

# Epidemics with Behavior

Christoph Carnehl\*      Satoshi Fukuda<sup>†</sup>      Nenad Kos<sup>‡</sup>

April 19, 2022

First Version: February 2021

## Abstract

We study social distancing in an epidemiological model. Distancing reduces the individual's probability of getting infected but comes at a cost. Equilibrium distancing flattens the curve and decreases the final size of the epidemic. We examine the effects of distancing on the outset, the peak, and the final size of the epidemic. First, the prevalence increases beyond the initial value only if the transmission rate is in the intermediate region. Second, the peak of the epidemic is non-monotonic in the transmission rate. A reduction in the transmission rate can increase the peak. However, a decrease in the cost of distancing always flattens the curve. Third, both a reduction in the transmission rate as well as a reduction in the cost of distancing decrease the final size of the epidemic. Our results suggest that public policies that decrease the transmission rate can lead to unintended negative consequences in the short run but not in the long run. Therefore, it is important to distinguish between interventions that affect the transmission rate and interventions that affect contact rates.

Keywords: epidemics; equilibrium distancing; transmission rate; interventions; SIR

JEL Classification Numbers: I12; I18; C73

---

\*Bocconi University, Department of Economics and IGIER. Email: [christoph.wolf@unibocconi.it](mailto:christoph.wolf@unibocconi.it).

<sup>†</sup>Bocconi University, Department of Decision Sciences and IGIER. Email: [satoshi.fukuda@unibocconi.it](mailto:satoshi.fukuda@unibocconi.it).

<sup>‡</sup>Bocconi University, Department of Economics, IGIER and CEPR. Email: [nenad.kos@unibocconi.it](mailto:nenad.kos@unibocconi.it).

# 1 Introduction

When faced with the possibility of contracting a hazardous disease, people undertake protective measures. They reduce social interactions due to the risk of meeting an infected person. Such behavior is not novel. During the plague pandemics, citizens would flee affected areas and wear costumes to protect themselves from the infection; long-beaked masks worn by physicians in the 17th century achieved particular notoriety. Public authorities eventually began to coordinate the response to epidemics. Famously, Venice required that the passengers on ships from affected areas confine themselves for forty days; thus, the term “quarantine” was minted.<sup>1</sup> Such behavior calls for the explicit incorporation of human behavior in epidemiological models. Yet, the standard SIR model of epidemics, introduced by [Ross and Hudson \(1917\)](#) and [Kermack and McKendrick \(1927\)](#), assumes that individuals engage in as many interactions at the height of the epidemic as they do when the disease is barely present.

We study a tractable model of epidemics that incorporates social distancing and show that explicitly modeling human behavior has important consequences on the predicted trajectory of an infectious disease.<sup>2</sup> Susceptible individuals non-cooperatively decide to which extent to reduce interactions at each point in time. Such distancing is costly but reduces the probability of getting infected. The cost of getting infected is fixed; building on the work of [Chen \(2012\)](#). We show that an equilibrium exists and that it is unique. If the disease spreads, the epidemic has a single peak: it propagates through the population until it reaches the peak prevalence, then it recedes and eventually dies out. Susceptible individuals distance throughout the epidemic, though the intensity of their distancing varies with the amount of actively infected individuals. Distancing affects three crucial and commonly discussed features: the conditions for an epidemic to start, its peak, and its final size.

First, we define a basic reproduction number taking distancing into account—the *behavioral basic reproduction number*. It consists of the classical, epidemiological basic reproduction number,  $R_0$ , multiplied by a behavioral term; a similar concept was introduced in [Fenichel et al. \(2011\)](#).<sup>3</sup> We show that the disease propagates itself if and only if the behavioral basic reproduction number is larger than one. The novelty is that the

---

<sup>1</sup>After the Italian word *quaranta* for forty; see [Snowden \(2019\)](#).

<sup>2</sup>Ours is not the first model of behavior during an epidemic. An account of the related literature follows below.

<sup>3</sup>One can derive an analogous *behavioral effective reproduction number*. That the basic reproduction number without distancing may be misleading when trying to understand epidemic dynamics has been recognized before. For example, [Caley et al. \(2008\)](#) find that the observed attack rate of the 1918-1919 influenza pandemic was substantially lower than the one expected based on the basic reproduction number and attribute this discrepancy to social distancing.

behavioral basic reproduction number is concave in the transmission rate and that the disease spreads only for intermediate values of the transmission rate. If the transmission rate is too high, individuals distance with such fervor that the prevalence never rises above the initial seed of infection. This finding stands in stark contrast with the predictions offered by the SIR model without distancing where the infection spreads if the transmission rate is high enough; see for example [Brauer and Castillo-Chavez \(2012\)](#).

Second, we derive results about the peak prevalence of the disease. The peak prevalence is crucial to understand whether a disease might cause the health system to reach its capacity. For example, the 1918 influenza pandemic hit an unprepared health system which soon became overwhelmed; see [Jester et al. \(2018\)](#) and [Schoch-Spana \(2001\)](#). In March 2020—less than a month after the coronavirus erupted in Italy—the healthcare system in Northern Italy was under such severe pressure that some pneumonia patients could not be treated.<sup>4</sup> In order to avoid the active number of infected individuals exceeding the health care system’s capacity, the goal became to *flatten the curve*. We show that an increase in the cost of distancing unequivocally leads to a reduction in distancing and, therefore, to a higher peak prevalence of the disease. However, peak prevalence is non-monotonic in the transmission rate. If the transmission rate is high enough for the disease to spread but not too high, an increase in the transmission rate leads to an increase in the peak prevalence. In contrast, when the transmission rate is sufficiently high, an increase in the transmission rate decreases the peak prevalence and causes flattening of the curve.

The above comparative statics lend themselves to two interpretations:

- i) A change in the transmission rate could be interpreted as a comparison of equilibrium trajectory for two different diseases. In this context, our results imply that a disease with a higher transmission rate can lead to a lower peak prevalence.
- ii) A change in either the transmission rate or the cost of distancing could be interpreted as a public health policy.<sup>5</sup> Our results suggest that a policy that decreases the transmission rate could lead to a higher peak prevalence.<sup>6</sup> In addition, the fact that peak prevalence is monotonic in the cost of distancing and non-monotonic in the transmission rate has important implications on how interventions should be modeled.

---

<sup>4</sup>See <https://www.nytimes.com/2020/03/12/world/europe/12italy-coronavirus-health-care.html>.

<sup>5</sup>Our comparative statics, strictly speaking, correspond to public policies that are implemented at the beginning and held until the end. The analysis can easily accommodate a policy that is enacted later as long as it persists indefinitely. Nevertheless, the underlying driving forces of our results—that a change in the transmission rate has two opposing effects and a change in the cost of distancing has only one effect—hold for any public policy that affects one of those variables.

<sup>6</sup>The idea that a policy which is meant to protect can lead to more risky behavior, known as *risk compensation*, was first documented by [Peltzman \(1975\)](#).

A large body of work that studies non-pharmaceutical interventions models these either as reductions in the transmission rate (see, for example, [Kruse and Strack, 2020](#); [Rachel, 2020a](#)) or as directly imposing restrictions on the activity level of (a fraction of) individuals (see, for example, [Acemoglu et al., 2021](#); [Alvarez et al., 2021](#); [Farboodi et al., 2021](#))—which are equivalent approaches in the SIR dynamics without behavior. Our results suggest that modeling individual distancing choices explicitly requires a careful choice of modeling interventions because their qualitative implications differ through the behavioral channel. On the one hand, those interventions affecting the rate at which the disease propagates conditional on meetings, e.g., mandatory mask mandates, should be modeled as a decrease in the transmission rate.<sup>7</sup> On the other hand, those interventions that directly affect the incentives to distance, e.g., restaurant, bar or museum closures, should be modeled as a decrease in the cost of distancing; e.g., [Fenichel et al. \(2011\)](#) model a public policy intervention as a change in the payoff structure of contacts.

Third, we find that the possibly detrimental short-run effects of a decrease in the transmission rate disappear in the long run. Despite the non-monotonicity of the peak prevalence in the transmission rate, the total number of infected individuals throughout the epidemic is monotonically increasing in both the cost of distancing and the transmission rate. Our model predicts a smaller final size of the epidemic (i.e., less total infections) than the standard SIR model due to distancing. Indeed, the model converges to the SIR model without distancing when the cost of distancing grows and so does the final size of the epidemic.

With these findings, we highlight an important trade-off between short-run mitigation, i.e., flattening the curve to avoid an overburdened health system, and long-run size of epidemics when considering the transmission rate. This trade-off arises due to the varying degree to which behavior matters during an epidemic. At the peak, the infection risks are high and individuals’ distancing decisions have a strong impact on the dynamics of the epidemic. When an epidemic fades out, however, behavior is of less importance as individual risks are low and the standard SIR mechanics dominate the behavioral effects. However, the trade-off disappears once policies are considered that directly affect distancing incentives of individuals and both short-run mitigation and long-run size of the epidemic are obtained with similar policies, i.e., lowering the cost of distancing.

In the final section, we present an environment in which the cost of infection is endo-

---

<sup>7</sup>Note that this result does not necessarily imply that mandating mask-wearing in public spaces will worsen an epidemic; it may flatten the curve as well. However, we want to highlight the *possibility* of this perverse effect arising. While [Yan et al. \(2021\)](#) document evidence that mask mandates lead to risk compensation behavior, they also argue that it is unclear whether such risk compensation behavior would lead to a net increase or decrease in transmission. [Chernozhukov et al. \(2021\)](#) show that mask mandates have reduced the number of COVID-19 cases and deaths in the US.

genized. We derive the cost of infection and show numerically that the non-monotonicity of peak prevalence in the transmission rate extends to that environment.

**Related Literature.** [Capasso and Serio \(1978\)](#) introduced non-linear contact rates into the standard SIR model as a reduced form of modeling behavior. For a more recent overview of literature on non-linear contact rates, see [Funk et al. \(2010\)](#) and [Verelst et al. \(2016\)](#). A strong point for explicitly modeling behavior was made by [Ferguson \(2007\)](#).

[Reluga \(2010\)](#) and [Fenichel et al. \(2011\)](#) introduced preventive behavior into SIR models explicitly and provided numerical analyses of equilibrium trajectories.<sup>8</sup> [Fenichel \(2013\)](#) studied the differences in incentives for distancing between a decentralized economy and an economy governed by a social planner. [Chen \(2012\)](#) introduced an SIR model with a constant cost of infection and derived conditions on the contact functions that deliver uniqueness of the Nash equilibrium in each period for a given prevalence of the disease; yet, he did not establish whether that leads to uniqueness of equilibrium trajectories. These papers do not provide analytical results about the equilibrium trajectories under social distancing.

Analytical results about equilibrium trajectories of the SIR model with distancing are few. Closest to ours is work by [Dasaratha \(2020\)](#), [McAdams \(2020\)](#), [McAdams et al. \(2021\)](#), [Rachel \(2020a\)](#) and [Toxvaerd \(2020\)](#). The last two analyze a model of behavior with a linear cost of distancing and an endogenous time-varying cost of getting infected. They derive the necessary conditions for an equilibrium and offer two different paths that satisfy the necessary conditions, but stop short from proving that these are indeed equilibria.<sup>9</sup> Characterization of equilibria in their model, therefore, remains an open question. It remains unresolved whether an equilibrium in their model is unique and thus comparative statics are non-ambiguous.<sup>10</sup> [Dasaratha \(2020\)](#) analyzes a model with a constant cost of infection (like ours), but where the infected individuals do not necessarily know whether they are infected. The complexity of his model requires that he mostly focuses on local results rather than on the entire path.<sup>11</sup> [McAdams \(2020\)](#) and [McAdams et al. \(2021\)](#) propose a model in which an individual's benefit of social activities depends on the actions of other individuals and shows that complementarities in distancing choices

---

<sup>8</sup>Behavior has also been studied in other (non-SIR) models of diseases such as HIV/AIDS. An account of that literature is beyond the scope of our paper.

<sup>9</sup>They never establish the existence of a co-state variable that supports the equilibrium behavior (in particular, they never verify that the agents are not distancing when the equilibrium prescribes that they should not) and satisfies the transversality condition.

<sup>10</sup>[Toxvaerd \(2020\)](#) claims uniqueness, but his argument treats the problem as a single person decision problem rather than a non-cooperative game.

<sup>11</sup>[Engle et al. \(2021\)](#) study a behavioral SIR model with a constant cost of infection but with a different meeting rate. They empirically analyze the incidence data for the 2009 Swine Flu and the COVID-19 pandemic.

may lead to multiplicity of equilibria. [Gans \(2022\)](#) explores the implications of directly imposing on a model of behavior that the effective reproduction number satisfies  $R_t = 1$ . [Budish \(2020\)](#) studies the optimal interventions (e.g., mask mandates and closure of large indoor gatherings) subject to  $R_t \leq 1$  as a constraint in a static model. [Avery \(2021\)](#) studies the interaction between social distancing and vaccination. An excellent account of the rapidly growing literature is provided by [McAdams \(2021\)](#).

A wide pool of papers studies how policy interventions affect distancing and, through that, the spread of a disease. [Farboodi et al. \(2021\)](#) and [Rachel \(2020b\)](#) build on the work discussed above to study lockdown effectiveness and the possibility of a second wave occurring. [Toxvaerd and Rowthorn \(2020\)](#) compare individuals' and a planner's decisions to apply treatments and vaccinations as pharmaceutical interventions during an epidemic. [Giannitsarou et al. \(2021\)](#) provide numerical projections for the COVID-19 pandemic under waning immunity, based on a model with endogenous distancing. [Acemoglu et al. \(2021\)](#) and [Brotherhood et al. \(2020\)](#) study the importance of age composition in the COVID-19 pandemic. With the exception of [Rachel \(2020b\)](#), these papers focus on numerical solutions of rather involved models without establishing either the existence or the uniqueness of equilibria. While one can argue that the numerical algorithms are bound to lead to an equilibrium, or something close to it, the lack of uniqueness of equilibria reduces the credibility of welfare comparisons of various policies in those models.

## 2 The Model

We study behavior in an otherwise standard SIR model. A continuum of individuals, indexed by  $i$  and normalized to unity, are infinitely lived with time indexed by  $t \in [0, \infty)$ . Each individual can be in one of three states: susceptible, infected (and infectious), or recovered. Susceptible individuals might get infected, in which case they transition into the infected state. Infected individuals can recover but cannot become susceptible again.<sup>12</sup> Recovered individuals acquire permanent immunity. This model is suitable for viral diseases which are transmitted directly from human to human.<sup>13</sup> We denote the share of the population that is susceptible at time  $t$  by  $S(t)$ , infected by  $I(t)$  and recovered by  $R(t)$ .

At each moment in time, susceptible individual  $i$  chooses how much activity to engage

---

<sup>12</sup>We abstract from the issues of testing for infection as studied in [Deb et al. \(2022\)](#) and [Ely et al. \(2021\)](#).

<sup>13</sup>[Avery et al. \(2020\)](#) provide an excellent assessment of the SIR model from the economics point of view.

in, denoted by  $\varepsilon_i(t) \in [0, 1]$ . The individuals enjoy the activity, but it exposes them to the danger of infection; hence, termed exposure. The converse,  $d_i(t) := 1 - \varepsilon_i(t)$ , is the measure of distancing. While susceptible, an individual incurs a flow payoff,  $\pi_S$ . Distancing is uncomfortable and comes at a cost  $\frac{c}{2}(d_i(t))^2$ . The cost of getting infected is  $\eta > 0$ . Later in the paper, we explore the model where the cost of infection is determined endogenously and might vary over time.

Individuals meet through a pairwise-matching technology where each individual has an equal chance of meeting any other individual—regardless of which state they are in. The only matches with an infection risk are the ones between a susceptible and an infected individual. The rate at which a susceptible individual who chooses exposure level  $\varepsilon_i(t)$  meets an infected individual and gets infected at time  $t$  is  $\beta \varepsilon_i(t) I(t)$ , where  $\beta > 0$  is the transmission rate.<sup>14</sup> Finally, infected individuals recover at rate  $\gamma > 0$ .

