**TECHNISCHE
UNIVERSITÄT
WIEN**
**Vienna University of Technology**

# Diplomarbeit

zum Thema

# Non-Linear Incidence Functions in Mathematical Epidemiology

ausgeführt am
Institut für Wirtschaftsmathematik
der Technischen Universität Wien

unter der Anleitung von
Univ. Prof. Dr. Vladimir Veliov

durch
Michaela Killian, BSc.
0625190
Hauptstraße 35
2000 Stockerau

Vienna, on May 13, 2012

**Statutory Declaration**

I declare in lieu of an oath that I have written this master thesis myself and that I have not used any sources or resources other than stated for its preparation. This master thesis has not been submitted elsewhere for examination purposes.

Vienna, on May 13, 2012                                    _____

# Contents

# Introduction

This study considers how non-linearities in the transmission of HIV/AIDS affect population dynamics. Since the early 1980s, when AIDS was first recognized, there has been uncertainty about the future trends and ultimate dimensions of this epidemy. This uncertainty persits because of the difficulties in measuring with any substantial degree of precision, the prevalence and more particularly the incidence of HIV/AIDS infections in any given population. As a result, many varieties of HIV models have been developed in an attempt to understand the dynamics and interrelationships of the major determinants of HIV transmission and the developing reliable estimates of the HIV/AIDS epidemy.

So first of all we are interested in what epidemiology generally is about. Epidemiology is the subject that studies the spread of diseases in populations, and primarily the human population. Mathematical epidemiology is concerned with the quantative aspects of the subject. Often the work of a mathematical epidemiologist consists of

1. model building,

2. estimation of parameters,

3. investigation of the sensitivity of the model to changes in the parameters,

4. simulations and interpretations.

All these activities are expected to tell us something about the spread of the disease in the population, the possibility to control this spread and maybe how to make the disease disappear from the population. The diseases which are modelled most often are the so called infectious diseases, that is, diseases that are contagious and can be transferred from one individual to another through contact. Examples of such diseases are the childhood diseases: measles, rubella, chicken pox and mumps and the sexually transmitted diseases: HIV/AIDS, gonorrhea, syphilis and others. Other examples of infectious diseases are hepatitis, tuberculosis, and one everybody is familiar with: influenza. What constitutes a contact so that the disease can be transmitted in each case is different. For example, we know that the common flu can be obtained merely by being physically close to a person who already has it. In most cases of sexually transmitted diseases, a sexual contact is necessary for transmission.

When a disease spreads in a population it splits the population in nonintersecting classes. These typically are:

- The group of people who can contract the disease under appropriate conditions. These people are called susceptible individuals or susceptibles. The size of this class is usually denoted by S.

- The group of people who have contracted the disease and are now ill with it. What is more important is that they can transmit the disease after a contact to a susceptible individual. These people are called infective individuals or infectives. The size of the class of infective individuals is usually denoted by I. The number of infected individuals in the population is called prevalence of the disease.

- The group of people who have recovered and cannot contract the disease again are called recovered individuals. The class of recovered individuals is usually denoted by R.

The number of individuals in each of these classes changes with time, that is, S(t), I(t), R(t). The epidemiological models consist of a system of ordinary differential equations (ODEw) which describe the dynamics in each class, for example, the rate of change of susceptible individuals $S'$(t) = - (number of individuals who become infected per unit of time). In epidemiology, the number of individuals who become infected per unit of time is called incidence, [4, 6].

In this work we analyse a kind of SI Model with no recovery (because we analyse HIV) and with a non-linear transmission function, which is able to prognose the expansion of the HIV prevalence and provides reasonable predictions for the expansion of the HIV disease based on scarce data. The interplay between epidemiology and population dynamics has been studied in various models. In this study, we discuss different non-linear transmission function in the epidemiology of HIV/AIDS and its relevance to prediction with less given data. We assume, however, that our population is not constant. So we presuppose the whole time that the sum of S and I is dynamical. So one of the main aspects we will discuss in this work is a model for HIV/AIDS with the feature of a contact rate that depends non-linearly on the total population, while this non linear function is the so called transmission or transition rate, which gives the changeover from S to I.

Finally, note that the autors of [15, 26, 22, 19, 12, 5] have studied cases where incidence rates depended non-linearly on the number of infectives or susceptibles. We are going to use some of those transmission functions for experiments with time series data from two different devoloping countries, taken from `www.who.int` (World Health Organization). This corresponds to a non-linear response to disease prevalence, not to a response to population density. Here we start to consider the effects of a nonlinear contact rate. Therefore, population dynamics with density dependent mortality is coupled with epidemics with a general shape of density-dependent transmission rate, which includes both sexually transmitted and environmentally transmitted diseases. Age structure and other

delays are neglected, so that the models result in ordinary differential equations. The epidemic model itself is, of the simplest kind, an SI. It is assumed that infection is permanent and that the infective individuals have a higher mortality than the susceptible. In detail, we have made all experiments with original time series data from Botswana and Swaziland, because these two had an almost increasing prevalence beween 1990 and 2010.

In the first chapter we are going to find a good non linear transmission function for increasing prevalences. Then we will use this found function in the second chapter for different experiments and we will also compare four different transition functions. This happend in experiments with data from Botswana and Swaziland. It will be shown that the transition function $\rho$ we chose in the first chapter approximates the prevalence the best. In the last chapter, Chapter 3, we start with a primitve control of the mortality rate, which you will not find effective in the long run for such highly infectious diseases as HIV/AIDS, but in the short run it is, which definitely makes sense as will be shown later. We would also like to thank Prof. Dr. Veliov, who gave me his notes as a motivation for the last section. You will find his notes about existence of a treshold line for the optimal treatment in a simple epidemic SI model in Section 3.3.

# Chapter 1

# On the Effect of Population Heterogeneity on Dynamics of Epidemic Diseases

This chapter investigates a class of SIS models of the evolution of an infectious disease in a heterogeneus population. The heterogeneity reflects in individual differences in the susceptibility or in the contact rates and leads to a distributed parameter system, requiring corresponding aggregated results of the homogenous population model (ODE) and the distributed one, at least in the expansion phase of the disease. However, this ODE model involves a nonlinear "prevalence-to-incidence" function which is not constructively defined. Based on several established properties of this function, a simple class of approximating function is proposed, depending on three free parameters that could be estimated from scarce data.

How the behaviour of a population depends on the level of heterogeneity–all other parameters kept equally–will be shown in detail later in this chapter. It turns out that for both the short run and the long run behaviour there exist threshold values, such that more heterogeneity is advantageous for the population if and only if the initial prevalence is above the treshold. Essentially, this chapter reprodues considerations and results from [26], which are necessary to justify the subsequent analyses.

## 1.1   The homogeneous and the heterogeneous models

The model we are going to analyse in the following chapter is one of SIS-type, that is such a model, which involves only susceptible and infected individuals. This is the heart of many more detailed epidemiological models. Furthermore, within this easy structure, the model is a rather general one, because nonlinear dependence of fertility, mortality and recovery rates on the population size is allowed, [26, 24].

### 1.1.1 The homogeneous model

We have two main variables in the model, which depend on time. These two variables are the size of the susceptible population S(t), and the size I(t) of the infected population. The dynamics of the disease is described by two differential equations

$$\dot{S}(t) = -\sigma p y S + \lambda(S,I)S + \gamma(S,I)I, \qquad S(0) = S_0, \qquad (1.1)$$

$$\dot{I}(t) = \sigma p y S - \delta(S,I)I, \qquad\qquad I(0) = I_0, \qquad (1.2)$$

where y is the prevalence of the disease, defined as

$$y(t) = \frac{I(t)}{S(t) + I(t)}.$$

The involved parameters and functions are described below:

$\lambda = \eta$ - $\mu$, where

$\eta$ is the birth rate of the susceptible individuals, where it is assumed that the new born of the susceptible individuals are all susceptible as well;

$\mu$ is the mortality rate of the suscebtible individuals;

$\gamma = \nu + \epsilon\tilde{\eta}$ is the inflow rate of susceptible individuals resulting from the infected population;

$\nu$ is the revovery rate from the infection;

$\tilde{\eta}$ is the fertility rate of the infected individuals;

$\varepsilon$ is the fraction of susceptible "'babies"' of infected mothers and

$\delta = \tilde{\mu} + \nu - (1 - \varepsilon)\tilde{\nu}$ is the net out-flow rate of infected individuals, where $\tilde{\mu}$ is the mortality rate of the infected individuals.

The rate of infection $\sigma p y(t)$ consists of three multipliers:

$\sigma$ represents the infectiousness

p is the average level of risk of the population, depending on the average intensity of participation in risky interactions, on the average immunity, etc. and

y(t) is the prevalence of the disease at time t.

The model that includes the above one as a base component has a drawback that could be significant if the population was heterogeneous with respect to risky behaviour of the individuals, in which case the individual values of p may be rather different from the average. In reality, individuals who are vulnerable become infected with higher chance than average. For many diseases, individuals who are more vulnerable as other susceptible are more infective when they become infected. As a result, in the early stage of the disease the rate of infection quests to be higher than $\sigma p y(t)$. Otherwise, if the mortality of the infected one is higher than that of the susceptible individuals, and the recovered ones do not increase their level of risk after recovery, then the average vulnerability to risk in the susceptible population decreases with the time. As a result, the value

$\sigma p y(t)$ overestimates the rate of infection in the late stage of the expansion phase. In the rest of this chapter, we are going to support these observations, and try to avoid the resulting distortion in predicting the evolution of the disease, see [3], which is about heterogeneity in host contact patterns profoundly shapes population-level disease dynamics. They find that human contact patterns are indeed more heterogeneous than assumed by homogeneous-mixing models, but are not as variable as some have speculated. Concluding, they mention that the homogeneous-mixing compartmental model is appropriate when host populations are nearly homogeneous, and can be modified effectively for a few classes of non-homogeneous networks.

### 1.1.2 Modeling the Heterogeneity

In this subsection, we explicitly take into account the heterogeneity of the population, supposing that the value of p is specific for each individual. In the same way as in [26] we introduce a variable $\omega$ that characterizes individual features that are relevant to the disease. As in [26] the variable $\omega$ will be called heterogeneity state, or shortly h-state. We assume that $\omega \in \Omega$, where $\Omega$ is measurable subset of a finite dimensional space.

Now we specify the following two variables on $\Omega$:

- $\bar{S}(t, \cdot)$ is the density of the susceptible population at time t,

- $\bar{I}(t, \cdot)$ is the density of the infected population at time t.

Moreover, $p(\omega) \geq 0$ will denote the level of risk at h-state $\omega$. In principle, $p(\omega)$ could combine the individual susceptibility and the individual contact rate. However, as explained later, more consistent with the other suppositions below is the interpretation of p as the individual rate of potentially risky contacts.

To avoid confusion we note that $\dot{\bar{S}}(t, \cdot), \dot{\bar{I}}(t, \cdot) : \Omega \to \mathbb{R}$ are not probability densities, since their integrals over $\Omega$, denoted furhter by S(t) and I(t), give the total size of the susceptible and of the infected population at t.

The dynamics of the heterogeneous population is described by the following model where "dot" means differentiation with respect to t:

$$
\begin{aligned}
\dot{\bar{S}}(t, \omega) = &-\sigma p(\omega) z(t) \bar{S}(t, \omega) - \mu \bar{S}(t, \omega) \\
&+ \eta \int_{\Omega} \psi_0(\bar{S}(t, \omega), \omega, \omega') \bar{S}(t, \omega') d\omega' \\
&+ \gamma \int_{\Omega} \psi(\bar{S}(t, \omega), \omega, \omega') \bar{I}(t, \omega') d\omega', \\
\dot{\bar{I}}(t, \omega) = &\ \sigma p(\omega) z(t) \bar{S}(t, \omega) - \delta \bar{I}(t, \omega).
\end{aligned}
\tag{1.3}
$$

9

The meaning of the rates $\mu, \eta, \gamma$ and $\delta$ are the same as in the previous subsection. It is supposed that these rates are independent of $\omega$. The density $\psi_0(\bar{S}(t,\omega), \omega, \omega')$ represents the "'probability"' that an offspring of an individual of h-state $\omega'$ is of h-state $\omega$. Similarly, $\psi(\bar{S}(t,\omega), \omega, \omega')$ represents the "probability" that an infected individual of h-state $\omega'$ passes to an h-state $\omega$ after recovery. These probabilities are allowed to depend on the current size of the susceptible population of h-state $\omega$. As before, the rates $\mu, \eta, \gamma$ and $\delta$, but also the densities $\psi_0$ and $\psi$, may depend on the total susceptible and infected population S(t) and I(t), which is not explicitly indicated in the above formulae.

In this chapter $\omega$ has more behavioural than purely biological meaning. That is, it represents habits or vulnerability to risk, rather than natural immunity or frailty. We are going to investigate HIV in Botswana and Swaziland. Therefore, we assume that the newborn individuals have the same h-distribution as the current susceptible population. However, this assumption holds in the case of genetically determined factors of the risk level, here of the susceptible, provided that only the suscepitble individuals give birth to non-infected children. Alternatively, if p is related to behaviour, then one can argue that the susceptible individuals represent a more attractive group to follow than the infected ones. Therefore the newborn individuals accept the behaviour of the former. Similar assumptions can also be made for the recovered individuals. The latter is certainly fulfilled if there is no recovery from the disease, as in our case – AIDS.

We assume further that

$$\psi_0(S, I, \bar{S}(t,\omega), \omega, \omega') = \psi(S, I, \bar{S}(t,\omega), \omega, \omega') = \frac{\bar{S}(t,\omega)}{S}. \qquad (1.4)$$

We got Formula (1.4) from [26], where you can look for the proof of these equalities. However, the term z(t) in (1.3) represents the infectivity of the environment in which the susceptible individuals live. It is called weighted prevalence and is defined as

$$z(t) = \frac{J(t)}{J(t) + R(t)},$$

where

$$R(t) = \int_\Omega p(\omega)\bar{S}(t,\omega)d\omega, \quad \text{and} \quad J(t) = \int_\Omega q(\omega)\bar{I}(t,\omega)d\omega.$$

That is, z(t) is the probability to randomly pick up an infected individual out of pool of all indivuduals, if the individuals are counted according to their weights $p(\omega)$ for the susceptible, and $q(\omega)$ for the infected ones. So the first term in (1.3) assumes the so called separable mixing: z(t) applies to all individuals, rather than allowing mixing preferences depending on p.

Furthermore, we aussume that

$$q(\omega) = \kappa p(\omega), \tag{1.5}$$

which seems reasonable, because if the level of risk is determined by the frequency – with which the individual is involved in a risky interaction – and if this frequency does not change, or changes proportionally, when the individual becomes infected.

The general heterogeneous model, under the simplifying assumption (1.4)-(1.5) and with the notations $\lambda$ and $\delta$ from the previous section, becomes

$$\dot{\bar{S}}(t,\omega) = -\sigma p(\omega)\frac{J(t)}{J(t)+R(t)}\bar{S}(t,\omega) + \lambda(S,I)\bar{S}(t,\omega) + \gamma(S,I)\frac{I}{S}\bar{S}(t,\omega), \tag{1.6}$$

$$\dot{\bar{I}}(t,\omega) = \sigma p(\omega)\frac{J(t)}{J(t)+R(t)}\bar{S}(t,\omega) - \delta(S,I)\bar{I}(t,\omega), \tag{1.7}$$

$$S(t) = \int_{\Omega} \bar{S}(t,\omega)d\omega, \tag{1.8}$$

$$I(t) = \int_{\Omega} \bar{I}(t,\omega)d\omega, \tag{1.9}$$

$$R(t) = \int_{\Omega} p(\omega)\bar{S}(t,\omega)d\omega, \tag{1.10}$$

$$J(t) = \kappa \int_{\Omega} p(\omega)\bar{I}(t,\omega)d\omega, \tag{1.11}$$

with initial conditions

$$\bar{S}(0,\omega) = \varphi_0^S(\omega)S_0,$$

and

$$\bar{I}(0,\omega) = \varphi_0^I(\omega)I_0.$$

The initial conditions are given in terms of the intitial size of the suceptible and of the infected sub-populations, $S_0$ and $I_0$ and the probabilistic densities $\varphi_0^S(\cdot)$ and $\varphi_0^I(\cdot)$ of their h-distributions, see [26].

The measurable mapping $(\bar{S}(\cdot,\cdot), \bar{I}(\cdot,\cdot), S(\cdot), I(\cdot), R(\cdot), J(\cdot))$ is a solution of (1.6)-(1.11) if for almost every $\omega$ the functions $\bar{S}(\cdot,\omega)$ and $\bar{I}(\cdot,\omega)$ are absolutely continuous and (1.6)-(1.11) hold almost everywhere. Below we suppose that the functions $\lambda, \gamma$ and $\delta$ are at

least continuous, or several times continuously differentiable, wherever apropriate. More-over, we assume that $\Omega$ is a closed subset of $\mathbb{R}^r$ with positive Lebesgue measure, that p is a nonnegative measurable function on $\Omega$, and that

$$\text{meas}\left\{(\omega_1, \omega_2) \in \Omega \times \Omega : p(\omega_1) = p(\omega_2)\right\} = 0. \tag{1.12}$$

This condition is obviously fulfilled if $\Omega = [0, 1]$ and p is strictly monotone, or if $\Omega = [0, 1]^r$ and p is continuous and strictly monotone along every line through the origin.

The above heterogeneous model is designed to avoid the shortcomings of the homogeneous models mentioned in the end of the previous subsection. However, it has its own disadvantages:

1. it involves distributed parameter integro-differential equations, therefore is more complicated for calculation;

2. it requires knowledge of the initial densities $\varphi_0^S(\cdot)$ and $\varphi_0^I(\cdot)$ which are usually not available.

The first disadvantage is not critical in simulation tasks but could be an obstacle for the design of optimal treatment or prevention strategies. The second disadvantage makes the direct use of the model impossible even for simulation tasks, due to the lack of data. For the HIV disease in Swaziland or Botwana, even the initial data $S_0$ and $I_0$ are vague, while for the densities $\varphi_0^S(\cdot)$ and $\varphi_0^I(\cdot)$ there is no data available at all, [26, 24].

## 1.2 Approximating the heterogeneous system by an infinite system of ODEs

Now we are aiming to approximate the heterogeneous model by a homogeneous one. This should work in such a way, that the results obtained by the homogeneous model are similar to those that could be achieved by the heterogeneous model if the relevant data were available. Integrating formulae (1.6) and (1.7) with respect to $\omega$ we obtain the following system of differential equations:

$$\dot{S} = -\sigma \rho^*(t) S + \lambda(S, I) S + \lambda(S, I) I, \qquad S(0) = S_0, \tag{1.13}$$

$$\dot{I} = \sigma \rho^*(t) S - \delta(S, I) I, \qquad I(0) = I_0, \tag{1.14}$$

where

$$\rho^*(t) = z(t) \frac{R(t)}{S(t)} = \frac{J(t)}{J(t) + R(t)} \frac{R(t)}{S(t)}. \tag{1.15}$$

Notice that the only information needed to solve the system of differential equations (1.13) and (1.14), to get the aggregated solution $S(\cdot)$ and $I(\cdot)$ of the heterogeneous system, is the function $\rho^*(t)$.

Define the normalized moments

$$m_k^S(t) = \int_\Omega (p(\omega))^k \frac{\bar{S}(t,\omega)}{S(t)} d\omega, \tag{1.16}$$

$$m_k^I(t) = \int_\Omega (p(\omega))^k \frac{\bar{I}(t,\omega)}{I(t)} d\omega, \qquad k = 0, 1, \dots. \tag{1.17}$$

Obviously

$$R(t) = m_1^S(t)S(t), \qquad J(t) = \kappa m_1^I(t)I(t). \tag{1.18}$$

From here we obtain

$$z(t) = \frac{\kappa m_1^I(t)I(t)}{m_1^S(t)S(t) + \kappa m_1^I(t)I(t)}, \tag{1.19}$$

$$\rho^*(t) = z(t)m_1^S(t) = \frac{\kappa m_1^S(t)m_1^I(t)y(t)}{(1 - y(t))m_1^S(t) + \kappa y(t)m_1^I(t)}. \tag{1.20}$$

We use for the prevalence $y(t)$ the notation $y(t) = \frac{I(t)}{S(t)+I(t)}$.

