



Machine learning-based mortality rate prediction using optimized hyper-parameter

Y.A. Khan^{a,d}, S.Z. Abbas^{b,d}, Buu-Chau Truong^{c,*}

^aSchool of Statistics, Jiangxi University of Finance and Economics, Nanchang, China

^bSchool of Mathematics and Statistics, Beijing Institute of Technology, Beijing 100081, China

^cFaculty of Mathematics and Statistics, Ton Duc Thang University, Ho Chi Minh City, Vietnam

^dDepartment of Mathematics and Statistics, Hazara University, Mansehra, Pakistan

ARTICLE INFO

Article history:

Received 8 May 2020

Accepted 6 August 2020

Keywords:

Prediction

Mortality rate

Hyper-parameter

Optimization

Covid-19 deaths rate

ABSTRACT

Objective and background: The current scenario of the Pandemic of COVID-19 demands multi-channel investigations and predictions. A variety of prediction models are available in the literature. The majority of these models are based on extrapolating by the parameters related to the diseases, which are history-oriented. Instead, the current research is designed to predict the mortality rate of COVID-19 by Regression techniques in comparison to the models followed by five countries.

Methods: The Regression method with an optimized hyper-parameter is used to develop these models under training data by Machine Learning Technique.

Results: The validity of the proposed model is endorsed by considering the case study on the data for Pakistan. Five distinct models for mortality rate prediction are built using Confirmed cases data as a predictor variable for France, Spain, Turkey, Sweden, and Pakistan, respectively. The results evidenced that Sweden has a fewer death case over 20,000 confirmed cases without observing lockdown. Hence, by following the strategy adopted by Sweden, the chosen entity will control the death rate despite the increase of the confirmed cases.

Conclusion: The evaluated results notice the high mortality rate and low RMSE for Pakistan by the GPR method based Mortality model. Therefore, the mortality rate based MRP model is selected for the COVID-19 death rate in Pakistan. Hence, the best-fit is the Sweden model to control the mortality rate.

© 2020 Elsevier B.V. All rights reserved.

1. Introduction

In the 1960s, the first identification of Coronaviruses occurred; their origination is still a mystery. The shape resemblance with crown-like proposed the name. They could infect humans as well as animals [1]. By functionality division, it is that kind that affects the sinuses, nose, and throat. The types NL63, 229E, HKU1, and OC43 of coronaviruses, like the common cold, usually cause illness of the upper respiratory tract. Many of the human population victimized by these types of viruses in the entire course of their lives. Such diseases lost in a minimal period. The general symptoms could be a headache, runny nose, dry cough, fever, sore throat, etc. [2].

It was December 31, 2019, when the world health organization announced several reasons of pneumonia cases in Wuhan city of Hubei province in China. This virus was noticed to be different from any other known type of viruses. For a new originated virus, we do not know the ways how it affects the peoples around, which raised significant concern. A few days later, the concerned authorities in China declared that they had identified a virus with a new shape. It was the coronavirus that causes the common cold like the MERS and SARS. The scientific name 2019-nCoV is suggested [3]. Now it is confirmed that new coronaviruses initiate the secretive respiratory sickness in Wuhan city. It is now clear that the secretive respiratory illness in Wuhan is undertaken by this virus vaguely associated with the SARS coronavirus abbreviated as SARS-CoV [4,5]. In humans, these viruses normally instigate through the surrounding of an infected entity by sneezing and coughing openly, close physical contact, touching of the objects with virus presence [6].

* Corresponding author.

E-mail addresses: abbassyedzaheer@tdtu.edu.vn (S.Z. Abbas), truongbuuchau@tdtu.edu.vn (B.-C. Truong).

The emerging trends of machine learning techniques possess a key feature in several fields of intelligence. It is based on the optimization of data by different algorithms to take preemptive measures. In data Sciences, it has the primary role for Data analytics. It makes us understand data and its processes better, make predictions based on historical data/experience data and categorize a group of data automatically called classification. It is observed that the Gaussian Based models have been commonly used in optimization applications [7]. In [8], an extensive comparative study was carried out between several surrogate models, comprising GPR, using simulation-optimization methodology with uncertainty parameters, in which it is concluded that the GPR models and their ensemble were efficient methods concerning the accuracy in prediction.

Similarly, the classical as supervised and unsupervised learning techniques feature as supervised learning techniques including Regression, Classification and Regression Trees (CART), and naive Bayes use labeled data to train the algorithms where input and output are known. The unsupervised learning techniques use unlabeled data to train the algorithms where the input of raw data given directly to these algorithms without knowing the output of that data [9,10]. The mortality models could be made more efficient by using Machine learning techniques. One such application is illustrated in [11], which is fully machine learning-oriented. Further, in [12], the authors extended the Lee-Carter model to multiple populations using neural networks.

In this research, the Gaussian Process Regression model with optimized hyper-parameter is used to develop the mortality models regarding COVID-19 for five different countries (Turkey, Spain, Sweden, France and Pakistan). Regression processes countered the flaws in these models. This model is fully featuring the machine learning techniques, which is capable of holding pieces of information which are not covered by standard models. We evaluate the enactment of these models, both in estimation and forecasting mortality rates, considering the available data for Pakistan.

The remaining paper is organized as follows:

Section 2 describes the proposed methodology for Predicting deaths due to COVID-19 for Pakistan by utilizing updated dataset samples. The discussion of the empirical study presented in Section 3. Section 4 concludes this work with possible enhancement as future work. All technical support is shown in Appendix A.

