



# Incentives for prosocial behavior under reputation persistence and policy lags

Francisco Candel-Sánchez<sup>1</sup> · Juan Perote-Peña<sup>2</sup>

Received: 19 June 2021 / Accepted: 20 March 2023  
© The Author(s) 2023

## Abstract

In this paper we show that a policy based on incentives to promote prosocial activities can be counterproductive in a context where the agents' reputation exhibits persistence over time and there exists a time lag between announcement of the policy and implementation. Reputation persistence in our model means that the reputation gained in past periods constrains the possibilities of changing reputation in the future. We present a two-period model in which agents use prosocial activities to signal their degree of altruism. If a subsidy is established for the second period, the set of agents that undertake social activities in that period enlarges. This worsens the reputation of the most altruistic agents, some of whom then react by decreasing their involvement in prosocial activities in the first period. We identify a condition under which subsidies cause a decrease in the global supply of pro-social activities.

**Keywords** Signaling · Incentives · Pro-social behavior · Corporate social responsibility

**JEL Classification** D82 · H25 · M14

## 1 Introduction

In recent years, companies are increasingly engaged in activities such as environmental protection, social inclusion or the development of business ethics. As part of the sustainable development agenda, governments are concerned with the conduct of a company's administration with respect to its stakeholders and society as a whole. Governments also support social behavior in NGOs that seek to promote actions for the welfare of others, and design policies to incentivize the philanthropic behavior

---

✉ Francisco Candel-Sánchez  
fcandel@um.es

<sup>1</sup> Universidad de Murcia, Murcia, Spain

<sup>2</sup> Universidad de Zaragoza, Zaragoza, Spain

of individuals who donate to charities or volunteer. The effectiveness of material incentives as a device to induce prosocial behavior has been questioned on several grounds (see Gneezy et al 2011). This paper deals with the effects of incentives on social activities in a dynamic environment where agents are concerned with reputation. Our main finding is that, if reputation exhibits some degree of persistence and there exist policy lags, establishing incentives to encourage social behavior may prove ineffective.

The formal framework depicted here is especially suited for (but not restricted to) the analysis of companies' Corporate Social Responsibility (CSR hereinafter). For the sake of exposition, we will conduct our analysis through this particular example throughout the paper, and hence the agents to be incentivized will be named as firms and the prosocial actions will be referred to as CSR activities. The nature of CSR activities as privately provided public goods might justify governmental intervention to mitigate the free-rider problem (the public good aspects of CSR are analyzed in Bagnoli and Watts (2003), Besley and Ghatak (2007), and Kotchen (2006)). However, the role of public policy with regard to CSR is still far from having reached a clear path of action.

Governments and public institutions in many countries are firmly committed to create a favorable environment for national companies to engage in CSR activities. For instance, the British Government uses tax incentives to ensure the socially responsible investment of public pension funds, and also disburses funds for CSR projects and organizations. The Brazilian National Economic Development Bank provides funds for entities that comply with national ethical labor codes. In the same fashion, the Ministry of Economics in Mexico created in 2004 the "Support Fund for Micro, Small, and Medium Enterprises", one of whose goals is to implement CSR programs in small and medium enterprises (SMEs). The Swedish Business Development Agency has also disbursed funds to SMEs to encourage CSR activities. The Dutch Government conditions the disbursement of funds to companies being familiar with the OECD guidelines for Multinational Enterprises. These are just a few examples of ongoing policy actions that leave little doubt about the interest of governments to encourage CSR activities. But, on which theoretical grounds are based such public policies? We claim here that, in a dynamic environment with reputation persistence and policy lags, the strategic response that subsidies trigger on firms may reduce companies' aggregate engagement in CSR.

Companies implement CSR policies, *motu proprio*, to show a positive attitude toward society and the environment. To the extent that social activities improve the image of the company and its reputation, CSR is an instrument to attract and retain consumers. In fact, a two decades survey for American consumers reveals that CSR has a positive impact on brand reputation, loyalty and affinity (see <http://www.conecomm.com/2017-CSR-Study>). In 2017, respondents of this survey "have a more positive image (92%), are more likely to trust (87%) and are more loyal (88%) to companies that support social and environmental issues". Consumers are willing to pay a premium for the provision of credence attributes, and, as pointed out in Baron (2011), firms have an incentive to form an organization to assure that credence attributes of goods are being supplied.

A substantial strand of the literature on CSR has been devoted to answer the question of why do firms engage in socially responsible behavior (see Kitzmueller and Shimshack (2012) for a synthesis of the role of preferences and economic justification of CSR). Here we adopt the view that firms engage in “profit maximizing” CSR, as they use their pro-social actions instrumentally as a means to maximize profits. Consumers do not know firms’ degree of altruism (types), and firms undertake CSR activities to send a credible signal of their types.

Investing in CSR is worthwhile as long as this investment increases the company’s profits in the future. For instance, “green” consumers prefer to buy “environmentally friendly” goods, provided by socially responsible companies (Arora and Gangopadhyay 1995). In a similar vein, McWilliams and Siegel (2011) argue that CSR contributes to firms’ sustainable comparative advantage. Thus, firms’ reputation can be viewed as an asset whose returns persist over time. In this respect, Roberts and Dowling (2002) find a positive impact of reputation on the path of future financial performance. Studies by Ang and Wight (2009) and Schultz et al (2001) show that reputation is “sticky”, meaning that it is “durable and tends to reproduce itself over time”. Based on the evidence provided in the mentioned articles, we assume that reputation exhibits some degree of persistence over time.

A standard assumption in the literature of CSR is that investments in social activities are valuable to stakeholders, who might reward this pro-social behavior with a preference for firms that engage in these activities, therefore having a positive impact on firms’ profits (for a survey of the most relevant approaches on this issue, see Bénabou and Tirole (2010), for instance). We adopt a slightly different view. In particular, we assume that consumers are concerned about the estimated degree of altruism of the firms they interact with rather than with the CSR actions themselves. Many experiments in social psychology (see Berg et al 1995) show that agents tend to reciprocate with each other. Consumers value positively those companies that are seen as altruistic and hence are more prone to buy from them. In this line, a study by Sen and Bhattacharya (2001) analyzes consumers’ reactions to companies’ behavior related to CSR. Firms are aware of the effect that CSR actions cause on consumers’ decisions and then use these actions as signals to inform stakeholders about the degree of altruism of the company. We abstract here from any other strategic or market consideration to isolate this effect and focus on its implications.

The effectiveness of material incentives to stimulate pro-social behavior has been questioned in a number of papers dealing with the so-called Motivation Crowding-Out Effect. In a prominent contribution, Bénabou and Tirole (2006) (B & T henceforth) consider a multidimensional signaling model to analyze the agents’ response to material incentives on pro-social activities. They show that, under a certain range of parameters, and assuming that types are normally distributed, rewards on the social activity reduce its supply. The mechanism that operates in our paper is very different to the one considered in B & T. In B & T there is a trade-off between the direct effect of rewards and their indirect effects on the perception of the agent’s intrinsic motivation and his/her reputation. Therefore, crowding out occurs in B & T because the reputation loss associated to rewards is higher than the direct utility gained from receiving these rewards. In contrast, in our setting, incentives affect the supply of CSR activities in the standard way. There is no crowding out of intrinsic

motivation, and the subsidy does not introduce noise in the signal extraction problem. A subsidy attracts companies to CSR activities, as expected. However, firms who are already undertaking CSR activities anticipate a future decline in the reputation attached to social activities and may react by curtailing their current involvement in CSR. The key to the decline in social activities comes from the dynamic nature of the model, not from the structure of multidimensional uncertainty considered in B & T.

In a similar fashion, Seabright (2009) puts forward a screening game where the public authority announces a price  $p$  that will be paid to all agents who choose to participate in a pro-social activity. Social image is very important for truly altruistic individuals because it leads to future profitable relationships with other altruistic agents. Seabright shows that truly altruistic individuals are more likely to choose a civic activity when there is no reward for this activity than when the activity is rewarded. The reason is that the reward spoils the signaling power of the civic activity. In our model, instead, rewards induce a reallocation of companies' involvement in social activities across periods. Some companies will embark on CSR activities attracted by rewards, while others abandon their current engagement in CSR in view of future reputation losses.

We present a two-period signaling model that takes into account explicitly the interplay between material incentives and reputation persistence. The model incorporates policy lags, meaning that it takes some time for announced incentives to be effectively implemented. Firms are heterogeneous in their degree of social concern (their type), and the more socially concerned a firm is, the lower the subjective cost it faces when undertaking a social activity. Each firm has only two actions available in each period: action 1 (the costly social activity), and action 0 (doing nothing). The reputation earned by a company is the posterior expectation on its type after the company has chosen an action. A separating equilibrium occurs when different types of firms choose different actions. We characterize a semi-separating equilibrium as a partition of the set of agents' types: The subset including the most cooperative firms choose the social activity and the rest of the firms choose action 0. In a first step, we characterize the companies' equilibrium choices without the presence of external incentives. After that, we analyze how this equilibrium is modified by incentives in a context of reputation persistence and policy lags.

