

A COMPARATIVE STUDY OF MACHINE LEARNING ALGORITHMS FOR PREDICTING CRIME HOTSPOTS

¹Dasyam Venus, ²K.Praveen Kumar, ³R.TamilKodi

¹PG student, ²Assistant Professor, ³Professor, ¹²³Department of Computer Application
¹²³Godavari Institute of Engineering and Technology(Autonomous), Rajahmundry, AP,
¹dasyamvenus@gmail.com, ²praveenkumar@giet.ac.in, ³tamil@giet.ac.in

ABSTRACT

Crime is one of the greatest and overwhelming issue in our general public and its counteraction is a significant undertaking. Every day there are enormous quantities of crimes carried out as often as possible. This require monitoring every one of the crimes and keeping an information base for same which might be utilized for future reference. The current issue confronted are keeping up of appropriate dataset of crime and breaking down this information to help in predicting and tackling crimes in future. The goal of this undertaking is to examine dataset which comprise of various crimes and predicting the sort of crime which may occur in future relying on different conditions. In this venture, we will utilize the method of AI and information science for crime expectation of crime informational collection. The crime information is separated from the authority entryway of areal police. It comprises of crime data like area description, kind of crime, date, time, scope, longitude. Prior to preparing of the model information pre handling will be finished after this component determination and scaling will be done with the goal that exactness acquire will be high. The K-Nearest Neighbor (KNN) arrangement and different calculations will be tried for crime expectation and one with better precision will be utilized for preparing. Perception of dataset will be done regarding graphical portrayal of numerous cases for instance at which time the criminal rates are high or at which month the criminal exercises are high. The spirit reason for this undertaking is to give a joke thought of how AI can be utilized by the law authorization offices to identify, foresee and settle crimes at a lot quicker rate and subsequently diminishes the crime rate. It not confined to particular area, this can be utilized in different states or nations relying on the accessibility of the dataset.

Keywords: *Nearest Neighbours Provision, Vector Machinery stirring average, persistent neuronal system, National Crime Records Bureau.*

1. INTRODUCTION

Crimes occur from little town, town to huge metropolitan regions moreover in tremendous metropolitan networks. India as we likely to be mindful is one of the metropolitan territories stacked with hostility, corruptions by administrative issues. Encroachment are of various sorts[1] – Burglary, murder, assault, trap, battery, false repression, appropriating, manslaughter, theft, youth abuse[2], Child Trafficking, Molestations, etc[3]. The Crime rehearses have been reached[4] out at a speedier rate and it is the responsibility of police division[5] to control and reduce the infringement works out[6]. Since Crimes are stretching out there is a need to illuminate[7] the cases and methodology in a lot snappier way. Crimes guess furthermore[8], criminal indisputable confirmations are the critical issues to the police office as there are gigantic extents of Crimes occasion information that exist[9]. So in this manner it is required to examination the[10] Crimes of each state and regions of India by months, a long time and season of present to thwart the advancing crimes[11]. AI is the investigation of having PCs make decisions without human intervention[12]. Lately, Machine Learning has been applied in self-driving vehicles, talk affirmation, [13]web search, and an improved perception of the

human genome[14]. It has in like manner made expecting Crimes [15]subject to referred to data achievable[16]. Game plan is a directed gauge framework which considers apparent class marks[17]. Gathering has been used in various regions including environment expecting, therapeutic thought, assets and banking, country security, and business knowledge. AI based Crimes assessment generally incorporates[18] data variety, request, plan recognizing confirmation, conjecture, and portrayal. Standard data mining techniques - alliance assessment, gathering and assumption,[19] pack examination, and exemption examination - perceive plans in coordinated data while fresher systems recognize plans from both coordinated and unstructured data. The fundamental objective of this work is to create an estimate model that can definitely expect Crimes. This paper talks about the field of law essential has framed into a totally unpredictable calling with limitless zones of strength and limit. Crimes assessment could be seen as the most present advancement to the field. This endeavour[20] executes a way to deal with check the Crimes practices and find the proportion of consistent things and close by that it predicts whether there will be any harsh Crimes occasion will be there or not all that that it might be hindered. In this paper, we first attempt to describe the past works done on the said subject with their benefits and negative marks. At thatpoint the paper pushes ahead to describe the proposed strategy and about the diverse calculation utilized, programming prerequisites and informational collection to accomplish the said theme. At that point we describe the procedure that we took up in settling the said task, we additionally furnish the watcher with legitimate square graph of the undertaking for better arrangement. The paper further gives the similar examination and its discoveries and afterward it at long last finishes up the point of the task with giving the future perspectives.