Using equations (1.6) and (1.13) in the expression for $\dot{m}_k^S$ one can obtain the following equations:

$$\dot{m}_1^S = -\sigma z(t)(m_2^S - m_1^S m_1^S),$$
$$\dot{m}_2^S = -\sigma z(t)(m_3^S - m_2^S m_1^S),$$
$$\dots = \dots\dots\dots\dots\dots\dots$$
$$\dot{m}_k^S = -\sigma z(t)(m_{k+1}^S - m_k^S m_1^S),$$
$$\dots = \dots\dots\dots\dots\dots\dots,$$

and similarly, using (1.7) and (1.14),

$$\dot{m}_1^I = \sigma z(t) \frac{1 - y(t)}{y(t)} (m_2^S - m_1^I m_1^S),$$

$$\dot{m}_2^I = \sigma z(t) \frac{1 - y(t)}{y(t)} (m_3^S - m_2^I m_1^S),$$

$$\ldots = \ldots\ldots\ldots\ldots\ldots\ldots\ldots$$

$$\dot{m}_k^I = \sigma z(t) \frac{1 - y(t)}{y(t)} (m_{k+1}^S - m_k^I m_1^S).$$

$$\ldots = \ldots\ldots\ldots\ldots\ldots\ldots\ldots$$

We know the following initial condition:

$$m_k^S(0) = \int_\Omega (p(\omega))^k \varphi_0^S(\omega) d\omega, \qquad m_1^I(0) = \int_\Omega (p(\omega))^k \varphi_0^I(\omega) d\omega.$$

Having in mind equation (1.19) for $z(t)$ and the one for $\rho^*(t)$ (1.20) we established that the above infinite system of differential equations, together with (1.13) and (1.14), determines the solution $(S(\cdot), I(\cdot))$. It can be used for numerical approximation of $(S(\cdot), I(\cdot))$ in a version of the method of Poincare by a truncation of the infinite system. Such an approximating procedure, however, still makes use of the data $\varphi_0^S(\cdot)$ and $\phi_0^I(\cdot)$, which should be avoided, as we argued in the end of Section 1.1, therefore, we do not discuss the details.

Another consequence of the above reformulation of the heterogeneous sytem is the following.

**Theorem 1.2.1** *Assume that $J(0) > 0$. Then for each $k \geq$ the moments $m_k^S(\cdot)$ is strictly monotone decreasing. For $k > 1$ also the normalized moment $\frac{m_k^S(\cdot)}{m_1^k(\cdot)}$ is strictly monotone decreasing. Moreover, if $m_k^S(t) \leq m_k^I(t)$ for $t = 0$, then this inequality holds for all $t \geq 0$.*

We are not going to proof this theorem here, details are provided in [26].

**Theorem 1.2.2** *The heterogeneous system (1.6)-(1.11) does not have a periodic solution with $J(0) > 0$, whatever is the functions $\lambda, \gamma$ and $\delta$.*

The strict monotonicity in Theorem 1.2.1 and the above Theorem 1.2.2 are obtained under the standing assumption (1.12), which obviously excludes the case of a homogeneous system. We mention that the homogenous system (1.1),(1.2) may have a periodic

solution of appropriate $\lambda, \gamma$ and $\delta$ depending on S and I.

The value

$$\frac{m_2^S(t)}{m_1^S(t)} - m_1^S(t)$$

can be viewed as a measure of heterogeneity of the current susceptible population: this fraction is just the $\frac{\text{variance of } p(\cdot)}{\text{mean of } p(\cdot)}$. If the two populations have the same mean risk $m_1^S(0)$ at time 0, then the one with the higher value of $m_2^S(0)$ is more heterogeneous, [26, 24].

## 1.3  Encapsulating heterogeneity in a homogeneous SI-model

In this section, we continue the analysis of the heterogeneous model (1.6)-(1.11). Our goal is to create such a homogeneous system of the form (1.1)-(1.2), which simulates the heterogeneous one. The key point is to replace the multiplier y in (1.1)-(1.2) with a nonlinear function of the prevalence, $\rho(y)$. As a result we will get that there exists such a function $\rho(y)$ for which the solution of (1.1)-(1.2) coincides with the (S,I)-part of the solution of the heterogenous system (1.6)-(1.11) in the expansion phase of the disease, where y increases. The advantage of knowing the appropriate function $\rho(y)$ is that one would not need to know the distributions $\varphi_0^S$ and $\varphi_0^I$ in order to simulate the evolution of the disease in a heterogenous population. However, the exact function $\rho$ will not be constructively defined.

In the next chapter we try to approximate the function $\rho(y)$ by measuring the prevalence y(t) at several moments t and applying a standard identification technique. For this purpose, one has to restrict the search of the function $\rho(y)$ to a class of functions $\Gamma$ depending on a few parameters. In order to justify a choice of the classe $\Gamma$ we first establish some qualitative properties of the function $\rho(y)$.

We stress that the approach below is appropriate for the expansion phase of the disease, where the prevalence y(t) is increasing. Let $(\bar{S}(\cdot, \cdot), \bar{I}(\cdot, \cdot), S(\cdot), I(\cdot), R(\cdot), J(\cdot))$ be a solution of the heterogeneous system (1.6)-(1.11), let $\rho^*(t)$ be the function defined in (1.15), and $y(t) = \frac{I(t)}{S(t)+I(t)}$ be the prevalence. We suppose that $\dot{y}(0) > 0$, thus there is an interval $[0, t^*)$, where $t^*$ may be $+\infty$, such that $\dot{y}(t) > 0$ on $[0, t^*)$ and $\dot{y}(t^*) = 0$, (this equation should be disregarded for $t^* = +\infty$). Let $[y_0, y^*)$ be the set of values of y(t) when t runs in $[0, t^*)$.

Because of the strict monotonicity of $y(\cdot)$, the equation

$$\rho(y(t)) = \rho^*(t) \tag{1.21}$$

defines (in a unique way) the function $\rho : [y_0, y^*) \mapsto \mathbb{R}$. Then the solution of the system

$$\dot{S} = -\sigma\rho\left(\frac{I}{I+S}\right)S + \lambda(S,I)S + \gamma(S,I)I, \qquad S(0) = S_0 \qquad (1.22)$$

$$\dot{I} = \sigma\rho\left(\frac{I}{I+S}\right)S - \delta(S,I)I, \qquad I(0) = I_0, \qquad (1.23)$$

coincides with the (S,I)-part of the solution of the heterogeneous model (1.6)-(1.11). So the definition of the nonlinear prevalence-to-incidence function $\rho$ is not constructive, since it requires knowledge of the solution of (1.6)-(1.11). Hence, the next step to constructive approximation will be to establish some properties of the function $\rho(y)$.

If the heterogeneity condition (1.12) holds, then according to

$$m_p m_s > m_q m_r$$

(which holds for every nonnegative integers $p < q < r < s$ with p+s=q+r) we have for the normalized dispersion

$$d = \frac{m_2^S(0)}{m_1^S(0)} - m_1^S(0) > 0.$$

Now we assume that

$$m_1^I(0) \leq m_1^S(0) + ed, \qquad (1.24)$$

where $e \in (0,1)$. For a homogeneous population $d = 0$ the assumption (1.24) reduces to an equality, which is automatically satisfied in this case. If the population is heterogeneous, then (1.24) allows $m_1^I(0)$ to be somewhat larger than $m_1^S(0)$, which is usually the case, since individuals of higher level of risk get infected with a higher probability.
Solving the differential equation for $m_1^I$ and denoting $q(t) = \frac{\sigma z(t)(1-y(t))}{y(t)}$ we have

$$m_1^I(t) = m_1^I(0)e^{-\int_0^1 q(\theta)m_1^S(\theta)d\theta} + \int_0^t e^{-\int_\xi^1 q(\theta)m_1^S(\theta)d\theta} q(\xi)m_2^S(\xi)d\xi$$

and integrating by parts and rearranging the terms we obtain

$$m_1^I(t) - \frac{m_2^S(t)}{m_1^S(t)} =$$

$$\left(m_1^I(0) - \frac{m_2^S(0)}{m_1^S(0)}\right)e^{-\int_0^1 q(\theta)m_1^S(\theta)d\theta}$$

$$-\int_0^t \frac{d}{d\xi}\left(\frac{m_2^S}{m_1^S}\right)(\xi)e^{-\int_\xi^1 q(\theta)m_1^S(\theta)d\theta}d\xi.$$

Substituting this in the equation for $m_1^I$ we have

$$\dot{m}_1^I(t) = q(t)m_1^S(t)e^{-\int_0^t q(\theta)m_1^S(\theta)d\theta} \cdot \left[\left(\frac{m_2^S(0)}{m_1^S(0)} - m_1^I(0)\right)\right.$$

$$+ \int_0^t \frac{d}{d\xi}\left(\frac{m_2^S}{m_1^S}\right)(\xi)e^{\int_0^\xi q(\theta)m_1^S(\theta)d\theta}d\xi].$$

The first term in the last brackets is constant and positive, according to (1.24). The derivative of $\frac{m_2^S}{m_1^S}$ is strictly negative according to Theorem 1.2.1. Therefore, the second term is zero at t= 0 and monotonically decreasing. We come to the following result.

**Theorem 1.3.1** *If $y(0) \in (0,1)$ and $t^* > 0$, then there exists $t_I \in [0, t^*]$ such that $\dot{m}_1^I(t) > 0$ on $[0, t_I)$ and $\dot{m}_1^I(t) < 0$ on $(t_I, t^*)$.*

It may happen that $t_I = t^*$, that is, $m_1^I$ is increasing in the whole expansion phase.

From the definition of the function $\rho$ we have

$$\rho'(y(t)) = \frac{\dot{\rho}^*(t)}{\dot{y}(t)}.$$

Taking into account (1.20) we obtain that

$$\rho^*(t) = \frac{\kappa m_1^S m_1^I y}{(1-y)m_1^S + \kappa y m_1^I}(t),$$

and

$$\rho'(y(t)) = \kappa\frac{m_1^I(m_1^S)^2 + \frac{\dot{m}_1^I(m_1^S)^2y(1-y)}{\dot{y}} + \frac{\kappa \dot{m}_1^S(m_1^I)^2y^2}{\dot{y}}}{[(1-y)m_1^S + \kappa y m_1^I]^2}(t). \tag{1.25}$$

**Theorem 1.3.2** *Denote $c = \frac{m_2^S(0)}{(m_1^S(0))^2}$. If $y(0) = y_0 \in (0, \frac{1}{1+c\sqrt{\frac{\kappa}{1-e}}})$, then the function $\rho(\cdot) : [y_0, y^*) \mapsto \mathbb{R}$ is strictly increasing close to $y = y_0$. If $t^* < +\infty$ and either $y^* = 1$ or $t_I < t^*$, then $\rho(\cdot)$ is strictly decreasing close to $y^*$.*

**Theorem 1.3.3** *$\rho(\cdot)$ has a bounded derivative in every compact subinterval $[y_0, y^*)$, but $\rho'(y)$ may converge to $-\infty$ at $y^*$.*

For the proofs of the above theorems we refor to [26].

Now we are interested in how to appropriately define a simple "approximation" of the nonlinear prevalence-to-incidence function $\rho$. We will see in Chapter 2 that different transition functions give a different fit of goodness to a given data set and approximate

different prevalence less or better than another transition function. One can look for a simple class of functions that has the above properties and could be used as an approximation of the prevalence-to-incidence function $y \to \rho(y)$. First of all, we chose the form $ay\psi(y)$, where this is a deformation of the linear function $a$ which corresponds to a homogeneous population. Second, the function $\rho$ can be defined only on a subinterval $(0, y^*) \subset [0, 1]$, in which case $\rho'(y^*) = -\infty$. Therefore, we chose the functional form $ay\psi(1 - \frac{y}{y^*})$, where now $\psi$ is a nonnegative differentiable function defined on $(0, 1)$, with $\psi'(0) = \infty$. If, in addition, $\psi(1) > 0$ then the property in Theorem 3.3.2 is automatically satisfied. A simple class of such functions is $\psi(x) = x^\beta : \beta \in [0, 1]$.

Based on the above phenomenological argument we propose the following class of functions depending on four parameters $a, b, \alpha, \beta$ to be used for approximations of the function $y \to \rho(y)$:

$$\rho(y) \equiv ay^\alpha(1 - by)^\beta, \tag{1.26}$$

where

$$a \geq 0, \qquad \alpha \in [0, 1], \qquad b \geq 1, \qquad \beta \in [0, 1].$$

Here, in fact $b = \frac{1}{y^*}$, but the last value is not known in advance. Notice that according to Theorem 1.3.3 we may fix $\alpha = 1$. We deliberately keep a redundancy by allowing values $\alpha \leq 1$, but we shall see later in Chapter 2 that estimating the four parameters from a given data set of prevalences of HIV in different countries leads to $\alpha$ converges to 1. Notice that no data about the heterogeneity of the population is required. Provided that the other parameters in (1.23)-(1.24) are konwn, we only need measurements of $y(t)$ for at least three moments of time.

Concluding, we can say that heterogeneous models with a single parameter of heterogeneity are numerically tractable, but nevertheless are often inapplicable due to the lack of h-distributed data. Section 2–4 propose a bridge from heterogeneous to homogeneous models of the dynamics of infectious diseases. The analysis shows that while a linear prevalence-to-incidence function is reasonable for a homogeneous population, the presence of heterogeneity tends to deform this function in a way qualitatively described in Section 4. We show that the evolution of the heterogeneous population in the expansion phase of the disease can be, in principle, exactly described by a non distributed homogeneous type model which, however, is not constructively defined. Based on the analysis in Section 4 we propose an appropriate class of transition functions which could be used to obtain reasonable approximations to the evolution of the disease within a non-distributed model. The functions from this class depend on four free parameters, one of which we know to be close to 1. The approach is especially appropriate for a disease such as HIV in Africa (where distributed data are scarce, while on the other hand, the heterogeneity plays an essential role). So these four free parameters could be estimated from a modest amount of non distributed data for the history of the prevalence, as we will see later on

in Chapter 2, where we estimate these four paramaters with historical HIV data from Botswana and Swaziland.

## 1.4   The SI Model

As mentioned aboved, since the approach is especially appropriate for a disease such as HIV in Africa we are going to observe models only with a set of susceptible S and infected I in the whole work. If we look at Figure 1.1 we recognize that the paramater $\gamma$, is zero, because individuums who are infected with HIV will not recover.
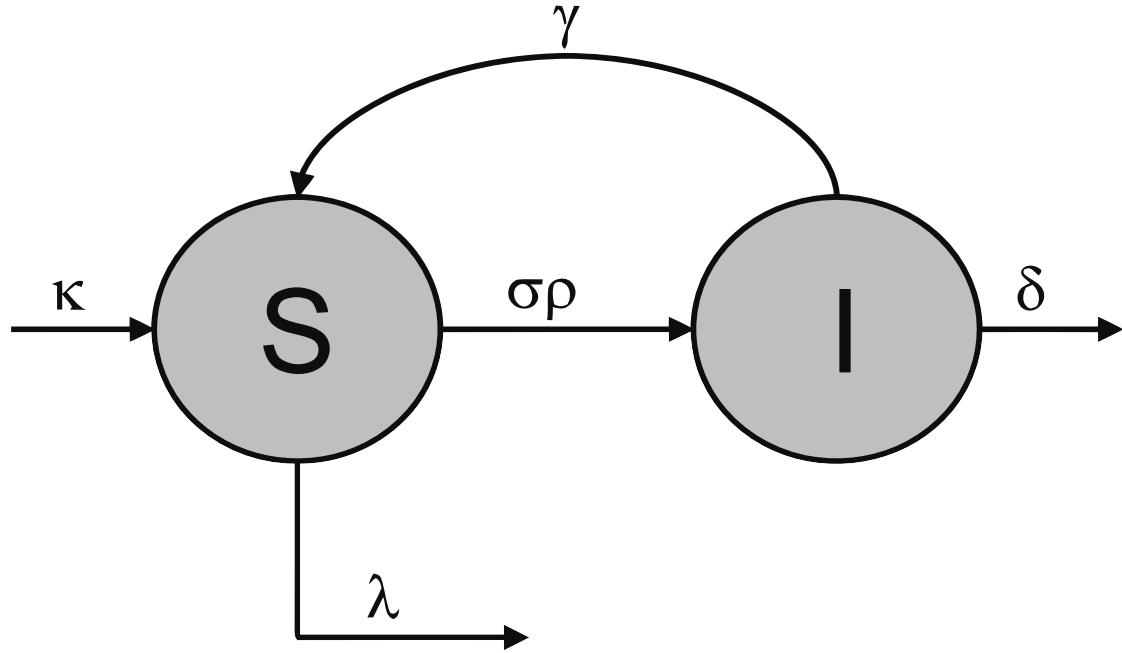


Figure 1.1: This figure shows the transfer diagram of our SI model.

We notice that $\sigma$ is the infectiousness of HIV, $\rho$ is the transmission function, so $\sigma\rho$ represents the crossover between susceptible and infected individuals, $\kappa$ is the birthrate of new individuals, e.g. the income rate into our system, $\delta$ is the mortality rate due HIV, $\lambda$ is the mortality rate of the susceptible and finally $\gamma$ is the recovery rate, which is equal to zero in the whole work.

Now our SI model, which we are going to analyse in the whole work, becomes:

$$\dot{S} = -\sigma\rho(y)S + \lambda S + \gamma I + \kappa, \tag{1.27}$$

$$\dot{I} = \sigma\rho(y)S - \delta I, \tag{1.28}$$

$$y = \frac{I}{S+I}. \tag{1.29}$$

In chapter 2 we will define three other transition functions $\rho_2$, $\rho_3$ and $\rho_4$ ($\rho$ is equal to $\rho_1$) to see if we are able to predict the future prevalence of HIV only with 4 given data points as good as with this $\rho$ for increasing prevalences. For $\rho$ we can take $\alpha$ equal to one in this comparison, because we want to make a fair comparison, where only 4 data points are used for indentification at each $\rho$.

# Chapter 2

# The Transmission Function in our SI Model

The main question that is raised is how pathogen transmission should be modelled. In this Chapter we are going to look at our $\rho$ from Chapter 1 and we want to know if our transmission function describes the prevalence obtained from real data from Botswana and Swaziland in a good way. In Section 2.2 we briefly discuss the phenomenon of mass action, before we will analyse our transmission function. But first of all I will give a short review on how I estimated the free parameters in the different transmission functions.

## 2.1 Least Squares Estimation

The method of least squares assumes that the best-fit curve of a given type is the curve that has the minimal sum of the deviations squared (least square error) from a given set of data.

Suppose that the data points are $(x_1, \theta_1), (x_2, \theta_2), \ldots, (x_n, \theta_n)$ where x is the independent variable, in our case the data set, and $\theta$ is the dependent variable vector, dependent on our sytem. The fitting curve f($\theta$) has the deviation (error) d from each data point, i.e., $d_1 = x_1 - f(\theta_1), d_2 = x_2 - f(\theta_2), ..., d_n = x_n - f(\theta_n)$. According to the method of least squares, the best fitting curve has the property that:

$$d_1^2 + d_2^2 + \ldots + d_n^2 = \sum_{i=1}^{n} d_i^2 = \sum_{i=1}^{n} (x_i - f(\theta_i))^2 \rightarrow \min_{\theta_i}$$

When a solution to the first order condition of the nonlinear least squares minimization (NLS) problem cannot be obtained analytically, the NLS estimates must be computed using numerical methods. To optimize a nonlinear function, an iterative algorithm starts from some initial value of the argument at that function and then repeatedly calculates the next available value according to a particular rule until an optimum is reached approximately. It should be noted that when there are multiple optima, an iterative

algorithm may not be able to locate the global optimum. In fact, it is more common that an algorithm gets stuck at a local optimum, except in some special cases, e.g., when optimizing a globally concave (convex) function. Several new methods, such as the simulated annealing algorithm, have been proposed to find the global solution. These methods have not been standardized yet because they are particularly difficult to implement and computationally very intensive. We will therefore confine ourselves to those commonly used "local" methods. In the following parameter estimations we used the implementation `fminsearch` in MATLAB R2007B. All solvers and MATLAB R2007B implementations are given in Appendix A.

## 2.2 Density-dependent Transmission Efficiency

In this section we give an example of the so called "mass action". We will see that a density-dependent transmission with fixed population size is not as good for epidemiological models as our non-linear $\rho$ from Chapter 1 and the other transmission functions $\rho_2$ to $\rho_4$ we are going to define in Section 2.4.

