ISTITUTO DI ANALISI DEI SISTEMI ED INFORMATICA

CONSIGLIO NAZIONALE DELLE RICERCHE

E. Pourabbas, A. d'Onofrio, M. Rafanelli

A METHOD TO ESTIMATE THE INCIDENCE OF COMMUNICABLE DISEASES UNDER SEASONAL FLUCTUATIONS WITH APPLICATION TO CHOLERA

R. 485  Novembre 1998

**Elaheh Pourabbas -** Istituto di Analisi dei Sistemi ed Informatica del CNR, viale Manzoni 30 - 00185 Roma, Italy. Email : pourabbas@iasi.rm.cnr.it.

**Alberto d'Onofrio -** Dottorato in Informatica Medica, Università degli Studi di Roma "La Sapienza", presso l'Istituto di Analisi dei Sistemi ed Informatica del CNR.

**Maurizio Rafanelli -** Istituto di Analisi dei Sistemi ed Informatica del CNR, viale Manzoni 30 - 00185 Roma, Italy. Email : rafanelli@iasi.rm.cnr.it.

2.

## Abstract

This paper describes a method for estimating the seasonal variation of infection rate (or contact rate) and the trajectories of the number of susceptible, infectious and removed individuals in a deterministic *SIRS* model. The key idea of the proposed method is that the number of periodically varying infectives at time $t$ can be represented as a sum of functions of the form $b_1/(1 + b_2(t - kT - b_3)^2)$, $k = \ldots, -1, 0, 1, \ldots$, where $b_1$, $b_2$ and $b_3$ are parameters to be estimated from the incidence data, and $T$ is the period. Given the infective trajectory, the other trajectories and the contact rate can be estimated via the model definition. The method is illustrated using cholera incidence data from three developing countries. Finally, an analysis of sensitivity of parameter estimation for validating the obtained results is made.

## 1. Introduction

The spread of some rare diseases, such as the plague, can be considered an isolated phenomenon and therefore, in a strict sense, such an event is defined as an epidemic. On the other hand, some diseases, such as measles, cholera, poliomyelitis, diphtheria, influenza, etc., tend to be permanently present in a given territory, with periodical oscillations that can be annual, poli-annual or chaotic. This behaviour depends on the effect of the seasonal fluctuations of the contact rate on the incidence of the disease.

The regular oscillations of the number of cases around the average endemic level have aroused a great interest in epidemiologists and mathematicians. In particular, the study of the biennial oscillation of measles in some large communities has been an object of inquiry and explanation attempts by means of deterministic and stochastic models [1, 3].

In the literature, different authors have discussed the estimation problem regarding the infectious transmission rate, or simply the contact rate ($\beta$), when this rate is constant [2, 4, 5, 8]. In the case of *variable* contact rates, the main contribution was given by London and Yorke [13]. They obtained an estimation of the contact rate using a delay equation model, which includes a latent period between the time of infection and the beginning of the infectious period of the disease.

In [12] the time dependent contact rate for measles in the *SEIR* model is reconstructed by means of an extended Kalman filter from the incidence data of the disease in the city of New York. The authors conclude that, although these data through the years show an irregular change in the number of infected people, the contact rate is periodic and follows the season.

The contact rate parameter $\beta$ in *SIR* and *SIRS* models is, in general, subject to remarkable short term oscillation due, for example, to seasonal variations of weather and/or to the school vacations. This results in a one-year period variation of $\beta$. The one-year periodic variation of parameters is a phenomenon common to many ecological models. This variation of $\beta(t)$ results either in periodic oscillations of one year (for $S$, $I$, $R$) and in poli-annual or chaotic behaviour.

In the papers [7, 9, 10, 14], where the model properties are studied (and where the problem of determining $\beta$ has not been solved), it is assumed that the contact rate $\beta$ undergoes a simple harmonic oscillation with a period of one year:

$$\beta(t) = \beta_o + \beta_1 Cos(\omega t), \ 0 \le \beta_1/\beta_o \le 1$$

This approximation is suitable to study the geometrical properties of the system, but it is not suitable for a quantitative study. In the present paper, instead, we consider many harmonics for the variable contact rate in order to obtain a more realistic value.

Our purpose is to develop, by mixing geometrical and quantitative considerations, a method to estimate both the curves of infectious $I(t)$, susceptible $S(t)$, and removed individuals $R(t)$, and the seasonally varying contact rate, for the *SIRS* epidemic model . To obtain realistic results, we refer to data on the incidence of cholera in three developing countries.

This study is part of a project in progress for planning health resources on the basis of the epidemiologic situation present in the territory itself [15]. The main characteristics of the data on cholera are discussed below.

In this paper we use the periodically forced *SIRS* model, and we assume that:

a) $I(t)$, $S(t)$, $R(t)$ undergo annual periodic oscillations;

4.

    b) the death rate of infectious individuals is not different from the death rate of susceptibles and removed;

    c) the total population size $N$ is constant, i.e. susceptibles are assumed to be added to the population at a constant rate that equals the total loss.

The method we use to analyze the seasonal variation of the incidence of cholera is based on an approximation of $I(t)$ by means of a series of functions similar to Agnesi's witch. A sensitivity analysis to give validity to the proposed method and results, is also performed.

