



Ministry of Higher Education and
Scientific Research
University of Diyala
Department of Computer Science



Development of Prediction Model for Population Concentration using LSTM and GRU Algorithms

A thesis

**Submitted to the Department of Computer Science \ College
of the Science \ University of Diyala in a Partial Fulfillment of the
Requirements for the Degree of Master in Computer Science**

By

Naseer Salah Abbas

Supervised By

Prof. Naji M. Sahib

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

وَالَّذِي هُوَ يُطْعِمُنِي وَيَسْقِينِ ﴿٧٩﴾

وَإِذَا مَرِضْتُ فَهُوَ يَشْفِينِ ﴿٨٠﴾

صدق الله العظيم

❧ Acknowledgment ❧

First of all, praise is to GOD, the lord of the whole creation, on all the blessing was the help in achieving this research to its end.

I wish to express my thanks to my college (college of science), my supervisor, prof. Naji Motar Sahib for supervising this research and for the generosity, patience and continuous guidance throughout the work. It has been my good fortune to have advice and guidance from him. My thanks to the academic and administrative staff at the Department of the computer sciences.



Naseer Salah Abbas

❧ Dedication ❧

I would like to dedicate this work to:

*To my candle that light my life My
Mother.*

*To My father may God have mercy on
him.*

For his unlimited love, his support,

And Encouragement for me,

About everything in my life.

He is my first and last success.

To My Brothers and Sister.

To All My Friends.

I produce this work with all my love....



Naseer Salah Abbas

Abstract

The health crisis that attributed to the quick spread of the COVID-19 has impacted the globe negatively in terms of economy, education and transport and led to the global lockdown. The risk of the COVID-19 infection has been increased due to a lack of a successful cure for the disease. Thus, social distancing is considered the most appropriate precaution measure to control the viral spread throughout the world.

In this study, a model was proposed for predicting the people's movement to find out the proportion of social distancing in the short term (one day) using deep learning algorithms to take the necessary measures and precautions to control the COVID-19 infection. The proposed model consists of four phases: data collection, data pre-processing phase, prediction and evaluation stage and comparison phase. The dataset is obtained from 428 mobility reports, collected based on data from users selected for their Google Account location history for a country such as Iraq from 2020-02-15 to 2021-04-17 (428 days) stored in a comma-separated values file (CSV). Then, deep learning algorithms Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU) and hybrid model (GRU & LSTM) are applied to pre-processed data to predict the people's movement. They are compared using statistical measures: Mean absolute error (MAE) and root mean square error (RMSE) for performance measurement of these machine learning algorithms, where the lower the error rate, the better and more accurate the prediction performance of the model.

The results of the GRU are the sum of MAE = 0.4284 and sum of RMSE = 0.6479 for predict people's movement with training time equal to 35.951 sec, while the results of the hybrid model are the sum of MAE = 0.4360 and sum of RMSE = 0.6558 for prediction and the training time equal to 71.190 sec and the results of the LSTM are the sum of MAE = 0.4418 and sum of RMSE = 0.6618 for prediction and the training time equal to 123.799 sec.

These statistical measurement values indicate proposed model GRU outperformed all other models, it showed a solid performance to predict person path and movement in coronavirus pandemic and took little time to train compared to other algorithms, while the hybrid algorithm showed good performance and a short time in training compared with the LSTM model.

Social distancing and mixing between people play a vital role in the rapid spread of the COVID-19 pandemic. Therefore, our forecasted results of future trends of people's movement in society are very helpful for the country to control the pandemic and for purposes of social distancing guidance.

List of Contents

<i>Subject</i>	<i>Page No.</i>
<i>Abstract</i>	<i>I</i>
<i>List of Contents</i>	<i>III</i>
<i>List of Abbreviations</i>	<i>VII</i>
<i>List of Figures</i>	<i>VIII</i>
<i>List of Tables</i>	<i>X</i>
<i>List of Algorithms</i>	<i>X</i>
<i>Chapter One: Introduction</i>	
<i>1.1 Overview</i>	<i>1</i>
<i>1.2 Related works</i>	<i>3</i>
<i>1.3 Problem Statement</i>	<i>6</i>
<i>1.4 Aim of Thesis</i>	<i>7</i>
<i>1.5 Outline of Thesis</i>	<i>7</i>
<i>Chapter Two: Theoretical Background</i>	
<i>2.1 Introduction</i>	<i>8</i>
<i>2.2 Check of missing data</i>	<i>8</i>
<i>2.3 Data Normalization</i>	<i>8</i>

