

where $W_1$ is the "relaxed" value function

$$W_1(r \mid q) \equiv q - c + \alpha \Lambda(r \mid q) + (1 - \alpha)(q - c), \tag{9}$$

which denotes P1's value given a continuation belief $r$ and a private belief $q$, *assuming* that P2 consumes and P3 consumes if they do not see a signal; that is, $V_1(r \mid q) = \mathbb{I}_{r \geq c} W_1(r \mid q)$.

In words, the value $q_v(p)$ tracks the lowest posterior at which P1 is indifferent between disclosing and not, assuming posteriors in $[q_v(p), c)$ are concealed in equilibrium. In proving Theorem 2, we show that $q_v$ is increasing, a result that constitutes yet another expression of the experimentation-accuracy trade-off. Intuitively, when $p$ is close to 1, the continuation belief $r$ is also close to 1, *regardless* of P1's disclosure strategy. As such, non-disclosure plays little role in altering P2's experimentation incentives, and in fact would only induce P3, being less well-informed, to make more mistakes in their consumption choice. Thus, the range of beliefs $q$ just below $c$ at which P1 would prefer to induce experimentation vanishes.

Combining the two constraints, the non-disclosure region is $[\underline{q}(p), c]$, where $\underline{q}(p) \equiv \max\{q_c(p), q_v(p)\}$. Intuitively, if $q_v(p) < q_c(p)$, then signals just greater than $x(p, q_v(p))$ cannot be concealed in equilibrium, despite there being an incentive to do so. For if they were, the resultant non-disclosure belief $p^\emptyset$ would be below $c$, meaning experimentation would fail. Conversely, if $q_v(p) > q_c(p)$, then even though the signal $x(p, q_c(p))$ could be concealed in equilibrium and keep $p^\emptyset \geq c$, there is no incentive to do so.

16

Note that the result relies on incentives that are distinct from the classic unraveling mechanism (Milgrom, 1981; Grossman, 1981; Dye, 1985a). In particular, the classic result does not depend on the prior distribution, whereas it does in our analysis. The key distinction is that in our setting, the sender's preferences over induced posterior beliefs are type-dependent, whereas in the classic setting, each type strictly prefers to induce the highest possible posterior. This type-dependence stems from the misaligned preferences at the heart of the model: for an intermediate range of signals, P1 strictly prefers to induce P2 to hold the lowest belief at which P2 still plays R and thus experiments, as described above. Outside of this range, P1 strictly prefers to have P2 share their belief. This complex pattern of experimentation is critically linked to the ex-post optimality of disclosure; we will show that it fails to hold when agents can commit to a disclosure policy prior to receiving their signal (Kamenica and Gentzkow, 2011).

## 5  Subsidizing Disclosure

Theorem 2 uncovers the complex relation between the degree of non-disclosure and the prior belief $p$, due to the two key constraint functions $q_c, q_v$ working against each other. But how are these constraints themselves determined by the ability to disclose, $\alpha$? The following result answers this question.

**Theorem 3.** *Fix $p \in (c, 1)$. Then there exist $0 < \tilde{\alpha}(p) < 1$ such that $c - \underline{q}$ is strictly increasing for all $\alpha \in [0, \tilde{\alpha}(p))$ and strictly decreasing for all $\alpha \in [\tilde{\alpha}(p), 1)$.*

Intuitively, when $\alpha$ is high, disclosure opportunities abound. This makes it harder for P1 to strategically non-disclose, as there is no room to "hide" behind a lack of disclosure opportunities. In contrast, when $\alpha$ is low, disclosure opportunities are rare, so less information is transmitted and thus beliefs are more persistent. As such, inducing incorrect beliefs comes with a greater chance of consumption errors by P3, which in turn makes non-disclosure less desirable to P1. Thus, for high $\alpha$, P1 is *unable* to induce experimentation through non-disclosure, whereas at low $\alpha$, they are not *willing* to. These counteracting forces result

in the experimentation region being non-monotone in $\alpha$, converging to the empty set when $\alpha$ tends to 0 or 1.

Formally, we first show that the belief constraint $q_c$ is increasing in $\alpha$, converging pointwise to $q_c(p) = c$ for all $p < 1$ as $\alpha \to 1$. In contrast, the incentive constraint $q_v$ is decreasing in $\alpha$, converging pointwise to $q_v(p) = c$ for all $p < 1$ as $\alpha \to 0$.

Taking the view that $1 - \alpha$ captures the fraction of consumers who find it prohibitively costly to leave feedback, Proposition 3 speaks to a commonly proposed intervention, namely that review platforms should attempt to subsidize costly feedback in order to generate more information and help consumers discover high-quality products more efficiently (Marinescu et al., 2021). Our finding cautions against a broad-brush approach to such an intervention, as there exist situations in which such a subsidy (increasing $\alpha$) paradoxically *reduces disclosure* (when $\alpha < \tilde{\alpha}(p)$), as well as others where it stimulates disclosure, but leads to *less experimentation* (when $\alpha \in [\tilde{\alpha}(p), 1)$).

## 6 DISCUSSION

### 6.1 MODEL DISCUSSION

Our model differs from standard social-learning settings in some important ways. First, whereas typically agents' private signals are hidden while their consumption choices are public (Banerjee, 1992; Bikhchandani et al., 1992), here it is the reverse, in the sense that (some) private signals are publicly disclosed.[16] Second, whereas typically agents receive a private signal prior to making a consumption choice, here our agents can only receive their signal if they consume. In our setting, the link between actions and private signal acquisition is crucial; were future agents to receive (and then disclose) private signals regardless of their predecessor's action choices, current agents would never seek to distort these choices

---

[16]Wolitzky (2018) also studies a social learning model with unobservable actions, but with observable outcomes and without strategic behavior. Bowen et al. (2023) study social learning from signals shared on social media, but here again the sharing is non-strategic. Closest is Acemoglu et al. (2024), in which agents decide optimally whether to share pieces of information they receive, but aiming here to maximize re-shares and minimize dislikes by others.

by withholding information, and thus full revelation would be optimal. We thus view our paper as belonging more to the literature on experimentation.

## 6.2 Cheap Talk and Persuasion

As discussed in the introduction, our modeling of feedback as verifiable disclosure is motivated by: (i) the absence of commitment on the part of consumers as to when to leave or not leave feedback, and the large proportion who elect not to; (ii) the fact that many online reviews left by consumers, and especially those that platforms make most salient, convey credible information; (iii) the selection and polarization biases in the reviews that are left.

To understand the role of the lack of commitment we analyze, in Online Appendix B, communication under persuasion. Then, to recognize the fact that many fraudulent or even purchased fake reviews still escape the platforms' scrutiny and regulators' efforts to curb them,[17] we analyze in Online Appendix C the case of cheap-talk, and also compare it with that of disclosure. We summarize below the main results of both cases.

**Persuasion** – We first examine the benchmark wherein P1 can commit to an arbitrary messaging rule prior to receiving their private signal $x$ (Kamenica and Gentzkow, 2011). We take $\alpha = 1$ for simplicity. Recently developed techniques in the persuasion literature allow us to completely characterize the solution (Dworczak and Martini, 2019).

Summarizing the results, communication under persuasion is again both polarized (pooling takes place on an interior interval) and positively selected (the average belief conditional on pooling is $c$, which is less than the prior $p$). In contrast to the disclosure benchmark, however, this pooling interval remains even when the prior $p$ is close to $c$. We refer the reader to Online Appendix B for full details.

**Cheap-Talk** – We next contrast our baseline results with the case of cheap-talk communication. To summarize our findings, assuming again $\alpha = 1$ for simplicity, we find that all equilibria are partitional. This property follows almost directly from the single-crossing properties of $\Lambda$, as summarized in Lemma 1. Furthermore, we show that the partition

---

[17]See for example https://tinyurl.com/2p9s5ehz in the case of the United Kingdom.

of any equilibrium admits finitely many cells on $[x(p,c),1]$, so that signals that generate posteriors above $c$ are pooled into finitely many intervals; see Proposition C.1 in Online Appendix C for details. This characterization is in stark contrast to Theorem 1, as well as to the one under commitment (Lemma B.1 in Online Appendix B), both of which exhibit a "separate-pool-separate" reporting structure.