Transmission is the key process in a pathogen interaction. In most models of such systems, transmission is assumed to occur via so-called "mass action": if we have the density of susceptible S and of infected I, then the number of new infected individuals per unit of time is $\beta$SI, where $\beta$ is the so called transmission coefficient. This $\beta$SI model assumes that infected and susceptible individuums mix completely with each other and move randomly within an arena of fixed size. If this is the case, there is a direct analogy between densities of susceptible and infected people. The question is if the mass-action model is a good approximation of our HIV SI Model. However, in 1995, de Jong published a paper that has been widely interpreted as claiming that $\beta$SI did not represent "true mass action": rather it was a model of "pseudo mass action", and transmission following "true mass action" should be represented by $\frac{\beta SI}{S+I}$. Since then, models have appeared that use either form of transmission, and terminology has been confused.
Sometimes $\beta$SI is described as "mass action", sometimes it is called "density-dependent transmission"; similiarly $\frac{\beta SI}{S+I}$ is in some instances called "mass action", in others it is called "frequency-dependent transmission". Empirical studies comparing models of transmission are only just beginning to appear.

Now the question is raised how pathogen transmission is modelled. If S and I represent densities, rather than numbers, $\beta$SI does represent "true mass action". However, $\frac{\beta SI}{S+I}$, might still give a better representation of the rate of transmission. Several more complex relationships between the densities of both susceptible and infected populution and pathogen transmission have also been proposed in Section 2.3 and 2.4 via $\rho = \rho_1$, $\rho_2$ to $\rho_4$. For a directly transmitted pathogen, the rate at which new infections occur in a population is the product of three things: (1) the contact rate; (2) the proportion of those contacts that are with susceptibles; and (3) the proportion of such "appropriate" contacts that actually result in infection. The assumption underlying mass action is that

the contact rate is directly proportional to density. At the other extreme, the contact rate might be independent of the density. Assuming that susceptible and infected individuums were randomly mixed, this would lead to transmission following $\frac{\beta SI}{S+I}$: on average, each susceptible S would make the same number of contacts regardless of density, and a proportion $\frac{I}{S+I}$ of these would be with infected people, the so called prevalence. This model of transmission is often called "frequency-dependent" or "density-independet". This is exactly the case we were studying in Chapter 1. We are going to make experiments with this density-independet transmission in Section 2.4, because – as we will see – it is going to represent epidemiological models better than density-dependent transmissoned models. The main reason why it is better for our sexualley tranmitted disease is because the number of sexual partners of an individual usually depends on the mating system of the species and is weakly related to density, [7, 19].

Various authors have proposed an asymptotic relationship between the contact rate and density. By contrast, the standard mass action transmission function, with contact rate proportional to density, is equivalent to a Holling Type I functional response. There are obvious and important parallels between contact rates in pathogen transmission and functional responses in predator – prey systems. The proportion of all contacts that are between susceptible and infected people might differ from the random-mixing assumption that underlies both the mass action and frequency-dependent transmission models for several reasons. In both of these models, the assumption is that a proportion $\frac{I}{S+I}$ of all contacts made by a susceptible individual are with an infected person. Alternatively, there might be physiological heterogeneity in susceptibility, which produces a nonlinear relationship between time and number of new infections acquired, as we have seen in Chapter 1. This is the case because the highly susceptible individuals tend to acquire infection first, with resistant individuals acquiring infection later, and at a slower rate, [15, 19].

The transmission coefficient is the most difficult parameter to estimate in any pathogen model. Some attempts have been made to establish it "bottom up" from a priori knowledge of population and disease behaviour, to predict probable disease dynamics and control in a host in which it had not yet been established. There are two approaches which are commonly used to estimate the transmission rate. One is to deduce the transmission coefficient and the form of the transmission function from results of experiments. The second is to deduce it from observations of disease behaviour in the field, in particular prevalence and dynamic responses to perturbations such as control of the population. We are going to estimate all free parameters due to an HIV prevalence time-series, [7, 15, 19].

Several laboratory studies have found that the $\beta$SI model is inadequate for describing pathogen transmission. Assuming mass action, the estimated $\beta$ increased with susceptible person density and decreased with the density of infected people. For the aggressive virus of HIV, either negative binomial transmission, or a power relationship – these transmission functions are defined in the Section 2.4 as $\rho_3$ for the power relationship and

$\rho_4$ for the negative binomial – were markedly superior to density-independent or rather frequency-dependent transmission. D'Amico [25] in [19] fitted a mass-action model and found that the estimated transmission coefficient declined with both infected and susceptible densities, showing that the mass-action model was inadequate to describe the transmission process. In [19] there are many descriptions of different experiments, but each of these small-scale experiments showed that simple mass action does not describe transmission adequately. A more fundamental problem is that it is difficult, if not impossible, to translate estimated rates, or even functional forms of transmission dynamics, from small-scale, homogeneous enclosures to large-scale, heterogeneous landscapes. An alternative approach to deducing the nature of transmission dynamics is to compare the fit of alternative transmission models to observed disease dynamics.

Begon and co-workers [29, 30] in [19] concluded that $\frac{\beta SI}{S+I}$ is a better descriptor of transmission dynamics than is density-dependent transmission $\beta SI$. Dobson and Meagher [31] in [19] compared models with density-dependent and frequency-dependent dynamics, and showed that frequency-dependence more accurately predicted the observed level of disease prevalence, although both transmission models captured the qualitative dynamics adequately. Increasingly, the weight of evidence is that simple mass action is not an adequate model in many situations.

Now we are going to give an example of the mass action transmission in our SI Model (1.27)-(1.29). Our $\rho$ is now equal to $\beta$SI. So, for this experiment, we used an HIV prevalence-time-series from Swaziland between 1990 and 2010, which we got from `www.who.int`. However, we can see the free parameters of the transmission function $\rho$ for two devoloping countries, Swaziland and Botswana, which are examples of those countries with the highest HIV prevalence, after estimating the mass action coefficient $\beta$ in Figure 2.1. The time series data we needed for the estimation above was provided by `www.who.int`. All data in this and the next section has been gathered between 1990 and 2010.
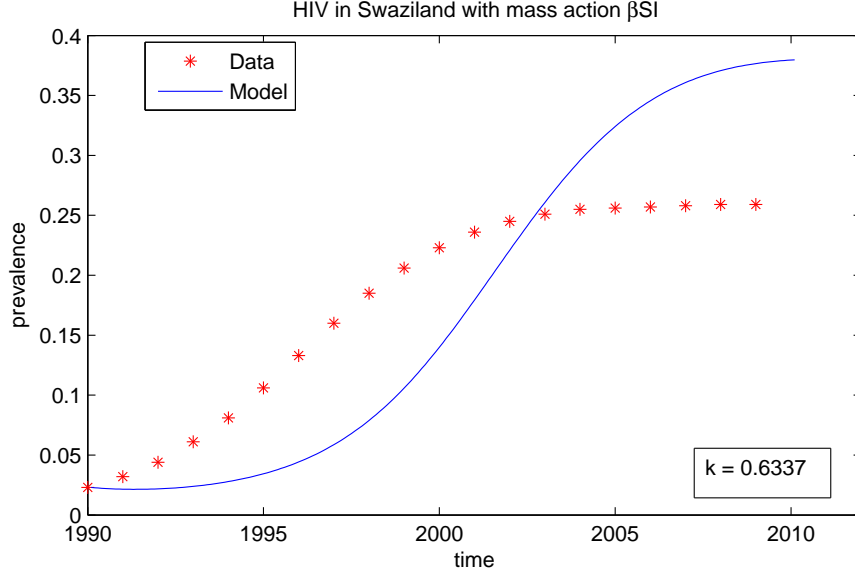
Figure 2.1: This figure shows the prevalence data and the calculated prevalence based on parameter identification using 20 data measurements for the mass action coefficent. The used transmission function is $\beta SI$.

In Figure 2.1, the red stars show the real prevalence data points of Swaziland. The blue line shows the Model with estimation coefficent k. As a result we can see that the mass action is not able to give us a good fitting curve on our real prevalence data from Swaziland, although, all data points were given as an input.

In comparison, Sections 2.3 and 2.4 will show in which excellent way the density-independent functions describe this epidemiological model of ODEs (1.27)-(1.29), especially with all given data points. In addition, we are going to see how well some of these functions are to predict the future prevalence only with 4 given data points.

## 2.3   Using $\rho = ay^\beta(1 - by)^\beta$ for Botswana and Swaziland

Now, as in Section 2.2 we are going to estimate the free parameters of the transmission function $\rho = ay^\alpha(1 - by)^\beta$, which we got in Chapter 1. The progress of estimation is going on in the way as it is described in Section 2.1. However, we did this estimation for two developing countries, Swaziland and Botswana. These are two of those countries with the highest HIV prevalence. The time series data of Swaziland and Botswana, we needed for the estimation, we got from www.who.int, as in Section 2.2. All data in this and the next Section is from the time period between 1990 and 2010.
For a start, we will have a look at the estimation of the four free parameters for Swaziland. In order to estimate the unknown parameters a, b, $\alpha$ and $\beta$ for the transmission function $\rho$, we minimize the square error:

$$\sum_{t=1990}^{2010} \left(P_{[data]}(t) - P_{[a,b,\alpha,\beta]}(t_1)\right)^2 \tag{2.1}$$

In the minimization $P_{[data]}(t)$ means the prevalence data points which came from the data set, $P_{[a,b,\alpha,\beta]}(t_1)$ is the solution of the ODE system given in our SI-Model (1.27)-(1.29) where $t_1$ is the adequate timestep dependent on step size h, which we needed in the ODE solver (please see Appendix A).

As we have seen in Section 1.3, the parameter $\alpha$ can be taken equal to 1 by theoretical resums. We use it as a free parameter just to show that the estimation from the data-fitting give a result close to one, which supports the theoretical conclusion. Fixing $\alpha = 1$ does not essentially change the results, therefore the transmission function $\rho$ can be viewed as depending on 3 free parameters only.

We now need the other parameters, which are constants for each country. We need them to calculate the solutation of the ODE System ($\dot{S}, \dot{I}$). Those constants for Swaziland we also got from http://www.who.int.
Our given parameter values are

- $\sigma = 0.43$ ,

- $\lambda = $ -0.153,

- $\gamma = 0$,

- $\delta = 0.259$,

- $\kappa = 0.2825$.

After calculating the minimal square error, for all 20 given data points, the minimal error was achieved for a= 1.4732, b= 3.1972, $\alpha$ =0.9555 and $\beta = 0.3698$. Figure 2.2 shows the resulting fit to the data.
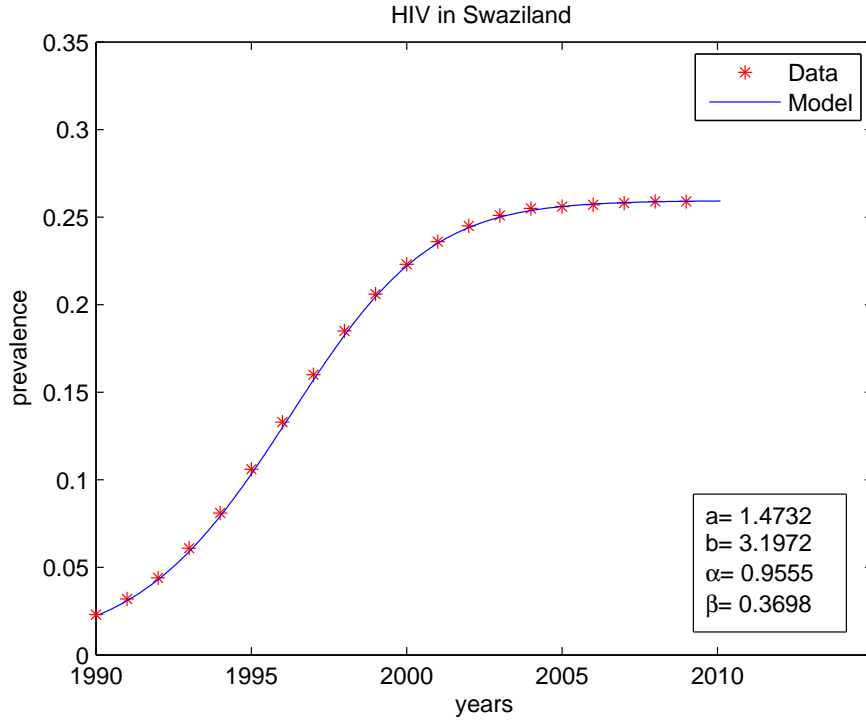
Figure 2.2: Prevalence data (red stars) and calculated prevalence based on parameter identification using 20 data measurements (blue line).

As we can see in Figure 2.2, the transmission function $\rho$ gives nearly an optimal approximation of the real data points. So what we want to know now is if our model with this $\rho$ is also able to predict the prevalence between 1997 to 2010, when we only have 10 given data points. After parameter identification we got for a= 1.5387, b= 3.3355, $\alpha$=0.9686 and $\beta$ =0.3729, as is shown in Figure 2.3. As we described in Chapter 2 and above, the value for $\alpha_{10\text{given points}}$ is higher than those of $\alpha_{20\text{given points}}$. We also see that our transmission function is able to predict the prevalence in a very good way, if only 7 data points were known.
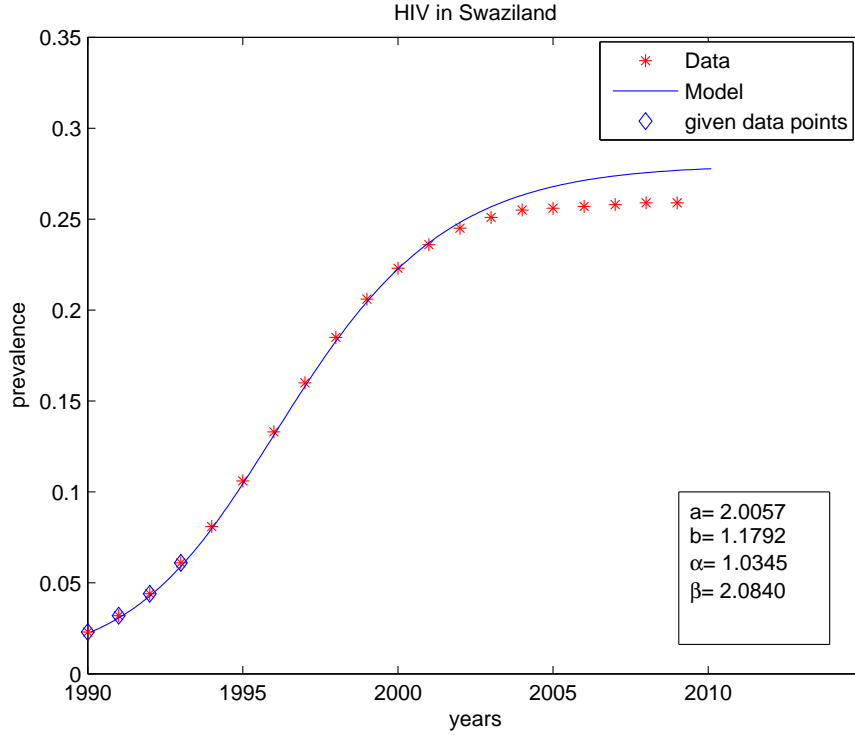
Figure 2.3: Prevalence data and calculated prevalence based on parameter identification using 10 data measurements.

Notice that the blue diamonds in the figures represent the used given data points.
The last figure, Figure 2.4, shows the prediction for HIV in Swaziland for 1994 to 2010. In the last experiment we only have the lowest limit of given data points, namely four. After paramter identification we got values for the four free parameters about a=2.0057, b= 1.1792, $\alpha = 1.0345$ and $\beta = 2.0840$. The transmission function $\rho$ with parameter identification based only on 4 points give a good prediction of the evolution of the prevalence in Swaziland between 1994 and 2003.

Figure 2.4: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements.

In particular, we see that $\alpha$ is close to 1, as predicted by the theory in Section 1.3.

Now we will take Botswana as an example for a country with an increasing prevalence, but we will only choose the data points up to the maximum. Swaziland's maximal element of the prevalence between 1990 and 2010 was the last data point, so its prevalence was monotonically increasing, but Botswana's prevalence between 1990 and 2010 has the maximal prevalence of 26,3% in year 2003. From 2004 to 2010 the prevalence is slowly decreasing. However, what we want to know is if our transmission function is able to estimate the real data points of the expansion phase in a good way and, if it does, we are also interested in the prediction. So the next experiment will show that $\rho$ is able to make a good prediction with minimal data points for the phase of increasing.

In order to estimate the unknown parameters a, b, $\alpha$ and $\beta$ for the transition function $\rho = ay^\alpha(1 - by)^\beta$, we minimize the square error, for given start values for a, b, $\alpha$ and $\beta$, in the same way as we did above for Swaziland:

$$\sum_{t=1990}^{2003} \left( P_{[data]}(t) - P_{[a,b,\alpha,\beta]}(t_1) \right)^2. \tag{2.2}$$

29

The other parameters were also taken from an HIV data time serie between 1990 to 2010, provided by http://www.who.int.

Our given constants are

- $\sigma = 0.38$ ,

- $\lambda$ = -0.108,

- $\gamma = 0$,

- $\delta = 0.23$,

- $\kappa = 0.2397$.

After calculating the minimal square error for all 13 given data points, the minimal error was achieved for a=1.6169 , b=3.5152 , $\alpha$=0.9182 and $\beta$= 0.2812.
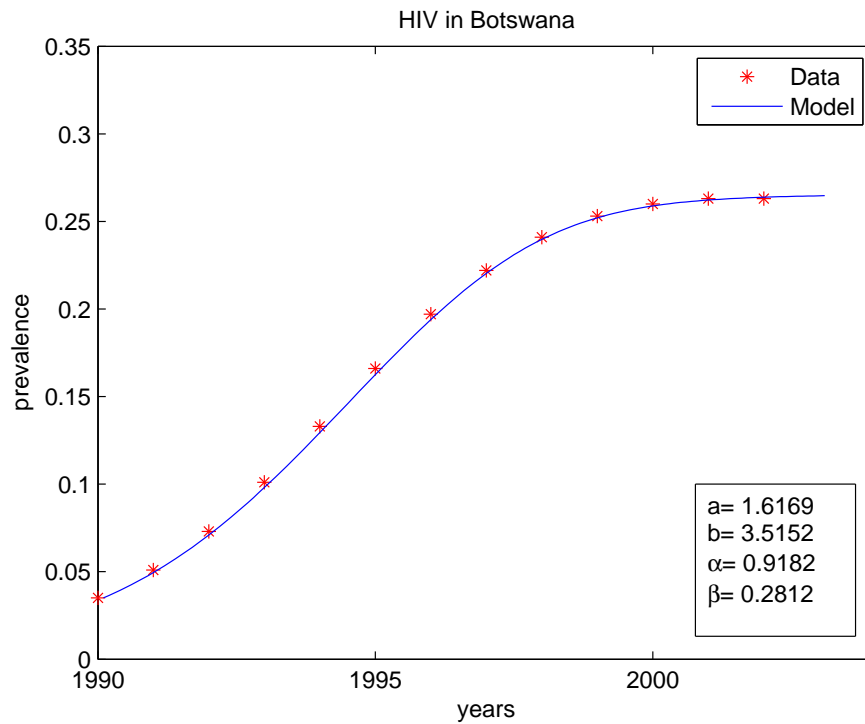


Figure 2.5: Prevalence data and calculated prevalence based on parameter identification using 13 data measurements.

