

Optimal control of deterministic epidemics

Horst Behncke*

FB Mathematik/Informatik, University of Osnabrück, 49069 Osnabrück, Germany

SUMMARY

Various deterministic optimal control models for SIR-epidemics are investigated in this paper. The epidemics are governed by a rather general interaction, which covers most cases studied in the literature. Vaccination, quarantine, screening or health promotion campaigns as forms of control are considered. In all cases one finds a maximum effort control on some initial time interval. In addition, uniqueness and monotonicity properties of these models are studied. The results are also extended to the infinite time-horizon situation. Copyright © 2000 John Wiley & Sons, Ltd.

KEY WORDS: optimal control; deterministic epidemics; infinite time-horizon problems; preventative health control

1. INTRODUCTION

There is a vast amount of literature on epidemics; see, for example, the books of Anderson May [1] and Bailey [2] or surveys of Hethcote [3] and Wickwire [4]. Only a few of these papers, however, deal with the problem of controlling epidemics. Most of these use simulations or model the control via intensity parameters. Other types of approximations are also used. Thus, there are very few papers with analytical results on the optimal control problems for epidemics. This paper takes the works of Wickwire and Morton [5] and Wickwire [6, 7] as its starting point to analyse a number of control models in epidemics. A closely related problem has been studied by Hethcote and Waltman [8] using simulations. While Wickwire analyses Kermack–McKendrick-type epidemics by means of the value function and the Bellman equation, we employ the maximum principle throughout and study the control problem qualitatively, thus avoiding differentiability assumptions on the value function. In this paper models with control by vaccination, quarantine and screening, or health campaigns are studied. The results on vaccination and quarantine extend those of Wickwire [6, 7] and Morton and Wickwire [5], because we use a general type of interaction, allow more general control and cost set-ups and deal with finite and infinite time-horizon situations. In addition, a gap in the paper [5] is closed. As in the above-named papers, we show that in all cases the optimal strategy is maximum effort on some initial time interval. While this result is plausible and natural, it has not yet been shown analytically, except

* Correspondence to: Horst Behncke, FB Mathematik/Informatik, University of Osnabrück, 49069 Osnabrück, Germany.

by Morton and Wickwire [5] for Kermack–McKendrick-type epidemics assuming piecewise continuity of the control and differentiability of the value function.

Even though most papers on epidemics use the mass-action-type interaction βxy or its probabilistically motivated counterpart $\beta xy/(x + y)$, we find it necessary to prove our results for a rather general type of interaction $f(x, y)$, because the assumptions leading to either form are rarely satisfied. These assumptions are that susceptibles and infected which move freely and independently throughout society. The two major objections are that most people only have a limited number of acquaintances and a limited territory in which they move. Secondly, it is well known, in particular for sexually transmitted diseases, that a high prevalence of the epidemic reduces the number of risky contacts. Sometimes this has been taken into account by using density-dependent transmission coefficients and saturation incidence [3, 9, 10]. Further deviations from the mass action law in epidemics may be caused by intermediate hosts in the transmission or interference effects on the part of the susceptibles or agents [1], (Section 6.5.6) and pair formation in sexually transmitted diseases. The following forms for f have been used in the literature $f(x, y) = \beta xy/(x + y)$ [9], $f(x, y) = x \cdot g(y)$ with $g(y) = y^p(1 + \alpha y^q)^{-1}$ [10].

Epidemiological models are usually classified by a simple notation which involves the letters S , E , I and R . S stands for susceptible and refers to the class of people who are not infected. I denotes the infected and E the infected who are not yet infectious. R simply stands for removed and means either the dead or immune. An SIR model thus describes a three-class model in which a susceptible becomes infected through contact with an infected person and eventually becomes immune or dies.

For such a general type of SIR and SEIR epidemics we show that the optimum is given by a maximum effort control at the beginning of the epidemic. This of course is intuitively quite obvious. The proof is based on the analysis of the switching function. This method can also be applied to a number of other more general vaccination models. In addition, a number of further results are shown for such epidemics, such as the uniqueness of the optimum or the monotonicity of the switching point with respect to the time horizon T . These results are new, even for simple Kermack–McKendrick epidemic models. The results on health campaigns are new, although they are rather similar in spirit to models in marketing; see, for example, the book by Feichtinger and Hartl [11]. The dynamics, however, differ substantially from the marketing situation. In as much as the model is linear in the campaign effort control we again have a maximum effort control initially, while the model with non-linear control allows a gradual phasing-out of effort as the epidemic subsides.

This paper is divided into five sections:

1. Introduction
2. SIR epidemics with control by vaccination
3. Quarantine and screening
4. Health-promotion campaigns
5. Infinite time-horizon problems

The notation is largely standard. Thus, g' denotes the derivative of g and partial derivatives will be written as subscripts such as $f_x, f_y \dots$ or $\partial_x f$ in order to avoid confusion. Throughout we assume that most of the functions $f, g \dots$ are sufficiently smooth. Thus, we will always assume that the derivatives that appear in some expressions exist and are continuous. The state space of all problems is $\Omega_n = \{x | x_i \geq 0\}$, the positive orthant in \mathbb{R}^n , where n denotes the number of dependent

variables. Ω_n^0 will then denote the interior of Ω_n . The characteristic function of a set A will be denoted by χ_A .

2. SIR EPIDEMICS WITH CONTROL BY VACCINATION

If we denote the susceptible, infected and removed in the population by $x(t)$, $y(t)$ and $z(t)$, respectively, the dynamics of the epidemics with interaction f can be written as

$$x' = -f(x, y), \quad y' = f(x, y) - \gamma y, \quad z' = \gamma y, \quad \gamma > 0 \quad (1)$$

Here γ denotes the relative rate of removal of the infected. Thus, γ^{-1} can be interpreted as the average duration of the disease. For the interaction f we assume

$$\begin{aligned} f(x, y) &= 0 \quad \text{for } x = 0 \quad \text{or } y = 0, \quad f(x, y) > 0 \quad \text{for } x, y > 0 \\ f_{xx}, f_{yy} &\leq 0 \quad \text{on } \Omega_2 \quad \text{and} \quad f_x, f_y, f_{xy} > 0 \quad \text{for } x, y > 0 \end{aligned} \quad (2)$$

The latter assumptions imply that we assume an almost first-order mass action type of interaction, at least for low densities. Equation (1) covers most cases considered in the literature.