While the equilibrium outcome is quite different across the three forms of communication, all three share a similar driving force, namely a preference toward biasing "upward". In Theorem 1, this force is what drives types with posteriors $q < c$ to prefer inducing the (higher) belief $c$. Under commitment, it is the same force that induces pooling (Lemma B.1). In contrast, under cheap talk, this same force generates a ripple effect for all higher types, whereby to preserve incentives, information transmission must necessarily be coarse throughout $[x(p,c),1]$ (Proposition C.1). Put differently, we see here an alternative manifestation of the same underlying accuracy-experimentation trade-off identified in Section 3: in order to foster experimentation by P2, equilibrium almost always induces consumption errors by P3. These represent complementary insights to those obtained under disclosure and persuasion; allowing for biased, ex-post optimal reviews greatly undermines the informativeness of consumer feedback, as almost all reviews induce inaccurate beliefs.

Our results on both cheap-talk and persuasion also feature prominently in Smirnov and Starkov (2024), who independently study the same model with both cheap-talk and persuasion. Over three-periods, our results (Proposition C.1 and Lemma B.1) and theirs align. Smirnov and Starkov (2024) also extend their analysis to the infinite horizon; they show how their analog of Lemma B.1 (pooling at intermediate values under persuasion) extends, but are not able to extend that of Proposition C.1 (partitional structure under cheap-talk).

### 6.3 LONGER HORIZONS

The three-period horizon on which we focus provides the simplest, most transparent setting in which strategic (non-)disclosure for purposes of inducing experimentation will arise. At

least three agents are required for the preference structure we employ to induce intertemporal conflict. Furthermore, our three-period model ensures that full disclosure is employed by P2 *independent* of their belief. P1's preferences over what belief to endow P2 with are thus shaped exclusively by their desire to induce experimentation by P2, rather than their desire to shape subsequent disclosure by P2. This feature of equilibrium disclosure greatly simplifies our analysis, and it is unlikely to hold over longer periods. For instance, were we to introduce a fourth agent, Theorem 2 tells us that now, P2's disclosure rule will be non-monotone in (and in particular, not independent of) their belief, implying that P1's disclosure rule must account for both P2's consumption choice as well as their subsequent disclosure rule. This added complexity is not an artifact of finite horizon, non-stationary analyses, and is likely to persist were we to extend our analysis to infinite horizon models.[18] That said, we are confident that certain results, for instance, polarized disclosure and the non-monotone comparative results regarding $\alpha$ will extend to the infinite horizon.[19]

One potential source of this complexity is the richness of the signal space. This was central to developing a theory of polarized disclosure, which requires at least three signal values. Adopting a signal space with only three values might afford sufficiently tractability to extend our model to longer horizons. Of course, these thoughts are speculative, and we leave a rigorous treatment of these avenues for other interested researchers.

## 7  Conclusion and Future Work

We studied a model of strategic information transmission, driven by a tension between selfish consumption and pro-social disclosure. Our analysis sheds light on an important question: when might consumers choose not to leave feedback in order to improve overall welfare? We showed that equilibrium disclosure is necessarily *polarized* and *positively*

---

[18]One could consider a model where each agent discounts exponentially all future payoffs, or an overlapping-generations model in which only the two following agents matter. In such models, it is likely that any stationary Markov perfect equilibrium will feature disclosure that is again non-constant in current beliefs. A full analysis of such extensions is beyond the scope of the current paper.

[19]As long as the sender values accuracy, they will disclose signals that are sufficiently informative. It also seems likely that full disclosure will obtain in the limits as $\alpha \to 0, 1$ over longer horizons, while full disclosure cannot be an equilibrium for interior $\alpha$.

*biased*, two well-established empirical regularities found in consumer reviewing behavior. We further showed that disclosure is hump-shaped with respect to both agents' prior and their opportunities for leaving reviews. This latter result cautions that making feedback less costly could potentially reduce experimentation. Of course, we do not claim that helping future consumers is the sole motivation for leaving feedback. For instance, the survey quoted in the introduction also suggests that expressing one's opinion is another leading factor. Combined with disclosure coming at a cost, this would readily deliver polarized disclosure. Further combining with an exogenous bias toward expressing positive feelings would then deliver positively selected disclosure.[20] Formulating such a theory is undoubtedly an interesting avenue for future research and would complement our paper nicely. We view our analysis as an important first step toward understanding a very natural mechanism for the existence and biases of reviews — that consumers might be pro-socially motivated when leaving feedback.

---

[20]Hui et al. (2024) find empirical evidence for polarized reviews, but that bias can in fact be either positive or negative and depends on the age and reputation of the firm

## REFERENCES

**Abeler, J., A. Becker, and A. Falk**, "Representative evidence on lying costs," *Journal of Public Economics*, 2014, *113*, 96–104.

__ , **D. Nosenzo, and C. Raymond**, "Preferences for Truth-Telling," *Econometrica*, 2019, *87* (4), 1115–1153.

**Acemoglu, D., A. Ozdaglar, and J. Siderius**, "A Model of Online Misinformation," *The Review of Economic Studies*, 2024, *91* (6), 3117–3150.

**Banerjee, A. V.**, "A simple model of herd behavior," *Quarterly Journal of Economics*, 1992, *107* (3), 797–817.

**Bénabou, R. and J. Tirole**, "Self-Confidence and Personal Motivation," *Quarterly Journal of Economics*, 2002, *117* (3), 871–915.

__ **and** __ , "Willpower and Personal Rules," *Journal of Political Economy*, 2004, *112* (4), 848–886.

**Bikhchandani, S., D. Hirshleifer, and I. Welch**, "A theory of fads, fashion, custom, and cultural change as informational cascades," *Journal of Political Economy*, 1992, *100* (5), 992–1026.

__ , __ , **O. Tamuz, and I. Welch**, "Information Cascades and Social Learning," *Working Paper*, 2022.

**Bowen, T.R., D. Dmitriev, and S. Galperti**, "Learning from Shared News: When Abundant Information Leads to Belief Polarization," *The Quarterly Journal of Economics*, 2023, *138* (2), 955–1000.

**Cabral, L. and L.I. Li**, "A Dollar for Your Thoughts: Feedback-Conditional Rebates on eBay," *Management Science*, 2015, *61* (9).

**Carrillo, J. D. and T. Mariotti**, "Strategic Ignorance as a Self-Disciplining Device," *Review of Economic Studies*, 2000, *67* (3), 529–544.

**Che, Y-K. and J. Hörner**, "Recommender Systems as Mechanisms for Social Learning," *Quarterly Journal of Economics*, 2018, *133*, 871–925.

**Dworczak, P. and G. Martini**, "The Simple Economics of Optimal Persuasion," *Journal of Political Economy*, 2019, *127* (5).

**Dye, R.A.**, "Disclosure of Nonproprietary Information," *Journal of Accounting Research*, 1985a, *23*, 123–145.

**Grossman, S. J.**, "The Informational Role of Warranties and Private Disclosure about Product Quality," *The Journal of Law and Economics*, 1981, *24* (3), 461–483.

**He, S., B. Hollenbeck, and D. Proserpio**, "The Market for Fake Reviews," *Marketing Science*, 2022, *41* (5), 461–483.

**Hui, X., T.J. Klein, and K. Stahl**, "Learning from Online Ratings," *working paper*, 2024.

**Jung, W. and Y. Kwon**, "Disclosure When the Market is Unsure of Information Endowment of Managers," *Journal of Accounting Research*, 1988, *26*, 146–153.

**Kamenica, E. and M. Gentzkow**, "Bayesian Persuasion," *American Economic Review*, 2011, *101*, 2590–2615.

**Kremer, I., Y. Mansour, and M. Perry**, "Implementing the Wisdom of the Crowd," *Journal of Political Economy*, 2014, *122* (5), 988–1012.

**Laibson, D.**, "Golden Eggs and Hyperbolic Discounting," *The Quarterly Journal of Economics*, 1997, *112* (2), 443–477.

**Marinescu, Ioana, Andrew Chamberlain, Morgan Smart, and Nadav Klein**, "Incentives can reduce bias in online employer reviews," *Journal of Experimental Psychology: Applied*, 2021.

**Milgrom, P. R.**, "Good News and Bad News: Representation Theorems and Applications," *The Bell Journal of Economics*, 1981, *12* (2), 380–391.

**Nosko, C. and S. Tadelis**, "The Limits of Reputation in Platform Markets: An Empirical Analysis and Field Experiment," *NBER working paper*, 2015.

**Rockafellar, R.T.**, *Convex Analysis*, New Jersey: Princeton University Press, 1997.

**Schoenmueller, V., O. Netzer, and F. Stahl**, "The Polarity of Online Reviews: Prevalence, Drivers and Implications," *Journal of Marketing Research*, 2020, *57* (5), 853–877.

**Slivkins, A.**, "Exploration and Persuasion," in V.V. Vazirani F. Echenique, N. Immorlica, ed., *Online and Matching-Based Market Design*, Cambridge University Press, 2022.

