

## STOCHASTIC MODELS IN THE PROBLEMS OF PREDICTING THE EPIDEMIOLOGICAL SITUATION\*

O. V. Bogdanov<sup>1</sup> and P. S. Knopov<sup>1†</sup>

UDC 519.21

**Abstract.** *The paper investigates some stochastic models with discrete and continuous time to solve important problems of predicting the spread of epidemiological diseases in the population. Various factors of epidemic spread and the main parameters influencing the forecast assessment are taken into account. Some test calculations based on the proposed methods have been performed.*

**Keywords:** *methods, optimization, modeling, stochastic equations, estimation, epidemic, discrete and continuous time.*

### INTRODUCTION

The COVID-19 pandemic has become a global challenge for humanity in the 21st century. It needs adequate methods and means of its control. In the absence of herd immunity and coronavirus drugs, as well as unequal access to vaccines, the epidemic threatens human life and health. But at the same time, long-term quarantine and measures to limit the pandemic cause economic damage and hamper the economic development. Therefore, decisions in disease spread control need special consideration: on the one hand, it is about the lives and health of a large number of people; on the other, there are significant economic losses and potential impoverishment.

Under these conditions, there is a growing need for modeling and decision-making support tools based on accurate calculations of their consequences. Such tools include various models of predicting the epidemiological situation and medical assistance needs, models of predicting the economic consequences of governmental (regional) decisions to limit the epidemic, etc. It is also necessary to take into account various risks and uncertainties that occur in modeling such complex processes with the stochastic (uncertain) nature of their components. This requires appropriate mathematical methods, including the use of random processes and fields, stochastic differential equations, regressions of special kind, modern apparatus of risk measures, etc. In what follows, we will present some approaches to solving the above problems.

As initial models, we took SIR (SEIR) and similar epidemiological models, which allow predicting the impact of restrictive measures on the dynamics of the spread of the disease. The main factor in these models is the virus replication rate (reproduction coefficient), which significantly depends on such measures. The study [1] analyzes the impact of school and workplace closures, public event cancellations, prohibition of public transportation, restrictions on domestic and international movement on the daily rate of spread of the disease.

The main problem for such models is the difficulty of their set up (identification) by real data. More detailed models require more complete data on the disease profile and its prevalence. In deterministic models, most of their parameters calculate average values. They take into account the stochasticity of processes on average, successfully

---

\* The study was partially supported by the National Research Foundation of Ukraine. Grant # 2020.02/0121.

---

<sup>1</sup>V. M. Glushkov Institute of Cybernetics, National Academy of Sciences of Ukraine, Kyiv, Ukraine, <sup>†</sup>*knopov1@yahoo.com*. Translated from *Kibernetika ta Systemnyi Analiz*, No. 1, January–February, 2022, pp. 70–76. Original article submitted September 8, 2021.

approximating the situation for large populations in a homogeneous environment. Stochastic models, in contrast to their deterministic analogs, more adequately reflect the course of the processes, especially in local or transient processes.

The present article describes some models for predicting the epidemiological situation and attempts to adjust them. Note that more detailed models require access to more complete data as to the disease profile and its prevalence. Also, one of the issues is incomplete testing of the population and latent course of diseases with mild symptoms in some people.

## DETERMINISTIC EPIDEMIOLOGICAL MODELS

According to [2], by the SIR model, population is divided into categories (compartments):  $S$  (Susceptible),  $I$  (Infected (detected illness cases)), and  $R$  (Recovered), whose dynamics can be described by the corresponding systems of differential equations:

$$\begin{aligned}\frac{dS}{dt} &= -\beta \frac{1}{N_0} SI, \\ \frac{dI}{dt} &= \beta \frac{1}{N_0} SI - \gamma I, \\ \frac{dR}{dt} &= \gamma I,\end{aligned}$$

where  $\beta$  and  $\gamma$  are morbidity rate and speed of recovery, respectively, and  $N_0$  is the number of people at risk of infection.

This system is substantiated as follows: increase in the number of infected people is directly proportional to the number of infected people (because the greater the number of infected people, the more sources of infection) and to the number of susceptible people. The number of people who recovered with the acquisition of immunity and the number of deaths are directly proportional to the number of infected people detected. This model clearly describes the beginning of an epidemic, when the density of infected people among susceptible ones is low.