2. Mathematical scheme

Gaussian Process Regression (GPR) is a non-parametric kernel-based probabilistic model that can handle complex non-linear relations between response and predictor variables [7]. The Gaussian process is random and is considered as a set of random variables with a Gaussian joint multivariate distribution [13]. GP mainly based on a mean and covariance function. GP can achieve non-parametric regression function learning from noisy data, and it has Gaussian distributions over the data [13,17]. The predictable mean value is a linear combination by GP computation of the covariance function [13]. Among many others, one of the essential applications of GPs is Gaussian process regression (GPR). GPR is a probabilistic and robust non-parametric Bayesian model that defines a priori distribution of the likelihood over function space [8]. It is one of the most significant Bayesian machine learning methods that estimate the subsequent deterioration of non-linear regression by restricting the previous distribution to match the available training data [13].

The productivity of the prediction is a Gaussian distribution of probability and is characterized by its mean and variance. Variance is the confidence factor for the output's expected mean value [18]. Usually, a GPR model is provided with training data, and weighting

targets calculate its performance in terms of error between training and test input [19].

For an in-depth understanding of the underline mechanism, consider the training set $\{(u_i, v_i); i = 1, 2, \dots, m\}$, where $u_i \in \mathbb{R}^d$ and $v_i \in \mathbb{R}$, drawn from an unknown distribution. A GPR model addresses the question of predicting the value of a response variable v_{new} , (in our case no. of deaths due to COVID-19) given the new input vector u_{new} , (which is the number of confirmed COVID cases) and the training data. A linear regression model is of the form

$$v = \mathbf{u}^T \boldsymbol{\beta} + \varepsilon, \quad (1)$$

where $\varepsilon \sim N(0, \sigma^2)$. The error variance σ^2 and the coefficients $\boldsymbol{\beta}$ are estimated from the data. A GPR model explains the response by introducing latent variables, $f(u_i), i = 1, 2, \dots, m$, from a Gaussian process (GP), and explicit basis functions, \mathbf{g} . The covariance function of the latent variables captures the smoothness of the response, and basis functions project the inputs u into a p -dimensional feature space.

A GP is a set of random variables, such that any finite number of them have a joint Gaussian distribution. If $f(u), u \in \mathbb{R}^d$ is a GP, then given m -observations $u_1, u_2, u_3, \dots, u_m$, the joint distribution of the random variables $f(u_1), f(u_2), \dots, f(u_m)$ follows Gaussian distribution. A GP is defined by its mean function $m(u)$ and covariance function, $l(u, u')$. That is, if $f(u), u \in \mathbb{R}^d$ is a Gaussian process, then $\mathbf{E}(f(\mathbf{u})) = \mathbf{m}(\mathbf{u})$ and

$$\text{cov}[f(\mathbf{u}), f(\mathbf{u}')] = \mathbf{E}\{[f(\mathbf{u}) - \mathbf{E}(f(\mathbf{u}))][f(\mathbf{u}') - \mathbf{E}(f(\mathbf{u}'))]\} = \mathbf{l}(\mathbf{u}, \mathbf{u}').$$

Now consider the following model

$$\mathbf{g}(\mathbf{u})^T \boldsymbol{\beta} + f(\mathbf{u}),$$

where $f(u) \sim GP(0, l(u, u'))$, that is $f(u)$ are from a zero-mean GP with covariance function, $l(u, u')$. $\mathbf{g}(u)$ is a set of basis functions that transform the original feature vector u in \mathbb{R}^d into a new feature vector $\mathbf{g}(u)$ in \mathbb{R}^p . $\boldsymbol{\beta}$ is a p -by-1 vector of basis function coefficients. This model represents a GPR model. An instance of response v can be modeled as

$$P(v_i | f(\mathbf{u}_i), \mathbf{u}_i) \sim N(v_i | \mathbf{g}(\mathbf{u}_i)^T \boldsymbol{\beta} + f(\mathbf{u}_i), \sigma^2). \quad (2)$$

Hence, a GPR model is a probabilistic model. There is a latent variable $f(u_i)$ introduced for each observation u_i , which makes the GPR model non-parametric. In vector form, this model is equivalent to

$$P(v | \mathbf{f}, \mathbf{U}) \sim N(v | \mathbf{G}\boldsymbol{\beta} + \mathbf{f}, \sigma^2 \mathbf{I}), \quad (3)$$

where

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_m^T \end{pmatrix}, \mathbf{v} = \begin{pmatrix} v_1 \\ \vdots \\ v_m \end{pmatrix}, \mathbf{G} = \begin{pmatrix} \mathbf{g}(\mathbf{u}_1^T) \\ \vdots \\ \mathbf{g}(\mathbf{u}_m^T) \end{pmatrix}, \mathbf{f} = \begin{pmatrix} f(\mathbf{u}_1) \\ \vdots \\ f(\mathbf{u}_2) \end{pmatrix}.$$

The joint distribution of latent variables $f(u_1), f(u_2), \dots, f(u_m)$ in the GPR model is as follows:

$$P(\mathbf{f} | \mathbf{U}) \sim N(\mathbf{f} | \mathbf{0}, \mathbf{l}(\mathbf{U}, \mathbf{U}')).$$

Which is similar to a linear regression model, where $l(\mathbf{U}, \mathbf{U}')$ looks as follows:

$$\mathbf{l}(\mathbf{U}, \mathbf{U}') = \begin{pmatrix} l(\mathbf{u}_1, \mathbf{u}_1) & \cdots & l(\mathbf{u}_1, \mathbf{u}_m) \\ \vdots & \ddots & \vdots \\ l(\mathbf{u}_m, \mathbf{u}_1) & \cdots & l(\mathbf{u}_m, \mathbf{u}_m) \end{pmatrix}$$

The covariance function $l(u, u')$ is usually parameterized by a set of kernel-parameters or hyper-parameters ϑ , and often written as $l(u, u' | \vartheta)$ to explicitly indicate the dependence on ϑ . It is used to represents the covariance between pairs of random variables in GPR and can be written as

$$L_{ij} = l(\mathbf{u}_i, \mathbf{u}_j) = \alpha \exp \left\{ - \frac{\| \mathbf{u}_i - \mathbf{u}_j \|^2}{2\sigma_1^2} \right\}, \quad (4)$$

where

σ_1 – Characteristics length scale
 α – Single variance

Here σ_1, α are hyper-parameters.