At each point in time  $t$ , a susceptible individual  $i$  solves the problem

$$\max_{\varepsilon_i(t) \in [0, 1]} \pi_S - \frac{c}{2}(1 - \varepsilon_i(t))^2 - \beta I(t) \varepsilon_i(t) \eta. \quad (1)$$

Let  $\varepsilon(t) := \frac{1}{S(t)} \int_i \varepsilon_i(t) di$  be the average exposure of susceptible individuals at time  $t$ . Analogously, define  $d(t) := 1 - \varepsilon(t)$  as the average distancing at time  $t$ . Then, the model is governed by the following dynamics

$$\dot{S}(t) = -\beta \varepsilon(t) I(t) S(t), \quad (2)$$

$$\dot{I}(t) = \beta \varepsilon(t) S(t) I(t) - \gamma I(t), \quad (3)$$

$$\dot{R}(t) = \gamma I(t), \quad (4)$$

with the assumption that there is a seed of infected,  $I(0) = I_0 \in (0, 1)$ , and susceptible individuals,  $S(0) = S_0 = 1 - I_0$ . Since  $S$ ,  $I$  and  $R$  are the only three states  $S(t) + I(t) + R(t) = 1$  at each instance of time.

**Definition 1.** An equilibrium is a tuple of functions  $(S, I, R, (\varepsilon_i)_i)$  with the following two properties: (i)  $(S, I, R)$  follow (2), (3) and (4) with the initial condition  $(S(0), I(0), R(0)) = (S_0, I_0, 0)$ , where  $\varepsilon$  is the average exposure; and (ii) each  $\varepsilon_i$  solves (1), that is,  $\varepsilon_i$  is a best-response to  $(S, I, R)$ . An equilibrium is symmetric if  $\varepsilon = \varepsilon_i$  for all  $i$ .

The first-order condition to the individual's problem yields the individual's optimal

---

<sup>14</sup>We implicitly assume that infected individuals choose full exposure. Though strong, the assumption is not as stark as it might at first seem. It is straightforward to accommodate exposure of infected with some parameter  $e$ , as long as it is fixed over time. Then, the same model as ours can be obtained by defining  $\tilde{\beta} = e\beta$ .

distancing choice

$$d_i(t) := 1 - \varepsilon_i(t) = \min \left( \frac{\eta\beta}{c} I(t), 1 \right). \quad (5)$$

When  $\frac{\eta\beta}{c} I(t)$  exceeds unity, individuals fully distance. Distancing at time  $t$  depends only on the infected population at time  $t$ —up to constants  $\beta$ ,  $c$  and  $\eta$ .<sup>15</sup> In equilibrium,  $\varepsilon_i = \varepsilon$  for all  $i$ , that is, every equilibrium is symmetric. By equation (5), exposure in a symmetric equilibrium is

$$\varepsilon(t) = \max \left( 1 - \frac{\eta\beta}{c} I(t), 0 \right).$$

Plugging (5) into the SIR dynamics yields a system of differential equations with an initial condition. All the proofs are collected in Appendix A.

**Proposition 1.** *An equilibrium exists, is unique and symmetric.*

The above result establishes uniqueness of the equilibrium in our model. In a similar model, Chen (2012) established uniqueness of a Nash equilibrium in each period for every state of the model. However, he did not study whether that leads to uniqueness of equilibrium trajectories.

### 3 Analysis

We establish several qualitative properties of the equilibrium and the resulting epidemic dynamics. First, observe that if  $\varepsilon(\tilde{t}) > 0$  for some  $\tilde{t}$ , then  $\varepsilon(t) > 0$  for all  $t > \tilde{t}$ . This follows from the observation that as long as  $\varepsilon(t) \in (0, \gamma/\beta)$ ,  $\dot{I}(t) < 0$  and thus  $\dot{\varepsilon}(t) > 0$ .  $\varepsilon(t)$  can, therefore, be 0 only at the beginning. To avoid this tedious contingency, we will often assume  $\varepsilon(0) > 0$ , or equivalently,  $I_0 < c/(\beta\eta)$ . Next, we establish that the number of active cases peaks at most once.

**Proposition 2.** *If  $\hat{t}$  is such that  $\dot{I}(\hat{t}) = 0$ , then  $\ddot{I}(\hat{t}) < 0$ .*

Proposition 2 implies that if  $I$  has a critical point, this critical point has to be a local maximum. If  $I$  does not have a critical point, it decreases throughout. Together with the continuous differentiability of  $I$ , this implies that  $I$  can have at most one peak. The

---

<sup>15</sup>A body of literature studies an extended SIR model where distancing behavior instead depends on deaths; see Weitz et al. (2020), Atkeson (2021) and Atkeson et al. (2021).



infection either immediately dies out or becomes an epidemic with a single peak.<sup>16</sup>

In the standard SIR model, the infection propagates itself ( $\dot{I}(0) > 0$ ) only if the basic reproduction number,  $R_0 := \frac{\beta}{\gamma}S_0$ , is larger than 1; see [Heesterbeek and Dietz \(1996\)](#).<sup>17</sup> However, the observed and measurable variable is how many secondary infections have been caused, given an individual's behavior. To capture this, we define the *behavioral basic reproduction number* as:

$$R_0^b := \frac{\beta}{\gamma}S_0\varepsilon(0). \quad (6)$$

Notice that  $R_0^b = \varepsilon(0)R_0$ ; the concept of a basic reproduction number that depends on the individuals' behavior was introduced in [Fenichel et al. \(2011\)](#). Equation (3) at  $t = 0$  can now be rewritten as  $\dot{I}(0) = \frac{I_0}{\gamma}(R_0^b - 1)$ . Therefore, the infection spreads,  $\dot{I}(0) > 0$ , if and only if  $R_0^b > 1$ , paralleling a similar result in the model without distancing. However, while in the standard SIR model  $R_0$  is increasing in  $\beta$ , the behavioral basic reproduction number  $R_0^b$  is non-monotonic and, in particular, concave. This finding has important implications on which types of an infection will spread.

**Proposition 3.** *Fix  $I_0 \in (0, 1)$ . Then,  $\dot{I}(0) > 0$  if and only if  $R_0^b > 1$ . Moreover:*

- (i) *if  $I_0 < \frac{1}{1+\frac{4\eta\gamma}{c}}$ , there exist  $\underline{\beta}$  and  $\bar{\beta}$ , with  $\frac{\gamma}{1-I_0} < \underline{\beta} < \bar{\beta} < \frac{c}{\eta I_0}$ , such that  $\dot{I}(0) > 0$  if and only if  $\beta \in (\underline{\beta}, \bar{\beta})$ .*
- (ii) *If  $I_0 \geq \frac{1}{1+\frac{4\eta\gamma}{c}}$ , then  $\dot{I}(t) \leq 0$  for all  $t$ .*

In the standard SIR model,  $\beta$  must be high enough ( $\beta > \frac{\gamma}{S_0}$ ) for the infection to spread. In the model with distancing, instead, the transmission rate has to be large enough to also overcome the initial distancing:

$$\beta > \frac{\gamma}{(1-d(0))S_0} > \frac{\gamma}{S_0}.^{18}$$

Our model predicts that a higher transmission rate is needed for the prevalence to increase from the start than in the standard SIR model. What differentiates the model with distancing even more starkly is that the prevalence starts to decrease immediately if the transmission rate is too high. If the disease is highly contagious, individuals are much more cautious, up to the point where their resolute distancing alone is sufficient to stop

<sup>16</sup>Eventually, the disease inevitably dies out;  $I_\infty := \lim_{t \rightarrow \infty} I(t) = 0$ . To see this, note that the fraction  $R(t)$  of recovered individuals is increasing and bounded above. Hence,  $\lim_{t \rightarrow \infty} R(t)$  exists. Together with equation (4), this implies that  $I_\infty = 0$ .

<sup>17</sup>Depending on the source  $R_0$  is defined either as  $\beta/\gamma$  or  $\beta S_0/\gamma$ . We use the latter definition as it allows for an easier presentation of results.

<sup>18</sup>It should be noted that  $d(0)$  depends on  $\beta$  as well.

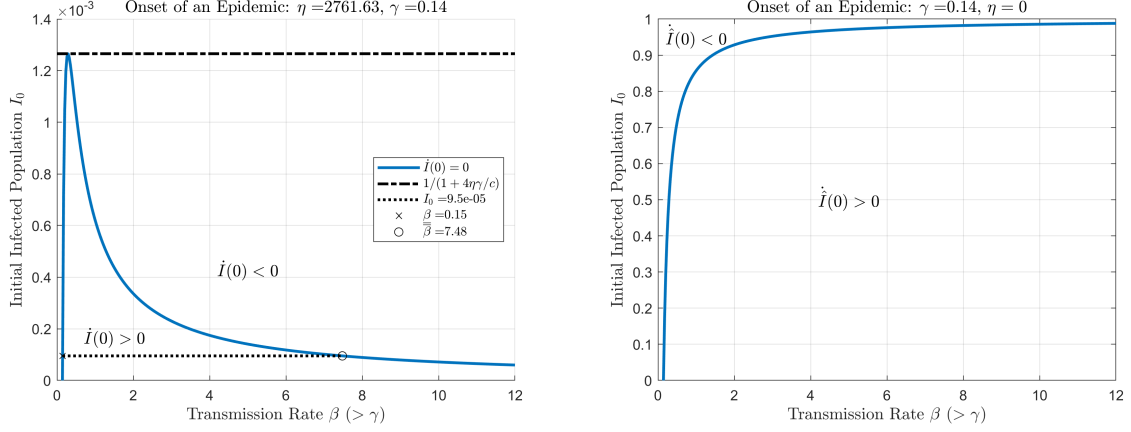


Figure 1: *The Onset of an Epidemic. Left: Our SIR model; Right: Standard SIR model.* The solid line in each panel depicts the combination of  $(\beta, I_0)$  with  $\dot{I}(0) = 0$ .

the disease in its tracks from the onset.<sup>19</sup> Due to preventive behavior, the prevalence increases beyond the initial seed of infection only if its transmission rate is large enough, but not too large, as illustrated in the left panel of Figure 1.<sup>20</sup>

Proposition 3 derives conditions on the transmission rate for the prevalence to take off. The same question can be analyzed along other dimensions. For example, the CDC has adopted a categorization for influenza viruses along the severity-transmissibility dimensions (see Reed et al., 2013). In our model, this can be interpreted as categorizing the combination of the cost of infection,  $\eta$ , and the transmission rate,  $\beta$ .

Recall that in the SIR model without behavior—which is nested in our model as the case  $\eta = 0$ —the epidemic takes off whenever  $\beta S_0 > \gamma$ . As the cost of infection,  $\eta$ , increases, individuals’ distancing incentives start to matter for the onset of an epidemic. In particular, for a fixed  $\beta$ , the higher the cost of infection, the more individuals engage in distancing. If the cost of infection becomes very large, it prevents the prevalence from rising altogether:  $\dot{I}(0) < 0$ . There is a cutoff cost of infection such that the prevalence will never increase when  $\eta > \frac{c}{4\gamma} \frac{S_0}{I_0}$  as getting infected is so costly for individuals that their distancing behavior compensates for any transmission rate  $\beta$ . By implication, the prevalence increases from the beginning only if the  $(\beta, \eta)$ -combination is intermediate. For a given  $\beta$ , the infection cost must not be too high; while for a given  $\eta$ , the transmission rate must neither be too high nor too low. The existence of an upper and a lower bound for  $\beta$  follows the same intuition as the one applying for Proposition 3.

<sup>19</sup>An informal discussion of the role of disease-intrinsic parameters and its effect on the outbreak of an epidemic can be found in Christakis (2020).

<sup>20</sup>We use parameters for COVID-19 in our simulations. A summary and justification of the parameters chosen can be found in Appendix B. We also describe our numerical algorithm there.

### 3.1 Solution Paths

A solution path in the  $(S, I)$ -space is a graph showing how the number of infected individuals changes with the number of susceptible individuals. To find the solution path  $(S, I) := (S(t), I(t))_{t \geq 0}$  in the phase space, one derives the quotient differential equation

$$\frac{dI}{dS} = -1 + \frac{\gamma}{\beta} \frac{1}{S \max\left(1 - \frac{\beta\eta}{c}I, 0\right)} \quad (7)$$

by dividing equation (3) by equation (2) and using (5) for  $\varepsilon$ .

**Proposition 4.** Suppose  $d(0) < 1$ .<sup>21</sup> The solution path  $(S, I)$  is implicitly determined by

$$S = \frac{\exp\left(-\frac{\beta^2\eta}{2\gamma c} \left(S + I - \frac{c}{\beta\eta}\right)^2\right)}{\exp\left(-\frac{\beta^2\eta}{2\gamma c} \left(1 - \frac{c}{\beta\eta}\right)^2\right) \frac{1}{S_0} + 2\beta\sqrt{\frac{\eta}{2\gamma c}} \int_{\beta\sqrt{\frac{\eta}{2\gamma c}}\left(S+I-\frac{c}{\beta\eta}\right)}^{\beta\sqrt{\frac{\eta}{2\gamma c}}\left(1-\frac{c}{\beta\eta}\right)} e^{-v^2} dv} \quad (8)$$

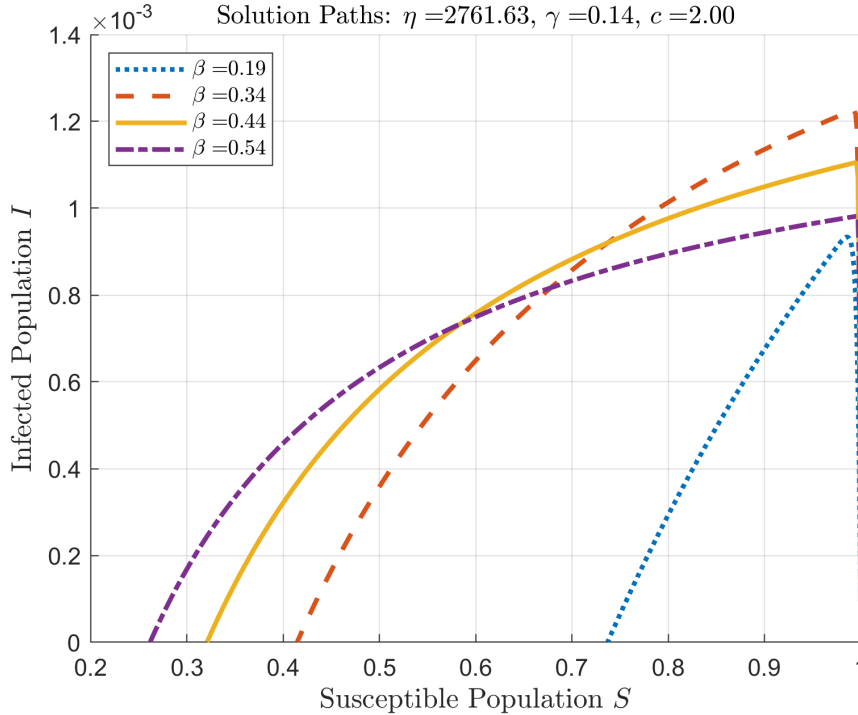


Figure 2: *Solutions Paths for Different Transmission Rates.*

Figure 2 depicts the solution paths in the phase space for different transmission rates.

<sup>21</sup>The assumption is made for ease of exposition directly on  $d(0)$ ; Formula (5) provides the corresponding assumptions on primitives. If  $d(0) = 1$ , then individuals engage in full distancing up to some point, after which an equation analogous to (8) determines the dynamics of the epidemic.

Notice that the solution paths are not monotonically ordered.

We show that the solution path  $(S, I)$  moves upwards as the cost of infection  $\eta$  decreases. The SIR model without distancing can be recovered as the special case of our model with no cost of infection ( $\eta = 0$ ).<sup>22</sup> Hence, for any level of the susceptible population, the corresponding number of active infections is lower in the SIR model with distancing. In particular, the peak prevalence in our model is below that of the standard model. This comparison is depicted in Figure 3.

