

Machine Learning model to predict the number of cases contaminated by COVID-19.

Allae Erraissi (✉ erraissi.allae@gmail.com)

Universite Hassan II Casablanca <https://orcid.org/0000-0002-2826-5495>

Mohamed Azouazi

Universite Hassan II Casablanca

Abdessamad Belangour

Universite Hassan II Casablanca

Mouad Banane

Universite Hassan II Casablanca

Case study

Keywords: Machine Learning, Spark ML model, COVID-19 Artificial Intelligence model, coronavirus, predicting COVID-19

Posted Date: May 13th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-23330/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.
[Read Full License](#)

Version of Record: A version of this preprint was published at International Journal of Computing and Digital Systems on November 25th, 2021. See the published version at <https://doi.org/10.12785/ijcds/100189>.

Abstract

Introduction:

This paper presents a dedicated machine learning model to predict the number of cases infected by the Corona Virus; the case of Morocco was chosen to validate this study.

Case description:

Completely realized in Spark ML with the 'Scala' language and tested for a certain number of algorithms generated on datasets coming from dedicated sources to gather Covid19 data in the world.

Discussion and Evaluation:

The results show the possibility of achieving better scores prediction after using the proposed method. We tested our model on the case of China and the results were relevant.

Conclusion

The proposed Machine Learning model can be applied to data from any country in the world. We have applied it in this paper to the case of Morocco and China. We are sending this work to the world to help them fight this 2019 Corona Virus pandemic.

1. Introduction

The field of artificial intelligence is undergoing a very big evolution these days and it risks to categorically change the world through us [1]. Fraud detection, recommendation systems, facial recognition, and decision-making automation are just a few of a variety where artificial intelligence has already taken place.

As its name suggests, artificial intelligence is intelligence (a set of skills) acquired by machines in order to function and react like humans. Machine learning [2] is a branch of artificial intelligence that concerns the design and development of algorithms that allow a computer (a machine in the broad sense) to learn to perform very complex tasks without having been explicitly programmed.

So this study aims to propose a model based on artificial intelligence, especially on machine learning [2]. This model uses a series of predefined instructions which leads a machine to generate a prediction of the number of cases infected with the Coronavirus COVID-19 in Morocco per day. Then we will test this model in the case of China since it has recorded enough data that can be used for verification. As part of this study, we used the library dedicated to distributed learning methods Spark ML [3]. This paper is detailed

as follows: Sect. 2 presents the related work on which we are concerned to carry out this study. Then, in Sect. 3 we talk about the Apache Spark ML solution dedicated to the use of Machine Learning algorithms and methods on Big Data [4, 5]. Section 4 presents our proposed Machine Learning model to predict the number of cases infected with the coronavirus (COVID-19). Finally, Sect. 5 presents the results obtained after the application of our model on the case of Morocco and China as well as a summary discussion of the work carried out.

2. Related Work

The use of machine learning techniques better known by its Anglicism "Machine Learning" is becoming more and more popular. It applies to a wide variety of fields, from image processing to processing human language, including computer security and the video game industry [2]. One of the main purposes of applying machine learning techniques is prediction. Any machine learning process can be subdivided into the following steps: data preprocessing, model design and training, and model evaluation. The preprocessing phase is still known as "Pre-processing" very often consists of choosing the most informative variables from the set of incoming vectors and normalizing them if necessary. This technique is also called "feature selection / feature engineering". The purpose of the selection of informative variables is to eliminate nonessential variables which can be considered as noise by the model and thus reduce its accuracy of prediction. Several research studies have been carried out in recent decades to propose either algorithms or criteria for the selection of informative variables.

In [6], authors are offered a technology for extracting informative variables based on the ontology of the field studied. The first phase of this technology consists of a preliminary extraction by the domain expert, of the variables which he deems useful without any other restriction than the variables chosen must be objects or attributes of the ontology. The second phase, still described in [6], consists of aggregating and filtering the preselected variables during the previous phase. The mathematical translation of this procedure is as follows: either $\Omega = \{ \omega_1, \dots, \omega_n \}$ - the class labels, $X = \{ X_1, \dots, X_n \}$ - the set of variable identifiers. x_s^i - a value taken by X_i . N_i - the definition domain of $X_i (x_s^i \in N_i)$. Let Σ be the set of learning vectors.

For a x_s^i value of X_i : $x_s^i \in N_i$ and in the class ω_k an aggregate $N_i(\omega_k) \in N_i$ is defined as follows:

$$x_s^i \in N_i(\omega_k) \text{ if and only if for } \forall \omega_v \in \Omega, v \neq k : p(\omega_k / x_s^i) > p(\omega_v / x_s^i) + \Delta, (1)$$

Where Δ a positive real number defining the dominance threshold.

At the end of phase 2, the author introduces the unitary predicates $B_i(\omega_k)$ which take the value 'true' (1) if and only if $x_s^i \in N_i(\omega_k)$. The final results of the second phase will, therefore, be the aggregates $N_i(\omega_k)$, the unitary predicates $B_i(\omega_k)$ and $i \in I(\omega_k)$: $I(\omega_k)$ - a subset of the indexes of the variables X_i that passed the test of inequality (1).

During Phase 3 a probabilistic approach was introduced to determine cause-and-effect dependencies in the conjunction of the predicates $B_i(\omega_k)$ and $i \in I(\omega_k)$, $\Omega = \{ \omega_1, \dots, \omega_n \}$ and $\omega_j \in \Omega$.

Let $\{ N_i(\omega_j) \} | i \in I(\omega_k), j = 1, \dots, m$ - the set of aggregates where $p(N_i(\omega_j)) = |N_i(\omega_j)| / |N_i|$, where $| |$ represents the standard of the corresponding assembly. It is obvious that: $p(B_i(\omega_j)) = p(N_i(\omega_j))$.

In the following three filters will be introduced to select the predicates:

Filter 1 :

$$I(B_i(\omega_j), \omega_k) = |p(B_i(\omega_j)\omega_k) - p(B_i(\omega_j))p(\omega_k)| / [p(B_i(\omega_j))p(\omega_k)] \geq \delta_{\min} > 0-a \quad (2)$$

Filter 2:

$$R(\widehat{B}_i(\omega_j), \omega_k) = |p(\omega_k / \widehat{B}_i(\omega_j)) - p(\omega_k / \widehat{B}_i(\omega_j))| \{ p(\widehat{B}_i(\omega_j)) [1 - p(\widehat{B}_i(\omega_j))] \} = |p(\widehat{B}_i(\omega_j) \omega_k) - p(\widehat{B}_i(\omega_j)) p(\omega_k)| / \{ p(\widehat{B}_i(\omega_j)) [1 - p(\widehat{B}_i(\omega_j))] \} \geq \delta_{\min}, \delta_{\min} > 0-a \quad (3)$$

Filter 3:

$$p(\widehat{\omega}_k / \widehat{B}_i(\omega_j)) = p(\widehat{B}_i(\omega_j) \widehat{\omega}_k) / p(\widehat{B}_i(\omega_j)) \geq \gamma_{\min} \quad (4)$$

The authors of [7] in Chap. 6 entitled "Dimensionality reduction" discuss the PCA (Principal Component Analysis [8]) algorithm based on the analysis of linear dependence. The idea is to replace the redundant variables with a new one that sums up the information contained in the original vector space. One of the first concepts to appear is SVD (Singular Value Decomposition) [9]. Let X be a matrix of dimension $n * d$; n - being the number of lines and d - the dimension of the initial vector space to be reduced. Let x be a column vector (transposed from one of the rows of X) and let v - be one of the new vector entities (main components) that we are looking for. According to the SVD theorem, matrix X can be decomposed sum follows $X = U \Sigma V^T$, U and V being orthogonal matrices such that $U^T U = I$ et $V^T V = I$. Σ a diagonal matrix containing the singular values of X In the following, the author develops the main stages of PCA [8]. The first step is to refocus the data in the initial vector space. The second step is the linear projection of the initial data vector x into new vectors v . The third step is to maximize the variance of the coordinates. The equation of the coordinate vector after projection is $z = Xv$.