The major finding of this paper is that a subsidy can lower the aggregate level of CSR. Because of policy lags, a subsidy announced today (first period) becomes effective tomorrow (second period). The choices made by firms in the first period, under the assumption of reputation persistence, affect the range of values for reputation (expected type) available in the second period, and firms are aware of this effect. If a subsidy is established in the second period, there will be more firms willing to carry out the social decision in that period. This will lead to a decrease in reputation of the most altruistic companies. As this reputation loss is discounted back, the benefits of behaving as an altruist in the first period fall because the subsidy lowers second period reputation. This causes some firms (those that were altruist, but not too much) not to embark in CSR activities in the first period. Reputation persistence plays a crucial role here because it is the link that connects the decisions intertemporally. In the new equilibrium, the level of CSR activities prior to the subsidy

(first period) decreases, whereas the level of CSR activities after the subsidy (second period) increases. If a specific condition holds, the former effect is larger than the latter, leading in this case to a global decrease in the level of companies' CSR.

It is worth mentioning that this result is not driven by a displacement of intrinsic motivation on the part of the companies' managers, caused by material incentives. Although corporate philanthropy may be the reason behind firms' differential costs of undertaking pro-social activities, CSR activities here are used instrumentally to signal the degree of altruism and ethical behavior of the firm. Our approach corresponds to the vision "doing well by doing good" whereby a socially responsible behavior is likely to make the firm more profitable (Orlitzky et al (2003) provide substantial evidence for this claim). Each firm trades-off the cost of sending the signal against the reputation benefits the signal entails and finds an optimal compromise between these two opposing forces. The novelty of the paper is to consider the introduction of material incentives in a dynamic setting where the stickiness of reputation plays a crucial role. A newly established subsidy provides incentives for companies' engagement in CSR, but also reduces the benefits that some formerly cooperative companies enjoyed from reputation persistence, to the extent that the subsidy may lead to a global decrease in the supply of CSR.

The paper is organized as follows. In Sect. 2 we introduce the model. In Sect. 3 we analyze the effects of setting up a subsidy on the social activity for the static case. In Sect. 4 we perform the same exercise for the two-period game. Section 5 concludes. All proofs are in the Appendix.

## 2 Model

We consider a two-period model to analyze firms' behavior in a simple dynamic setting. We first introduce a set of definitions and concepts valid for any period. In order to keep notation simple, we do not use time indicators at this stage. Later on, we will refer to a specific period using superscript  $t = 1, 2$ . Let  $\Theta \equiv [0, 1]$  be the set of all firms. Each firm  $i \in \Theta$  is indexed by an element of the unit interval, and we write as  $\theta_i \in [0, 1]$  the type of firm  $i$ . We interpret the parameter  $\theta_i$  as the true degree of firm  $i$ 's social concern or altruism. Types are private information, constant across periods, and distributed in interval  $[0, 1]$  according to a continuous cumulative distribution function  $F(\cdot)$  with associated density function  $f(\cdot)$ .

Let us consider the set of actions  $\{a_0, a_1\}$ . The costly action  $a_1$  represents a social activity undertaken by the firm. Imagine, for instance, a pro-environmental policy or an action that improves labor conditions. Choosing the costless action  $a_0$  means that the firm does not carry out any social activity. We are modeling, in a stylized way, dichotomous decisions like: The firm develops a "community involvement program" or not; the firm launches a campaign to promote fair trade or not; the firm takes an action to reduce its carbon footprint or not, etc. Firms may use social activities as instruments to signal their types. We consider that companies seek to increase their profits by undertaking social activities that improve their reputation and social image. Firms are aware that a good reputation exerts a positive influence on profits,

and also know that reputation can be built through their engagement in activities related to CSR.

Consumers are concerned about the firms' true types because they prefer to engage in a relationship with firms that are somehow altruistic or cooperative. One reason for this may be the existence of significant incomplete contracts in the individual interactions of consumers (stakeholders in general) with companies. Contracts between firms and consumers include prices and some verifiable conditions on quality and additional services, but these contracts cannot include all relevant aspects of the good or service bought, like, for instance, the number of phone calls needed for the customer service to respond or how efficient and decent the company will be when facing unforeseeable problems with the product. These are aspects of stakeholders' interactions with the firm that cannot be easily observable, verifiable and included in the contract but are still important in estimating the quality and value of the good or service traded. Such aspects might well be positively correlated with the kind of "moral capital" that makes some firms devote resources to observable public goods and charities. By sacrificing part of the profits through CSR investments, firms also prove that they are less prone to fail in fulfilling their non-contracted obligations with their stakeholders.

The firm acquires reputation when engaging in social activities. The reputation of a firm with type  $\theta_i$  depends on the information about the firm's actions in period  $t = 1, 2$ . Let  $a^t$  be the vector of actions the firm has undertaken until period  $t = 1, 2$ . The reputation function of a firm in period  $t$  depends on the consumers' system of beliefs given  $a^t$ , and its expression is given by:

$$R^t(a^t) = E[\theta_i | a^t] = \int_0^1 \theta_i f(\theta_i | a^t) d\theta_i, \quad (1)$$

where  $f(\theta_i | a^t)$  is the conditional density function of the types that represents consumers' beliefs on firms' types after firm's actions until period  $t$  have been publicly observed. As in Bénabou and Tirole (2006) and Seabright (2009), we assume that a company's reputation affects positively its payoff. A meta-analysis of 52 studies performed by Orlitzky et al (2003) identifies a positive correlation between the companies' social responsibility and their financial performance.

On the other hand, carrying out a pro-social activity is costly for the firm. A necessary condition for the signals to reveal information about the types is that their cost is different across types. In particular, we assume that the higher the firm's true social concern, the lower the (subjective) cost it faces for any given pro-social activity. Formally, we denote by  $c(a_j, \theta_i)$  the cost of action  $a_j$  ( $j = 0, 1$ ) for firm with type  $\theta_i$  in each period  $t$ . We assume that cost function  $c(\cdot)$  is such that: (i)  $c(a_1, \theta_i) \geq 0$  is continuous, differentiable and strictly decreasing in  $\theta_i$ , with  $c(a_1, 0) = c > 0$ , and  $0 \leq c(a_1, 1) < c(a_1, \theta_i)$  for all  $\theta_i < 1$ ; and (ii)  $c(a_0, \theta_i) = 0$  for all  $\theta_i \in \Theta$ . Parameter  $c$  must be high enough for signals to have an informative value. We denote by  $c'(a_1, \theta_i) < 0$  the derivative of function  $c(\cdot)$  with respect to  $\theta_i$ .

To keep the model tractable, the firm’s payoff in period  $t$  is modeled as a “black box” with a tradeoff between the benefit of reputation and its cost, with the following linear form:

$$\pi^t(a^t, \theta_i) = \gamma R^t(a^t) - c(a_j, \theta_i). \tag{2}$$

$\gamma R^t(a^t)$  represents the income associated with the reputation acquired until period  $t$ , and parameter  $\gamma > 0$  measures the relative importance of reputation with respect to the cost of carrying out a social activity. Parameter  $\gamma$  can be related to the visibility of the social activities carried out by the firm, or it may represent the impact on firm’s profits of socially concerned (or “ethical”) consumers, the ones who enjoy higher utility buying products from “ethical” corporations. The payoff function in Eq. (2) represents the expected net return of the company’s investment in CSR. Observe that, while the value for reputation in period  $t$  depends upon the pattern of actions across periods until period  $t$ , the cost function is evaluated only on the action that is undertaken in period  $t$ .

### 3 The effects of incentives in the one-shot game

In the static case, information about the firm is simply given by the action chosen, either  $a_0$  or  $a_1$ . Then,  $a^t$  in Eqs. (1) and (2) is equal to  $a_j$  with  $j = 0, 1$ , and we write the reputation as  $R(a_j)$  and the payoff function as  $\pi(a_j, \theta_i)$ . A strategy for firm  $i \in \Theta$  is a mapping from the space of types to the set of actions:

$$\sigma_i : [0, 1] \rightarrow \{a_0, a_1\}.$$

After firm  $i$  learns its own type, and knowing the prior distribution of all firms’ types, firm  $i$  either selects  $a_1$  (the pro-social action), or  $a_0$ . The choices of all firms are publicly observed, and then consumers form expectations on firms’ types. A system of beliefs in this scenario is a rule that associates each action with a subset in interval  $[0, 1]$ . For instance, if a firm  $i$  carries out action  $a_j$ , consumers believe that type  $\theta_i$  belongs to set  $\Theta_j$ . Therefore, after observing  $a_j$ , beliefs are updated and represented by the conditional density  $f(\theta_i|a_j)$ , where  $f(\theta_i|a_j) = 0$  for all  $\theta_i \notin \Theta_j$ .

A profile of strategies and a system of beliefs form a Perfect Bayesian Equilibrium (PBE) when the strategy of each firm is optimal given the system of beliefs, and beliefs can be updated using Bayes’ rule from the equilibrium play of all firms. Note that, since the actions set is discrete and the space of types is continuous, any action choice only provides partial information about the type.

We focus our attention on PBE where some information is revealed. Specifically, as the cost of action  $a_1$  is strictly decreasing in the firm’s type, we are able to characterize a semi-separating equilibrium in which higher types select the social activity while lower types do not engage in any social activity. Because there are only two actions available, a semi-separating equilibrium necessarily implies a partition on the type space consisting of two subsets. Our discrete action model allows to understand the process whereby information accrues along both periods and is correctly interpreted by the agents, who can gain a more precise assessment on firms’ types

as time goes by. If we considered instead a continuum of actions, as in Seabright (2009), total separation could be obtained in the first period. Then, prosocial behavior in the second period would not add any information that can be relevant to the agents.

