

# AI Based Algorithms in the Emergency Room for Predicting Patient Waiting Time in the Queue System

Arvind Kumar Kachhi<sup>1</sup>, Prof. Vikash Verma<sup>2</sup>, Prof. Saurabh Sharma<sup>3</sup>, Prof. Vishal Paranjape<sup>4</sup>

<sup>1,2,3,4</sup>Global Nature Care Sangathan Group of Institutions, Jabalpur(M.P), India

**Abstract--** Many hospitals use the amount of time patients spend in lines to gauge how crowded their emergency rooms (ERs) are. Many ER departments have exorbitant wait times, which makes it difficult to adequately treat patients and raises overall expenditures. Modern methods like deep learning (DL) and machine learning have been widely used in queuing system applications. In addition to, or instead of, queuing theory, this work intends to utilise DL algorithms for historical queuing variables to estimate patient waiting times in a system (QT). Four optimization algorithms—SGD, Adam, RMSprop, and AdaGrad—were used. To select the model with the lowest mean absolute error, the algorithms were compared (MAE). For further comparisons, a conventional mathematical simulation was employed. The outcomes demonstrated that by activating a lowest MAE of 10.80 minutes (24% error reduction) to estimate patient waiting times, the DL model is applicable when employing the SGD method. By achieving the highest performing model to better prioritise patients in the queue, this work makes a theoretical contribution of estimating patients' waiting times with alternative methodologies. Additionally, this work provides a useful contribution by utilising actual ER data. In addition, we suggested methods that would forecast patient wait times more accurately than a conventional mathematical approach. Using information from electronic health records (EHRs), the queue system in the healthcare industry can quickly adopt our method.

**Keywords—**Healthcare Management, Patient Priority, Waiting Time, Deep Learning, Queuing Theory

## I. INTRODUCTION

Most hospitals' emergency rooms (ER) are severely overcrowded with patients since they receive more than 50% of all hospital admissions. Due to the importance of the ER to hospitals, most departments need a lot of resources to accommodate the lengthy patient lines (Mor et al. 2015). Queuing is a danger in a context like healthcare since downtime may be costly for staff members and uncomfortable for patients. Additionally, it could have an impact on a patient's life or health circumstances (Gupta and Denton 2008). Traditional queuing theory (QT) is a

Historically, queuing systems have been studied using a mathematical technique (Gupta 2013). However, due to methodology limitations, such as unrealistic assumptions about the time distribution needed to perform queuing analysis, the typical QT technique may not be enough in real-world applications (Mahadevan 2015; Pianykh and Rosenthal 2015). Alternative methods, such as deep learning (DL) algorithms, are therefore thought to considerably increase ER effectiveness. Another category of machine learning technique is DL algorithms. Additionally, recent research revealed that the approach used to estimate patient wait times in emergency rooms had a limited degree of accuracy (Pak et al. 2020). In addition, DL algorithms are more accurate than conventional techniques while also reducing human error (Shafaf and Malek 2019). The purpose of this study was to create a brand-new, more precise model for predicting waiting times as well as a crucial tool for quick reactions in the event that emergency rooms report lengthy wait times. Due to prior studies' significant error rates, this objective was prompted. Compared to earlier research on this subject, the unique model used in this study minimises error prediction.

From a practical standpoint, DL was used to create a novel method that will increase the ER queuing predictor variables' ability to accurately estimate waiting times for patients with low acuity. The DL methodology was contrasted with conventional mathematical methods. Between January and December of 2018, realistic data from the triage monitoring system at an ER in Saudi Arabia with 30,909 patients was used.

According to recent studies, client dissatisfaction levels and waiting times have relationships (Abe 2019). They are therefore urged to think about allocating sufficient resources to reduce line waiting times. Customer service needs to be improved in the healthcare industry in particular if general happiness and successful health outcomes are to increase. Extreme wait times are a gauge of access to healthcare facilities and are associated with worse healthcare outcomes (Liang 2010).

To maximise queueing and resource utilisation, many strategies are commonly used (such as mathematical analysis) (Bittencourt et al. 2018). By analysing wait times in hospital pharmacies and other multiple points of service, queue models, for instance, are routinely used to handle excessive demand. Similar to this, queueing models are used in other service sectors that need security controls, such airports (Abe 2019). Additionally, the length of the line is used as a gauge for traffic management technique effectiveness. For instance, more than 90% of the delays in travel time and traffic congestion at the airport are caused by queuing delays (Peterson et al. 1995). Queuing models can also be used in daily life, such as when people wait in line for food at the grocery store or a restaurant. Longer wait times in any system may result in higher consumption, according to studies (Dong et al. 2019; Ülkü et al. 2020). Slow-moving lines increase waiting times and their prominence, which calls for the use of more resources.

The following is the study's contribution to the literature at the moment: First, using real data on the low patient acuity obtained from electronic health records (EHR) at an ER in Saudi Arabia, DL models were created alongside or in instead of queuing theory to estimate waiting time in a queue. The second improvement brought about by DL was a 24% decrease in prediction error as measured by the MAE metric. Third, taking into account model understandability and the feature extraction procedure, this study offers guidelines for waiting time analysis in the queue, not only in the healthcare industry but also in other sectors. These guidelines are based on trials carried out during the research. The outcomes, in our opinion, will be useful to practitioners and researchers who tackle related issues in other domains.