2.LITERATURE REVIEW

Previously, there have been various such systems, where bad behaviour data is penniless down using different estimations, basically K-Means, K-Medoids, KNN, etc A segment of the models and technique are clarified: The makers, Jain et al. in their paper "Bad behaviour Prediction using K-Means Algorithm" [2], have used K-implies gathering computation [21] find plans from the bad behaviour dataset. K-Means clustering computation [22] is distance-based estimation. The Euclidean distance metric is used to find the distance of a point from the nearest centres and picks if that point should have a spot with the pack or not[23]. The amount of packs can't be settled at the start of the estimation. Therefore, various patterns of K-Means should be performed. The makers have used Rapid Miner gadget for examination [24]in light of its flexibility and adaptability. The essential place of the assessment was to fathom which year was the wrongdoing rate generally significant and least. Supporting this data,[25] visual graphs are plotted for each gathering. The paper was circulated in 2013, after which, a lot of better results [26] using better headways were introduced. Bad behaviour dataset is an England based record of infringement from the year 1990 to 2012. Bogahawatte and Adikari [27] proposed an approach in which they highlighted the usage of data mining strategies, gathering what's more, gathering for effective assessment of bad behaviours and criminal distinctive verification by developing a structure named Intelligent Crime Examination System (ICSIS) that could perceive a criminal put together [28] up with respect to the confirmation accumulated from the bad behaviour region. They used bundling to perceive the bad behaviour plans which are used to complete bad behaviours realizing the way that each bad behaviour has certain models. The informational collection is set up with an oversaw learning estimation, Naïve Bayes [29] to expect possible suspects from the criminal records. His methodology consolidates developing a multi-expert for bad behaviour configuration recognizing evidence. There are experts for the spot, time, work brand name what's more, substance of gangsters what secludes the piece of the criminals in fragments. [30] The system is a multi-expert structure and made with administered Java Beans. It simplifies it to

embody the referenced components in the work into things and returns it to the bean for revealing properties. Gathering the culprits/associates is in see with the Naïve Bayes classifier for recognizing commonly possible suspects from bad behaviour data. Packing the culprits relies upon the model to help with recognizing instances of completing bad behaviours. Agarwal et al. [3] used the quick digger instrument for analysing the wrongdoing rates and assumption for wrongdoing rate using unmistakable data mining methodologies. Their work done is for bad behaviour assessment using the K-Means Clustering estimation. The guideline objective of their bad behaviour assessment work is to eliminate the bad behaviour plans, anticipate the bad behaviour reliant upon the spatial scattering of existing data and disclosure of bad behaviour. Their assessment consolidates the following homicide wrongdoing rates beginning with one year then onto the following Kiani et al. [4] played out a bad behaviour examination work subject to the gathering and game plan systems. Their work fuses the extraction of bad behaviour plans by bad behaviour assessment subject to available criminal information, assumption for bad behaviours reliant upon the spatial movement of existing data and bad behaviour affirmation. They proposed a model where the examination and assumption for bad behaviours are done through the progression of special case area overseer limits which is performed through the Genetic Algorithm. The features are weighted in this model and the low-regard features were deleted through picking a proper cut-off. After which the gatherings are bundled by the k-implies packing computation for gathering of bad behaviour dataset. Satyadevan et al. [5] has achieved a work which will show high probability for bad behaviour occasion and can picture bad behaviour slanted districts. Maybe than just focusing in on the bad behaviour occasions, they are focusing in essentially on the bad behaviour factors of consistently. They used the Guileless Bayes, Logistic Regression and SVM classifiers for plan of bad behaviour models and bad behaviour parts of consistently. Their technique includes a model distinctive confirmation stage which can recognize the examples and models in bad behaviour using the Apriori Calculation. The assumption for bad behaviour spots is done with the help of Choice Tree estimation which will perceive the bad behaviour possible zones and their models

2.1 Disadvantages

The rank order of the location with

- 1 being the location with the most incidents,
- 2 being the location with the next most incidents,
- 3 being the location with the third most incidents, and so forth until those locations that have only one incident each;
4. The frequency of incidents at the location. This is the number of incidents occurring at that location;
5. The X coordinate of the location;
6. The Y coordinate of the location.