## 2. Overview of data

The data on which we applied our method, details of which are given in the following section, refer to the cumulative cases of hospitalized persons who have developed clinical symptoms of cholera. These data, informally provided by the World Health Organization in Rome-Italy, gave the weekly number of cases of cholera in various countries where such a disease was present between 1993 and 1994.

    Attention was given to cholera data for two main reasons. The first is because such a disease is endemic in the majority of the developing countries and, hence, these data are a good starting point for studying the periodic behaviour of the contact rate in the *SIRS* mathematical model [6, 11]. The second reason deals with the difficulty, and sometime the impossibility, of acquiring data on the time-scale of weeks, as experienced by the health organizations of developing countries for different diseases.

    From all the data received, we have chosen the data of 3 countries (El Salvador, Nicaragua, and Somalia) on the basis of a criterion of relative completeness of the surveys (see Figure 1). Nevertheless, it is evident from the figure that there are no records for some weeks in the period considered. The total population, for the three countries considered, together with the life expectancy, are given in Table 1.

Table 1

Demographic data for the populations considered

| Country | Population | Life expectancy |
| --- | --- | --- |
| El-Salvador | 5.500.000 | 67.50 |
| Nicaragua | 4.206.000 | 64.52 |
| Somalia | 7.300.000 | 55.74 |

One of the conditions to take into account in applying a *SIR* or *SIRS* model refers to the hypothesis of homogeneous mixing of the population. This hypothesis appears reasonable for all the three countries chosen, because the populations under study are prevalently rural, with rather homogeneous social and economic conditions. In the literature, for example in [2] and [12], where the urban populations (Leningrad and New York) were considered, the homogeneous mixing was assumed for total population sizes of 3.5 millions and, respectively, 9 millions of individuals.

    Furthermore, we note that the proposed method is not constrained to the whole country, but it can be applied even to a single region or area, which is not subject to immigration and emigration, where a disease that can be described by means of the *SIRS* model is present.
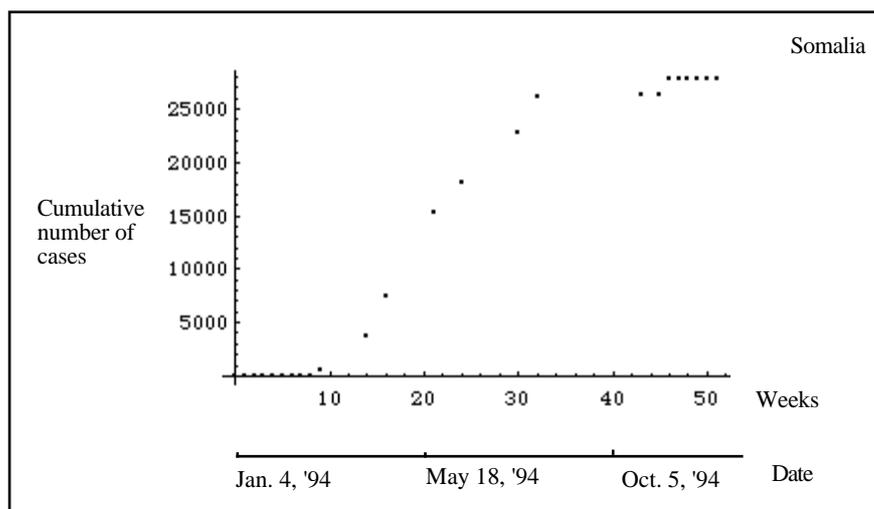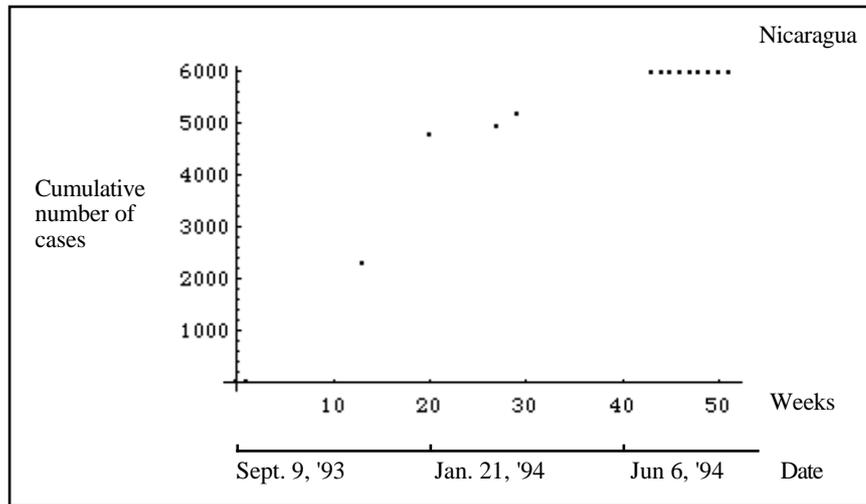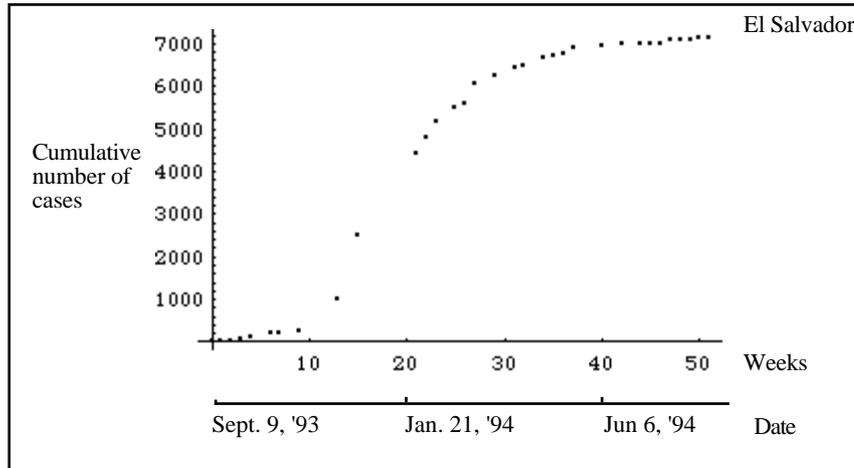
Figure 1. Cumulative number of cholera cases