In literature, the model given under Eq. (2) is called a surrogate for the objective function. The proxy is more comfortable with optimizing than the objective function. GP methods find the next set of hyper-parameters to evaluate the actual objective function by selecting the best hyper-parameters that perform on this surrogate function.

2.1. Parameter optimization

In regression models, the objective of parameter optimization is to find the parameters of a given algorithm that return the best performance on a validation set while training and testing the model [13]. It is mathematically represented as

$$\beta^* = \arg_{\min_{\mathbf{x} \in X} \mathbf{f}(\mathbf{x})}. \quad (5)$$

Here

$f(x)$ – represents an objective score to minimize the root mean squared error (RMSE) evaluated on the validation set

β^* – is the set of hyper-parameters which yields the lowest score of RMSE,

and

x – is any value in the problem domain X .

Even though hyper-parameter optimization is time-consuming, it yields good prediction accuracy than traditional regression models.

Proposed algorithm for hyper-parameter optimization

1. Initialize hyper-parameters for GPR model based on the problem.
2. An objective function of GPR model which takes in these hyper-parameters and outputs a RMSE score that has minimal value.
3. Define alternate model of the objective function.
4. Specify the selection criteria for evaluating hyper-parameters which have to choose next from the substitute model.
5. Maintain the history of (score, hyper-parameter) pairs used by the GPR algorithm to update the substitute model.
6. Repeat steps 2–5 until maximum iterations or time is reached.

2.2. Advantages of Gaussian process regression

- Gaussian process regression is probabilistic and robust non-parametric Bayesian model that defines a priori distribution of the likelihood over function space [14].
- It is one of the most significant BML-methods that evaluates the successive deterioration of non-linear regression by limiting the previous distribution to match the available training data [13].
- It has high flexibility and accurate prediction for processing small data set and also for high-dimensional data [13,15].
- It can have trained from noisy data using non-parametric regression function and, sidestepping simple parametric assumptions [16].

2.3. Polynomial regression model

Polynomial regression is a particular form of multiple linear regression models in which the maximum degree of the predictor variable is more than 1. In this technique, the best fit line is in the curve shape. It can be used to approximate a complex non-linear relationship [17].

The P th-order polynomial model in one variable is given by

$$y_i = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_k x^p + \varepsilon_i. \quad (6)$$

Where

y_i – the response in the i th trial

β_0 – intercept

β_i – regression coefficients $i = 1, 2, \dots, k$

x^i – values of predictor variables

ε_i – error term

A model with k explanatory variables x_1, x_2, x_3 is called multiple linear regression.

If a polynomial regression model is express in the form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}. \quad (7)$$

Then the methods of linear regression model estimation can be easily assumed for fitting the polynomial regression model.

2.4. Root mean square error (RMSE)

Root mean square error is a statistic used for accounting the average error size. It is the square root of the squared differences measured between estimated and actual observation and can be stated as

$$RMSR = \sqrt{\frac{1}{N} \sum (\hat{\theta} - \theta)^2}. \quad (8)$$

Here

N = Number of observations

$\hat{\theta}$ = Estimated value

θ = Actual value

In our case $N = 20$ days, and $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n)'$ are point forecasts of $\theta = (\theta_1, \theta_2, \dots, \theta_n)'$.

In this study, the GPR method is examined as a potential regression model for handling the non-linear variable in predicting the COVID-19 death rate for Pakistan. For comparison purposes, a polynomial regression method is employed. The final results indicate that the GPR method, although it is more time consuming, it proved to be more efficient in terms of the root mean square error (RMSE), a standard performance measure of regression models. Fig. 1 below clearly depicts the flow of the proposed methodology.

The best and worst-case scenarios for COVID-19 spread across the globe are taken to model the COVID-19 mortality rate model for five countries using COVID-19 daily confirmed cases time series data of those countries. These data are taken from www.ourworldindata.org, the most trusted website worldwide. Countries Sweden and Turkey are considered for the best-case scenario, whereas Spain and France are made for the worst-case scenario. The models for the mortality rate of COVID-19 are developed using Polynomial Regression Model and GP Regression model for these countries. RMSE used as a qualitative performance indicator to choose the best model. Finally, predict the mortality rate for Pakistan using the best death rate model.