The next lemma establishes that this partition comprises two sub-intervals of interval  $[0, 1]$ , separated by threshold  $z \in [0, 1]$ , where  $z$  can be interpreted as the type of the firm indifferent between choosing  $a_1$  or  $a_0$ . (In case where  $\theta_i = z$ , we assume that firm  $i$  chooses  $a_1$ ).

**Lemma 1** *For any given system of beliefs  $f(\theta_i|a_j)$ , if a semi-separating equilibrium of the one-shot signaling game exists, it must induce a partition on the type space consisting of sub-intervals  $[0, z)$  and  $[z, 1]$ , where all types  $\theta_i \geq z$  choose  $a_1$ , all types  $\theta_i < z$  choose  $a_0$ , and threshold  $z$  is uniquely determined by the equality  $\pi(a_1, z) = \pi(a_0, z)$ .*

It is worth mentioning that Lemma 1 holds regardless the system of beliefs considered. Besides, a remarkable consequence of Lemma 1 is that  $R(a_1) > R(a_0)$ . Next, we define a notion of equilibrium for this game.

**Definition 1** A  $z$ -equilibrium is a semi-separating PBE consisting of strategies  $\sigma_i^*$  and beliefs  $f(\theta_i|a_j)$  such that, for all  $i \in \Theta$ ,

$$\sigma_i^*(\theta_i) = \begin{cases} a_1 & \iff \theta_i \geq z \\ a_0 & \iff \theta_i < z, \end{cases}$$

and

$$\begin{cases} f(\theta_i|a_1) > 0 & \iff \theta_i \in [z, 1] \\ f(\theta_i|a_0) > 0 & \iff \theta_i \in [0, z) \end{cases}$$

where  $z$  is such that  $\pi(a_1, z) = \pi(a_0, z)$ .

An interesting implication of Lemma 1 is that, if a  $z$ -equilibrium exists, it must be unique. The reason is that function  $\pi(a_j, z)$  is monotone in  $z$ , provided that  $R(a_1) > R(a_0)$  and  $c'(a_j, \theta_i) < 0$ . To put it simply, if  $\pi(a_1, z)$  and  $\pi(a_0, z)$  cross, they cross once. Under our assumption that  $c$  is high enough, there could be also a pooling equilibrium where all firms selected  $a_0$  and  $f(\theta_i|a_0) = f(\theta_i)$ . However, this pooling equilibrium would not survive Cho and Kreps' Intuitive Criterion (Cho and Kreps 1987), that restricts out of equilibrium beliefs to be somehow "reasonable". If any firm deviates from  $a_0$ , it must be a high type firm for sure (as long as  $c$  is high enough,  $a_1$  is equilibrium dominated for low types). Note that a firm with type  $\theta_i = 1$  always bears the lowest possible cost from the social action, so this argument holds even for high values of  $c$ .



For the analysis that follows, it is convenient to write the equilibrium reputation of a firm as a function of the interval where the firm’s type lies. That is, for any  $a, b \in [0, 1]$ , with  $a \leq b$ ,  $r(a, b) = \frac{\int_a^b \theta_i f(\theta_i) d\theta_i}{\int_a^b f(\theta_i) d\theta_i}$ . Thus,  $r(a, b)$  can be interpreted as the average type (reputation) of firms whose types belong to interval  $[a, b]$ . Now, based on the notion of a  $z$ -equilibrium in Definition 1, we write the equilibrium reputation of a firm with type  $\theta_i$  as:

$$R(\sigma_i^*(\theta_i)) = \begin{cases} r(z, h) & \iff \theta_i \geq z \\ r(l, z) & \iff \theta_i < z \end{cases}$$

For any threshold  $z \in [l, h] \subseteq [0, 1]$ , we define function  $\Delta r(z;l, h) = r(z, h) - r(l, z)$ . This function represents the difference in reputation between companies belonging to interval  $[l, z]$  (the low interval) with respect to firms whose types are in interval  $[z, h]$  (the high interval). Let  $\Delta r'(z;l, h)$  stand for the derivative of function  $\Delta r$  with respect to  $z$ , that is,  $\Delta r'(z;l, h) = r_z(z, h) - r_z(l, z)$ . The derivative  $\Delta r'(z;l, h)$  measures how the difference between “high” and “low” reputation changes in response to a marginal change in  $z$ . In the next lemma, we derive specific properties of functions  $r(\cdot)$  and  $\Delta r(\cdot)$  that we shall use later.

**Lemma 2** (i) For any  $a, b \in \Theta$ , function  $r(a, b)$  is continuous and strictly increasing in both arguments. Therefore, provided that  $h \geq z \geq l$ , we have  $\Delta r(z;l, h) > 0$ ; (ii) If cdf  $F(\cdot)$  is C2 and concave, we have that  $\Delta r'(z;l, h) \geq 0$ .

The approach followed here to describe reputation is formally similar to the one developed in Candel-Sánchez and Perote-Peña (2020) to analyze the effectiveness of subsidies when agents can use multiple (substitute) prosocial activities to build reputation. Let us denote by  $x$  the amount of a subsidy that is set up by the government with the intention to encourage CSR activities. Each firm must select either  $a_0$  or  $a_1$ , and action  $a_1$  is rewarded with a monetary subsidy  $x \geq 0$ . In the next proposition we characterize a  $z$ -equilibrium in the presence of subsidy  $x$ .

**Proposition 1** If  $c$  is high enough, a  $z$ -equilibrium always exists where threshold  $z(x)$  is implicitly defined by

$$\pi(a_0, z) = \pi(a_1, z) + x,$$

and  $f(\theta_i/a_1) > 0 \iff \theta_i \geq z$ .

The aggregate level of CSR-related activities in this case is given by  $1 - F[z(x)]$ , where  $z(x)$  is implicitly defined by equation  $\Delta r(z;0, 1) + x = c(a_1, z)$ . We are now prepared to analyze the effect of  $x$  on aggregate CSR. Note that, to rule out trivial situations where  $x$  is so high that it always induces all types of firms to choose the social action, we consider that  $x$  must be small enough relative to  $c$ . The following

corollary to Proposition 1 establishes that a subsidy is effective at incentivizing CSR.

**Corollary 1** *A higher subsidy  $x$  induces higher amount of CSR activities, that is*

$$\frac{d[1 - F[z(x)]]}{dx} > 0.$$

The proof of this corollary is immediate, since

$$z'(x) = -\frac{1}{\Delta r'(z; 0, 1) - c'(a_1, z)} < 0.$$

The conclusion is that, in a static setting, the supply of social activities can be effectively incentivized by establishing a subsidy on them. As we will see below, this rather obvious result no longer holds in a framework where reputation is persistent across periods.

#### 4 The effects of incentives in the two-period game with reputation persistence

Along this section, we assume a policy lag, so that a subsidy that is announced at the beginning of the game becomes effective in period 2. Time lags in the implementation of policies are pervasive in real life. For instance, the Dutch government published 2020 tax plan on September 17, 2019. This plan included tax exemptions or reliefs on activities related to companies' CSR. Likewise, the Disabled Access Credit provided a non-refundable credit of up to \$5,000 for small businesses that incur expenditures for the purpose of providing access to persons with disabilities. This law was introduced July 25, 2019, and was effective after December 31, 2019.<sup>1</sup> The lags associated to environmental policy are particularly long (see Di Maria et al., 2012). In dynamic contexts, implementation lags cannot be overlooked because the strategic behavior of firms depends critically on the sequence of events in the game.

Next, we adapt our basic model in Sect. 2 to the two-periods case to account for reputation persistence, and analyze the firms' equilibrium strategies in the light of this fact. In this dynamic case, action  $a_j(j = 0, 1)$ , undertaken in period  $t = 1, 2$ , is denoted by  $a_t^j$ . The action profile until period  $t$ ,  $a^t$ , is given by the first period action if  $t = 1$  and by the vector of actions undertaken in both periods if  $t = 2$ . The payoff earned in period  $t = 1$  by a firm of type  $\theta_i$  that carries out action  $a_j^1$  is denoted by  $\pi^1(a_j^1, \theta_i)$ , and given by Eq. (2) where  $a^1 = a_j^1$ . The second period payoff for a firm

<sup>1</sup> See the Disabled Credit Act at: <https://www.congress.gov/bill/116th-congress/house-bill/4045/text?r=7&s=1>.

with type  $\theta_i$  that chooses  $a_j^2$ , given the action chosen in the first period  $a_j^1$ , is denoted by  $\pi^2(a_j^2, a_j^1, \theta_i)$ . The total payoff for the firm is  $\pi^1(\cdot) + \pi^2(\cdot)$ . For simplicity, we assume that there is no discount factor for the payoffs realized in the second period.

As in the static case,  $x$  represents the amount of the subsidy that is announced by the government, and, as before, we assume that the subsidy  $x$  is not too large relative to cost  $c$ . The subsidy is made effective in period 2, so if a firm of type  $\theta_i$  chooses  $a_1^2$ , it receives a monetary transfer of  $x$ , and its payoff in the second period becomes  $\pi^2(a_1^2, a_j^1, \theta_i) + x$ . If the firm chooses  $a_0^2$  instead, its second period payoff remains as  $\pi^2(a_1^2, a_j^1, \theta_i)$ . The cost function is  $c(a_j^t, \theta_i)$  in each period  $t$ , and has the same functional form as in the static case. For notational simplicity we shall omit superscript  $t$  where obvious.