The following lemma is rather elementary and probably known in one form or another. However, it will be needed subsequently. Thus we will only sketch the proof.

Lemma 2.1

For $(x(0), y(0)) \in \Omega_2^0$ the solution $(x(t), y(t))$ of (1) with (2) satisfies $x(t) \searrow x^*$ and $y(t) \rightarrow 0$ for $t \rightarrow \infty$. Moreover, one has $f_y(x^*, 0) < \gamma$ if f is twice continuously differentiable near $(x^*, 0)$ for $x^* > 0$ and $f_{xy}(x^*, 0) > 0$.

Proof. First, we will show that the solution never leaves the positive quadrant, once it starts there. By (2) we have $x' \geq -Kx$ or $(x \cdot e^{Kt})' \geq 0$. Thus $x(t)e^{Kt} \geq x(0)$ or $x(t) > 0$ if $x(0) > 0$. Argue similarly for y . This type of argument will be used repeatedly in the remainder. The assumptions on f and (1) imply $x(t) \searrow x^*$ because the right-hand side of the x -equation is negative for $x, y > 0$. Similarly, (1) implies $(x + y)(t) \searrow$. Thus $y(t)$ also converges for $t \rightarrow \infty$ to y^* . Using (1) again finally shows $f(x^*, y^*) = 0, f(x^*, y^*) - \gamma y^* = 0$. This implies $y^* = 0$. In order to determine the asymptotics of $(x(t), y(t))$, expand f up to second order around $(x^*, 0)$. If $f_y(x^*, 0) < \gamma$, $(x^*, 0)$ is a possible ω -limit point to which the solution converges exponentially because near $(x^*, 0)$ $y' \approx (f_y(x^*, 0) - \gamma)y$. If $f_y(x^*, 0) = \gamma$, consider $dy/dx = -(f - \gamma y)/f \approx f_{xy}(x^*, 0)\gamma^{-1}(x - x^*)$ in order to deduce $y \leq C(x - x^*)^2$ for some C by means of comparison results for differential equations. This shows that the dominant term in the expansion of $f - \gamma y$ is the one associated with f_{xy} . This in turn implies $y' \geq 0$ near $(x^*, 0)$ if $f_{xy}(x^*, 0) > 0$. This, however, is impossible, because y has to approach 0. \square

Remark

The lemma shows that the epidemic dies out exponentially for large t . This, however, will not be needed in the sequel. But we will make use of the fact that y is integrable. This follows from (1) and

the fact that $x + y + z$ is constant. The integrability of y and the equation $dy/dx = -(f - \gamma y)/f$ also show that system (1) is Lyapunoff stable. This means that small changes in the initial values cause only small changes uniformly in t for $x(t)$, $y(t)$, $z(t)$. This fact will be needed when these results are extended to the infinite time-horizon situation. The lemma is also valid if $f_{xy}(x^*, 0) = 0$ and if f is three times continuously differentiable. This latter assumption, however, would be unreasonable from a modelling point of view.

In the case of carrier-borne epidemics [7], where a fraction p of the infected is infective without showing symptoms, equations (1) have to be replaced by

$$x' = -f, \quad y' = pf - \gamma y, \quad 0 < p \leq 1$$

In this case Lemma 2.1 remains valid in as much as it implies the integrability of $y(t)$ and $y(t) \rightarrow 0$.

If vaccination is introduced as a means to control the epidemic, (1) has to be replaced by

$$x' = -f(x, y) - g(x)u, \quad y' = pf(x, y) - \gamma y, \quad g(x) \geq 0, g'(x) \geq 0, g'(0) > 0, \quad 0 < p \leq 1 \quad (3)$$

If $g(0) > 0$ the condition $g'(x) \geq 0$ suffices. Here the term u has to be interpreted as vaccination effort and $g = g(x)$ describes the efficiency or effectiveness of vaccination. In this sense $g' \geq 0$ is reasonable, because vaccination at higher densities is cheaper and easier. For $g \equiv 1$, u is simply the rate of vaccination.

Since technical and financial aspects limit the vaccination rate, one has

$$0 \leq u(t) \leq u_0 \quad (4)$$

In addition, we have to take into account the constraints

$$x(t), \quad y(t) \geq 0$$

Of these only the x -constraint can be active, because it follows from (1), (2) and the first argument in Lemma 2.1 is that $y(t) > 0$ for all $t \geq 0$ if $y(0) > 0$. The cost of an epidemic results from care of the infected and economic losses, as well as from the vaccination. Normalizing the cost of an infected person per unit time to 1, the cost of the vaccination programme in $[0, T]$ can be written as

$$\mathcal{C} = \int_0^T (y(t) + c(x(t))g(x(t))u(t) + du(t)) dt + ay(T), \quad c(x), d \geq 0, \quad \gamma^{-1} \geq a \geq 0, \quad c'(x) \leq 0 \quad (5)$$

Here we consider the cost of vaccination as x -dependent, with $c'(x) \leq 0$, because vaccination at higher densities may be cheaper and easier. The term $ay(T)$ describes additional cost resulting from the remaining infected. $\gamma^{-1} \geq a$ means that only those which are actually present at T are considered.

In general, the relative cost of vaccination c will be rather small, $c \ll 1$. We do not assume this here though, because regardless of the size of c we have:

Theorem 2.1

Assume f satisfies (2), then the above optimal control problem has the solution

$$u(t) = u_0 \cdot \chi_{[0, t^*]} \quad \text{with } 0 \leq t^* < T$$

Proof. We may assume $x(0), y(0) > 0$, because otherwise $t^* = 0$. In the first part of the proof we will also assume $x(t) > 0$ for all $t \leq T$. It follows from the y equation in (3) that $y(t) > 0$ if $y(0) > 0$. The existence of an optimal solution follows easily from Filippov's Theorem [12] as presented in the book by Cesari [13, Section 9.3]. Filippov's proof is based on Ascoli's theorem and a selection result. The necessary conditions for Ascoli's theorem follow from the boundedness of the solutions, here $x(t) + y(t) < x(0) + y(0) + z(0)$, and the compactness of the time interval. The selection theorem requires convexity of the orbit ranges. For linear controls restricted to a finite interval, $0 \leq u(t) \leq u_0$, this is trivially satisfied. It should be noted, however, that Filippov's theorem only guarantees a measurable control. The necessary conditions for an optimum are given in the classical book by Pontryagin *et al.* [14] and in a very general form in Neustadt's book [15]. We will, however, use the article by Hartl *et al.* [16] mainly, because state constraints play an important role here.