2.5. Data and computational environment

The data used in this research is obtained from <https://ourworldindata.org/coronavirus>¹. Since January 21, 2020, the data is updated daily, with an increment of the number of infected people, the number of recovered and the number of deaths due to the coronavirus in Covid-19 infected countries across the world. In this study, we used recent data by taking daily observations into account from March 21, 2020, to May 10, 2020, for countries under consideration. The experimentation platform is a laptop with 2.7-GHz Intel CORE i5 and 8 GB of memory running 64-bits OS of MS Windows 10. All results quoted in this research paper were performed in MatLab, and Figures were construed in R, which are a

¹ European Center for Disease Prevention and Control (ECDC)

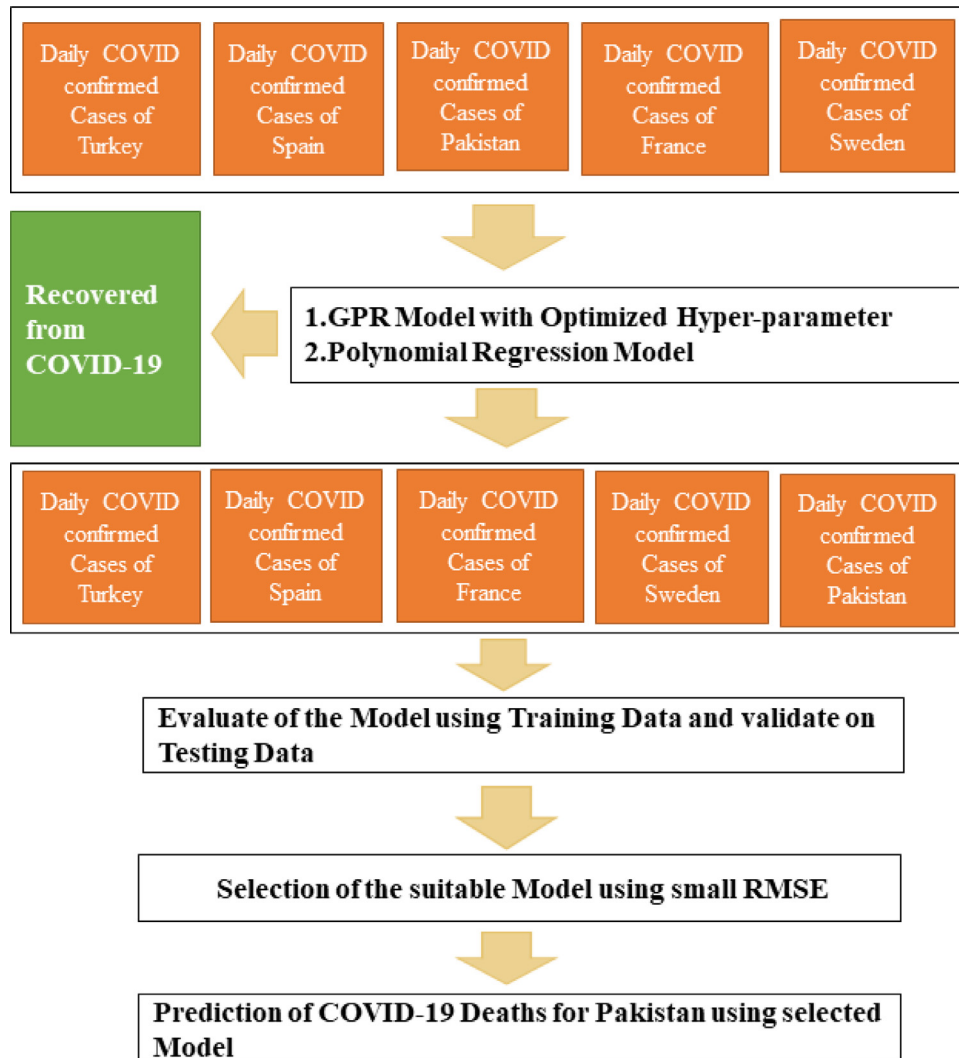


Fig. 1. Hierarchy of the anticipated methodology.

user-friendly software and publically available online. For hyper-parameter optimization, using the specified Algorithm takes about 56 s to reach the maximum iteration.

3. Application and discussion

An experiment is carried out to evaluate the efficiency of the proposed methodology for Deaths Prediction of the COVID-19 pandemic in Pakistan. Recent daily data of 51 days of the chosen countries were used for this purpose, details of which are given in the methodological section. To achieve the accuracy of the estimation and to validate the obtained result, we further subdivide the same filtered sample data and constitute training and testing sets of data. The training data contain 36 observations, while the testing set includes 15 days' data. A mortality rate prediction model was built for each selected countries using training data set where confirmed COVID cases are considered as the predictor variable and number of death due to COVID correspond to the response variable. All model was validated on testing data set. By the proposed algorithm for optimization, the parameter of the mortality rate prediction model was optimized, which was used for prediction purposes. Fig. 2 below represents the components of COVID-19 confirmed cases. A list of countries having a higher rate of COVID-19 confirmed cases are presented in Fig. 3. Fig. 4 represents

the number of COVID-19 death cases for countries such as Spain, Turkey, Pakistan, France, and Sweden.

Kernel Parameters such as σ_M , σ_F and σ of the best objective function for GPR models of various countries using Squared Exponential Method as kernel function for Gaussian process regression model are presented in Table 1.

Each country has a different tuning hyper-parameter value, which reflects that each country has a distinct trend in the COVID-19 spread and has the rate of increased confirmed cases also different. The regression loss for predicted value using the GPR model for five countries is also given in Table 2. Root Mean Square Error for the same are tabulated and are presented in the same table. Table 3 shows the RMSE value for the Polynomial Regression model of five countries.

Fig. 5 provides the graphical comparison of RMSE value for Polynomial Regression and GPR model, which helps in the selection of the best model among the two. It indicates that the GPR model has very low RMSE as compare to the Polynomial Regression model for five countries. Hence, Gaussian process regression with an optimized hyperparameter model is chosen as the best model which also supported by literature and used to predict deaths due to COVID-19 for Pakistan after lifting lockdown in the country where the chances of infection are very high.