## II. LITERATURE REVIEW

Previous studies have demonstrated that prolonged waiting times cause patients to become frustrated, angry, anxious, and dissatisfied (Curtis et al. 2018; Sun et al. 2000; Ward et al. 2017). Numerous research have used various approaches to examine forecasts of ER waiting times. For instance, Kuo et al. (2020) used systems thinking and machine learning to forecast waiting times in emergency rooms. Arha (2017) employed many machine learning techniques, such as Elastic Net and Random Forest, to forecast the waiting time for low patient acuity in ER. Stage (2020) implemented a variety of approaches, including machine learning and a simulation, to predict patient waiting time.

Lastly, Curtis et al. (2018) created a number of machine learning techniques, including neural networks, to estimate patient waiting times while taking into account a variety of factors, including patient arrival time, service completion time, and examination. Additionally, research have created forecasting models using algorithms like quantile regression to estimate the length of time before therapy for patients with low acuity (Pak, Gannon, and Staib 2020). Our study is distinct from earlier studies on this subject since we used many DL optimization strategies to increase accuracy. Additionally, we took into account other predictors by gaining access to fresh data from the patient's entry into the queue (such as the minute, hour, and day), the length of the patient's wait in the queue, and the time of departure.

In their emergency rooms, many hospitals around the world regularly experience excessive wait times and crowding. Every year in the United States, there are steadily more visits to ERs (Di et al. 2015). The National Center for Health in 2016

According to statistics, there are roughly 145.6 million ER visits per year (Kea et al. 2016). Not only have ER visits climbed, but so have ER wait times. For instance, according to a 2017 report from the Canadian Institute for Health Information, ER wait times have significantly increased since 2015. A workable solution to these issues is to assess the effectiveness of emergency rooms (Rasouli et al. 2019).

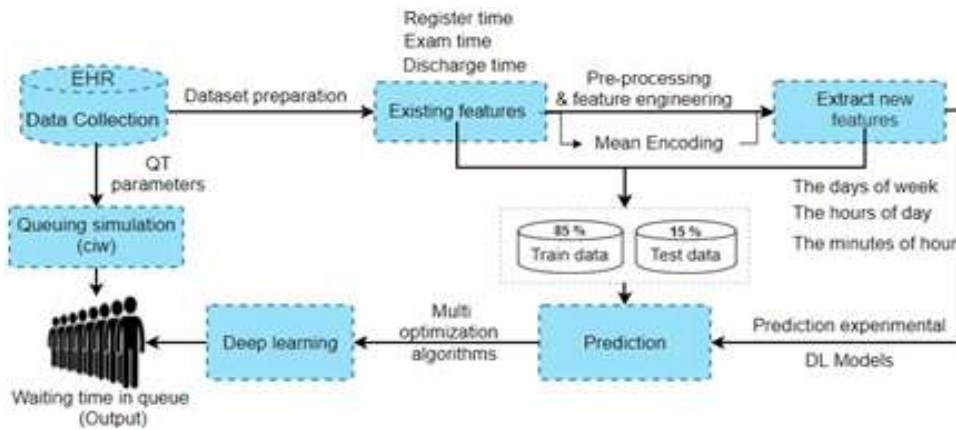
By analysing patient arrival times, some hospitals are utilising queueing models to improve staffing levels and optimise patient care (Kaushal et al. 2015; Sasanfar et al. 2020). The medical sector is finding increasing usage for predictive models. Seasonal arrival and waiting times can be reduced by using historical data to estimate future patient wait times (Ruben et al. 2010; Cai et al. 2016). The information stored in the EHR is essential for analysing and resolving healthcare issues that could have hidden components. Other research has concentrated on improving the healthcare queueing system, particularly how it might be applied to the development of predictive models for the analysis of future behaviour (Eiset et al. 2019). Additionally, the concept of machine learning has been used to analyse the projection of queuing behaviour (Srivastava 2016; Stagge 2020). The two research projects rely on a predictive modelling strategy, although their work on time series analysis on queue data prediction is flawed. Historical waiting times were evaluated in the multi-hospital study by Dong et al. (2019), and the results demonstrated that patients take into account ER waiting times when determining where to seek medical care.

The previously provided data is assisting in operational decisions that reduce waiting times and crowding in the emergency room (Abir et al. 2019). In order to forecast queueing behaviours in businesses, Stintzing and Norrman (2017) compared optimization using queueing theory with artificial neural networks (ANN) as a prediction method. The results using ANN, according to the scientists, were encouraging and might be applied to forecast the ideal level of service each day. Numerous forecasting methods have been used in queue analysis to reduce wait times (Moreno-Carrillo et al. 2019, however our model with multiple optimization algorithms can be used to assess ER wait times for low-acuity patients.

As a result, by using EHR data, the suggested model in this study can be utilised to inform ER medical staff about how long patients will wait in line.