The objective function for the main components is $\max_w W^T W$, where $W^T W = 1$. It appears that the optimal form of W is the eigenvector of $X^T X$.

In [10] the authors based on Bayes' theorem:

$$P(A_j | B) = \frac{P(A_j)P(B|A_j)}{\sum_{j=1}^n P(A_j)P(B|A_j)} \quad (5)$$

and on Information theory [11], in particular, the notion of self-information $I(x_i) = -\log_2 P(x_i)$ proposes an algorithm supposed to improve the forecasting accuracy of the NBC (Naive Bayesian Classifier) model [12].

The algorithm divided into two stages consists of its first phase of calculating the weight of each variable $W_{F_i} = -\log_2 P(C|F_i)$. The second phase consists of selecting the variables whose weight W_{F_i} would have exceeded a certain threshold δ determined by the user.

A new method called LFE (Learning Feature Engineering) is proposed in [13]. At the heart of the LFE [14], we find a set of multilayer classifying perceptrons. Each perceptron corresponding to a transformation. LFE takes as input a dataset and recommends a set of paradigms allowing reconstructing a subset of the informative variables of the initial dataset. Each paradigm is made up of a transformation and the ordered list of variables for which the transformation was most efficient.

In the following, we will propose a model based on the machine learning methods already mentioned to predict the number of COVID-19 cases in Morocco per day. Through the use of the library dedicated to distributed learning methods Spark ML [3].

3. Spark ML

Spark ML is an essential software brick of the Apache Spark platform [3]. The Spark ML API is dedicated to the implementation of learning methods. Spark ML offers services covering data preparation, enrichment, development of learning methods, and deployment.

The Spark ML API offers several transformers and algorithms for developing and fine-tuning methods of:

- Classification
- Regression
- Clustering
- Recommendation

3.1. Concepts for implementing Spark ML algorithms

Here we cite the key concepts for implementing a Spark ML algorithm [15]:

- **DataFrame:** is a distributed, column-oriented data structure suitable for statistical data or learning methods.
- **Estimators:** these are algorithms aimed at adjusting a function between the explanatory variables and the target variable. The algorithms are based on the data of a DataFrame.
- **Evaluators:** measures to assess the performance of a learning method.
- **Transformers:** these are algorithms used to transform one DataFrame into another.
- **Parameters:** Transformers and Estimators share the same API to specify parameters.

- **Pipeline:** a pipeline consists of several Transformers and Estimators for the implementation of a ML workflow.
- **Cross-Validation:** a technique to develop a learning method.

3.2. Pipeline & PipelineStage

A learning method involves processing and learning phases. A Pipeline, in Spark MLIB [16], is a sequence of elements of type pipelineStage. Each PipelineStage consists of a Transformer or an Estimator. Therefore, the order of steps (PipelineStage) is important.

3.3. Transformation Stage

A transformer transforms data. In Spark ML, for a transformation step, the transform() method is called on a DataFrame and returns a DataFrame. Transformer is an abstract class that extends the PipelineStage abstract class.

Some examples of Transformers:

- **Tokenization:** segment text into a list of words.
- **Binarization:** transform numerical values into binary values.
- **Normalizer:** normalize vectors using p-norm.
- **N-gram:** calculate the next character from an article sequence q
- **MinMaxScaler:** transform a quantitative variable to a quantitative variable with values in a given interval.
- Etc.

3.4. Estimation Stage

An estimator allows you to evaluate or approach an unknown parameter. A statistical estimator aims to evaluate an unknown parameter via a sample of data. In Spark ML, for an estimation step, the fit() method of an Estimator is called on a DataFrame to produce a Transformer.

Some examples of estimators:

- **StringIndexer:** code categories.
- **IndexToString:** decode indexes.
- **VectorIndexer:** identify categorical variables according to a frequency threshold.
- **QuantileDiscretizer:** transform a quantitative variable into a qualitative variable using the quantile.
- **StandardScaler:** normalize a vector using the variance.
- Etc.

3.5. Evaluate a learning method

To evaluate the performance of a learning method, we use dedicated measures. For a given individual, an evaluation measure uses the adjusted value and the reference value. We have measures for regression methods and classification methods.

Evaluate a Learning method: Classification

It aims to measure the quality of a classification method. It indicates whether the method succeeds in correctly classifying individuals.

Table 1
Confusion matrix [17].

		Predicted	
		Positive	Negative
Real	Positive	TP	FN
	Negative	FP	TN
The evaluation measures are:			

- **Sensitivity:** the rate of correctly identified positive cases.
- **Specificity:** the rate of correctly identified negative cases.
- **Accuracy:** the rate of correctly identified positive or negative cases.
- **F-measure:** the rate that combines sensitivity and specificity.
- **Evaluate a learning method: regression**

It aims to measure the quality of a regression method by calculating a global error. An error close to zero indicates the relevance of the model. The main measures are [15]:

- **MSE:** Mean Squared Error.
- **RMSE:** this is the square root of MSE (root Mean Squared Error).
- **R:** Coefficient of determination.
- **MAE:** mean Absolute Error.
- **MAPE:** Mean Absolute Percentage Error.

3.6. API for the evaluation of a learning algorithm

The Spark ML API provides three implementations for the Evaluator abstract class. The classes are: RegressionEvaluator for regression methods, BinaryClassificationEvaluator for binary classification methods, and MulticlassClassificationEvaluator for multi-class classification methods.

4. Case Description

In order to better propose a Machine Learning model to predict the number of cases of the Corona COVID-19 [19] virus in Morocco over the coming weeks.

As it is established that, during epidemics, the daily number of infected follows a bell-shaped diffusion curve which gives, in cumulative, an S curve, the choice of a logistic function was obvious.

The data we had captured the results of tests by health authorities in different countries around the world. But we knew without question that these data did not reflect reality. Some countries have chosen not to be over-tested either by choice of public health policy or because they simply did not have the means. And even for the countries that tested more significantly, it can be reasonably assumed that, in the first periods, the number of tests concerning the number of real cases was lower than in the following periods. The awareness of the dangers increases with time.

In order to verify the relevance of our proposed ML model, we conducted a series of tests on data from the site of the Moroccan Ministry of Health (www.covidmaroc.ma) [20].

So the dataset first gathered was of the type (CSV file [21]).

```
Date;02/03/2020;03/03/2020;04/03/2020;05/03/2020;06/03/2020;07/03/2020;08/03/2020;09/03/2020;10/03/2020;11/03/2020;12/03/2020;13/03/2020;14/03/2020;15/03/2020;16/03/2020;17/03/2020;18/03/2020;19/03/2020;20/03/2020;21/03/2020;22/03/2020;23/03/2020;24/03/2020;25/03/2020;26/03/2020;27/03/2020;28/03/2020;29/03/2020;30/03/2020;31/03/2020;01/04/2020;02/04/2020;03/04/2020;04/04/2020;05/04/2020;06/04/2020;07/04/2020;08/04/2020;09/04/2020;
Number of cases;1;0;0;1;0;0;0;1;0;3;0;2;10;10;9;7;10;9;16;17;19;28;27;55;50;70;57;77;77;61;37;54;83;128;102;99; 64;91; 99;
```

This CSV type file must be transformed into a DataFrame, because we are going to use Spark ML which is based on the DataFrame API. To do this, we used the version 3.0.0 of Spark. And as for programming language, we used version 2.12.10 of the Scala programming language [22]. You need to make sure that Java is installed on your machine. All of these tools were installed on a machine with a 2.70 GHz Intel (R) Core (TM) i7 processor. With 1 TB storage space and 16 GB RAM memory. The Java version "jre1.8.0_241" is used for the development of our Machine Learning model.