Variables may often not be the absolute values  $S$ ,  $I$ , and  $R$  but their proportions with respect to  $N_0$ , i.e.,  $S_* = S / N_0$ ;  $I_* = I / N_0$ ; and  $R_* = R / N_0$ . Such a system can be represented in a more compact form.

Let us use the regular notation:  $T_{\text{inf}} = 1 / \gamma$ ,  $R_0 = \beta / \gamma$ , where  $T_{\text{inf}}$  is the average period of time when an infected person  $y_j^*$  spreads the disease,  $R_0$  is respectively the reproductive rate, i.e., average number of infecting by a sick person during their infectivity period. Then  $\gamma = 1 / T_{\text{inf}}$  and  $\beta = R_0 \gamma = R_0 / T_{\text{inf}}$ . Thus, we get the equations

$$\begin{aligned}\frac{dS}{dt} &= -\frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} SI, \\ \frac{dI}{dt} &= \frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} SI - \frac{1}{T_{\text{inf}}} I, \\ \frac{dR}{dt} &= \frac{1}{T_{\text{inf}}} I.\end{aligned}\tag{1}$$

Note that the distribution of the population by the respective categories satisfies the overall balance  $N_0(t) = S(t) + I(t) + R(t)$ . Therefore, medium- and long-term modeling should take into account the dynamics of the population  $N_0(t)$  (demographics), given the number of births  $b$  per unit time and mortality rate  $\mu$  in the population:

$$\begin{aligned}\frac{dS}{dt} &= -\frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} SI + b - \mu S, \\ \frac{dI}{dt} &= \frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} SI - \frac{1}{T_{\text{inf}}} I - \mu I, \\ \frac{dR}{dt} &= \frac{1}{T_{\text{inf}}} I - \mu R.\end{aligned}\tag{2}$$

The SEIR model includes an additional category  $E$  (Exposed) for people in the incubation period, when a person has only become infected but does not infect others. In addition to the dynamics of the category  $E(t)$ , the previous equations are supplemented with the parameter  $T_{\text{inc}}$ , which is average incubation period, and the system becomes

$$\begin{aligned}\frac{dS}{dt} &= -\frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} SI, \\ \frac{dE}{dt} &= \frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} SI - \frac{1}{T_{\text{inc}}} E, \\ \frac{dI}{dt} &= \frac{1}{T_{\text{inc}}} E - \frac{1}{T_{\text{inf}}} I, \\ \frac{dR}{dt} &= \frac{1}{T_{\text{inf}}} I.\end{aligned}\tag{3}$$

A variable  $D$ , which determines the mortality caused by infection [3] can be introduced in this model. Then the third equation of system (3) is replaced with the following:

$$\frac{dI}{dt} = \frac{1}{T_{\text{inc}}} E - \frac{1}{T_{\text{inf}}} (1 + \delta)I$$

and the fifth equation is added:

$$\frac{dD}{dt} = \frac{1}{T_{\text{inf}}} \delta I.$$

Here,  $\delta = Cfp / (1 - Cfp)$ , where  $Cfp$  describes the average proportion of fatal (lethal) cases for infected people.

The models described are then specified by adding new categories, such as people to be hospitalized and those who require the use of ventilators.

However, detailed models require appropriate settings. Even the simple model (3) requires data on the distribution of infected people into categories  $E$  and  $I$ , i.e., into those who are in the incubation period and those who have already left it.

Since testing cannot be considered sufficient, another parameter can be introduced for model (1) to adjust it. Because the official data only show a certain share  $P_{\text{inf}}$  of infected people (the number of detected cases), and instead of real values  $I$  only  $I_* = P_{\text{inf}} I$  are known, introducing the notation  $I = I_* / P_{\text{inf}}$ , we get

$$\begin{aligned}\frac{dS}{dt} &= -\frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} S \frac{1}{P_{\text{inf}}} I_*, \\ \frac{dI_*}{dt} &= \frac{R_0}{T_{\text{inf}}} \frac{1}{N_0} SI_* - \frac{1}{T_{\text{inf}}} I_*, \\ \frac{dR}{dt} &= \frac{1}{T_{\text{inf}} P_{\text{inf}}} I_*.\end{aligned}\tag{4}$$

In this case, the parameters  $R_0$ ,  $T_{\text{inf}}$ , and  $P_{\text{inf}}$  of system (4) can be estimated by the proximity of trajectories  $I_*(t)$  and  $S(t)$  to the respective observations. Since they can only be estimated by the trajectory  $I_*(t)$ , parameter  $P_{\text{inf}}$  does not affect the solution. Therefore, system (4) does not differ from system (1).