<i>2.4 Machine learning (ML)</i>	<i>10</i>
<i>2.5 Deep learning model</i>	<i>11</i>
<i>2.6 Recurrent Neural Networks (RNNs)</i>	<i>12</i>
<i>2.6.1 Activation Function</i>	<i>14</i>
<i>2.6.2 Loss Function</i>	<i>15</i>
<i>2.6.3 Backpropagation</i>	<i>16</i>
<i>2.6.4 RNN architectures</i>	<i>17</i>
<i>2.6.5 Exploding and vanishing gradients</i>	<i>19</i>
<i>2.6.6 Long Short-Term Memory (LSTM) RNNs</i>	<i>19</i>
<i>2.6.7 Gated Recurrent Units (GRU) RNNs</i>	<i>26</i>
<i>2.6.8 Optimization Algorithm</i>	<i>28</i>
<i>2.6.9 Overfitting</i>	<i>29</i>
<i>2.6.10 Accuracy Metrics</i>	<i>30</i>
<i>Chapter Three: The Proposed Model</i>	
<i>3.1 Introduction</i>	<i>31</i>
<i>3.2 Structure of the Proposed System</i>	<i>32</i>
<i>3.2.1 Dataset Information</i>	<i>32</i>
<i>3.2.2 Dataset Pre-Processing</i>	<i>33</i>

<i>i. Cleaning Dataset and Rename Features</i>	<i>33</i>
<i>ii. Handling Missing Values</i>	<i>33</i>
<i>iii. Data Normalization</i>	<i>35</i>
<i>iv. Reshape Input Data</i>	<i>36</i>
<i>3.2.3 Deep Learning Model</i>	<i>36</i>
<i>i. LSTM Algorithm</i>	<i>37</i>
<i>ii. GRU Algorithm</i>	<i>41</i>
<i>iii. Hybrid (LSTM and GRU) models</i>	<i>44</i>
<i>3.2.4 Evaluation and Comparison Phase</i>	<i>47</i>
<i>3.3 The Prediction Systems</i>	<i>48</i>
<i>Chapter Four: Experimental Results and Evaluation</i>	
<i>4.1 Introduction</i>	<i>50</i>
<i>4.2 Implementation Environment</i>	<i>50</i>
<i>4.3 Pre-processing Phase</i>	<i>50</i>
<i>4.3.1 Cleaning Dataset and Rename Features</i>	<i>51</i>
<i>4.3.2 Check of Missing Data</i>	<i>52</i>
<i>4.3.3 Exploratory Data Analysis (EDA)</i>	<i>53</i>
<i>4.3.4 Data Normalization</i>	<i>55</i>

<i>4.3.5 Reshape Input Data</i>	<i>55</i>
<i>4.4 The Prediction Systems Results</i>	<i>55</i>
<i>4.4.1 First System: LSTM Results</i>	<i>56</i>
<i>4.4.2 Second System: GRU Results</i>	<i>60</i>
<i>4.4.3 Third System: Hybrid model Results</i>	<i>64</i>
<i>4.5 Performance Comparison</i>	<i>68</i>
<i>Chapter Five: Conclusions and Future Works</i>	
<i>5.1 Conclusions</i>	<i>72</i>
<i>5.2 Suggestions for Future Works</i>	<i>73</i>
<i>References</i>	<i>74</i>
<i>Appendix A</i>	<i>79</i>

List of Abbreviations

<i>Abbreviation</i>	<i>Description</i>
<i>Adam</i>	<i>Adaptive Moment Optimizer</i>
<i>AI</i>	<i>Artificial Intelligence</i>
<i>ANN</i>	<i>Artificial Neural Network</i>
<i>BPTT</i>	<i>Backpropagation Through Time</i>
<i>COVID</i>	<i>Coronavirus</i>
<i>CSV</i>	<i>Comma Seprated Values</i>
<i>CT</i>	<i>Computed Tomography</i>
<i>DL</i>	<i>Deep Learning</i>
<i>GRU</i>	<i>Gated Recurrent Unit</i>
<i>IoT</i>	<i>Internet of Things</i>
<i>LSTM</i>	<i>long Short-Term Memory</i>
<i>MAE</i>	<i>Mean Absolute Error</i>
<i>ML</i>	<i>Machine Learning</i>
<i>MSE</i>	<i>Mean Square Error</i>
<i>RMSE</i>	<i>Root Mean Square Error</i>
<i>RNN</i>	<i>Recurring Neural Network</i>