3. PROPOSED METHOD

In this paper, different algorithm, for example, Naïve Bayes, KNN, Gradient boosting, Logistic regression, Random forest, Genetic Algorithm are utilized on same dataset. From the crime dataset, it tends to be seen that a characterization of crime relies upon different factors, for example, Date, Category, Description of the crime, Day of the Week on which crime happened, Police station close to

nearest to the crime scene, Resolution of the crime and co-ordinates (for example longitude and scope) of the crime. Various components impact the order of the crime in an unexpected way. A few factors like crime co-ordinates and season of the crime has high co-connection with the crime type, while different factors, for example, "Day of the week" has low co-connection with the crime. Henceforth, utilizing those components which has high connection with crime type will improve the precision of any algorithm. Results fluctuates for various classifier, consequently genetic algorithm is utilized to sort out the highlights for a specific classifier with which it gives great exactness to characterization. With the assistance of dataset having a place with various urban areas, pictorial portrayal of crime is acquired for the better comprehension of the crime type and what piece of the city struggles with what sort of crime. Crime shrewd "Heat maps" are produced to stick point the piece of the city on the guide and the degree of specific crime in that part. Genetic algorithm regards different highlights as "chromosomes" and treats them thusly. It is recursive in nature and incorporates the cycle of "change" and "get over". At first, it will make different stochastic arrangements of "highlights", at that point discover precision for a specific classifier with different sets. It keeps a vector of "sets of highlights" which gives best exactness's. In the wake of discovering the exactness, it does transformation and hybrid of various arrangements of highlights. At long last, it furnishes with the arrangement of highlights which will give the best exactness for specific classifier.

4.METHODOLOGY

The Hotspot detection is defined by the following equation:

$$F^k_l = (I_{x,y} * K^k) \text{ Eqn (4.1)}$$

It summarizes comparable data in the place of the open field and yields the main reaction inner this community district by following.

$$Z_i = f_p(F^l_{x,y}) \text{ Eqn(4.2)}$$

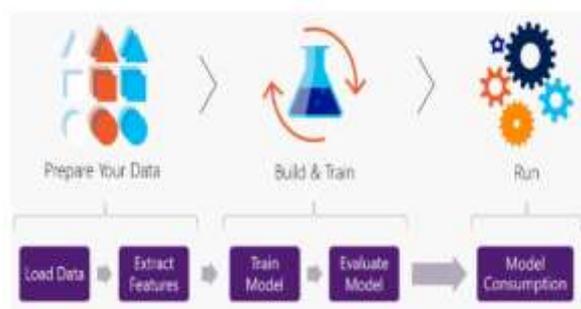


Fig : Flow Diagram of implementation process by ensemble of algorithms

THE METHODOLOGY IMPLEMENTATION IS CONDUCTED IN FOLLOWING STEPS

Step :a) Dataset type and resources: This part defines that it can permit the utilization of table labels and names of state where the crime Occurring, client reference and URLs. These require unexpected preparing in comparison to different data in all set. The amount of vast dataset or base can be gathered of local datasets from any political or local government organization for records or any cloud based database such as tinyDB from thinkable platform or live fetching environment using python libraries (e.g. -tweepy) for the provided terms and data.

Step: b)Data tokenization: Information Tokenization is the cycle used to supplant touchy information with remarkable recognizable proof images that hold all the fundamental data about the information without trading off its security. It will give in numeric state whether the information is in one of the wrongdoing class

Step: c)Machine Learning Algorithm: Data gathered from data warehouse will go through processing of ensemble algorithms and coded in python language and provide outputs of individuals. This will show how efficient can one platform is. In end the terminal of console shows the accuracy of supervised and unsupervised categorized program to the developer and compare of all algorithms.

Step: d) Data Diagrams: Client will see the process through chart and graphs about the state of the crimes amount in individual areas. Take the visualization in particular axis and precautions can be taken. In this way predicted areas can be revealed and in future the crimes can be avoided.

Step e) Heat map: Implementation time and processing will return a graph and chart of affected areas along with that it will provide the heat map prediction of effected with the accuracy of individual algorithms and areas according to months and years of data gathered to prevent crime happenings later on.

Step: f) Comparison: The supervised and unsupervised algorithms will be compared with each other for checking in efficiency. The more the accuracy come across according the dataset given the more the result will be up to the mark. The comparison gives an output of one outcome so clients don't have to go with individual checking and come up with one answer which is more time convenient and less efficient.