## 3. Estimation method

A simple deterministic mathematical model suitable for the study of infectious diseases (such as cholera) is the *SIRS* model with temporary immunity. The vital dynamics are described by the following set of differential equations [11]:

$$\frac{dI}{dt} = \beta(t)S(t)I(t) - (\gamma + \mu)I$$

$$\frac{dR}{dt} = \gamma I - (\mu + \theta)R \tag{1}$$

$$S(t) + I(t) + R(t) = N$$

where $\gamma$ is the removal rate, $\mu$ is the death rate and $\theta$ is the loss of immunity rate. The rates $\gamma$, $\mu$, $\theta$ are taken constant in time and known a priori.

The available data, as mentioned above, give the cumulative amount of people who were hospitalized during one year. In our model, this amount will be expressed by the function $u(t)$ defined as follows:

$$u(t) = \int_0^t \gamma I(\eta)d\eta, \qquad 0 \le t \le T \tag{2}$$

where $T = 52$ weeks (time will be measured in weeks in the following) .

By inspection of real data (see Figure. 1), it is easy to note that these data qualitatively seem to lie on a translated arc-tangent function. This suggests that $I(t)$ can be approximated by functions similar to the derivative of the arc-tangent function, translated with respect to the starting time (corresponding to the time of initial observed data) and suitably rescaled. Let $v(t)$ be the derivative of $Arctan(t)$:

$$v(t) = \frac{1}{1 + t^2} \tag{3}$$

Then $I(t)$, that we assume to have annual periodicity, will be expressed by the following formula

$$I(t; A, c, t_M) = (\frac{1}{\gamma}) \sum_{k=-\infty}^{+\infty} Acv(c(t - t_M - kT))$$

$$= (\frac{1}{\gamma}) \sum_{k=-\infty}^{+\infty} \frac{Ac}{1 + c^2(t - t_M - kT)^2} \tag{4}$$

where $T = 52$ weeks is the period, $A > 0$ is the parameter which controls the maximal amplitude of a single function, and $t_M < T$ is the time in which the maximum occurs for the function with $k = 0$. At $t = t_M \pm \frac{1}{c}$ this function has the value $0.5Ac$. We note that for a particular choice of parameters ($A = 4a^2$, $c = 1/2a$, $t_M = 0$) the basic function in series (4) is known as Agnesi's witch.

If the series in Eq.(4) converges, then, for construction, it will converge to a periodic function of period $T$. Therefore, we can study the convergence only in $[0, T]$.

Let us denote by $v_k(t)$ the function:

$$v_k(t) = (\frac{1}{\gamma}) \frac{Ac}{1 + c^2(t - t_M - kT)^2} \qquad k = \ldots, -1, 0, 1, \ldots \tag{5}$$

and by $s_{p,q}(t)$ the following partial sum:

$$s_{p,q}(t) = \sum_{k=-p}^{q} v_k(t) \tag{6}$$

The series in Eq.(4) is totally convergent. In fact, by considering the $\infty$-norm:

$$\|f(t)\|_\infty = \operatorname{Sup}_{t \in [o, T]} |f(t)| \tag{7}$$

it results:

$$
\begin{aligned}
\left\| s_{p,q}(t) \right\|_\infty &= \left\| \sum_{k=-p}^{q} v_k(t) \right\|_\infty \leq \sum_{k=-p}^{q} \left\| v_k(t) \right\|_\infty \\
&= \sum_{k=-p}^{-1} |v_k(0)| + v_0(t_M) + \sum_{k=1}^{q} |v_k(T)| \\
&= \left( \frac{Ac}{\gamma} \right) \left( \sum_{k=1}^{p} \frac{1}{1 + c^2(kT - t_M)^2} + 1 + \sum_{k=1}^{q} \frac{1}{1 + c^2(T - t_M - kT)^2} \right)
\end{aligned}
\tag{8}
$$

The convergence is obvious.