List of Figures

<i>Figure No.</i>	<i>Figure Title</i>	<i>Page No.</i>
(2.1)	<i>Techniques of Machine learning</i>	10
(2.2)	<i>Architecture of Neural Network</i>	11
(2.3)	<i>General RNN model</i>	13
(2.4)	<i>Standard Sigmoid function</i>	15
(2.5)	<i>Standard Hyperbolic tangent function</i>	15
(2.6)	<i>The different types of RNNs</i>	17
(2.7)	<i>The repeating module in a standard RNN contains a single layer</i>	19
(2.8)	<i>The structure shows forget, input and output gates</i>	20
(2.9)	<i>The structure of an LSTM block</i>	21
(2.10)	<i>The cell state</i>	22
(2.11)	<i>The forget gate layer</i>	23
(2.12)	<i>Input gate layer and candidate memory layer</i>	24
(2.13)	<i>Output gate layer</i>	24
(2.14)	<i>The update of the old cell state</i>	25
(2.15)	<i>The structure of a GRU block</i>	27
(3.1)	<i>General Block Diagram of the proposed prediction system</i>	32
(3.2)	<i>The flowchart for check and process the missing data</i>	34
(3.3)	<i>Timesteps for Data</i>	36
(3.4)	<i>Proposed LSTM model</i>	40

(3.5)	<i>The proposed GRU model</i>	44
(3.6)	<i>Schematic diagram of the hybrid model</i>	47
(3.7)	<i>The Proposed Prediction Systems Flowchart</i>	48
(4.1)	<i>Sample of a dataset for Iraq country before preprocessing</i>	51
(4.2)	<i>Check the type of dataset</i>	51
(4.3)	<i>Dataset after cleaning</i>	52
(4.4)	<i>a- Dataset before processing the missing data. b- Dataset after processing the missing data.</i>	52
(4.5)	<i>Exploratory data analysis for features samples</i>	54
(4.6)	<i>Data normalization between 0 and 1</i>	55
(4.7)	<i>LSTM model layers</i>	58
(4.8)	<i>Accuracy of predicted results in LSTM model</i>	59
(4.9)	<i>The loss for the LSTM model</i>	60
(4.10)	<i>GRU model layers</i>	61
(4.11)	<i>Accuracy of predicted results in the GRU model</i>	63
(4.12)	<i>The loss for the GRU model</i>	64
(4.13)	<i>Hybrid Model layers</i>	65
(4.14)	<i>Accuracy of predicted results in the Hybrid model</i>	67
(4.15)	<i>The loss for the Hybrid model</i>	68
(4.16)	<i>MAE for prediction models</i>	69
(4.17)	<i>RMSE for prediction models</i>	70
(4.18)	<i>The sum of test errors for prediction models</i>	71

List of Tables

Tables No.	Tables Title	Page No.
(4.1)	<i>Model summary (LSTM)</i>	58
(4.2)	<i>Errors for prediction in LSTM model</i>	59
(4.3)	<i>Model summary (GRU)</i>	62
(4.4)	<i>Errors for prediction in GRU model</i>	62
(4.5)	<i>Model summary (Hybrid)</i>	65
(4.6)	<i>Errors for prediction in Hybrid model</i>	66
(4.7)	<i>MAE for prediction models</i>	68
(4.8)	<i>RMSE for prediction models</i>	69
(4.9)	<i>The sum of test errors for prediction models</i>	70

List of Algorithms

Algorithm No.	Algorithm Title	Page No.
(3.1)	<i>Pre-processing Algorithm</i>	34
(3.2)	<i>LSTM Cell Network Algorithm</i>	38
(3.3)	<i>GRU cell Network Algorithm</i>	41
(3.4)	<i>Hybrid Model Algorithm</i>	44
(3.5)	<i>The Proposed Prediction System Algorithm</i>	49

Chapter One

Introduction

Chapter one

Introduction

1.1 Overview

Coronavirus (COVID-19) is a member of the Coronaviridae family of viruses. In every country around the world, a virus with no cure is wreaking havoc on people's lives, as well as financial and economic losses. The epidemic was declared a public health emergency and pandemic on January 30, 2020, by the World Health Organization (WHO). This virus causes exhaustion, dry cough, fatigue and respiratory problems among other things. Maintaining social distance is the safest way to avoid and slow down transmission. As a result, an automated detection system must be implemented to prevent the virus from spreading among people. The role of computer technologies in discovering Corona disease is large and effective. In the battle against the COVID-19 crisis, artificial intelligence is a key tool. Machine Learning (ML) and Deep Learning (DL) are two subdomains of artificial intelligence (AI). It has numerous uses in the areas of Computer Vision it aids in the diagnosis and prediction of that virus. Deep learning and machine learning techniques can be used to create alerts to maintain social distance, diagnose and treat COVID-19, track COVID-19 events, create dashboards, forecast and for another potential mechanism of control [1].