### III. RESEARCH METHODOLOGY

Deep learning techniques are implemented in this study to predict patient waiting time in queueing system alongside, on the other hand, queueing theory (QT) using EHR data. Next, we compare the DL algorithms to find the best model with the lowest MAE. The model is presented, and a flowchart of the proposed methodology is shown in Figure 1. Each step is illustrated in the following subsections.



**Figure 1. Proposed methodology for the study**

#### 3.1 Data Description and Preparation

The Saudi Arabian Ministry of Health created the Triage Monitoring system, a national database, to guarantee the standard of patient treatment. The information was taken from the Triage Monitoring system, which also has information on how lines formed and were served in ERs at hospitals between January and December of 2018. From the time of registration until the patient departed the hospital, it tracked and recorded the patient flow. These are the main references used in the use of machine learning, particularly in the estimation of the waiting time for a new patient who joins a queue. These data included wait times, arrival and registration times, wait times in the queue, wait times at the server point, wait times for doctor examinations, and the total amount of time spent on all system activities (length of stay).

The primary information taken from the EHR was inserted at random. The data was cleaned, analysed, and finalised in several processes. In the first step, we translated our data utilising the arrival time/register time into weeks.

Step 2 created daily statistics using information from Step 1. Step 3 involved sorting the data according to arrival time to arrange the entries in order of arrival. Step 4: We removed the data's high values and missing values (which were caused by manual data entry errors) from our analysis. Only patients with level 3 through 5 acuities were included in step 5 because they made up more than 70% of the data we collected, which after data cleaning contained about 30,909 patients who were employed in the training model. Additionally, non-relevant variables from the dataset were eliminated, including patient ID and names.

The triage service time was combined with the waiting time because the dataset only reported one server (a doctor's examination) (for time from arrival to first time being seen by doctor). These patients are regarded as less urgent or non-urgent. These patients often receive care according to the order in which they arrive and do not require immediate attention. The waiting time that the machine learning algorithms attempt to anticipate is thus the output variable in this model.

The dataset's mean waiting time was 44.76 minutes, the median was 39.0 minutes, and the standard deviation was 20.23 minutes. To provide preliminary insights into the data from our model, a variety of input variables, including service time, waiting time, and individuals waiting vs days of the week, were examined. The service period in our data is the span from the beginning of the patient's medical care and its conclusion. In this instance, the dataset (new characteristics collected) is utilised to determine the number of patients in the line as well as the number of patients who join the queue. Every time a patient left the queue, we added up the waiting time and the arrival time to find the total number of individuals in the queue, and then we counted the number of people who remained when a new patient joined the queue. Data preparation and feature selection are often employed techniques in machine learning.

### 3.2 Pre-processing and Feature Engineering

The feature selection (selection of predictors) is an essential element in the machine learning model structure that determines the model's performance (Chandrashekar and Sahin 2014). There were main features extracted (e.g., minute, hour, and day) from the patient joining the queue in this study. Also, the patient's waiting time in the queue and leaving time were extracted. The following three are their main features:

1. Day was in the range of Monday (0) to Sunday (6).
2. Time in hours from 0 to 23rd hour.
3. Time in minutes starting at 0 minutes and continuing through 59th minute.

We applied different optimizer algorithms, including Adam, Adagrad, RMSprop, and SGD for the iterative update of network weights based on our data training and to describe the math behind the algorithms; equations (2) to (12) below are cited and summarized from Ruder (2016). Stochastic gradient descent (SGD) optimization algorithm does not change during training for all

$$\theta = \theta - \eta \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)}) \quad (2)$$

Adam (adaptive moment estimation) is a combination of Root Mean Square Propagation (RMSprop) and momentum. It can also be used instead of the standard stochastic gradient descent procedure to update network weights iteratively based on training data (Kingma et al. 2014).

The categorical features were encoded using mean target encoding and extracting the new features from current features in the dataset. We adopted this method (feature extraction) as presented by Kyritsis and Michel (2019), which was applied in the bank. The mean target encoding was used to encode our data with the new features because it is a fast way to get most of the categorical variables encoded and gives higher cardinality features for regression problems (Pargent et al. 2019).

### 3.3 Prediction Experimental

The experiment on machine learning in this study used TensorFlow version 2.0.0-beta1 and Python version 3.7.3. Also, different libraries were used to prepare and pre-process the data, such as Matplotlib, Date Time and Pandas. Accordingly, to validate and test the sensitivity of the model's performance, we split our dataset into two factions: the test set was 15%, and the training set was 85%. The test set was kept hidden throughout the training process. Moreover, by validating our model, it means that we used a test harness was used to give a fair estimation of the model's performance for making predictions on new data because it shows how sensitive the method is to applied data or new data that can be introduced to the model. Different optimization algorithms were used for this model in order to find the best with the lowest MAE. MAE is one of the metrics used to measure the machine learning model performance accuracy; it gives an idea of the magnitude of the error. It calculates all recorded means for the absolute errors by subtracting the prediction value from the actual value, as shown in the Eq. (1):