In Figure 2.5, we see that our transmission function also represents the prevalence of Botswana in a very good. In the Figures 2.6 and 2.7 we will see the predictions with 7 given data points and 4 given points. We also see that $\alpha$ is in all cases close to 1.
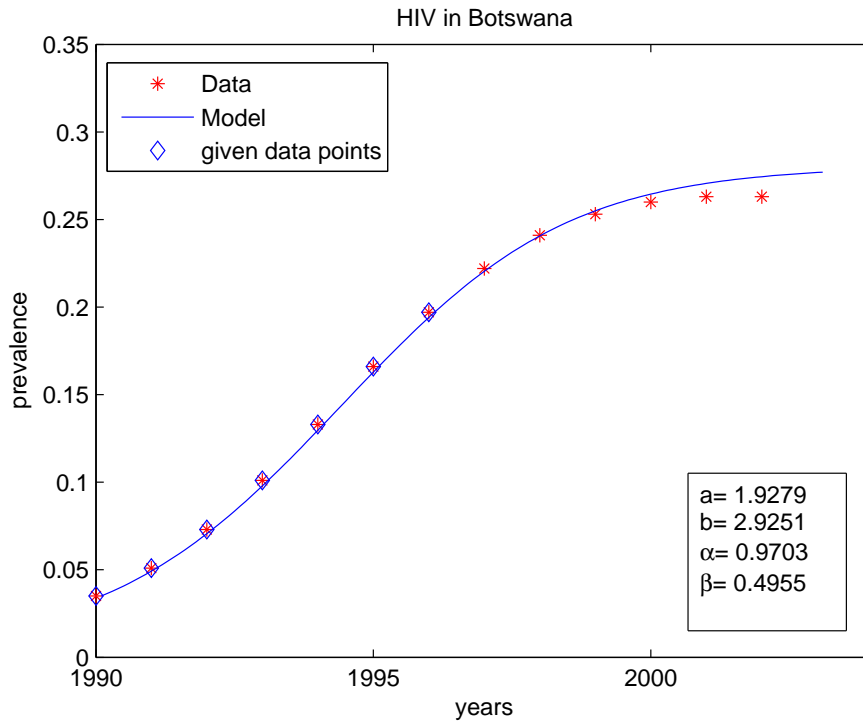


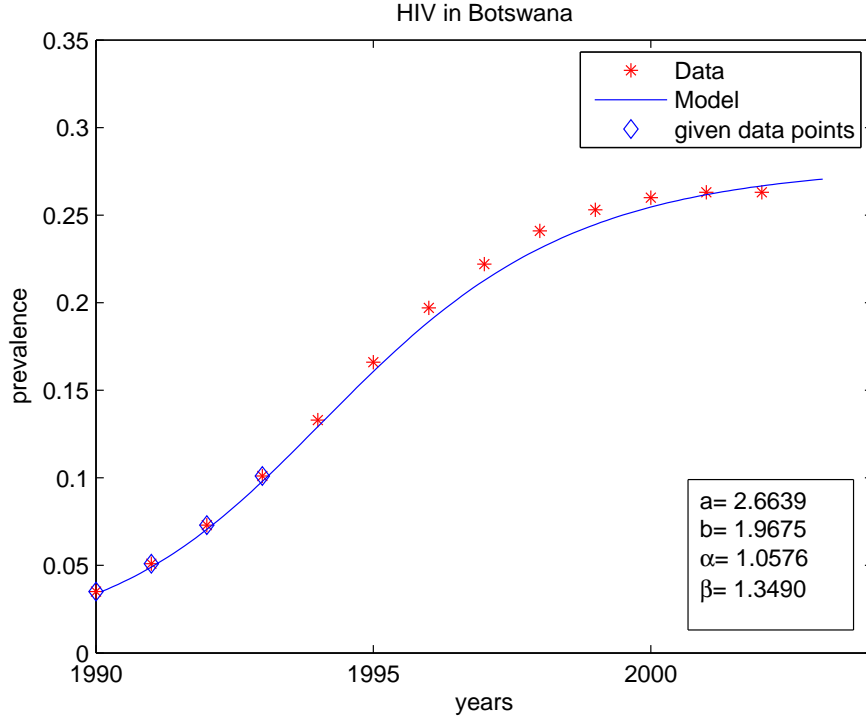Figure 2.6: Prevalence data and calculated prevalence based on parameter identification using 7 data measurements.

Figure 2.7: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements.

In Figure 2.6, after the parameter identification we got a= 1.9279, b= 2.9251, $\alpha$= 0.9703 and $\beta$= 0.4955. In Figure 2.7 we see in a very pleasing way how good our transmission function describes the increasing phase of the prediction, and for the parameter values we got a= 2.6639, b= 1.9675, $\alpha$=1.0576 and $\beta$= 1.3490. In comparison, Figure 2.7 gives a very good prediction of the real data points.

Concluding, this section has shown that the theoretical considerations in Chapter 1 are rather efficient in the context of experiments with real HIV-time-series.

## 2.4 Comparison of four different non-linear transmission functions

In this section we use three other non-linear incidence functions from [19] and [12] with our $\rho$ for an experimental comparison. As a consequence, we are interested in the question if our used transmission function approximates the real data better or worse than other functions.

First of all, we are going to define the new functions. Notice that $\rho$ is in this sec-

tion equal to $\rho_1$. The first of the new transmission is from [12]. The function $\rho_2$ has three free parameters. Firstly, we have to estimate values for, p, q and $\alpha$. This function is defined by

$$\rho_2 = \frac{I^p}{1 + \alpha I^q}. \tag{2.3}$$

The third one is called Power relationship function and is taken from [19]. This function has three free parameter values, p, q, $\beta$ and can be defined as

$$\rho_3 = \beta S^p I^I. \tag{2.4}$$

The last transmission function we used for comparison is also taken from [19]. It has two free parameters k and $\beta$, while it shall be mentioned that small k corresponds to highly aggregated infection. This function is presented by

$$\rho_4 = kS \ln(1 + \frac{\beta I}{k}). \tag{2.5}$$

The first experiments will be for Swaziland. On the one hand, we will show a comparison where every transmission function got all 20 given data points from the HIV-time serie to estimate the free parameters and choose them in the best way for a good fit, and on the other hand, we will show a comparison of the predictions. Due the fact that we want to show a fair comparison in this experiment, all functions got 4 given data points for their parameter approximations. In the first step you can see the figures with all given data points, which are shown in fFigure 2.8 to Figure 2.11.
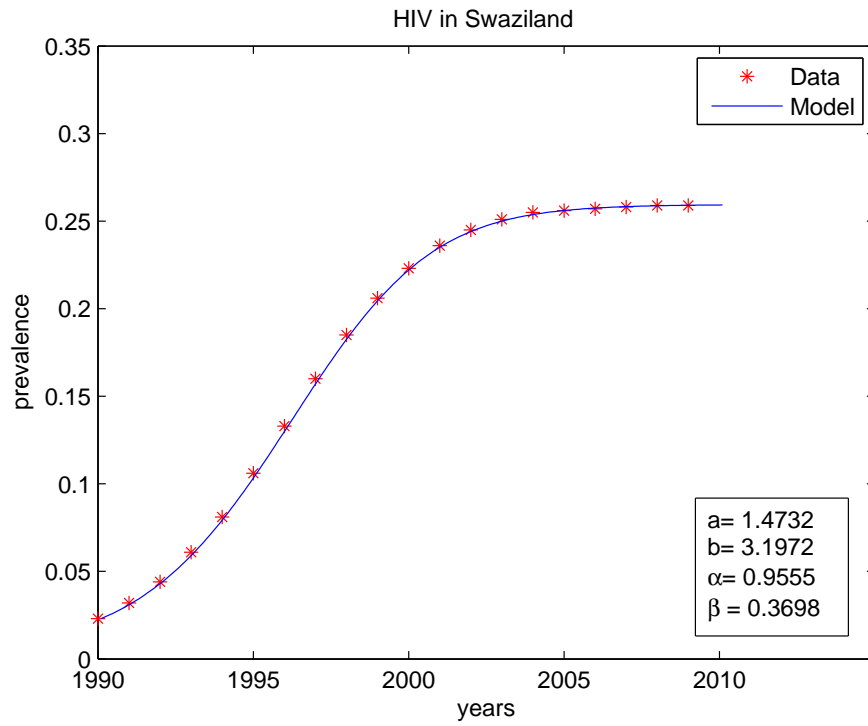
Figure 2.8: Prevalence data and calculated prevalence based on parameter identification using 20 data measurements for the transmission function $\rho_1 = ay^\alpha(1 - by)^\beta$.
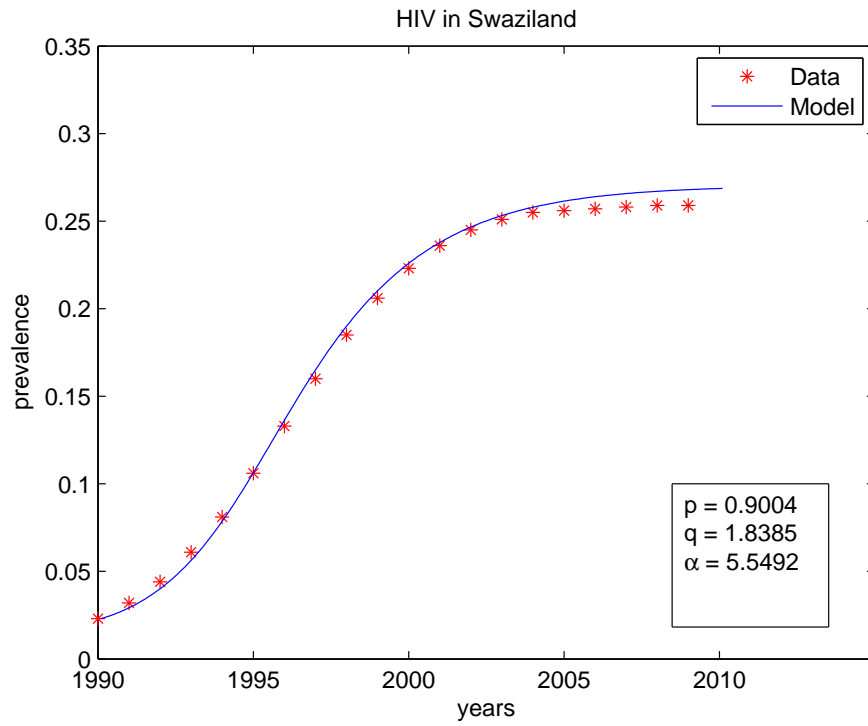
Figure 2.9: Prevalence data and calculated prevalence based on parameter identification using 20 data measurements for the transmission function $\rho_2 = \frac{I^p}{1+\alpha I^q}$.
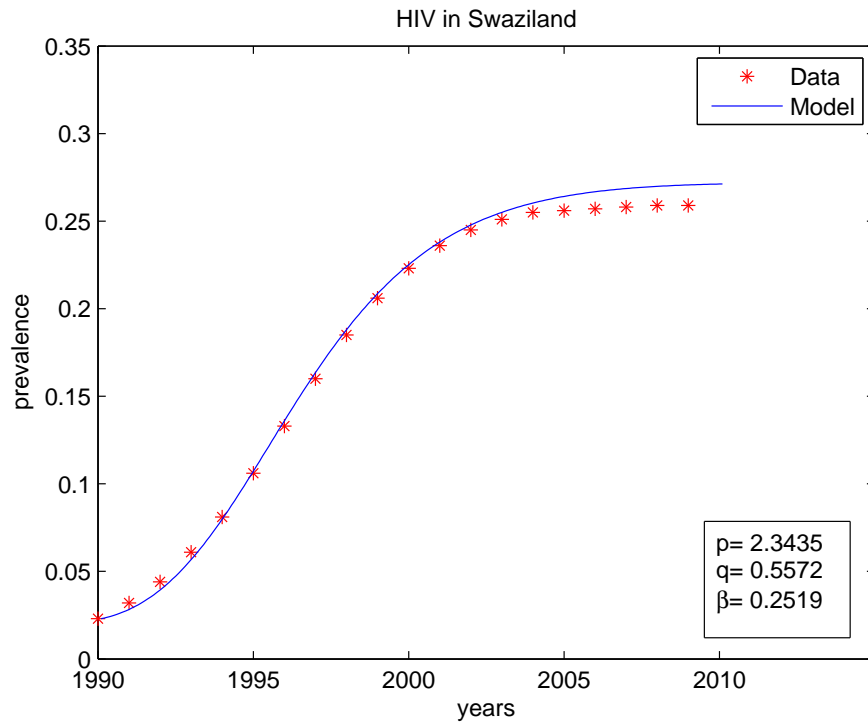
Figure 2.10: Prevalence data and calculated prevalence based on parameter identification using 20 data measurements for the power relationsship $\rho_3 = \beta S^p I^I$.
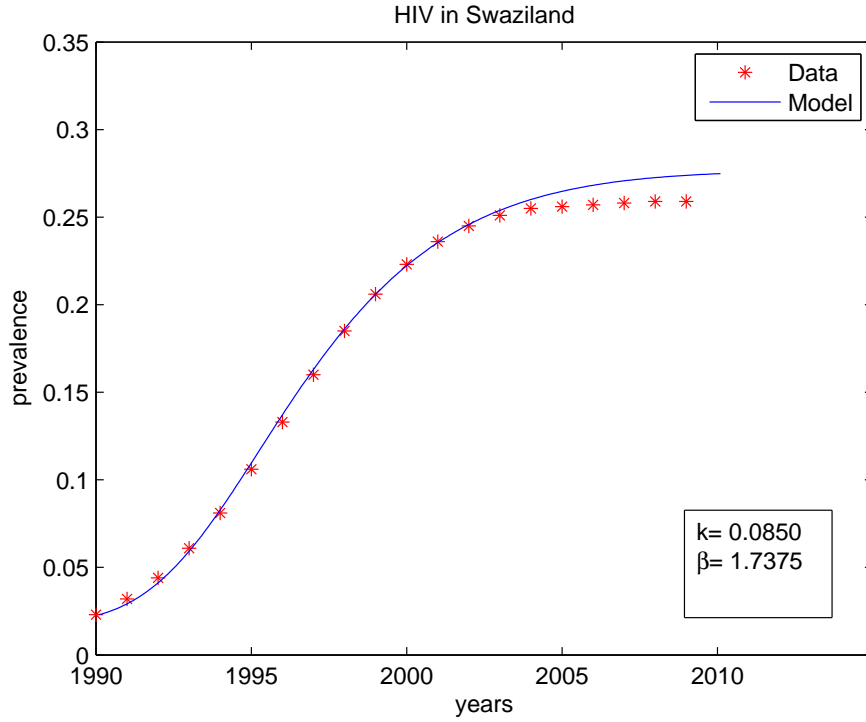
Figure 2.11: Prevalence data and calculated prevalence based on parameter identification using 20 data measurements for negative binomial $\rho_4 = kS\ln(1 + \frac{\beta I}{k})$

As a result, we can see that our function $\rho_1$ gives the best fit for the HIV prevalence time series data for 20 given data points, but also $\rho_2, \rho_3$ and $rho_4$ give a remarkably good fit to the data. Therefore, when we look at the prediction with only 4 given data points, we usually know from above that our function gives a very good prediction to the prevalence data from Swaziland's HIV prevalence. But we are interested in how good the other transmission functions are able to fit the data with less given points. This is shown in Figure 2.12 to Figure 2.15 .
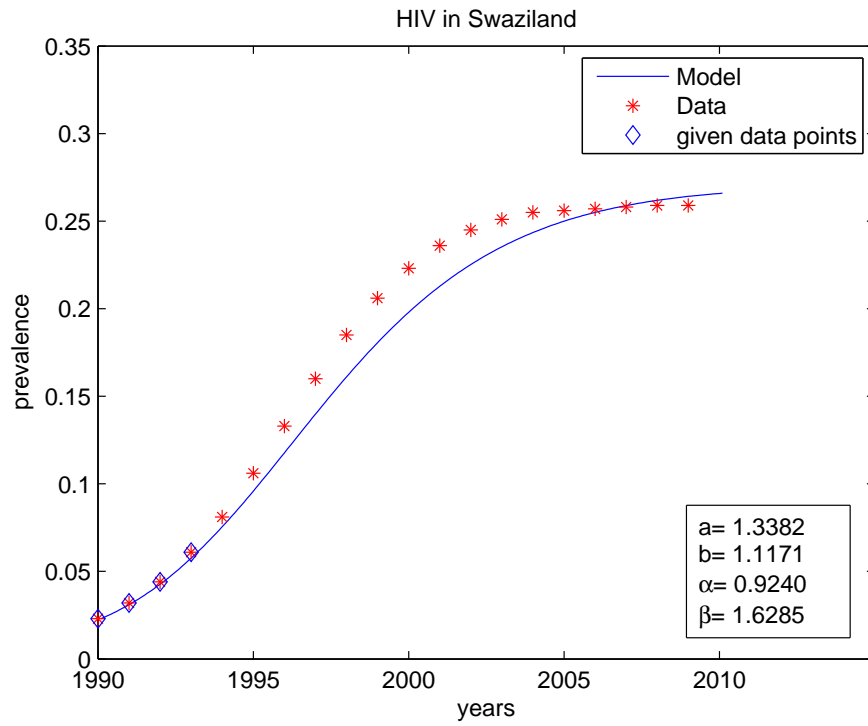
Figure 2.12: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for the transmission function $\rho_1 = ay^\alpha(1 - by)^\beta$.
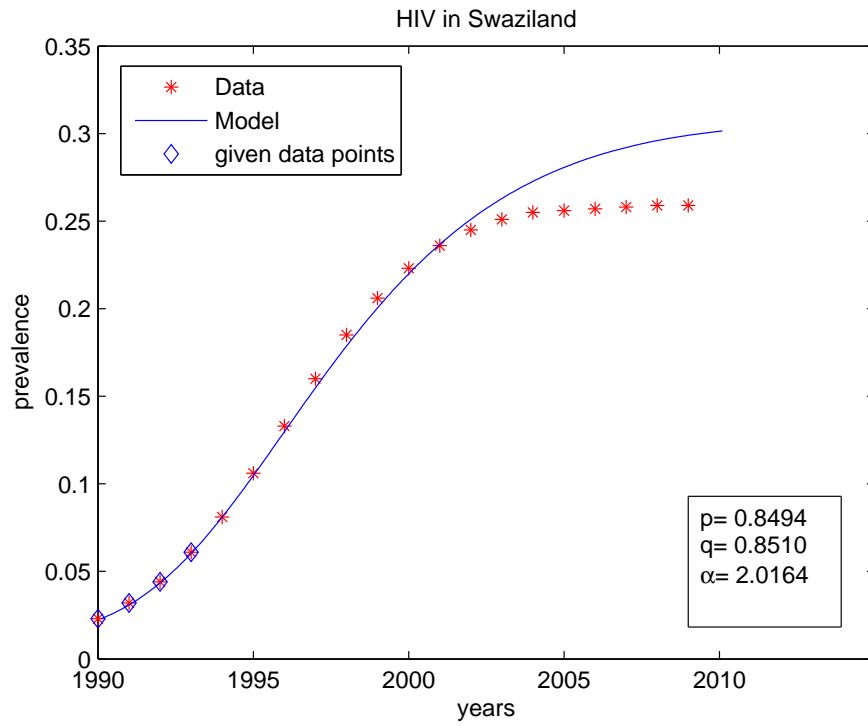
Figure 2.13: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for the transmission function $\rho_2 = \frac{I^p}{1+\alpha I^q}$.
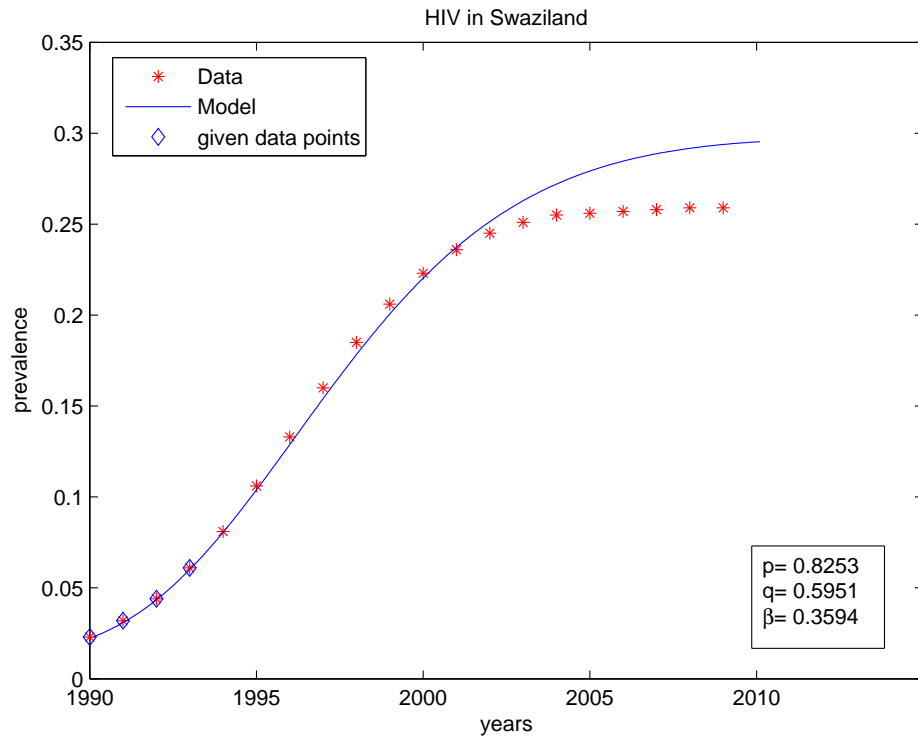
Figure 2.14: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for the power relationsship $\rho_3 = \beta S^p I^I$.
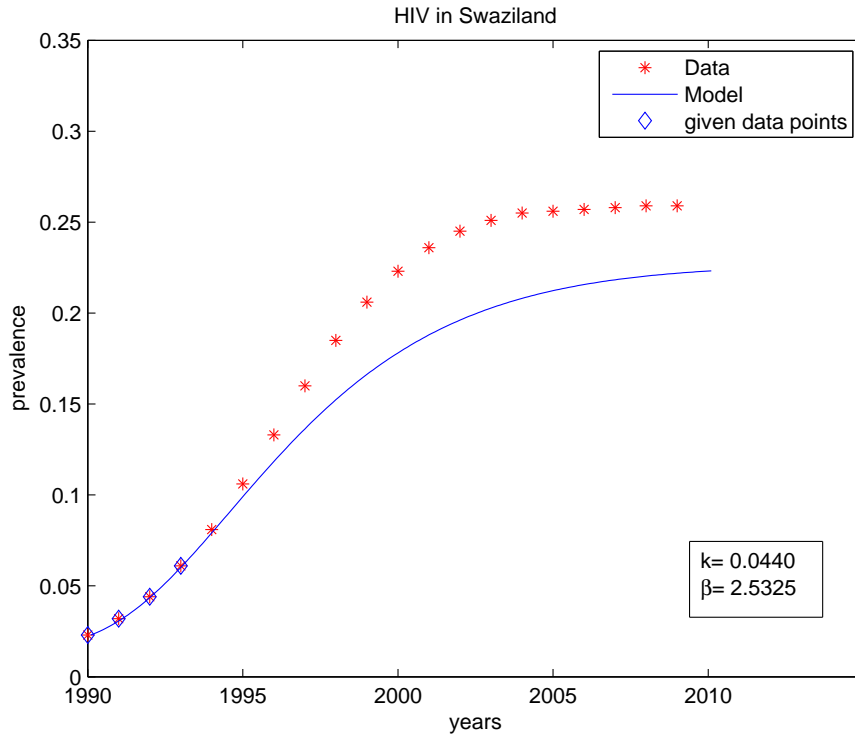
Figure 2.15: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for negative binomial $\rho_4 = kS \ln(1 + \frac{\beta I}{k})$