**Smirnov, A. and E. Starkov**, "Designing Social Learning," *working paper*, 2024.

**Smith, L. and P.N. Sørensen**, "Pathological Outcomes of Observational Learning," *Econometrica*, 2000, *68* (2), 371–398.

**_ , _ , and J. Tian**, "Informational Herding, Optimal Experimentation, and Contrarianism," *Review of Economics Studies*, 2021, *8* (5), 2527–2554.

**Wolitzky, A.**, "Learning from Others' Outcomes," *American Economic Review*, 2018, *108*, 2763–2801.

# A  Proofs

## A.1  Proof of Lemma 1

To prove the first part of the lemma, it suffices to note that $\frac{\partial \Lambda(r|q)}{\partial r}$ is proportional to $q^{x(r,c)} - c$, which has the sign of $q - r$, by the MLRP. To prove the first part of the lemma, it suffices to note that $\frac{\partial \Lambda(r|q)}{\partial r}$ has the sign of $q^{x(r,c)} - c$, (since $\partial x(r,c)/\partial r < 0$ by the MLRP), which itself has the sign of $q - r$ by the MLRP and is thus negative for $q < r$. For the second part, basic algebra confirms that:

$$\Lambda(r \mid q) = q(1 - F_H(x(r,c)))(1 - c) + (1 - q)(1 - F_L(x(r,c)))(-c). \tag{A.1}$$

When P1 holds belief $q$ and P2 belief $r$, P1 believes the state is high with probability $q$ and that P3 will consume with probability $1 - F_H(x(r,c))$, receiving a payoff $1 - c$, and similarly when the state is low. The third part follows from

$$\Lambda(c \mid c) = \int_{\hat{x}}^{\bar{x}} (c^z - c) \, f_c(z) \, dz > 0,$$

since $c^z > c$ for $z > \hat{x}$.

## A.2  Proof of Theorem 1

We first leverage Lemma 1 to characterize $V_1$. To do so, it is often convenient to study the "relaxed" value function $W_1(r \mid q)$. From equation (9), it is straightforward to demonstrate that $W_1$ obtains the same properties as $\Lambda$, since $u$ also preserves these properties.

**Lemma A.1.**  *1.  $r \mapsto W_1(r \mid q)$ is strictly increasing on $[c, q)$ and strictly decreasing on $(q, 1]$.*

*2.  $q \mapsto W_1(r \mid q)$ is strictly increasing (and affine) on $[0, 1]$.*

*3.  $W_1(c \mid c) > 0$.*

To establish the result, note first that following a signal leading P1 to hold a posterior $q$, their disclosure decision hinges on the sign of $V_1(q \mid q) - V_1(p^\emptyset \mid q)$. First, we show that disclosure occurs in equilibrium for all signals such that $p^x \geq c$, i.e. for sufficiently high signals.

**Lemma A.2.**  *(Positive Selection) If $p^x \geqslant c$, then $d(x) = 1$ is a strictly dominant strategy.*

*Proof.* The proof proceeds in two cases. Let $q = p^x$. First, suppose that $p^\emptyset < c < q$, so that non-disclosure causes consumption to stop. Then, using (2),

$$V_1(p^\emptyset \mid q) = 0 < u(q \mid q) + \alpha \Lambda(q \mid q) + (1 - \alpha)u(q \mid q) = V_1(q \mid q),$$

since by Lemma 1,

$$\Lambda(q \mid q) \geq \Lambda(c \mid q) \geq \Lambda(c \mid c) = \int_{\hat{x}}^{\bar{x}} (c^z - c) \, f_c(z) \, dz > 0.$$

Next, suppose that $p^\emptyset \geqslant c$, so that non-disclosure leads to consumption (and subsequent disclosure by P2) in spite of a lower belief. In this case,

$$V_1(q \mid q) - V_1(p^\emptyset \mid q) = \alpha \left[ \Lambda(q \mid q) - \Lambda(p^\emptyset \mid q) \right] \geq 0,$$

since the first part of Lemma 1 showed that that $r \mapsto \Lambda(r \mid q)$ is maximized at $q$, for $q \geqslant c$. $\qquad\square$

Lemma A.1 then implies that if non-disclosure occurs in equilibrium, it must take an interval form; $D(d) = [\underline{q}, c]$. Finally, that $\underline{q} > 0$ follows from the fact that $V_1(r \mid 0) < 0 = V_1(0 \mid 0)$ for $r \geq c$; thus, by continuity, revealing is strictly preferred to inducing experimentation for sufficiently low $q$.

## A.3 PROOF OF THEOREM 2

To establish Theorem 2, we will use a series of lemmas characterizing beliefs following non-disclosure and the functions $q_c(p), q_v(p)$. We start with properties of $\phi(p, q)$, which as defined in (7) denotes the continuation public belief if signals in the range $[x(p, q), x(p, c))$ are not disclosed. We showed in Theorem 1 that, if non-disclosure happens, it is over an interval of exactly this type, so $\phi(p, q)$ is indeed the relevant computation for the equilibrium belief $p^\emptyset$ following non-disclosure.

**Lemma A.3.** *1. For $p \geqslant c$, $q \mapsto \phi(p, q)$ is strictly increasing and differentiable on $[0, c]$, with $\phi(p, 0) \in (0, p)$ and $\phi(p, c) = p$.*

*2. For $q \leq c$, $p \mapsto \phi(p, q)$ is strictly increasing on $[c, 1]$.*

*Proof.* For part i), differentiability is clear, and $\partial \phi(p, q) / \partial q$ has the sign of

$$-\frac{\partial x(p, q)}{\partial q} p^{x(p,q)} \left[ (1 - \alpha) + \alpha \int_{x(p,q)}^{x(p,c)} dF_p(z) \right] + \frac{\partial x(p, q)}{\partial y} \left[ (1 - \alpha)p + \alpha \int_{x(p,q)}^{x(p,c)} p^z dF_p(z) \right].$$

27

Given that $q \mapsto x(p, q)$ is increasing by the MLRP and $p^{x(p,q)} \equiv q$, that sign is also that of

$$(1 - \alpha)(p - q) + \alpha \int_{x(p,q)}^{x(p,c)} (p^z - q) \, dF_p(z) > 0,$$

since $q \leq c \leq p$ and $p^z \geq q$ for $z > x(p, q)$. The bounds on $q \mapsto \phi(p, q)$ follow immediately.

For part ii), let us first re-write $\phi(p, q)$ as

$$\phi(p, q) = \frac{(1 - \alpha)p + \alpha \int_{x(p,q)}^{x(p,c)} p^z f_p(z) \, dz}{(1 - \alpha) + \alpha \int_{x(p,q)}^{x(p,c)} f_p(z) \, dz} = \frac{(1 - \alpha)p + \alpha \int_q^c r \, dG_p(r)}{(1 - \alpha) + \alpha \int_q^c dG_p(r)}.$$

Let $a(p) \equiv \int_q^c r \, dG_p(r)$ and $b(p) \equiv \int_q^c dG_p(r)$. By Proposition 4 in Smith et al. (2021), Assumption 2 implies that

$$\frac{d}{dp} \left( \frac{a(p)}{b(p)} \right) > 0.$$

In our case, $P_2$'s not having received a signal may also be due to $P1$ not having had the opportunity to leave feedback, which occurs with probability $1 - \alpha$. As a result, $\partial \phi(p, q) / \partial p$ has the sign of

$$[(1 - \alpha) + \alpha b(p)][(1 - \alpha) + \alpha a'(p)] - [(1 - \alpha)p + \alpha a(p)][\alpha b'(p)]$$

$$= (1 - \alpha)^2 + (1 - \alpha)\alpha a'(p) + \alpha b(p)(1 - \alpha) - \alpha(1 - \alpha)p b'(p) + \alpha^2 \underbrace{(b'(p)a(p) - b(p)a'(p))}_{>0}$$

$$\geqslant \alpha(1 - \alpha)(a'(p) - p b'(p) + b(p)).$$

But

$$a'(p) - p b'(p) + b(p) = \frac{\partial}{\partial p} \int_q^c r \, dG_p(r) - p \frac{\partial}{\partial p} \int_q^c dG_p(r) + \int_q^c dG_p(r)$$

$$= \int_q^c \left[ r \frac{\partial g_p(r)}{\partial p} - p \frac{\partial g_p(r)}{\partial p} + g_p(r) \right] dr = \int_q^c r(g_H(q) - g_L(q)) + g_L(q) \, dr$$

$$= \int_q^c r g_H(q) + (1 - r)g_L(q) \, dr = \int_q^c g_r(r) \, dr > 0,$$

thus proving the claim.