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \lambda(x_i; \eta)| \quad (1)$$

weight updates and the learning rate and maintains a single learning rate (termed alpha). A learning rate is maintained for each network weight (parameter) and separately adapted as learning unfolds. In contrast SGD performs a parameter update for each training example  $x^{(i)}$  and label  $y^{(i)}$ :



$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \tag{3}$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \tag{4}$$

Where:  $m_t$  is estimates of the first moment (the mean) and  $v_t$  is the second moment (the uncentered variance) of the gradients, respectively, hence the technique's name.  $m_t$  and  $v_t$  are initialized as vectors of zero (as the author of Adam observed that  $m_t$  and  $v_t$  are biased towards zero, especially during

the initial time steps and when the decay rates are low (e.g.,  $\beta_1$ , and  $\beta_2$  are close to one)).  $m_t$  and  $v_t$  counteract the biases by computing bias corrected first and second-moment estimates:

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \tag{5}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t} \tag{6}$$

Then, to update the parameters, we use these as shown in RMSprop, which yields the Adam update rule:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t}} \hat{m}_t + \epsilon \tag{7}$$

RMSprop maintains per parameter learning rates which are adapted based on the average of recent magnitudes of the gradients for weight, such as how quickly it is changing. RMSprop is very effective but an unpublished, adaptive learning

rate method proposed by Geoff Hinton in Lecture 6 slide 29 of his class (McMahan and Streeter 2014). In fact, RMSprop is identical to the first update vector derived from AdaDelta, which is an extension of AdaGrad optimization algorithm:

$$E[g^2]_t = 0.9 E[g^2]_{t-1} + 0.1 g_t^2 \tag{8}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t}} g_t + \epsilon \tag{9}$$

Where:  $E[g^2]_t$  is the decaying average over past squared gradients. Adaptive Gradient (AdaGrad) maintains a per-parameter learning rate that improves performance on problems with sparse gradients, such as computer vision and natural language problems (Brownlee 2020). For each parameter  $\theta_{t,ii}$  at every time step  $t$ , AdaGrad uses a

different learning rate. First, AdaGrad's per-parameter is updated, which then is vectorized, for brevity;  $g_t, i, ii$  is set to be the gradient of objective function w.r.t. to the parameter  $\theta_{t,ii}$  at time step  $t$ :

$$g_{t,ii} = \nabla_{\theta_{t,ii}} J(\theta_{t,ii}) \tag{10}$$

Stochastic gradient descent updates for each parameter  $\theta_{t,ii}$  at each time step  $t$  then becomes:

$$\theta_{t+1,ii} = \theta_{t,ii} - \eta \cdot g_{t,ii} \tag{11}$$

AdaGrad modifies its update rule, the general learning rate  $\eta$  at each time step  $t$  for every parameter  $\theta_{t,ii}$  based on the past gradients that have been computed for  $\theta_{t,ii}$ :

$$\theta_{t+1,ii} = \theta_{t,ii} - \frac{\eta}{\sqrt{\sum_{\tau=0}^{t-1} g_{\tau,ii}^2}} g_{t,ii} + \epsilon \tag{12}$$

Where:  $G_t \in \mathbb{R}^{d \times d}$  is a diagonal matrix where each diagonal element  $y_{ii}$  is the sum of the squares of the gradients w.r.t.  $\theta_{t,ii}$  up to time step  $t$ . And  $\epsilon$ : is a smoothing term that avoids division by zero (usually on the order of  $1e-8$ ).

Then, The Rectified Linear Unit (ReLU) was used in hidden layers. The ReLU activation function is a linear function that will output the input directly if it is positive ( $x$ ); otherwise, it will output zero.

(that is, if it receives any negative output it will return zero. (Hara et al. 2015)). It is used in this model because it achieves better performance and is more comfortable to

$$ff(x) = \max(0, x)$$

train when compared with other optimization functions (e.g., Sigmoid Function). ReLU can be written as Eq.(13):