Let x, y, u be an optimal solution of our problem. In order to derive the necessary conditions we introduce the adjoint variables $\lambda_0, \lambda_1, \lambda_2$ and the Hamiltonian

$$\mathcal{H} = \lambda_0(y + cgu + du) - \lambda_1(f + gu) + \lambda_2(pf - \gamma y) \quad (6)$$

The adjoint variables satisfy

$$\lambda'_1 = -\lambda_0(c'gu + cg'u) + (\lambda_1 - \lambda_2 p)f_x + \lambda_1 g'u \quad (7)$$

$$\lambda'_2 = -\lambda_0 + (\lambda_1 - \lambda_2 p)f_y + \lambda_2 \gamma$$

The terminal conditions are [16, Section 4.10]

$$\lambda_1(T) = 0, \quad \lambda_2(T) = a\lambda_0, \quad \lambda_0 \leq 0 \quad (8)$$

The singular case $\lambda_0 = 0$ cannot be realized, however, because then the existence and uniqueness theorem for differential equations implies $\lambda_1 = \lambda_2 \equiv 0$, which is impossible.

Thus we may assume $\lambda_0 = -1$. The maximum principle shows that

$$\varphi = -cg - d - g\lambda_1 \quad (9)$$

is a switching function, i.e. $u(t) = u_0$ if $\varphi(t) > 0$, while $\varphi(t) < 0$ implies $u(t) = 0$. Equations (7) suggest to introduce the function

$$\psi = \lambda_1 - p\lambda_2$$

and to rewrite everything in terms of φ and ψ . Then

$$\mathcal{H} = \varphi u - \psi f - y - \lambda_2 \gamma y$$

and

$$\psi' = -\frac{g'u}{g}\varphi - \frac{dg'}{g}u + c'gu + \psi(f_x - pf_y) - p - \lambda_2 \gamma p \quad (10)$$

The shadow price interpretation of λ_1 and λ_2 suggests ψ to be positive for $t < T$, because infected persons are more costly than susceptibles. This is indeed the case near T , because $u = 0$ near T and either $\psi(T) = ap > 0$ for $a > 0$ or $\psi'(T) = -p$ if $a = 0$. Assume there is a maximal $t_1 < T$ with $\psi(t_1) = 0$. Then the slope of ψ at t_1 is positive and (10) and the assumptions on g and c imply

$$0 \leq \psi'(t_1) \leq -p - \lambda_2(t_1)\gamma p$$

Thus $y(t_1) > 0$ results in $\mathcal{H}(t_1) \geq y(t_1)\psi'(t_1) \geq 0$. The Hamiltonian, however, is constant, because the problem is autonomous

$$\mathcal{H}(T) = [-apf - y(1 - a\gamma)](T) < 0$$

This contradiction shows $\psi(t) > 0$ for $t < T$. It follows from (7) that

$$\varphi' = c'gf - g'g^{-1}fd - gf_x\psi - g'g^{-1}f\varphi = A - g'g^{-1}f\varphi$$

Since $A < 0$ the switching function can have at most one zero in $[0, T]$. This must be the switching time t^* , because $\varphi(T) < 0$.

It may happen that $x(t_2) = 0$ for some minimal $0 < t_2 \leq T$, though this is unlikely in realistic situations. In this case, when the state constraint for x is active, the equation for λ_1 has to be modified by an additional distributional term [16, 4.16]. In addition, we note that this can only occur if $g(0) > 0$. If $t_2 = T$ the proof can be completed as above. Thus assume $t_2 < T$. Then $u(t) = 0$ for $t \geq t_2$ or $\varphi(t) \leq 0$ there. In order to reach $x(t_2) = 0$ one needs $\varphi(t) \geq 0$ for $t \leq t_2$ near t_2 . Hence $\varphi(t_2 -) \geq 0$. It follows from [16, Theorems 4.2, 4.16] that λ_1 can at most jump downward. $\varphi(t_2 -) \geq 0$ and $\varphi(t_2 +) \leq 0$, however, preclude any jumps and φ is continuous at t_2 with $\varphi(t_2) = 0$. If $\psi(t_2) > 0$ or $\psi(t_2) = 0$ but $\psi'(t_2 -) < 0$, the proof can be completed as above. On the other hand, $\psi(t_2) < 0$ gives $\varphi(t) < 0$ left of t_2 , which is impossible. The only remaining case is $\psi(t_2) = 0$, but $\psi'(t_2 -) \geq 0$ leads to a contradiction as above.

Remark 1

The theorem shows that the full vaccination programme should be applied as early as possible. This result is, of course, intuitively quite obvious. Likewise, it is apparent that the cost \mathcal{C} is a decreasing function of u_0 . It follows similarly that \mathcal{C} is a monotonically increasing and t^* is a monotonically decreasing function of c if c is constant. On the basis of the shadow price interpretation of the conjugate variables, the positivity of ψ follows, because the cost arises from the infected. The positivity of ψ also shows that an early vaccination programme is more

advantageous. In this proof we have only used the first two conditions of (2), but not that f is twice continuously differentiable.

Theorem 2.1 has been shown by Morton and Wickwire for $f(x, y) = \beta xy$, g, c constant and $T = \infty$ with the assumptions of the differentiability of the value function and piecewise continuity of the controls. This, however, presupposes the existence and uniqueness of the problem [17].