So, we start with the command: spark-shell. The following figure shows the execution result and the configuration used:

Then we transform the input file of type CSV to a DataFrame with the following command:

Here is the result of the creation of the DataFrame from a data source which is of type CSV in this work:

Now that the DataFrame is well created, it will be used by the Spark Machine Learning library to predict the number of cases of the Corona Covid-19 virus in Morocco per day. This part shows an extract from the code of the model that we have developed for this study:

```

import org.apache.spark.ml.{Pipeline, PipelineModel}

import org.apache.spark.ml.classification.{DecisionTreeClassifier, DecisionTreeClassificationModel}

import org.apache.spark.ml.feature.{StringIndexer, VectorAssembler}

// load CovidMorocco data and rename columns

val MoroccoCases = spark.read.option("header", true).option("inferSchema", "true") .option("delimiter";
").csv("C:/Users/Allae Erraissi/Desktop/CoronoVirusCasesMorocco.txt").toDF("Date", "Number")

// Create assemble to group the columns

val assembler = new VectorAssembler().setInputCols(Array("MorningDate", "EveningDate"))
.setOutputCol("ResultOfDay")

// transform the Species variable into an indexer

val indexer = new StringIndexer().setInputCol("Number").setOutputCol("NumberIndex")

// Create the learning model DT

val dt = new DecisionTreeClassifier().setLabelCol("NumberIndex").
.setPredictionCol("NumberPredictCol").setFeaturesCol("features")

// create a pipeline: assembler, indexer, dt

val pipeline = new Pipeline().setStages(Array(assembler, indexer, dt))

// adjust the model

val model = pipeline.fit(MoroccoCases)

// evaluate the model on learning data

val predictions = model.transform(MoroccoCases)

// adjust results

predictions.show

```

The execution of the COVID 19 model on the sample of data from Morocco gives us the results of this form:

Then we present a part of code which presents the techniques relating to the evaluation of our learning method:

```

// Import of class for evaluation
import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator

// Create the assessment instance

// Fix the label column

// Fix the prediction column

// Set the evaluation measure

val evaluation = new MulticlassClassificationEvaluator().setLabelCol("NumberIndex")
.setPredictionCol("NumberPredictCol").setMetricName("accuracy")

// calculate the value of Accuracy

val ACCURACY = evaluation.evaluate(predictions)

// display the value of ACCURACY

print(s" Accuracy= ${ACCURACY}")

```

To evaluate our proposed Machine Learning model, we followed all of these steps :

- Load data.
- Prepare the data.
- Adjust the ML model.
- Initiate the evaluator and position the params.
- Evaluate the method via the Accuracy measure.

5. Discussion And Evaluation

In order to verify the relevance of the proposed methods, we had conducted a series of tests on data from the site of the Moroccan Ministry of Health (www.covidmaroc.ma) [20]. The classification algorithms used are: 'Decision Tree' (DTC) [23], 'Gaussian Naive Bayes' (NBC) [24], 'Support Vector Machine' (SVM) [25], 'Logistic Regression' (LR) [26], 'Random Forest'(RFC) [27] and finally 'Voting Classifier' (VC) [28]. The characteristics of the input data are shown in Fig. 9. These data show the number of COVID-19 cases in Morocco per day.

To make predictions, whether we have the partial differential equation (PDE) model [29] which determines how the different parameters involved in the phenomenon interact with each other or ordinary differential equations (ODE) [30] when we only have two quantities, one in function on the other which act on the phenomenon. In this case, the data collected is injected into the equations and it is the resolution of the equations which gives the forecasts. The examples illustrating this are numerous and diverse, we can cite the example of the weather. What you also need to know is that the models are approximate and their resolution is done by approximation methods. If the forecast in the case of the weather is more and more

precise. This is because the models keep getting better and the resolution methods are more and more efficient and fast so they give fewer errors. Another thing to clarify is that the phenomena for which we have PDE or ODE equation models are phenomena that have a certain regularity, even if the phenomenon is complex, the regularity can be concealed. If we don't have the equation model, can we make predictions? The answer is yes if the phenomenon is random and if we can find the law of probability that governs it, we can use it to make forecasts. Now, if we don't have an equation model or a probability law can we forecast anyway? In this case, we use extrapolation methods to extrapolate the history. The least-squares method consists in looking for a polynomial that best approaches the data represented in the form of a point cloud. When we say who "approaches the best" we would need a means of measurement, a criterion to select the best approximation, this is what is called standard. Following the course of evolution, this polynomial is considered as the mathematical model that governs the evolution of the phenomenon, we speak of the polynomial model. The polynomial model does not work for all phenomena.

Case of COVID-19:

We are faced with a completely random phenomenon. The number of people infected depends on the trajectory of each infected individual and the people they met on their way and also depends on the people who have touched anywhere the virus is found. Assuming that the model follows a Gaussian law is not meaningless, without doubt, the evolution continues its growth, relatively fast, it will not grow indefinitely, it will go through a peak then it will decrease to 0 Maybe not with the same speed. The whole problem is to approach this Gaussian. We are going to show on a theoretical example that polynomials cannot "efficiently" approach a Gaussian using a reduced number of data. And therefore the polynomials cannot be used as predictors in this case.

We can clearly see in the graph that the long term cannot be predicted by polynomials. To remedy this, the use of our Machine Learning model with Spark ML to give satisfactory results. We have given as input to our model several parameters related to the situation of our paid in this confinement period.

We have taken into consideration the percentage of illiteracy in Morocco and its relationship with the application of the instructions of the various ministries of the country which oblige people to isolate themselves and stay at home. Also, consideration has been given to the approaching of the holy month of Ramadan for Muslims. It is, for this reason, our Machine Learning model predicted two possible scenarios for the COVID-19 [31] in Morocco.

As already mentioned, we have applied this model to data extracted from the site: www.marocovid.com. The result is shown in the following figure:

We applied our model to detect the number of cases of the Coronavirus in Morocco on the DataFrame. We have defined two possible scenarios for our study. The first scenario consists in that the Moroccan population has respected the instructions of the ministries responsible for fighting against the COVID-19 [31].

With the application of our Covid-19 model, our estimate is as follows for the first scenario:

We can see that the growth rate has decreased significantly which has considerably influenced the peak and also the date of attenuation.

The second scenario was defined based on the place of the holy month Ramadan among Muslims in general, and especially for Moroccans. So during this month, people will not respect the instructions of the Moroccan ministries. To take this information into account, we have changed the weighting percentages for the country's illiteracy, as well as the percentage of non-compliance with the instructions of the Moroccan authorities.

With the application of our COVID-19 model, our estimate is as follows for the second scenario:

Verification of our model on the case of China:

China has recorded enough data that can be used for verification. The data on which we do this verification is extracted from the site: <https://www.worldometers.info/coronavirus/country/china/> [32].

It goes without saying that the accuracy of forecasts depends on the "method" approach first, on the number of data and the distance from the future to be predicted.

6. Conclusion

This work presents a dedicated machine learning model to predict the number of cases infected with the Corona Virus; the case of Morocco was chosen to validate this study. Completely realized in Spark ML with the 'Scala' language and tested for a certain number of algorithms generated on datasets from dedicated sources to group COVID-19 data to the world, the results show the possibility of achieving better scores prediction after using the proposed methods. We are sending this work to the world to help them in our way to fight this 2019 Coronavirus pandemic.

Abbreviations

PCA

Principal Component Analysis.

SVD

Singular Value Decomposition.

NBC

Naive Bayesian Classifier.

LFE

Learning Feature Engineering.

MSE

Mean Squared Error.

RMSE

this is the square root of MSE (root Mean Squared Error).

R

Coefficient of determination.

MAE

mean Absolute Error.

MAPE

Mean Absolute Percentage Error.

DTC

Decision Tree.

NBC

Gaussian Naive Bayes.

SVM

Support Vector Machine.

LR

Logistic Regression.

RFC

Random Forest.

VC

Voting Classifier.

PDE

Partial Differential Equation (PDE).

ODE

Ordinary Differential Equations.

Declarations

Availability of data and materials

All data generated or analysed during this study are included in this published article.

Competing interests

The authors declare that they have no competing interests.

