

Received January 28, 2021, accepted March 30, 2021, date of publication April 21, 2021, date of current version May 26, 2021.

Digital Object Identifier 10.1109/ACCESS.2021.3074791

A Multi-Agent Approach for Personalized Hypertension Risk Prediction

SUNDUS ABRAR¹, CHU KIONG LOO¹, (Senior Member, IEEE),
AND NAOYUKI KUBOTA², (Member, IEEE)

¹Faculty of Computer Science and Information Technology, University of Malaya, Kuala Lumpur 50603, Malaysia

²Faculty of System Design, Tokyo Metropolitan University, Hachioji-shi 192-0397, Japan

Corresponding author: Chu Kiong Loo (ckloo.um@um.edu.my)

This work was supported in part by the UM Partnership Grant from the University of Malaya under Project RK012- 2019, in part by the IF017-2018 Office of Naval and Research Global (ONRG) Grant, U.K., under Project ONRGNICOP-N62909-18-1-2086, in part by the International Interfaculty Initiative in Computational Systems Care at Tokyo Metropolitan University, Japan, and in part by the University of Malaya, Impact Oriented Interdisciplinary Research Grant (IIRG)- IIRG002C-19HWB.

ABSTRACT Hypertension is a global health problem and a leading factor in severe and life-threatening cardiovascular diseases (CVD) and stroke. The onset is dependent on individual lifestyle choices, and no single root cause of the condition exists. Various machine learning solutions are proposed for the early diagnosis of hypertension and its prediction, but they are based on standard guidelines and do not provide personalized solutions. Current models mainly rely on batch learning methods and do not readily learn the new incoming data. There is also a lack of an intelligent technique for handling anomalies in data, which leads to unreliable prediction results. In this paper, an integrated multi-agent-based hypertension risk prediction system is proposed that detects and computes missing values in the time series and provides personalized hypertension risk predictions. The proposed solution incorporates Gaussian mixture models for enhancing the input data, and an Online Infinite Echo State Gaussian Process (OIESGP) is used to obtain real-time prediction distribution of blood pressure. The prediction system readily absorbs new incoming data, and the model is updated to learn any new patterns in the data. The hypertension risk score is estimated using the Framingham hypertension risk estimator, and a 4-year hypertension risk is computed. The prediction performance of the proposed model is evaluated on blood pressure data gathered from the Malaysian population using mean absolute error, mean square error, and root-mean-square errors. The experimental results indicate that the proposed prediction model exhibits greater prediction accuracy than existing state-of-the-art online prediction methods.

INDEX TERMS Blood pressure, Gaussian mixture models, hypertension, online infinite echo state Gaussian process, personalised prediction model.

I. INTRODUCTION

High blood pressure is a global health problem. It directly contributes to various chronic diseases such as stroke, cardiac arrest, memory loss, renal function failure, and multiple other disabilities [1]–[3]. In the year 2008, around the globe, an estimated 1 billion adults suffered from hypertension. According to statistics, the number of hypertensive patients is expected to rise to 1.56 billion by 2025 [4]. The number of hypertensive adults in Malaysia alone rose from 33.6% in 2011 to 35.3% in 2015. 36% of mortality in Malaysia is due to cardiovascular diseases (CVD) and is a significant factor in premature death [5], [6]. This growing prevalence

can be associated with varying lifestyle choices that act as behavioral risk factors, such as smoking, unhealthy eating habits, inadequate physical activity, harmful use of alcohol, and stress. Even though high blood pressure causes almost 6% fatalities worldwide, it is in fact a preventable disease and can be controlled with preventive measures [4]. Table 1 shows a brief classification of blood pressure value ranges (in mm Hg) for optimal and hypertension scenarios.

There is no single root cause of hypertension due to multiple unique risk factors related to genetics, environment and lifestyle for each individual [7]. This makes early diagnosis and prevention of the disease difficult for clinicians. One solution for controlling the disease is early diagnosis, which depends on constant monitoring of blood pressure. Constant monitoring helps in defining an individualized blood pressure

The associate editor coordinating the review of this manuscript and approving it for publication was Gang Wang¹.

TABLE 1. Blood pressure classification.