Coronavirus infection is the world's most global epidemic, since there is no effective vaccine or cure for this virus, non-pharmaceutical treatments such as contact tracing, hospitalization or self-isolation, quarantine, group lockout and social distancing are the only ways to reduce the spread of infection [2].

Machine Learning (ML) can be used to manage vast amounts intelligently and data forecast disease spread. To monitor the virus, forecast the epidemic's progress and devise strategies and policies to control its spread, cloud computing and machine learning can be used efficiently. Researchers will use Machine Learning (ML) and Artificial Intelligence (AI) to predict when and where the virus will spread and then alert those provinces to make the necessary preparations [3].

This thesis presents a model to predict the movement of individuals in society by using Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) algorithms, In addition to creating a hybrid algorithm from the two previous algorithms.

1.2 Related works

It is known that the Coronavirus appeared at the end of 2019 and continues until this thesis was written and there are many studies and research that appeared during this period dealing with the use of (ML) in contributing to finding solutions in the discovery and method of dealing with this virus or pandemic. As for finding the relationship between the people's movement and time, we did not find any direction for this type of research.

Many types of research that have used machine learning and artificial intelligence to contribute to finding solutions to this pandemic or virus as reviewed:

- **Wang, B. et al. (2019)[4]:** They used the Internet of Things (IoT) technology to monitor the acquired data, process the data, and predict the next data using a neural network. One algorithm proposed is a two-layer model prediction algorithm based on Long Short Term Memory Neural Network and Gated Recurrent Unit (LSTM & GRU). set a double-layer Recurrent Neural Network to predict the PM2.5 value.

This model is an improvement and enhancement of the existing prediction method Long Short Term Memory (LSTM). The experiment integrates data monitored by the IoT node and information released by the national environmental protection department. First, the data of 96 consecutive hours in four cities were selected as the experimental samples. The experimental results are close to the true value. Then, selected daily smog data from 2014/1/1 to 2018/1/1 as a train and test dataset. It contains smog data for 74 city sites. The first 70% of the data was used for training and the rest for testing. The results of this experiment show that the hybrid model can play a better prediction.

- **Ardabili, S. F. et al. (2020) [5]:** They compared machine learning and soft computing models to Susceptible-Exposed-Infectious-Recovered (SEIR) and susceptible-infected-recovered (SIR) models to predict the COVID virus outbreak. Two models out of a wide range of (ML) investigated showed promise. There has been strong generalization potential for longer-term predictions in the results of the adaptive network fuzzy (ANFIS) and multi-layered perceptron (MLPs) deference scheme. This study shows that (ML) is an efficient method to model an outbreak based on findings and the highly complex nature of the viral outbreak and its behavior variations from nation to nation. Work also shows that the integration of machine learning and SEIR models will achieve true innovation in outbreak prediction.
- **Tuli, S. et al. (2020) [4]:** To evaluate and forecast how the epidemic will grow to control the disease in the future, predict epidemic progression and develop strategies and policies to combat its spread and an improved mathematical model was used. To forecast the possible danger of the Coronavirus in countries around the world, an improved model based on machine learning (ML) has been used. They demonstrate that iterative weighting can be better suited to the

development of a prediction System utilizing a Generalized Inverse Weibull. This was used on a cloud computing platform to forecast the growth behavior of the epidemic more accurately and in real-time. A more accurate data-driven method can be very helpful for government and citizen proactiveness. Finally, they proposed a range of study possibilities and premises for further use.

- **Yadav, M. et al. (2020) [6]:** They suggested that different tasks connected to the current COVID-19 be investigated using a new SVM method. They also use the supported vectors in this work to obtain greater classification accuracy rather than just a regression line. The approach is evaluated and contrasted with other well-known regression models on the normal datasets available. The positive results reflect both success and accuracy.
- **Kasilingam, D. et al. (2020) [2]:** They used infrastructures, the environment, policies and independent variables associated with infections to predict early containment to create predictive supervised machine learning models. Data about infection with coronaviruses was used in 42 countries. Logistic regression findings indicate a positive significant relationship and signs of early containment between healthcare infrastructure and lock-down policies. Logistical regression, decision tree, random forest and SVM machine learning models are being built to demonstrate precision from 76.2% to 92.9% to forecast early symptoms of infection containment.
- **Lars, L. et al. (2020) [7]:** Used a successful sampling algorithm, they employed a system of temporary point processes and their model to quantify the impacts on the path of business restrictions, measures of social distancing and the disease of various tracing and testing methods. Based on this algorithm, Bayesian optimization was used to estimate the rate of transmission from infectious persons at the sites

they visited and their homes, as well as the reduction of mobility because of the social distancing from longitudinal case statistics.