Funding

Not applicable.

Authors' contributions

All authors contributed to the study conception and design. Material preparation, data collection and analysis were performed by Allae Erraissi, Mouad Banane, Abdessamad Belangour and Azouazi

Mohamed. The first draft of the manuscript was written by Allae Erraissi and all authors commented on previous versions of the manuscript. All authors read and approved the final manuscript.

Acknowledgements

Not applicable.

References

1. Kulkarni S, Seneviratne N, Baig MS, Khan AH. Artificial Intelligence in Medicine: Where Are We Now?. *Academic radiology*. 2019 Oct 19.
2. 10.1016/j.artmed.2020.101822
Waring J, Lindvall C, Umeton R. (2020, April 1). Automated machine learning: Review of the state-of-the-art and opportunities for healthcare. *Artificial Intelligence in Medicine*. Elsevier B.V. <https://doi.org/10.1016/j.artmed.2020.101822>.
3. Karau H, Warren R. High performance Spark: best practices for scaling and optimizing Apache Spark. " O'Reilly Media, Inc.", 2017.
4. « Meta-Modeling of Big Data Management Layer »
<https://doi.org/10.30534/ijeter/2019/01772019>
Erraissi A. and Abdessamad Belangour. « Meta-Modeling of Big Data Management Layer ». *International Journal of Emerging Trends in Engineering Research* 7, no 7, 36–43, 2019. <https://doi.org/10.30534/ijeter/2019/01772019>.
5. « Hadoop Storage Big Data Layer: Meta-Modeling of Key Concepts and Features »
<https://doi.org/10.30534/ijatcse/2019/49832019>
Erraissi A. and Abdessamad Belangour. « Hadoop Storage Big Data Layer: Meta-Modeling of Key Concepts and Features ». *International Journal of Advanced Trends in Computer Science Engineering* 8, n 3, 64653, 2019. <https://doi.org/10.30534/ijatcse/2019/49832019>.
6. Gorodetsky V, Samoylov V. 'Feature Extraction for Machine Learning: Logic-Probabilistic Approach', *JMLR: Workshop and Conference Proceedings* 10: 55–65 The Fourth Workshop on Feature Selection in Data Mining.
7. Zheng A, Casari A. 'Feature Engineering for Machine Learning', O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.
8. Wold S, Esbensen K, Geladi P. Principal component analysis. *Chemometrics intelligent laboratory systems*. 1987;2(1–3):37–52.
9. Golub GH, Reinsch C. Singular value decomposition and least squares solutions. In: *Linear Algebra*. Berlin: Springer; 1971. pp. 134–51.
10. Mani K, Kalpana P. 'An Efficient Feature Selection based on Bayes Theorem, Self-Information and Sequential Forward Selection'. *IJ Information Engineering Electronic Business*. 2016;6:46–54.

11. Blahut RE, Blahut RE. Principles and practice of information theory. Vol. 1. Reading: Addison-Wesley; 1987.
12. Kononenko I. March). Semi-naive Bayesian classifier. In: European Working Session on Learning. Berlin: Springer; 1991. pp. 206–19.
13. Fatemeh Nargesian H, Samulowitz U, Khurana EB, Khalil. Deepak Turaga, 'Learning Feature Engineering for Classification', Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI-17).
14. Nargesian F, Samulowitz H, Khurana U, Khalil EB, Turaga DS. (2017, August). Learning Feature Engineering for Classification. In IJCAI (pp. 2529–2535).
15. Spark A. "Apache spark." Retrieved January 17 (2018): 2018.
16. Meng X, et al. "Mllib: Machine learning in apache spark.". The Journal of Machine Learning Research. 2016;17(1):1235–41.
17. Townsend JT. Theoretical analysis of an alphabetic confusion matrix. Perception Psychophysics. 1971;9(1):40–50.
18. Bifet A, et al. "Streamdm: Advanced data mining in spark streaming." 2015 IEEE International Conference on Data Mining Workshop (ICDMW). IEEE, 2015.
19. Lai CC, Shih TP, Ko WC, Tang HJ, Hsueh PR. (2020). Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) and corona virus disease-2019 (COVID-19): the epidemic and the challenges. International journal of antimicrobial agents, 105924.
20. Morocco Covid-19 cases. <http://www.covidmaroc.ma/pages/Accueil.aspx>.
21. Mahmud SH, Hossin MA, Jahan H, Noori SRH, Bhuiyan T. (2018, May). CSV-ANNOTATE: Generate annotated tables from CSV file. In 2018 International Conference on Artificial Intelligence and Big Data (ICAIBD) (pp. 71–75). IEEE.
22. Odersky M, Lex Spoon, and Bill Venners. Programming in scala. Artima Inc, 2008.
23. Safavian S, Rasoul, Landgrebe D. A survey of decision tree classifier methodology. IEEE transactions on systems man cybernetics. 1991;21(3):660–74.
24. Ontivero-Ortega M, et al. "Fast Gaussian Naïve Bayes for searchlight classification analysis" Neuroimage. 2017;163:471–9.
25. Ben-Hur A, Weston J. "A user's guide to support vector machines." Data mining techniques for the life sciences. Humana Press, 2010. 223–239.
26. Kleinbaum DG, et al. Logistic regression. New York: Springer-Verlag; 2002.
27. Liaw A. and Matthew Wiener. "Classification and regression by randomForest." R news 2.3 (2002): 18–22.
28. Ruta D, Gabrys B. "Classifier selection for majority voting." Information fusion 6.1 (2005): 63–81.
29. Strikwerda JC. (2004). Finite difference schemes and partial differential equations (Vol. 88). Siam.
30. Hale JK. Functional differential equations. In: Analytic theory of differential equations. Berlin: Springer; 1971. pp. 9–22.

31. World Health Organization. "Coronavirus disease 2019 (COVID-19): situation report, 72." (2020).

32. Coronavirus Cases in China. <https://www.worldometers.info/coronavirus/country/china/>.

Figures

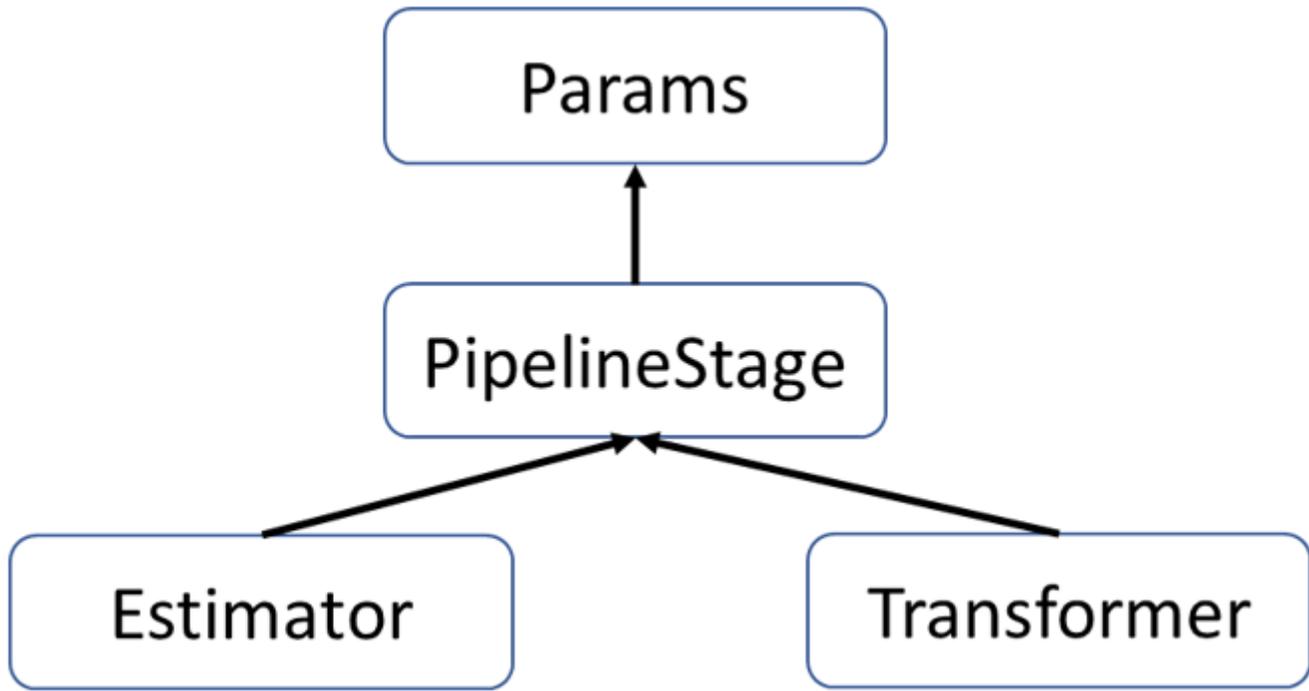


Figure 1

Hierarchy: Estimator and Transformer [16].