In practice, it is possible to use the following approximated form of $I$, with $Z$ large enough:

$$I(t; A, c, t_M) \approx (\frac{1}{\gamma}) \sum_{k=-Z}^{+Z} \frac{Ac}{1 + c^2(t - t_M - kT)^2} \tag{9}$$

Therefore, we obtain the following function to represent the notified cases:

$$\tilde{u}(t; A, c, t_M) = A \sum_{k=-Z}^{+Z} [Arctan(c(t - t_M - kT)) + Arctan(c(t_M + kT))] \qquad 0 \leq t \leq T \tag{10}$$

For determining the parameters $A$, $t_M$, and $c$ which optimally fit the function $\tilde{u}$ to the observed data, the least square method was applied using the non-linear optimization program *NonlinearFit* function of *Mathematica*®. The value of $Z$ was set equal to 6. Higher values did not change the results in all the considered cases.

Let $\hat{I}(t)$ be the estimated function for $I(t)$ given by

8.

$$\hat{I}(t) = (\frac{1}{\gamma}) \sum_{k=-\infty}^{+\infty} \frac{\hat{A}\hat{c}}{1 + \hat{c}^2 (t - \hat{t}_M - kT)^2} \tag{11}$$

where $\hat{A}$, $\hat{t}_M$, $\hat{c}$ are the parameter estimates obtained as described above. The function $\hat{I}(t)$ can be written in Fourier's series form as

$$\hat{I}(t) = \sum_{k=0}^{+\infty} \hat{I}_k Cos(k\omega t - \alpha_k), \quad \omega = \frac{2\pi}{T}, \quad \alpha_0 = 0 \tag{12}$$

From the second equation of system (1), we can obtain the estimate $\hat{R}(t)$ for the removed individuals as the solution of the equation:

$$\frac{d\hat{R}}{dt} + (\mu + \theta)\hat{R} = \gamma\hat{I}(t) \tag{13}$$

In the periodic condition, $\hat{R}(t)$ can be represented by:

$$\hat{R}(t) = \gamma \sum_{k=0}^{+\infty} \frac{\hat{I}_k}{|(\theta + \mu) + jk\omega|} Cos(k\omega t - \alpha_k - Arctan(\frac{k\omega}{(\theta + \mu)})) \tag{14}$$

where $j = \sqrt{-1}$.

In a similar way, the estimate $\hat{S}(t)$ for the number of susceptible individuals is given by:

$$\hat{S}(t) = N - \hat{I}(t) - \hat{R}(t) \tag{15}$$

From the approximating functions $\hat{R}$, $\hat{I}$ and $\hat{S}$, and from the first equation of (1), we can obtain the following estimate of $\beta(t)$:

$$\hat{\beta}(t) = \frac{\gamma + \mu}{\hat{S}(t)} + \frac{1}{\hat{S}(t)\hat{I}(t)} \frac{d\hat{I}}{dt} \tag{16}$$

## 4. A case study: results for cholera

The previous method was applied to the data of the above mentioned countries to determine $S(t)$, $I(t)$ and $R(t)$, and to estimate the parameter $\beta(t)$. The theoretical curves representing the best fitting of cumulative cholera cases of the three considered countries are shown in Figure 2. As an example of the complete results, in Figures 3, 4, and 5 the time-course of infectious, susceptible, and removed individuals, the contact rate and the phase diagram for El Salvador, are shown. The values for parameters $\gamma$ and $\theta$ commonly used in literature are: $\gamma = 1.4 \ week^{-1}$ (the average length of infectious period for cholera is 5 days) and $\theta = 0.0034 \ week^{-1}$ (the temporary immunity period is about 5 years). Concerning the total number of individuals in the population and the mortality rate, we used the demographic data reported in Table 1, estimating the mortality rate as the inverse of the life expectancy.

In Table 2, the estimated values of parameters $A$, $t_M$, and $c$ are reported for each country considered. From the parameter $c$, the epidemic length was calculated. This length is defined as the time interval in which $I(t) \geq 0.5 I_{max}$. We have approximated this value as $2/c$ (see Eq. (4)), and this quantity is also reported in Table 2.

<div align="center">

Table 2

Estimated parameter values

</div>

| Country | $A$ | $t_M$ | $c$ | $2/c$ | $Max - Min(of\ infectious)$ |
|---|---|---|---|---|---|
| El Salvador | 2353.9 | 18.4 | 0.21 | 9.64 | 50.0 - 350.0 |
| Nicaragua | 1968.0 | 13.9 | 0.25 | 7.94 | 35.0 - 370.0 |
| Somalia | 9177.0 | 20.4 | 0.22 | 9.09 | 100.0 - 1500.0 |

Finally, in Table 3 the mean, minimum, and maximum values of contact rate and the relative amplitude of $\beta(t)$ are listed. We note that the values for $A$, $t_M$, and $c$ so obtained are in convincing physical ranges. As regards $\beta(t)$, it should be remarked that the peak-peak variation of $\beta$ with respect to its mean value is also in an acceptable physical range (29% - 34%, as reported in Table 3).