5.RESULT ANALYSIS

Dissecting crime information utilizing different AI algorithm assists with understanding what algorithm is proficient for the grouping for various crimes. On the off chance that arrangement of crime is done appropriately, it will make it conceivable to envision which region is caused by which kind of crime and what strategy change is needed in a specific region to battle a particular crime. It very well may be effortlessly seen that the random forest classifier is the awesome order of a crime dataset followed by Extra Tree classifier, Gradient Boosting algorithm. Genetic algorithm depends on the idea of change and consequently is utilized to sort out highlights that helps in improving the precision of a specific classifier. The issue of genetic algorithm is that it is iterative in nature, along these lines classifiers which have high execution time or are themselves iterative in nature such "Logistic Regression", "Slope Boosting Classifier", "Random Forest Classifier" ought not be utilized with Genetic Algorithm. Utilizing genetic algorithm with classifiers, for example, KNN or Naïve Bayes results in effectively increment of precision.

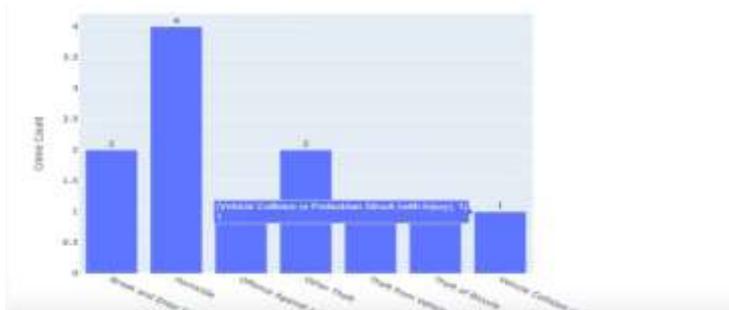


Figure 5.1 Crime committed analysed by the system

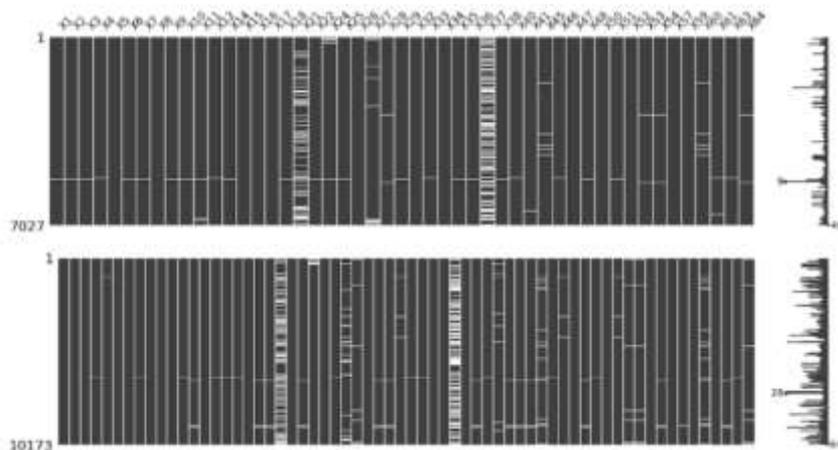


Figure 5.2 heat Graph of crimes

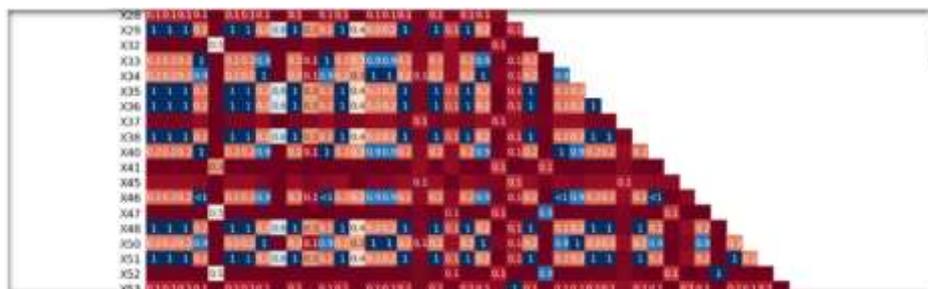


Figure5.3 Corelation graphs

The above figure is a result of the all the steps that are being followed in the procedure the data is separated. Each meta data with its related data was extracted. These data point will be taken to be the reference points for the algorithm. The data points are crime data, area of the crime committed and the nature of the crime that is committed. From these data points the nearest neighbours for the crime data were found and the crime datapoints are feed into the random forest algorithm for the classification of hotspots as shown in the figure. The results were compared with real world data and found to be 90% accurate. Every hotspot was tested to precision on every iteration.