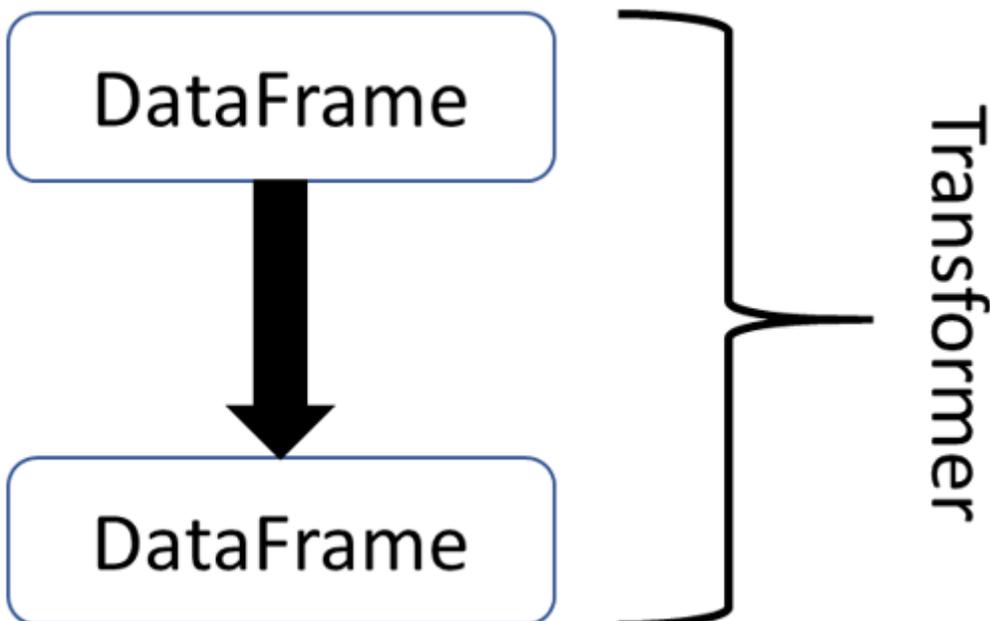


Figure 2

The role of Transformer [16].

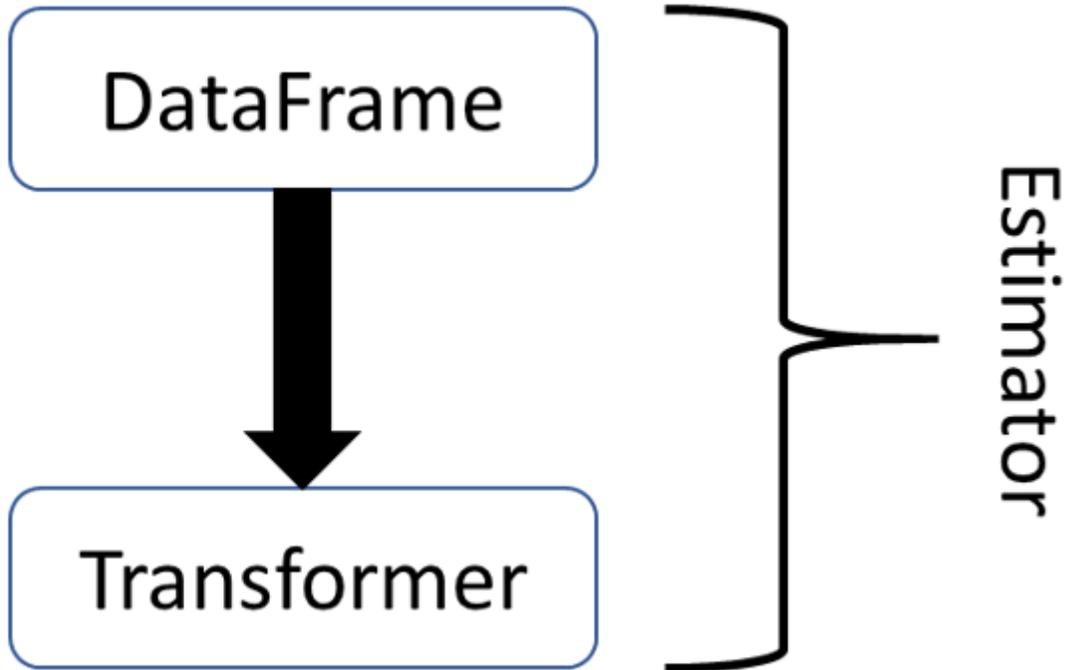


Figure 3

The role of Estimator [16].

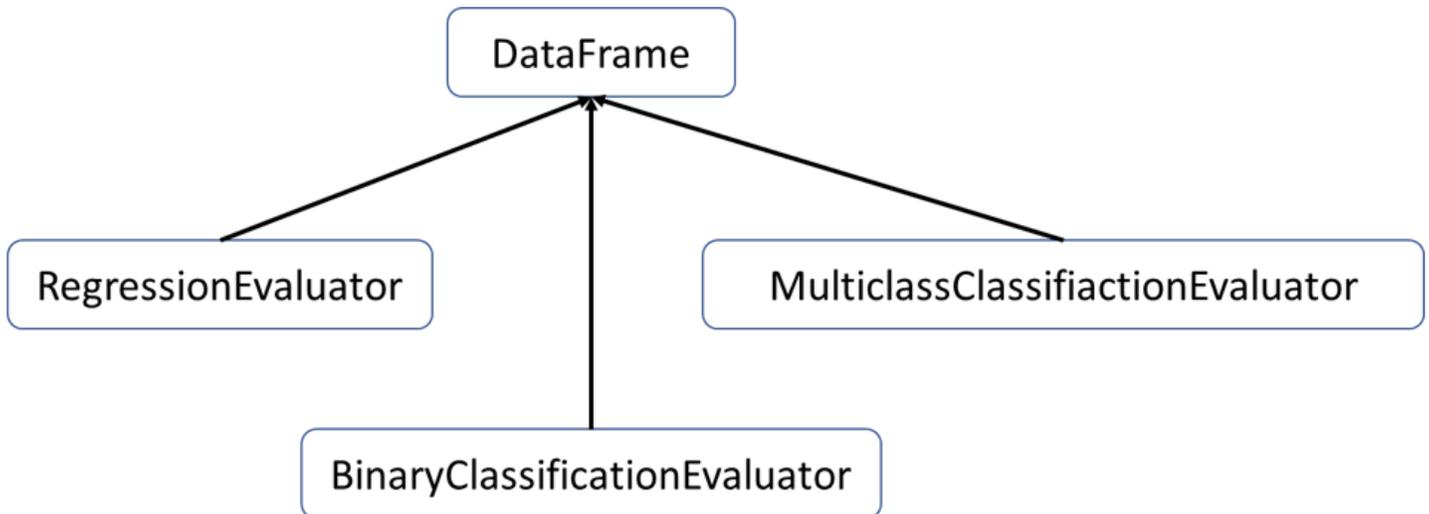


Figure 4

Implementations for the Evaluator abstract class [18].