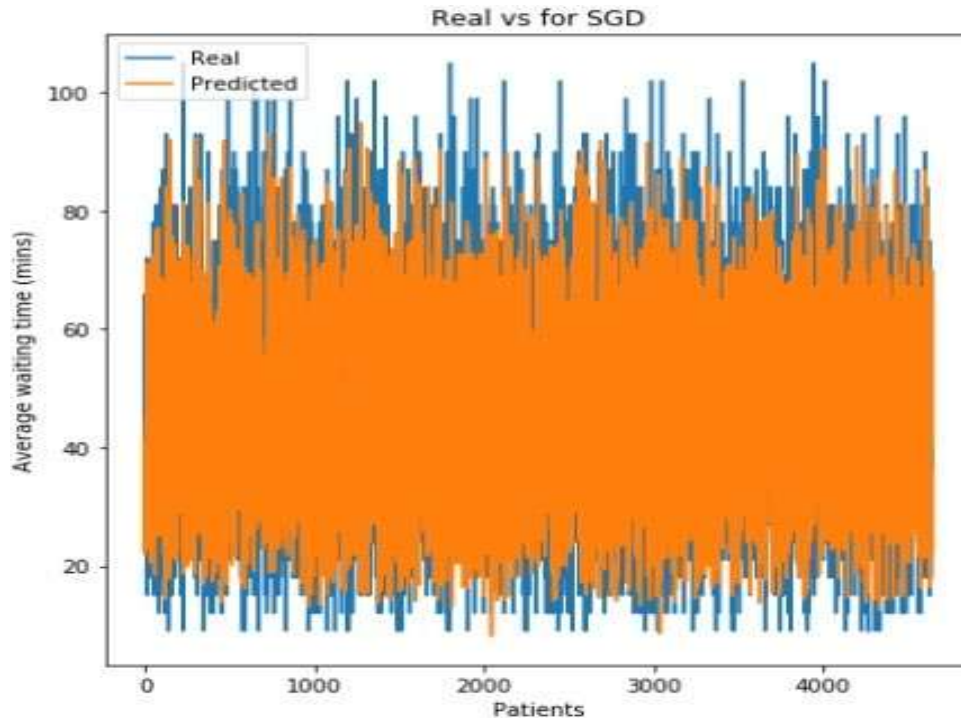
$$(13)$$

#### IV. RESULTS

In this study, we aimed to apply a DL approach with queueing theory. DL is one of the machine learning methods based on artificial neural networks. DL's power is the libraries built, such as Keras, which help to create extensive networks quickly and easily. Also, the simulation model for the queueing system was built to compare with the DL model using the Ciw library. Ciw is a discrete event simulation (DES) library supported by Python for queue networks (Palmer et al. 2019).

#### 4.1 DL Models

Keras library by Python was used to apply DL mode and was trained with four input visible layers, 25 neurons for the first hidden layer, 18 neurons in the next hidden layer, and one output in the output layer. After 150 epochs of model training. Figure 2 shows the model predicted average waiting time against actual waiting time for the best out perform optimization algorithms (SGD). The blue color represents the real (actual) waiting time, and the orange color represents the predicted waiting time. It shows the predicted waiting time as being closest to the actual waiting time. The idea of what score a good/poor model can achieve only makes sense when it is interpreted in the situation of the skill scores of other models and trained on the same data. For this purpose, different optimization algorithms are compared and trained on the same data.



**Figure 2. Waiting time predicted vs. actual for SGD algorithms**

In our experiment, the four optimization algorithms are listed in order from high to low MAE in Table 1. The stochastic gradient descent (SGD) had the lowest MAE with 10.80 minutes, followed by RMSprop, and then Adam and AdaGrad optimization algorithm with around 12 minutes.

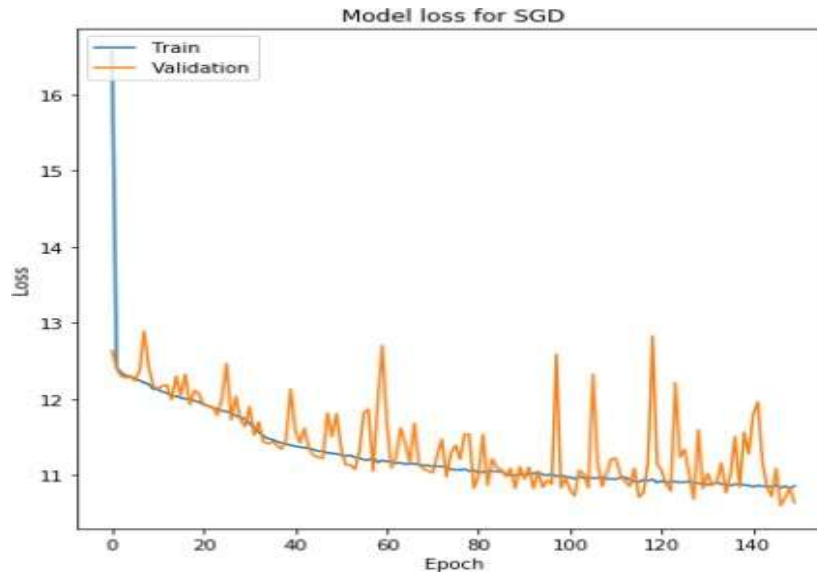
After exhaustive tuning of all related hyperparameters used in the model, [25-18-1] values of the architecture were found to be suitable for the neural network for this model.

**Table 1.**  
**Summary of the DL model (MAE results)**

Optimization Algorithms	Network Architecture	Mean Absolute Error
AaGrad	[25-18-1]	12.78 minutes
Adam	[25-18-1]	11.16 minutes
RMSprop	[25-18-1]	11.14 minutes
SGD	[25-18-1]	10.80 minutes

On the plot of loss, the model has comparable performance on both training and validation datasets for all optimization algorithms as shown in Figure 3. The loss on both datasets may use this as a sign to stop training at an earlier epoch if these parallel plots start to depart consistently.

Also, it shows the comparable skill on both train and validation data sets between the different optimization algorithms. The goal of using different optimizer algorithms is to change the attributes of our DL model, such as weights, learning rate, and to reduce the losses and reach the lowest MAE.