6.COMPARATIVE STUDY

Comparing the accuracy in different machine learning algorithms are here shown in a tabular format.

CLASSIFIER	F1	F2	F3	F4	F5	F6	F7
NBB	24.20	23.8	24.3	21.03	12.63	21.56	21.11
LR	23.81	23.8	23.8	33.16	33.6	22.31	26.2
KNN	12.82	14.59	21.61	59.83	13.62	21.45	23.12
RF	24.87	27.27	24.86	58.98	36.12	29.35	31.33
ETC	25.15	23.89	24.16	56.35	38.64	26.37	27.56
NBG	23.78	21.77	17.77	40.85	32.65	24.21	24.51

By considering the various features, the comparison of machine learning algorithms can be done. F1- Hour, F2- Day of the week, F3-Police station, F4- Crime description, F5-Crime resolution, F6- Crime address, F7- Crime co-ordinates.

The above table ,the algorithms that are mentioned :-NBB: Bernoulli Naïve Bayes, LR: Logistic Regression, KNN: K Nearest neighbour, RF: Random Forest, ETC: Extra Tree Classifier, NBG: Gaussian Naïve Bayes. From the above table, It is clear that Crime description and Crime resolution among the mentioned features are having good accuracy, and Day of week and Hour are having less accuracy when compared. The classification can be done with good correlation features. This improves the accuracy nearly each algorithm to a great extent.

7.CONCLUSION

From this paper, it tends to be presumed that due to high number of classes where the crime has been grouped (37 classes), precision of different AI algorithm is low. There are numerous ML algorithm, for example, "Inclination Boosting" and "Random forest" which can characterize crime with high exactness. On the off chance that kinds of crime is gathered and new classes are shaped for arrangement, the exactness of other algorithm, for example, Naïve Bayes and logistic regression can be improved. The utilization of "Genetic algorithm" with different classifier can prompt improve in the exactness of the classifier to extraordinary degree. From this paper, it tends to be noticed for KNN. With Genetic algorithm, precision of the classifier helped to 97.21 from 49.91. The crime investigation part is additionally holding the forecast interaction in this manner it gives aftereffects of individual algorithms and yield will be appeared. In a concise this specific task is adaptable enough where we can include ne new highlights alongside execution on Cloud Platform, for example, AWS (Amazon Web Service) Microsoft Azure and so forth, in future if necessary. The specific application can likewise be included not many IoT related ventures, for example, Image preparing by utilizing a robot or vehicle and investigate the climate exhaustive reconnaissance camera.

REFERENCES

1. J. R. Kling and J. Ludwig, "Is Crime Contagious?" *The Journal of Law and Economics*, vol. 50, no. 3, pp. 491–518, 2007.
2. C. F. Manski, "ID of endogenous social impacts: the reflection issue," *Review of Economic Studies*, vol. 60, no. 3, pp. 531–542, 1993.
3. K. A. Goss and P. Cook A, *A Selective Review of the Social-Contagion Literature*, Terry Sanford Institute Working Paper, Duke University, Durham, NC, USA, 1996.
4. C. F. Manski, "Monetary examination of social communications," *Journal of Economic Perspectives (JEP)*, vol. 14, no. 3, pp. 115–136, 2000.
5. W. J. Wilson, *The Truly Disadvantaged: The Inner City, the Underclass, and Public Policy*, The University of Chicago Press, Chicago, Illinois, 1987.
6. R. J. Sampson, S. W. Raudenbush, and F. Lords, "Neighborhoods and rough wrongdoing: a staggered investigation of aggregate viability," *Science*, vol. 277, no. 5328, pp. 918–924, 1997.
7. C. Jencks and S. E. Mayer, "The Social Consequences of Growing Up," in *Inner-City Poverty in the United States*, L. Lynn and M. McGeary, Eds., National Academy of Sciences, Washington, DC, USA, 1990.
8. E. L. Glaeser, B. Sacerdote, and J. A. Scheinkman, "Wrongdoing and social cooperations," *The Quarterly Journal of Economics*, vol. 111, no. 2, pp. 507–548, 1996.