<div align="center">

Table 3

Mean, minimum, and maximum values of contact rate and its relative amplitude

</div>

| Country | $\beta(t) \times 10^7$ | | | $\dfrac{Max - Min}{\overline{\beta}(t)}$ |
|---|---|---|---|---|
| | $\overline{\beta}(t)$ | $Max$ | $Min$ | |
| El Salvador | 2.56 | 2.93 | 2.19 | 0.29 |
| Nicaragua | 3.35 | 3.93 | 2.78 | 0.34 |
| Somalia | 1.95 | 2.26 | 1.65 | 0.31 |

Since the estimates of $\beta(t)$ were computed in practice by truncating the series in Eq. (14), the effect of such a truncation was checked by integrating numerically system (1) using for $\beta(t)$ its estimate. The function $I(t)$ so computed was compared, for $t \in [0, T]$, to the function $I(t; \hat{A}, \hat{c}, \hat{t}_M)$ as given by (9), and truncations were adjusted to make differences negligible.
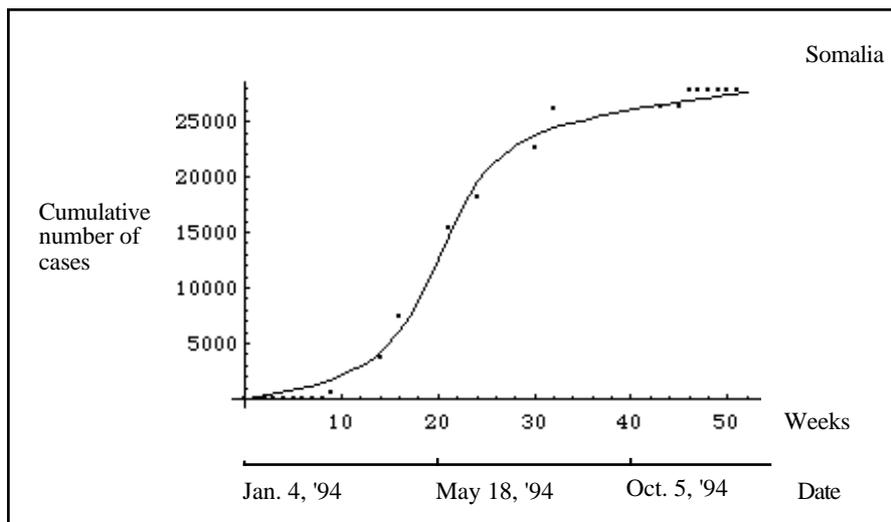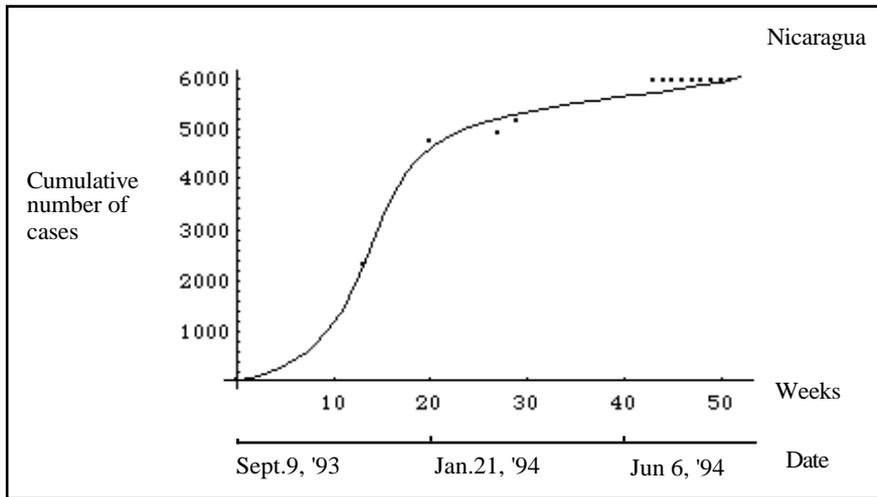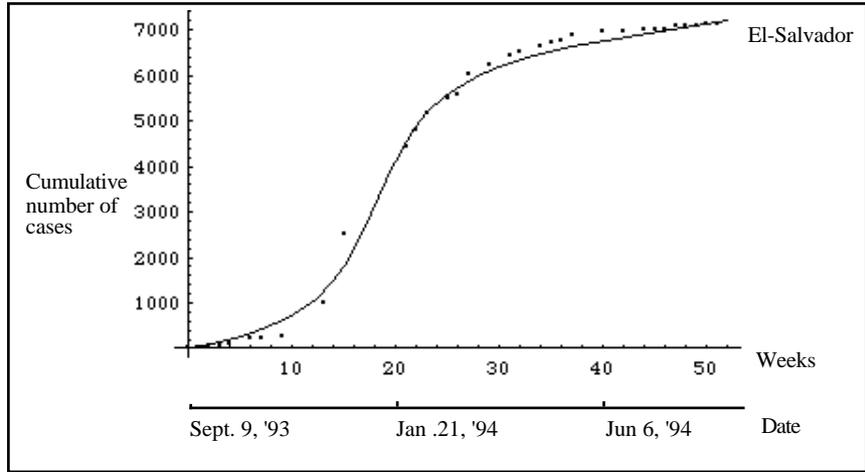
10.