**Figure 3. Plot of model loss on training and validation dataset**

#### 4.2 QT Models

The classic approach of queue theory was used to simulate the model in this study. The arrival rate ( $\lambda$ ) and the service rate ( $\mu$ ) have been calculated using the same data applied in the DL model. Data was used for one day as shown in (Appendix A Table A), and is assumed the probability distribution of service time as an exponential distribution. The number of arriving patients per unit of time follows the Poisson distribution.

Because patients with level 3 through 5 acuties were used in this study, the queueing model of M/M/1 system for simulating the analysis was used to reflect the ER system. The model initially runs for 1,440 minutes (one day). To ensure that the simulation reflects reality, the model runs for ten simulations in a loop and takes warm-up time and cool-down times of 100-time units. Different seeds were considered every time, so each trial yielded different results.

Then, to achieve a more confident answer, the mean effect was taken over the trial results (68.24, 87.07, 59.377, 61.08, 63.64, 63.01, 70.75, 88.66, 63.64, 58.29 in minutes). Consequently, with each trial, the model ran for one day + 200 minutes (1,640 minutes). Patient mean waiting time resulting from the simulation model was 58.29 minutes, and the service time was 53.27 minutes. Comparing QT, the results to DL models in the dataset, the mean waiting time was about 44.74 minutes which is close to DL model prediction.

## V. DISCUSSION

Urgent and stochastic processes in the ER establish a challenge to waiting time prediction. For example, the ER provides high acuity patients with effective emergency care but is typically not as efficient with patients seeking attention for non-urgent ailments. This leads to increased ER occupancy. While non-urgent patients wait in the ER, patients requiring highly urgent attention bypass waiting times, which may increase the waiting times for those who are non-urgent. Also, patients with non-urgent cases may vary from case to case due to patient-level attention needs (e.g., level, 3 to 5). However, it was established that low acuity care has a significant impact on overall ER waiting and service times for high acuity patients (Arha 2017).

The goal of this study was to deliver a more accurate model for waiting time prediction and create an essential tool for reactive actions if ERs report long waiting time. For instance, this model compares other similar models for predicting waiting time in ER. Kuo et al. (2020) developed models with a mean-square error accuracy between 0.15 to 0.22. The model included different significant variables, such as Patient's triage categories, arrival time, number of doctors (within three hours of the Patient's arrival), number of patients in a queue (for triage, consultation, and departure upon the Patient's arrival). The proposed model, Kuo et al.'s model is limited by implementing a local triage system (Hong Kong). Also, a regional triage system by Arha (2017) estimated patient waiting time in an ER in Tennessee using a simple regression model. Arha used similar predictor variables (e.g., time of day, day of week, and month of year), and a mean square error as predictive accuracy (Arha 2017). To calculate this model variable and compare it with the proposed model requires collected clinical data.

Pak et al. (2020) also developed a waiting time prediction model for low acuity patients assigned to the waiting room with an overall accuracy of 20% mean squared prediction error; the proposed models with SGD and RMS prop algorithms reduced the prediction errors by 24% compared to model improvement in Pak et al. (2020).

This study has some limitations, including data availability; not all ER information was included, such as a patient type of injury, X-ray process time, and laboratory test time. What is available in the dataset was extracted. Also, DL is known as data hunger; in this case, data was collected for only one year. As shown in the results, 10.80 minutes was reached as the lowest MAE, but this could decrease if the amount of data increases. In the experiment, 30,909 patients (level 3 through 5 acuity) were used in training after removing other levels (level 1 to 2 acuity) and missing data. Significant improvement was shown in waiting time prediction with available data when compared with predicted average waiting time. Also, the model is simple enough to be implemented into an EHR system using relative information.

The second limitation is the patient levels in the local triage system (assigned as levels 1 to 5) may differ geographically. For example, this data, levels 3 to 5 were set as low urgent, and levels 1 and 2 as high critical, but this may be different in other ER triage systems globally. The third limitation is this is a single location study, which could potentially impact the accuracy of the model, requiring more work to validate the model by using data from other ERs in different locations and with different populations.

## VI. CONCLUSION

This paper proposed a novel model to improve the accuracy of waiting time prediction for low acuity patients using DL techniques and ER data. The study used historic queuing variables to predict patient waiting time in a queuing system alongside, or in place of, traditional approaches (queueing theory). The traditional methods may not be sufficient in real-life applications due to the limitations of the method, such as unrealistic assumptions of the time distribution required to do queueing analysis.

In the current literature, research reported that the methodology applied to predict patient waiting time in ERs has limited accuracy.



Furthermore, DL algorithms can reduce human error and achieve better accuracy, when compared with traditional methods. Thus, alternative techniques, such as DL algorithms, must be used to significantly improve efficiency. For this purpose, a novel model for waiting time prediction was created and as an essential tool for reactive actions if ERs report long waiting times. Furthermore, four optimization algorithms, including SGD, Adam, RMSprop, and AdaGrad, were compared to find the best accuracy considering MAE metrics. Also, algorithms were compared with traditional mathematical approaches and data was utilized from the triage monitoring system in Saudi Arabia. The results show that the DL model achieved better prediction accuracy than the traditional approach. Moreover, the novel model produced in this study resulted in a 24% error reduction when compared to prior work on this topic. The theoretical contribution of this paper is to predict patient waiting times with alternative techniques by achieving the highest performing model to better prioritize patient waiting in the queue. Also, this study offers a practical contribution by using real-life data from ERs. Furthermore, model have been proposed to predict patient waiting times with more accuracy than traditional mathematical models.