Remark 2

The above solution is unique if f is twice continuously differentiable. In order to see this one defines the cost functional $W(s, x_0, y_0)$ for problem (1), (2) with the control $u_s = u_0 \chi_{[0, s]}$. Then W is twice continuously differentiable in all variables if $s < T$ and $x(t) > 0$ for all $t \leq T$. For the space variables (x_0, y_0) this follows from the smooth dependence of solutions of differential equations on initial values, while the case for s can be shown by direct expansion. Then $\partial_s W = -\phi u_0$. This latter fact can also be used numerically in some gradient-type algorithms. Now formally define with u_s and the corresponding solution x_s the Hamiltonian \mathcal{H}_s as in (6). Then \mathcal{H}_s is piecewise constant, with a possible jump in s . Following the lines of the proof above with this pseudo-Hamiltonian, \mathcal{H}_s , then shows that the second variation of W is positive. As it turns out, $\partial_s^2 W > 0$ is intimately connected with the positivity of ψ , which was also crucial in the proof of Theorem 2.1. Thus $s \rightarrow W(s, x_0, y_0)$ has a unique minimum $s_0 = s_0(x_0, y_0)$ and $\mathcal{C} = W(s, x_0, y_0)$. This also shows that $(x_0, y_0) \rightarrow \mathcal{C}$ is twice continuously differentiable. If the state $x(t) = 0$ is attained, the solution is clearly unique

Remark 3

Even though the above problem has such a simple and intuitively obvious solution, the extension of this method to problems with births in a stable population presents a serious challenge. In this case the equations are

$$x' = -f - u + \mu - \mu x, \quad y' = f - (\gamma + \mu)y \quad (11)$$

where μ is the birth—as well as the death rate. In this case the optimal strategy for large T seems to drive the epidemic just barely below the threshold.

Remark 4

The average rate of expenditure of vaccination during the time of illness is given by γc , since γ^{-1} is approximately the duration of illness. Compared to this we had set the running expenditure of sick people to 1. Assume c and g to be constant, $g = 1$. Common sense tells us that it is advantageous to vaccinate, if vaccination is cheap and if the epidemic is expanding. This is indeed the case, because $\gamma c < 1$ and $f(x(0), y(0)) - \gamma y(0) > 0$ implies $t^* > 0$. Conversely, vaccination should be stopped if it is expensive, $\gamma c \geq 1$, and if the epidemic is subsiding, i.e. $f(x(0), y(0)) - \gamma y(0) < 0$. Indeed, in this case $t^* = 0$ follows easily by evaluating \mathcal{H} at t^* and T . In the first case it turns out that t^* is also a monotonically increasing function of T .

The results of Wickwire and Morton [5] show that these results are optimal, because for $f = \beta xy$, $\gamma c = 1$ and $T = \infty$, t^* is defined by the threshold.

The above result also extends to the corresponding vaccination model with latency. This model is defined by

$$x' = -f - gu, \quad w' = f - \delta w, \quad y' = \delta w - \gamma y \quad (12)$$

with cost

$$\mathcal{C} = \int_0^T (y + cgu + du) dt + ay(T) \quad (13)$$

Here the latently infected are denoted by w . The period of latency will then last about δ^{-1} . $f = f(x, y)$, $g = g(x)$ and $c = c(x)$ are assumed to have the same properties as above.

Theorem 2.2

Assume $f = yh(x)$, g, c, d, a are defined as above. Then the optimal control problem (12), (13) has the optimal solution $u = u_0 \cdot \chi_{[0, t^*]}$ with $0 \leq t^* < T$.

Proof. Since the proof follows the steps of Theorem 2.1 we will only sketch it. The Hamiltonian is given by

$$\mathcal{H} = -(y + cgu + du) - \lambda_1(f + gu) + \lambda_2(f - \delta w) + \lambda_3(\delta w - \gamma y)$$

As above, the switching function is

$$\varphi = -cg - d - \lambda_1 g$$

It satisfies $\varphi' = c'gf - dg'g^{-1}f - \psi gf_x - \varphi g'g^{-1}f$, where $\psi = \lambda_1 - \lambda_2$. Since

$$\psi' = -g'g^{-1}u\varphi - g'g^{-1}ud + c'gu - \delta\phi + \psi f_x, \quad \text{where } \phi = \lambda_2 - \lambda_3$$

one shows first with an argument as in Theorem 2.1 that $\phi(t) > 0$ for $t < T$. Here one needs

$$\phi' = \delta\phi - \psi f_y - 1 - \lambda_3\gamma$$

and

$$\mathcal{H} = \varphi u - \phi\delta w + \psi(f_y y - f) + \phi'y - \delta\phi y$$

The equation for ψ' and the positivity of ϕ imply that $\psi(t) > 0$ for $t < T$. Now the proof can be completed as above. It follows as in Theorem 2.1 that only the state constraint for x can be active. But this case is treated as above. \square

A number of results and methods used above in the proof of Theorem 2.1 will be applied quite often in the same spirit. In order to avoid unnecessary repetitions, we note that Filippov's existence theorem applies to all finite horizon problems, considered here. Similarly, we use $\lambda_0 = -1$ and the constancy of the Hamiltonian.

3. QUARANTINE AND SCREENING

For severe and contagious diseases screening and isolation of the infected represents an efficient means of controlling the epidemic. For humans, however, this meets with severe ethical and moral problems and thus, in general, only restricted versions of this, such as self-reporting and voluntary or anonymous screening are practised. The model below (14) describes a situation in which the infected are treated such that their infectiousness is reduced either due to treatment or because of isolation. The first results in this section extend those of Wickwire. They describe a situation in which self-reporting or isolation of the infected showing symptoms are the dominant modes of control. The final model in this section treats the effect of screening with intensity u and efficiency $g \cdot u$. In all these cases the group of infected is usually much smaller than the total population, i.e. $y \ll x$, $f_x \ll f_y$ and $f(x, y) > \gamma y$, unless x is comparable to y .

If quarantine is not used as a preventive measure to isolate people with an unknown disease status, but only to isolate the infected, one is led to the following model:

$$x' = -f, \quad y' = f - \gamma y - ug \quad (14)$$

where g is a function of x and y , which is rather similar to $a + y/(x + y + b)$ if u is interpreted as an isolation effort or screening efficiency. However, we will only use

$$g > 0 \quad \text{and} \quad -g_x, g_y \geq 0 \quad \text{and} \quad g - yg_y + yg_x \geq 0 \quad \text{in} \quad \Omega_2^0 \quad (15)$$

Since the average duration of the disease is γ^{-1} , the cost can be written as

$$\mathcal{C} = \int_0^T (y + cug) dt + ay(T) \quad \text{with } c > 0 \quad (16)$$

As before, we assume $0 \leq a < c$, $a\gamma < 1$ and $0 \leq u \leq u_0$. In addition, we have the constraints

$$x(t), y(t) \geq 0$$

Of these only $y(t) \geq 0$ can be active only if $g(x, 0) > 0$. This requires the analysis of two cases. But the result is valid also with these constraints.