In the Figures 2.12-2.15 we see that $\rho_2$ and $\rho_3$ predict the used data noticably good until the year 2011, but then they are not able to further give an appropriate approximation with only 4 given points . The function $\rho_1$ does not give such a good prediction in the first years as $\rho_2$ and $\rho_3$, but in summary, it nevertheless gives the best prediction on the whole time interval. The transmission function $\rho_4$ only gives a good fit until 1995 and afterwards it gives an unsatisfactory prediction of the true prevalence. The transmission function $\rho = ay^\alpha(1 - by)^\beta$ gives the best fitting curve with 20 given dates, because it describes exactly the given prevalence, and for 4 given dates our function also predicts the given HIV data set the best, regarding the whole time interval.

Now we are going to use the HIV-time-series prevalence data of Botswana. As we have see before, Botswana's prevalence includes a maximal data point after which the prevalence is decreasing. Since we consider only the expansion phase of the disease, we use only 13 data points in the next considerations. Figures 2.16 to 2.19 show the fitting curve on Botswana's prevalence with 13 given data points for the four different transmission functions.

41

Figure 2.16: Prevalence data and calculated prevalence based on parameter identification using 13 data measurements for the transmission function $\rho_1 = ay^\alpha(1 - by)^\beta$.

Figure 2.17: Prevalence data and calculated prevalence based on parameter identification using 13 data measurements for the transmission function $\rho_2 = \frac{I^p}{(1+\alpha I^q)}$.

Figure 2.18: Prevalence data and calculated prevalence based on parameter identification using 13 data measurements for the transmission function $\rho_3 = \beta S^p I^q$.
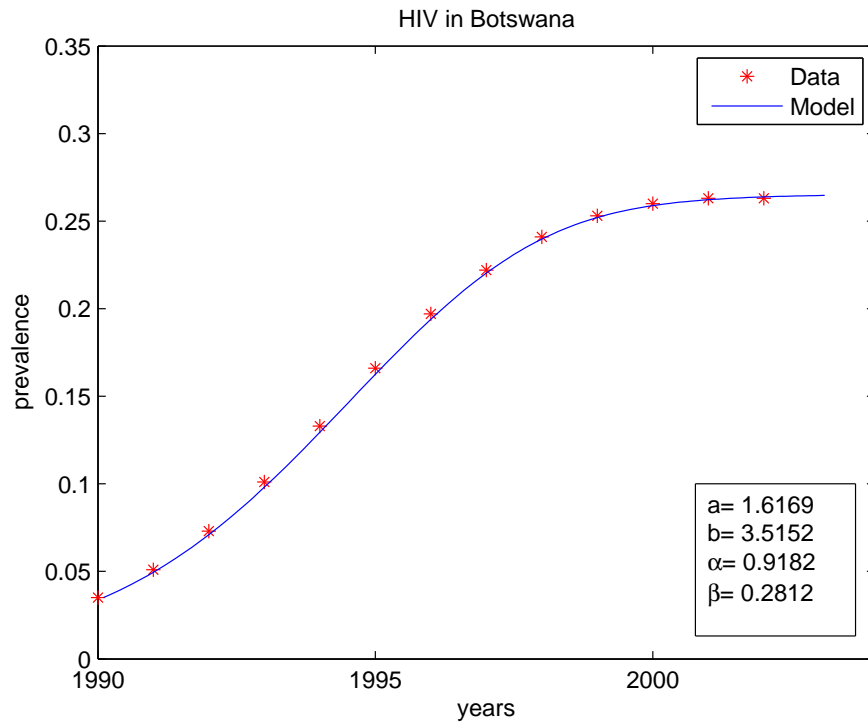
Figure 2.19: Prevalence data and calculated prevalence based on parameter identification using 13 data measurements for the transmission function $\rho_4 = kS\ln(1 + \frac{\beta I}{k})$.

Now we are interested in the prevalence's prediction with 4 points. In the same context as before, we are only interested in a fair comparison of the four transmission functions, so all of them get the same input data. These experimental plots are shown in Figures 2.20 to 2.23.

Figure 2.20: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for the transmission function $\rho_1 = ay^\alpha(1 - by)^\beta$.
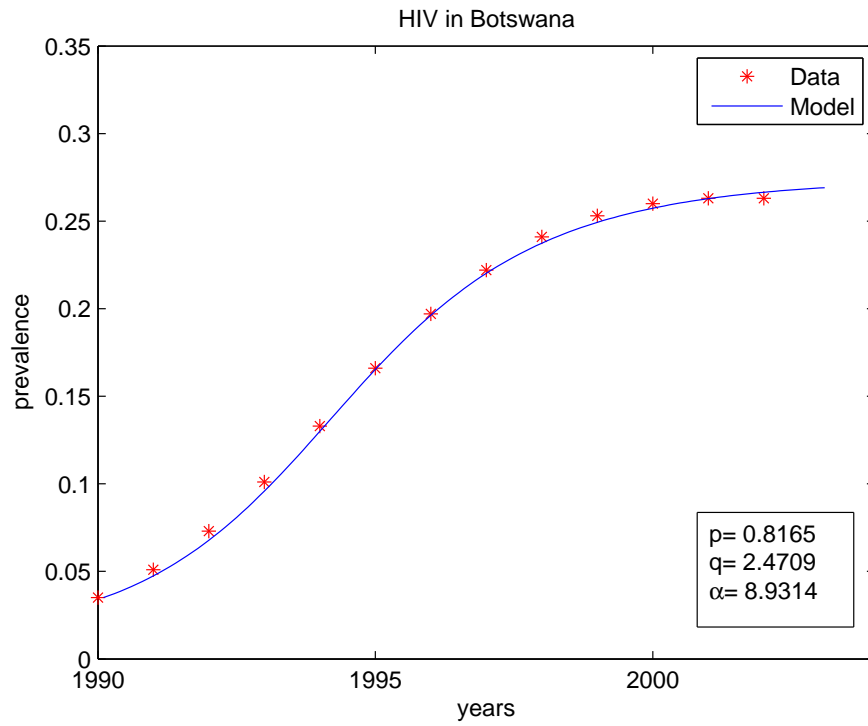
Figure 2.21: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for the transmission function $\rho_2 = \frac{I^p}{(1+\alpha I^q)}$.
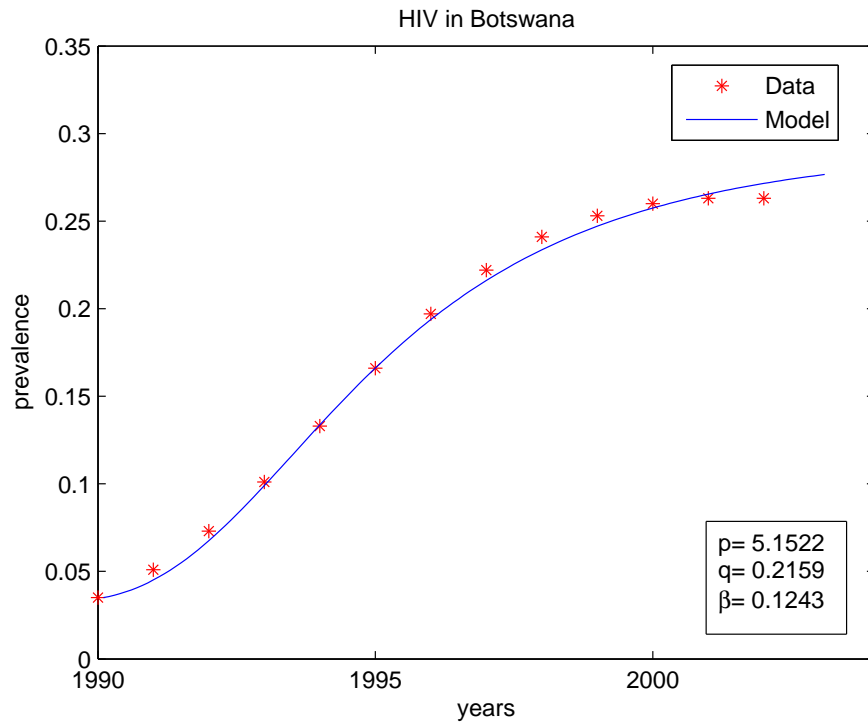
Figure 2.22: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for the transmission function $\rho_3 = \beta S^p I^q$.
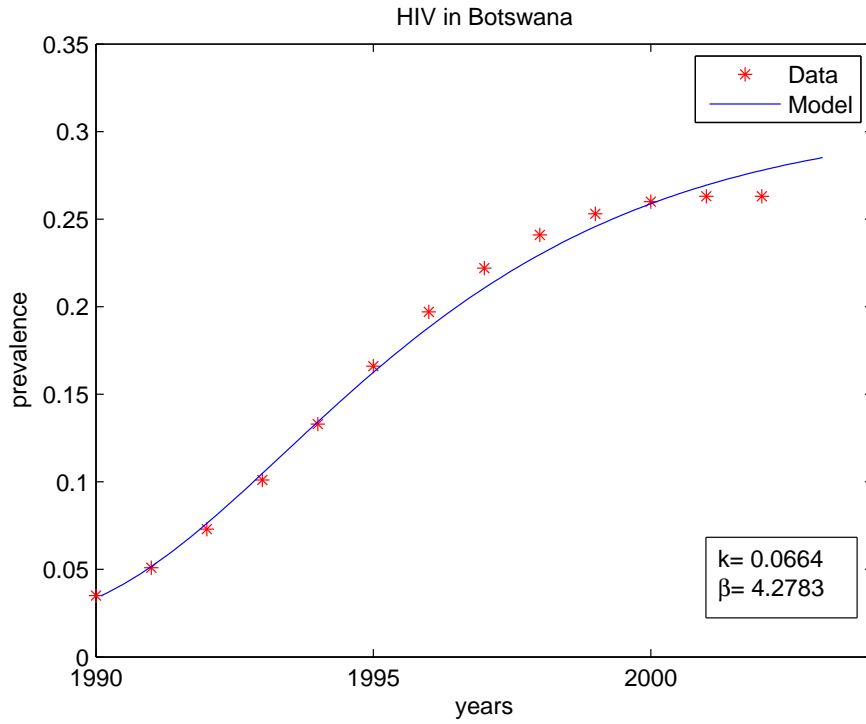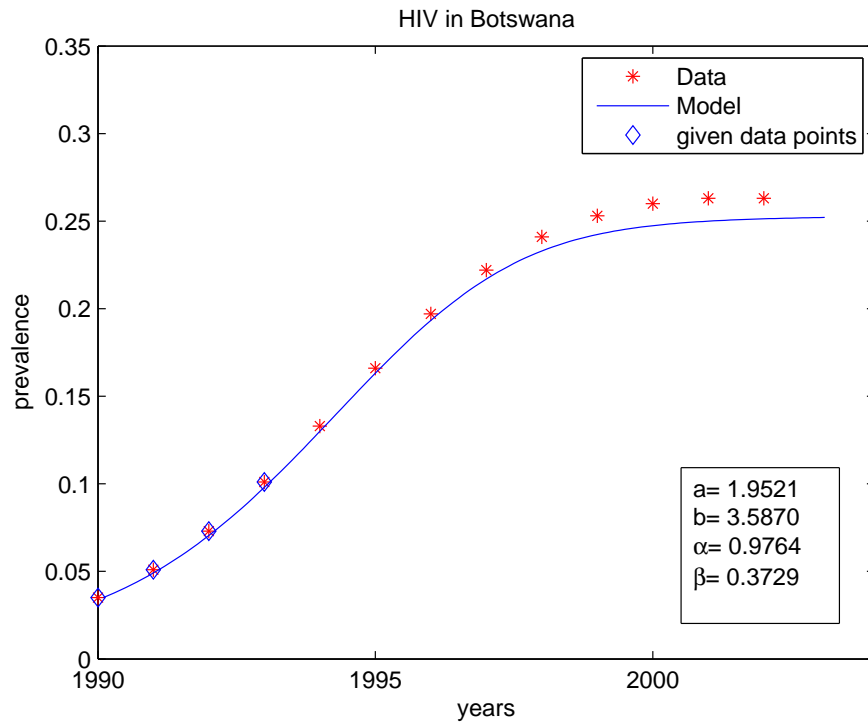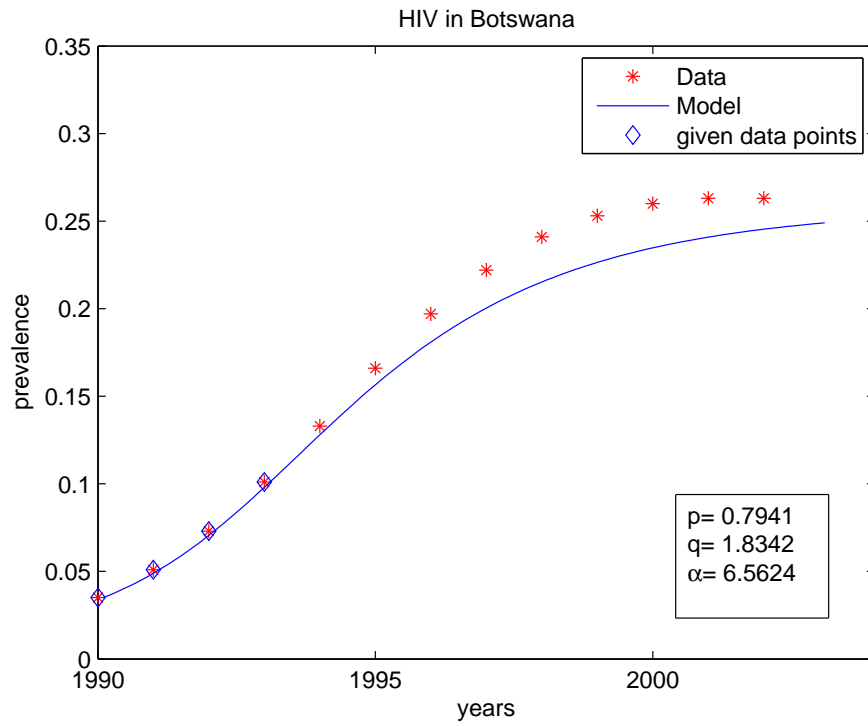
Figure 2.23: Prevalence data and calculated prevalence based on parameter identification using 4 data measurements for the transmission function $\rho_4 = kS\ln(1 + \frac{\beta I}{k})$.

The results of all the above experiments are sumarized in the following tables giving the values of the $L^1$-norm and the $L_\infty$-norm at the deviations of the results obtained by the fitted ODE model from the real data. We need the values for comparison of the four transmission functions, the norms are defined as

$$x_{(y,y^*)} = \left| P_{[data]}(i) - P_{[a,b,\alpha,\beta]}(t_i) \right|$$

$$l_1(x_{(y,y^*)}) = \left\| x_{(y,y^*)} \right\|_1 = \frac{1}{n-1} \sum_{i=1}^{n} |x_i|$$

$$l_\infty(x_{(y,y^*)}) = \left\| x_{(y,y^*)} \right\|_\infty = \max(|x_1|, |x_2|, \ldots + |x_n|)$$

where

$$y = P_{[data]}(i) \qquad y^* = P_{[a,b,\alpha,\beta]}(t_i)$$

The following table gives the values for Botswana with 13 used data points. All tables are calculated with MATLAB7.6.0 R2007b.

49

| used $\rho_i$ for $Botswana_{13}$ | $l_2(y, y^*)$ | $l_\infty(y, y^*)$ |
|---|---|---|
| $\rho_1$ | 0.0004 | 0.001 |
| $\rho_2$ | 0.0019 | 0.0038 |
| $\rho_3$ | 0.0037 | 0.0091 |
| $\rho_4$ | 0.0062 | 0.0156 |

As you can see in the table above, $\rho_1$ has the least error of all four fitting functions. Concluding, $\rho_1$ gives the best fit for Botswana with all 13 data points.

The next table gives the values for Botswana with 4 used data points for the parameter identification. So now we want to now wich function predicts the values in the best way for the future.

| used $\rho_i$ for $Botswana_4$ | $l_1(y, y^*)$ | $l_\infty(y, y^*)$ |
|---|---|---|
| $\rho_1$ | $9.6849 \cdot 10^{-6}$ | 0.0016 |
| $\rho_2$ | $1.2178 \cdot 10^{-4}$ | 0.0197 |
| $\rho_3$ | $1.633 \cdot 10^{-4}$ | 0.02348 |
| $\rho_4$ | $2.2641 \cdot 10^{-4}$ | 0.02496 |

We can see, that $\rho_1$ gives the best prediction for the future. So our transmission function gives the best fit of the four functions.

Now we are going to analyse the norm values for Swaziland with 20 used data points, in the same way as we did for Botswana with 13 measurements.

| used $\rho_i$ for $Swaziland_{20}$ | $l_1(y, y^*)$ | $l_\infty(y, y^*)$ |
|---|---|---|
| $\rho_1$ | 0.0003 | 0.0017 |
| $\rho_2$ | 0.0021 | 0.0040 |
| $\rho_3$ | 0.0022 | 0.0037 |
| $\rho_4$ | 0.0033 | 0.0063 |

The table above shows the norm values for Swaziland with all used data points, as we saw in the case of Botswana, here gives $\rho_1$ gives the best fit here, too. However, we are more interested in the prediction of the future, so the last table gives the values for Swaziland with 4 used data points for the parameter identification.

| used $\rho_i$ for $Swaziland_4$ | $l_1(y, y^*)$ | $l_\infty(y, y^*)$ |
|:---:|:---:|:---:|
| $\rho_1$ | $6.0251 \cdot 10^{-5}$ | 0.0018 |
| $\rho_2$ | $2.0962 \cdot 10^{-4}$ | 0.0328 |
| $\rho_3$ | $2.4087 \cdot 10^{-4}$ | 0.0320 |
| $\rho_4$ | $3.2160 \cdot 10^{-3}$ | 0.0462 |

Concluding, we take a look at the tables above, so we can declare that the transmission function $\rho_1$ predicts the true prevalence in the best way of all four functions. The fitting curve with all data points and also the curve for 4 data points with $\rho_1$ gave in both cases, Botswana and Swaziland, the best fit of the four used transmission functions. So the main statement of this chapter is, that our transmission function $\rho_1$ is the best fitting function of those four and predicts the values for the future in the best way.