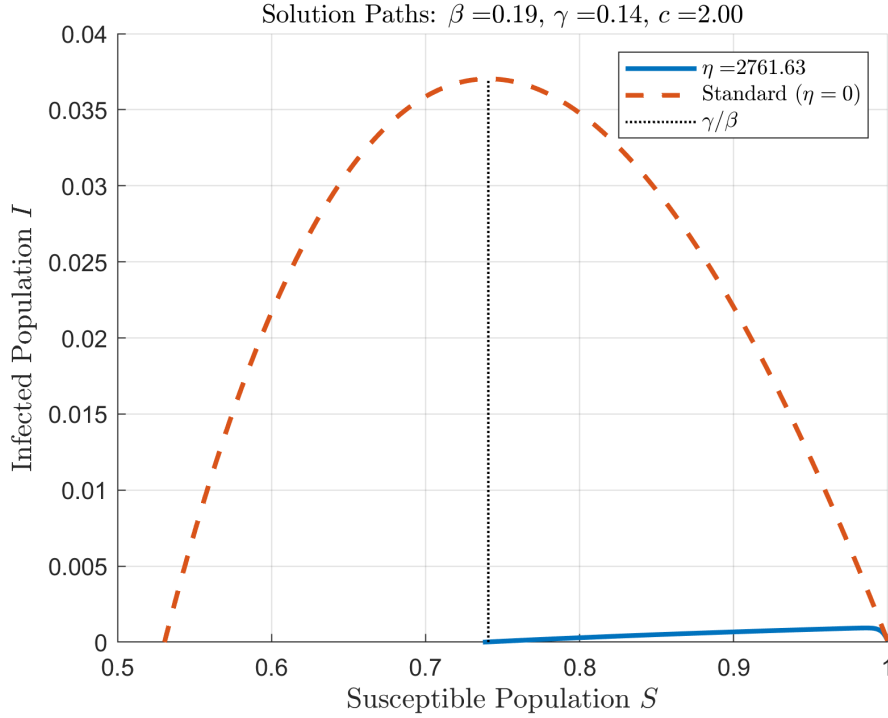


Figure 3: *Solution Paths for the SIR Models With and Without Behavior.* The solid curve depicts the solution path for our SIR model with behavior. The dashed curve depicts the solution path of the SIR model without behavior.

**Proposition 5.** *Assume that  $d(0) < 1$ . Then, if  $\eta$  decreases, the solution path  $(S, I)$  moves upwards (lies above the original solution path). In particular, the solution path  $(\hat{S}, \hat{I})$  of the standard SIR model, which is associated with  $\eta = 0$ , lies above the original solution path.*

Increasing the cost of infection,  $\eta$ , raises the incentives to distance and pushes the solution path down, that is, it decreases the infected population at any level of susceptibles

<sup>22</sup>We denote by  $(\hat{S}, \hat{I}, \hat{R})$  the proportion of susceptible, infected, and recovered individuals in the standard SIR model. The dynamics of the standard SIR model is obtained by replacing  $(S, I, R, \varepsilon)$  with  $(\hat{S}, \hat{I}, \hat{R}, 1)$  in equations (2), (3) and (4). The solution path  $(\hat{S}, \hat{I})$  of the standard SIR model is captured by  $\frac{d\hat{I}}{d\hat{S}} = -1 + \frac{\gamma}{\beta} \frac{1}{\hat{S}}$ .

in the phase space. Since the SIR model without distancing is the special case with  $\eta = 0$ , its solution path is above the path for any  $\eta > 0$ .

### 3.2 Peak Prevalence

The peak prevalence of an epidemic has profound consequences on the overall provision of health care services. A large number of infected individuals may lead to an overwhelming demand for personal protective equipment, such as face masks, and for medical devices, such as ICU beds and ventilators. The shortage of medical resources, in turn, may cause a suboptimal treatment and health care coverage; see, for example, [Schoch-Spana \(2001\)](#) for the 1918 influenza pandemic, [Ferguson et al. \(2020\)](#) for the COVID-19 pandemic, and [Reed et al. \(2013\)](#) for influenza epidemics. The high demands of the epidemic on the health system also divert medical resources from other important activities. What is more, healthcare workers themselves are at high risk of infection.<sup>23</sup> The peak prevalence is, therefore, of paramount interest for epidemic preparedness and optimal policy responses.

When the epidemic takes off ( $\dot{I}(0) > 0$ ), Proposition 2 implies that the prevalence is maximized when  $\dot{I}(t) = 0$ , that is, when

$$\varepsilon(t)S(t) = \gamma/\beta.$$

Denote by  $I^* := \max_t I(t)$  the peak prevalence. In the standard SIR model with  $R_0 > 1$ , the peak prevalence  $\hat{I}^* := \max_t \hat{I}(t)$  is given by  $\hat{I}^* = 1 - \frac{\gamma}{\beta} + \frac{\gamma}{\beta} \log\left(\frac{\gamma}{\beta S_0}\right)$ ; see, for example, [Hethcote \(2008\)](#) or [Brauer and Castillo-Chavez \(2012\)](#). The peak prevalence is attained when the population  $\hat{S}(t)$  of susceptibles reaches the threshold of herd immunity  $\frac{\gamma}{\beta}$ . When the peak prevalence  $I^*$  of our model is attained, the population  $S(t)$  of susceptibles is larger than  $\frac{\gamma}{\beta}$ . Since the solution path  $(S, I)$  is below the path  $(\hat{S}, \hat{I})$ , our model predicts a smaller peak prevalence than the SIR model without behavior,  $I^* < \hat{I}^*$ .

We study how the peak prevalence changes with the parameters  $\beta$  and  $c$ . To focus on the case in which the infection can take place, we assume  $I_0 < \frac{1}{1 + \frac{4\eta\gamma}{c}}$ ; see Proposition 3. When this assumption fails, the prevalence decreases from the start irrespective of the transmission rate:  $I^* = I_0$ .

**Proposition 6.** *The following holds:*

---

<sup>23</sup>[Elston et al. \(2017\)](#) survey the health impact of the 2014-15 Ebola outbreak in West Africa. For Sierra Leone, they report a 20 % decrease in measles coverage, an overall 20-23 % decrease in deliveries and Caesarian sections. 10.7 % of the healthcare workforce were infected and 6.9 % died from Ebola virus disease.

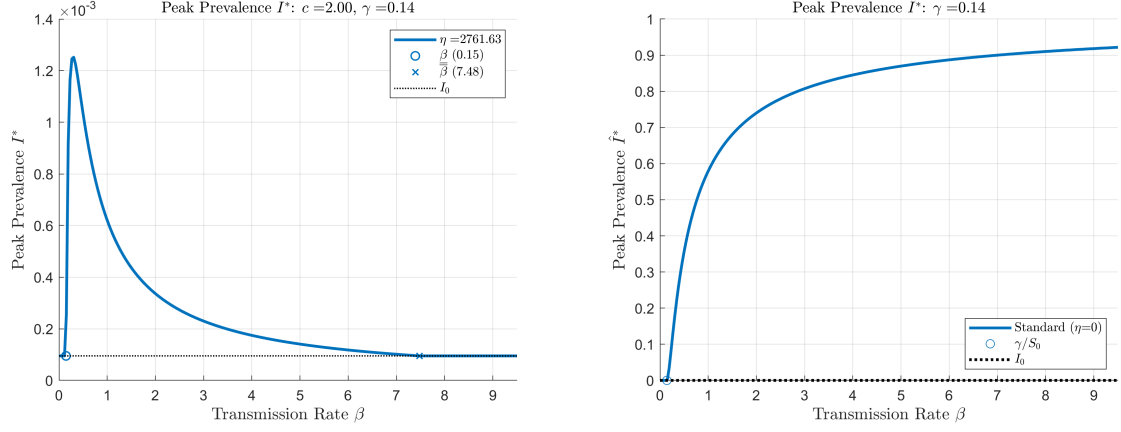


Figure 4: *Peak Prevalence as Function of Transmission Rate for the SIR Models with and without Behavior.* The left panel depicts the peak prevalence of the SIR model with behavior. The right panel depicts the peak prevalence of the SIR model without behavior.

- (i) Fix  $\gamma$ ,  $c$  and  $\eta$  and let  $I_0 < \frac{1}{1+\frac{4\eta\gamma}{c}}$ . Then, there exist  $\beta_1 < \beta_2$  such that  $I^*$  is increasing in  $\beta$  for  $\beta \in (\underline{\beta}, \beta_1)$  and decreasing in  $\beta$  for  $\beta \in (\beta_2, \bar{\beta})$ .
- (ii) The peak prevalence  $I^*$  is non-decreasing in  $c$ . It is strictly increasing in  $c$  whenever  $\dot{I}(0) > 0$ .

In the SIR model without behavior, the peak prevalence  $\hat{I}^*$  is monotonically increasing in the transmission rate  $\beta$ , as illustrated in the right panel of Figure 4. In contrast, in our model, a higher transmission rate leads to, ceteris paribus, more distancing. This effect can be so strong that a higher transmission rate reduces the peak prevalence and flattens the infection curve. Indeed, peak prevalence is non-monotonic in  $\beta$ ; see the simulation in the left panel of Figure 4. When the transmission rate is low, the peak of the infection increases in  $\beta$ . In contrast, when the transmission rate is high, the peak prevalence decreases in  $\beta$ .<sup>24</sup> A measure imposed to fight the epidemic through a reduction in  $\beta$  could, therefore, have a daunting short-run effect; for example, if the potential resulting increase in prevalence leads to stress of the health care system.<sup>25</sup> We want to emphasize that this

<sup>24</sup>Avery (2021) constructs an upper bound on the prevalence (that is never attained) and shows that it is non-monotonic in the transmission rate. Since the bound is not tight, its implications on the behavior of peak prevalence are unclear.

<sup>25</sup>This indirect effect of a measure reducing individual risk on taking fewer precautions is reminiscent of risk compensation introduced by Peltzman (1975); for a survey see Hedlund (2000). The importance of risk compensation has been recognized to play an important role in the economic epidemiology literature at least since Philipson and Posner (1993). Kremer (1996) and Geoffard and Philipson (1996) showed in an SI model that a policy intended to decrease HIV prevalence can lead to exactly the opposite due to the riskier behavior of individuals. Greenwood et al. (2019) construct a search model and quantify such behavior. It should be pointed out that these findings were developed in models very different from the SIR model under consideration here and that those of the above-mentioned papers by no means imply our results.

effect arises only for a subset of potential parameters. In particular, [Chernozhukov et al. \(2021\)](#) show that the introduction of mask mandates—a  $\beta$ -reducing policy—reduced the number of active cases in 2020 during the COVID-19 pandemic in the US. At the same time, [Knotek II et al. \(2020\)](#) and [Yan et al. \(2021\)](#) report evidence that some individuals view mask-wearing as a substitute for physical distancing. For a more comprehensive discussion of mask-wearing, see [Howard et al. \(2021\)](#).

In contrast, an increase in the cost of distancing always increases the peak prevalence. A higher cost of distancing leads to less distancing, all else equal. The disparity in effects of changes in  $c$  and  $\beta$  on peak prevalence can be seen by studying how the slope of the solution path at a fixed point in the phase space varies with changes in the two parameters. Differentiating the slope with respect to the cost of distancing parameter yields

$$\frac{\partial}{\partial c} \left( \frac{dI}{dS} \right) = -\frac{\gamma}{\beta \varepsilon^2 S} \frac{\partial \varepsilon}{\partial c} < 0,$$

where the inequality follows from the observation that, for a fixed  $I$ , the exposure increases if the cost of distancing increases. Importantly, the only effect an increase in the cost has on the solution path is through the change in distancing. By implication, the slope of a solution path with a higher cost is smaller than the slope of a solution path with a smaller cost of distancing at any point of intersection. The fact that they start from the same point,  $(S_0, I_0)$ , then implies that everywhere else the solution path corresponding to a higher cost must be above the one with the lower cost.

The change in the transmission rate, though, has a more nuanced effect. Differentiating the slope of the solution path at a fixed point yields

$$\frac{\partial}{\partial \beta} \left( \frac{dI}{dS} \right) = -\frac{\gamma}{\beta^2 \varepsilon S} - \frac{\gamma}{\beta \varepsilon^2 S} \frac{\partial \varepsilon}{\partial \beta}.$$

An increase in  $\beta$  has two effects. Holding everything else fixed, it results in more secondary infections from each infected individual, thereby increasing the speed of the spread of the disease. Such a direct effect is absent from changes in the cost of distancing. The second, indirect, effect is due to the distancing response to the change in the transmission rate. A more infectious disease results in more distancing and thus dampens the evolution of the epidemic. The two effects run in opposite directions. Depending on which of the two dominates, an increase in  $\beta$  can lead to either a smaller or a larger slope of the solution path.

**Public Policies.** A change in  $\beta$  can be interpreted in two ways. First, models with

two different  $\beta$ s can be thought of as comparing the trajectories of two different diseases. Our result implies that a disease with a higher transmission rate might, indeed, lead to a lower peak prevalence due to the effect of distancing on the epidemic path.

Second, the above finding has an important implication on how various preventive policies should be studied in models with an epidemiological component. Such models commonly adopt one of two apparatuses: behavior is either modeled implicitly by changes in  $\beta$  in the standard SIR model (see, for example, [Capasso and Serio, 1978](#); [Brauer, 2019](#); [Kruse and Strack, 2020](#)) or by directly imposing behavioral changes (see, for example, [Acemoglu et al., 2021](#); [Alvarez et al., 2021](#); [Farboodi et al., 2021](#); [Rachel, 2020a](#)). Our results highlight the importance of differentiating between policies that change the transmission rate and policies that change the cost of distancing.<sup>26</sup> For example, if a government imposes temporary restaurant closures to slow the spread of the disease, this gives individuals fewer reasons to go out and should be modeled as a decrease in the cost of distancing, and not as a decrease in the transmission rate directly. Holidays or vacation periods can be seen as *increases* in the cost of distancing. Closest to this, ([Fenichel et al., 2011](#)) model a public policy intervention as a change in the payoff accruing from contacts.

Moreover, our analysis highlights an additional consideration that is rarely taken into account. The effect of a particular measure to decrease the cost of distancing depends on the current epidemic state. If there are few actively infected individuals, the same measure has a much lower effect on the individuals' distancing decision than when there are many actively infected individuals. This implies that if one models policy as directly affecting the level of exposure, the cost of doing so should also depend on the epidemic's state.

### 3.3 Final Size of the Epidemic

We have established the effect of changes in the transmission rate and the cost of distancing on the peak prevalence. Turning to the long-run effects of behavior, an important characteristic of a disease is  $S_\infty := \lim_{t \rightarrow \infty} S(t)$ , the number of remaining susceptible individuals once the epidemic is over. The converse,  $1 - S_\infty$ , is referred to as the final size of the epidemic.

---

<sup>26</sup>Strictly speaking, our comparative statics correspond to a public policy that is implemented from the outset of the epidemic to the end. However, our model can accommodate a policy that takes effect later during the epidemic if it remains in effect indefinitely. In particular, a policy introduced at time  $\hat{t}$  can be analyzed by re-parametrizing our model with the new initial conditions  $(I_{\hat{t}}, S_{\hat{t}}, R_{\hat{t}})$ . The basic idea extends even to public policies that are enacted for a limited amount of time: a change in  $\beta$  has two opposing effects while a change in  $c$  has one effect.



A lesson from the SIR model without distancing is that once  $S$  falls below the threshold  $\gamma/\beta$ , the prevalence decreases. Indeed, it decreases faster than the number of susceptible individuals,  $\hat{S}_\infty := \lim_{t \rightarrow \infty} \hat{S}(t) \in (0, \gamma/\beta)$ . Therefore, at the end of the epidemic, a strictly positive fraction (but smaller than  $\gamma/\beta$ ) of the population remains susceptible. The following proposition establishes properties of the final size of an epidemic when individuals can undertake preventive behavior.

**Proposition 7.** *The final size of the epidemic,  $S_\infty$ , satisfies  $\hat{S}_\infty \leq S_\infty < \frac{\gamma}{\beta}$  and is decreasing in  $\beta$ , for  $\beta \in [0, \frac{c}{\eta I_0}]$ , and decreasing in  $c$ .*

The model with distancing predicts a smaller size of the epidemic than without. More importantly, as long as  $\dot{I}(0) > 0$ , the size of the epidemic is monotone in  $\beta$  and  $c$ .<sup>27</sup> A higher  $\beta$  leads to a larger amount of individuals that contract the disease during the epidemic, as does an increase in the cost of distancing. This result is somewhat surprising in light of the result that the peak prevalence is non-monotonic in the transmission rate. Together, the two results establish the existence of a region of transmission rates where an increase in the transmission rate leads to a decrease in the peak prevalence but an increase in the final size of the epidemic.

While Proposition 3 established that the prevalence decreases throughout if  $\beta$  is too large (larger than  $\bar{\beta}$ ), Proposition 7 establishes that the size of the epidemic nevertheless increases in the transmission rate—even at such large  $\beta$ . In such a case, the disease is spreading very slowly through the population. From a practical viewpoint, if the spread of infection is sluggish, a vaccine or at least a treatment is likely to be developed that would significantly decrease, if not eliminate, the cost and discomfort brought upon by the infection.

Propositions 6 and 7 can be interpreted through two different prisms. First, through the prism of public policies: while policies that affect  $\beta$  might have perverse effects in the short run—e.g., a decrease in the transmission rate,  $\beta$ , may lead to an increase in peak prevalence—, in the long run, they will unequivocally decrease the number of infected individuals. Therefore, one needs to be circumspect if the medical capabilities are at or close to the capacity in the short run.