The cost can be written in this way, because the dominant part of the cost arises from the treatment of those infected and in quarantine, as well as from the corresponding economic loss. In this case both costs are rather similar or $1 \approx c\gamma$. As before, $ay(T)$ measures the detrimental effect the remaining infectives have.

Before we study the general quarantine and screening system let us briefly look at a simplified version, which can be solved explicitly. This system is determined by

$$y' = -\gamma y - u, \quad y(0) = y_0 > 0, \quad 0 \leq u \leq u_0 \quad \text{and} \quad \mathcal{C} = \int_0^T (y + cu) dt. \quad (17)$$

It may be considered as an approximation to a system, where x is very large and y is small. We note, however, that γ^{-1} does not necessarily represent the average duration of the infection any more and $\gamma < 0$ should be allowed for strong epidemics. From the equation $\lambda' = 1 + \gamma\lambda$, $\lambda(T) = 0$ for the adjoint variable one finds that $t^* = 0$ if $c\gamma \geq 1$ and $t^* = T + \gamma^{-1} \ln(1 - c\gamma)$ for the unique switching point t^* . This solution is valid if $y(t) > 0$. Taking this into account one arrives at

$$t^* = \min [T + \gamma^{-1} \ln(1 - c\gamma), \gamma^{-1} \ln(1 + \gamma y_0/u_0)]$$

Theorem 3.1

Assume f satisfies (2), g satisfies (15) and $a < c$ and $c\gamma \leq 1$. Then the above optimal control problem (14) and (16) has the solution $u = u_0 \chi_{[0, t^*]}$ with $t^* < T$. This result also remains valid for $c\gamma > 1$ if $f(x, y) = yh(x)$.

Proof. As before, we treat the case $x(t), y(t) > 0$ first. Again it should be clear that the state $x(t) = 0$ cannot be reached, while $y(t) = 0$ can only be attained if $g(x, 0) > 0$. In this case the Hamiltonian is given by

$$\mathcal{H} = \varphi u - \psi f + \varphi \gamma y g^{-1} + (c\gamma - 1)y$$

where $\varphi = -g(\lambda_2 + c)$ and $\psi = \lambda_1 - \lambda_2$. The transversality conditions give $\lambda_1(T) = 0$ and $\lambda_2(T) = -a$. Thus $\mathcal{H}(T) = -y(T)(1 - a\gamma) - af < 0$. From the adjoint equations one obtains

$$\psi' = g^{-1}\varphi[ug_y - ug_x + \gamma] + \psi(f_x - f_y) + (c\gamma - 1)$$

$$\varphi' = g^{-1}[f(g_y - g_x) + \gamma(g - yg_y)]\varphi - g[1 - c\gamma + \psi f_y]$$

As above, one first shows $\psi(t) > 0$ for all $t < T$. If this were false, there would be a maximal $t_1 < T$ with $\psi(t_1) = 0$ and

$$\psi'(t_1) = [\varphi g^{-1}(ug_y - ug_x + \gamma) + (c\gamma - 1)](t_1) \geq 0$$

One derives from this

$$\mathcal{H}(t_1) \geq g^{-1}\varphi u_0[g - g_y y + g_x y](t_1) + \psi' y(t_1) \geq 0$$

which contradicts $\mathcal{H}(t_1) = \mathcal{H}(T) < 0$. The claim now follows from the fact that $\varphi'(t_0) < 0$, whenever $\varphi(t_0) = 0$. This, however, follows from (18), because $\psi f_y > 0$ and $1 \geq c\gamma$. Now we will treat the case with active constraints.

Assume that there is a minimal $t_1 \leq T$ with $y(t_1) = 0$. Then $\varphi(t) \leq 0$ for $t > t_1$ while $\varphi(t_1) \geq 0$ is necessary to reach $y(t_1) = 0$. Arguments as in the proof of Theorem 2.1 now show that φ is continuous at t_1 with $\varphi(t_1) = 0$. Since $\lambda'_1(t) = 0$ for $t \geq t_1$ one gets $\lambda_1(t) = 0$ for $t \geq t_1$ and $\psi(t_1) = c > 0$. The proof can then be completed as above for $c\gamma \leq 1$. The proof, however, is a little more involved since $\mathcal{H} = 0$ only.

If $c\gamma > 1$ one shows as above $\varphi(t) < 0$ near T and $\psi(t) > 0$ for all $t < T$. Assume there is a maximal t_2 with $\varphi(t_2) = 0$ and $\varphi'(t_2) \geq 0$. Then $\varphi'(t_2) = \mathcal{H}(t_2)g/y \leq 0$, which is a contradiction. \square

Even though the method of isolation is rarely used today, a form of quarantine is employed if infected people are restricted by sanctions or law to mix freely with others and thus spread the disease. This, in fact, seems to be the most widespread form of isolation and is largely applied to sexually transmitted diseases.

Theorem 3.1 can also be extended to the situation with latency by similar methods.

Theorem 3.2

Assume $c\gamma \leq 1$ and that $f(x, y) = yh(x)$ and $g = g(y)$ satisfy the conditions of the previous theorem. Then the corresponding model with latency, which is defined in Theorem 2.2 by

$$x' = -f, \quad w' = f - \delta w, \quad y' = \delta w - \gamma y - ug \quad (18)$$

and costs (16), has the following optimal solutions:

(i) $u = u_0\chi_{[0, t^*]}$ with $0 \leq t^* < T$ if $y(t) > 0$ for all $t \leq T$

$$(ii) \quad u(t) = \begin{cases} u_0, & 0 \leq t \leq t_1, \\ \delta w(t)/g(0), & t_1 < t \leq t_2, \\ 0, & t_2 < t \leq T, \end{cases} \quad y(t) = \begin{cases} > 0, & 0 \leq t \leq t_1 \\ 0, & t_1 \leq t \leq t_2 \\ > 0, & t_2 < t \leq T \end{cases}$$

Since the proof is somewhat more involved than the previous one, but uses essentially the same techniques, we will not present it here.

In this last model we will consider the case of screening the population with an effort u and the removal and treatment of those with positive tests. The model for this is given by (14), (15) and costs

$$\mathcal{C} = \int_0^T (y + cug + du) dt \quad (19)$$

Thus gu has to be interpreted as the amount of positive tests. The cost term cug thus contains the cost for possible further tests and treatment.