Future and extended work of this research could be as follows: more information from EHR could be implemented to the model such as different queueing predictor parameters. Moreover, different datasets from other hospitals and locations could be implemented. The service time of patients with the same acuity levels could be predicted. In addition, different machine learning algorithms could be applied to this model including linear and nonlinear regression. The model could be implemented on similar problems in different fields or sectors, including services and customer queueing. As part of future work, the model could be deployed as a web application to allow patients to join the queue prior to using EHR data.

#### REFERENCES

- [1] Abe, Y., Designing educative passenger journey by utilizing queueing and waiting times, Masters Theses Available: <https://www.theseus.fi/handle/10024/265246>, 2019.
- [2] Abir, M., Goldstick, J.E., Malsberger, R., Williams, A., Bauhoff, S., Parekh, V.I., Steven, K., and Jeffrey, S., Evaluating the impact of emergency department crowding on disposition patterns and outcomes of discharged patients, *International Journal of Emergency Medicine*, vol. 12, no. 1, pp. 1-11, 2019.
- [3] Arha, G., Reducing wait time prediction in hospital emergency room: a lean analysis using a random forest model, Masters Theses, Available [https://trace.tennessee.edu/utk\\_gradthes/4722/](https://trace.tennessee.edu/utk_gradthes/4722/), 2017.
- [4] Bittencourt, O., Vedat, V., and Morty, Y., Hospital capacity management based on the queueing theory, *International Journal of Productivity and Performance Management*, vol. 67, no. 2, pp. 224-38, 2018.
- [5] Brownlee, J., Gentle introduction to the Adam optimization algorithm for deep learning, machine learning mastery. Available: <https://machinelearningmastery.com/adam-optimization-algorithm-for-deep-learning/>, 2020.
- [6] Chandrashekar, G., & Sahin, F. (2014). A survey on feature selection methods. *Computers and Electrical Engineering*, 40(1), 16-28.
- [7] Cai, X., Oscar, P., Enrico, C., Fernando M., Richard D., David R., and Blanca G., Real-time prediction of mortality, readmission, and length of stay using electronic health record data, *Journal of the American Medical Informatics Association*, vol. 23, no. 3, pp. 553-61, 2016.
- [8] Chandrashekar, G., and Ferat, S., A survey on feature selection methods, *Computers and Electrical Engineering*, vol. 40, no. 1, pp. 16-28, 2014.
- [9] Curtis, C., Chang, L., Thomas, J.B., and Oleg, S.P., Machine learning for predicting patient wait times and appointment delays, *Journal of the American College of Radiology*, vol. 15, no. 9, pp. 1310-1316, 2018.
- [10] Dong, J., Elad, Y., and Galit, B. Y., The impact of delay announcements on hospital network coordination and waiting times, *Management Science*, vol. 65, no. 5, pp. 1969-1994, 2019.
- [11] Di S. S., Paladino, L. V., Lalle, I., Magrini, L., and Magnanti, M., Overcrowding in emergency department: an international issue, *Internal and emergency medicine*, vol. 10, no. 2, pp. 171-175, 2015.
- [12] Eiset, A.H., Hans, K., and Mogens, E., Crowding in the emergency department in the absence of boarding - a transition regression model to predict departures and waiting time, *BMC Medical Research Methodology*, vol. 19, no. 1, pp. 68, 2019.
- [13] Gupta, D., *Queueing Models for Healthcare Operations*, handbook of healthcare operations management, Springer New York LLC, vol. 184, pp. 19-44, 2013.
- [14] Gupta, D., and Brian, D., Appointment scheduling in healthcare: challenge and opportunities, *IIE Transactions*, vol. 40, no. 9, pp. 800-819, 2008.
- [15] Hara, K., Daisuke, S., and Hayaru, S., Analysis of function of rectified linear unit used in deep learning, *Proceedings of the International Joint Conference on Neural Networks*, Killarney, Ireland, 12-17 July 2015.
- [16] Kaushal, A., Yuancheng, Z., Qingjin P., Trevor, S., Erin, W., Michael, Z., and Aleks, C., Evaluation of fast-track strategies using agent-based simulation modeling to reduce waiting time in a hospital emergency department, *Socio-Economic Planning Sciences*, vol. 50, pp. 18-31, 2015.
- [17] Kea, B., Rochelle, F., Robert, A. L., and Benjamin, C. S., Interpreting the national hospital ambulatory medical care survey: United States Emergency Department Opioid Prescribing, *Academic Emergency Medicine*, vol. 23, no. 2, pp. 159-165, 2006-2010.
- [18] Kuo, Y. H., Nicholas, B. C., Janny, M. Y. L., Helen, M., Anthony, M. C. S., Kelvin, K. F. T., and Colin, A. G., An integrated approach of machine learning and systems thinking for waiting time prediction in an emergency department, *International Journal of Medical Informatics*, vol. 139, pp. 104-143, 2020.
- [19] Kyritsis, A.I. and Michel, D., A machine learning approach to waiting time prediction in queueing scenarios, *Proceedings of 2<sup>nd</sup> International Conference on Artificial Intelligence for Industries*, pp. 17-21, 2019.