□

We now define and characterize the belief constraint: for $p \geqslant c$, let

$$q_c(p) = \inf \{ q \in [0, 1] \mid \phi(p, q) \geqslant c \}. \tag{A.2}$$

Lemma A.3 and Corollary A.3 tell us that $p \mapsto q_c(p)$ is well-defined on $[c, 1]$. We now establish key properties of the function $q_c(p)$

**Lemma A.4.** *The map $p \mapsto q_c(p)$ is everywhere continuous, with $q_c(c) = c$. Furthermore, there exists $\hat{p}_c \in (c, 1)$ such that: (i) on $[c, \hat{p}_c]$, $q_c(p)$ is strictly decreasing, differentiable and solves $\phi(p, q_c(p)) = c$; (ii) on $[\hat{p}_c, 1]$ $q_c(p) = 0$ .*

*Proof.* Note that $p \mapsto q_c(p)$ is defined as the minimum of a continuous function, $(q \mapsto \phi(p, q))$, on a compact set. Therefore, the infimum is attained by Weierstrass' Theorem, and continuity follows from Berge's Theorem (note that the constraint $\phi(p, q) \geqslant$ defines an upper-hemicontinuous correspondence, since $\phi(p, q)$ is continuous). That $q_c(c) = c$ follows from Corollary A.3, part 2.

Next, that there exists a $\hat{p}_c \in (c, 1)$ such that $q_c(p) = 0$ for all $p \in [\hat{p}_c, 1]$ follows from the definition of $\phi(p, q)$, since as $p \to 1$,

$$\phi(p, 0) \to \frac{(1 - \alpha) \cdot 1 + \alpha \cdot 1}{1 - \alpha + 0} = 1.$$

Finally, that $p \mapsto q_c(p)$ is strictly decreasing on $[c, \hat{p}_c]$ follows directly from Lemma A.3. $\qquad \square$

Having characterized the "belief constraint" $q_c(p)$ bearing on P1's disclosure rule, we next turn to the "incentive constraint" $q_v(p)$.

To begin, we demonstrate that disclosure is strictly optimal after sufficiently extreme signal realizations. We do so by proving a property of the relaxed value function $W_1(r \mid q)$ defined in Section 3.1.

**Lemma A.5.** *For all $p \in [c, 1)$,*

$$\lim_{q \to 0, 1} [W_1(\phi(p, q) \mid q) - V_1(q \mid q)] < 0.$$

*Proof.* The lower limit follows immediately since $\Lambda(r \mid q) < 0$ and $q - c < 0$ for sufficiently small $q$. The upper limit is obtained by noting that as $q \to 1$, $V_1(q \mid q)$ achieves the upper bound on $V_1$. $\quad \square$

Away from these limits, note that the minimization defining $q_v(p)$ is well defined:

**Lemma A.6.** *The map $p \mapsto q_v(p)$ is well-defined, with $q_v(p) < c$ for all $p \in (c, 1)$ and $q_v(1) = c$.*

*Proof.* Note that $W_1(r \mid c) > 0$ for all $r \in (c, 1)$, since $\Lambda(r \mid c) > 0$. Furthermore, for $q$ sufficiently close to $c$, $\phi(p, q) \geq c$, so that $W_1(\phi(p, q) \mid c) > 0$. Thus by continuity, $W_1(\phi(p, q) \mid q) > 0$ for $q$

in some neighborhood below $c$. Lemma A.5 combined with the Intermediate Value Theorem then implies there exists $q' \in (0, c)$ such that $W_1(\phi(p, q') \mid c) = V_1(q' \mid c) = 0$, thus proving that $q_v(p) < c$ for all $p \in (c, 1)$. On the other hand, $W_1(1 \mid c) = 0$, so that by Lemma A.1, $W_1(1 \mid q) < 0$ for all $q < c$, and thus $q_v(1) = c$. $\qquad\square$

**Lemma A.7.** *The map $p \mapsto q_v(p)$ is continuous and strictly increasing.*

*Proof.* For $p \geqslant c$, let $q(p) < c$ be any solution to the equation $W_1(\phi(p, q), q) = V_1(q \mid q)$. From Lemmas A.1 and A.3 and the chain rule, it follows that if $q(p) < c$, then $q'(p) > 0$. Therefore, $p \mapsto q_v(p)$ must be strictly increasing. $\qquad\square$

Taken together, these lemmas immediately show that the two loci $q_c$ and $q_v$ cross at a unique interior point:

**Lemma A.8.** *There exists $\tilde{p} \in (c, 1)$ such that $q_v(p) \leqslant q_c(p)$ if and only if $p \leqslant \tilde{p}$.*

*Proof.* Follows from Lemma A.4 and Lemma A.7 part ii). $\qquad\square$

To complete the proof of the theorem, we proceed in two cases:

1. If $q_v(p) \in [0, q_c(p))$, then setting $\underline{q} = q_c(p)$ defines the MEE. To see this, note first that the equilibrium belief condition (3) is satisfied by definition. Next, we will verify the IC condition (4), which in this case amounts to $V_1(c \mid q) \geq V_1(q \mid q)$ for all $q \in [q_c(p), c)$. But if $q_v(p) \leqslant q_c(p)$ then $\phi(p, q_v(p)) \leqslant \phi(p, q_c(p)) = c$ by (A.3), and so for all $q \in [q_c(p), c)$,

$$V_1(c \mid q) \geq V_1(\phi(p, q_c(p)) \mid q) = W_1(\phi(p, q_c(p)) \mid q) \geq W_1(\phi(p, q_v(p)) \mid q) = 0 = V_1(q \mid q),$$

with the first equality holding since $V_1(r \mid q) = W_1(r \mid q)$ for all $r \geq c$, and the second one holding by Lemma A.1. This verifies incentive compatibility. That $\underline{q}$ defines an EE is then immediate. To verify that this is a MEE, note that were $\underline{q} < q_c(p)$, then one would have $\phi(p, \underline{q}) < c$ and thus no experimentation by P2 could be supported.

2. If $q_v(p) \in [q_c(p), c)$, then set $\underline{q} = q_v(p)$. Again, (3) is satisfied immediately since $q_v(p) \geq q_c(p)$. Next, note that $\underline{q} = q_v(p) \geq q_c(p)$ implies that $\phi(p, \underline{q}) \geq \phi(p, q_c(p))$, and so $W_1(\phi(p, \underline{q}) \mid q) = V_1(\phi(p, \underline{q}) \mid q) \geq 0$ for all $q \in [\underline{q}, c)$. Thus, (4) is verified. Furthermore, since (4) is binding, this must also be a MEE (setting $\underline{q} < q_v(p)$ would violate (4)).

## A.4 Proof of Lemma 2

Consider first any $d \in \mathcal{E}$ that is not the MEE, $d^*$, and which has an associated threshold $q_d$ (characterized by Theorem 1). By definition of the MEE, it must be that $q_d > \underline{q}$ (we suppress the relation of $q$ on $p$ throughout this proof for convenience). We then have that

$$
\begin{aligned}
\mathcal{W}(d^*; p) - \mathcal{W}(d; p) &= \int_{\underline{q}}^{q_d} V_1(\phi(p, \underline{q}) \mid q) g_p(q) \, dq + \int_{q_d}^{c} \left[ V_1(\phi(p, \underline{q}) \mid q) g_p(q) - V_1(\phi(p, q_d) \mid q) \right] g_p(q) \, dq \\
&> \int_{q_d}^{c} \left[ V_1(\phi(p, \underline{q}) \mid q) g_p(q) - V_1(\phi(p, q_d) \mid q) \right] g_p(q) \, dq \\
&> \int_{q_d}^{c} \left[ V_1(\phi(p, \underline{q}) \mid q) g_p(q) - V_1(\phi(p, \underline{q}) \mid q) \right] g_p(q) \, dq \\
&= 0,
\end{aligned}
$$

where the first inequality holds since $\underline{q} \geq q_v(p)$ by construction, and hence $V_1(\phi(p, \underline{q}) \mid q) > 0$ for $q \in (\underline{q}, q_d]$, and the second inequality holds since $r \mapsto V_1(r \mid q)$ is strictly decreasing on $[0, c]$ for $q \in [0, c]$ (Lemma A.1) and $q \mapsto \phi(p, q)$ is strictly increasing (Lemma A.3).

Finally, note that any equilibrium $d$ that is not in $\mathcal{E}$ supports precisely the welfare under full disclosure, as experimentation is not induced with positive probability. Hence,

$$
\mathcal{W}(d^*; p) - \mathcal{W}(d; p) = \int_{\underline{q}}^{1} V_1(\phi(p, \underline{q}) \mid q) g_p(q) \, dq
$$

$$
> 0
$$

by the same reasoning.