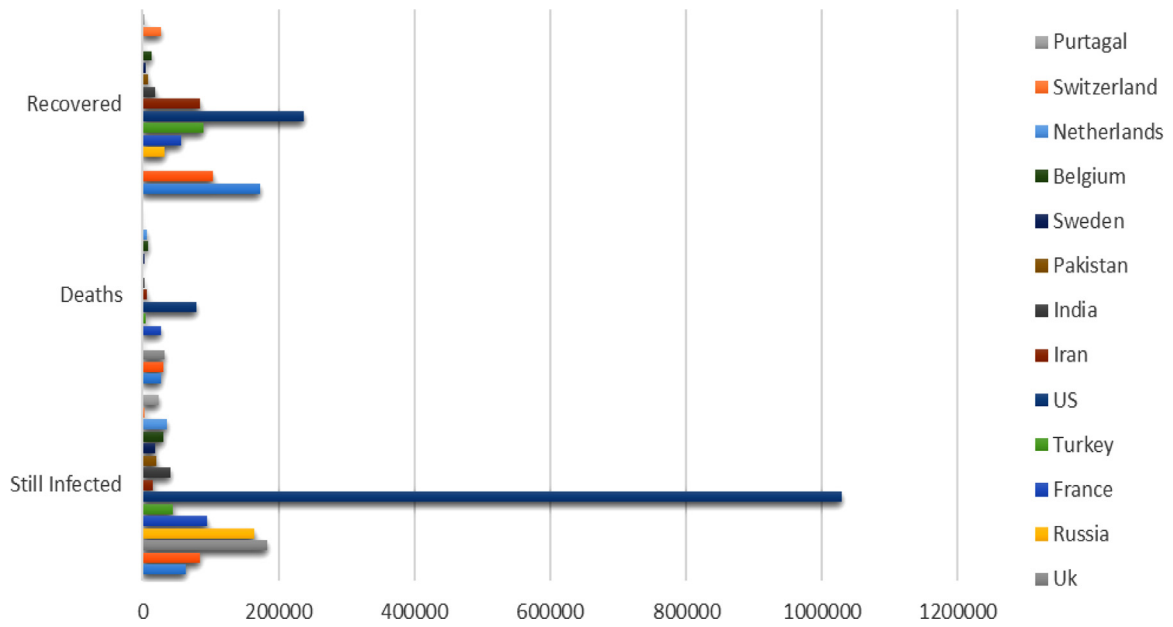


Fig. 2. Percentage distribution of COVID-19 infected cases, deaths and recovered from coronavirus as on May 10, 2020.

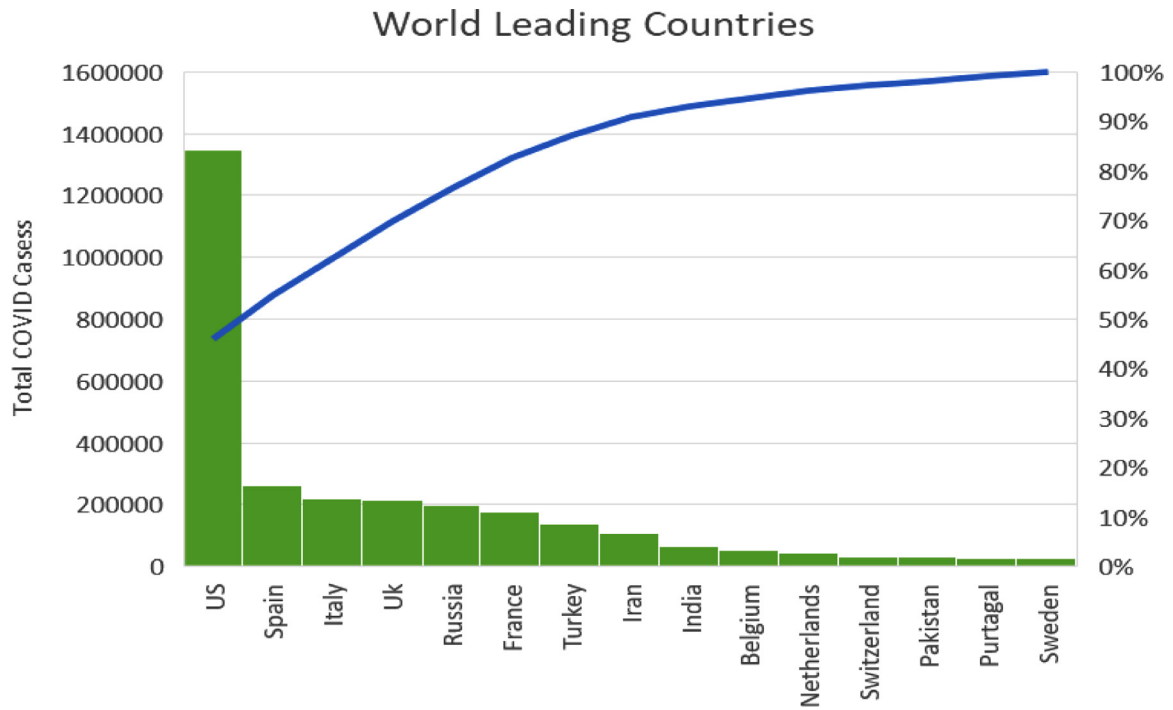


Fig. 3. Leading countries in the world with COVID-19 confirmed cases as on May 10, 2020.

Table 1
Estimated optimized hyper-parameters of GPR model for different countries.

Estimated kernels parameters	Turkey	Sweden	France	Spain	Pakistan
$\hat{\sigma}_m$	256876.3	589.8239	22374.43	79290.37	2667.455
$\hat{\sigma}_f$	7008.68	20.46376	332.0723	10241.42	50.8837
$\hat{\sigma}$	37.07613	1.110905	4.428102	87.57416	1.654951

Table 2
Loss and error value for GPR model.