In the next Section we are looking at Zimbabwe, a country that had a high prevalence in the 90s, then had a maximum prevalence, and after that it was monotonically decreasing until 2010. The question we ask now is of our transmission function $\rho$ is able to make a reasonbaly good fit for the given HIV prevalence time series.

## 2.5 A Country with Increasing and Decreasing Prevalence: Zimbabwe

In this section we are going to do on experiment with prevalences which are not monotonically increasing for the whole given time interval. Our selected country is also in Africa, namely Zimbabwe. Firstly, we are going to see that our transmission function is going to fail on the whole estimation.

Now, as in Section 2.3, we are going to estimate the free parameters of the transmission function $\rho = ay^\alpha(1 - by)^\beta$. The parameters $\alpha$, $\beta$, a, b are estimated in the same way as it is described in Section 2.1. However, we did this estimation for two developing countries, Swaziland and Botswana. These are two of those countries with the highest HIV prevalence. The time series data of Swaziland and Botswana we needed for the estimation were taken, as in Section 2.2, from `www.who.int`. This time series data also lie between 1990 and 2010.

We will start with the estimation of the four free parameters for Zimbabwe. In order to estimate the unknown parameters a, b, $\alpha$ and $\beta$ for the transmission function $\rho$ we minimize the square error:

$$\sum_{t=1990}^{2010} (P_{[data]}(t) - P_{[a,b,\alpha,\beta]}(t_1))^2 \tag{2.6}$$

Our given parameter values for Zimbabwe are

- $\sigma = 0.27$ ,

- $\lambda = -0.1358$,

- $\gamma = 0$,

- $\delta = 0.24$,

- $\kappa = 0.3265$.

After calculating the minimal square error for all 20 given data points, the minimal error was achieved for a= 0.1231, b= 0.0725, $\alpha$= -0.9176 and $\beta$=-0.0172. Figure 2.24 shows the resulting fit to the data.



Figure 2.24: This figure shows the expansion phase for the prevalence data (red stars) and the calculated prevalence based on parameter identification using 20 measurements (blue line).

As a result, which we have already known before this experiment, $\rho$ does not give a good fitting curve of non-monotonically increasing prevalences. It is not possible to balance the values of the free parameters, so not even the increasing phase is well fitted. However, if we only use the data points from the expansion phase, we get Figure 2.25. This figure shows the model which only used the first 13 data points from the time series, i.e. those points until the maximum prevalence value of Zimbabwe. Now the minimal error was

achieved for a= 3.5682, b= 3.5600, $\alpha$= 0.9929 and $\beta$= 0.1221.
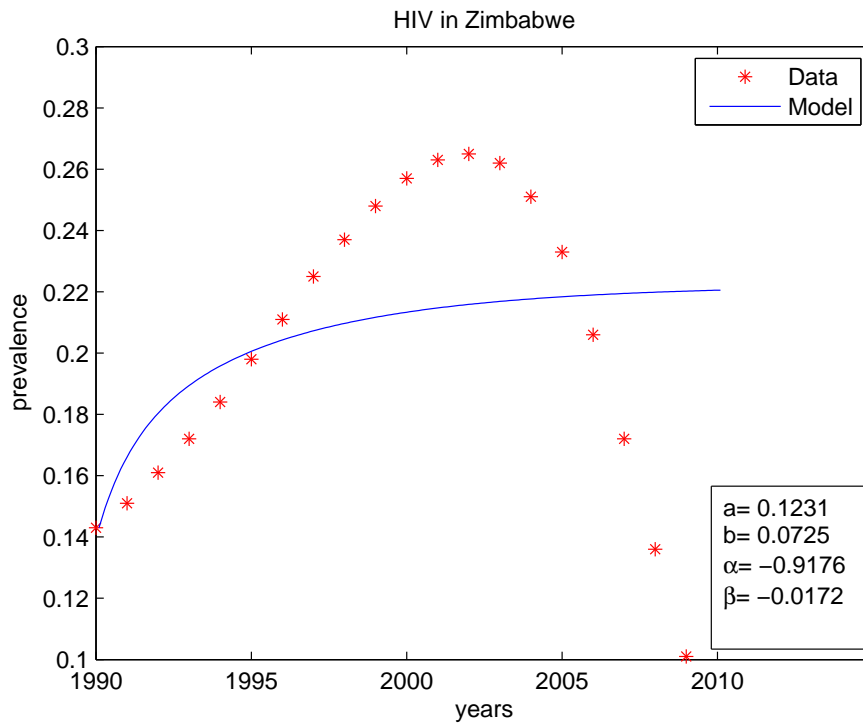


Figure 2.25: This figure shows the expansion phase for the prevalence data and the calculated prevalence based on parameter identification using only 13 measurements, which reflect the expansion phase of Zimbabwe's prevalence.

Concluding, we see that the fitting curve is very good for the increasing phase of the prevalence, especially we only use the data points of until the maximum for the parameter identifaction. Further resarch is needed to obtain a transmission function $\rho$ from analyses of a model taking into account the heterogeneity of the population (as in Chapter 1) that will relevant fot the descending phase of the disease.

# Chapter 3

# The Optimal Control

## 3.1 Analysis of the primitive control of the mortality rate

In this section we analyze the short-term and the long-term effects of the mortality rate $\delta$ depending on the SI-Model (1.27)-(1.29) with transmission function $\rho = a y^\alpha (1 - by)^\beta$ in the case of Swaziland. This means we have for the fixed parameters the following values:

- $\sigma = 0.43$ ,

- $\lambda = $ -0.153,

- $\gamma = 0$,

- $\delta = 0.259$,

- $\kappa = 0.2825$.

The first question we want to answer is which effects a smaller mortality rate $\delta_1 < \delta$ has depending on I. Therefore, we choose $\delta_1 = 0.2 < \delta = 0.259$, in Figure 3.1 we see the effects of this minor modification.

Figure 3.1: Short and long run effect of a small change of the mortality rate.

We see in Figure 3.1 that for the short run the smaller mortality rate $\delta_1$ gives a small set of infected, but after 13 years there exists an intersection between the set of I with a $\delta = 0.259$ and the smaller one. So after 13 years and also in the long run, the set depending on $\delta_1$ is higher than those of $\delta$. This fact represents that if we minimize the mortality rate in such a country like Swaziland it is counterproductive in the long run, if the people in Swaziland will not change there sexual behaviour, concluding they will infect more and more people with HIV in the long term. Of course, in the short run it gives a smaller set of infected than a higher mortality rate, but we are interested in the prediction of HIV and in the minimization of the virus, so a smaller $\delta$ gives the converse effect of what we are searching for.

We also take a short look at the set of the susceptible and the set of the susceptible plus the set of the infected. So, Figure 3.2 shows the set of susceptible with $\delta$ and $\delta_1$.

Figure 3.2: Set of susceptible with a change of the mortality rate.

In the Figure 3.2 above, we see that there is no significant difference between the red and the green line in the short run, but after 5 years we see that the set of susceptible with the smaller $\delta_1$ has a higher number of susceptible than those with the original $\delta$. We also see that the red line is only decreasing after 5 years and seems to be constant after 25 years. In comparison, the green line is decreasing after 5 years and gets increasing after 7 more years. At last we take a short look at Figure 3.3.

Set of Susceptible plus Infected in Swaziland

Figure 3.3: Set of susceptible plus set of infected with a change of the mortality rate.

Figure 3.3 just shows that the summation of the two sets gives a higher number of individuals if the mortality rate is smaller, which is quite logical. In comparison, the most interesting fact of this section is Figure 3.1, because it shows that the long run behaviour of a smaller $\delta_1$ is counterproductive for the whole population and only results in more people getting infected with HIV.

In the last figure, Figure 3.4, the comparison of the prevalences with different values for $\delta$ is shown.

Figure 3.4: Prevalence data, calculated from the SI Model (1.27)-(1.29) with a change of the mortality rate.

In Figure 3.4 is shown that in the first 25 years of the disease, the prevalence with $\delta = 0.259$ is a bit higher than those with $\delta_1 = 0.2$, but between 25 and 26 years there exists an intersection between the two prevalence functions.

To see the intersection in more detail, Figure 3.5 shows Figure 3.4 with $x_{min} = 20$ years, $x_{max} = 30$ years, $y_{min} = 0.256$ and $y_{max} = 0.261$.

Figure 3.5: Prevalence data, calculated from the SI Model (1.27)-(1.29) with a change of the mortality rate, in more detail.

In the next sections we are going to define and analyse the basics of optimal control theory, which is relevant after such a primitve control of the mortality rate. Finally we want to know if we are able to finde controls $u_1$ and $u_2$, which represent the effect of prevention and the effect of treatment, to minimize the number of people who are infected with HIV.

## 3.2 Controlled ODEs

Controlled ordinary differential equations (ODEs) are defned by:

$$\dot{x}(t) = f(t, x(t), u(t)), \qquad x(0) = x^0 \in \mathbb{R}^n, \qquad t \in [0, T], \qquad (3.1)$$

with the terminal constraint

$$x(t) \in M. \tag{3.2}$$

where $M \subset \mathbb{R}^n$ is a given closed target set.

Admissible controls are in $\mathcal{U}[0,T] = \mathcal{U}^{pc}[0,T]$ or $\mathcal{U}[0,T] = \mathcal{U}^{L}[0,T]$, $U \subset \mathbb{R}^m$. Notice that:

1. $\mathcal{U}^{pc}[0,T] = \{u : u \text{ is piecewise continuous and } u(t) \in U \text{ for all } t \in [0,T]\}$

2. $\mathcal{U}^{L}[0,T] = \{u \in L_1[0,T] : u(t) \in U \text{ for almost every } t \in [0,T]\}$

Thus, we have to assume that $f : [0,T] \times \mathbb{R}^n \times U \mapsto \mathbb{R}^n$ and $h : [0,T] \times \mathbb{R}^n \times U \mapsto \mathbb{R}$ are measurable and bounded in t, continuously differentiable in x and continuous in u. Let us define an $x : [0,T] \mapsto \mathbb{R}^n$, which is called a solution of $\dot{x}(t)$, where x is continuous and

$$x(t) = x^0 + \int_0^t f(s, x(s), u(s)) ds \qquad \forall t \in (0,T].$$

If $u \in \mathcal{U}^{pc}[0,T]$, then the solution x, if it exists, is differentiable at the points of continuity of $f(t, x(t), u(t))$ and the derivative $\dot{x}$ is piecewise continuous. If $u \in \mathcal{U}^{L}[0,T]$, then x is absolutely continuous. In both cases (3.1) is satisfied for almost every $t \in (0,T)$. We notice that x is absolutely continuous, by definition, iff it is a.e. differentiable and $x(t) = x(0) + \int_0^t \dot{x}(s) ds$.

The definition concerning the solution x, which is mentioned above, is based on the following theorem.

**Theorem 3.2.1** *(Caratheodory) Based on the assumptions made above, if u is measurable, then $f(t, x(t), u(t))$ is measurable in t.*

We also assume that for every $u \in \mathcal{U}[0,T]$ a unique solution, $x[u]$, of (3.1) exists on $[0.T]$. A solution may fail to exist only if $|x(t)|$ escapes to infinity for $t \leq T$.

Now we want to give important sufficient optimality conditions, the so called Letimann-Stalford sufficient optimality conditions. Consider the problem

$$\max \int_0^T h(t, x(t), u(t)) dt \quad \text{subject to} \quad (3.1) \quad \text{and} \quad u(t) \in U$$

and under the additional condition that $x(T) \in M$, where $M \subset \mathbb{R}^n$ is a given closed target set. Equivalently,

$$\max_{u \in \mathcal{U}} \int_0^T h(t, x[u](t), u(t)) dt \qquad (3.3)$$

under the condition that $x[u](T) \in M$.

Notice that h has the same properties as f above, while $U \subset \mathbb{R}^m$ and $M \subset \mathbb{R}^n$ are arbitrary. Let us define the Hamiltonian:

60

for $\lambda \in \mathbb{R}^n$

$$H(t, x, u, \lambda) = h(t, x, u) + \langle \lambda, f(t, x, u) \rangle,$$

the meaning as which will be explained in the next Section.

**Theorem 3.2.2** *Let $u^* \in \mathcal{U}[0, T]$ and $x^* = x[u^*]$. Let $\lambda : [0, T] \mapsto \mathbb{R}^n$ be absolutely continuous. Assume that the following two conditions are satisfied:*

*1. $\forall u \in \mathcal{U}$ and a.e. $t \in [0, T]$*

$$H(t, x^*, u^*, \lambda(t)) - H(t, x[u](t), u(t), \lambda(t)) + \left\langle \dot{\lambda}, x^* - x[u](t) \right\rangle \leq 0;$$

*2. for every $x \in M$*

$$\langle \lambda(T), x^*(t) - x \rangle \leq 0.$$

*Then $(u^*, x^*)$ is an optimal solution.*

Now for Arrow-type sufficient optimality conditions we denote the maximized Hamiltonian by $\mathcal{H}$:

$$\mathcal{H}(t, x^*(t), \lambda(t)) := \max_{u \in U} H(t, x^*(t), u, \lambda(t)).$$

**Theorem 3.2.3** *Assume that for every $\lambda \in \mathbb{R}^n$ the function $\mathcal{H}(t, \cdot, \lambda)$ is concave and continuously differentiable. Let for $u^* \in \mathcal{U}[0, T]$ and $x^* = x[u^*]$, then there exists an absolutely continuous function $\lambda : [0, T] \mapsto \mathbb{R}^n$ such that for a.e. $t \in [0, T]$*

$$\dot{\lambda}(t) = -H_x(t, x^*(t), u^*(t), \lambda(t)), \qquad \lambda(T) \in -N_M^{\#}(x^*(T))$$
$$H(t, x^*(t), u^*(t), \lambda(t))) = \max_{u \in U} H(t, x^*(t), u, \lambda(t)).$$

*Then the pair $(u^*, x^*)$ is an optimal solution of problem (3.1) and (3.3).*

### 3.2.1 Pontryagin's Maximum Principle

Now, we are able to discuss Pontryagin's Maximum Principle, which is used in optimal control theory to find the best possible control for taking a dynamical system from one state to another, especially in the presence of constraints for the state or input controls. It was formulated in 1956 by the Russian mathematician Lev Semenovich Pontryagin and his students.

On a fixed interval $t \in [0, T]$, consider the Bolza problem

$$\max_{u\in\mathcal{U}[0,T]} \left\{ g(x(t)) + \int_0^T h(t,x(t),u(t))dt \right\} \tag{3.4}$$

The function $g : \mathbb{R}^n \mapsto \mathbb{R}$ is continuously differentiable. The sets U and M are closed. For every $u \in \mathcal{U}$ the solution $x[u]$ exist on $[0,T]$.

To solve such a problem you need the *Hamiltonian*, so for $\lambda \in \mathbb{R}^n$ define

$$\mathcal{H}(t,x,u,\lambda,\lambda_0) = \lambda_0 h(t,x,u) + \langle \lambda, f(t,x,u) \rangle .$$

**Theorem 3.2.4** *Assume that $u^* \in \mathcal{U}$ is an optimal control and let $x^* = x[u^*]$. Furthermore assume that the functions $f(t,x^*(t),u^*(t)), h(t,x^*(t),u^*(t))$ and the corresponding derivatives $f'_x, h'_x$, are bounded. Then there exist a $\lambda_0 \in (0,1)$ and $\xi \in N_M(x^*(T))$ with $\lambda_0 + |\xi| > 0$, such that the **adjoint equation***

$$\dot{\lambda} = -\mathcal{H}_x(t,x^*(t),u^*(t),\lambda) \tag{3.5}$$
$$\lambda(T) = \alpha g'(x^*(T)) + \xi \tag{3.6}$$

*has a unique solution $\lambda$ on $[0,T]$, and for a.e. $t \in (0,T)$*

$$\mathcal{H}(t,x^*(t),u^*(t),\lambda(t),\lambda_0) = \max_{u\in U} \mathcal{H}(t,x^*(t),u,\lambda(t),\lambda).$$

The meaning of the set $N_M(x)$ is given by the *transversality condition*:
For a general closed set $M \subset \mathbb{R}^n$ and $x \in M$ the appropriate definition of $N_M(x)$ is given by F. Clarke. The maximum principle remains true.

1. For a convex M: $N_M(x) = I \in \mathbb{R}^n : \langle I, y - x \rangle \leq 0 \forall y \in M$. In particular, if $M = \mathbb{R}^n$, then $N_M(x) = 0$, hence the transversality condition becomes the one we know:

$$\lambda(T) = \lambda_0 g'(x^*(T)).$$

   In this case, the theorem holds with $\lambda_0 = 1$.

2. If $M = x \in \mathbb{R}^n : \varphi(x) = 0$, where $\varphi : \mathbb{R}^n \mapsto \mathbb{R}$ with $r \leq n$, then

$$N_M(x) = \varphi_x(x)p : p \in \mathbb{R}^r.$$

   The transversality condition becomes

$$\lambda(T) = \lambda_0 g'(x^*(T)) + \varphi_x(x^*(T))p$$

   for some undetermined $p \in \mathbb{R}^r$. These r unknowns should be determined using the r equations $\varphi(x^*(T)) = 0$.

3. If $M = x \in \mathbb{R}^n : \varphi(x) \geq 0$, then

$$N_M(x) = - \left\{ \sum_{i \in I(x)} \varphi_x^i(h) p_i : p_i \geq 0, i \in I(x) \right\},$$

where $I(x) = i \in i, \ldots, r : \varphi^i(x) = 0$. The transversality condition becomes

$$\lambda(T) = \lambda_0 g'(x^*(T)) + \sum_{i \in I(x^*(T))} \varphi_x^i(x^*(T)) p_i.$$

The unknown $p_i, i \in I(x^*(T))$ should be determined from the equations $\varphi^i(x^*(T)) = 0, i \in I(x^*(T))$.

The last thing we need is the definition of the **canonical system**. So since $u^*$ and $x^*$ are not known in advance, the maximum principle can be viewed as a system of equations for x, u and $\lambda$:

$$\dot{x} = f(t, x, u(t)),$$
$$\dot{\lambda} = -(f_x'(t, x, u(t)))^T \lambda - h_x'(t, x, u(t))$$
$$\mathcal{H}(t, x(t), u(t), \lambda(t)) = \max_{u \in U} \mathcal{H}(t, x(t), u, \lambda(t)),$$

with the additional boundary condition for $M = (R)^m$ $x(0) = x^0$ and $\lambda(T) = g'(x(T))$. Shortly you get

$$\dot{x} = \mathcal{H}_\lambda(t, x, u(t), \lambda), \qquad\qquad x(0) = x^0 \qquad (3.7)$$
$$\dot{\lambda} = -\mathcal{H}_x(t, x, u(t), \lambda), \qquad\qquad \lambda(T) = g'(x(T)), \qquad (3.8)$$
$$\mathcal{H}(t, x(t), u(t), \lambda(t)) = \max_{u \in U} \mathcal{H}(t, x(t), u(t), \lambda(t)). \qquad\qquad (3.9)$$

This is a full system for determinung u, x, and $\lambda$.