The second prism corresponds to a comparison of various diseases. Our results imply that a more transmissible disease always infects a larger number of individuals. However, it may result in a smaller peak prevalence than a less transmissible disease.

---

<sup>27</sup>Note that  $\bar{\beta} < \frac{c}{\eta I_0}$ .

## 4 Endogenous Cost of Infection

In this section, we present a model with an endogenous cost of infection, develop the formula for the cost, show how the model with fixed cost can be used to bound the model with endogenous cost, and provide numerical support for the non-monotonicity of peak prevalence in the transmission rate in the model with an endogenous cost of infection.

As before, the individuals at each point in time decide to which extent to distance, which determines how likely they are to get infected. An individual's flow payoff from being in state  $\theta \in \{S, I, R\}$  is  $\pi_\theta$ . We assume  $\pi_S \geq \pi_R \geq \pi_I$ .<sup>28</sup> The endogeneity of the cost of infection results from differences in the flow payoff across the states and the individuals taking future infection risks into account. The individual discounts the future at rate  $\rho > 0$ .

A susceptible individual  $i$  with exposure  $\varepsilon_i(t)$  enjoys the instantaneous payoff  $\pi_S - \frac{c}{2}(1 - \varepsilon_i(t))^2$ . Let  $1 - p_i(t)$  be the probability of being susceptible at time  $t$  and, thus,  $p_i(t)$  the probability that an individual has become infected in the past. Then,  $\dot{p}_i(t)$  represents the rate at which susceptible individuals become infected

$$\dot{p}_i(t) = \varepsilon_i(t)\beta I(t)(1 - p_i(t)), \quad (9)$$

with  $p_i(0) = 0$ ; since we model the behavior of susceptible individuals, the probability that they are infected at the outset is zero. Once an individual gets infected, her progression to recovery is independent of her behavior. Her continuation payoff from the moment she became infected,  $V_I$ , is:

$$V_I = \frac{1}{\rho + \gamma} \left( \pi_I + \frac{\gamma}{\rho} \pi_R \right). \quad (10)$$

See Remark 1 in Appendix A for the derivation.

A susceptible individual who faces average exposure  $\varepsilon$  from her peers solves the prob-

---

<sup>28</sup>Models with endogenous cost of infection have been presented in Reluga (2010); Fenichel et al. (2011); Fenichel (2013); McAdams (2020); Rachel (2020a); Toxvaerd (2020), among others. Yet, analytical characterizations of equilibria are rather elusive.

lem

$$\begin{aligned}
& \max_{\varepsilon_i(\cdot) \in [0,1]} \int_0^\infty e^{-\rho t} \left\{ (1 - p_i(t)) \left[ \pi_S - \frac{c}{2} (1 - \varepsilon_i(t))^2 \right] + p_i(t) \rho V_I \right\} dt \\
& \text{s.t.} \\
& \dot{p}_i(t) = \beta \varepsilon_i(t) I(t) (1 - p_i(t)), \\
& p_i(0) = 0,
\end{aligned} \tag{11}$$

and the underlying dynamics given by equations (2), (3) and (4) with the initial condition  $(S(0), I(0), R(0)) = (1 - I_0, I_0, 0)$  and  $I_0 \in (0, 1)$ . The individual's payoff can be thought of as the expected value of being susceptible or infected at each point in time where the flow payoff of an infected individual is  $\rho V_I$ . An individual's behavior affects her probability of infection directly, but none of the population dynamics as she is small.

We study symmetric equilibria (equilibria for short).

**Definition 2.** A symmetric equilibrium is a tuple of functions  $(S, I, R, (\varepsilon_i, p_i)_i)$  with the following three properties: (i)  $(S, I, R)$  follow (2), (3) and (4) with the initial condition  $(S(0), I(0), R(0)) = (S_0, I_0, 0)$ , where  $\varepsilon$  is the average exposure; (ii) each  $\varepsilon_i$  solves (11), that is,  $\varepsilon_i$  is a best-response to  $(S, I, R)$ , where the average exposure  $\varepsilon$  is induced by  $(\varepsilon_j)_{j \neq i}$ ; and (iii)  $\varepsilon_i = \varepsilon$  for all  $i$ .

In equilibrium, each  $p_i$  is determined by the average exposure  $\varepsilon$  and  $I$ , and thus  $p = p_i$  for each  $i \in [0, 1]$ . For ease of exposition, we denote an equilibrium by  $(S, I, R, \varepsilon, p)$ .

**Assumption 1.**  $\pi_S - \frac{c}{2} > \rho V_I$ .

Even if a susceptible individual is fully distancing, her flow payoff of being susceptible is greater than the flow payoff of being infected. The current-value Hamiltonian of problem (11) is

$$\mathcal{H}_i = (1 - p_i(t)) \left[ \pi_S - \frac{c}{2} (1 - \varepsilon_i(t))^2 \right] + p_i(t) \rho V_I - \eta_i(t) \beta \varepsilon_i(t) I(t) (1 - p_i(t)),$$

where  $\eta_i(t)$  is the current-value co-state variable.<sup>29</sup> It represents the marginal cost of an increase in the probability of being infected at time  $t$ . The optimality condition with respect to exposure  $\varepsilon_i(t)$  at time  $t$  is

$$\frac{\partial \mathcal{H}_i}{\partial \varepsilon_i(t)} = (1 - p_i(t)) [c(1 - \varepsilon_i(t)) - \beta \eta_i(t) I(t)] = 0.$$

---

<sup>29</sup>Note that we define the co-state as the negative of the usual co-state to interpret it as a cost of infection rather than as benefit of being susceptible to relate it directly to our constant cost of infection model.

It can be verified that  $p_i(t) < 1$ ; see Remark 2 in Appendix A. Thus, the optimality condition delivers equilibrium distancing

$$d_i(t) = \frac{\beta}{c} \eta_i(t) I(t), \quad (12)$$

provided that the entire distancing path admits an interior solution, i.e., that  $d_i(t) \in [0, 1]$  for all  $t$ . One should keep in mind that the marginal cost of an increased probability of infection,  $\eta_i(t)$ , is positive due to Assumption 1. The extent to which an individual distances is, ceteris paribus, increasing in the infection rate,  $\beta$ , and the size of the infected population,  $I(t)$ , and the co-state,  $\eta_i(t)$  and decreasing in the cost parameter,  $c$ . Importantly, the decisions today influence the probability of getting infected both today and in the future, which in turn affects the distancing decisions today—a fact that is captured by the co-state  $\eta_i(t)$ . The current-value co-state variable  $\eta_i$  follows the adjoint equation

$$\begin{aligned} \dot{\eta}_i(t) &= \rho \eta_i(t) + \frac{\partial \mathcal{H}_i}{\partial p_i(t)} \\ &= \eta_i(t) (\rho + \varepsilon_i(t) \beta I(t)) + \left( \pi_S - \frac{c}{2} (1 - \varepsilon_i(t))^2 - \rho V_I \right). \end{aligned} \quad (13)$$

The transversality condition is  $\lim_{t \rightarrow \infty} e^{-\rho t} \eta_i(t) = 0$ . In equilibrium,  $\eta = \eta_i$  for all  $i$ . Using the adjoint equation and the transversality condition, we solve for  $\eta$ .

**Lemma 1.** *Suppose that the rest of the population is following the strategy  $\varepsilon$ , and  $\varepsilon_i$  is the individual  $i$ 's best response. Then*

$$\eta_i(t) = \int_t^\infty e^{-\rho(s-t)} \frac{1 - p_i(s)}{1 - p_i(t)} \left( \pi_S - \frac{c}{2} (1 - \varepsilon_i(s))^2 - \rho V_I \right) ds. \quad (14)$$

Let  $(S, I, R, \varepsilon, p)$  be an equilibrium. Then

$$\eta(t) = \int_t^\infty e^{-\rho(s-t)} \frac{S(s)}{S(t)} \left( \pi_S - \frac{c}{2} (1 - \varepsilon(s))^2 - \rho V_I \right) ds. \quad (15)$$

We term  $\pi_S - \frac{c}{2} (1 - \varepsilon(t))^2 - \rho V_I$  the *susceptibility premium* at time  $t$ . It is the difference in flow payoffs between being susceptible and being infected. The cost of getting infected,  $\eta(t)$ , is the discounted value of the susceptibility premium over time weighted by the conditional probability of being susceptible at each time in the future,  $s \geq t$ ,  $\frac{S(s)}{S(t)}$ . Distancing over a period of time reduces the quality of life and, thus, the susceptibility premium. However, it also decreases the probability that the individual will get infected and rewards her with the premium for a longer period of time. The functional form of  $\eta_i$  demonstrates the difficulty of the dynamic problem. Optimal exposure at time  $t$  depends on the exposure of the remaining individuals through the effect it has on the spread of the infection, as well as on the exposure of individual  $i$  at each instance in the

future.

Alternatively, one can decompose  $\eta$  in two parts

$$\eta(t) = (V_S(t) - V_I(t))$$

where

$$V_S(t) = \int_t^\infty e^{-\rho(s-t)} \left( \frac{S(s)}{S(t)} \left( \pi_S - \frac{c}{2}(1 - \varepsilon(s))^2 \right) + \left( 1 - \frac{S(s)}{S(t)} \right) \rho V_I \right) ds$$

is the continuation payoff of being susceptible and

$$V_I(t) = V_I,$$

is the continuation payoff of being infected.

The above discussion implies that analytically characterizing the set of equilibria is untenable. To verify whether a distancing function  $\varepsilon$  can be part of an equilibrium, one needs to posit that the individuals use it, derive  $S$ ,  $I$ ,  $R$  and  $\eta$ , and then verify that  $\varepsilon$  is indeed a best reply given the dynamics. This task is made more challenging by the fact that even the SIR model without distancing does not have a tractable closed-form solution and that  $\eta$  is pinned down only in the limit rather than at any point.

However, we can make use of the model with an endogenous cost of infection to inform our parameter choices in the constant cost of infection model. The following lemma provides bounds for  $\eta$ , which enable us to connect the two models.

**Lemma 2.** *Let  $(S, I, R, \varepsilon, p)$  be an equilibrium. Then*

$$\frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho + \beta} \leq \eta(t) \leq \frac{\pi_S - \rho V_I}{\rho}, \quad (16)$$

and

$$\lim_{t \rightarrow \infty} \eta(t) = \frac{\pi_S - \rho V_I}{\rho}. \quad (17)$$

If  $\dot{\eta}(0) < 0$ , then

$$\eta(t) \geq \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho}. \quad (18)$$

As time passes,  $\eta$  eventually converges to the upper bound. The bound is attained when individuals choose full exposure in perpetuity without facing any risk of becoming infected. This is the scenario in which getting infected would be most costly as there is

no need to distance and no risk of future infection. The convergence to this bound is intuitive: as time goes to infinity, the disease dies out and obviates the need for distancing.

The above lemma also provides a lower bound on  $\eta$ . This bound applies even if  $\eta$  is locally increasing at time 0. When  $\eta$  is decreasing at the onset, which occurs if  $I_0$  is sufficiently small, the lower bound  $\frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho}$  is approximately tight. This bound corresponds to the cost of infection when individuals are fully distancing from now until eternity.

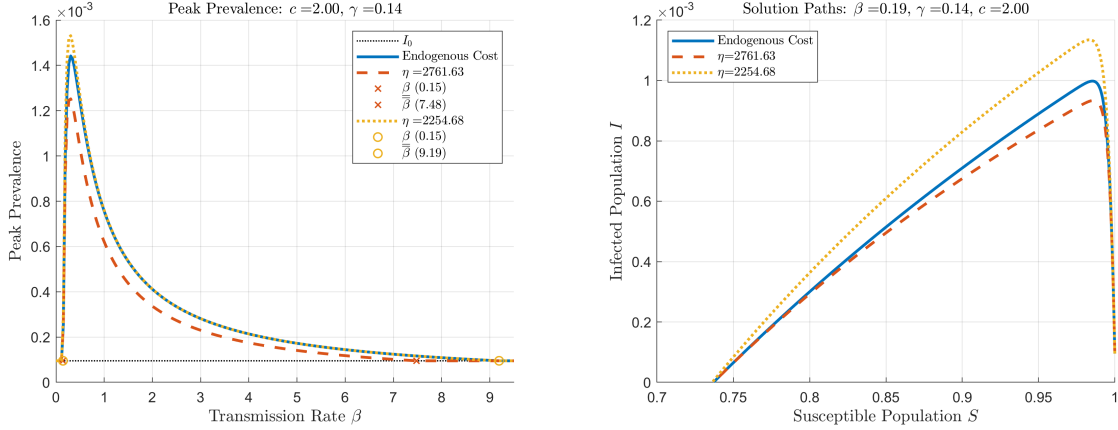


Figure 5: *Peak Prevalence and Solution Path of the Endogenous Cost of Infection Model.* In the left panel, the solid blue curve shows the peak prevalence in the endogenous cost of infection model as a function of the transmission rate. The dashed and dotted curves reproduce the constant cost of infection model's peak prevalence using the derived bounds on  $\eta$ . In the right panel, the solid blue curve represents the solution path of the endogenous cost of infection model. The dashed and dotted curves represent the constant cost of infection model's solution paths using the derived bounds on  $\eta$ .

Lemma 2 connects the solution paths of the model analyzed here and the model with a fixed cost of infection. Towards that, let  $(S, I, R, \varepsilon, p)$  be an equilibrium of the endogenous infection cost model with  $\eta$  being the corresponding co-state given by (15). Let  $\eta_L$  and  $\eta_H$  be the lower and the upper bound on  $\eta$  as given by Lemma 2. Finally, let  $(S_j, I_j, R_j, \varepsilon_j)$ , for  $j \in \{L, H\}$ , be the equilibria of the model with the constant cost of infection corresponding to the lower and upper bounds of  $\eta$ . The following result shows how the model with a constant cost of infection can be used to bound the model with the endogenous cost of infection.

**Proposition 8.** *In the phase space, the graph of  $(S_L, I_L)$  is above that of  $(S, I)$ , which, in turn, is above that of  $(S_H, I_H)$ .*

Finally, we numerically solve the endogenous cost of infection model using commonly

used parameters for COVID-19 following [Farboodi et al. \(2021\)](#) with our objective function and show how peak prevalence varies with the transmission rate; see Figure 5. The non-monotonicity of the peak prevalence in  $\beta$  persists in the environment with the endogenous cost of infection.

## 5 Conclusion

We analyze an epidemiological model with human behavior. We establish that changes in the cost of distancing have a markedly different effect on the progression of a pandemic than changes in the transmission rate. This result has important implications on how public interventions should be modeled. While the existing literature invariably models public interventions as direct changes in the transmission rates, our paper demonstrates that interventions affecting distancing incentives directly—e.g., via bar and restaurant closures—rather than the likelihood of transmission conditional on a meeting—e.g., via mask mandates—should be more appropriately modeled as changes in the cost of distancing. Similarly, changes in the transmission rate should be modeled *together* with a behavioral response rather than only as a direct change in the transmission rate in a standard SIR model because the indirect behavioral effects may outweigh the direct effect on transmission.

Our model is a stylized depiction of reality. Many details require a more thorough investigation. In future work, we plan to study in more detail how the variations in the cost of distancing affect the model’s predictions. Two types of changes are of particular interest. First, distancing becomes more costly over time due to increasing distancing fatigue. This effect can be captured by having the cost of distancing depend on past distancing. Second, sudden and significant increases in distancing cost may occur that alter the course of the epidemic. A salient example is holidays around which families and friends would gather normally. Distancing in such circumstances is much more difficult to sustain as the opportunity cost is high.

Prudence and further study of the spread of infectious diseases are of utmost importance. Not a quarter way into the century, COVID-19 is already the third coronavirus outbreak—after SARS and MERS—, not to mention other outbreaks of diseases like Ebola virus and swine flu, to name a few.