Content of the DataFrame.

```
>res 0: the prediction of the day (07/04/2020) is : 65 cases
>res 1: the prediction of the day (08/04/2020) is : 92 cases
>res 2: the prediction of the day (09/04/2020) is : 98 cases
>res 3: the prediction of the day (10/04/2020) is : 97 cases
```

Figure 8

The result of the execution of our model.

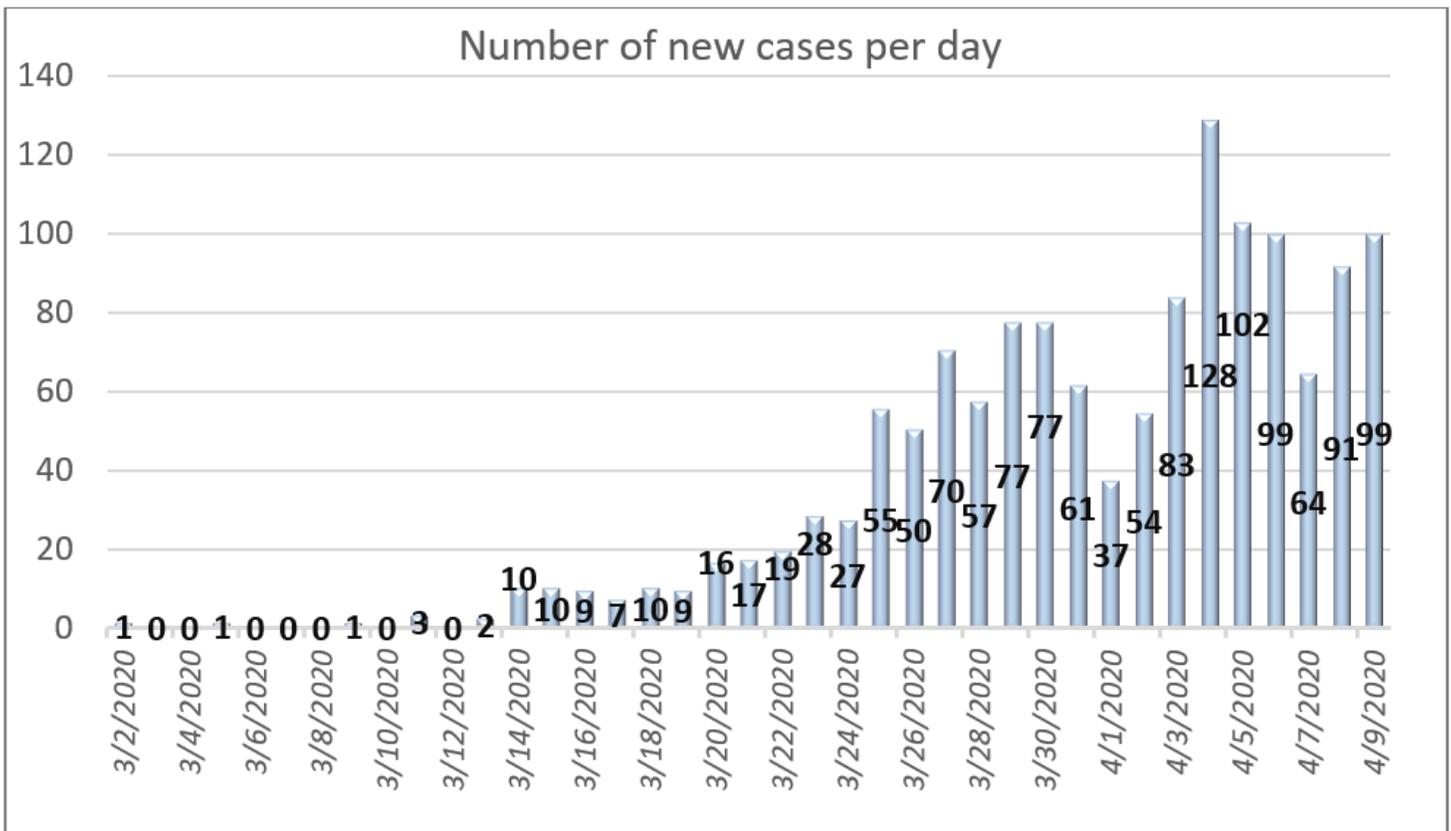


Figure 9

Number of COVID-19 cases per day in Morocco.

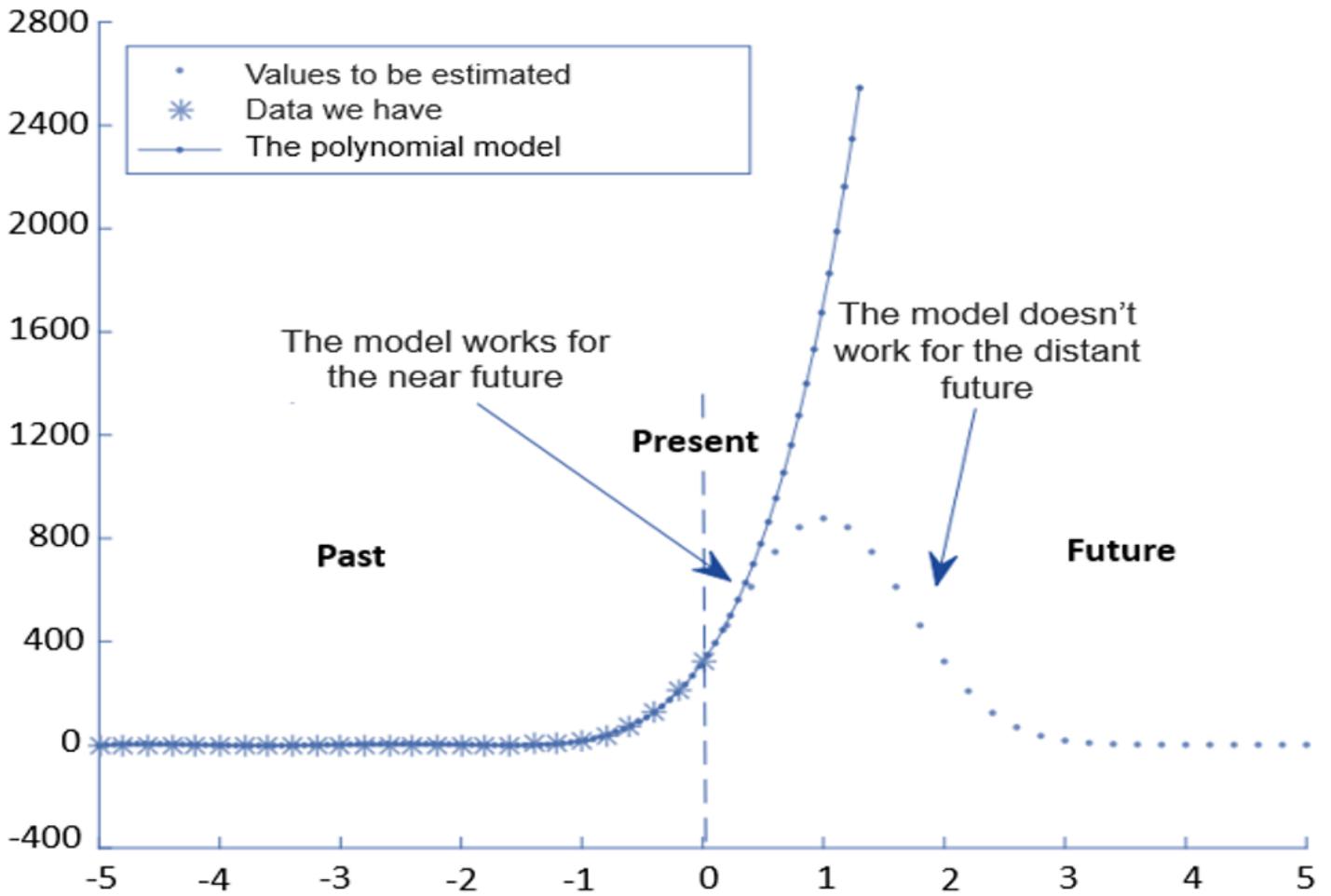


Figure 10

Result of the polynomial approach.