If we assume that for any given vectors x and $\lambda$ in $\mathbb{R}^n$ the Pontryagin maximization problem 3.2.6, arising from (3.9).

$$\max_{u \in U} \mathcal{H}(t, x, u, \lambda) \qquad\qquad (3.10)$$

has a unique solution, denoted by $\hat{u}(t, x, \lambda)$, that is

$$\mathcal{H}(t, x, \hat{u}(t, x, \lambda), \lambda) = \max_{u \in U} \mathcal{H}(t, x, u, \lambda), \qquad\qquad (3.11)$$

then (3.7) and (3.8) become what is called the **canonical system**:

$$\dot{x} = \mathcal{H}_\lambda(t, x, \hat{u}(t, x, \lambda), \lambda), \tag{3.12}$$

$$\dot{\lambda} = -\mathcal{H}_x(t, x, \hat{u}(t, x, \lambda), \lambda) \tag{3.13}$$

with the boundary conditions

$$x(0) = x^0, \qquad \lambda(T) = g^{'}(x(T)). \tag{3.14}$$

This is a boundary value problem for a system of ODEs. It has to be mentioned that $\hat{u}(t, x, \lambda) = \hat{u}(x, \lambda)$ is independent of t, if all data are independent of t. In general, the canonical system does not have a unique soltuion, as it happens for many economic models.

**Theorem 3.2.5** *(unique solution) Assume that the considered optimal control problem has a solution. Let the Pontryagin system (3.7)-(3.9) has a unique solution $(u^*, x^*, \lambda)$. Then $u^*$ is the unique optimal control.*

Let a unique maximizer $\hat{u}(t, x, \lambda)$ of the Hamiltonian (3.10) exist. On the assumption of the last theorem, if the canonical system with boundary condition (3.14) has a unique solution $(x^*, \lambda)$, then $u^*(t) := \hat{u}(t, x^*(t), \lambda(t))$ is the unique optimal solution.

So we can briefly formulate the Pontryagin's Maximum Principle as:

**Proposition 3.2.6** *(Pontryagin's Maximum Principle)*

$$\mathcal{H}(t, x^*(t), u^*(t), \lambda(t)) = \max_{u \in U} \mathcal{H}(t, x^*(t), u, \lambda(t)) \tag{3.15}$$

*and at every point t where $u(\cdot)$ is continuous*

$$\dot{\lambda} = -\mathcal{H}_x(t, x^*(t), u^*(t)) \tag{3.16}$$

*Furthermore, the transversality condition*

$$\lambda(T) = g_x(x^*(T), T) \tag{3.17}$$

*holds, where*

$$\mathcal{H}(t, x, u, \lambda) := g(t, x, u) + \lambda f(t, x, u). \tag{3.18}$$

## 3.3 Existence of a Treshold Line for the Optimal Treatment in a Simple Epidemic Model

This Section is explicitly based on notes from Prof. Vladimir Veliov, who was so supportive as to provide them to me for writing them in this work.

Consider the following model with infinite horizon of infectious disease:

$$\max \int_0^{+\infty} e^{-rt}[S(t) + \alpha I(t)]dt$$

$$\dot{S} = -\sigma \frac{SI}{S+I} - \mu S, \qquad (3.19)$$

$$\dot{I} = \sigma \frac{SI}{S+I} - u(t)I, \qquad (3.20)$$

$$u(t) \in [\mu_0, \mu_1] \qquad (3.21)$$

Here

- S - susceptible individuals;

- I - infected individuals;

- $\mu$ - net mortality rate of the non-infected individuals;

- $\mu_1$ - mortality rate of the non-medicated infected individuals;

- $\mu_0$ - mortality rate of the fully-medicated infected individuals;

- u - medication control; u = $\mu_1$ - no medication, u = $\mu_0$ - full medication, $0 \le \mu < \mu_0 < \mu_1$;

- $\sigma$ - strength of infection;

- $\alpha$ - relative productivity of the infected individuals, $\alpha \in [0, 1]$.

We suppose that the following inequalities related to the reproduction number of the disease hold:

$$\frac{\sigma + \mu}{\mu_1} < 1 \quad \text{and} \quad \frac{\sigma + \mu}{\mu_0} > 1. \qquad (3.22)$$

We suppose also that

$$r > 0, \quad \mu_0 < \sigma - r, \quad 2\mu_0 > \sigma + \mu - r. \qquad (3.23)$$

To avoid infinitely expanding population we suppose below that $\mu \geq 0$. The case $\mu = 0$, in which the population would be in equilibrium if the disease was not present, is especially interesting since the system with $u = \mu_1$ has no stable equilibrium, as we shall see later on.

We mention that in this problem there exists an optimal solution since the equations depend linearly on u, the objective function is independent of (thus convex with respect to) u and $[\mu_0, \mu_1]$ is convex.

For a fixed constant control value u and initial data $S_0, I_0$ the solution of the system (3.19)-(3.20) is given by

$$S(t) = e^{-rt}(p_0 e^{(\sigma + \mu - u)t} + 1 - p_0)^{\frac{\sigma}{u - \sigma - \mu}} S_0,$$

$$I(t) = e^{(\sigma - u)t}(p_0 e^{(\sigma + \mu - u)t} + 1 - p_0)^{\frac{\sigma}{u - \sigma - \mu}} I_0,$$

where $p_0 = \frac{I_0}{S_0 + I_0}$. The formula for the prevalence $p(t) = \frac{I(t)}{S(t) + I(t)}$ is

$$p(t) = \frac{p_0}{p_0 + (1 - p_0)e^{(u - \sigma - \mu)t}}.$$

We can verify this by substituting the expression in the equations. The derivation is also not difficult. First we can show that the prevalence satisfies the equation

$$\dot{p} = (\sigma + \mu - u)p(1 - p), \quad p(0) = p_0. \tag{3.24}$$

This is a Ricati equation that can be solved by changing the variable $p = \frac{1}{z}$. Once the prevalence p(t) is obtained, equation (3.19) becomes a linear homogeneous (non stationary) equation that can be solved explicity to obtain S.

From (3.20) we conclude that $\dot{I}(t) \geq -\mu I(t)$, therefore $I(t) \geq e^{-\mu t} I_0$. In particular, if $I(0) > 0$, then $I(t) > 0, \forall t$. Additionally, from the formula for S(t) we know that $S(t) > 0$ if $S(0) > 0$. Further, we shall consider only nonzero initial points $S_0$ and $I_0$.

We stress the next properties that follow from the above formulas and (3.23):

(P1) p(t) is strictly increasing for $u = \mu_0$ and strictly decreasing for $u = \mu_1$.

(P2) The trajectory S(t) for $u = \mu_1$ converges to zero with $t \to +\infty$ if $\mu > 0$, but if $\mu = 0$, then
$$\lim_{t \to +\infty} S(t) = (1 - p_0)^{\frac{\sigma}{\mu_1 - \sigma}} S_0, \quad \lim_{t \to +\infty} I(t) = 0.$$

Moreover, from the first order homogeneity of the problem we easily establish the following:

(P3) If $u(\cdot)$ is an optimal control for initial point $(S_0, I_0)$, then it is an optimal control also for any other initial point $(S_1, I_1)$ with the same initial prevalence, that is, with $\frac{I_1}{I_1 + S_1} = \frac{I_0}{I_0 + S_0}$.

Now we shall prove the following property:

**Claim 3.3.1** *Either there exists an optimal trajectory along which the prevalence is constant, or for any optimal trajectory the corresponding prevalence is strictly monotone.*

To prove this we notice first that the set of trajectories $(S(\cdot), I(\cdot))$ starting from a given initial compact set is compact in the topology of uniform convergence on every compact interval. This follows from the uniform boundedness of the trajectories, the linearity of the equations with respect to the control, and the convexity and the compactness of the control set $[\mu_0, \mu_1]$. Moreover, the limit of optimal trajectories is an optimal trajectory, since the objective function is independent [1] of u.

Take an arbitrary optimal solution $(S(\cdot), I(\cdot))$ and consider the corresponding prevalence $p(\cdot)$. If the claim is false for this solution there exist $t_1 < t_2$ such that $p(t_1) = p(t_2)$. We consider two cases:

(i) There is no nondegenerated subinterval of $[t_1, t_2]$ on which p is constant. Then p is not constant also in $[t_1, t_2]$. Suppose that $p(t) > p(t_1)$ for some $t \in [t_1, t_2]$ (the alternative case is analogous). Then p achieves its maximum $\hat{p}$ on $[t_1, t_2]$. From the supposition of this case and from elementary topological properties of R it follows that every $\epsilon > 0$ there is some nondegenerate $[\tau_1^\epsilon, \tau_2^\epsilon] \subset [t_1, t_2]$ such that $p(\tau_1^\epsilon) = p(\tau_2^\epsilon) = p^\epsilon$, $\tau_2^\epsilon - \tau_1^\epsilon \leq \epsilon$, and $|p^\epsilon - p| \leq \epsilon$. Then the dynamic programming principle and property (P3) imply that for the initial points with prevalence $p^\epsilon$ there is an optimal trajectory for which the prevalence $p^\epsilon(\cdot)$ equals $p^\epsilon$ at least once in every interval of length $\epsilon$. From the compactness of the set of trajectories we conclude that for initial prevalence $\hat{p}$ there is an optimal trajectory along which the prevalence is constant.

(ii) There is a nondegenerate interval on which $p(t) = \hat{p}$ is constant. In this case, the dynamic programming principle implies that for initial points with prevalence p there is an optimal solution with constant prevalence. Thus, the claim is also true.

The maximum principle, in current values terms, claims that for any optimal control $\hat{u}(\cdot)$ the adjoint sytem

$$\dot{\xi} = r\xi + \sigma p^2 (\xi - \eta) + \mu \xi - 1,$$
$$\dot{\eta} = r\eta + \sigma(1-p)^2(\xi - \eta) + u\eta - \alpha,$$

---

[1] In fact, the objective function may depend on u in a concave way

has a solution $(\xi, \eta)$ such that $\hat{u}(t)$ satisfies, after cancellation of $I(t) > 0$,

$$\eta(t)\hat{u}(t) = \min_{u \in [\mu_0, \mu]} \eta(t)u \tag{3.25}$$

for a.e. $t > 0$,

It is remarkable that the adjoint equations depend on the state only through p, therefore we come up with a 3-dimensional primal/dual sytem

$$\dot{p} = (\sigma + \mu - u)p(1 - p), \tag{3.26}$$
$$\dot{\xi} = r\xi + \sigma p^2(\xi - \eta) + \mu\xi - 1, \tag{3.27}$$
$$\dot{\eta} = r\eta + \sigma(1 - p)^2(\xi - \eta) + u\eta - \alpha. \tag{3.28}$$

Now we shall investigate the possibility of singular solutions. Let us take an arbitrary optimal solution and consider the corresponding p, $\xi$ and $\eta$ for which (3.25)-(3.28) hold almost everywhere. Since $u(t) = \mu_0$ whenever $\eta(t) > 0$ and $u(t) = \mu_1$ whenever $\eta(t) < 0$, the functions $\xi$ and $\eta$ have Lipschitz continuous derivatives and satisfy (3.27)-(3.28) everywhere. We also mention that in view of (3.26) and (3.23) Claim 3.3.1 can be reformulated in the follwing way: for every optimal control u either the set of points where $u(t) = \mu_0$ or the set of points where $u(t) = \mu_1$ is of measure zero.

Suppose that the set $\Omega$ of all points t for which $\eta(t) = 0$ has a positive measure. Then $\dot{\eta}(t) = 0$ at every non-isolated points of $\Omega$ for which the derivative exists (and almost all points of $\Omega$ are such). Then from (3.28) we obtain that $\sigma(1 - p)^2\xi = \alpha$ a.e. on $\Omega$. Differentiating this equality and using (3.26) and (3.27) we have

$$2(\sigma + \mu - u)p\xi = r\xi + \sigma p^2\xi + \mu\xi - 1.$$

Excluding $\xi$ we obtain that the following relation between p and u holds almost everywhere in $\Omega$:

$$\sigma + \mu - u = \frac{1}{2p}(r + \sigma p^2 + \mu) - \frac{\sigma}{2\alpha p}(1 - p)^2. \tag{3.29}$$

We consider the following two cases.

*Case 1*
Suppose that $\eta(t) < 0$ on a set of positive measure. Then $u(t) = \mu_1$ on a set of positive measure, therefore $\eta(t) > 0$ never happens.

Let $\tau$ be a Lebesgue point of both the set $\Omega$ and the set $\Delta = [0, +\infty)$. We shall call such

68

points *contact points*. We shall estimate $\ddot{\eta}(t)$ for $t \in \Delta$ ( in a neighborhood of $\tau$, and for t for which this derivative and also $\dot{p}(t)$ exists). Below, the notation $O(\cdot)$ is used in the standard way. Having in mind that $\eta(\tau) = \dot{\eta}(\tau) = 0$ we have from (3.28)

$$\ddot{\eta}(t) = \sigma[-2(1 - p(t))\dot{p}(t)\xi(t) + (1 - p(t))^2\dot{\xi}(t)] + O(|t - \tau|).$$

Since $\tau$ is a Lebesgue point also for $\Omega$, there is a sequence $t_k \to \tau$ such that $t_k \in \Omega$ and the derivatives $\dot{p}(t_k)$ and $\dot{\eta}(t_k) = 0$ exist. Then from the equality $\sigma(1 - p)^2\xi = \alpha$ wich holds a.e. in $\Omega$ we obtain

$$0 = \sigma[-2(1 - p(t_k))\dot{p}(t_k)\xi(t_k) + (1 - p(t_k))^2\dot{\xi}(t_k)].$$

Combining the last two formulas we have

$$\ddot{\eta}(t) = -2\sigma(1 - p(t))(\dot{p}(t) - \dot{p}(t_k))\xi(t) + O(|t - \tau|) + O(|t_k - t|).$$

Using (3.26) and the equality $u(t) = \mu_1$ we obtain

$$\ddot{\eta}(t) = 2\sigma(\mu_1 - u(t_k))p(t)(1 - p(t))^2\xi(t) + O(|t - \tau|) + O(|t_k - t|). \tag{3.30}$$

Assume that

$$\lim_k u(t_k) = \mu_1 \tag{3.31}$$

does not hold. Since $\xi(t_k) = \frac{\alpha}{\sigma(1-p(t_k))^2} > 0$ uniform in k, (3.30) implies that some $\epsilon > 0$ in such that for almost every $t \in \Delta$

$$\ddot{\eta} \geq \epsilon - O(|t - \tau|) - O(|t_k - t|).$$

Since $\eta(\tau) = \dot{\eta}(\tau) = 0$ and since $\tau$ is a Lebesgue point for the set of t for which the above inequality holds, it is implied that $\eta(t) > 0$ in a left or right neighborhood of $\tau$. This contradicts the assumption of Case 1. Thus (3.31) holds.

Since $\dot{p}(t_k)$ and $\dot{\eta}(t_k)$ exist and $t_k \in \Omega$, equality (3.29) is fulfilled for $t_k$. Taking the limit in k we obtain that

$$\sigma + \mu - \mu_1 = \frac{1}{2p(\tau)}(r + \sigma p(\tau)^2) + \mu) - \frac{\sigma}{2\alpha p(\tau)}(1 - p(\tau))^2.$$

This equation holds for every contact point $\tau$. Viewing it as an equation $p = p(\tau)$ we rewrite it as

$$\sigma(1 - \alpha)p^2 - 2p(\sigma + \alpha(\sigma + \mu - u)) - \alpha(r + \mu) + \sigma = 0, \tag{3.32}$$

69

where $u = \mu_1$, but in Case 2 the same equation arises with $u = \mu_0$. In both cases condition (3.23) implies that the above quadratic function changes values in $[0, 1]$, which means that exactly one solution $p_0^* \in (0, 1)$ ($p_1^* \in (0, 1)$ in Case 2, respectively) exists. For later use we also introduce the solution $p^*$ of (3.32) with $u = \sigma + \mu$, which is an admissible value thanks to condition (3.22). Clearly, $p^* \in (p_1^*, p_0^*)$.

Since the function p(t) is strictly monotone, this implies that most often one contact point $\tau$ may exist. Then $\Omega$ is (modulo a set of measure zero) either an interval $[0, \tau)$ with finite or infinite $\tau$, or an interval $(\tau, +\infty)$.

Now we assume that the interval $\Omega$ is infinite. The control singular $u(\cdot)$ on $\Omega$ is determined from the equation (3.29), or equivalently, by (3.32) in a feedback from $u = v(p)$. Differentiating (3.32) we obtain that v is strictly decreasing. Since $p(\cdot)$ is strictly increasing in the case that we consider, $u(t) = v(p(t))$ is strictly decreasing. Since the value of $\sigma + \mu - u(t)$ is non-negative (otherwise $p(\cdot)$ would not be increasing) and $u(\cdot)$ is increasing in $\Omega$, we conclude that $\sigma + \mu - u(t) \geq \epsilon$ for some $\epsilon > 0$ and all $t \in \Omega$, eventually excluding an interval at the left side of $\Omega$. This implies that $p(t) \to 1$. But for values $p > p_0^*$ we have $v(p) < \mu_0$, therefore $u(t) = v(p(t))$ is not admissible for all sufficiently large t. Thus $\Omega$ cannot be unbounded.

*Case 2*
The case where $\eta(t) > 0$ on a set of positive measure can be considered in exactly the same way. Thus we obtain the following:

**Claim 3.3.2** *Any optimal control is either the constant $u = \sigma + \mu$, or has the following structures:*
*there is a finite $\tau \geq 0$ so that the control u is determined from the feedback system*

$$u(t) = \sigma + \mu - \frac{1}{2p(t)}(r + \sigma p(t)^2 + \mu) + \frac{\sigma}{2\alpha p(t)}(1 - p(t))^2, \qquad (3.33)$$

$$\dot{p}(t) = (\sigma + \mu - u)p(1 - p), \qquad p(0) = p_0, \qquad (3.34)$$

*while on $(\tau, +\infty)$ the control equals either $\mu_0$ or $\mu_1$. Moreover, $\tau > 0$ may only happen if $p(0) = p_0 \in [p_1^*, p_0^*]$ and $\tau$ is uniquely determined by the initial prevalence $p_0$: this is the first moment $t = \tau_1(p_0)$ at which $p(t) = p_1^*$ (this happens if $p_0 < p^*$), or the first moment $t = \tau_1(p_1)$ at which $p(t) = p_0^*$ (this happens if $p_0 > p^*$).*

Denote by $\hat{u}[p_0](\cdot)$ the control determined by (3.33) and (3.34). We introduce the following three control functions depending on the initial prevalence $p_0$:

$$u_0[p_0](t) = \mu_0,$$
$$u_1[p_0](t) = \mu_1,$$

$$u_s[p_0](t) = \begin{cases} \begin{cases} \bar{u}[p_0](t) & \text{for } t \in [0, \tau_1(p_0)], & \text{if } p_0 < p^*, \\ \mu_1 & \text{for } t > \tau_1(p_0), & \text{if } p_0 < p^*, \end{cases} \\ \sigma + \mu \quad \forall t & \text{if } p_0 = p^*, \\ \begin{cases} \bar{u}[p_0](t) & \text{for } t \in [0, \tau_0(p_0)], & \text{if } p_0 > p^*, \\ \mu_0 & \text{for } t > \tau_0(p_0), & \text{if } p_0 > p^*. \end{cases} \end{cases}$$

Notice that each of the mappings $p_0 \to u_0[p_0](\cdot)$, $p_0 \to u_1[p_0](\cdot)$, $p_0 \to u_s[p_0](\cdot)$ is continuous in the topology of uniform convergence on every compact interval. This is obvious for the first two mappings and can be easily checked for the third one, taking into account that $\tau_0$, $\tau_1$ and $p_0 \to \bar{u}[p_0](p_0)$.

Denote by $J_0(p_0)$, $J_1(p_0)$ and $J_s(p_0)$ the objective values corresponding to the controls $u_0$, $u_1$ and $u_s$. According to Claim 3.3.2, any of the above three controls for which the objective value equals max $J_0(p_0), J_1(p_0), J_s(p_0)$ is optimal, and only these controls are optimal. We denote by $\Omega_0$, $\Omega_1$ and $\Omega_s$ the sets of those $p_0 \in (0, 1)$ for which $u_0$, $u_1$, respectively $u_s$ is optimal. We denote the continuity of the above controls with respect to $p_0$ and $r > 0$ and we obtain the sets $\Omega_0$, $\Omega_1$ and $\Omega_s$ as closed.