## References

- ACEMOGLU, D., V. CHERNOZHUKOV, I. WERNING, AND M. D. WHINSTON (2021): “Optimally Targeted Lockdowns in a Multi-group SIR Model,” *American Economic Review: Insights*, 3, 487–502.
- ALVAREZ, F., D. ARGENTE, AND F. LIPPI (2021): “A Simple Planning Problem for COVID-19 Lock-down, Testing, and Tracing,” *American Economic Review: Insights*, 3, 367–82.
- ATKESON, A. (2021): “Behavior and the Dynamic of Epidemics,” *Brookings Papers on Economic Activity*, 67–88.
- ATKESON, A. G., K. KOPECKY, AND T. ZHA (2021): “Behavior and the Transmission of COVID-19,” *AEA Papers and Proceedings*, 111, 356–60.
- AVERY, C. (2021): “A Simple Model of Social Distancing and Vaccination,” Working paper.
- AVERY, C., W. BOSSERT, A. CLARK, G. ELLISON, AND S. F. ELLISON (2020): “An economist’s guide to epidemiology models of infectious disease,” *Journal of Economic Perspectives*, 34, 79–104.
- BRAUER, F. (2019): “The Final Size of a Serious Epidemic,” *Bulletin of Mathematical Biology*, 81, 869–877.
- BRAUER, F. AND C. CASTILLO-CHAVEZ (2012): *Mathematical Models in Population Biology and Epidemiology*, Springer, second ed.
- BROTHERHOOD, L., P. KIRCHER, C. SANTOS, AND M. TERTILT (2020): “An Economic Model of the Covid-19 Pandemic with Young and Old Agents: Behavior, Testing and Policies,” Working paper.
- BUDISH, E. (2020): “Maximize Utility subject to  $R \leq 1$ : A Simple Price-Theort Approach to Covid-19 Lockdown and Reopening Policy,” Working paper.
- CALEY, P., D. J. PHILP, AND K. MCCracken (2008): “Quantifying Social Distancing Arising from Pandemic Influenza,” *Journal of the Royal Society Interface*, 5, 631–639.
- CAPASSO, V. AND G. SERIO (1978): “A Generalization of the Kermack-McKnedrick Deterministic Epidemic Model,” *Mathematical Biosciences*, 42, 43–61.
- CHEN, F. (2012): “A Mathematical Analysis of Public Avoidance Behavior during Epidemics Using Game Theory,” *Journal of Theoretical Biology*, 302, 18–28.



- CHERNOZHUKOV, V., H. KASAHARA, AND P. SCHRIMPF (2021): “Causal Impact of Masks, Policies, Behavior on Early Covid-19 Pandemic in the US,” *Journal of Econometrics*, 220, 23–62.
- CHRISTAKIS, N. (2020): “Nicholas Christakis on Fighting Covid-19 by Truly Understanding the Virus,” *The Economist*.
- DASARATHA, K. (2020): “Virus Dynamics with Behavioral Responses,” Working paper.
- DEB, R., M. PAI, A. VOHRA, AND R. VOHRA (2022): “Testing alone is insufficient,” *Review of Economic Design*, 26, 1–21.
- ELSTON, J., C. CARTWRIGHT, P. NDUMBI, AND J. WRIGHT (2017): “The Health Impact of the 2014-15 Ebola Outbreak,” *Public Health*, 143, 60–70.
- ELY, J., A. GALEOTTI, O. JANN, AND J. STEINER (2021): “Optimal test allocation,” *Journal of Economic Theory*, 193, 105236.
- ENGLE, S., J. KEPPO, M. KUDLYAK, E. QUERCIOLO, L. SMITH, AND A. WILSON (2021): “The Behavioral SIR Model, with Applications to the Swine Flu and COVID-19 Pandemics,” Working paper.
- FARBOODI, M., G. JAROSCH, AND R. SHIMER (2021): “Internal and external effects of social distancing in a pandemic,” *Journal of Economic Theory*, 196, 105293.
- FENICHEL, E. P. (2013): “Economic Considerations for Social Distancing and Behavioral Based Policies during an Epidemic,” *Journal of Health Economics*, 32, 440–451.
- FENICHEL, E. P., C. CASTILLO-CHAVEZ, M. G. CEDDIA, G. CHOWELL, P. A. G. PARRA, G. J. HICKLING, G. HOLLOWAY, R. HORAN, B. MORIN, C. PERRINGS, M. SPRINGBORN, L. VELAZQUEZ, AND C. VILLALOBOS (2011): “Adaptive Human Behavior in Epidemiological Models,” *Proceedings of the National Academy of Sciences*, 108, 6306–6311.
- FERGUSON, N. (2007): “Capturing Human Behavior,” *Nature*, 446, 733.
- FERGUSON, N. M., D. LAYDON, G. NEDJATI-GILANI, N. IMAI, K. AINSLIE, M. BAGUELIN, S. BHATIA, A. BOONYASIRI, Z. CUCUNUBÁ, G. CUOMO-DANNENBURG, A. DIGHE, I. DORIGATTI, H. FU, K. GAYTHORPE, W. GREEN, A. HAMLET, W. HINSLEY, L. C. OKELL, S. VAN ELSLAND, H. THOMPSON, R. VERITY, E. VOLZ, H. WANG, Y. WANG, P. G. WALKER, C. WALTERS, P. WINSKILL, C. WHITTAKER, C. A. DONNELLY, S. RILEY, AND A. C. GHANI (2020): “Report 9: Impact of Non-Pharmaceutical Interventions (NPIs) to Reduce COVID19 Mortality and Healthcare Demand,” Imperial college covid-19 response team.

- FUNK, S., M. SALATHÉ, AND V. A. A. JANSEN (2010): “Modelling the Influence of Human Behavior on the Spread of Infectious Diseases: a Review,” *Journal of the Royal Society Interface*, 7, 1247–1256.
- GANS, J. S. (2022): “The Economic Consequences of  $\hat{\mathcal{R}} = 1$ : Towards a Workable Behavioral Epidemiological Model of Pandemics,” *Review of Economic Analysis*, 14, 3–25.
- GEOFFARD, P.-Y. AND T. PHILIPSON (1996): “Rational Epidemics and Their Control,” *International Economic Review*, 37, 603–624.
- GIANNITSAROU, C., S. KISSLER, AND F. TOXVAERD (2021): “Waning Immunity and the Second Wave: Some Projections for SARS-CoV-2,” *American Economic Review: Insights*, 3, 321–38.
- GREENWOOD, J., P. KIRCHER, C. SANTOS, AND M. TERTILT (2019): “An Equilibrium Model of the African HIV/AIDS Epidemic,” *Econometrica*, 87, 1081–1113.
- HALL, R. E., C. I. JONES, AND P. J. KLENOW (2020): “Trading Off Consumption and COVID-19 Deaths,” Working paper.
- HEDLUND, J. (2000): “Risky business: safety regulations, risk compensation, and individual behavior,” *Injury Prevention*, 6, 82–89.
- HEESTERBEEK, J. A. P. AND K. DIETZ (1996): “The Concept of  $R_0$  in Epidemic Theory,” *Statistica Neerlandica*, 50, 89–110.
- HETHCOTE, H. W. (2008): “The Basic Epidemiology Models: Models, Expressions for  $R_0$ , Parameter Estimation, and Applications,” in *Mathematical Understanding of Infectious Disease Dynamics*, ed. by S. Ma and Y. Xia, World Scientific, 1–61.
- HOWARD, J., A. HUANG, Z. LI, Z. TUFEKCI, V. ZDIMAL, H.-M. VAN DER WEST-HUIZEN, A. VON DELFT, A. PRICE, L. FRIDMAN, L.-H. TANG, V. TANG, G. L. WATSON, C. E. BAX, R. SHAIKH, F. QUESTIER, D. HERNANDEZ, L. F. CHU, C. M. RAMIREZ, AND A. W. RIMOIN (2021): “An Evidence Review of Face Masks against COVID-19,” *Proceedings of the National Academy of Sciences*, 118.
- JESTER, B. J., T. M. UYEKI, A. PATEL, L. KOONIN, AND D. B. JERNIGAN (2018): “100 Years of Medical Countermeasures and Pandemic Influenza Preparedness,” *American Journal of Public Health*, 108, 1469–1472.
- KERMACK, W. O. AND A. G. MCKENDRICK (1927): “A Contribution to the Mathematical Theory of Epidemics,” *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 115, 700–721.

- KNOTEK II, E., R. S. SCHOENLE, A. M. DIETRICH, G. F. MÜLLER, K. O. R. MYRSETH, AND M. WEBER (2020): “Consumers and COVID-19: Survey Results on Mask-Wearing Behaviors and Beliefs,” *Economic Commentary*, 2020-20.
- KREMER, M. (1996): “Integrating Behavioral Choice into Epidemiological Models of AIDS,” *Quarterly Journal of Economics*, 111, 549–573.
- KRUSE, T. AND P. STRACK (2020): “Optimal Control of an Epidemic through Social Distancing,” Working paper.
- MCADAMS, D. (2020): “Nash SIR: An Economic-Epidemiological Model of Strategic Behavior During a Viral Epidemic,” Working paper.
- (2021): “The Blossoming of Economic Epidemiology,” *Annual Review of Economics*, 13.
- MCADAMS, D., Y. SONG, AND D. ZOU (2021): “Equilibrium Social Activity during an Epidemic,” Working paper.
- PELTZMAN, S. (1975): “The Effects of Automobile Safety Regulation,” *Journal of Political Economy*, 83, 677–725.
- PHILIPSON, T. J. AND R. A. POSNER (1993): *Private choices and public health: The AIDS epidemic in an economic perspective*, Harvard University Press.
- RACHEL, L. (2020a): “An Analytical Model of Covid-19 Lockdowns: Equilibrium Mitigation Flattens the Epidemic Curve, Optimal Lockdown Does Not,” Working paper.
- (2020b): “Second Wave,” Working paper.
- REED, C., M. BIGGERSTAFF, L. FINELLI, L. M. KOONIN, D. BEAUVAIS, A. UZICANIN, A. PLUMMER, J. BRESEE, S. C. REDD, AND D. B. JERNIGAN (2013): “Novel Framework for Assessing Epidemiologic Effects of Influenza Epidemics and Pandemics,” *Emerging Infectious Diseases*, 19, 85–91.
- RELUGA, T. C. (2010): “Game Theory of Social Distancing in Response to an Epidemic,” *PLoS Computational Biology*, 6.
- ROSS, R. AND H. P. HUDSON (1917): “An Application of the Theory of Probabilities to the Study of a Priori Pathometry.—Part III,” *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 93, 225–240.
- SCHOCH-SPANNA, M. (2001): ““Hospital’s Full-Up”: the 1918 Influenza Pandemic,” *Public Health Reports*, 116, 32–33.

- SNOWDEN, F. M. (2019): *Epidemics and Society: From the Black Death to the Present*, Yale University Press.
- TOXVAERD, F. (2020): “Equilibrium Social Distancing,” Working paper.
- TOXVAERD, F. AND R. ROWTHORN (2020): “On the Management of Population Immunity,” Working paper.
- VERELST, F., L. WILLEM, AND P. BEUTELS (2016): “Behavioural Change Models for Infectious Disease Transmission: a Systematic Review (2010-2015),” *Journal of the Royal Society Interface*, 13, 20160820.
- WALTER, W. (1998): *Ordinary Differential Equations*, Springer.
- WEITZ, J. S., S. W. PARK, C. EKSIN, AND J. DUSHOFF (2020): “Awareness-driven behavior changes can shift the shape of epidemics away from peaks and toward plateaus, shoulders, and oscillations,” *Proceedings of the National Academy of Sciences*, 117, 32764–32771.
- YAN, Y., J. BAYHAM, A. RICHTER, AND E. P. FENICHEL (2021): “Risk compensation and face mask mandates during the COVID-19 pandemic,” *Scientific Reports*, 11, 3174.

## A Appendix

**Proof of Proposition 1.** An individual’s problem (1) is concave; therefore, the first-order condition (5) is also sufficient. This pins down the individual’s optimal distancing in the SIR dynamics.

Using the exposure obtained from (5) in the SIR dynamics yields

$$\dot{S}(t) = -\beta S(t)I(t) \max \left( 1 - \frac{\eta \beta I(t)}{c}, 0 \right), \quad (19)$$

$$\dot{I}(t) = \beta S(t)I(t) \max \left( 1 - \frac{\eta \beta I(t)}{c}, 0 \right) - \gamma I(t), \quad (20)$$

$$\dot{R}(t) = \gamma I(t). \quad (21)$$

Thus, in any equilibrium  $(S, I, R)$  is characterized by the system of differential equations  $\frac{d}{dt}(S, I, R) = F(t, S, I, R)$ , where  $F$  is defined by (19), (20), and (21). The initial condition is  $(S(0), I(0), R(0)) = (S_0, I_0, 0)$ . Then, the initial value problem admits

a unique solution  $(S, I, R)$  on  $[0, \infty)$ , as the system satisfies the standard conditions. Namely, the function  $F$  is continuous on the domain  $D = [0, \infty) \times [0, 1]^3$ , and  $F$  is uniformly Lipschitz continuous in  $(S, I, R)$ : there exists a Lipschitz constant  $L$  satisfying  $\|F(t, S, I, R) - F(t, \tilde{S}, \tilde{I}, \tilde{R})\| \leq L\|(S, I, R) - (\tilde{S}, \tilde{I}, \tilde{R})\|$  for each  $t \in [0, \infty)$ . See, for example, [Walter \(1998\)](#). Now,  $\varepsilon = \varepsilon_i$  is uniquely determined, and hence the model admits a unique and symmetric equilibrium.  $\square$

**Proof of Proposition 2.** Let  $\hat{t}$  be as in the supposition of the proposition. We first show  $\varepsilon(\hat{t}) \in (0, 1)$ . Since  $I(t) > 0$  for all  $t$ , evaluating  $\dot{I}(t) = 0$  at  $\hat{t}$  yields  $\beta S(\hat{t})\varepsilon(\hat{t}) = \gamma$ . Hence,  $\varepsilon(\hat{t}) \in (0, 1)$ .

Next, since  $\dot{I}$  is differentiable at  $\hat{t}$ , it follows that  $\ddot{I}(\hat{t})$  exists. We show:

$$\begin{aligned}\ddot{I}(\hat{t}) &= \beta \left( \dot{S}(\hat{t})I(\hat{t})\varepsilon(\hat{t}) + S(\hat{t})\dot{I}(\hat{t})\varepsilon(\hat{t}) + S(\hat{t})I(\hat{t})\dot{\varepsilon}(\hat{t}) \right) - \gamma\dot{I}(\hat{t}) \\ &= \beta I(\hat{t}) \left( \dot{S}(\hat{t})\varepsilon(\hat{t}) + S(\hat{t})\dot{\varepsilon}(\hat{t}) \right) \\ &= \beta S(\hat{t})I(\hat{t}) \left( -\beta I(\hat{t})\varepsilon^2(\hat{t}) + \dot{\varepsilon}(\hat{t}) \right) \\ &= -\beta S(\hat{t})I^2(\hat{t})\varepsilon^2(\hat{t}) < 0.\end{aligned}$$

The second equality follows from  $\dot{I}(\hat{t}) = 0$ , the third from equation (2), and the fourth from

$$\dot{\varepsilon}(\hat{t}) = -\frac{\eta\beta}{c}\dot{I}(\hat{t}) = 0, \quad (22)$$

which, in turn, follows from optimality condition (5) and  $\varepsilon(\hat{t}) \in (0, 1)$ .  $\square$

**Proof of Proposition 3.** The first statement follows from the equality  $\dot{I}(0) = \frac{I_0}{\gamma} (R_0^b - 1)$ .

Part (i): From (3) it follows that

$$\dot{I}(0) > 0 \text{ if and only if } I_0 \left( \beta \left( 1 - \frac{\beta\eta}{c}I_0 \right) (1 - I_0) - \gamma \right) > 0.$$

Therefore,  $\dot{I}(0) > 0$  if and only if  $\beta \in (\underline{\beta}, \bar{\beta})$  where  $\underline{\beta}$  and  $\bar{\beta}$  are solutions to the quadratic equation

$$\beta \left( 1 - \frac{\eta I_0}{c}\beta \right) (1 - I_0) - \gamma = 0. \quad (23)$$

The discriminant of the quadratic equation is positive if and only if  $I_0 < \frac{1}{1 + \frac{4\eta\gamma}{c}}$ ; the solid curve in the left panel of Figure 1 corresponds to equation (23). Since  $\varepsilon(0) = 1 - \frac{\eta I_0}{c}\beta < 1$

and  $I_0 > 0$ , the left-hand side of the above equation is negative at  $\beta = \frac{\gamma}{1-I_0}$ . Thus,  $\underline{\beta} > \frac{\gamma}{1-I_0}$ . If  $\beta = \frac{c}{\eta I_0}$ , then  $\varepsilon(0) = 0$  and  $\dot{I}(0) < 0$ . Thus,  $\bar{\beta} < \frac{c}{\eta I_0}$ .