9. E. L. Glaeser, B. I. Sacerdote, and J. A. Scheinkman, "The social multiplier," *Journal of the European Economic Association*, vol. 1, no. 2-3, pp. 345–353, 2003.
10. J. Crane, "The pandemic hypothesis of ghettos and neighborhood impacts on exiting and high school childbearing," *American Journal of Sociology*, vol. 96, no. 5, pp. 1226–1259, 1991.
11. D. A. Toll and P. R. Nail, "Virus: a hypothetical and observational audit and reconceptualization," *Genetic Social and General Psychology Monographs*, vol. 119, no. 2, pp. 233–284, 1993.
12. J. Fagan, D. L. Wilkinson, and G. Davies, "Social virus of savagery," in *The Cambridge Handbook of Violent Behavior and Aggression*, D. Flannery, A. T. Vazsonyi, and I. D. Waldman, Eds., Cambridge University Press, Cambridge, UK, 2007.
13. R. S. Burt, "Social Contagion and Innovation: Cohesion versus Structural Equivalence," *American Journal of Sociology*, vol. 92, no. 6, pp. 1287–1335, 1987.
14. L. O. Gostin, "The interconnected pandemics of medication reliance and AIDS," *Harvard Civil Rights Civil Liberties Law Review*, vol. 26, no. 1, pp. 113–184, 1991.
15. J. L. Rodgers and D. C. Rowe, "Social infection and juvenile sexual conduct: A formative EMOSA model," *Psychological Review*, vol. 100, no. 3, pp. 479–510, 1993.
16. S. B. Patten, "Pandemics of viciousness," *Medical Hypotheses*, vol. 53, no. 3, pp. 217–220, 1999.
17. E. Massad, A. F. Rocha, F. A. B. Coutinho, and L. F. Lopez, "Demonstrating the spread of images: How innovations are communicated from one mind to another," *Applied Mathematical Sciences*, vol. 7, no. 45-48, pp. 2295–2306, 2013.
18. K. Kirkpatrick, "The Social Contagion of Violence: A Theoretical Exploration of the Nature of Violence in Society," in *Social Science*, vol. 461, 462, California Polytechnic State University, San Luis Obispo, CA, USA, 2018, <http://digitalcommons.calpoly.edu/socssp/78>.
19. S. B. Patten and J. A. Arboleda-Flórez, "Plague hypothesis and gathering viciousness," *Social Psychiatry and Psychiatric Epidemiology*, vol. 39, no. 11, pp. 853–856, 2004.
20. C. C. N. Dias, *PCC: Hegemonia nas prisões e monopólio da violência (Portuguese)*, Editora Saraiva, São Paulo, Brazil, 2013.
21. S. Adorno and F. Salla, "Criminalidade organizada nas prisões e os ataques do PCC," *Estudos Avançados*, vol. 21, no. 61, pp. 7–29, 2007.
22. R. M. Anderson and R. M. May, *Population Biology of Infectious Diseases*, Springer-Verlag, Berlin, Heidelberg, New York, 1982.
23. R. M. Anderson and R. M. May, *Infectious Diseases of Humans: Dynamics and Control*, Oxford University Press, Oxford, England, second version, 1991.
24. S. Machin, M. Olivier, and V. Suncica, "The wrongdoing decreasing impact of schooling," *Economic Journal*, vol. 121, no. 552, pp. 463–484, 2011.
25. M. Maguire, "Cheerful Morash: Understanding Gender, Crime and Justice," *Critical Criminology*, vol. 16, no. 3, pp. 225–227, 2008.
26. J. K. Sound, *Ordinary Differential Equations*, Krieger, Basel, Switzerland, 1980.
27. H. W. Hethcote, "The math of irresistible illnesses," *SIAM Review*, vol. 42, no. 4, pp. 599–653, 2000.
28. V. Lakshmikantham, S. Leela, and A. A. Martynyuk, *Stability Analysis of Nonlinear Systems*, Marcel Dekker, New York, NY, USA, 1989.
29. E. Massad and F. A. B. Coutinho, "Vectorial limit, essential generation number, power of disease what not: Formal documentation to finish and change their old style ideas and conditions," *Memórias do Instituto Oswaldo Cruz*, vol. 107, no. 4, pp. 564–567, 2012.

30. O. Sharomi, C. N. Podder, A. B. Gumel, E. H. Elbasha, and J. Watmough, "Job of occurrence work in antibody prompted in reverse bifurcation in some HIV models," *Mathematical Biosciences*, vol. 210, no. 2, pp. 436–463, 2007.