The timing of the two-period game is as follows:

1. In period 1, the government announces a subsidy of amount  $x$  on CSR activities, to be implemented in period 2. Then, each firm chooses either  $a_1^1$  or  $a_0^1$ .
2. In period 2, having observed the first period choices of all firms, each firm chooses either  $a_1^2$  or  $a_0^2$ .
3. The payoffs are realized.

Reputation persistence in our context means that the types are constant, and that the agents do not forget what they have learned about these types in period one. Therefore, for the inference on firms' types in period two, agents use both the action taken in the second period and the first period action choice (which led to the inference on the type made in the first period). The term "persistence" then refers to the fact that the information inferred in the first period constrains what can be inferred in the second.

The scenario of reputation persistence affects critically the way how agents form the second period beliefs. Note that, because there are only two actions available in each period, and the firm's payoff is monotonic in its type, the result established in Lemma 1 holds for any period  $t = 1, 2$ . Therefore, there exists a unique partition of interval  $[0, 1]$  that separates the firms' types in each period. For this dynamic case, we call  $\alpha$  the separating threshold in the first period. The reputation earned by firms in the second period depends on all information available at the time the belief is formed. This includes the first period actions and also the threshold  $\alpha$ , that separates firms in two subsets. When firms in the first period make their action choices they take into account how these choices will influence beliefs about their types both in the present and in the future. In particular, the value of threshold  $\alpha$  determines the interval in which the type of each firm (reputation) is expected to lie in the second period. Then, although second period decisions affect the inference made on the firms' types, these are constrained to belong either to interval  $[\alpha, 1]$ , if the firm chose  $a_1^1$ , or to interval  $[0, \alpha)$  if the firm chose  $a_0^1$  instead. In other words, firms earn some baseline reputation in the first period, and, in the second period, they have the possibility to modify, to a certain extent, their baseline reputation. For instance, a firm that carried out CSR

activities in the first period can further improve its reputation by choosing again the social action in the second period, or it can avoid the cost of CSR activities in that period and still enjoy a reputation above  $\alpha$ .

A system of beliefs in this dynamic setting with reputation persistence is given by posterior density  $f(\theta_i|a_j^1)$  in the first period, and by posterior density  $f(\theta_i|a_j^1, a_j^2)$  in the second period. The value for the reputation acquired in the second period, after a pair of actions  $(a_j^1, a_j^2)$  has been chosen and threshold  $\alpha$  has been determined, is given by:

$$R^2(a_j^1, a_j^2) = E[\theta_i|a_j^1, a_j^2, \alpha] = \int_0^1 \theta_i f(\theta_i|a_j^1, a_j^2) d\theta_i.$$

Formally, the constraint on the second period reputation introduced by reputation persistence is expressed as:  $R^2(a_1^1, a_2^2) \geq \alpha$  and  $R^2(a_0^1, a_2^2) < \alpha$  for all  $a_2^2$ .

When choosing the first period actions, firms only know their own type, the distribution of the types of all firms and the value of subsidy  $x$ . With this information, each firm can deduce the location of the threshold that separates “good” firms from the “bad” ones in the first period. Moreover, since the types are constant and its distribution does not change across periods, firms are also able to infer the equilibrium actions of all firms in the second period. Therefore, for any given subsidy  $x$ , the only relevant information that affects each firm’s strategy is its own type.

A strategy for firm  $i \in \Theta$  is a pair  $\sigma_i = (\sigma_i^1, \sigma_i^2)$  that associates each type  $\theta_i \in [0, 1]$  with a pattern of actions, where:

$$\sigma_i^1 : [0, 1] \rightarrow \{a_0^1, a_1^1\},$$

and

$$\sigma_i^2 : [0, 1] \times \{a_0^1, a_1^1\} \rightarrow \{a_0^2, a_1^2\}.$$

There are four possible combinations of actions, namely:  $(a_0^1, a_0^2)$ ,  $(a_1^1, a_0^2)$ ,  $(a_0^1, a_1^2)$ ,  $(a_1^1, a_1^2)$ . Note that, whatever action was taken by the firm in the first period, by Lemma 1 there exists a unique partition of the type space also in period two, where a certain threshold, say  $\lambda$ , separates the firms that carry out  $a_1^2$  from those that carry out  $a_0^2$ . The next lemma establishes that the only possible relationship between thresholds  $\alpha$  and  $\lambda$  at equilibrium is  $\lambda \geq \alpha$ .

**Lemma 3** *In any semi-separating equilibrium of the two-period signaling game, under the assumption of reputation persistence, it holds that  $\lambda \geq \alpha$ .*

Notice that, for Lemma 3 to hold, it is necessary that subsidy  $x$  is small enough so as to guarantee that the cost of the social action in the second period minus  $x$  is positive. Given the characteristics of the cost function, this holds trivially true in a positive neighborhood of zero. This assumption is consistent with the empirical observation that performance varies non-monotonically with incentives (Gneezy

and Rustichini 2000). In general, the maximum possible size for the subsidy  $x$  that guarantees  $\lambda \geq \alpha$  depends on the particular distribution function assumed for the firms' types, the cost function and other parameters in the model, such as  $c$  and  $\gamma$ . The relationship between all these elements becomes apparent in the specific example developed below, after Proposition 3.

The phenomenon of reputation persistence has an interesting implication on the equilibrium strategies of firms: information accumulates and is made more precise on time, but losing the opportunity for an early disclosure cannot be compensated in the future with higher investments, since a rational agent with a high type will always invest in it as soon as she has the opportunity. Therefore, only firms that chose  $a_1^1$  can be interested in improving their reputation by doing the costly action  $a_1^2$  in the second period. In other words, if a firm chose  $a_0^1$ , it will never be interested in getting "redeemed" by carrying out CSR activities in the second period. Observe that the "productivity" of early investment in reputation is higher, since it anchors the range of values possible for future information release through new investments. As a consequence, in any semi separating equilibrium, some types above  $\alpha$  will choose the social action in the second period while some other types above  $\alpha$  will prefer to choose  $a_0^2$  instead. Besides, all types below  $\alpha$  will choose  $a_0^2$ . We focus our attention now on a semi-separating  $(\alpha, \lambda)$ -equilibrium, that we define as follows:

**Definition 3** A  $(\alpha, \lambda)$ -equilibrium for the dynamic case is a set of strategies  $\sigma_i^* = (\sigma_i^{*1}, \sigma_i^{*2})$  and beliefs  $[f(\theta_i|a_j^1), f(\theta_i|a_j^1, a_j^2)]$  such that, for all  $i \in \Theta, j = 0, 1$ ,

$$\sigma_i^{*1}(\theta_i) = \arg \max_{\{a_j^1\}} \pi^1(a_j^1, \theta_i) + \pi^2(\sigma_i^{*2}(\theta_i, a_j^1), a_j^1, \theta_i) = \begin{cases} a_1^1 & \Leftrightarrow \theta_i \geq \alpha \\ a_0^1 & \Leftrightarrow \theta_i < \alpha \end{cases}$$

$$\sigma_i^{*2}(\theta_i, a_j^1) = \arg \max_{\{a_j^2\}} \pi^2(a_j^2, a_j^1, \theta_i) = \begin{cases} a_1^2 & \Leftrightarrow \{a_j^1 = a_1^1 \& \theta_i \geq \lambda\} \\ a_0^2 & \Leftrightarrow \text{otherwise} \end{cases}$$

and

$$\left\{ \begin{array}{l} f(\theta_i|a_1^1) > 0 \iff \theta_i \in [\alpha, 1] \\ f(\theta_i|a_0^1, a_0^2) > 0 \iff \theta_i \in [0, \alpha] \\ f(\theta_i|a_1^1, a_0^2) > 0 \iff \theta_i \in [\alpha, \lambda] \\ f(\theta_i|a_1^1, a_1^2) > 0 \iff \theta_i \in [\lambda, 1] \end{array} \right.$$

where  $\alpha$  and  $\lambda$  are the thresholds that "separate" firms that choose the social action from those that do not, in periods 1 and 2, respectively.

This definition of equilibrium basically extends to the two-period case the concept of a  $z$ -equilibrium already defined for the static case. However, the definition includes explicit dynamic features and sophisticated (forward-looking) behavior on the part of firms. In particular, the inference on the firm's type made in period 2 depends on both the action that is currently being taken in that period, but also on the action chosen in the first period. Besides, at the time of choosing the first period action, the firm anticipates and takes into account her second period optimal behavior.

We are now ready to characterize the firms' equilibrium choices for the dynamic case. Our next proposition establishes conditions that define implicitly the equilibrium thresholds  $\alpha$  and  $\lambda$ .