Part (ii): Let  $I_0 \geq \frac{1}{1+\frac{4\eta\gamma}{c}}$ . If  $\varepsilon(0) > 0$ , then the quadratic equation (23) has at most one solution. Thus,  $\dot{I}(0) \leq 0$ . Proposition 2 then implies that if  $\dot{I}(t) \leq 0$  for some  $t$  (take  $t = 0$ ) then  $\dot{I}(s) < 0$  for all  $s > t$ . If  $\varepsilon(0) = 0$ , then there exists a  $\underline{t}$  such that  $\varepsilon(t) = 0$  on  $t \in [0, \underline{t}]$  and  $\varepsilon(t) > 0$  on  $t \in (\underline{t}, \infty)$ . By implication  $\dot{I}(t) = -\gamma I(t) < 0$  for  $t \in [0, \underline{t}]$ , while  $\varepsilon(t) > 0$  for  $t \in (\underline{t}, \infty)$  by the same reasoning as in the  $\varepsilon(0) > 0$  case. Finally,  $\dot{I}(\underline{t}) = -\gamma I(\underline{t})$ .  $\square$

**Proof of Proposition 4.** By Assumption  $d(0) < 1$ , it can be seen that  $\varepsilon(t) \in (0, 1)$  for all  $t$ . Then, we have

$$\frac{dS}{d(S+I)} = \frac{\beta}{\gamma} S \varepsilon = \frac{\beta}{\gamma} \left( \frac{\beta\eta}{c} S^2 + S - \frac{\beta\eta}{c} (S+I)S \right),$$

where the first equality follows from dividing (2) by the sum of (2) and (3), and the second uses (5) and simple manipulations. The above expression can be rewritten as

$$\frac{d}{d(S+I)} \left( \frac{1}{S} \right) + \left( \frac{\beta}{\gamma} - \frac{\beta^2\eta}{\gamma c} (S+I) \right) \frac{1}{S} = -\frac{\beta^2\eta}{\gamma c},$$

which is a linear first-order differential equation with respect to  $\frac{1}{S}$  and  $(S+I)$ . For ease of notation, let  $y = \frac{1}{S}$  and  $x = S+I$ . Then,

$$\frac{dy}{dx} + \left( \frac{\beta}{\gamma} - \frac{\beta^2\eta}{\gamma c} x \right) y = -\frac{\beta^2\eta}{\gamma c}. \quad (24)$$

Let  $\mu(x) := \exp \left( \int \left( \frac{\beta}{\gamma} - \frac{\beta^2\eta}{\gamma c} x \right) dx \right)$  be the integrating factor. We have

$$\mu(x) = k \cdot \exp \left( -\frac{\beta^2\eta}{2\gamma c} \left( x - \frac{c}{\beta\eta} \right)^2 \right), \quad (25)$$

where  $k$  is the constant of integration. Then, equation (24) reduces to

$$\frac{d}{dx} [\mu(x)y] = \mu(x) \left[ \frac{d}{dx} y + \left( \frac{\beta}{\gamma} - \frac{\beta^2\eta}{\gamma c} x \right) y \right] = -\mu(x) \frac{\beta^2\eta}{\gamma c}. \quad (26)$$

Integrating the outer most sides of Expression (26) and using (25) yield

$$\left[ \exp \left( -\frac{\beta^2\eta}{2\gamma c} \left( x - \frac{c}{\beta\eta} \right)^2 \right) y \right]_{S+I}^1 = \frac{-\beta^2\eta}{\gamma c} \int_{S+I}^1 \exp \left( -\frac{\beta^2\eta}{2\gamma c} \left( x - \frac{c}{\beta\eta} \right)^2 \right) dx. \quad (27)$$

The left-hand side of (27) reduces to

$$\exp\left(-\frac{\beta^2\eta}{2\gamma c}\left(1-\frac{c}{\beta\eta}\right)^2\right)\frac{1}{S_0}-\exp\left(-\frac{\beta^2\eta}{2\gamma c}\left(S+I-\frac{c}{\beta\eta}\right)^2\right)\frac{1}{S}.$$

For the right-hand side of (27), let  $v = \beta\sqrt{\frac{\eta}{2\gamma c}}\left(x - \frac{c}{\beta\eta}\right)$ . Since  $\frac{dv}{dx} = \beta\sqrt{\frac{\eta}{2\gamma c}}$ , the right-hand side of (27) reduces to

$$-\beta\sqrt{\frac{2\eta}{\gamma c}}\int_{\beta\sqrt{\frac{\eta}{2\gamma c}}\left(S+I-\frac{c}{\beta\eta}\right)}^{\beta\sqrt{\frac{\eta}{2\gamma c}}\left(1-\frac{c}{\beta\eta}\right)}e^{-v^2}dv.$$

Hence, we can rewrite equation (27) as

$$\exp\left(-\frac{\beta^2\eta}{2\gamma c}\left(S+I-\frac{c}{\beta\eta}\right)^2\right)\frac{1}{S}=\exp\left(-\frac{\beta^2\eta}{2\gamma c}\left(1-\frac{c}{\beta\eta}\right)^2\right)\frac{1}{S_0}+\beta\sqrt{\frac{2\eta}{\gamma c}}\int_{\beta\sqrt{\frac{\eta}{2\gamma c}}\left(S+I-\frac{c}{\beta\eta}\right)}^{\beta\sqrt{\frac{\eta}{2\gamma c}}\left(1-\frac{c}{\beta\eta}\right)}e^{-v^2}dv,$$

and finally we obtain (8), as desired.  $\square$

**Proof of Proposition 5.** We prove the assertion with respect to  $\frac{\eta}{c}$ . Denote by  $(S, I(S))$  a point on the solution path. Differentiating the quotient differential equation  $\frac{dI}{dS}$  with respect to  $\frac{\eta}{c}$  at a fixed point  $(S, I(S))$  yields

$$\frac{\partial}{\partial \frac{\eta}{c}} \frac{dI}{dS} = \frac{\gamma I}{S} \frac{1}{(1 - \beta I \frac{\eta}{c})^2} > 0.$$

Now, take  $\frac{\eta}{c}$  and  $\tilde{\eta}$  with  $\frac{\eta}{c} > \tilde{\eta}$ . Denote by  $(\tilde{S}, \tilde{I})$  the solution path associated with  $\tilde{\eta}$  and  $\tilde{c}$ . By the above inequality  $\frac{dI}{dS} > \frac{d\tilde{I}}{d\tilde{S}}$  at any point of intersection. It follows that there exists a  $\delta_1 > 0$  such that  $\tilde{I}(S_0 - \delta) > I(S_0 - \delta)$  for every  $\delta < \delta_1$ . Now, it is sufficient to show that two curves  $I$  and  $\tilde{I}$  do not intersect at any other point. Suppose to the contrary that  $I$  and  $\tilde{I}$  did intersect. Let,  $\bar{S} := \sup\{S \in (0, S_0 - \delta_1] \mid \tilde{I}(S) = I(S)\}$ . Since the solution curves are continuous they intersect at  $\bar{S}$  and therefore  $\frac{d\tilde{I}}{d\tilde{S}}(\bar{S}) < \frac{dI}{dS}(\bar{S})$ . But now we have that  $\frac{d\tilde{I}}{d\tilde{S}} < \frac{dI}{dS}$  both at  $\bar{S}$  and  $S_0$  and that the two curves do not intersect anywhere in between, a contradiction.  $\square$

**Proof of Proposition 6.** We prove the result with respect to  $c$  first, then with respect to  $\beta$ .

Part (ii): The proof of Proposition 3 has established that

$$\dot{I}(0) > 0 \text{ if and only if } I_0 \left( \beta \left( 1 - \frac{\beta\eta}{c} I_0 \right) (1 - I_0) - \gamma \right) > 0.$$

Therefore,  $\dot{I}(0) > 0$  if and only if  $c > \underline{c}$ , where  $\underline{c}$  can be recovered from the above

inequality. The peak prevalence when  $c \leq \bar{c}$  is  $I_0$ . If  $c > \bar{c}$ , then the peak prevalence is strictly greater than  $I_0$ . We show that the peak prevalence is strictly increasing in  $c > \bar{c}$ . Differentiating equation (7) with respect to  $c$ , while holding  $S$  and  $I$  fixed, yields<sup>30</sup>

$$\frac{\partial}{\partial c} \left( \frac{dI}{dS} \right) = -\frac{\gamma\eta I}{c^2 S (1 - \beta I_c^\eta)^2} < 0.$$

If two solution paths corresponding to  $c$  and  $c' > c$  intersect at some point, the solution path corresponding to  $c'$  has a smaller slope. A certain point of intersection is the beginning of the infection  $(S_0, I_0)$ . At this point in the graph with  $S$  on the horizontal and  $I$  on the vertical axis the solution path corresponding to  $c'$  is steeper; the solution paths are decreasing at  $(S_0, I_0)$ . Just below  $S_0$ , then, the solution path corresponding to  $c'$  is above the one corresponding to  $c$ . If they were to intersect at some other  $S < S_0$ , the solution path corresponding to  $c'$  would have to intersect the solution path corresponding to  $c$  from above and stay below it. This would contradict the finding that the solution path corresponding to  $c'$  is above the one corresponding to  $c$  for  $S$  slightly below  $S_0$ . Finally, given that the solution path under  $c'$  is above the solution path under  $c$ , the peak of infection under  $c'$  must be higher than the peak of infection under  $c$ .

Part (i): We break up the proof for  $\beta$  into two steps.

Step 1:  $I^*$  is decreasing in  $\beta$  for  $\beta \in [\frac{c}{2I_0\eta}, \bar{\beta}]$ . The derivative of the quotient differential equation (7) at a given point  $(S, I(S))$  with respect to  $\beta$  is

$$\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right) = -\frac{\gamma}{\beta^2 S} \frac{1 - 2\frac{\beta\eta}{c} I(S)}{(1 - \frac{\beta\eta}{c} I(S))^2}. \quad (28)$$

The above derivative evaluated at  $(S_0, I_0)$  is greater or equal to 0, for  $\beta \geq \frac{c}{2I_0\eta}$ . This means that at  $(S_0, I_0)$ , a higher  $\beta$  leads to a slower spread of the infection when the starting  $\beta$  is high enough. At  $(S_0, I_0)$  solution paths are decreasing, thus the positive derivative with respect to  $\beta$  means that the solution path becomes flatter as  $\beta$  increases. That is, around  $(S_0, I_0)$  the solution path corresponding to a higher  $\beta$  is, therefore, below the one with the lower  $\beta$ .

Moreover,  $\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right) \geq 0$  at  $(S_0, I_0)$ , for  $\beta \geq \frac{c}{2I_0\eta}$ , implies that the same is true for all  $(S, I)$  with  $I > I_0$ . This means that if two solution paths corresponding to some  $\beta$  and  $\beta' > \beta$  in  $[\frac{c}{2I_0\eta}, \bar{\beta}]$  intersect, then the solution path corresponding to  $\beta'$  must have a larger slope. One such point of intersection is  $(S_0, I_0)$ . Therefore, a solution path for  $\beta'$  is below the one of  $\beta$  just below  $S_0$  and it cannot intersect it anymore as long as  $I \geq I_0$ . In other words, the solution path of  $\beta'$  is strictly below the solution path of  $\beta$  for all  $I > I_0$ .

---

<sup>30</sup>Recall that when  $\dot{I}(0) > 0$ , exposure is interior for all  $t$ .



The maximum of  $I$  for  $\beta'$  is, therefore, strictly below the maximum of  $I$  for  $\beta$ .

Step 2: There exists a  $\beta_1$  such that  $I^*$  is increasing in  $\beta$  on  $(\underline{\beta}, \beta_1)$ . We divide this step into three substeps. First, we show that the peak  $I^*$  is continuous in  $\beta$ . Then, we show that  $\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right) < 0$  along the entire solution path. Finally, we combine these two insights to show that for  $\beta > \underline{\beta}$  but sufficiently close,  $\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right) < 0$  implying that the peak is increasing in  $\beta$  for  $\beta \in (\underline{\beta}, \underline{\beta} + \delta)$  for some  $\delta > 0$ .

Step 2.1: We argue that  $I^*$  is continuous in  $\beta \in (0, \bar{\beta})$ . For  $\beta \in (0, \underline{\beta})$ ,  $I^* = I_0$ .

For  $\beta \in (\underline{\beta}, \bar{\beta})$ , in the  $(S, I)$ -phase space,  $(S, I) = (S^*, I^*)$  satisfies equation (8) and  $\frac{dI}{dS} = 0$ , i.e.,  $S^* = \frac{\gamma}{\beta} \frac{1}{1 - \frac{\beta \eta}{c} I^*}$ . Substituting the latter equation into the former and rearranging, we obtain

$$\begin{aligned} & \exp \left( -\frac{\eta}{2\gamma c} \left( \beta - \frac{c}{\eta} \right)^2 \right) \frac{1}{S_0} + \beta \sqrt{\frac{2\eta}{\gamma c}} \int_{\sqrt{\frac{\eta}{2\gamma c} \left( \frac{\gamma}{1 - \frac{\beta \eta}{c} I^*} + \beta I^* - \frac{c}{\eta} \right)}}^{\sqrt{\frac{\eta}{2\gamma c} \left( \beta - \frac{c}{\eta} \right)}} e^{-v^2} dv \\ &= \frac{1}{\gamma} \left( \beta - \frac{\beta^2 \eta}{c} I^* \right) \exp \left( -\frac{\eta}{2\gamma c} \left( \frac{\gamma}{1 - \frac{\beta \eta}{c} I^*} + \beta I^* - \frac{c}{\eta} \right)^2 \right). \end{aligned}$$

This implies that  $I^*$  is differentiable and thus continuous in  $\beta$  for  $\beta \in (\underline{\beta}, \bar{\beta})$ .

Finally, if  $\beta = \underline{\beta}$  then  $(S^*, I^*) = (S_0, I_0)$  satisfies the above implicit equation. Thus,  $I^*$  is also continuous at  $\beta = \underline{\beta}$ .

Step 2.2:  $\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right) < 0$ , along the entire solution path whenever it holds along the path that  $1 - 2\frac{\beta \eta}{c} I(S) > 0$ . This is satisfied for  $\beta = \underline{\beta}$ .

Recall that the sign of  $\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right)$  at each  $(S, I(S))$  is determined by the negative of the sign of  $1 - 2\frac{\beta \eta}{c} I(S)$ . Thus, it is sufficient for the derivative to be negative along the entire path that  $1 - 2\frac{\beta \eta}{c} I^* > 0$ .

Observe that the solution  $\underline{\beta}$  of equation (23) is given by

$$\underline{\beta} = \frac{c}{2\eta I_0} \left( 1 - \left( 1 - 4\frac{\eta \gamma}{c} \frac{I_0}{S_0} \right)^{\frac{1}{2}} \right).$$

Therefore,

$$1 - 2\frac{\beta \eta}{c} I_0 = \left( 1 - 4\frac{\eta \gamma}{c} \frac{I_0}{S_0} \right)^{\frac{1}{2}} > 0,$$

where the inequality follows due to the assumption on  $I_0$  in the statement of the result.

Consequently,

$$\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\underline{\beta}) \right) < 0.$$

Step 2.3:  $I^*$  is increasing in  $\beta$  for  $\beta \in (\underline{\beta}, \underline{\beta} + \delta)$  for some  $\delta > 0$ .

Since  $1 - 2\frac{\beta\eta}{c}I_0 > 0$ , there exists a  $\delta_1 > 0$  such that  $1 - 2\frac{\beta\eta}{c}I_0 > 0$  for all  $\beta \in [\underline{\beta}, \underline{\beta} + \delta_1)$ .

By continuity of the peak, for every  $\delta_2 > 0$ , there exists a  $\delta_3 > 0$ , such that  $\beta \in [\underline{\beta}, \underline{\beta} + \delta_3)$  implies  $I^*(\beta) < I_0 + \delta_2$ . Choose  $\delta_2$  to correspond to the  $\delta_1$  argued above Step 2.2, and let  $\delta_3$  corresponds to such  $\delta_2$ . This guarantees that we consider  $\beta$  to lie in a range such that the peak is sufficiently low to ensure that the slopes of the solution paths can be ordered by comparing  $\beta$ .