Theorem 3.3

Assume f is of the form $f(x, y) = yh(x)$ and $g(x, 0) = 0$ with $g_y, g_x, g_{xy} > 0$ in Ω_2^0 ; $g_{yy} \leq 0$ and $g_{xx}h \leq -g_xh_x$. Then the optimal screening strategy is again of the form $u(t) = u_0\chi_{[0, t^*]}$ with $0 \leq t^* < T$.

Proof. We may assume $y(0) > 0$. Since $g(x, 0) = 0$ one has $y(t) > 0$ for all $t \leq T$. Moreover, the problem can be shown to be piecewise continuous. We will now apply Green's method [18] to this case. Assume $(x(t), y(t))$ is an optimal solution. Let $P_0 = (x(0), y(0))$ and $P_1 = (x(T), y(T))$ and

let $\Gamma = \{(x(t), y(t)) | 0 \leq t \leq T\}$ be the optimal path connecting P_0 and P_1 . Then, using (1), \mathcal{C} may be rewritten as

$$\mathcal{C} = - \int_T \left\{ \left[h^{-1} + \left(c + \frac{d}{g} \right) \left(1 - \frac{\gamma}{h} \right) \right] dx + \left(c + \frac{d}{g} \right) dy \right\}$$

Let Γ' be another admissible solution from P_0 to P_1 with cost \mathcal{C}' . Then Green's theorem allows $\mathcal{C} - \mathcal{C}'$ to be expressed as

$$\mathcal{C}' - \mathcal{C} = \iint_B \frac{d}{g^2} \left(-g_x + g_y \left(1 - \frac{\gamma}{h} \right) \right) ds$$

where B is the area enclosed by Γ and Γ' .

From this one can deduce that maximum effort control should be applied as long as

$$g_y \left(1 - \frac{\gamma}{h} \right) - g_x > 0 \quad \text{or} \quad F = g_y(h - \gamma) - g_x h > 0$$

while no control is optimal if this expression is negative. The conditions on g then imply that the curve defined by $F = 0$ is intersected non-tangentially with positive slope by any solution of (14). Again this shows that the control is maximum effort initially. \square

Remark

If g is independent of x , $t^* > 0$ if $h(x(0))y(0) > \gamma y(0)$. In addition, we note that the solutions for Theorems 3.1 and 3.2 are unique. This is shown as in Section 2.

4. HEALTH-PROMOTION CAMPAIGNS

Health-promotion campaigns are likewise an efficient means of controlling epidemics. The aim of health campaigns is to reduce the spread of epidemics by inducing people to hygienic or risk-averse behaviour. Health campaigning has long been in use in connection with sexually transmitted diseases and with epidemics in developing countries. More recently, it has been widely used in the prevention of AIDS. Health-promotion-campaigns aim directly at reducing the interaction term f responsible for the spreading. Its methods are those of marketing. Thus, the health campaign model introduced here resembles well-known models of marketing [11]. This similarity is rather superficial, however, because the underlying dynamics are quite different.

In health campaigns u measures the amount of funds used for advertising, television and other forms of campaigns, such as counselling or hygienic aid. In a general setting the interaction should be described by $f(x, y, u)$ if we assume an immediate influence of the campaign. For simplicity, however, we model the effect of the campaign multiplicatively, as has been done in connection with models of marketing [11, 19]. Thus we write

$$x' = -w(u)f(x, y), \quad y' = w(u)f(x, y) - \gamma y \quad (20)$$

where f satisfies the same conditions as above. In this context we thus assume an instantaneous action of health campaign measures, the efficiency of which is described by w . This might well apply to campaigns using TV or advertising. Since campaign expenditures are increasingly less effective, we assume $w(0) = 1$ and

$$0 < w(u) \leq 1, \quad w'(u) < 0 \quad \text{and} \quad w''(u) > 0, \quad w''' \leq 0 \quad (21)$$

The cost is given by

$$\mathcal{C} = \int_0^T (y + u) dt \quad (22)$$

and the Hamiltonian becomes

$$\mathcal{H} = \lambda_0(y + u) - \lambda_1 wf + \lambda_2 wf - \lambda_2 \gamma y$$

where λ_1 and λ_2 are, respectively, again the covariables of x and y . It follows from (20) that the constraints $x(t), y(t) \geq 0$ can never be active. Thus $\lambda_1(T) = \lambda_2(T) = 0$ and $\lambda_0 = -1$ follow easily. Arguing as above one sees $\psi(t) > 0$ for all $t < T$. The maximum principle requires maximizing

$$\mathcal{H} = -\psi fw - y - u - \lambda_2 \gamma y = -\psi(fw - \gamma y) - u - y - \lambda_1 \gamma y$$

with respect to u . The unique solution is obviously given by

$$u = 0 \quad \text{if} \quad -\psi fw' < 1, \quad \text{and} \quad \min(u^*, u_0)$$

where $0 \leq u^*$ is the solution of $0 = -\psi fw' - 1$. This shows in particular that $u(t) = 0$ near T , which is quite plausible. By an appropriate choice of w the constraint on u could be disposed of. In order to interpret these conditions determining u , we note that ψ represents the shadow price difference of the susceptible versus the infected. While f is a measure of the force of the epidemic, $-w'(u)$ describes the relative efficiency of advertising. Thus, there will not be any health campaigns if the epidemic is weak, advertising is inefficient or the shadow price for the infected is low. If

$$-\psi fw'(0) > 1 \quad (23)$$

u is non-zero and the health campaign level is a monotonic function of this quantity. From the equations of the covariables and Hamiltonian one deduces

$$(\psi f)' = \psi[\gamma f - \gamma f_y y] - f(1 + \lambda_1 \gamma) = \gamma \psi(f - f_y y) + f/y[-y(T) + \psi(fw - \gamma y) + u]$$

Since the first two summands of this expression are small and negative for large times T compared to γ^{-1} , the behaviour of $(\psi f)'$ is largely determined by $(fw - \gamma y)$. For strong epidemics the campaign level will thus be high initially and phase out slowly as the epidemic subsides.