Category	Systolic		Diastolic
Optimal	< 120	&	< 80
Prehypertension	120 - 139	&/or	80 - 89
Hypertension			
Stage 1	140 - 159	&/or	90 - 99
Stage 2	160 - 179	&/or	100 - 109
Stage 3	>= 180	&/or	>= 110

profile and develop a personalized treatment plan for that patient. However, in real life, this is impractical and quite inconvenient for the patients to visit hospitals so frequently, and the traditional cuff-based measurement method using mercury is not easy to operate at home. Hence, the clinicians are left with incomplete or missing blood pressure readings. Due to the complexity associated with disease prediction, clinicians are prone to make errors, especially in these cases where the data is incomplete or noisy.

Recent advancements in telemedicine and neuro-computing have enabled researchers to bridge the gap between computer science and medicine. Machine learning, in particular, is being applied in the medical field to accurately analyze medical images and serve as medical diagnostic systems [8], [9]. A combination of genetic algorithms (GA) and k-nearest neighbours (KNN) is utilised in the prediction of paroxysmal atrial fibrillation [10] and feedforward neural networks are used in the diagnosis of diabetes mellitus [11]. Classifier systems are developed that allow early diagnosis of congestive heart failure using k-nearest neighbors, linear discriminant analyses, multilayer perceptron, support vector machines, and radial basis function artificial neural network [12]. Linear regression models, recurrent neural networks (RNN) and artificial neural networks (ANN) are extensively applied in blood pressure analysis [13]–[15] however, most of these solutions use a standard universal approach for all individuals, and there is no consideration for personalization. The models often require a large amount of input data for training, need a continuous input time series in order to make accurate predictions, and do not update the model as new input data is received. The effect of missing values in the input data is not considered, and no intelligent technique is used to handle missing values in the time series. Furthermore, no complete system exists for hypertension risk prediction in the Malaysian population.

Our study focuses on the implementation of a multi-agent architecture for real-time risk prediction of hypertension in the Malaysian population. The proposed system handles missing records in the blood pressure timeseries, besides identifying and removing abnormalities from the data set. A comprehensive data set is gathered from the Malaysian population, and their hourly blood pressure and heart rate are collected. A user agent is designed that is responsible for managing blood pressure and risk factors from the user. A data processing agent employs Gaussian mixture models (GMM) for cleaning the data and converting it into a complete-time series. A prediction agent uses an enhanced recurrent neural network combined with Gaussian

estimation known as the Online infinite echo state Gaussian process (OIESGP) for hourly predictions. The blood pressure status agent evaluates the risk of hypertension with the help of the Framingham classifier and presents a 4-year risk of Hypertension. There is no research focusing on the real-time risk prediction of hypertension in a multi-agent environment to the best of our knowledge,. The prediction performance of the proposed model is evaluated using blood pressure data gathered under the supervision of medical practitioners from the Malaysian population at public hospitals.

The rest of the paper is organized as follows: in section II, we briefly discuss the existing state of the art work on hypertension prediction. Section III describes the proposed multi-agent system and our methodology. In section IV, we explain our experiments for the system's evaluation. Section V presents the results, followed by the discussions in section VI, and finally, the conclusion and future work in Section VII.

II. RELATED WORK

Artificial intelligence is playing a significant role in digitalizing health care, and existing studies extensively discussed hypertension risk prediction.

Ambika *et al.* [16] developed a personalised decision support system based on a support vector machine (SVM) and fuzzy association rule mining (ARM) to predict the probability of acquiring hypertension. The missing values in the data are substituted using mean and mode value substitution and the interquartile range (IQR) technique is used to remove the outliers. The model enhanced the data using AdaBoost and predicted various stages of hypertension using boosted SVM. The model also took into account personal behavioral factors along with medical history for prediction. The model reported a prediction accuracy of 91.8%.

Mohammadi *et al.* [17] developed models based on logistic regression and recurrent neural networks (RNN), namely the long short term memory (LSTM), to predict the risk of uncontrollable hypertension in the upcoming 3-months. The developed model was evaluated on electronic health data from 17,000 patients, and any missing values in the data set were replaced with the average value of that variable. Patient data containing fewer than two entries were excluded from the analysis. The model achieved an area under the curve of 0.714 and 0.696 precision.