**Proposition 2** *If the distribution function of the types,  $F(\cdot)$ , is concave, a unique  $(\alpha, \lambda)$ -equilibrium exists, where thresholds  $\alpha$  and  $\lambda$  are implicitly defined by the following conditions:*

$$\gamma \Delta r(\lambda; \alpha, 1) + x = c(a_1, \lambda) \tag{3}$$

$$\gamma[\Delta r(\alpha; 0, 1) + \Delta r(\alpha; 0, \lambda)] = c(a_1, \alpha) \tag{4}$$

Condition (3) is the incentive condition (with equality) that defines threshold  $\lambda$ . That is, a firm with type strictly higher than (equal to)  $\lambda$  is better off (indifferent about) choosing the prosocial action (or not) in the second period. Similarly, Condition (4) is the incentive condition that determines the value of first period threshold  $\alpha$ . Notice that Condition (3) also depends on  $\alpha$  and Condition (4) also depends on  $\lambda$ . Therefore, both conditions jointly determine equilibrium thresholds  $\alpha$  and  $\lambda$ . The equilibrium thresholds in Proposition 2 fulfill the inequality  $0 \leq \alpha \leq \lambda \leq 1$ . This means that firms with types  $\theta_i \geq \lambda$  will choose the strategy  $(a_1^1, a_1^2)$ , firms with types  $\theta_i \in [\alpha, \lambda)$  will select  $(a_1^1, a_0^2)$  and firms with types  $\theta_i < \alpha$  will choose  $(a_0^1, a_0^2)$ . Therefore, the total supply of CSR at equilibrium is given by  $2 - F(\alpha) - F(\lambda)$ .

The values of these thresholds, and the difference between them, depends on several elements, namely: (i) the shape of the distribution function of firms' types; (ii) the parameter  $\gamma$ ; and (iii) the shape of cost function  $c(a_j, \theta_i)$ . For instance, for the case of a uniform distribution of the types and cost function  $c(a_1, \theta_i) = (1 - \theta_i^\beta)c$ , it is easy to show that, if parameter  $\beta$  is low enough, the gap between  $\lambda$  and  $\alpha$  increases with  $c$  and decreases with  $\gamma$ . The difference  $\lambda - \alpha$  is a measure of the proportion of agents that follow the strategy  $(a_1^1, a_0^2)$ . This strategy is optimal for agents that, by carrying out  $a_1$  in the first period, secure a reputation above threshold  $\alpha$  in the second period, but for whom it might not be worth to continue carrying out the social action in the second period, as the cost of this action outweighs the gain in reputation. In general, this will occur for the types that are slightly above  $\alpha$ .

Now we turn to the primary question of the section: Is the subsidy  $x$  effective at incentivizing the supply of CSR activities? To answer this question, we must determine how does the total supply of CSR activities change in response to a variation in subsidy  $x$ .

**Corollary 2** *If the distribution function of the types,  $F(\cdot)$ , is concave, a higher subsidy  $x$  induces higher amount of CSR activities in the second period, but also provokes a decrease in CSR activities in the first period, namely*

$$\alpha'(x) = \frac{\gamma r_\lambda(\alpha, \lambda)}{D(\alpha, \lambda)} > 0 \tag{5}$$

and

$$\lambda'(x) = -\frac{\gamma [\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda)] - c'(a_1, \alpha)}{D(\alpha, \lambda)} < 0. \tag{6}$$

Denominator  $D(\alpha, \lambda)$  is a complex expression of  $\alpha$  and  $\lambda$  whose value is provided in the Appendix. The derivatives in (5) and (6) have been computed using implicit differentiation in Eqs. (3) and (4) together. A detailed explanation of the derivation of (5) and (6) is included in the Appendix.

The negative value of derivative  $\lambda'(x)$  does not come as a surprise. If the subsidy on CSR activities increases, the proportion of firms that carry out such activities in the second period must increase as well. However, the behavior of threshold  $\alpha$  is, at first view, surprising: if the amount of the subsidy  $x$  increases, firms decrease participation in CSR activities in the first period! Therefore, the response of the total supply of CSR to a higher subsidy depends on the relative magnitude of these two derivatives, and also on the shape of the distribution function of the types.

**Proposition 3** *The total amount of CSR activities decreases in response to an increase in subsidy  $x$  if:*

$$[\gamma [\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda)] - c'(a_1, \alpha)]f(\lambda) < \gamma r_\lambda(\alpha, \lambda)f(\alpha). \tag{7}$$

Condition (7) means that the reduction in the amount of CSR activities in the first period, caused by an increase in  $x$ , is higher than the second period increase in such activities. The key of this result can be found in Expressions (5) and (6). According to them, the equilibrium values of  $\alpha$  and  $\lambda$  move in opposite directions. Hence, if  $x$  increases, more firms will engage in CSR in the second period, but less will do it in the first period.

What is the rationale of this result? Note that subsidizing CSR activities in the second period lowers the reputation (average type) of the most cooperative firms (the highest types), because threshold  $\lambda$  decreases. This reputation loss comes from the fact that the set of cooperative firms enlarges. As in any positional game, what is important here is the expectation on the own type *relative* to others' types. Hence, if more firms are included in the category of "altruistic", the reputation of all firms within this category must fall. At the same time, firms that chose  $a_1$  in the first period and in the second period select  $a_0^2$  (that is, those firm that are cooperative, but not "too much") also suffer reputation losses due to the subsidy. Therefore, the positive signal about the firm's type sent by a firm that is engaged in CSR activities loses part of its value after a monetary subsidy on such activities is set up. In other words,

the benefits of being cooperative in the first period are lower after the subsidy. Consider a firm whose type is slightly above  $\alpha$  before the subsidy. This firm will choose  $a_1$  in the first period. After a subsidy is introduced in the second period, part of the benefits from choosing  $a_1$  in the first period disappear. It can then be the case that the firm now prefers to choose  $a_0$ . This means that the indifferent type,  $\alpha$ , has increased.

On the one hand, the subsidy discourages some firms to undertake CSR activities in the first period because they anticipate that the subsidy will cause reputation losses in the second period. On the other hand, the proportion of firms that choose CSR activities in the second period grows with respect to a no-subsidy scenario. The final amount of CSR activities carried out depends on the relative magnitude of these two opposite effects. If the condition in Proposition 3 holds, incentives on CSR turn out to be counterproductive.

We can illustrate Proposition 3 with a simple example. Let  $F(\cdot)$  be the uniform distribution, and consider the following family of (convex) cost functions:  $c(a_1, \theta_i) = ce^{-\beta\theta_i}$ , with  $\beta < 1$ . The uniform distribution of the types implies that  $f(\alpha) = f(\lambda) = 1$ ,  $\Delta r'(\alpha; 0, 1) = \Delta r'(\alpha; 0, \lambda) = \Delta r'(\lambda; \alpha, 1) = 0$ , and  $r_\lambda(\alpha, \lambda) = \frac{1}{2}$ . Under these assumptions, Condition (7) in Proposition 3 remains as  $\beta ce^{-\beta\alpha} < \frac{\gamma}{2}$ , and the equations in Proposition 2 that characterize thresholds  $\lambda$  and  $\alpha$  are as follows:

$$\gamma \frac{1 - \alpha}{2} + x = ce^{-\beta\lambda}, \tag{3}$$

$$\gamma \frac{1 + \lambda}{2} = ce^{-\beta\alpha}. \tag{4}$$

Now we consider the following set of parameters:  $\beta = 0.4$ ,  $\gamma = 0.937$ ,  $c = 1$  and  $x = 0.5$ . We solve (numerically) for  $\alpha$  and  $\lambda$  the equations system given by Eqs. (3) and (4), and obtain:  $\alpha \cong 0.5$ , and  $\lambda \cong 0.75$ . Condition (7) holds for  $\alpha = 0.5$  and the parameters considered in this example. Namely,  $0.4e^{-0.2} = 0.327 < \frac{0.937}{2} = 0.468$ . Therefore, in this example, a subsidy on activity  $a_1$  reduces the total amount of CSR.

Our model implicitly considers that both periods are of the same length. However, one can think of a context where the effects of the subsidy remain in the future (possibly along several periods), so the long-term payoff following the enacting of the subsidy might be greater than the payoff earned between the announcement of the subsidy and the enacting of it. A simple way to model this would be to put a higher weight on the second period reputation. In this case, as the reputation loss caused by the subsidy lasts more periods, the proportion of agents who do not wish to behave prosocially in period one increases. Although from the point of view of agents' reputation, the longest the length of the second period, the stronger the incentive to abandon prosocial behavior in the first period, it is also true that if the subsidy is active from period 2 on, and anything else changes from period 2 on, then the global amount of prosocial behavior must increase proportionally to the length of the second period (and also the budgetary cost of the subsidy).



In our model, the assumption of reputation persistence is critical. If we alternatively considered that agents do not recall the actions taken in period one, the subsidy would be effective. In this case, agents do not update beliefs and we would be dealing with the static framework analyzed in Sect. 3, played twice. However, we think it is natural to assume that (at least) part of the information accrued by the agents in the first period can be used in the second to form expectations about the types. In this context, even if there is only partial recall, our qualitative results with respect to the ineffectiveness of the subsidy still hold.

The shape of the distribution function of the types influences the results in the paper. Without concavity, we cannot guarantee that the same results will hold. The assumption of concavity of function  $F(\cdot)$  is not, however, a necessary condition. It is just sufficient to obtain a clear sign of the effects of  $x$  on thresholds  $\alpha$  and  $\lambda$  and it can also be used to prove uniqueness of the equilibrium. In particular, assuming that the distribution function of the types is concave ensures that the sign of  $\Delta r'(z;l, h)$  is positive. Therefore, the difference in reputation between types above a threshold and those below it increases as the threshold increases, which is a sufficient condition for Eqs. (5) and (6) to hold. Other shapes for the distribution function  $F(\cdot)$  could also be compatible with our results.