## A.5 Proof of Theorem 3

**Lemma A.9.** *1. For fixed $p \in (c, 1)$, there exists $\hat{\alpha}(p) \in (0, 1)$ such that $q_c(p)$ is strictly increasing in $\alpha$ for $\alpha \in [\alpha(\hat{p}), 1]$ and $q_c(p) = 0$ otherwise. Furthermore, $\lim_{\alpha \to 1} q_c(p) = c$.*

*2. For fixed $p \in (c, 1)$, $q_v(p)$ is strictly decreasing in $\alpha$. Furthermore, $\lim_{\alpha \to 0} q_v(p) = c$.*

*Proof.* For part i), note that by Lemma A.3, $q \mapsto \phi(p, q)$ is strictly increasing. Furthermore, from equation (7), $\phi(p, q)$ is strictly decreasing in $\alpha$. Since $q_c(p)$ solves $\phi(p, q_c(p)) = c$, this proves the first claim. For the second part, note that from equation (7), $\lim_{\alpha \to 1} \phi(p, q) = \mathbb{E}(r \mid r \in [q, c])$, and hence $\lim_{\alpha \to 1} \phi(p, c) = c$, while $\lim_{\alpha \to 0} \phi(p, q) = p$, and hence $\lim_{\alpha \to 0} \phi(p, c) = 0$.

31

For part ii), note that $\partial W_1(\phi(p,q) \mid q)/\partial q > 0$, as asserted in Lemma A.7. Next, by Lemma A.3, $q \mapsto \phi(p,q)$ is strictly increasing. Thus, the first part of the claim obtains provided that $\partial W_1(\phi(p,q) \mid q)/\partial \alpha > 0$. To see that such is the case, note that $\partial W_1(r \mid q)/\partial r > 0$ as argued in Lemma A.7, and from equation (7), $\phi(p,q)$ is strictly decreasing in $\alpha$. Finally,

$$\frac{\partial W_1(r \mid q)}{\partial \alpha} = \Lambda(r \mid q) - (q - c) = \int_{x(r,c)}^{\bar{x}} (q^z - c) f_r(z)\, dz \geq 0,$$

since $q^z > c$ for $z > x(r,c)$. Hence

$$\frac{\partial W_1(\phi(p,q) \mid q)}{\partial \alpha} = \frac{\partial W_1(r \mid q)}{\partial \alpha} + \frac{\partial W_1(r \mid q)}{\partial r} \frac{\partial \phi(r,q)}{\partial \alpha} > 0.$$

To prove the second part, we proceed as in Lemma A.6; note that $W_1(p \mid c) > 0$ for all $p \in (c,1)$ and all $\alpha \in (0,1)$, and thus $q_v(p) < c$, while for $\alpha = 0$, $W_1(p \in c) = 0$, so that $q_v(p) = c$. □

Finally, we can put these results together to prove the theorem. For $p \in (c,1)$, let $\tilde{\alpha}(p)$ be that value of $\alpha$ such that $q_v(p) = q_c(p)$. Such a value exists and lies in $(0,1)$ by Lemma A.9, as we have that $q_v(p) = c > q_c(p)$ at $\alpha = 0$. Finally, $\tilde{\alpha}(p) < 1$ since $\tilde{p} < 1$ for all $\alpha \in (0,1)$.

# Online Appendix
## "(Pro)-Social Learning and Selective Disclosure"

Roland Bénabou, Nikhil Vellodi

## A  Heterogeneous Payoffs

We now extend the analysis to allow for heterogeneous payoffs, by introducing an idiosyncratic component to utility. We also relax the restriction to equilibria in which P2 is assumed to fully disclose. Besides adding realism this will serve to show that, under general conditions on the form of this heterogeneity, disclosure is still polarized and positively biased, and that all equilibria are *necessarily* EE's in which P2 strictly prefers to disclose.

Let the payoff to agent $t$ from receiving signal $x$ now be $x\epsilon_t$, where each $\epsilon_t$ is drawn from a distribution $H$, independently from $x$. Without loss of generality, we assume that $\mathbb{E}(\epsilon) = 1$ and $H$ has full support on $[0, \infty)$, with a density $h$ that is everywhere positive. We further assume that the realization of their own $\epsilon_t$ is observable to an agent prior to their consumption decision –e.g., it represents the intensity of their need for such a product– whereas the value $x\epsilon_t$ (or, equivalently, $x$ itself) is revealed only when consumption occurs. Thus $\epsilon_t$ guides the experimentation decision $a_t$, but when $a_t = 1$ the relevant information for the disclosure decision $d_t$ remains $x$ itself, since $\epsilon_t$ is irrelevant to any successor. Formally, consumption rules now map both from beliefs and shocks, i.e. $a_t : [0, 1] \times [0, \infty) \to \{0, 1\}$, while disclosure rules remain as before.

The expected values, from P1's perspective, of subsequent players' consumptions are now:

$$u(r \mid q) \equiv \int_{c/r}^{\infty} (q\epsilon - c) \, dH(\epsilon),$$

$$\Lambda(r \mid q) \equiv \mathbb{E}_{\epsilon,z}(u\left(r^z \epsilon \mid q^z \epsilon\right)) = \int_0^1 \int_{\frac{c}{r^z}}^{\infty} (q^z \epsilon - c) \, dH(\epsilon) dF_q(z).$$

We start with some basic properties of $u$ and $\Lambda$.

**Lemma A.1.** *1. Both maps $r \mapsto u(r \mid q)$ and $\Lambda(r \mid q)$ are strictly maximized at $q$.*

*2. $\Lambda(r \mid q) \geq (>)u(r \mid q)$ for all $r \leq (<)q$.*

*Proof.* Direct calculation verifies that $\frac{\partial u}{\partial r} = -\frac{c^2}{r^2} \left( \frac{q}{r} - 1 \right) h \left( \frac{c}{r} \right)$, which is equal to zero if and only if $q = r$. Since $\Lambda(r \mid q) = \mathbb{E}_{\epsilon,z}(u\left(r^z \epsilon \mid q^z \epsilon\right))$, point 1 is verified. To verify point 2, note that

$$
\begin{aligned}
\Lambda(r \mid q) - u(r \mid q) &= \int_0^1 \int_{\frac{c}{r^z}}^\infty (q^z \epsilon - c) \, dH(\epsilon) dF_q(z) - u(r \mid q) \\
&= \int_0^{\hat{x}} \int_{\frac{c}{r^z}}^\infty (q^z \epsilon - c) \, dH(\epsilon) dF_q(z) \\
&\quad + \underbrace{\int_{\hat{x}}^1 \int_{\frac{c}{r^z}}^\infty (q^z \epsilon - c) \, dH(\epsilon) dF_q(z)}_{\geq u(r|q)} - u(r \mid q) \geq 0,
\end{aligned}
$$

where the last inequality holds because $z \mapsto q^z \epsilon - c$ is positive on the range $[c/r^z, \infty)$ by the MLRP, since by assumption $r \leq q$. $\qquad\square$

Note that since $V_2(r \mid q) = u(r \mid q)$, Lemma A.1 implies that full disclosure by P2 is a *strictly* dominant strategy. In Section 2, P2 was indifferent over posterior beliefs that induce the same action by P3. Now, greater accuracy leads to a strictly lower chance of erroneous consumption choices by P3 due to idiosyncratic shocks.

As before, this allows us to simplify player 1's value function,

$$
V_1(r \mid q) = u(r \mid q) + \left( \alpha C(r)\Lambda(r \mid q) + (1 - \alpha + \alpha \left(1 - C(r)\right) u(r \mid q) \right), \tag{A.1}
$$

where

$$
C(r) \equiv \int_{c/r}^\infty dH(\epsilon)
$$

is the probability of consumption given a prior belief $r$, prior to the realization of $\epsilon$.

2

### A.0.1   SELECTED DISCLOSURE

In order to draw comparison to the results in Sections 3 and 4, we first adapt the definition of experimentation equilibria in the most natural manner. Now, let

$$X_E(\sigma) = \{x \in N_1(d) \mid a_2(p^\emptyset, \epsilon) > a_2(p^x, \epsilon) \quad \forall \epsilon \in [0, \infty)\}$$

denote the experimentation set for an equilibrium $\sigma$. First, we recover the result of polarized disclosure.