Figure 2. Theoretical curves representing the best fitting of cumulative cholera cases
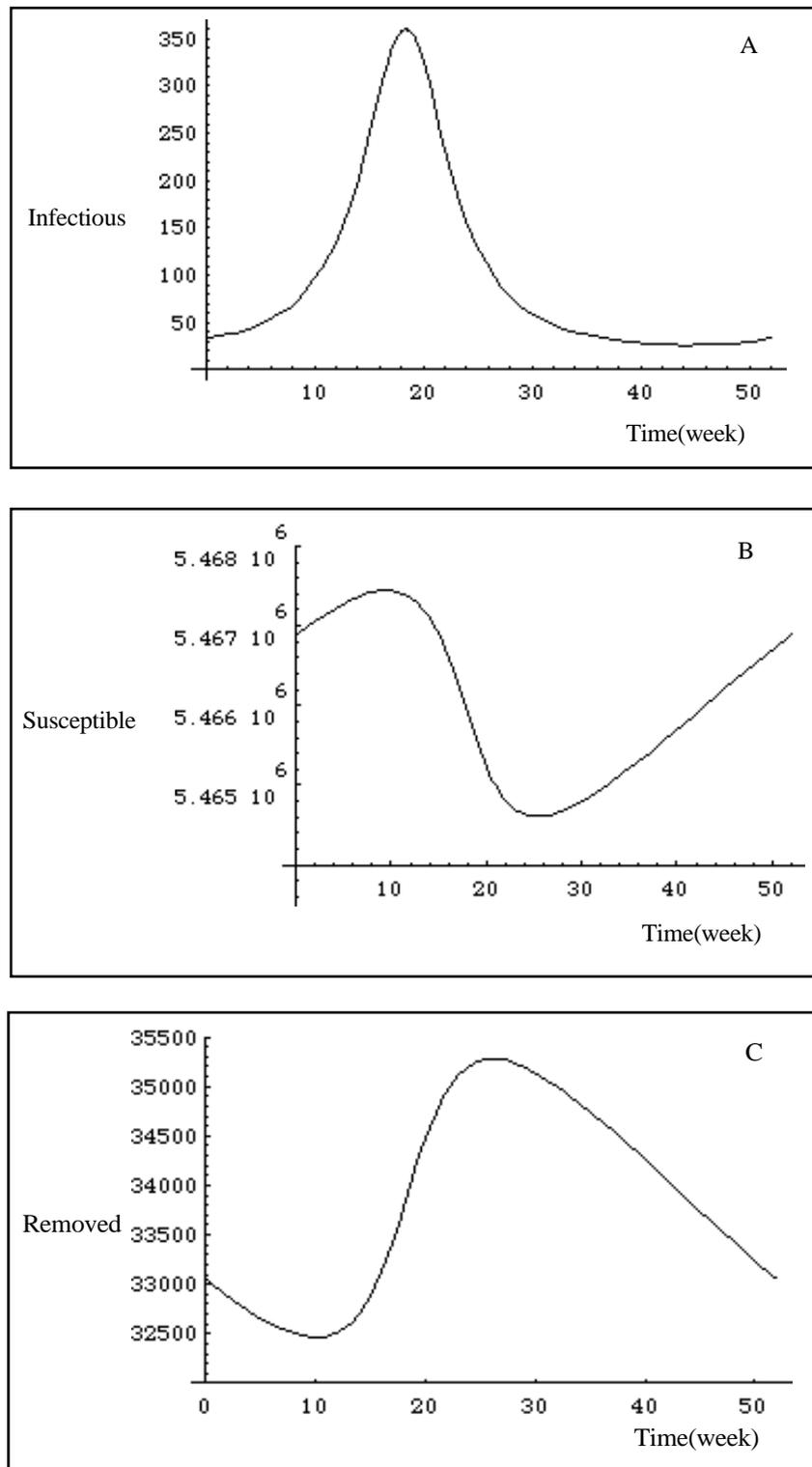
Figure 3. Estimated number of infectious (A), susceptible (B), and removed individuals (C) in El Salvador. Time is counted from September 9, 1993.
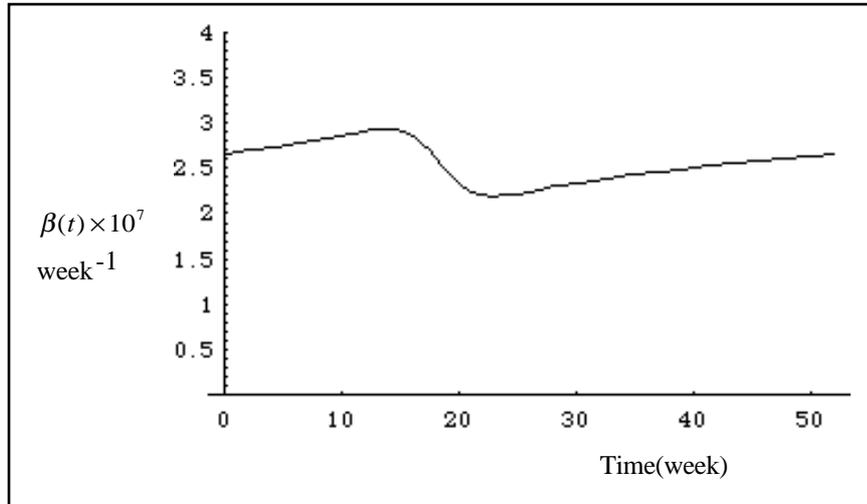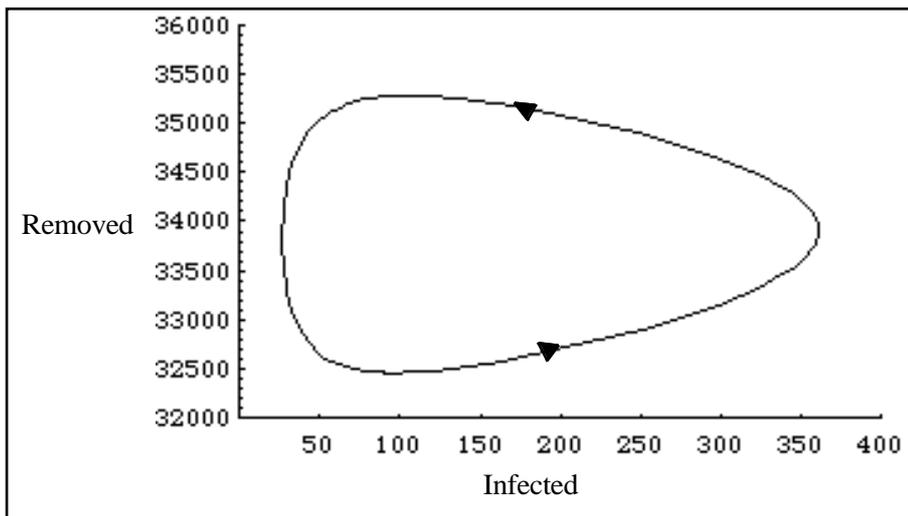
12.