Although Proposition 3 presents a mainly negative result, our model also admits situations where incentivizing CSR through subsidies can be effective. For instance, if the effect of reputation persistence is sufficiently small (because consumers are short sighted or because visibility of firms' CSR is low), the appropriate modeling would be the one-shot game in the paper, and incentives would have the desired effect. It can also be the case that reputation persistence is important, but the variables in the economy are such that Condition (7) does not hold.

In general, the model presented here applies to situations where (rational) agents use social activities as a signal of altruism and reputation exhibits some degree of persistence. In this context, if a subsidy aimed to promote prosocial behavior is announced and implemented with a policy lag, the resulting total amount of prosocial activities may fall. The reputation formation process is linked to the subsidy policy. The announcement of a subsidy gives agents the opportunity to update expectations on firms' types. The fact that reputation exhibits inertia is critical to obtain the (possible) ineffectiveness of the subsidy. By carrying out a prosocial activity today, agents can secure that future reputation will be above a certain level. This fact introduces strategic dynamic considerations in the model: A subsidy incentivizes prosocial behavior in the second period (as expected), but it may also trigger a strategic (negative) response from the agents in the first period. Provided that the average reputation of agents who carry out the social activity in the second period falls after the subsidy is implemented, being prosocial tomorrow entails lower benefits for the cooperative agents. This effect is discounted back by the agents who are cooperative, but not "too much" (the ones whose type is slightly above the value of threshold  $\alpha$  in the scenario before the subsidy scheme). After the subsidy, some of these agents will be better off not behaving prosocially in period one. In general, the total effect

of the subsidy is undetermined: although it increases the proportion of agents carrying out prosocial activities in period two ( $\lambda$  decreases), it also lowers this proportion in period one ( $\alpha$  increases). Our model should be seen as a counterexample to the general assertion that establishing a subsidy is an effective means to incentivize social activities. Considering the effectiveness of material incentives in a dynamic context with reputation persistence adds a novel element to the theoretical explanations in Bénabou and Tirole (2006) and Seabright (2009), among others.

## 5 Conclusions

The main message of this paper is that setting up material incentives for social activities can be counterproductive in a dynamic environment where reputation exhibits persistence over time. We develop our argument using a signaling model in which prosocial (CSR) activities are instruments that increase the agents' payoffs via reputation. A policy based on future subsidies for social activities may, paradoxically, induce some altruistic agents to abandon these activities. The subsidy attracts more agents towards prosocial activities, and hence lowers the average reputation of the most altruistic agents. This reputation loss provokes that some agents reverse their decisions regarding prosocial behavior, to the extent that the aggregate supply of social activities may decrease.

The use of economic incentives in this paper is embedded in a dynamic framework, where the agents' strategic behavior depends crucially upon two features of the model: reputation persistence and policy lags. We show that the phenomenon of reputation persistence is important for the agents as it constitutes a link between the decisions taken in the present and the payoffs earned in the future. On the other hand, the time lag between announcement and implementation of the subsidy allows the agents to react before the policy is in force. Both effects together may induce a global reduction in the activity that was intended to incentivize.

If the subsidy was implemented without any delay, or if, in general, the length from announcement to implementation was shortened in such a way that no decision could be made in between, we would face a one-shot game, and our conclusions would be similar to the ones established for the static case, also analyzed in the paper. However, as long as policies are subject to unavoidable regulatory and implementation lags, and considering that the sole announcement of any measure triggers a quick response from the agents, the dynamic approach taken in the paper depicts a realistic framework to assess the effectiveness of subsidies to promote social behavior. Any policy aimed at promoting prosocial activities by means of subsidies should be designed taking into account the dynamic features of the framework in which the policy is to be developed. Otherwise, the economic incentive will interact with the agents' reputation concerns in such a way that their strategic response may lead to a global reduction of the activity that was intended to incentivize.

## Appendix

### Proof of Lemma 1

The system of beliefs  $f(\theta_i|a_i)$  induces reputation  $R(a_1)$  for firms that choose  $a_1$  and reputation  $R(a_0)$  for firms that choose  $a_0$ . First, we prove that if some type  $\tilde{\theta}_i$  chooses  $a_0$ , then, all types below  $\tilde{\theta}_i$  must also choose  $a_0$ . We prove it by contradiction. If a firm with type  $\tilde{\theta}_i$  chooses  $a_0$ , this firm obtains reputation  $R(a_0)$ . Since it bears zero cost from taking the action  $a_0$ , its payoff is  $\gamma R(a_0)$ . Because in a separating equilibrium every firm maximizes payoffs given beliefs, it must be true that  $\gamma R(a_0) \geq \gamma R(a_1) - c(a_1, \tilde{\theta}_i)$ . We show now that any other type  $\hat{\theta}_i \in [0, \tilde{\theta}_i]$  must also choose  $a_0$ . If type  $\hat{\theta}_i$  chose  $a_1$  in equilibrium, the following inequality should hold:  $\gamma R(a_1) - c(a_1, \hat{\theta}_i) \geq \gamma R(a_0)$ . However, this would imply that  $c(a_1, \tilde{\theta}_i) \geq c(a_1, \hat{\theta}_i)$ , which is not possible because the cost function is decreasing in the type and  $\tilde{\theta}_i \geq \hat{\theta}_i$ . We conclude that, if some type  $\tilde{\theta}_i$  selects  $a_0$ , all other types below  $\tilde{\theta}_i$  must also select  $a_0$ . It is easy to see that this conclusion implies that, if some type, say  $\check{\theta}_i$ , chooses  $a_1$ , then all types above  $\check{\theta}_i$  must also choose  $a_1$ . Suppose that  $\theta_i > \check{\theta}_i$  and  $\theta_i$  chooses  $a_0$ . Then, this would imply that  $\check{\theta}_i$  should also choose  $a_0$ , which is a contradiction to the initial hypothesis. So, when the cost of the social action is strictly decreasing in the type, the only possible partition of the types in a separating equilibrium is given by sub-intervals  $[0, z)$  and  $[z, 1]$ , where the value of  $z$  depends on the system of beliefs  $f(\theta_i|a_1)$  and is implicitly defined by the following equation:

$$\gamma \left[ \int_0^1 \theta_i f(\theta_i|a_1) d\theta_i \right] - c(a_1, z) = \gamma \left[ \int_0^1 \theta_i f(\theta_i|a_0) d\theta_i \right] - c(a_0, z).$$

### Proof of Lemma 2

Statement (i) follows immediately from our assumptions on function  $F(\cdot)$  (see Kotz et al (2004)). A version of Statement (ii) was proven, in a different setup, in Candel-Sánchez and Perote-Peña (2020). We replicate here this proof with the adequate changes in notation and introducing the specific features of the present model. Let us denote by  $r_z(z, h)$  and  $r_z(l, z)$  the derivatives of functions  $r(z, h)$  and  $r(l, z)$  with respect to  $z$ , respectively. We now prove that  $\Delta r'(z; l, h) = r_z(z, h) - r_z(l, z) \geq 0$ . First of all, we compute

$$r_z(z, h) = \frac{f(z)}{F(h) - F(z)} [r(z, h) - z] \geq 0,$$

and

$$r_z(l, z) = \frac{f(z)}{F(z) - F(l)} [z - r(l, z)] \geq 0.$$

We prove that

$$\frac{r(z, h) - z}{F(h) - F(z)} \geq \frac{z - r(l, z)}{F(z) - F(l)}.$$

If  $l \rightarrow h$ , in the limit, we would have  $r_z(l, z) \rightarrow r_z(z, h)$ . For the inequality above to be true, it is sufficient that  $r_z(l, z)$  is increasing in  $l$ , that is,

$$\frac{dr_z(l, z)}{dl} = \frac{f(l)}{(F(z) - F(l))^2} [z + l - 2r(l, z)] \geq 0.$$

Clearly, function  $r_z(l, z)$  is increasing in  $l$  if and only if  $r(l, z) \leq \frac{z+l}{2}$ . Now, we consider the truncation of cdf  $F(\cdot)$  to interval  $[l, z]$ . Let  $f(\theta|l \leq \theta \leq z)$  be the associated density function. On the other hand, let us consider the truncation of the uniform distribution to interval  $[l, z]$ . The associated density function is constant and equal to  $\frac{1}{z-l}$ . Provided that the cdf  $F(\cdot)$  is concave in all its domain (the interval  $[0, 1]$ ), it holds that  $F(\theta) \geq \theta$  for all  $\theta$ . This is also true if we restrict the domain of  $\theta$  to the interval  $[l, z]$ . Therefore, the cdf of the uniform distribution on this interval is first order stochastic dominant over the cdf  $F(\cdot)$ , which implies that

$$\frac{z+l}{2} = \int_l^z \theta \frac{1}{z-l} d\theta \geq \int_l^z \theta f(\theta|l \leq \theta \leq z) d\theta = r_2(l, z).$$

The above expression proves that  $\Delta r'(z;l, h) \geq 0$ .