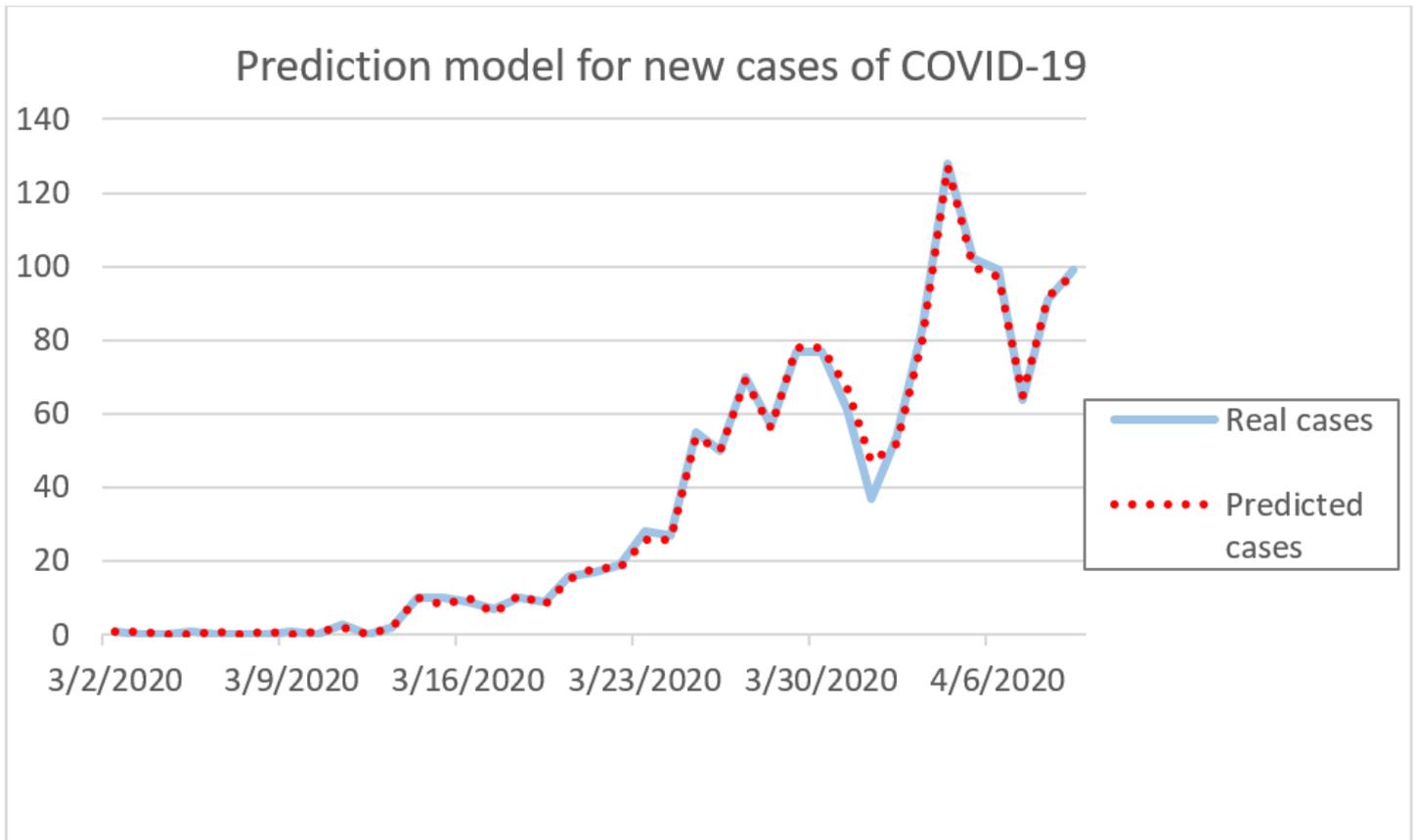


Figure 11

Prediction model for new cases of COVID-19 in Morocco.

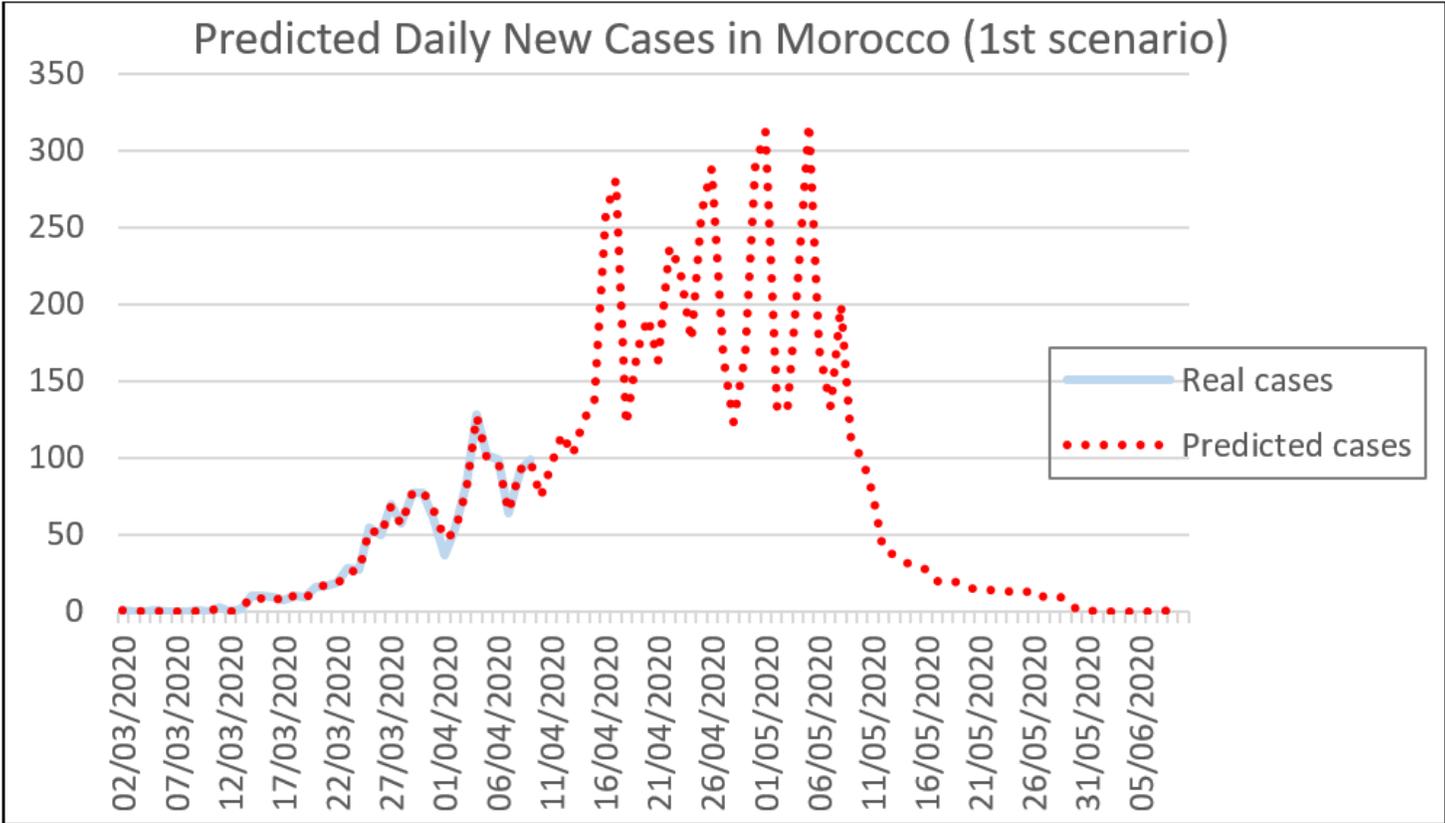


Figure 12

Predicted daily new cases in Morocco (1st scenario).