By Steps 2.1 and 2.2, for any such  $\beta$ ,  $1 - 2\frac{\beta\eta}{c}I(S) > 0$  and therefore  $\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right) < 0$ . This implies that whenever two solution paths corresponding to different  $\beta$  in  $(\underline{\beta}, \underline{\beta} + \delta_3)$  intersect at a point, the one with the higher  $\beta$  has the smaller slope. Indeed, one point of intersection is  $(S_0, I_0)$ . The solution path with a higher  $\beta$  must be steeper than the other path; the two are decreasing at  $(S_0, I_0)$ . Suppose the two solution paths were to intersect at some  $S < S_0$  and let  $\tilde{S}$  be the largest such  $S$ . Then due to  $\frac{\partial}{\partial \beta} \left( \frac{dI}{dS}(\beta) \right) < 0$ , the solution path with the higher  $\beta$  would need to intersect the solution path with the lower  $\beta$  from above, and fall below it for  $S > \tilde{S}$ . But this contradicts the fact that at  $S_0$  the solution path corresponding to a higher  $\beta$  is above the one with the lower  $\beta$ .  $\square$

**Proof of Proposition 7.** We begin by proving a lemma establishing bounds for  $S_\infty$ .

**Lemma 3.** *The following chain of inequalities holds:*

$$0 < S_0 e^{-\frac{\beta}{\gamma}} \leq \hat{S}_\infty \leq S_\infty < \frac{\gamma}{\beta}.$$

**Proof of Lemma 3.** We start with showing  $S_0 e^{-\frac{\beta}{\gamma}} \leq \hat{S}_\infty$ . Let  $\hat{R}_\infty := \lim_{t \rightarrow \infty} \hat{R}(t)$ . It follows from the SIR dynamics that

$$\hat{S}_\infty = S_0 \exp \left( -\beta \int_0^\infty \hat{I}(s) ds \right) = S_0 \exp \left( -\frac{\beta}{\gamma} \hat{R}_\infty \right) \geq S_0 e^{-\frac{\beta}{\gamma}}.$$

The first equality follows from integrating both sides of (2) with  $S$  and  $I$  replaced by  $\hat{S}$  and  $\hat{I}$ , respectively, and with  $\varepsilon = 1$ . The second equality follows from integrating (4) (precisely, with  $R$  and  $I$  replaced by  $\hat{R}$  and  $\hat{I}$ , respectively). The inequality follows

because  $\hat{R}_\infty \leq 1$ .

Second, we show  $\hat{S}_\infty \leq S_\infty$ . It suffices to show that  $\hat{S}(t) \leq S(t)$ , as letting  $t \rightarrow \infty$  yields the desired result. Suppose to the contrary that there exists some  $\tilde{t}$  such that  $S(\tilde{t}) < \hat{S}(\tilde{t})$ . At time 0,  $S(0) = \hat{S}(0)$  and  $\dot{S}(0) > \dot{\hat{S}}(0)$ . Thus, there exists an interval in which  $S(\cdot) > \hat{S}(\cdot)$ . Then there would have to exist  $t_0$  such that  $S(t_0) = \hat{S}(t_0)$  and  $\dot{S}(t_0) < \dot{\hat{S}}(t_0)$ . However, it follows from  $S(t_0) = \hat{S}(t_0)$  and the previous argument that  $I(t_0) \leq \hat{I}(t_0)$ , and thus

$$\dot{S}(t_0) = -\beta\varepsilon(t_0)S(t_0)I(t_0) > -\beta S(t_0)I(t_0) > -\beta\hat{S}(t_0)\hat{I}(t_0) = \dot{\hat{S}}(t_0),$$

which is impossible.

Third, we show  $S_\infty < \frac{\gamma}{\beta}$  in two steps. The first step establishes  $S_\infty \leq \frac{\gamma}{\beta}$ . Suppose not. As  $S(t)$  is weakly decreasing throughout, there exists a  $\delta > 0$  such that  $S(t) \geq \delta + \frac{\gamma}{\beta}$  for all  $t \geq 0$ . Since  $\lim_{t \rightarrow \infty} \varepsilon(t) = 1$  and  $\delta > 0$ , for a given  $\kappa \in (0, \delta)$ , there exists  $t_1 \in [t_0, \infty)$  such that  $\delta\varepsilon(t) - \frac{\gamma}{\beta}(1 - \varepsilon(t)) > \kappa$  for all  $t \geq t_1$ . Then, for all  $t \geq t_1$ , we have

$$\dot{I}(t) = \beta I(t)(S(t)\varepsilon(t) - \frac{\gamma}{\beta}) \geq \beta I(t)((\delta + \frac{\gamma}{\beta})\varepsilon(t) - \frac{\gamma}{\beta}) > \beta I(t)\kappa,$$

that is,  $\frac{\dot{I}(t)}{I(t)} > \beta\kappa$  (note that, since  $\dot{I}(t) \geq -\gamma I(t)$ ,  $I(t)$  is always positive:  $I(t) \geq I(0)e^{-\gamma t} > 0$ ). Thus,  $I(t) \geq I(t_1)e^{\beta\kappa t}$ , which yields  $I_\infty = +\infty$ . This is a contradiction to  $I_\infty = 0$ .

The second step establishes  $S_\infty \neq \frac{\gamma}{\beta}$ . Suppose to the contrary  $S_\infty = \frac{\gamma}{\beta}$ . Then,  $\frac{dI}{dS}(S_\infty) = -1 + \frac{\gamma}{\beta} \frac{1}{S_\infty} = 0$  as  $\lim_{t \rightarrow \infty} \varepsilon(t) = 1$ . However, note that

$$\begin{aligned} \frac{d}{dS} \frac{dI}{dS}(S_\infty) &= -\frac{\gamma}{\beta} \frac{1}{\varepsilon(I(S))S} \left( \frac{1}{S} + \frac{1}{\varepsilon(I(S))} \frac{d\varepsilon(I(S))}{dI(S)} \frac{dI}{dS} \right) \\ &= -\frac{\gamma}{\beta} \frac{1}{\varepsilon(I(S))S^2} < 0 \end{aligned}$$

as  $\frac{dI}{dS}(S_\infty) = 0$ , where  $\varepsilon(I(S)) = 1 - \beta_c^{\frac{\eta}{c}} I(S)$ . Thus, there is a  $\delta > 0$  such that for  $S \in (S_\infty, S_\infty + \delta)$ ,  $\frac{dI}{dS}(S_\infty + \delta) < 0$  and, hence, that  $I(S_\infty + \delta) < 0$ , a contradiction. Thus,  $S_\infty < \frac{\gamma}{\beta}$ .  $\square$

We now move on to proving the comparative statics of  $S_\infty$ . It follows from Proposition 5 that  $S_\infty$  is increasing in  $c$ . Thus, we show that  $S_\infty$  is decreasing in  $\beta$  for the following three cases: (1)  $\beta \in [0, \underline{\beta}]$ ; (2)  $\beta \in [\underline{\beta}, \bar{\beta}]$ ; and (3)  $\beta \in [\bar{\beta}, \frac{c}{\eta I_0}]$ .

Case 1. Let  $\beta \in [0, \underline{\beta}]$ . In this case,  $I^* = I_0$ , and  $\dot{I}(t) < 0$  for all  $t \in (0, \infty)$ . Also,  $\varepsilon(\cdot) \in (0, 1)$  as  $\varepsilon(t) = 1 - \beta \frac{\eta}{c} I(t)$  is decreasing in  $I(t)$ ,  $I(t) > 0$  is decreasing, and  $\beta < \frac{c}{\eta I(0)}$ .

The derivative of the quotient differential equation with respect to  $\beta$  at  $(S, I(S))$  is

$$\frac{\partial}{\partial \beta} \frac{dI}{dS} = -\frac{\gamma}{\beta S} \frac{1}{\beta(1 - \frac{\beta \eta}{c} I(S))} \left( 1 - 2 \frac{\beta \eta}{c} I(S) \right) < 0.$$

This implies that, for any  $\beta, \beta' \in [0, \underline{\beta}]$  with  $\beta < \beta'$ , the solution path associated with  $\beta'$  has a flatter slope than the one associated with  $\beta$  at any point  $S \in (S_0, S_\infty(\beta))$ , where  $S_\infty(\beta)$  is  $S_\infty$  associated with  $\beta$ . Thus,  $I(S_\infty(\beta)) > 0$  for the solution path associated with  $\beta'$ , and hence  $S_\infty(\beta') < S_\infty(\beta)$ .

Case 2. Let  $\beta \in [\underline{\beta}, \bar{\beta}]$ . In this case,  $\dot{I}(0) \geq 0$  and  $\varepsilon(\cdot) \in (0, 1)$ . Substituting  $(S, I) = (S_\infty, 0)$  into (8) yields

$$S_\infty = \frac{\exp\left(-\frac{\eta}{2\gamma c} \left(\beta S_\infty - \frac{c}{\eta}\right)^2\right)}{\exp\left(-\frac{\eta}{2\gamma c} \left(\beta - \frac{c}{\eta}\right)^2\right) \frac{1}{S_0} + 2\beta \sqrt{\frac{\eta}{2\gamma c}} \int_{\sqrt{\frac{\eta}{2\gamma c}(\beta S_\infty - \frac{c}{\eta})}}^{\sqrt{\frac{\eta}{2\gamma c}(\beta - \frac{c}{\eta})}} e^{-v^2} dv}. \quad (29)$$

Rewriting Expression (29),

$$\exp\left(-\frac{\eta}{2\gamma c} \left(\beta - \frac{c}{\eta}\right)^2\right) \frac{S_\infty}{S_0} + 2\beta S_\infty \sqrt{\frac{\eta}{2\gamma c}} \int_{\sqrt{\frac{\eta}{2\gamma c}(\beta S_\infty - \frac{c}{\eta})}}^{\sqrt{\frac{\eta}{2\gamma c}(\beta - \frac{c}{\eta})}} e^{-v^2} dv = \exp\left(-\frac{\eta}{2\gamma c} \left(\beta S_\infty - \frac{c}{\eta}\right)^2\right). \quad (30)$$

For the right-hand side,

$$\frac{\partial}{\partial \beta}(\text{RHS}) = - \underbrace{\exp\left(-\frac{\eta}{2\gamma c} \left(\beta S_\infty - \frac{c}{\eta}\right)^2\right)}_{=(\text{RHS})} \frac{\eta}{\gamma c} \left(\beta S_\infty - \frac{c}{\eta}\right) \left(S_\infty + \beta \frac{\partial S_\infty}{\partial \beta}\right).$$

For the left-hand side, we obtain

$$\begin{aligned} \frac{\partial}{\partial \beta}(\text{LHS}) &= \exp\left(-\frac{\eta}{2\gamma c} \left(\beta - \frac{c}{\eta}\right)^2\right) \frac{S_\infty}{S_0} \left(-\frac{\eta}{\gamma c} \beta(1 - S_0) + \frac{1}{\gamma} + \frac{\partial S_\infty}{\partial \beta} \frac{1}{S_\infty}\right) \\ &\quad + 2\beta S_\infty \sqrt{\frac{\eta}{2\gamma c}} \int_{\sqrt{\frac{\eta}{2\gamma c}(\beta S_\infty - \frac{c}{\eta})}}^{\sqrt{\frac{\eta}{2\gamma c}(\beta - \frac{c}{\eta})}} e^{-v^2} dv \left(\frac{1}{\beta} + \frac{\partial S_\infty}{\partial \beta} \frac{1}{S_\infty}\right) \\ &\quad - \beta S_\infty \frac{\eta}{\gamma c} \left(S_\infty + \beta \frac{\partial S_\infty}{\partial \beta}\right) \exp\left(-\frac{\eta}{2\gamma c} \left(\beta S_\infty - \frac{c}{\eta}\right)^2\right). \end{aligned}$$

Equating the derivatives of the left-hand and right-hand sides and using Expression (30) and rearranging yield

$$\begin{aligned} & \exp\left(-\frac{\eta}{2\gamma c}\left(\beta S_\infty - \frac{c}{\eta}\right)^2\right) \left(\left(\frac{\beta}{\gamma} - \frac{1}{S_\infty}\right) \frac{\partial S_\infty}{\partial \beta} + \frac{1}{\gamma} \left(S_\infty - \frac{\gamma}{\beta}\right)\right) \\ &= \exp\left(-\frac{\eta}{2\gamma c}\left(\beta - \frac{c}{\eta}\right)^2\right) \frac{S_\infty}{S_0} \left(\frac{1}{\gamma} \left(1 - \frac{\eta\beta(1-S_0)}{c} - \frac{\gamma}{\beta}\right)\right). \end{aligned}$$

Thus,

$$\left(\frac{\beta}{\gamma} - \frac{1}{S_\infty}\right) \frac{\partial S_\infty}{\partial \beta} = \frac{\exp\left(-\frac{\eta}{2\gamma c}\left(\beta - \frac{c}{\eta}\right)^2\right)}{\exp\left(-\frac{\eta}{2\gamma c}\left(\beta S_\infty - \frac{c}{\eta}\right)^2\right)} \frac{S_\infty}{S_0} \frac{1}{\gamma} \left(\varepsilon(0) - \frac{\gamma}{\beta}\right) - \frac{1}{\gamma} \left(S_\infty - \frac{\gamma}{\beta}\right). \quad (31)$$

Since  $S_\infty < \frac{\gamma}{\beta}$  follows from Lemma 3,<sup>31</sup> it follows that

$$\frac{\partial S_\infty}{\partial \beta} = \frac{S_\infty}{\beta} \left( \frac{\exp\left(-\frac{\eta}{2\gamma c}\left(\beta - \frac{c}{\eta}\right)^2\right)}{\exp\left(-\frac{\eta}{2\gamma c}\left(\beta S_\infty - \frac{c}{\eta}\right)^2\right)} \frac{S_\infty}{S_0} \frac{\varepsilon(0) - \frac{\gamma}{\beta}}{S_\infty - \frac{\gamma}{\beta}} - 1 \right) < 0.$$

Case 3. The case with  $\beta \in [\bar{\beta}, \frac{c}{\eta_0}]$  is analogous to Case 1, and thus the proof is omitted.  $\square$

**Remark 1.** First, we derive equation (10). Suppose that an individual gets infected at time  $\tau$ . The (conditional) probability that the individual will have been recovered after time  $\tau + t$  is  $1 - e^{-\gamma t}$ . Therefore,

$$V_I(\tau) = \int_0^\infty e^{-\rho t} (e^{-\gamma t} \pi_I + (1 - e^{-\gamma t}) \pi_R) dt = \frac{1}{\rho + \gamma} \left( \pi_I + \frac{\gamma}{\rho} \pi_R \right),$$

which is independent of  $\tau$ ; see also Toxvaerd (2020).

Second, the payoff in (11) can be obtained from

$$\int_0^\infty e^{-\rho t} (1 - p_i(t)) \left[ \pi_S - \frac{c}{2} (1 - \varepsilon_i(t))^2 + \frac{\dot{p}_i(t)}{1 - p_i(t)} V_I \right] dt.$$

With probability  $1 - p_i(t)$  individual  $i$  has not been infected by time  $t$  and receives the flow payoff  $(\pi_S - \frac{c}{2} (1 - \varepsilon_i(t))^2) dt$ . In addition, with probability  $\dot{p}_i(t) dt$  she becomes infected

---

<sup>31</sup>In fact, Equation (31) itself yields  $S_\infty \neq \frac{\gamma}{\beta}$ . Since  $\dot{I}(0) \geq 0$ , we have  $\varepsilon(0) \geq \frac{\gamma}{\beta S_0} > \frac{\gamma}{\beta}$ . Since the first-term of the right-hand side of (31) is not zero, it cannot be the case that  $S_\infty = \frac{\gamma}{\beta}$ .

and receives the lump sum payoff  $V_I$ . The above payoff is obtained by integration by parts. This approach was previously used in [Toxvaerd \(2020\)](#); for the approach dealing with all three states ( $S$ ,  $I$  and  $R$ ) see [Rachel \(2020a\)](#).

**Remark 2.** We have assumed  $p_i(t) < 1$  in deriving equation (12). We show that the condition is satisfied in three steps. First, the proof of the inequality  $S_0 e^{-\frac{\beta}{\gamma}} \leq S_\infty$  in Lemma 3 holds for any SIR dynamics (19), (20) and (21) with  $\varepsilon(\cdot) \in [0, 1]$ . Specifically, it holds for the model with the endogenous cost of infection in which  $\eta$  evolves according to (13). Second,  $\frac{1-p_\infty}{1-p(0)} = \frac{S_\infty}{S_0} > 0$  holds, where  $p_\infty := \lim_{t \rightarrow \infty} p(t)$  and where the equality follows from observations in the proof of Lemma 1 and the inequality from the first step. Third,  $p_i$ , which follows (9), is weakly increasing and satisfies  $p_i = p$  in equilibrium. Then,  $p(t) \leq p_\infty < 1$ , as desired.