### Proof of Lemma 3

To conclude that  $\lambda \geq \alpha$ , we must prove that, under the assumption of reputation persistence, it does not exist any PBE where: (i) all types  $\theta_i \geq \alpha$  choose the strategy  $(a_1^1, a_1^2)$ ; or (ii) all types  $\theta_i \geq \alpha$  choose the strategy  $(a_1^1, a_0^2)$ ; or (iii) some types choose the strategy  $(a_1^1, a_0^2)$  and other types choose  $(a_0^1, a_1^2)$ . If there is no PBE such that either (i), or (ii), or (iii) hold, then, necessarily, any PBE of the signaling game must be such that  $\lambda \geq \alpha$ . We prove each one of the statements separately.

Statement (i): *There is no PBE such that all types higher than  $\alpha$  choose the social action in the second stage.* Suppose that there exists a PBE in which all types  $\theta_i \geq \alpha$  choose the strategy  $(a_1^1, a_1^2)$ . This means that, for the beliefs corresponding to this equilibrium, all firms with type  $\theta_i \geq \alpha$  maximize their payoff by choosing  $(a_1^1, a_1^2)$ . That is, all firms that chose the prosocial action in the first period prefer to choose the prosocial action also in the second period. Can this be an equilibrium? Notice that, given all actions taken by all agents in the first period, the value of threshold

$\alpha$  is fixed and publicly known in the second period. If all types above  $\alpha$  select the social action again, there is no updating of beliefs with respect to the first period. But this implies that the reputation of these firms in the second period holds constant with respect to the one already acquired in the first period, i.e.,  $R^1(a_1^1) = R^2(a_1^1, a_1^2)$ . Since there is no reputation gain to firm  $i$  by choosing  $a_1^2$ , and firms with types  $\theta_i \geq \alpha$  must pay a net cost of  $c(a_1^2, \theta_i) - x$ , it is not possible that all firms with types  $\theta_i \geq \alpha$  are maximizing their payoff when choosing action  $a_1^2$ . This contradicts the initial hypothesis.

Statement (ii): *There is no PBE such that all firms with types higher than  $\alpha$  choose  $a_0^2$ .* We prove that a situation where all firms above  $\alpha$  select  $a_0^2$  and beliefs are not updated does not constitute an equilibrium either, provided that out of equilibrium beliefs fulfill the Intuitive Criterion. The reason is that, for firms with low types, choosing  $a_1^1$  is equilibrium dominated, but this is not true for high types. If a company deviates from the (pooling) equilibrium where all types choose  $a_0^2$ , this firm must have a high type for sure, and the prior should be updated accordingly. Then, a situation where all types above  $\alpha$  choose  $a_0^2$  is not an equilibrium when the system of beliefs satisfies Cho and Kreps' Intuitive Criterion.

From statements (i) and (ii), we conclude that there is no possible PBE in which all types above  $\alpha$  are pooled in the second period (that is, either all choose  $a_1^1$  or all choose  $a_0^2$ ). Therefore, in any equilibrium there must be some firms above  $\alpha$  choosing  $a_1^1$  and other firms above  $\alpha$  choosing  $a_0^2$ . Next, we show that, if  $(a_1^1, a_0^2)$  is an equilibrium strategy for some types, then  $(a_0^1, a_1^2)$  cannot be an equilibrium strategy for any type. This result, together with claims (i) and (ii) allows us to conclude that  $\lambda \geq \alpha$ .

Statement (iii): *There is no PBE where some types choose the strategy  $(a_1^1, a_0^2)$  and other types choose  $(a_0^1, a_1^2)$ .* Recall that we are ruling out (extremely high) values for the subsidy that induce all firms, regardless of their type, to carry out the social action. Let us assume that there exists an equilibrium in which some agents choose the strategy  $(a_1^1, a_0^2)$  and other agents choose  $(a_0^1, a_1^2)$ . This hypothesis will lead us to a contradiction. Let us suppose that  $(a_1^1, a_0^2)$  is the equilibrium strategy for some firm, say  $\theta_h$ , and  $(a_0^1, a_1^2)$  is the equilibrium strategy for a firm with type  $\theta_l$ . Since  $(a_1^1, a_0^2)$  is the equilibrium choice for a firm with type  $\theta_h$ , the payoff from  $(a_1^1, a_0^2)$  for this type is higher than the payoff it gets from  $(a_0^1, a_1^2)$ , that is,

$$\gamma [R^1(a_1^1) + R^2(a_1^1, a_0^2)] - c(a_1^1, \theta_h) \geq \gamma [R^1(a_0^1) + R^2(a_0^1, a_1^2)] - c(a_1^2, \theta_h) + x.$$

By the same reason, the payoff from  $(a_0^1, a_1^2)$  to  $\theta_l$  is higher than the payoff this type gets from choosing  $(a_1^1, a_0^2)$ . Therefore:

$$\gamma [R^1(a_0^1) + R^2(a_0^1, a_1^2)] - c(a_1^2, \theta_l) + x \geq \gamma [R^1(a_1^1) + R^2(a_1^1, a_0^2)] - c(a_1^1, \theta_l).$$

Provided that the cost of the social action is the same across periods for any given type, the two inequalities above imply that:

$$\gamma [R^1(a_1^1) + R^2(a_1^1, a_0^2)] = \gamma [R(a_0^1) + R^2(a_0^1, a_1^2)] + x.$$

However, we know that  $R^1(a_1^1) > R^1(a_0^1)$  by Lemma 1, and  $R^2(a_1^1, a_0^2) \geq R^2(a_0^1, a_1^2)$  by reputation persistence. Hence, for small values of  $x$ , we have

$$\gamma [R^1(a_1^1) + R^2(a_1^1, a_0^2)] > \gamma [R(a_0^1) + R^2(a_0^1, a_1^2)] + x,$$

which is a contradiction to the initial hypothesis.

**Proof of Proposition 1**

Let us note that  $r(z, 1) > r(0, z)$  for any  $z \in [0, 1]$ . Given the equilibrium beliefs, we prove that  $\pi(a_0, \theta^i) < \pi(a_1, \theta^i) + x$  for  $\theta^i > z$  and  $\pi(a_0, \theta^i) > \pi(a_1, \theta^i) + x$  for  $\theta^i < z$ . We define the continuous function  $\Delta\pi(\theta_i) = \Delta r(\theta^i, 0, 1) + x - c(a_1, \theta_i)$ . Since  $c(a_1, \theta_i)$  is strictly decreasing in  $\theta_i$ , and  $\Delta r'(\theta^i, 0, 1) \geq 0$ , it is straightforward that  $\Delta\pi(\theta_i)$  is strictly increasing and continuous in  $\theta_i$ . We just need to show that  $\Delta\pi(\theta_i) > 0$  for  $\theta^i > z$  and  $\Delta\pi(\theta_i) < 0$  for  $\theta^i < z$ . Clearly, we have that  $\Delta\pi(1) = \Delta r(1; 0, 1) + x > 0$ . If  $c$  is high enough we have also that  $\Delta\pi(0) = \Delta r(0; 0, 1) + x - c < 0$ . Then, by Bolzano’s theorem, there exists a type, say  $z$ , for which  $\Delta\pi(z) = 0$ .

**Proof of Proposition 2**

We first obtain the characterizing conditions and then prove uniqueness of the equilibrium. In the second period, a firm with type  $\theta^i$  that chose  $a_1^1$ , will choose  $a_1^2$  instead of  $a_0^2$  whenever

$$\gamma r(\lambda, 1) - c(a_1, \theta_i) + x \geq \gamma r(\alpha, \lambda).$$

Then, if  $\theta_i = \lambda$  we have:

$$\gamma \Delta r(\lambda, \alpha, 1) + x = c(a_1, \lambda).$$

In the first period, and being aware of the second period threshold  $\lambda$ , a firm with type  $\theta^i < \lambda$  chooses  $a_1^1$  instead of  $a_0^1$  if

$$\gamma r(\alpha, 1) - c(a_1, \theta_i) + \gamma r(\alpha, \lambda) \geq 2\gamma r(0, \alpha).$$

Therefore,  $\alpha$  is the value of type  $\theta_i$  for which:

$$\gamma [\Delta r(\alpha; 0, 1) + \Delta r(\alpha; 0, \lambda)] = c(a_1, \alpha).$$

We now take the total differential in Eq. (3) with respect to  $\alpha$  and  $\lambda$  and reorder the resulting expression to obtain:

$$\frac{d\lambda}{d\alpha} = \frac{\gamma r_\alpha}{\gamma \Delta r'(\lambda; \alpha, 1) - c'(a_1, \lambda)},$$

where  $\Delta r'(\lambda; \alpha, 1) = r_\lambda(\lambda, 1) - r_\lambda(\alpha, \lambda)$ . Therefore, for  $\frac{d\lambda}{d\alpha} > 0$  it is sufficient (but not necessary) that  $\Delta r'(\lambda; \alpha, 1) \geq 0$ , provided that  $c'(a_1, \lambda) < 0$ . We call  $\lambda(\alpha)$  the implicit function derived from Eq. (3). From Lemma 2-(ii), the concavity of  $F(\cdot)$  implies  $\lambda'(\alpha) > 0$ .