**Lemma A.2.** *(Polarized disclosure) Fix $r \in (0, 1)$. Then, $\lim_{q \to 0,1} [V_1(q \mid q) - V_1(r \mid q)] > 0$.*

*Proof.* For the lower limit, note that

$$u(r \mid 0) = \int_{c/r}^\infty -c \, dH(\epsilon) < 0, \quad \Lambda(r \mid 0) = \int_0^1 \int_{c/r^z}^\infty -c \, dH(\epsilon) dF_q(z) < 0,$$

whereas $u(0 \mid 0) = \Lambda(0 \mid 0) = 0$. Thus, by the expression for $V_1(r \mid q)$ given in (A.1), $V_1(r \mid q) < 0 = V_1(q \mid q)$. For the upper limit, note that

$$\Lambda(r \mid 1) = \int_0^1 \int_{\frac{c}{r^z}}^\infty (\epsilon - c) \, dH(\epsilon) dF_q(z),$$

which is strictly increasing in $r$ by the MLRP, since the integrand is strictly positive. Similarly, $r \mapsto u(r \mid 1)$ is strictly increasing. Finally, $r \mapsto C(r)$ is also strictly increasing, and thus so is $r \mapsto V_1(r \mid 1)$. Therefore, the claim is verified. $\square$

Next, we demonstrate that for any prior $p \in (0, 1)$, any posterior $q \geqslant p$ (i.e. any signal $x \geqslant \hat{x}$) is disclosed by P1. Note that whereas in the baseline model (Lemma A.2) it was dominant for all posteriors $q \geq c$ to be disclosed, here this is no longer necessarily the case.

**Lemma A.3.** *If $p^x \geqslant p$, then $d_1(x) = 1$ is a strictly dominant strategy.*

*Proof.* Suppose not, so that there exists an $x > \hat{x}$ such that $d_1(x) = 0$. Take the largest

such $x$ and let $q = p^x$. By construction, to satisfy the equilibrium belief condition (3) it must be that $p^\emptyset < q$. But then by Lemma A.1, $u(p^\emptyset \mid q) < u(q \mid q)$ and $\Lambda(p^\emptyset \mid q) < \Lambda(q \mid q)$, while we also have $C(p^\emptyset) < C(q)$ as $r \mapsto C(r)$ is strictly increasing. Combining, we have that $V_1(p^\emptyset \mid q) < V_1(q \mid q)$. $\qquad\square$

Finally, we prove that non-disclosure of signals that convey marginally bad news (namely, such that the posterior $p^x$ is just below the prior $p$) is optimal. This result has no direct analog in the baseline model, insofar as non-disclosure now occurs at signals the revelation of which would have induced consumption with strictly positive probability ($c < p^x < p$).

**Lemma A.4.** *(Positive selection) Let $\tilde{V}_1(q) \equiv V_1(r \mid q)$. Then $\tilde{V}_1'(q) > \frac{\partial V_1}{\partial q}|_{r=q}$.*

*Proof.* Since $\tilde{V}_1'(q) = \partial V_1(r \mid q)/\partial r|_{r=q} + \partial V_1(r \mid q)/\partial q|_{r=q}$, the claim is equivalent to proving that $\partial V_1(r \mid q)/\partial r|_{p=q} > 0$. But

$$
\begin{aligned}
\frac{\partial V_1|(r \mid q)}{\partial r}|_{r=q} = & \underbrace{\frac{\partial u}{\partial r}_{r=q}}_{=0} + \frac{c}{q^2} h\left(\frac{c}{q}\right) \left[\Lambda(q \mid q) + (1 - C(q))u(q \mid q)\right] \\
& + C\left(\frac{c}{q}\right) \left[\underbrace{\frac{\partial \Lambda}{\partial r}|_{r=q}}_{=0} + (1 - C(q))\underbrace{\frac{\partial u}{\partial r}|_{r=q}}_{=0} - C'(q)u(q \mid q)\right] \\
= & \Lambda(q \mid q) - C\left(\frac{c}{q}\right) u(q \mid q)\frac{c}{q^2} h\left(\frac{c}{q}\right) > 0,
\end{aligned}
$$

where the last inequality holds because $C(q \mid q) < 1$ and $\Lambda(q \mid q) > u(q \mid q)$. $\qquad\square$

In particular, for $x = \hat{x} - \varepsilon$ where $\varepsilon$ is small, non-disclosure is optimal. Combining Lemmas A.3 and A.4 with a continuity argument yields that non-disclosure takes place in (at least) some interval $[\hat{x} - \varepsilon, \hat{x})$, and thus disclosure is positively biased. Furthermore, we have:

**Lemma A.5.** *Any equilibrium is an EE.*

*Proof.* To see that all equilibria admit a non-empty experimentation region, note that Lemmas A.2 and A.3 imply that in any equilibrium $\sigma$, for each $p$ there exists a minimal posterior

$\underline{q}(p) < p$ that is concealed. Continuity of $r \mapsto V_1(r \mid q)$ then ensures the existence of a $\delta > 0$ such that posteriors in the interval $[\underline{q}(p), \underline{q}(p) + \delta)$ are concealed. But for $\delta$ sufficiently small, $\underline{q}(p) + \delta < p$, and so $[\underline{q}(p), \underline{q}(p) + \delta) \subset X_E(\sigma)$. □

## B  OPTIMAL FEEDBACK – PERSUASION

We now turn to the benchmark wherein P1 can commit to an arbitrary messaging rule prior to receiving their private signal $x$ (Kamenica and Gentzkow, 2011). Formally, P1 chooses an *information structure*, consisting of a message space $\mathcal{S}$ along with a collection of conditional probabilities $(\pi(\cdot \mid x))_{x \in X}$, where $\pi(s \mid x)$ denotes the likelihood of P1 sending the message $s$ given that they received signal $x$. Let $\mathcal{M} = X \cup \{\emptyset\}$ denote the (rich) message space that naturally associates messages with outcomes, as well as a privileged message $\emptyset$ that denotes no signal reported. We may take $\mathcal{S} = \mathcal{M}$. Since communication is no longer constrained to be verifiable, we can set $\alpha = 1$ without loss of generality. Contrasting this case with that of hard-evidence disclosure will thus shed light on how ex-post IC constraints shape optimal feedback. Recently developed techniques in the persuasion literature allow us to completely characterize the solution (Dworczak and Martini, 2019). Denote $V_1(q \mid q)$ by $V_1(q)$ for simplicity. The following result is illustrated in Figure 2.

**Theorem B.1.** *There exist $\underline{q}^*(p) < c < \bar{q}^*(p)$ such that the solution to the persuasion problem takes the following form: reveal $x$ if either $p^x < \underline{q}^*(p)$ or $p^x \geqslant \bar{q}^*(p)$, and pool all $x$ such that $p^x \in [\underline{q}^*(p), \bar{q}^*(p))$. Furthermore, $\underline{q}^*(p), \bar{q}^*(p)$ solve*

$$\mathbb{E}_p(q \mid q \in [\underline{q}^*(p), \bar{q}^*(p))) \equiv \frac{\int_{\underline{q}^*(p)}^{\bar{q}^*(p)} q \, dG_p(q)}{\int_{\underline{q}^*(p)}^{\bar{q}^*(p)} dG_p(q)} = c, \tag{B.1}$$

*and*

$$\frac{V_1(\bar{q}^*(p))}{V_1(c)} = \frac{\bar{q}^*(p) - \underline{q}^*(p)}{c - \underline{q}^*(p)}. \tag{B.2}$$

*Proof.* Since $q \mapsto V_1(r \mid q)$ is affine, standard arguments imply that the problem faced by

P1 under commitment is to solve

$$v^*(p) = \max_{H \in \Delta([0,1])} \int_0^1 V_1(q) \, dH(q), \tag{B.3}$$

subject to the constraint that $H$ is a mean-preserving contraction of $G_p$ (Kamenica and Gentzkow, 2011). First, we prove that $V_1(q)$ is convex on $[c, 1]$. To see this, note that Lemma 1 implies that $V_1(q \mid q) = \sup_{r \in [0,1]} V_1(r \mid q)$ for $q \in [c, 1]$, and that $q \mapsto V_1(r \mid q)$ is affine. The convexity of $V_1(q)$ then follows from standard results in convex duality (Rockafellar, 1997, Theorem 13.2).