## ADJUSTING THE SIR MODEL

Let us return to system (1), which we will configure according to the observations of active infected people  $y_t$ ,  $t = 1, \dots, n$ , which are defined as the daily “total number of infected” minus “the total number of those who became ill.” Since the process under study is dynamically changing, as the criterion we will consider the summation of the weighted absolute value of deviation of the trajectory from the observation points. The values of the weights are selected

depending on the passage of time. For observations  $y_i, i=1, \dots, n$ , such a criterion has the form

$$\text{Crit} = \sum_{t=1}^n w_t |I(t) - y_t|, \quad w_t = 2t / (n(n+1)). \quad (5)$$

Let the population of Ukraine at the time of modeling be  $N_0 = 41$  million 858 thousand of people (according to some estimates as of March 1, 2020). The population dynamics can be taken into account according to model (2).

Note that the parameter  $T_{\text{inf}}$  determines the average period of time during which an infected person spreads the disease. This parameter is determined by the action of the virus and the human body, so it cannot change arbitrarily. According to the Oxford model for COVID-19,  $T_{\text{inf}} \approx 4.5$ . Therefore, we can limit the selection of the parameters by the condition  $4.4 \leq T_{\text{inf}} \leq 4.6$ , and selection of the parameters  $R_0$  and  $T_{\text{inf}}$  by the minimization of criterion (5), which can be performed using standard procedures for constrained global optimization.

A large number of models of such “deterministic” type have been developed by now; the behavior and properties of many of them have been studied. The main disadvantage of these models is the lack of stochastic approaches, although the process of epidemic development is essentially random.

## THE GUY KATRIEL MODEL AND ITS GENERALIZATIONS

This model is based on the results of [3]. It proposes a stochastic model of an epidemic with discrete time, in which the daily number of new diseases is binomially distributed depending on the number of diseases in the previous days. This model has the following advantages.

1. The model takes into account variations in the level of infectivity during the disease development, i.e., the probability of transmitting the disease on each day of the disease to individual patients.

2. The model is stochastic, which corresponds to the actual spread of infection among the population.

3. The model is easy to use, well-known formulas are used to calculate estimates by the method of maximum likelihood of the parameter (basic reproductive rate), which is equal to the average number of people infected by one ill person during the entire period of the disease. This parameter determines the epidemic spread rate. This assessment makes it possible to determine the parameter using the previous statistics of the daily number of new diseases to predict the further development of the epidemic.

The study [3] also considers the versions of models according to which the population is divided into subgroups, for example, by age or absence of acquired immunity.

An extended version of the model has also been developed [4–6].

1. An additional parameter is introduced: the probability of detecting the disease. Since in real life not all cases of the disease are detected or taken into account by statistics, the estimate of the parameter on the basis of previous data is not accurate; therefore, the parameter is used to adjust the statistics, taking into account a certain level of inaccuracy.

2. The ability to divide the epidemic duration into several periods with different values of the parameters at different stages is added. Estimates of the parameters at certain stages are not independent; therefore, it is necessary to maximize the approximation of the statistics for the entire epidemic. Division into stages is necessary when new quarantine measures are introduced (the parameter being changed in this case) or when the level of monitoring of the population is changed ( $DR$  is changed).

In cases of long-lasting epidemics (such as the COVID-19 pandemic), the disease spread dynamics may be seasonal (due to the effect of weather on the infectivity level and/or seasonal variation in the number of contacts among the population).

A program has been developed for parameter estimation and further simulation of the development of the epidemic.

## RANDOM MODELS OF EPIDEMIC SPREAD

Let us consider some models and methods used to determine the spread of epidemics as random rather than deterministic processes [7–10].

Let  $n$  be the number of people who fell ill. Every day, every sick person (regardless of other patients) may recover with probability  $\beta/n$  and die with probability  $\gamma/n$ . Also, patients receive a certain amount of medication  $x$  every day, which in our model is considered absolutely effective. The process ends when all patients either recover or die. The task is: for the given values of the parameters  $\gamma, \beta$ , and  $n$  find  $x$  for which the effectiveness of the provided drugs is the maximum.