We proceed in a similar way in Eq. (4) and obtain:

$$\frac{d\lambda}{d\alpha} = - \frac{\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda) - c'(a_1, \alpha)}{r_\lambda(\alpha, \lambda)},$$

where  $\Delta r'(\alpha; 0, 1) = r_\alpha(\alpha, 1) - r_\alpha(0, \alpha)$  and  $\Delta r'(\alpha; 0, \lambda) = r_\alpha(\alpha, \lambda) - r_\alpha(0, \alpha)$ . For  $\frac{d\lambda}{d\alpha}$  to be negative, it is sufficient that  $\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda) \geq 0$ , since  $c'(a_1, \alpha) < 0$  and  $r_\lambda(\alpha, \lambda) > 0$ .

The concavity of  $F(\cdot)$  is sufficient for  $\Delta r' \geq 0$ , which in turn is sufficient to have a positive relationship between  $\alpha$  and  $\lambda$  in Eq. (3) and a negative one in Eq. (4). Therefore, the functions implicit in both equations cannot cross more than once.

Let us substitute  $\lambda(\alpha)$  (the implicit function between  $\alpha$  and  $\lambda$  in Eq. (3)), in Eq. (4). Then, Eq. (4) is as follows:

$$\gamma[\Delta r(\alpha; 0, 1) + \Delta r(\alpha; 0, \lambda(\alpha))] - c(a_1, \alpha) = 0.$$

We call  $\varphi(\alpha)$  the function in the left-hand side of the above equation. Function  $\varphi(\alpha)$  is continuous and differentiable, under our assumptions. Furthermore, this function is monotone increasing since

$$\varphi'(\alpha) = \gamma[\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda(\alpha))\lambda'(\alpha)] - c'(a_1, \alpha) > 0.$$

By Bolzano's Theorem, we only need  $\varphi(0) < 0$  and  $\varphi(1) > 0$  to ensure that there exists  $\hat{\alpha}$  such that  $\varphi(\hat{\alpha}) = 0$ . Observe that

$$\varphi(0) = \gamma[r(0, 1) + r(0, \lambda(0))] - c(a_1, 0).$$

For  $c(a_1, 0) = c$  high enough, we have  $\varphi(0) < 0$ . On the other hand,

$$\varphi(1) = 2\gamma[1 - r(0, 1)] - c(a_1, 1) > 0.$$

Therefore, if the cdf of the types,  $F(\cdot)$ , is concave, and the cost  $c$  faced by a firm whose type is  $\theta_i = 0$  is high enough (i.e., the maximum subjective cost of a social action is high enough), then, thresholds  $\alpha$  and  $\lambda$  uniquely characterize a  $(\alpha, \lambda)$ -equilibrium as described in Definition 3.

### Proof of Corollary 2

The equations that characterize thresholds  $\alpha$  and  $\lambda$  can be rewritten as

$$\gamma r(\lambda, 1) - c(a_1, \lambda) + x = \gamma r(\alpha, \lambda) \tag{3}$$

$$\gamma r(\alpha, 1) - c(a_1, \alpha) + \gamma r(\alpha, \lambda) = 2\gamma r(0, \alpha) \tag{4}$$

In order to ascertain the sign of derivatives  $\alpha'(x)$  and  $\lambda'(x)$ , we resort to implicit differentiation of Eqs. (3) and (4):

$$\gamma [\Delta' r(\lambda; \alpha, 1) d\lambda - r_\alpha(\alpha, \lambda) d\alpha] + dx = c'(a_1, \lambda) d\lambda \tag{3'}$$

$$\gamma [\Delta' r(\alpha; 0, 1) d\alpha + \Delta' r(\alpha; 0, \lambda) d\alpha + r_\lambda(\alpha, \lambda) d\lambda] = c'(a_1, \alpha) d\alpha \tag{4'}$$

Now we solve for  $d\lambda$  in Eq. (3'):

$$d\lambda = \frac{\gamma r_\alpha(\alpha, \lambda)}{\gamma \Delta r'(\lambda; \alpha, 1) - c'(a_1, \lambda)} d\alpha - \frac{1}{\gamma \Delta r'(\lambda; \alpha, 1) - c'(a_1, \lambda)} dx.$$

We substitute  $d\lambda$  computed above in Eq. (4'), and rewrite the equation as:

$$\frac{d\alpha}{dx} = \frac{\gamma r_\lambda(\alpha, \lambda)}{[\gamma [\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda)] - c'(a_1, \alpha)] [\gamma \Delta r'(\lambda; \alpha, 1) - c'(a_1, \lambda)] + \gamma^2 r_\lambda(\alpha, \lambda) r_\alpha(\alpha, \lambda)}.$$

This is the expression that appears in Eq. (5) in the paper (with the value for the denominator denoted as  $D(\alpha, \lambda)$ ). Finally, we divide by  $dx$  the expression of  $d\lambda$  just computed, and substitute in it the expression for  $\frac{d\alpha}{dx}$ , yielding

$$\frac{d\lambda}{dx} = - \frac{\gamma [\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda)] - c'(a_1, \alpha)}{[\gamma [\Delta r'(\alpha; 0, 1) + \Delta r'(\alpha; 0, \lambda)] - c'(a_1, \alpha)] [\gamma \Delta r'(\lambda; \alpha, 1) - c'(a_1, \lambda)] + \gamma^2 r_\lambda(\alpha, \lambda) r_\alpha(\alpha, \lambda)}.$$

This is the same expression as in Eq. (6) in the paper, where the denominator is denoted by  $D(\alpha, \lambda)$ .

### Proof of Proposition 3

The marginal change in the supply of CSR caused by a marginal change in the amount of the subsidy is given by

$$\frac{d[2 - F(\alpha(x)) - F(\lambda(x))]}{dx} = -f(\alpha)\alpha'(x) - f(\lambda)\lambda'(x).$$

Using Eqs. (5) and (6), we find that the sign of the expression above is negative if and only if Condition (7) holds.

**Funding** Open Access funding provided thanks to the CRUE-CSIC agreement with Springer Nature. Funding was provided by ministerio de economía y competitividad (Grant No ECO2016-75631-P), Francisco Candel-Sánchez and Juan Perote-Peña (Grant No ECO2016-75631-P).

**Data availability statement** Data sharing not applicable to this article as no datasets were generated or analyzed during the current study.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Ang SH, Wight AM (2009) Building intangible resources: The stickiness of reputation. *Corp Reput Rev* 12(1):21–32
- Arora S, Gangopadhyay S (1995) Toward a theoretical model of voluntary overcompliance. *J Econ Behav Organ* 28(3):289–309
- Bagnoli M, Watts SG (2003) Selling to socially responsible consumers: competition and the private provision of public goods. *J Econ Manag Strategy* 12(3):419–445
- Baron DP (2011) Credence attributes, voluntary organizations, and social pressure. *J Public Econ* 95(11–12):1331–1338
- Bénabou R, Tirole J (2006) Incentives and prosocial behavior. *Am Econ Rev* 96(5):1652–1678
- Bénabou R, Tirole J (2010) Individual and corporate social responsibility. *Economica* 77(305):1–19
- Berg J, Dickhaut J, McCabe K (1995) Trust, reciprocity, and social history. *Games Econom Behav* 10(1):122–142
- Besley T, Ghatak M (2007) Retailing public goods: The economics of corporate social responsibility. *J Public Econ* 91(9):1645–1663
- Candel-Sánchez F, Perote-Peña J (2020) Optimal incentives on multiple prosocial activities when reputation matters. *Scand J Econ* 122(3):1207–1230
- Cho IK, Kreps DM (1987) Signaling games and stable equilibria. *Q J Econ* 102(2):179–221
- Di Maria C, Smulders S, Van der Werf E (2012) Absolute abundance and relative scarcity: environmental policy with implementation lags. *Ecol Econ* 74:104–119
- Gneezy U, Rustichini A (2000) Pay enough or don't pay at all. *Q J Econ* 115(3):791–810
- Gneezy U, Meier S, Rey-Biel P (2011) When and why incentives (don't) work to modify behavior. *J Econ Perspect* 25(4):191–210
- Kitzmueller M, Shimshack J (2012) Economic perspectives on corporate social responsibility. *J Econ Lit* 50(1):51–84
- Kotchen MJ (2006) Green markets and private provision of public goods. *J Polit Econ* 114(4):816–834
- Kotz S, Balakrishnan N, Johnson NL (2004) Continuous multivariate distributions, volume 1: models and applications, vol 1. Wiley, New York
- McWilliams A, Siegel DS (2011) Creating and capturing value: Strategic corporate social responsibility, resource-based theory, and sustainable competitive advantage. *J Manage* 37(5):1480–1495
- Orlitzky M, Schmidt FL, Rynes SL (2003) Corporate social and financial performance: a meta-analysis. *Organ Stud* 24(3):403–441
- Roberts PW, Dowling GR (2002) Corporate reputation and sustained superior financial performance. *Strateg Manag J* 23(12):1077–1093
- Schultz M, Mouritsen J, Gabrielsen G (2001) Sticky reputation: Analyzing a ranking system. *Corp Reput Rev* 4:24–41
- Seabright PB (2009) Continuous preferences and discontinuous choices: How altruists respond to incentives. *BE J Theor Econ*. <https://doi.org/10.2202/1935-1704.1346>
- Sen S, Bhattacharya C (2001) Does doing good always lead to doing better? Consumer reactions to corporate social responsibility. *J Mark Res* 38(2):225–243

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.