Loss value	Turkey	Sweden	France	Spain	Pakistan
Regression loss	179.1848	1.878231	3.009694	1811.484	0.634793
RMSE	15.38199	2.370486	3.731184	42.56153	1.796739

Cumulative COVID-Deaths

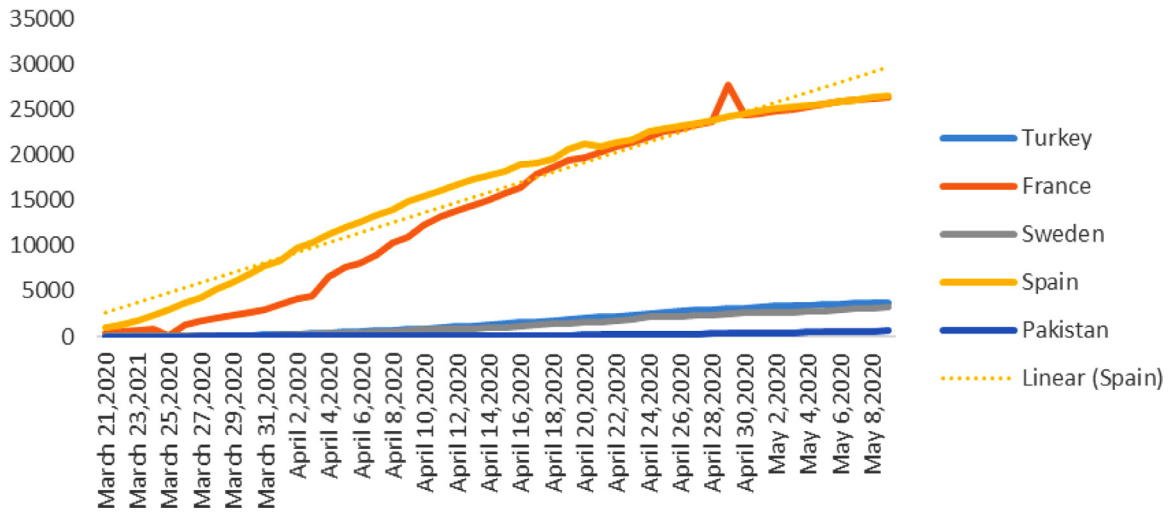


Fig. 4. Number of COVID-19 death cases for five countries.

RMSE COMPARISON

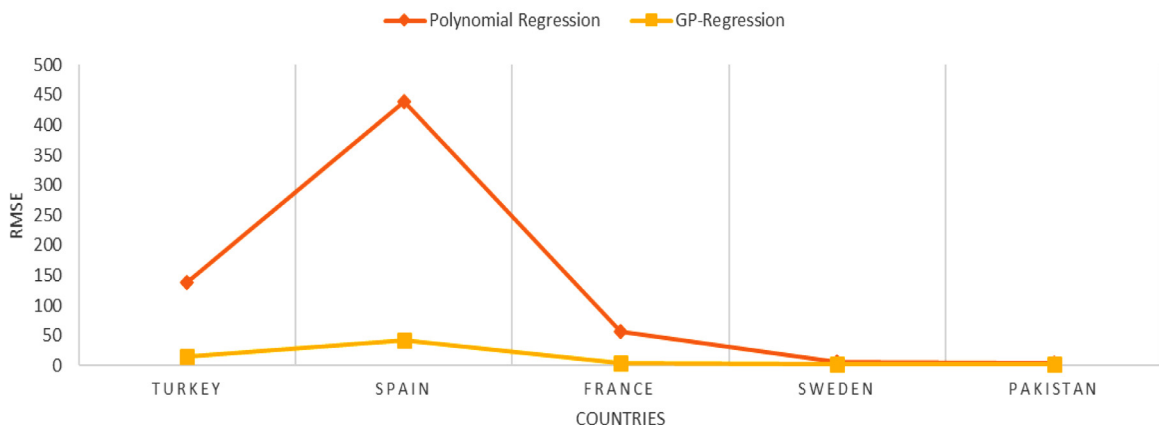


Fig. 5. Root mean square error comparison of the models.

Tables 4 and 5 represent the number of predicted COVID-19 death cases for countries like Turkey, France, Sweden, Spain, and Pakistan using the GPR and PR model, respectively. Fig. 6 illustrates the ranking of the COVID-19 mortality rate model of five countries.

It is known that Pakistan has a meager RMSE value because the magnitude of the confirmed cases is low due to strict lockdown. At the same time, the average death rate for Pakistan is the lowest

among the five countries taken for study. Sweden is the only country with having a shallow mortality rate without locking down the country. Sweden's RMSE value for the MRP model is the second-lowest in the table, which is selected as the best model. While the number of COVID-19 confirmed cases increases, the number of death cases for Sweden is deficient when compared to other countries in those tables. Although, for Turkey, the number of death

Table 4
COVID-19 deaths prediction using GPR model.