Clearly, it is not possible that there are $p_0 \in \Omega_0$ and $p_1 \in \Omega_1$ for which $p_0 < p_1$. Indeed in this case the corresponding trajectories $p_0(t)$ and $p_1(t)$ would cross at some moment $\tau$, which would give rise to a non-monotone optimal trajectory (switching from $\mu_0$ to $\mu_1$ at $\tau$, or vice versa). For the same reason there are no points $p_1 \in \Omega_1$ and $p_s \in \Omega_s$ such that $p_1 > p_s > p^*$. Similarly, there are no points $p_0 \in \Omega_0$ and $p_s \in \Omega_s$ such that $p_0 < p_s < p^*$.

Let us compare $u_0$ and $u_s[p]$ for $p > p^*$ for which these two controls are different. Assume that there is a pair $p_0$ and $p_s$ such that $p_0 < p_s$ and $u_0$ is optimal for $p_0$, while $u_s[p_s] \neq u_0$ is optimal for $p_s$. Then $p_s < p_0^*$, since $u_s[p] = u_0$ for $p \geq p_0^*$. Hence, for every $p \leq p_s$ it holds

$$u_s[p](0) \geq \mu_0 + \delta \tag{3.35}$$

with some $\delta > 0$.

We already know that the control $u_1$ cannot be optimal for any $p > p_0$. Then for every $\epsilon > 0$ there exist $p_0^\epsilon \geq p_0$ and $p_s^\epsilon \leq p_s$ so that $p_0^\epsilon < p_s^\epsilon$ and $u_s[p_s^\epsilon]$ is optimal for $p_s$. But according to (3.35) the opimtal trajectories $p_0^\epsilon(t)$ and $p_s^\epsilon(t)$ will cross each other if $\epsilon$ is sufficiently small, which leads to a contradiction, as before. The same argument is valid also for $p < p^*$.

Thus the following cases arise:

1. $\min\{J_0(p^*), J_1(p^*)\} \leq J_s(p^*)$. Then we consider two following possibilities.

(a) $J_1(p^*) \leq J_0(p^*)$. Then there are numbers $0 < p^* \leq \hat{p}_{1s} \leq \hat{p}_{0s} \leq 1$ such that the following control is optimal:

$$u_1 \text{ for } p_0 \in (0, \hat{p}_{1s}], \quad u_s[p_0] \text{ for } p_0 \in [\hat{p}_{1s}, \hat{p}_{0s}], \quad u_0 \text{ for } p_0 \in [\hat{p}_{0s}, 1).$$

Clearly in this case $\hat{p}_{1s}$ is a critical point, provided that $\hat{p}_{1s} < 1$. If $\hat{p}_{1s} < \hat{p}_{0s} < 1$ then, in fact $\hat{p}_{0s} \leq p_0^*$. If this inequality is strict, then $\hat{p}_{0s}$ is also critical.

If $\hat{p}_{1s} = 1$, there are certainly no critical points. If $\hat{p}_{0s} = \hat{p}_{1s}$, or $\hat{p}_{0s} = p_0^*$ then $\hat{p}_{0s}$ is not critical, because of the continuity of the mapping $p_0 \to u_s[p_0]$ and the fact that $u_s[p_0] = u_1$ for $p_0 \geq \hat{p}_0$.

(b) $J_1(p^*) \geq J_0(p^*)$, similarly to (a). In the particular case $J_1(p^*) \geq J_0(p^*)$ the point $p^*$ is the only critical point.

2. $\min\{J_0(p^*), J_1(p^*)\} > J_s(p^*)$. In this case there are points $0 < p_1^* \leq \hat{p}_{1s} \leq \hat{p}_{0s} \leq p_0^* < 1$ such that the following control is optimal:

$$u_1 \text{ for } p_0 \in (0, \hat{p}_{1s}], \quad u_s[p_0] \text{ for } p_0 \in [\hat{p}_{1s}, \hat{p}_{0s}], \quad u_0 \text{ for } p_0 \in [\hat{p}_{0s}, 1).$$

In this case $\hat{p}_{1s}$ is a critical point if and only if $p_1^* < \hat{p}_{1s}$, and $\hat{p}_{0s}$ is a critical point if and only if $\hat{p}_{0s} < p_0^*$.

We see that the conditions for the existence of a critical point in Case 2 are not easy to check, since they involve the unknown point $\hat{p}_{1s}$, resp. $\hat{p}_{0s}$. In Case 1, however, a sufficient condition for existence of a critical point is $\hat{p}_{1s} < 1$, resp. $\hat{p}_{0s} > 0$. Both would hold if we prove that $J_0(1) < J_1(1)$ and $J_1(p) < J_0(p)$ for all sufficiently small $p > 0$.

We know that for any fixed control value u the corresponding objective value is

$$J(p; u) = \int_0^{+\infty} e^{-rt} \left[ e^{-\mu t} (p_0 e^{(\sigma+\mu-u)t} + 1 - p_0)^{\frac{\sigma}{u-\sigma+\mu}} (1 - p) \right. \tag{3.36}$$
$$\left. + \alpha e^{(\sigma-u)t} (p_0 e^{(\sigma+\mu-u)t} + 1 - p_0)^{\frac{\sigma}{u-\sigma+\mu}} p \right] dt.$$

The functions $J_0(p) = J(p; \mu_0)$ and $J_1(p) = J(p; \mu_1)$ are differentiable in $(0, 1)$. Obviously, we have

$$J(1; \mu_1) = -\frac{\alpha}{r + \mu_1} > -\frac{\alpha}{r + \mu_0} = J(1; \mu_0).$$

Moreover, $J(0; \mu_0) = J(0; \mu) = -\frac{1}{r+\mu}$. Thus, the existence of a solution to the equation $J(p; \mu) = J(p; \mu_0)$ in $(0, p)$ will be implied by the relation $J'(0; \mu_0) > J'(0; \mu_1)$. We calculate

$$J_p(0; u) = -\frac{1}{r+\mu} + \frac{\alpha}{\sigma - r - u} + \frac{\sigma}{(r+\mu)(\sigma - \mu - u)} - \frac{\sigma}{\sigma - \mu - u} \int_0^{+\infty} e^{(\sigma - r - u)t} dt.$$

According to (3.23), we have $r < \sigma - \mu_0$. Then $J(0; \mu_0) = -\infty$ while $J(0; \mu_1)$ is always finite. Thus, $J'(0; \mu_0) > J'(0; \mu)$, which was what we needed.

From (3.36) it is easy to find that

$$J_s(p^*) = J(p^*; \sigma + \mu) = \frac{1 - (1 - \alpha)p^*}{r + \mu + \sigma p^*}.$$

We have proven so far that a critical point $p_0$ (at least one) exists in $(0, 1)$, provided that

$$\max\{J(p^*; \mu_0), J(p^*; \mu_1)\} \geq \frac{1 - (1 - \alpha)p^*}{r + \mu + \sigma p^*}, \tag{3.37}$$

where $p^*$ is the unique solution of the equation

$$\sigma(1 - \alpha)p^2 - 2\sigma p - \alpha(r + \mu) + \sigma = 0.$$

The corresponding control is constant $u^* = \sigma - \mu$. We need to compare the values $J(p^*; \mu_0).J(p^*; u^*)$ and $J(p^*; \mu_1)$, for which the explicite formula (3.36) can be used. However, to achieve this analytically is very difficult. This is the only instance where we involve a "numerical" argument in the consideration.

**Proposition 3.3.3** *Let the inequalities (3.22)–(3.23) hold. Assume that the inequality (3.37) is fulfilled. Then there exists at least one critical point $\hat{p} \in (0, 1)$.*

It shall now be studied what happens with the trajectories at a critical point $\hat{p} = \hat{p}_{1s}$ that always arises in the case (a). The most interesting aspect is the parameter value $\mu = 0$, therefore we concentrate on this case. For $p_0 \leq \hat{p}$ the control $u = \mu_1$ is optimal and the optimal trajectory asymptotically approaches $S = (1 - p_0)^{\frac{\sigma}{\mu_1 - \sigma}} S_0 > 0$ and $I = 0$. For $p_0 \geq \hat{p}$ the optimal control is either $u = \mu_0$, or $u_s[p_0]$. In the first case, the critical prevalence $\hat{p}$ is the unique solution of the equation $J(p; \mu_0) = J(p; \mu_1)$. The optimal trajectory converges to $S = 0$ and $I = 0$. In this way $u_s[p_0]$ is optimal, the prevalence reaches $p_1^*$ in a finite time $\tau_0(p_0)$, then the control takes the constant value $\mu_0$. Therefore, the population asymptotically becomes extinct in this case, too.
The qualitative behaviour is the same also in the case (b).

# Appendix A

# Matlab Implementation

## A.1 Solver

The following ODE Solver is written by Prof. Dr. Vladimir Veliov, who gave it to me for solving my system.

```
% Solves the DE with r.-h. side f on the interval [t0, t0+n*h], n >= 0.
% h - the discretization step
% X0 - the initial condition (row)
% f should have the format void Y = f(X,t,[PAR1],[PAR2],[PAR3])
%     where Y is a column
% X - the trajectory: at each moment - one row
% If h<0, then the end-value prolem is solved on [t0+n*h,t_0],
% with the same meaning of y.
% The standard 4-th order Runge-Kutta metchod is applied. For
% the t-variable, the values of f(.,t) at the 2n+1
% points t0, t0+0.5h, t0+h, ..., t0+nh are used.

function X = ode_rk4(t0,T,h,X0,fun1,PAR1)

X(1,:) = X0;
n = floor(T/h)+1;
if (n <= 0) return
else
   t = t0;
   h2 = 0.5*h;
   h6 = h/6.0;
   for (i=1:(n-1))
      X1 = X(i,:);
      Y1 = feval('fun1',X1,t,PAR1);
      W = X1 + h2*Y1.';
      t = t + h2;
```

```
        Y2 = feval('fun1',W,t,PAR1);
        W = X1 + h2*Y2.';
        Y3 = feval('fun1',W,t,PAR1);
        W = X1 + h*Y3.';
        t = t + h2;
        Y4 = feval('fun1',W,t,PAR1);
        X(i+1,:) = X1 + h6*(Y1 + 2.0*(Y2+Y3)+Y4).';
    end


end
%disp('RK-4');
```

## A.2   Right side of the ODE System

```
function fun1  = fun1(X0,t, PAR1)

%this function gives the right side of the ODE System
%we need this to compute the values in the ode_rk4 solver
%X0=[S(0) I(0)]
%PAR 1 is a 9-dim parameter vector
%t is the initial time

a = PAR1(1);
b = PAR1(2);
alpha = PAR1(3);
beta = PAR1(4);
sigma = PAR1(5);
lambda= PAR1(6);
gamma = PAR1(7);
delta = PAR1(8);
kappa = PAR1(9);

y = (X0(2)/(X0(2)+X0(1)));  %calculate the prevalence

rho = a*(y^alpha)*((1-b*y)^beta);    %transmission function for increasing prevalence

S = -sigma*rho*X0(1) + kappa + lambda*X0(1) + gamma*X0(2);
I = sigma*rho*X0(1) - delta*X0(2);


fun1 = [S;I];
```

```
end
```

## A.3   Parameter Identification

```
function main

%this program is used for the estimation of the transmission function's free paramters


[parameter,fval,exitflag, iterations]=fminsearch(@funestimation,[2 3.5 1 0.2])

%on the left side are those values we got after minimizing the function funestimation
% using fminsearch with 4 good start values for a, b, alpha and beta

function f = funestimation(para)
a = para(1);
b = para(2);
c = para(3);   %stands for alpha
d = para(3);   %stands for beta


data = [0.023 0.032 0.044 0.061 0.081 0.106 0.133 0.16 0.185 0.206 0.223 0.236 0.245 0.251

%data=[0.035 0.0510 0.0730 0.1010 0.1330 0.1660 0.197 0.222 0.241 0.253 0.26 0.263 0.263 ]

n = n(:,2);

I0 = data(1);
S0 = 1 - data(1); % I/S+I = y0 ; S*y0+I*y0 =y0 ; S + I = 1 --> I =y0 : S=1-y0
X0 =[S0 I0];
h=0.1;

%PAR1 = [para(1) para(2)  para(3) para(4) 0.38 -0.108 0 0.23 0.2397]; %Botswana
PAR1 =[para(1) para(2) para(3) para(4) 0.43 -0.153 0 0.259 0.2825]; %Swaziland


z = feval('ode_rk4',0,n,h,X0,feval('fun1',X0,0, PAR1),PAR1);  %z uses the solver and fun1

prev = (z(:,2)./(z(:,1)+z(:,2))));
prev = prev(:,1);
data = data';
```

```
f = sum((data - prev(1:1/(h):n*(1/h))).^2);

% PAR1 =[a, b, alpha, beta, sigma, lambda, gamma, delta, kappa]
```

## A.4   Time

Time is just a short function wich scales the x label in a correct way.

```
function t = Time(T,h)

n = floor(T/(h))+1;
   t = h*[1:n];

end
```

## A.5   Plotting the Prevalence

```
function prevalence = prevalence(T,h)

t0 = 0;

%I0 = 0.035; % initial values for the infected people of Botswana
%S0 = 1 - 0.035;  % initial values for the susceptible people of  Botswana
% calculated by  I/S+I = y0 ; S*y0+I*y0 =y0 ; S + I = 1 --> I =y0 : S=1-y0
I0 = 0.023; % initial values for the infected people of Swaziland
S0 = 1- 0.023; %initial values for the susceptible people of Swaziland

X0 =[S0 I0];


% PAR1 =[a, b, alpha, beta, sigma, lambda, gamma, delta, kappa]

%PAR1 = [1.9521    3.5870    0.9764    0.3729 0.38 -0.108 0 0.23 0.2397]; %Botswana
PAR1 =[ 1.4733    3.1971    0.9555    0.3698 0.43 -0.153 0 0.259 0.2825]; %Swaziland


z = feval('ode_rk4',0,T,h,X0,feval('fun1',X0,0, PAR1),PAR1);
```

```
z = (z(:,2)./(z(:,1)+z(:,2)));    %computes the prevalence (I/(S+I))


t = feval('Time', T, h);


prevalence= plot(t, z);
set(gca, 'XTick', 0:5:20);
set(gca, 'XTickLabel', [ 1990:5:2010 ] ); %scales the x lable into the correct years
```

## A.6  Plotting the Set of Susceptible and Infected People

```
function set_of_SI = set_of_SI(T, h)

t0 = 0;

I0 = 0.023; % initial values for the infected people of Swaziland
S0 = 1- 0.023; %initial values for the susceptible people of Swaziland

X0 =[S0 I0];


% PAR1 =[a, b, alpha, beta, sigma, lambda, gamma, delta, kappa]

PAR1 =[ 1.4733    3.1971    0.9555    0.3698 0.43 -0.153 0 0.259 0.2825]; %Swaziland

z = feval('ode_rk4',0,T,h,X0,feval('fun1',X0,0, PAR1),PAR1);
t = feval('Time', T, h);
S1 = z(:,1);
I1 = z(:,2);

set_of_SI= plot(t, S1, 'b');
hold on
set_of_SI = plot(t, I1, 'r');

set(gca, 'XTick', 0:5:20);
set(gca, 'XTickLabel', [ 1990:5:2010 ] );
```

## A.7 Plotting the Summation of Susceptible and Infected People

```
function set_of_SplusI = set_of_SplusI(T, h)

t0 = 0;

I0 = 0.023; % intital values of Swaziland's set of infected
S0 = 1- 0.023; %intital values of Swaziland's set of susceptible

X0 =[S0 I0];


% PAR1 =[a, b, alpha, beta, sigma, lambda, gamma, delta, kappa]
PAR1 =[ 1.4733    3.1971    0.9555    0.3698 0.43 -0.153 0 0.259 0.2825]; %Parmeter vector

z = feval('ode_rk4',0,T,h,X0,feval('fun1',X0,0, PAR1),PAR1);
t = feval('Time', T, h);
S1 = z(:,1);
I1 = z(:,2);
u = S1+I1;
set_of_SplusI= plot(t, u, 'g');

set(gca, 'XTick', 0:5:20);
set(gca, 'XTickLabel', [ 1990:5:2010 ] );
```

# Bibliography

[1] H. D'ALBIS & E. AUGERAUD-VERON(2010). The Optimal Prevention of Epidemics. Toulouse School of Economics, University of La Rochelle.

[2] C. ALMEDER, G. FEICHTINGER, W. C. SANDERSON & V. VELIOV(2007). Prevention and medication of HIV/AIDS: the case of Botswana. CEJOR, 15, 47-61.

[3] S. BANSAL, B. T. GRENFELL & L. A. MEYERS(2007). When individual behaviour matters: homogeneous and network models in epidemiology. The Journal of the Royal Society Interface,4, 879-891.

[4] R. BONITA, R. BEAGLEHOLE & R. BEAGLEHOLE(2006). Basic Epidemiology. World Health Organisation, Version 0002.

[5] J. CHIN & S.K. LWANGA(1991). Estimation and projection of adult AIDS cases: a simple epidemiological model. WHO Bulletin OMS. Vol 69.

[6] O. DIEKMANN & J. A. P. HEESTERBEEK(2000). Mathematical Epidemiology of Infectious Diseases: Model Building, Analysis and Interpretation. John Wiley & Sons, Vol 1.

[7] K. DIETZ & K. P. HADELER(1988). Epidemiological models for sexually transmitted diseases. The Journal of Mathematical Biology, 26, 1-25.

[8] P. VAN DEN DRIESSCHE & JAMES WATMOUGH(2000). A simple SIS epidemic model with a backward bifurcation. The Journal of Mathematical Biology, 40, 525-540.

[9] S. GOLDMAN & JAMES LIGHTWOOD(1995). The SIS Model of Infectious Disease with Treatment.

[10] H. W. HETHCOTE(1994). A Thousand And One Epidemic Models .Frontiers in Mathematical Biology, 504-515.

[11] H. W. HETHCOTE(1989). Three Basic Epidemiological Models. Biomathematics, Vol. 18, 119-144.

[12] H. W. HETHCOTE & P. VAN DEN DRIESSCHE(1990). Some epidemiological models with nonlinear incidence. The Journal of Mathematical Biology, 29, 271-287.

[13] H. W. HETHCOTE & S. A. LEVIN(1989). Periodicity in Epidemiological Models. Biomathematics, Vol. 18, 193-211.

[14] H. W. HETHCOTE & D. W. TUDOR(1980). Integral Equation Models for Endemic Infectious Diseases. The Journal of Mathematical Biology, 9, 37-47.

[15] M. E. HOCHBERG(1991). Non-linear Transmission Rate and the Dynamics of Infectious Disease. The Journal of theoretic Biology, 153, 301-321.

[16] M.H. R. KHOUZANI(2011). Optimal Control of Epidemic Evolution. IEEE INFO-COM 2011, 1683-1691.

[17] W. LIU, S. A. LEVIN & Y. IWASA(1986). Influence of nonlinear incidence rates upon the behavior of SIRS epidemiological models. The Journal of Mathematical Biology, 23, 187-204.

[18] P. MICHEL(1982). On the Transversality Condition in Infinite Horizon Optimal Problem. Econometrica, Vol. 50, 1975-1985.

[19] H. McCALLUM, N. BARLOW & J. HONE(2001). How should pathogen transmission be modelled? Trends in Ecology and Evolution, Vol. 16.

[20] J. MENA-LORCA & H. W. HETHCOTE(1992). Dynamic models of infectious diseases as regulators of population sizes. The Journal of Mathematical Biology, 30, 693-719.

[21] D. J. NOKES & R. M. ANDERSON(1988). The use of mathematical models in the epidemiological study of infectious diseases and in the design of mass immunization programmes. Epidemiology and Infection, 101, 120.

[22] A. PUGLIESE(1990). Population models for diseases with no recoverys. The Journal of Mathematical Biology, 28,65-82.

[23] R. E. ROWTHORN, R. LAXMINARYAN & C. A. GILLIGAN(2009). Optimal control of epidemics in metapopulations. Journal of the Royal Society, Vol.6, Issue 41, 1135-1144.

[24] H. RAHMANDAD & J. STERMAN(2004). Heterogeneity and Network Structure in the Dynamics of Diffusion: Comparing Agent-Based and Differential Equatin Models. MIT Sloan School of Management Working Paper 4512-04.

[25] R. ROWTHORNY & F. TOXVAERD(2008). The Optimal Control of Infectious Diseases via Prevention and Treatment. JEL Classification, C73, I18.

[26] V. M. VELIOV(2005). On the effect of population heterogeneity on dynamics of epidemic diseases. The Journal of Mathematical Biology, 51, 123-143.