We may now apply (Dworczak and Martini, 2019, Theorem 1). In particular, consider the function $\psi$ defined by

$$\psi(q) = \begin{cases} V_1(q) & \text{if} \quad p^x \geqslant \bar{q}^*(p) \\ V_1(c) \left( \frac{q - \underline{q}^*(p)}{c - \underline{q}^*(p)} \right) & \text{if} \quad p^x \in [\underline{q}^*(p), \bar{q}^*(p)) \\ V_1(q) & \text{if} \quad p^x < \underline{q}^*(p), \end{cases}$$
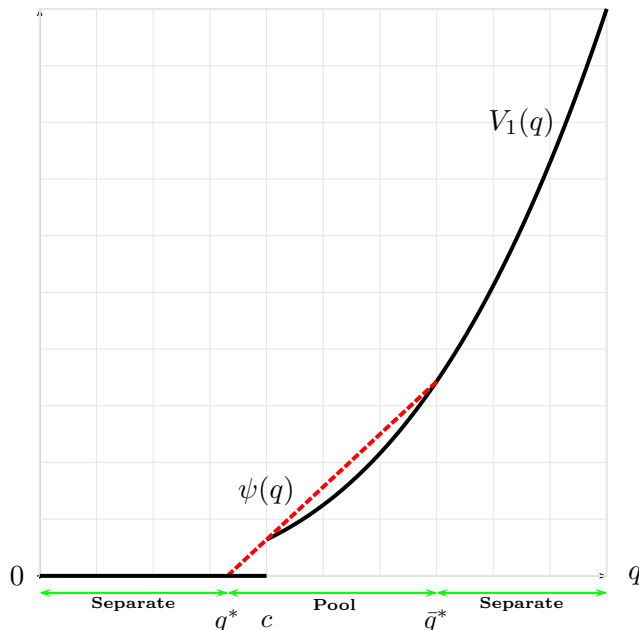
and the distribution $H_p : [0, 1] \to [0, 1]$ defined by

$$H_p(q) = \begin{cases} G_p(q) & \text{if} \quad p^x \geqslant \bar{q}^*(p) \\ G_p(c) + \mathbb{I}_{q \geq c}[G_p(\bar{q}^*(p)) - G_p(\underline{q}^*(p))] & \text{if} \quad p^x \in [\underline{q}^*(p), \bar{q}^*(p)) \\ G_p(q) & \text{if} \quad p^x < \underline{q}^*(p), \end{cases}$$

which reveals $q$ when either $q \geq \bar{q}^*(p)$ or $q \leq \underline{q}^*(p)$ and pools otherwise. It is readily verified that $\psi$ and $H$ together satisfy conditions 3.1-3.3 of (Dworczak and Martini, 2019, Theorem 1), and thus constitute a solution to the commitment problem. Finally, note that since $q \mapsto G_p(q)$ is continuous and strictly increasing, so too are $\underline{q}^*(p), \bar{q}^*(p)$. $\qquad\square$

Communication under persuasion is also both polarized (pooling takes place on an interior interval) and positively selected (the average belief conditional on pooling is $c$, which

6

Figure 2: Persuasion Solution

Disclosure under commitment. Value function $V_1(q) \equiv V_1(q \mid q)$: solid black lines. $\underline{q}^*, \bar{q}^*$ are determined by both $\mathbb{E}(q \mid q \in [\underline{q}^*, \bar{q}^*)) = c$ and lying on a straight-line segment $\psi(q)$ (dotted red) intersecting $V_1(q)$ at $\underline{q}^*, \bar{q}^*$ and $c$.

is less than the prior $p$). In contrast to the disclosure benchmark, however, this pooling interval remains even when the prior $p$ is close to $c$. Crucially, under persuasion, P1 can "pool down" by pooling posteriors above $c$ with those below $c$, while still averaging to $c$ (equation (B.1)). This allows them to maintain a positive-measure pooling interval as the prior $p$ converges to either $c$ or 1. In contrast, under disclosure such pooling down cannot occur, since P1 finds it ex-post optimal to disclose (separate) at all posteriors above $c$. This heavily constrains their ability to not disclose at posteriors below $c$. The logic presented here highlights the role of ex-post optimality (equation (4)) that disclosure rules must satisfy in shaping optimal feedback.

**Corollary B.1.** *Both $\underline{q}^*(p)$ and $\bar{q}^*(p)$ are strictly decreasing in $p$. Furthermore, $\lim_{p \to c,1} \underline{q}^*(p) < c < \lim_{p \to c,1} \bar{q}^*(p)$.*

*Proof.* Note that the constraint (B.2) is independent of $p$, whereas a simple application of

the posterior monotonicity property (Proposition 4 in Smith et al. (2021)) implies that for fixed $\underline{q}, \bar{q}, \mathbb{E}_p(q \mid q \in [\underline{q}, \bar{q}))$ is strictly increasing in $p$. Thus, to keep $\mathbb{E}_p(q \mid q \in [\underline{q}, \bar{q}))$ fixed, we must lower both $\underline{q}$ and $\bar{q}$. The final part of the corollary follows by noting that $V(q)$ is strictly increasing and convex for $q \geq c$ and strictly positive at $c$, and thus for all $p \in [c, 1]$ the line segment intersecting the three points $(\underline{q}^*(p), 0), (c, V_1(c)), (\bar{q}^*(p)$ and $V_1(\bar{q}^*(p)))$ can only exist if $\underline{q}^*(p) \neq \bar{q}^*(p)$, while the constraint that $\mathbb{E}_p(q \mid q \in [\underline{q}^*(p), \bar{q}^*(p))) = c$ further implies that $\underline{q}^*(p) < c < \bar{q}^*(p)$. $\qquad\square$

Finally, notice that the persuasion outcome — which did not assume information to be verifiable – can be implemented via commitment to the verifiable disclosure rule

$$
d(x) = \begin{cases} 1 & \text{if} \quad p^x \geqslant \bar{q}^*(p) \\[2mm] \emptyset & \text{if} \quad p^x \in [\underline{q}^*(p), \bar{q}^*(p)) \\[2mm] 1 & \text{if} \quad p^x < \underline{q}^*(p). \end{cases}
$$

This is due to the simple structure of optimal persuasion; it is not only monotone partitional (Dworczak and Martini, 2019), but includes only one pooling region (see Figure 2). Thus, the pooling region can be interpreted as non-disclosure and the separating regions as disclosure, satisfying the verifiability assumption. In this sense, the benefit of persuasion over (ex-post) verifiable disclosure comes directly from which posteriors (signals) are credibly concealed, rather than the communication language itself.

## C  Biased Feedback – Cheap Talk

We now consider a natural variant on our baseline model, by relaxing the requirement of hard evidence disclosure and instead permitting arbitrary message reporting (cheap talk). Such a variant is important for several reasons. First, in many applied settings, it might not only be feasible but strategically optimal for consumers to misreport their experiences. The hard-evidence baseline abstracts from this possibility, thus providing a useful benchmark;

even when fake reviews are impossible, might there be scope for strategic disclosure? In this section, we explore the extent of strategic information transmission when lying is both feasible and costless. Second, by studying an alternative, well-established form of equilibrium information transmission, we make clear the features of strategic disclosure that are invariant to the information-sharing technology available to agents.

Specifically, we endow each agent with a rich messaging space $\mathcal{M} = [0,1] \times \{\emptyset\}$ that allows not only for full separation but also for agents to send a privileged message that pools with non-arriving consumers, so that messaging rules (previously, disclosure rules) are now mappings $d_t : X \times [0,1] \to \mathcal{M}$.[21] Again, full transparency is dominant for P2, so we focus on P1's messaging strategy. Let $r^*(m)$ denote P2's equilibrium belief upon observing message $m$. Then the IC constraint (4) is replaced with the condition

$$d(x) \in \underset{m \in \text{supp}(d)}{\arg\min} V_1(r^*(m) \mid p^x). \qquad (\text{C.1})$$

We focus on the case where $\alpha = 1$. Combining various insights learned through the baseline analysis, we summize that all equilibria must admit a partitional structure. The proof of Theorem C.1 is constructive. First, we identify a lower-bound on the degree of experimentation possible; there exists a $q_{min}$ such that $V_1(c \mid q_{min})$, thus any type lower prefers to terminate experimentation, regardless of the continuation belief $r$. Each equilibrium is then determined by its associated $q \in [q_{min}, c]$, that is the lowest type whose message induces experimentation. See Figure 3 for a graphical illustration of this construction.

**Theorem C.1.** *All equilibria are partitional. That is, for all $r \in [0,1]$ induced in equilibrium, the set of $q$ in which $r$ is induced forms an interval in $[0,1]$. Furthermore, there must be at most finitely many such intervals on $[c,1]$.*

*Proof.* We proceed with a series of lemmas.