Number of confirmed COVID cases	No of deaths Turkey	No of deaths France	No of deaths Sweden	No of deaths Spain	No of deaths Pakistan
1000	29	56	3	89	55
2000	35	53	4	119	73
5000	65	61	10	264	99
10,000	141	79	39	710	560
20,000	350	530	45	1190	730
25,000	350	530	48	1190	735
30,000	410	370	53	1470	867
35,000	470	290	61	1750	961
40,000	530	284	72	2040	1010
45,000	590	280	81	2340	1021
50,000	650	280	102	2640	990

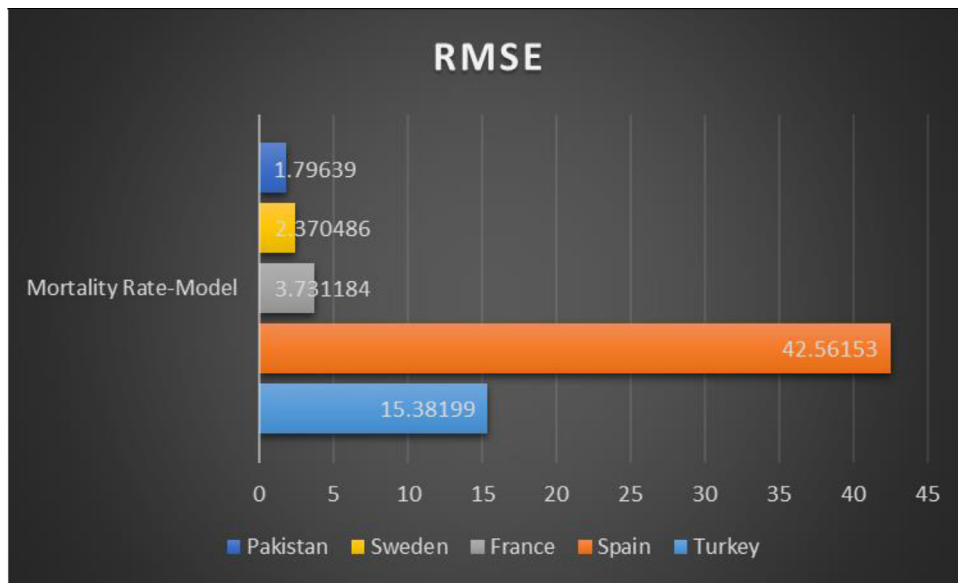


Fig. 6. RMSE value for five countries mortality rate model.

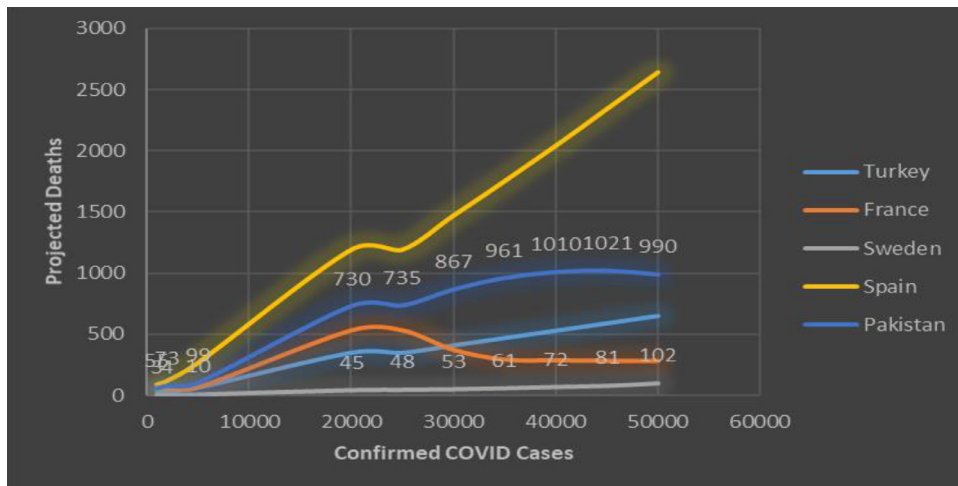


Fig. 7. Graphical presentation of COVID-19 death prediction using GPR model.

Table 3
RMSE value for polynomial regression model.

Loss value	Turkey	Sweden	France	Spain	Pakistan
RMSE	137.567	5.227635	57.21517	439.4404	3.445828

Moreover, the RMSE value for Sweden COVID- 19 Mortality model is lower than Turkey. Fig. 7 represents the number of predicted death cases in five countries. Table 5 shows the mortality rate prediction for France, Sweden, Turkey, Spain, and Pakistan, respectively. Fig. 7 represents the ranking of five countries based on the mean mortality rate. Based on the mean mortality rate, Sweden’s model is the best model for MRP.

cases increasing linearly but still low over the different levels of confirmed cases.

Table 5
COVID-19 deaths prediction using PR model.

Number of confirmed COVID cases	No of deaths Turkey	No of deaths France	No of deaths Sweden	No of deaths Spain	No of deaths Pakistan
20,000	200	760	50	460	650
25,000	200	760	50	460	650
30,000	300	920	90	1010	780
35,000	390	1080	140	1550	910
40,000	490	1230	180	2090	1040
45,000	580	1390	220	2640	1170
50,000	680	1540	270	3180	1310

4. Conclusion

In sighting the situation around the world, coronavirus becomes a biological bomb whose impact is more severe than a nuclear weapon. Although corona not only digests thousands of precious lives but also destroy the economy of the world. Almost all countries in the world practicing social distancing and observing lockdown from the last two months make the life of human being hell.

Although it is indeed essential to get an estimate of the financial losses to occur due to the deadly virus, which will be obtained years later, it is imperative to recognize the pattern of deaths. To minimize as much as possible, future losses of precious lives all over the world.

The proposed study investigated the advantages of Gaussian process regression using hyperparameter optimization, composed of the number of confirmed cases and deaths for the duration 21, March 2020 to May 10, 2020. A comparison is also given with the Polynomial Regression. Better performance is noticed for the Gaussian Process Regression model.

As outlined by Hong et al. [13] in his work, Gaussian process regression has the advantage of utilizing prior information to estimates the subsequent variation of a non-linear pattern under few assumptions, which is also validated in this study.