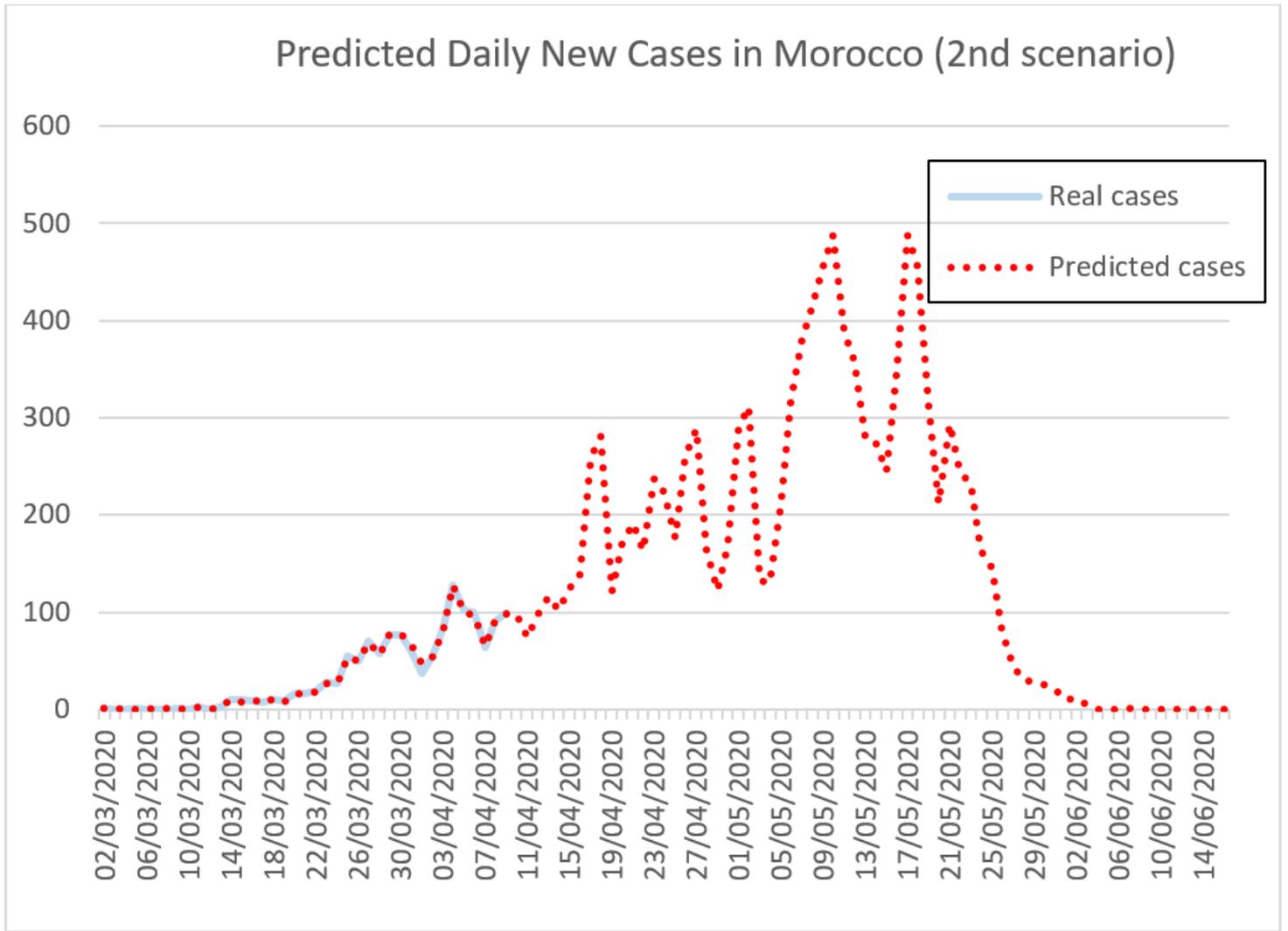


Figure 13

Predicted daily new cases in Morocco (2nd scenario).

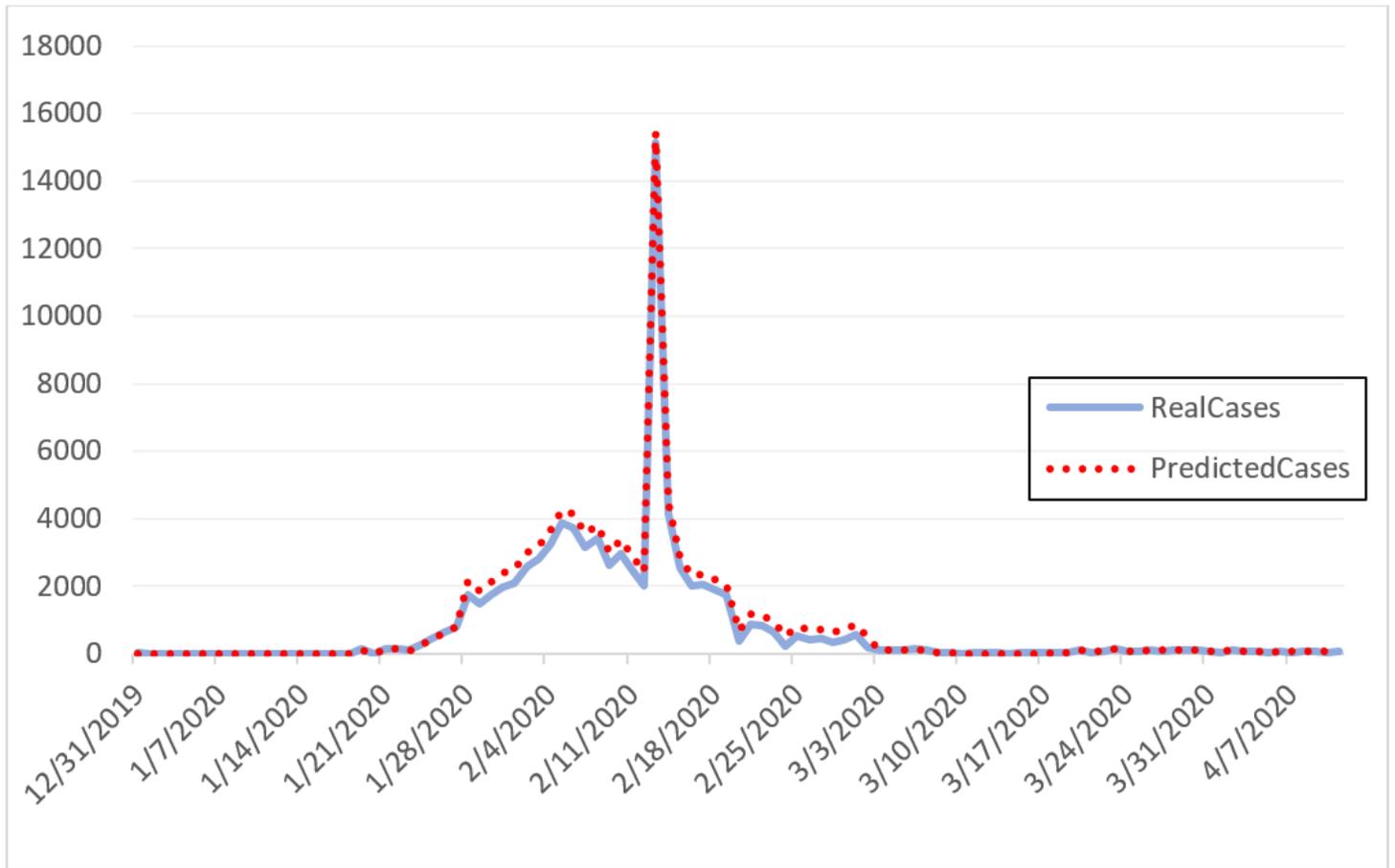


Figure 14

Application of our model on the case of China.