---

[21] We focus on pure-strategy equilibria for simplicity, noting the usual implementation via uniform randomization in cheap-talk games

**Lemma C.1.** *All equilibria are partitional. Furthermore, there must be at most countably infinitely many such intervals on $[c, 1]$.*

*Proof.* Lemma 1 tells us that $r \mapsto V_1(r \mid q)$ is maximized at $r = q$, and that $q \mapsto V_1(r \mid q)$ is strictly increasing, so that $\arg\max_{r \in [0,1]} V_1(r \mid q)$ is strictly increasing in $q$, which proves the first claim. To prove the second one, we argue that there can be no interval in $[c, 1]$ on which separation can occur. Suppose there were, and take the lowest such interval $[q_1, q_2], q_1 \leq q_2$. If $q_2 < 1$, then we claim that types $q \in (q_2 + \varepsilon]$ have an incentive to pool with $q_2$. For since this was the lowest separating interval, it must be that types $q \in (q_2 + \varepsilon]$ induce a belief $\hat{q} = q_2 + \delta, \delta > 0$. By Lemma A.1, $V_1(\hat{q}_2 \mid q_2 + \varepsilon) < V_1(q_2 \mid q_2 + \varepsilon) \approx V_1(q_2 \mid q_2) + \varepsilon V_1'(q_2 \mid q_2)$ for small enough $\varepsilon > 0$. If $q_2 = 1$, then we claim that $q \in (q_1 - \varepsilon]$ have an incentive to pool with $q_1$ by analogous reasoning. $\square$

**Lemma C.2.** *There exists $q_{min} < c$ such that full revelation is weakly dominant for all types $q \in [0, q_{min})$.*

*Proof.* $q_{min}$ is the unique root of $q \mapsto V_1(c \mid q)$ on $[0, c]$, which is well-defined since the map is continuous, strictly increasing with $V_1(c \mid 0) < 0 < V_1(c \mid c)$. $\square$

It is thus without loss to associate an equilibrium with a lowest type $q > 0$ that forms part of a pooling interval that itself induces experimentation. More specifically, combining with Lemma C.1, an equilibrium can be described by a (possibly infinite) sequence ($q \equiv q_0 < q_1 < q_2 < \ldots$ such that types in $[q_i, q_{i+1})$ pool and $\hat{q} \equiv \mathbb{E}(q \mid q \in [q, q_1)) \geq c$. More generally, we denote $\hat{q}_{i+1} = \mathbb{E}(q \mid q \in [q_i, q_{i+1}))$.

We next prove that the two first intervals $[q, q_1), [q_1, q_2)$ cannot be "too small" as types just below $q$ would then profitably deviate by pooling with $[q_1, q_2)$ to induce $\hat{q}_1$.

**Lemma C.3.** *For all $q \in [q_{min}, c]$ there exists $\hat{q}_{2,min} > c$ such that in any equilibrium, $\hat{q}_2 \geq \hat{q}_{2,min}$.*

*Proof.* If not, then for any $\varepsilon > 0$ there exists an equilibrium with $\hat{q}_2 \leq c + \varepsilon$. But since by definition $\hat{q} \geq c$, it must be that $\underline{q} > \underline{q}_{min}$ for sufficiently small $\varepsilon$, and by the sandwich theorem, $V_1(\hat{q} \mid \underline{q}) > 0$, violating the IC constraint at $\underline{q}$. $\square$

**Lemma C.4.** *All equilibrium partitions essentially admit at most finitely many intervals covering $[c, 1]$.*
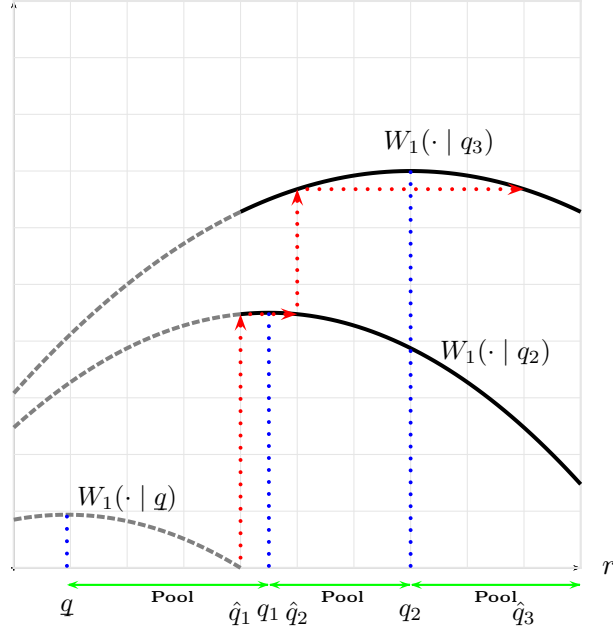
*Proof.* We proceed constructively, via the following algorithm:

1. Fix a $\underline{q} \geq \underline{q}_{min}$. Compute $\hat{q}_{max} \equiv \mathbb{E}_p(q \mid q \in [\underline{q}, 1])$.

   (a) If $V_1(\hat{q}_{max} \mid \underline{q}) > 0$, then $N^*(\underline{q}) = 0$ and $\underline{q}$ cannot be implemented in equilibrium.

   (b) If not, then there exists a unique $q_1 > c$ such that $V_1(\hat{q}_1 \mid \underline{q}) = 0$, where $\hat{q}_1 \equiv \mathbb{E}(q \mid q \in [\underline{q}, q_1])$. (Such a value exists by continuity and strict monotonicity of $r \mapsto V_1(r \mid \underline{q})$ on $[\underline{q}, 1]$, the Intermediate Value Theorem and because $V_1(c \mid \underline{q}) > V_1(c \mid \underline{q}_{min}) = 0$ by Lemma C.2).

2. Compute $V_1(\hat{q}_1 \mid q_1)$.

   (a) If $V_1(1 \mid q_1) \geq V_1(\hat{q}_1 \mid q_1)$, then $N^*(\underline{q}) = 1$.

   (b) If not, then there exists a unique $q_2 > q_1$ such that $V_1(\hat{q}_2 \mid q_1) = V_1(\hat{q}_1 \mid q_1)$, where $\hat{q}_2$ is analogously defined, and $q_2$ exists by the same reasoning as $q_1$.

3. Repeat from step 2.

Finally, we argue that this algorithm terminates in finitely many steps. Suppose not. Then for all $\varepsilon > 0$, there exists an equilibrium and an interval $[q_i, q_{i+1}) \subset [c, 1]$ such that $q_{i+1} - q_i \leq \varepsilon$. Without loss, assume equality, and further assume that $[q_i, q_{i+1})$ is the lowest such interval (this is possible due to Lemma C.3). Let $\hat{q}_{i+1} = \mathbb{E}(q \mid q \in [q_i, q_{i+1}))$. Then there exists $\delta(\varepsilon) < \varepsilon$ such that $\hat{q}_{i+1} - q_i = \delta(\varepsilon)$. The Mean Value Theorem implies that

$$V_1(\hat{q}_i \mid q_i) - V_1(q_i \mid q_i) = V_1'(\varphi_1 \mid q_i)(\hat{q}_1 - q_i),$$

11

Figure 3: Cheap-Talk Equilibria: Construction



An equilibrium with three pooling intervals covering $[c, 1]$. Relaxed value function $W_1(r \mid q)$. For $r \geq c$, $W_1(\cdot \mid q) = V_1(\cdot \mid q)$ (solid black lines). For $r < c$, $V_1(\cdot \mid q) = 0$.

for some $\varphi_i \in (\hat{q}_i, q_i)$. But since $r \mapsto V_1(r \mid q)$ has a global maximum at $q$, we know that

$$V_1(q_i + \delta(\varepsilon) \mid q_i) - V_1(q_i \mid q_i) \approx \frac{\partial^2 V_1}{\partial r^2}(q_i \mid q_i)\delta(\varepsilon)^2.$$

Combining these terms implies that $q_i - \hat{q}_i = \kappa\delta(\varepsilon)$, for some $\kappa > 0$, and so $\hat{q}_{i+1} - \hat{q}_i = (\hat{q}_{i+1} - q_i) + (q_i - \hat{q}_i) = \kappa_i\delta(\varepsilon)$, for some $\kappa_i > 0$. Now, since $\varepsilon > 0$, there exists a finite $I > 0$ such that $q_{i-I} \leq \underline{q}$ (if not, then Lemma C.3 is violated) and thus a simple inductive argument implies that $\hat{q}_1 - \hat{\underline{q}} = \kappa_{i-I}\delta(\varepsilon)$, for some $\kappa_{i-I} > 0$. Taking $\varepsilon$ (and thus $\delta(\varepsilon) < \varepsilon$) sufficiently small violates Lemma C.3. $\qquad\square$

$\square$

12