Figure 4. Estimated contact rate



Figure 5. Phase diagram

## 5. Sensitivity analysis of parameter estimation

The problem of characterizing the error in parameter estimation is particularly difficult. In particular, this is a recurrent problem in epidemiology, because no information on the state variables can be considered reliable. In our case, the sensitivity of the estimates of parameters to the errors affecting the data, or to other possible data corruption, was evaluated by suitable computer simulations.

Since there is no precise information on the type of error that affects the available data, we can only make some hypotheses. A diagnosis of cholera when an other disease has occurred can be considered improbable, since cholera is a disease thoroughly studied, and easy to identify because it has strong and characteristic symptoms. Some extent of under-reporting is instead possible, although the percentage of sick people escaping hospitalization should not be very large, due to the great visibility of cholera symptomatology.

Moreover, inspection of data shows that the number of data-points on the curve of cumulative number of cases is very variable, and that their location is far from uniform.

We started assuming that the "true" number of infectious individuals at time t, $I*(t)$, is given by Eq. (4) with parameter values $A*$, $t*_M$, and $c*$.

Since the real data sets contain a number of points much less than $52$, that is there is not always a report every week, we extracted at random $N$ values from the set $\{1,\ldots,52\}$ and let $\{t_1,\ldots,t_N\}$ be the set of these values. The set $\{t_1,\ldots,t_N\}$ will be the simulated instance of times of reporting.

The number of individuals that should be hospitalized between the $t_{i-1}$-th and $t_i$-th week can then be computed as

$$w*_i = \int_{t_{i-1}}^{t_i} \gamma I*(t)dt \tag{17}$$

where $t_0 = 0$.

To simulate the presence of under-reporting, we simulate the actual number of hospitalized people between the $t_{i-1}$-th and $t_i$-th week as follows:

$$w_i = w*_i - x_i \qquad i = 1,\ldots,N \tag{18}$$

where $x_i$ is a random variable uniformly distributed in the interval $[0, 0.1\,w*_i]$. The simulated cumulative incidence data are given by

$$z_i = \sum_1^i w_i, \qquad i = 1,\ldots,N \tag{19}$$

From these "new" data we have estimated the values of parameters $A$, $t_M$, and $c$. This procedure was repeated many times (exactly $100$ iteration was made) and the statistical means for parameters are calculated, as table 4 and 5 show. In these table, the mean, standard deviation (SD) and bias (i.e. mean value minus true value) of the parameter estimates are given.

Two sets of true parameter values were considered (close to the estimated values for El Salvador and Somalia, respectively), and $N = 15$ and $N = 30$ were assumed.

Table 4

Mean, standard deviation, and bias calculated for $A*=2500$, $c*=0.25$, $t*_M=15.0$

| Points | Mean | Standard Deviation | Bias |
|--------|--------|--------------------|-----------|
|        | 2250.45 | 66.69 | - 249.55 |
| 15     | 0.2526 | 0.0106 | + 0.0026 |
|        | 15.04 | 0.17 | + 0.04 |
|        | 2244.35 | 52.45 | - 255.65 |
| 30     | 0.2496 | 0.0075 | - 0.0004 |
|        | 14.99 | 0.11 | - 0.01 |

Table 5

14.