**Proof of Lemma 1.** We prove the equation (14) in two steps. First, it follows from equation (9) that

$$\frac{d}{dt} \log(1 - p_i(t)) = -\frac{\dot{p}_i(t)}{1 - p_i(t)} = -\varepsilon_i(t)\beta I(t).$$

Integrating both sides from some  $t_0$  to  $t_1 > t_0$  and taking the exponential yield

$$\frac{1 - p_i(t_1)}{1 - p_i(t_0)} = \exp\left(-\int_{t_0}^{t_1} \beta \varepsilon_i(t) I(t) dt\right). \quad (32)$$

Second, since equation (13) is a linear first-order differential equation, let

$$\mu(t) := e^{-\rho t} \exp\left(-\beta \int_0^t \varepsilon_i(\tau) I(\tau) d\tau\right)$$

be the integrating factor. Since  $\frac{d}{dt} [\mu(t)\eta_i(t)] = \mu(t) (\dot{\eta}_i(t) - (\rho + \beta \varepsilon_i(t) I(t))\eta_i(t))$ , it follows that

$$\frac{d}{dt} [\mu(t)\eta_i(t)] = \mu(t) \left( (\pi_S - \rho V_I) - \frac{c}{2}(1 - \varepsilon_i(t))^2 \right).$$

Integrating both sides on  $[t, \infty)$  and using the transversality condition give

$$\begin{aligned} & e^{-\rho t} \exp\left(-\beta \int_0^t \varepsilon_i(\tau) I(\tau) d\tau\right) \eta_i(t) \\ &= \int_t^\infty e^{-\rho s} \exp\left(-\beta \int_0^s \varepsilon_i(\tau) I(\tau) d\tau\right) \left( (\pi_S - \rho V_I) - \frac{c}{2}(1 - \varepsilon_i(s))^2 \right) ds. \end{aligned}$$

Thus,

$$\begin{aligned}\eta_i(t) &= \int_t^\infty e^{-\rho(s-t)} \frac{\exp\left(-\beta \int_0^s \varepsilon_i(\tau) I(\tau) d\tau\right)}{\exp\left(-\beta \int_0^t \varepsilon_i(\tau) I(\tau) d\tau\right)} \left( (\pi_S - \rho V_I) - \frac{c}{2}(1 - \varepsilon_i(s))^2 \right) ds \\ &= \int_t^\infty e^{-\rho(s-t)} \frac{1 - p_i(s)}{1 - p_i(t)} \left( (\pi_S - \rho V_I) - \frac{c}{2}(1 - \varepsilon_i(s))^2 \right) ds,\end{aligned}\quad (33)$$

where the last equality used (32).

Next, we derive equation (15) in two steps. First, observe that (2) can be rewritten as

$$\frac{d}{dt} \log(S(t)) = -\beta \varepsilon(t) I(t).$$

Integrating both sides from some  $t_0$  to  $t_1 > t_0$  and taking the exponential yield

$$\frac{S(t_1)}{S(t_0)} = \exp\left(-\int_{t_0}^{t_1} \beta(t) \varepsilon(t) I(t) dt\right). \quad (34)$$

Second, in an equilibrium, (33) reduces to

$$\begin{aligned}\eta(t) &= \int_t^\infty e^{-\rho(s-t)} \frac{\exp\left(-\beta \int_0^s \varepsilon(\tau) I(\tau) d\tau\right)}{\exp\left(-\beta \int_0^t \varepsilon(\tau) I(\tau) d\tau\right)} \left( (\pi_S - \rho V_I) - \frac{c}{2}(1 - \varepsilon(s))^2 \right) ds \\ &= \int_t^\infty e^{-\rho(s-t)} \frac{S(s)}{S(t)} \left( (\pi_S - \rho V_I) - \frac{c}{2}(1 - \varepsilon(s))^2 \right) ds,\end{aligned}$$

where the last equality used (34). □

**Proof of Lemma 2.** We first show (17). We rearrange (15) as

$$\eta(t) = \int_t^\infty e^{-\rho(s-t)} \frac{S(s)}{S(t)} (\pi_S - \rho V_I) ds - \int_t^\infty e^{-\rho(s-t)} \frac{S(s)}{S(t)} \frac{c}{2} (1 - \varepsilon(s))^2 ds. \quad (35)$$

For the first term of (35), since  $\pi_S - \rho V_I > 0$ ,

$$\frac{\pi_S - \rho V_I}{\rho} = \int_t^\infty e^{-\rho(s-t)} (\pi_S - \rho V_I) ds \geq \int_t^\infty e^{-\rho(s-t)} \frac{S(s)}{S(t)} (\pi_S - \rho V_I) ds \geq \frac{S(\infty)}{S(t)} \frac{\pi_S - \rho V_I}{\rho}.$$

As  $t \rightarrow \infty$ , the first term of (35) converges to  $\frac{\pi_S - \rho V_I}{\rho}$ . For the second term of (35), observe  $I_\infty = 0$ . This is because, if  $I_\infty > 0$ , then  $R$  is unbounded, which is impossible. By optimality condition (12),  $\lim_{t \rightarrow \infty} \varepsilon_i(t) = 1$ . Then, for any small number  $\kappa > 0$ , there

exists  $t_0 \in [0, \infty)$  such that if  $t \geq t_0$  then

$$0 \leq \int_t^\infty e^{-\rho(s-t)} \frac{S(s)}{S(t)} \frac{c}{2} (1 - \varepsilon(s))^2 ds \leq \int_t^\infty e^{-\rho(s-t)} \frac{c}{2} (1 - \varepsilon(s))^2 ds \leq \frac{c\kappa^2}{2\rho}.$$

Thus,

$$0 \leq \lim_{t \rightarrow \infty} \int_t^\infty e^{-\rho(s-t)} \frac{S(s)}{S(t)} \frac{c}{2} (1 - \varepsilon(s))^2 ds \leq \frac{c\kappa^2}{2\rho}.$$

Since  $\kappa$  is arbitrary, the second term of (35) converges to zero. Hence, we obtain (17), as desired.

As for the bounds, the upper bound is obtained by replacing  $\varepsilon(t) = 1$  and  $\frac{S(s)}{S(t)} = 1$  for all  $s \geq t$  in (15). For the lower bound, it follows from (13) that  $\dot{\eta}_i(t) < 0$  if and only if

$$\eta(t) < \frac{\pi_S - \rho V_I - \frac{c}{2}(1 - \varepsilon(t))^2}{\rho + \varepsilon(t)\beta I(t)}.$$

If  $\eta(t)$  satisfies  $\eta(t) < \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho + \beta}$ , then from time  $t$  on  $\eta$  is always decreasing, which contradicts the statement that  $\eta$  converges to its upper bound as time goes to infinity.

Next, assume  $\dot{\eta}(0) < 0$ . Observe that  $\eta$  is bounded because it is continuous and converges to the finite upper bound (17). Letting  $t_\eta$  be a time at which  $\eta$  attains a minimum, it follows from the assumption  $\dot{\eta}(0) < 0$  that  $\dot{\eta}(t_\eta) = 0$ . Thus,

$$\eta(t_\eta) = \frac{\pi_S - \rho V_I - \frac{c}{2}(1 - \varepsilon(t_\eta))^2}{\rho + \varepsilon(t_\eta)\beta I(t_\eta)}.$$

Substituting for  $\beta I(t)$  from equation (12) for optimal distancing and rearranging yield

$$\begin{aligned} \eta(t_\eta) &= \frac{\pi_S - \rho V_I}{\rho} - \frac{c}{2\rho}(1 - \varepsilon^2(t_\eta)) \\ &\geq \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho}. \end{aligned} \tag{36}$$

Finally, we show that the lower bound  $\frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho}$  is approximately tight when  $\dot{\eta}(0) < 0$ . Substituting (36) into optimality condition (12) yields the following quadratic equation with respect to  $\varepsilon(t_\eta)$ :

$$\frac{\beta I(t_\eta)}{2\rho} \varepsilon^2(t_\eta) + \varepsilon(t_\eta) - 1 + \frac{\beta}{c} I(t_\eta) \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho} = 0.$$



This quadratic equation admits a unique solution  $\varepsilon(t_\eta) \in [0, 1]$ :

$$\varepsilon(t_\eta) = -\frac{\rho}{\beta I(t_\eta)} \left( 1 - \sqrt{1 + 2 \frac{\beta I(t_\eta)}{\rho} \left( 1 - \frac{\beta}{c} I(t_\eta) \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho} \right)} \right).$$

Since  $1 - \sqrt{1 + 2x} \approx -x$  and  $1 - \sqrt{1 + 2x} \geq -x$ ,

$$\varepsilon(t_\eta) \approx \frac{\rho}{\beta I(t_\eta)} \left( \frac{\beta I(t_\eta)}{\rho} \left( 1 - \frac{\beta}{c} I(t_\eta) \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho} \right) \right) = 1 - \frac{\beta}{c} I(t_\eta) \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho}.$$

Comparing the last equation with optimality condition (12), we obtain  $\eta(t_\eta) \approx \frac{\pi_S - \rho V_I - \frac{c}{2}}{\rho}$ . □

**Proof of Proposition 8.** Recall that

$$\frac{dI}{dS} = -1 + \frac{\gamma}{\beta S \max(0, 1 - \frac{\beta \eta}{c} I)},$$

and denote the solution path of the model with endogenous cost of infection by  $(S_e, I_e)$  and its co-state by  $\eta_e$ . By construction,  $\eta_e(\cdot) \in [\eta_L, \eta_H]$ . Therefore, for any fixed values of  $S$  and  $I$ , the following chain of inequalities obtains:  $\frac{dI_L}{dS_L} \leq \frac{dI_e}{dS_e} \leq \frac{dI_H}{dS_H}$ . Finally, recall that all three solution paths go through  $(S_0, I_0)$ .

Consider first the solution path of the model with an endogenous cost of infection and the model with the fixed cost of infection  $\eta_H$ . Since  $\frac{dI_H}{dS_H} \geq \frac{dI_e}{dS_e}$ , at any point of intersection the solution path of the model with the fixed cost  $\eta_H$  intersects the model with the endogenous cost from below. Hence, for  $\delta > 0$  small enough  $I_H(S_0 - \delta) \leq I(S_0 - \delta)$ . But then there can be no intersection for any  $S < S_0$  as otherwise at such an intersection  $\frac{dI_H}{dS_H} < \frac{dI_e}{dS_e}$ . Thus,  $I_H(S) \leq I_e(S)$ .

The proof for the case with the fixed cost of infection  $\eta_H$  is analogous and, therefore,  $I_L(S) \geq I_e(S)$ . □

## B Parameters and Computational Algorithm

We simulate the model at a daily frequency. We follow Farboodi et al. (2021) for most model parameters as summarized in Table 1. We set  $\gamma = 1/7$ , assuming that the average length of disease is 7 days. For the transmission rate  $\beta$  for the baseline simulation of the endogenous cost of infection model, we assume that the initial growth rate  $\frac{\dot{I}(0)}{I(0)}$  without

behavior is 0.3. Since it is given as  $\beta - \gamma$  for the dynamics of the standard SIR model with  $S_0 = 1$ , we set  $\beta = 0.3 + \gamma = 0.443$ . This gives  $R_0 = 3.1$  without behavior. We vary  $\beta$  for various numerical simulations. For  $I_0$ , we match 194 people who died from COVID-19 in the US on or before March 18, a week after the pandemic declaration of the WHO on March 11, 2020. Given a population of 328 million and an IFR of 0.0062, we set  $I_0 = 0.95 \times 10^{-4}$ . We take  $\rho = \tilde{\rho} + \lambda = (0.05 + 0.67)/365$ , where  $\rho$  captures a 5 percent annual discount rate, and  $\lambda$  implies an expected time until the arrival of a cure of 1.5 years as in [Alvarez et al. \(2021\)](#) and [Farboodi et al. \(2021\)](#).

For the flow payoff, we normalize it to be  $-(1 - \varepsilon(t))^2$ . Thus, we set  $c = 2$  and  $\pi_S = 0$ . To compute the parameter  $\eta$  of the constant cost of infection model, we follow the same steps as in [Farboodi et al. \(2021\)](#). We assume the value of a statistical year of life to be US\$ 270,000 and an average remaining life expectancy of COVID-19 victims to be 14.5 years, which gives US\$ 3,915,000 where the numerical values are taken from [Hall et al. \(2020\)](#). Hence, to avoid a 0.1 percent probability of death an individual would be willing to pay US\$  $0.001 \times 3,915,000$ . Using the discount rate to translate this into flow units we obtain US\$  $\rho \cdot 3,915$  as the willingness to pay to avoid the 0.1 percent probability of death. To translate this into utils, we also use the US per capita consumption from [Hall et al. \(2020\)](#) of US\$ 45,000 per year implying that an individual is willing to give up  $\frac{3,915\rho \cdot 365}{45,000} = 31.755\rho$  in terms of annual consumption units, i.e.,  $\varepsilon = 1 - 31.755\rho$ , to avoid a 0.1 percent risk of death. Applying the assumed utility function, an individual, who is willing to give up  $31.755\rho$  units of consumption per period to avoid a 0.1 percent risk of death, is indifferent between this and full exposure with a 0.001 risk of death which has a utility cost of  $v$ :

$$-\frac{(1 - 1)^2}{\rho} - 0.001v = -\frac{(1 - 31.755\rho)^2}{\rho}.$$

Multiplying this value of life in utils by the death rate of 0.0062 (also from [Hall et al., 2020](#)) yields a cost of infection  $\eta = 2761.63$ .

For the endogenous cost of infection model, we set  $\pi_R = 0$  and  $\pi_I = -399.96$  so that  $V_I = \frac{\pi_I}{\rho + \gamma} = -\eta$  works as the upper bound of  $\eta(t)$  in the endogenous cost of infection model. The lower bound of  $\eta$  is  $2761.63 - \frac{c/2}{\rho} = 2254.68$ , which we also use in the constant cost of infection model.

We have solved the constant cost of infection model using the fourth-order Runge-Kutta method. For the endogenous cost of infection model, recall that the equilibrium of the model is characterized as follows. First,  $(S, I, R)$  follow (2), (3), and (4) with the initial condition  $(S(0), I(0), R(0)) = (S_0, I_0, 0)$ , where  $\varepsilon(t) = 1 - \frac{\beta\eta(t)}{c}I(t)$  is the average

Table 1: *Table of Baseline Parameters for Numerical Analysis.*

Parameter	Description	Value	Source
$\gamma$	Recovery Rate	$1/7$	<a href="#">Farboodi et al. (2021)</a>
$\beta$	Transmission Rate	$0.3 + \gamma$	<a href="#">Farboodi et al. (2021)</a>
$I_0$	Initial Seed of Infections	$0.95 \times 10^{-4}$	Based on death toll in the US before March 19, 2020
$\tilde{\rho}$	Discount Rate	$0.05/365$	<a href="#">Farboodi et al. (2021)</a>
$\lambda$	Arrival Rate of Cure	$0.67/365$	<a href="#">Farboodi et al. (2021)</a>
$c$	Cost of Distancing	2	Normalization
$\pi_S$	Flow Payoff of Susceptibles	0	Normalization
$\eta$	Cost of Infection	$\{2254.68, 2761.63\}$	<a href="#">Hall et al. (2020)</a>

exposure. Second,  $\eta$  follows equation (13) with  $\lim_{t \rightarrow \infty} \eta(t) = \frac{\pi_S - \rho V_I}{\rho}$  as in (17).

To numerically solve  $(S, I, R, \eta)$ , we set  $\eta(T) = \frac{\pi_S - \rho V_I}{\rho}$  at  $T = 400 \times 365$  (days). Then, given  $\eta$ , we solve for  $(S, I, R)$  with the initial condition. In turn, given  $(S, I, R)$ , we solve for  $\eta$  with the terminal condition  $\eta(T) = \frac{\pi_S - \rho V_I}{\rho}$ . We iterate the procedure until the sum of the distances of  $(S, I, R, \eta)$  in two successive iterations is below a threshold value. To facilitate the computation, at each iteration, when  $S(t) - S(t+1)$  and  $I(t+1)$  are below threshold values, we have terminated the simulation of  $(S, I, R)$  at  $t+1$ , and we start the computation of  $\eta$  with  $\eta(t+1) = \frac{\pi_S - \rho V_I}{\rho}$  and  $(S, I, R)$ . Once the iterations end, we have checked whether  $\varepsilon(\tau) \in [0, 1]$  for every time  $\tau$ . The right panel of Figure 5 depicts the peak prevalence when  $\varepsilon(\tau) \in [0, 1]$  for every time  $\tau$ .