The mortality model for Turkey, Spain, Sweden, France, and Pakistan had been build using GPR and PR model. Their performance had been analyzed using RMSE value. The evaluated results notice the high mortality rate and low RMSE for Pakistan by the GPR method based Mortality model. Therefore, the morality rate based MRP model is selected for the COVID-19 death rate in Pakistan. Hence, the best-fit is the Sweden model for Pakistan to control the mortality rate.

4.1. Suggestions based on this study

Coronavirus will lasted for a long time. Although many countries are observing lockdown, they not keep it for a long time. After all, they will re-open sooner or later, because it severely affected the economy of all the countries. The only solution to survive in this station is to re-open as Pakistan did on May 9, 2020. As the coronavirus is contagious and transmitted from person-to-person, therefore the risk of being infected is still very high.

The social distancing, COVID-19 awareness and best self-hygienic practices are crucial factors to constraint COVID-19 deaths in these circumstances. Sweden is the only country having the best MRP model, although Sweden did not observe lockdown for a single day. Due to the reasons, the Sweden Mortality model is selected as the best model for predicting COVID-19 death cases. If Pakistan also adopts Sweden's strategies, Pakistan's COVID-19 death cases, while working in the presence of the coronavirus, can be restricted up to a small number.

Declaration of Competing Interest

The authors declared no conflict of interest regarding this manuscript submitted to Computer Methods and programs in Biomedicine.

Appendix A: Technical assistance

MatLab syntax/command used in reproduction of this research are as follows:

Syntax	
<code>gprMdl = fitrgp(X,y)</code>	<i>:- You can train a GPR model using the function</i>
<code>L = loss(gprMdl, Xnew,Ynew)</code>	<i>:- Regression loss for Gaussian process regression mode</i>
<code>gprMdl=fitrgp(Xtrain, ytrain,FitMethod','exact' 'Predict-Method','exact','KernelFunction', ardsquaredexponenti-aal','Standardize',1)</code>	<i>:- For kernel parameter using Squared/exponential</i>
<code>Regression loss=loss(gprMdl,Xtest,ytest)</code>	<i>:- For regression loss</i>
<code>ypredtest=predict(gprMdl,xtest)</code>	<i>:- Prediction from test data</i>
<code>ypred = predict(gprMdl,Xnew)</code>	<i>:-GPR Prediction of response</i>

Note: y "response" is no of deaths, and X "predictor" is no of confirmed cases. R codes of parameter optimization the Algorithm will be provided on personal request.

References

- [1] Coronavirus, WebMD. www.webmd.com/lung/coronavirus#1.
- [2] Common Human Coronaviruses, CDC,2020. <https://www.cdc.gov/coronavirus/about/symptoms.html>
- [3] WHO, 2019. <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>.
- [4] New Virus Discovered by Chinese Scientists Investigating Pneumonia Outbreak. Available online: <https://www.wsj.com/articles/new-virus-discovered-by-chinese-scientists-investigating-pneumoniaoutbreak-11578485668> (accessed on January 21 2020)
- [5] New-Type Coronavirus Causes Pneumonia in Wuhan: Expert. Available online: <http://www.xinhuanet.com/english/2020-01/09/c138690570.htm> (2020)
- [6] Transmission of Coronavirus, CDC. www.cdc.gov/coronavirus/about/transmission.html
- [7] Ebden M., Gaussian Processes for Regression: A Quick Introduction, No. 2008.
- [8] C.E. Rasmussen, Gaussian processes in machine learning, Adv. Lect. Mach. Learn. 14 (2) (2004).
- [9] R. Richman, M.V. Wüthrich, A Neural Network Extension of the Lee-Carter Model to Multiple Populations, SSRN, Rochester, 2018.
- [10] B.A. Tolson, S. Razavi, D.H. Burn, Review of surrogate modeling in water resources, Water Resour. Res. 48 (7) (2012), doi:10.1029/2011wr011527.
- [11] P.C. Deka, N.S. Raghavendra, Multistep ahead groundwater level time-series forecasting using Gaussian Process Regression and ANFIS, in: R. Chaki, A. Cortesi, K. Saeed, N. Chaki (Eds.), *Advanced Computing and Systems for Security: Advances in Intelligent Systems and Computing*, 396, Springer, New Delhi, India, 2016.
- [12] B. Datta, D.K. Roy, Trained meta-models and Evolutionary Algorithm based multi-objective management of coastal aquifers under parameter uncertainty, J. Hydro Inform. 6 (2018) 1247–1267.
- [13] S. Hong, Z. Zhou, Application of Gaussian process regression for bearing degradation assessment, in: Proceedings of the 2012 Sixth International Conference on New Trends in Information Science and Service Science and Data Mining (ISSDM).
- [14] C.E. Rasmussen, Gaussian processes in machine learning., Adv. Lect. Mach. Learn. 14 (2), 2004
- [15] J. Kocijan, K. Ažman, Dynamical systems identification using Gaussian process models with incorporated local models, Eng. Appl. Artif. Intell. 24 (2) (2011).
- [16] M. Huber, Recursive Gaussian process: online regression and learning, Pattern Recognit. Lett. 45 (2014) 85–91.
- [17] C.K.I. Williams, C.E. Rasmussen, Gaussian Processes for Machine Learning, MIT Press, Cambridge, 2006.
- [18] A. Saxena, K. Goebel, B.Saha S.Saha, Distributed prognostic health management with gaussian process regression, in: Proceedings of the 2010 IEEE Aerospace Conference, 2010, pp. 1–8.
- [19] J. Kocijan, K. Ažman, Dynamical systems identification using Gaussian process models with incorporated local models, Eng. Appl. Artif. Intell. 24 (2) (2011).