Let us consider the features of the problem.

1. Unlike most epidemiological models, no new individuals are added to the class of sick people during the treatment process. This is the case, for example, when the disease is genetic or caused by a single catastrophe.
2. All the problem parameters ( $\gamma$ ,  $\beta$ ,  $n$ , and  $x$ ) are assumed to be positive numbers.
3. We assume that  $x$  does not change over time. In contrast to the parameters  $\gamma$ ,  $\beta$ ,  $c$ ,  $n$ , we regulate by the number of drug units  $x$ , but this number remains constant.

**Problem Solution.** Let  $N(t)$  be the number of sick people at time  $t$ . Consider the process given by the equation

$$M(t) = n - \sum_{i=0}^t (\xi(i) + x + \mu(i)), \quad M(0) = n, \quad (6)$$

where  $\xi(i)$  is the number of people who died at time  $i$  and  $\mu(i)$  is the number of those who recovered on their own at time  $i$ .

For any trajectory, we get

$$N(t) = \begin{cases} M(t), & M(t) > 0; \\ 0 & M(t) \leq 0. \end{cases}$$

Indeed,  $N(t)$  ceases to satisfy Eq. (6) only when the number of sick people becomes less than the daily number of drug units  $x$ . In what follows, we will use  $M(t)$  and will show that the results can be applied to  $N(t)$ .

Three lemmas were proved to find an efficient method of drug delivery [10].

**LEMMA 1.** For the mathematical expectation, the statement holds:

$$E[(M(t))] = \left(1 - \frac{\gamma}{n} - \frac{\beta}{n}\right)^t n \left(1 + \frac{x}{\gamma + \beta}\right) - \frac{nx}{\gamma + \beta}.$$

**LEMMA 2.** For the second moment  $M(t)$ , the statement takes place:

$$E[M^2(t)] = a_1^t \left(n^2 - \frac{a_4}{1-a_1}\right) + a_2 a_3 \frac{a_1^t - a_3^t}{a_1 - a_2} + \frac{a_4}{1-a_1},$$

where

$$\begin{aligned} a_1 &= \frac{n(n-\gamma-\beta) + (\gamma+\beta)^2}{n^2}, \\ a_2 &= \frac{-2xn^2 + n(2x+1)(\gamma+\beta) - (\gamma+\beta)^2}{n} \left(1 + \frac{x}{\gamma+\beta}\right), \\ a_3 &= 1 - \frac{\beta}{n} - \frac{\gamma}{n}, \\ a_4 &= \frac{2xn^2 - n(2x+1)(\gamma+\beta) + (\gamma+\beta)^2}{n} \left(\frac{x}{\gamma+\beta}\right) + x^2. \end{aligned}$$

**LEMMA 3.** Let  $a \in \mathbf{R}$ . Then

$$\frac{D[M(an)]}{(E[M(an)])^2} \rightarrow 0, \quad n \rightarrow \infty.$$

Using these lemmas, we obtain the statement [10].

**THEOREM 1.** Let  $T(n)$  be the total disease duration for  $n$  infected people, i.e.,  $T(n) = \min_{t \in N} \{t : N(t) = 0\}$ . Then

$$\forall \varepsilon > 0 \quad P(|T(n) - a_0 n| > \varepsilon n) \rightarrow 0, \quad n \rightarrow \infty,$$

where  $a_0 = \frac{\ln \frac{x+\gamma+\beta}{x}}{\gamma+\beta}$ .

**THEOREM 2.** For the mathematical expectation of the total number of drug units spent and the total number of deaths, the statement takes place:

$$(i) \eta(x, n) = xT(n) \Rightarrow E\eta(x, n) = xET(n) \approx a_0xn;$$

$$(ii) E[\xi(x, n)] = E\left[\sum_{i=1}^{T(n)} \xi(i)\right] = \frac{\gamma}{n} \sum_{i=1}^{T(n)} EN(t) \approx n \left[ (1 - e^{-a_0(\gamma+\beta)}) \frac{(x+\gamma+\beta)\gamma}{(\gamma+\beta)^2} - \frac{a_0x\gamma}{\gamma+\beta} \right].$$

The result of Theorem 2 makes it possible to find, by numerical methods, for the given values of the parameters  $\gamma$ ,  $\beta$ , and  $n$ , the value of  $x$  that maximizes the efficiency of providing drugs to patients.