In the language of marketing the above model might be termed a campaign response model, and it shows some similarity to the well-known Vidale–Wolf model. A model for health campaigns which corresponds to advertising capital models like the Nerlove–Arrow model will now be presented. Here advertising and other campaign measures, such as counselling or meetings, build up a capital stock w which is denoted by goodwill in marketing. In this setting w should be interpreted as concern or the willingness to abstain from risky behaviour. In order to keep the model reasonably simple, however, we will neglect the dynamics of the susceptible. Thus, this model describes a situation in which the number of those susceptible is and remains large compared to the number of infected. This is observed for sexually transmitted diseases in Europe and the U.S. In order to allow an extension to an infinite time horizon, we shall build in a discount rate r , as is common for control models in economics [11]. Following the notation above, this model is determined by

$$y' = f(y, w) - \gamma y, \quad w' = u - \delta w \quad 0 \leq u \leq u_0 \quad (24)$$

and

$$\mathcal{C} = \int_0^T e^{-rt}(y + cu) dt \quad (25)$$

Here f gives the rate of new infections from contact with the infected at an average level of concern w . The rate of forgetting or decay of concern is δ . For f the following assumptions seem natural:

$$f, f_y, -f_w, -f_{wy} > 0, \quad f_{ww}, -f_{yy} \geq 0 \quad \text{if } w, y > 0 \quad (26)$$

Since the number of infected is assumed to be much smaller than the number of the susceptible, f is almost a linear function of y . Later we will therefore assume more specifically

$$f(y, w) = yh(w) \text{ with } h(w) = \rho(\alpha + (1 - \alpha)e^{-\beta w}) \quad (27)$$

The current value Hamiltonian is then given by

$$\mathcal{H} = -(y + cu) + \lambda_1(f - \gamma y) + \lambda_2(u - \delta w)$$

Thus the costates λ_1 and λ_2 for y and w , respectively, satisfy

$$\lambda'_1 = \lambda_1(r + \gamma - f_y) + 1, \quad \lambda'_2 = \lambda_2(r + \delta) - \lambda_1 f_w \quad (28)$$

If $y(0) > 0$, the transversality conditions $\lambda_1(T) = \lambda_2(T) = 0$ imply $-\lambda_1(t), \lambda_2(t) > 0$ for $t < T$. Moreover, one has $u = 0$ near T , which is quite plausible. As above we denote by $\varphi = -c + \lambda_2$ the switching function. Then

$$\varphi' = \varphi(r + \delta) + c(r + \delta) - \lambda_1 f_w \quad (29)$$

In order to eliminate the λ_1 term, differentiate once more and use (29) to substitute λ_1 . One obtains

$$\begin{aligned}\varphi'' &= A(r + \delta)\varphi + ((r + \delta) - A)\varphi' + c(r + \delta)A - f_w \quad \text{with} \\ A &= -(r + \gamma - f_y + f_{yw}f_w^{-1}(f - \gamma y) + f_{ww}f_w^{-1}(u - \delta w))\end{aligned}\tag{30}$$

The substitution $z = \varphi \exp B$ with $B' = -\frac{1}{2}(r + \delta - A)$ transforms (30) into

$$z'' = \left[\frac{1}{4}(r + \delta + A)^2 + \frac{1}{2}A'\right]z + (c(r + \delta)A - f_w)\exp B\tag{31}$$

In short this can be written as

$$z'' = Fz + G \exp B$$

Assuming (27) one gets

$$\begin{aligned}A &= -r + \beta(u - \delta w), \quad F = \frac{1}{4}\delta^2 + \frac{1}{4}\beta^2(u - \delta w)^2 \\ G &= c(r + \delta)(\beta(u - \delta w) - r) + y\rho(1 - \alpha)\beta e^{-\beta\omega}\end{aligned}\tag{32}$$

The following result is now obvious.

Proposition 4.1

Assume F to be non-negative and G to be positive for $u = u_0$, then the control in this model is maximum effort control on some initial time interval $[0, t^*]$ with $0 \leq t^* < T$.

Proof. Contrary to the claim assume φ has a (t_1, t_2) positive arc with $0 < t_1 < t_2$, i.e. $\varphi(t_1) = 0 = \varphi(t_2)$ and $\varphi(t) > 0$ inbetween. This also defines a (t_1, t_2) positive z arc, which is impossible in view of (31). \square

Remark 5

It is obvious that the oscillation properties of z or φ are largely determined by the oscillatory behaviour of G . Thus φ has no negative arcs if G with $u = 0$ is negative and $G \equiv 0$ along a singular arc. Without specifying f any further, no general assertions can be made. For this reason assume (27) now and normalize γ to 1. Moreover, assume $\rho, \delta \geq \frac{1}{2}$, while r , the rate of interest, is assumed to be much smaller. For $w \approx 0$ a campaign effort of $u = \Delta$ increases the concern to $w \approx \Delta$. This in turn results in savings of $f_w(y, 0)\Delta \approx y\rho(1 - \alpha)\beta\Delta$. Taking the decay of concern into account, we can say that campaigning is efficient if $-f_w$ is large and inexpensive if $c \max(1, \delta) < y\rho(1 - \alpha)$. Thus in general we expect G to be positive if campaigning is cheap and efficient and the level of infection high. In those cases initial maximum effort control is optimal.

Remark 6

Since the above analysis is largely qualitative, this result also holds if f is slightly time-dependent. The influence of the susceptible on the dynamics can be modelled in this way.

Remark 7

This result can also be shown qualitatively for an advertising capital model, which takes the susceptible into account. The assumptions, however, will have to be much more specific in order to derive any useful results. For this reason we have abstained from elaborating on this.

5. INFINITE TIME-HORIZON PROBLEMS

In most cases the extension of an optimal control problem to an infinite time horizon causes considerable problems with respect to existence [20], in relation to finite time-horizon problems or the asymptotics of the adjoint variables. In this case these problems are considerably alleviated because of the structure of our solutions. Even though we will formulate the proof for vaccination controls only, it obviously extends to the quarantine models.