- **Arun, S. et al. (2020) [8]:** Evaluated the effects of various testing and tracing techniques, social distance measures and business constraints, using an efficient sampling algorithm for the temporal points method. Based on this algorithm, Bayesian optimization was used to estimate the rate of transmission from infectious persons at the sites they visited and their homes, as well as the reduction of mobility because of the social distancing from longitudinal case data.
- **Barstugan, M. et al. (2020) [9]:** Detected on abdominal Computed Tomography (CT) images were carried out utilizing machine learning methods. Four separate datasets were generated for the detection of the Coronavirus by taking patches from the 150 CT images in the dimensions 16x16, 32x32, 48x48, 64x64. Patches were used to improve classification efficiency in the process of feature extraction. The function extraction approach was used with Discrete Wavelet Transform (DWT), Gray Level Size Zone Matrix (GLSZM), Grey Level Run Length Matrix (GLRLM), Local Directional Pattern (LDP) and Gray Level Co-occurrence Matrix (GLCM). The extracted features were graded by SVM. To assess classification results, F-score metrics, accuracy, specificity and sensitivity were used. With 10-fold cross-validation and the GLSZ M feature extraction process, the best classification precision was achieved as 99.68 %.
- **ArunKumar, K. et al. (2021)[10]:** They suggested state-of-art deep learning Recurrent Neural Networks (RNN) models to predict the country-wise cumulative confirmed cases, cumulative recovered cases and cumulative fatalities. The Gated Recurrent Units (GRUs) and Long Short-Term Memory (LSTM) cells along with Recurrent Neural Networks (RNN) were developed to predict the future trends of the

COVID-19 and used publicly available data from John Hopkins University's COVID-19 database and emphasize the importance of various factors such as age, preventive measures, and healthcare facilities, population density, etc. that play a vital role in the rapid spread of COVID-19 pandemic.

- **Koç, E. and Türkoğlu, M. (2021) [11]:** Suggested a deep learning solution focused on a deep long-term memory network to anticipate requests for health kits and the cases of Coronavirus outbreak numbers. A regression layer, a fully connected layer, a Multilayer LSTM network, a dropout layer and a normalizing layer is used for the proposed scheme. This model is used to predict the number of cases, the number of intensive care units and the amount of respiratory equipment during the coming days. A dataset with 77-days of Coronavirus data for sixty-eight days for training and nine days for testing was used to check the suggested method validity. The experimental outcomes showed the suggested MAPE values of (4,80 %, 3,29 %, 2,89 %), respectively (99,72 %, 99,85 %, 99,90 %) for the estimate of the cases, the respiratory kit number and intensive care beds.
- **Bodapati, S. et al. (2020) [12]:** Used proper learning models on time-series analysis, LSTMs and RNN was applied to predict future patterns such as the number of people recovered from COVID-19, virus deaths and reported positive viral cases number in the coming days. The data collection used was used for encouraging experimental outcomes.

1.3 Problem Statement

One of the important problems that emerged with the emergence of Covid-19 is the spread's speed and relationship with the people's movement and social distancing and find the relationship between time and people's

movement to determine crowded places to control and limit coronavirus spread speed by taking the necessary measures and choosing an appropriate forecast model for predicting social movement.

1.4 Aim of Thesis

This thesis aims to find the relationship between time and people's movement to controls coronavirus spread speed by predicting the people's movement and identifying crowded places (standard, high, low) to take the necessary measures by using recurrent neural network algorithms (RNN).

1.5 Outline of Thesis

Besides this chapter, the remaining parts of this thesis include the following chapters:

Chapter Two: Theoretical Background

In this chapter, the theoretical tools and techniques that were used in this thesis are presented.

Chapter Three: The Proposed Model

This chapter introduces the steps of the proposed prediction system, with its design and implementation.

Chapter Four: Experimental Results and Evaluation

This chapter presents the results of the proposed model implementation, analysis, testing and evaluates these results.

Chapter Five: Conclusions and Suggestions for Future Work

This chapter presents the conclusions of this work. Furthermore, it provides suggestions for future work.