Mean, standard deviation, and bias calculated for $A^*=10000$, $c^*=0.25$, $t^*_M=20.0$

| Points | Mean | Standard Deviation | Bias |
|--------|------|--------------------|------|
|        | 9069.18 | 242.78 | - 930.82 |
| 15     | 0.2520 | 0.0104 | + 0.0020 |
|        | 19.97 | 0.15 | - 0.03 |
|        | 9017.53 | 548.52 | - 982.47 |
| 30     | 0.2413 | 0.0423 | - 0.0087 |
|        | 19.99 | 0.10 | - 0.01 |

As we can see, the standard deviation of the estimates is very small. A bias on the estimates of $A$ is instead present, as expected, since noise with negative values was added to $w^*_i$ and it follows very close, as percentage, the extent of this noise.

## 6.  Concluding remarks

In this paper, we have proposed a method which, starting from simple geometrical considerations on the shape of the cumulative curve of cholera incidence data, leads to an analytical development that allows us to obtain an acceptable time-behaviour both for the system of *SIRS* model variables and for the contact rate.

The main idea of the method is to represent the periodic time-course of the number of infectives by a rather flexible function defined by three parameters. The values of parameters are to be estimated by fitting the integral of this function, times the removal rate, to the cumulative incidence data. The sensitivity of parameter estimation to data corruption appears to be reasonably low, on the basis of the results obtained on computer-generated data.

The use of a smoothing function to represent the data curve may result in some loss of information in recovering $I(t)$. It can be noted, however, that this loss refers mainly to the high frequency components, if they are present, of the unknown functions. Thus, the estimate of $I(t)$ is expected to be reasonably accurate if this quantity has slow temporal variations (as it seems likely). If high frequency components are present in the time behaviour of $I$, they cannot be recovered by the described procedure. This problem is indeed ill-conditioned, since high-frequency components are found filtered in the output curve $u(t)$. Similar considerations are also valid for the estimation of $\beta(t)$.

## Acknowledgements

## References

[1]  N. T. J. Bailey, The Mathematical Theory of Infectious Diseases and Its Applications, Griffin, London, 1975.

[2]  N. Bailey, Macro-Modelling and Prediction of Epidemic Spread at Community Level, in The Population dynamics of infectious diseases: Theory and applications (R. M. Anderson, Ed.), Chapman and Hall, 1982.

[3] M. S. Bartlett, Deterministic and stochastic models for recurrent epidemics, Proc. Third Berkeley Symp. Math. Statist. & Prob., 4 (81), 1956.

[4] N. Becker, J. Angulo, On Estimation the Contagiousness of a Disease Transmitted from Person to Person, Math. Biosci., 54 (137), 1981.

[5] N. Becker, Estimation for an Epidemic Model, Biometrics, 32 (769), 1976.

[6] B. Cvietanovitz, B. Grab, K. Uemura, Dynamics of Acute Bacterial Diseases. Epidemiological Models and their Application in Public Health, Supplement No. 1 to Vol. 56 of the Bullettin of the World Health Organization, Geneve 1978.

[7] K. Dietz, The Incidence of Infectious Diseases Under the Influence of Seasonal Fluctuations, Mathematical Models in medicine, Workshop, Mainz, March , Ed. J. Berger, W. BÜhler, R. Repges, and P. Tautu, Lecture Notes in Biomathematics, 11(1), 1976.

[8] G. Di Lena, G. Serio, A Discrete Method for the Identification of Parameters of a Deterministic Epidemic Model, Math. Biosci. 60 (161), 1982.

[9] N. M. Ferguson, D. J. Nokes, R. M. Anderson, Dynamical Complexity in Age-Structured Models of the Transmission of the Measles Virus: Epidemiological Implications at High Levels of Vaccine Uptake, Math. Biosci., 13 (101), 1996.

[10] Z. Grossman, I. Gumowski, K. Dietz, The Incidence of Infectious Diseases Under the Influence of Seasonal Fluctuations in Analytical Approach, Nonlinear Systems and Applications, An International Conference, V. Lakshmikantham Ed. , Academic press, Inc. 1977.

[11] H. W. Hethcote, Qualitative Analyses of Communicable Disease Models, Math. Biosci. 28 (335), 1976.

[12] M. Kalivianakis, S. L. J. Mous, J. Grasman, Reconstruction of the seasonally varying contact rate for measles, Math. Biosci. 124 (225), 1994.

[13] W. P. London, J. A. York , Recurrent Outbreaks of measles, chickenpox and mumps, I: Seasonal variation in contact rates, Amer. J. Epidem., 98 (453), 1973.

[14] E. Pourabbas, A. d'Onofrio: A SIR Epidemic Model and the Parametric Resonance, 3rd European Conference on Mathematics Applied to Biology and Medicine, Heidelberg (Germany), 6-10 October, 1996.

[15] E. Pourabbas, Design and implementation of an Information System Integrated with Epidemic Models for Planning Health Care Resources, Ph. D. Theses, University of Bologna, 1997.

[16] H. E. Soper, Interpretation of periodicity in disease prevalence, J. R. Statist. Soc., 92 (34), 1929.