Since the cost \mathcal{C}_0 of the no-vaccination policy on $[0, \infty)$ is finite, the infimum \mathcal{C} of the cost on $[0, \infty)$ for (3)–(5) exists. For each finite T let $t^*(T)$ be the switching point of vaccination, which minimizes the cost of the corresponding $[0, T]$ -time-horizon problem. Clearly, $t^*(T) \leq K$ for a suitable K . Now let T_n be any sequence tending monotonically towards infinity. By selecting a subsequence, if necessary, we may assume that $t_n^* = t^*(T_n)$ converges to t^* . Let (x_n, y_n) or (\bar{x}, \bar{y}) be admissible solutions of (3) on $[0, \infty)$ corresponding to $u_n = u_0 \chi_{[0, t_n^*]}$ or $\bar{u} = u_0 \chi_{[0, t^*]}$. Then the stability of (1) shows that (x_n, y_n) converges uniformly towards (\bar{x}, \bar{y}) . Moreover, it follows from (3) that

$$\int_0^\infty y_n \, dt \rightarrow \int_0^\infty \bar{y} \, dt$$

Hence, the corresponding costs also satisfy $\mathcal{C}_n \rightarrow \bar{\mathcal{C}}$. Now let $\varepsilon > 0$ and let $(\tilde{x}, \tilde{y}, \tilde{u})$ be an admissible solution with cost $\tilde{\mathcal{C}} \leq \mathcal{C} + \varepsilon$. A $T_a \geq 0$ exists such that the costs of (\tilde{x}, \tilde{y}) and (\bar{x}, \bar{y}) on $[T_a, \infty)$ are less than ε and such that $|\mathcal{C}_n - \mathcal{C}| < \varepsilon$ whenever $T_n \geq T_a$. For such n one finds $\mathcal{C}_n \leq \tilde{\mathcal{C}} + 2\varepsilon$, because (x_n, y_n) is optimal on $[0, T_n]$. This in turn implies $\mathcal{C} \leq \mathcal{C} + 4\varepsilon$ or $\mathcal{C} = \bar{\mathcal{C}}$. Thus (\bar{x}, \bar{y}) is an optimal admissible solution.

Theorem 5.1

The infinite time-horizon vaccination problems in Section 2 and the corresponding quarantine problem possess optimal solutions. These solutions are given by a maximal effort control on some initial interval. They are unique.

Proof. The above proof also covers the case $g(x, 0) = 0$, provided the $t^*(T_n)$ remain bounded. If the $t^*(T_n)$ are unbounded, we replace (\bar{x}, \bar{y}) above by the solution of (3) with $u \equiv u_0$. Otherwise

we keep the notation and denote the costs for (x_n, y_n, u_n) on $[0, t^*(T_n)]$ by \mathcal{C}_n . For $\varepsilon > 0$ choose a T_a such that the costs for (\bar{x}, \bar{y}) and (\tilde{x}, \tilde{y}) on $[T_a, \infty)$ are less than ε . Then we have for all n with

$$t_n^*(T_n) \geq T_a: \bar{\mathcal{C}} - \varepsilon \leq \mathcal{C}_n \leq \tilde{\mathcal{C}} + \varepsilon \leq \mathcal{C} + 2\varepsilon \leq \bar{\mathcal{C}} + 2\varepsilon$$

because (x_n, y_n, u_n) is optimal on $[0, T_n]$ and because \mathcal{C} is minimal. Thus $\mathcal{C} = \bar{\mathcal{C}}$. \square

Theorem 5.1 does not cover the case with screening of Theorem 3.3, because state $y = 0$ cannot be attained and because $u(t) = u_0$ on an infinite interval leads to an infinite cost, unless the costs are discounted.

REFERENCES

1. Anderson RM, May RM. *Infectious Diseases of Humans: Dynamics and Control*. Oxford University Press: Oxford, 1991.
2. Bailey N. *The Mathematical Theory of Epidemics*, Griffin: London, 1957.
3. Hethcote HW. A thousand and one epidemic models. In *Frontiers. Math. Biology*, Levin SA (ed.), Lecture Notes in Biomathematics, vol. **100**. Springer: Berlin, Heidelberg, 1994.
4. Wickwire KH. Mathematical models for the control of pests and infectious diseases. *Theoretical Population Biology* 1977; **11**:182–238.
5. Morton R, Wickwire KH. On the optimal control of a deterministic epidemic. *Advances in Applied Probability* 1974; **6**:622–635.
6. Wickwire KH. Optimal isolation policies for deterministic and stochastic epidemics. *Mathematical Biosciences* 1975; **26**:325–346.
7. Wickwire KH. A note on the optimal control of carrier-borne epidemics. *Journal of Applied Probability* 1975; **12**:565–568.
8. Hethcote HW, Waltman P. Optimal vaccination schedules in a deterministic epidemic model. *Mathematical Biosciences* 1973; **18**:365–381.
9. Brauer F. Models for the spread of universal fatal diseases. *Journal of Mathematical Biology* 1990; **28**:451–462.
10. Hethcote HW, van den Driessche P. Some epidemiological models with nonlinear incidence. *Journal of Mathematical Biology* 1991; **29**:271–287.
11. Feichtinger G, Hartl RF. *Optimale Kontrolle ökonomischer Prozesse. Anwendungen des Maximumprinzips in den Wirtschaftswissenschaften*. de Gruyter: Berlin, 1986.
12. Filippov AF. On certain questions in the theory of optimal control. *SIAM Journal of Control* 1962; **1**:76–84.
13. Cesari L. *Optimization-Theory and Applications*. Springer: New York, 1983.
14. Pontryagin LS, Boltyanskii VG, Gamkrelidze RV, Mishchenko EF. *The Mathematical Theory of Optimal Processes*. Interscience, Wiley: New York, London, Sydney, 1962.
15. Neustadt L. *Optimization*. Princeton University Press: Princeton, NJ.
16. Hartl RF, Sethi SP, Vickson RG. A survey of the maximum principles for optimal control problems with state constraints. *SIAM Reviews* 1995; **37**:181–218.
17. Cannarsa P, Frankowska H. Some characterizations of optimal trajectories in control theory. *SIAM Journal of Control* 1991; **29**:1322–1347.
18. Hermes H, Haynes G. On the nonlinear control problem with control appearing linearly. *SIAM Journal of Control* 1963; **1**:85–108.
19. Sethi S. Dynamical optimal control models in advertising. A survey. *SIAM Review* 1977; **19**:685–725.
20. Baum RF. Existence theorems for Lagrange control problems with unbounded time domains. *Journal of Optimization Theory and Applications* 1976; **19**:89–115.