## REGRESSION MODELS OF EPIDEMIC SPREAD, SWITCHING REGRESSIONS

More and more attention is being paid to modeling, analysis, and forecasting of objects of variable structure and (or) with time-varying parameters. Such processes take place in medicine, economics and technology.

There are several approaches to creating models of such objects. One of the most promising, in our opinion, is the use of switching regressions, where the switch points are unknown, so they need to be evaluated.

The essence of switching regression is that the regression parameters are not constant throughout the observation interval. They are constant on its subintervals, which are separated from each other by switch points. By estimating the switch points, it is possible to determine the time intervals on which structural changes of the object took place. There are two forms of switching regressions: with continuous regression line and with regression line that has discontinuities at switch points.

In studies of such regressions with discrete time, the contribution by P. Perron and co-authors is significant [7–9]. They showed the possibility of applying switching regressions in the economy.

In [11], a new class of switching regressions in continuous time was introduced and a method of their construction was proposed. Based on this study, the article [12] proposes a method for constructing switching regressions in discrete time.

Article [13] provides a preliminary statistical analysis of the spread of coronavirus disease in Ukraine based on the use of switching regression. The calculation procedure described there can be automated, which will allow real-time data processing.

Switching regression can also be used to determine the duration of treatment in people infected by coronavirus, as well as to monitor the course of various epidemics.

## CONCLUSIONS

We have considered some approaches to creating stochastic models of epidemic prediction, as well as the models, mathematical methods, and software for their implementation. The further development of these models for prediction and assessment of the spread of epidemics is associated with the use of regression models with continuous time and stochastic diffusion equations.

## REFERENCES

1. I. Brovchenko, “Developing a mathematical model of the COVID-19 epidemic spread in Ukraine,” *Svitohlyad*, No. 2 (82), 2–14 (2020).
2. M. J. Keeling and P. Rohani, *Modeling Infectious Diseases in Humans and Animals*, Princeton University Press, Princeton, NJ (2007).
3. G. Katriel, “Stochastic discrete-time age-of-infection epidemic models,” *Intern. J. of Biomathematics*, Vol. 6, No. 1, 999–1005 (2013).
4. G. Chowell, J. M. Hyman, L. M. A. Bettencourt, and C. Castillo-Chavez, *Mathematical Models in Population Biology and Epidemiology*, Springer, Dordrecht (2009).

5. W. Kermack and A. McKendrick, "Contributions to the mathematical theory of epidemics. I," *Bulletin of Mathematical Biology*, Vol. 53 (1–2), 33–55 (1991).
6. P. S. Knopov and O. V. Bogdanov, "Using a stochastic model to predict long-term epidemics," *Problemy Upravl. Inform.*, No. 3, 50–57 (2021).
7. R. Garcia and P. Perron, "An analysis of the real interest rate under regime shifts," *Review of Economics and Statistics*, Vol. 78, No. 1, 111–125 (1996).
8. J. Bai and P. Perron, "Estimating and testing linear models with multiple structural changes," *Econometrica*, Vol. 66, No. 1, 47–78 (1998).
9. J. Bai and P. Perron, "Computation and analysis of multiple structural change models," *J. of Applied Econometrics*, Vol. 18, No. 1, 1–22 (2003).
10. P. S. Knopov and O. V. Bogdanov, "Modeling of epidemics," *Kibernetyka ta Komp. Tekhnologii*, No. 2, 30–44 (2020).
11. P. S. Knopov and A. S. Korkhin, "Continuous-time switching regression method with unknown switching points," *Cybern. Syst. Analysis*, Vol. 56, No. 1, 68–80 (2020). <https://doi.org/10.1007/s10559-020-00222-z>.
12. A. S. Korkhin, "An approximate method of constructing a switching regression with unknown switch points," *Cybern. Syst. Analysis*, Vol. 56, No. 3, 426–438 (2020). <https://doi.org/10.1007/s10559-020-00258-1>.
13. P. S. Knopov and A. S. Korkhin, "Statistical analysis of the dynamics of coronavirus cases using stepwise switching regression," *Cybern. Syst. Analysis*, Vol. 56, No. 6, 943–952 (2020). <https://doi.org/10.1007/s10559-020-00314-w>.