- [20] Liang, T.K., Queuing for health care, Article in Journal of Medical Systems, vol. 36, no. 2, pp. 541-547, 2010.
- [21] Mor, A., Shlomo, I., Avishai, M., Yariv, N.M., Yulia, T., Galit, B. Y., Onpatient flow in hospitals: A data-based queueing-science perspective, Stochastic Systems, vol. 5.1, pp. 146-194, 2015.
- [22] Moreno, Atilio, Lina A., Julián, F., Camilo, C., Sandra, T., and Oscar, M.M., Application of queueing theory to optimize the triage process in a tertiary emergency care (ER) department, Journal of Emergencies, Trauma and Shock, vol. 12, no. 4, pp. 268-273, 2019.
- [23] McMahan, B., and Streeter, M., Delay-tolerant algorithms for asynchronous distributed online learning. In Advances in Neural Information Processing Systems, pp. 2915-2923, 2014.
- [24] Mahadevan, B., Operations Management Theory and Practice, 3rd Edition, Pearson Education, India, 2015.
- [25] Pak, A., Brenda, G., and Andrew, S., Predicting waiting time to treatment for emergency department patients, International Journal of Medical Informatics, vol. 145, pp. 104303, 2020.
- [26] Palmer, G.I., Vincent, A.K., Paul, R.H., and Asyl, L.H., Ciw: an open-source discrete event simulation library, Journal of Simulation, vol. 13, no. 1, pp. 68-82, 2019.
- [27] Pargent, F., Bischl, B., and Thomas, J., A benchmark experiment on how to encode categorical features in predictive modeling, Master Thesis, 2019. Peterson, M.D., Dimitris, J.B., and Amedeo, R.O., Models and algorithms for transient queueing congestion at airports, Management Science, vol. 41, no. 8, pp. 1279-1295, 1995.
- [28] Pianykh, O.S. and Daniel, I.R., Can we predict patient wait time? Journal of the American College of Radiology, vol. 12, no. 10, pp. 1058-1066, 2015.
- [29] Rasouli, H.R., Esfahani, A.A., and Mohsen, A.F., Challenges, consequences, and lessons for way-out to emergencies at hospitals: a systematic review study, BMC Emergency Medicine, vol. 19, no. 1, pp. 1-10, 2019.
- [30] Ruder, S., An overview of gradient descent optimization algorithms, Available: <https://arxiv.org/abs/1609.04747>, 2016. Ruben, A., Billy, J.M., Ying, P.T., Mark, H.D., Christopher, A.C., Song, Z., Gary, R., Timothy, S.S., Ying, M., and Ethan, A.H., An automated model to identify heart failure patients at risk for 30-day readmission or death using electronic medical record data, Medical Care, vol. 48, No. 11, pp. 981-988, 2010.
- [31] Sasanfar, S., Morteza, B., and Afrooz, M., Improving emergency department: simulation-based optimization of patients waiting time and staff allocation in an Iranian hospital, International Journal of Healthcare Management, vol. 16, pp. 1-8, 2020.
- [32] Shafaf, N., and Hamed, M., Applications of machine learning approaches in emergency medicine: a review article, Archives of Academic Emergency Medicine, vol. 7, no. 1, pp. 34, 2019.
- [33] Srivastava, T., How to predict waiting time using queueing theory? Available: <https://www.analyticsvidhya.com/blog/2016/04/predict-waiting-time-queueing-theory/>, December 17, 2019.
- [34] Stagge, A., A time series forecasting approach for queue wait-time prediction, Thesis, Available: <https://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1458832&dswid=9120,2020>.
- [35] Stintzing, J., and Fredrik, N., Prediction of Queuing Behaviour through the Use of Artificial Neural Networks, Thesis, Available: <https://www.diva-portal.org/smash/record.jsf?pid=diva2%3A111289&dswid=9120,2017>.
- [36] Sun, B.C., Adams, J., Orav, E.J., Rucker, D.W., Brennan, T.A., and Burstin, H.R., Determinants of patients' satisfaction and willingness to return with emergency care, Annals of Emergency Medicine, vol. 35, no. 5, pp. 426-434, 2000.
- [37] Ülkü, Sezer, Chris, H., and Shiliang, C., Making the wait worthwhile: experiments on the effect of queueing on consumption, Management Science, vol. 66, no. 3, pp. 1149-1171, 2020.
- [38] Ward, P. R., Philippa, R., Clinton, C., Mariastella, P., Nicola, D., Simon, A.C., and Samantha, M., Waiting for 'public and private hospitals: a qualitative study of patient trust in south Australia, BMC Health Services Research, vol. 17, no. 1, pp. 1-11, 2017.