Kanegae *et al.* [18] proposed a prediction model for the onset of new hypertension. The model was tested on clinical data of hypertensive patients. Any missing values in the data were imputed using last observation carried forward, mean and mode substitutions. The model was a combination of logistic regression, random forest and XGBoost technique combined with the help of bagging technique. The model achieved 0.992 AUC.

Melin *et al.* [19] proposed a hybrid model based on a modular neural network with fuzzy inference for hypertension diagnosis. Each module in the neural network received systolic, diastolic, heart rate and age values respectively as

input variables and the neural network was trained using the backpropagation algorithm. Two fuzzy inference systems (FIS) handled nocturnal hypertension and heart rate, and a third FIS classified hypertension. The system classified blood pressure range and hypertension with 90% accuracy in each module.

Ye *et al.* [20] used a large dataset extracted from electronic health repositories (EHR) to predict essential hypertension for the upcoming year. Missing values in the dataset were imputed using k-nearest neighbours (KNN), and feature reduction is applied using the Cochran-Armitage trend test and logistic regression. The data is then fed to XGBoost, and the model is trained with the data of the previous 1-year consisting of chronic and medical history, health conditions, clinical utilization history, and social determinants. The prediction model achieved an accuracy of 0.971 and 0.87 in retrospective and validation, respectively.

LaFreniere *et al.* [15] developed a hypertension diagnosis system based on an ANN. The model used various risk factors affecting hypertension, which were identified based on the patient's health status, medical history, and geographical location. Eleven risk factors were identified and set as input nodes for the ANN with seven hidden nodes and two output nodes. The records with the majority missing entries were excluded from the analysis, and recordings with few missing lab entries were set to 0. The outliers were identified by evaluating z-scores, and data with extreme scores were excluded from the analysis. The model classified individuals as hypertensive or non-hypertensive with reported 82.3% accuracy and was trained on a sufficiently large data set, therefore capable of handling various test cases.

Wang *et al.* [21] the use of simple demographic data and lifestyle choices, contributing significantly to the onset of hypertension. They proposed a hybrid model for hypertension diagnosis based on logistic regression and ANN. A binary logistic regression model determines the factors affecting hypertension and identified 13 risk factors. The identified risk factors are then used by a neural network with back-propagation as input parameters to predict the onset of hypertension. The performance of the model was evaluated using questionnaire surveys. The records with missing data entries were excluded from the analysis. The model achieved an accuracy of 72.12% with minimal standard variations.

A brief comparison of the works mentioned above has been summarized in Table 2. To summarize the above discussion, existing approaches towards hypertension risk prediction:

- 1) Do not present a concrete system for personalized hypertension risk prediction.
- 2) Mostly ignore missing values in historical records, and noisy or anomalous data is not appropriately handled.
- 3) Do not perform online update of the prediction model and thus present potentially outdated models.

The main contributions of this research work are:

- 1) An integrated multi-agent-based architecture for personalized hypertension risk estimation in a mobile application.

- 2) A data preprocessing technique that can estimate missing values and remove outliers in the time series before the prediction process.
- 3) A personalized online prediction model that is capable of learning new patterns in the input data.
- 4) A 4-year hypertension risk predictor using the Framingham risk calculator

III. MULTI-AGENT FRAMEWORK

The goal of this research was to develop a personalized hypertension risk prediction system. The proposed solution can gather user blood pressure periodically and perform its prediction for the next 24 hours. It also estimates the 4-year risk of hypertension for the user to avoid the onset of hypertension. The system is developed using a multiagent approach such that the tasks are divided and distributed to four agents, with each agency responsible for a specific job. The agents communicate with each other on various steps to make the system more robust and scalable. The architecture diagram is given in Figure 1.

The user agent runs on the user device and collects all user-specific data (age, weight, height, gender, smoking status, and BMI) to build a personalized profile. It collects blood pressure and heart rate data hourly from the user and passes them on to the data processing agent. The data processing agent is responsible for preprocessing the data and saving it securely for further analysis. The blood pressure prediction agent requests the data processing agent for the hourly readings for the past 24 hours and the updated user profile to perform BP predictions for the next 24 hours. The prediction agent is also responsible for notifying the blood pressure status agent about the prediction results. The blood pressure status agent requests the results from the data processing agent and calculates the risk of hypertension. The results are communicated to the user agent, converting them into a user-readable format and presenting them.

The following section discusses details of each agent and its dedicated responsibilities. The algorithms used for computing missing values in the time series and the prediction of blood pressure are also explained below.

A. USER AGENT

The user agent is responsible for all interactions associated with the end-user. It has three main tasks: (1) collect user profile information, (2) collect hourly blood pressure and heart rate data, and (3) pass user data to the data processing agent.

The profile information consists of age, weight, height, gender, smoking status, and body mass index (BMI). The user agent also keeps track of updates in the profile data. It is also responsible for collecting the user's vitals (heart rate and blood pressure in this case) periodically and ensure this information is secured until it is passed on to the data preprocessing agent. The user agent achieves its tasks with the help of a mobile application, and Bluetooth low energy (BLE)

TABLE 2. Summary of existing research work on hypertension prediction.

Author	Year	Method	Personalised	Online	Advantages	Disadvantages
Ambika [16]	2020	Support vector machine	✓	✗	<ul style="list-style-type: none"> • Learns from imbalanced data • Introduces personalised prediction 	<ul style="list-style-type: none"> • Unintelligent schemes for data processing (mean/mode substitution)
Mohammadi [17]	2019	Logistic regression, recurrent neural networks	✗	✗	<ul style="list-style-type: none"> • Large dataset to train the prediction model • Includes blood profile along with BP data 	<ul style="list-style-type: none"> • Missing values replaced with averages • Ignores relevant clinical data during analysis
Kanegae [18]	2019	Logistic regression, Random Forest, XGBoost	✗	✗	<ul style="list-style-type: none"> • Large dataset used for model training • Identifies effective BP measurement for prediction • Age and BMI incorporated for BP prediction 	<ul style="list-style-type: none"> • Clinical BP unable to identify white-coat hypertension • Imputation of missing data introduces bias
Ye [20]	2018	XGBoost, KNN	✗	✗	<ul style="list-style-type: none"> • Used large feature set used for model training • Real time predictive model 	<ul style="list-style-type: none"> • KNN based for imputation can cause bias for patients with large number of missing data • Directly using EHR data did not capture all related risk factors
Melin [19]	2018	Artificial neural networks, Fuzzy System	✗	✗	<ul style="list-style-type: none"> • Fuzzy systems handles BP variability • BP variation during the night taken in account 	<ul style="list-style-type: none"> • Dataset assumed to be a complete time series • Insufficient data samples to capture blood pressure trend
LaFreniere [15]	2016	ANN	✗	✗	<ul style="list-style-type: none"> • Use large dataset to train model 	<ul style="list-style-type: none"> • Do not account for behavioural information
Wang [21]	2015	Logistic Regression and ANN	✗	✗	<ul style="list-style-type: none"> • Lifestyle factors effecting hypertension are identified and used for prediction 	<ul style="list-style-type: none"> • Removes records with missing entries, thus resulting in loss of data

enabled wrist band worn by the user. The mobile app provides a visual interface for the collection of profile data.

The wrist band measures the user blood pressure and heart rate periodically. After receiving the reading, the user agent then passes the complete data, including the user profile and the physiological data, to the data processing agent. The other agents use this data to perform prediction and estimate the risk of hypertension. The user agent receives the prediction results and the risk score, which displays them to the user in the mobile application. Fig. 2 shows the working of the mobile app. In this way, the user is provided with an overview of his physical health condition, enabling him to take preemptive measures to avoid any health risk.

B. DATA PROCESSING AGENT

The data processing agent manages the data set used for training the prediction model and storing prediction results. It receives data from all the other agents and stores it in a

database for further processing. The main tasks of the agent are below:

- 1) Preprocessing: This comprises estimating missing data entries and removal of outliers.
- 2) Storage: This comprises saving received data securely for further analysis.
- 3) Communication: This comprises transmitting the data to the relevant agent.

Data preprocessing is a computationally expensive task and requires a significant amount of time. Recent advancements in cloud technology enabled researchers to shorten this processing time by carrying out most of these tasks in the cloud. The large storage capacity and faster processing resources in the cloud system enable various machine-learning algorithms for data cleaning, knowledge discovery, and analysis to run simultaneously. The proposed data preprocessing agent runs on the cloud and carries out

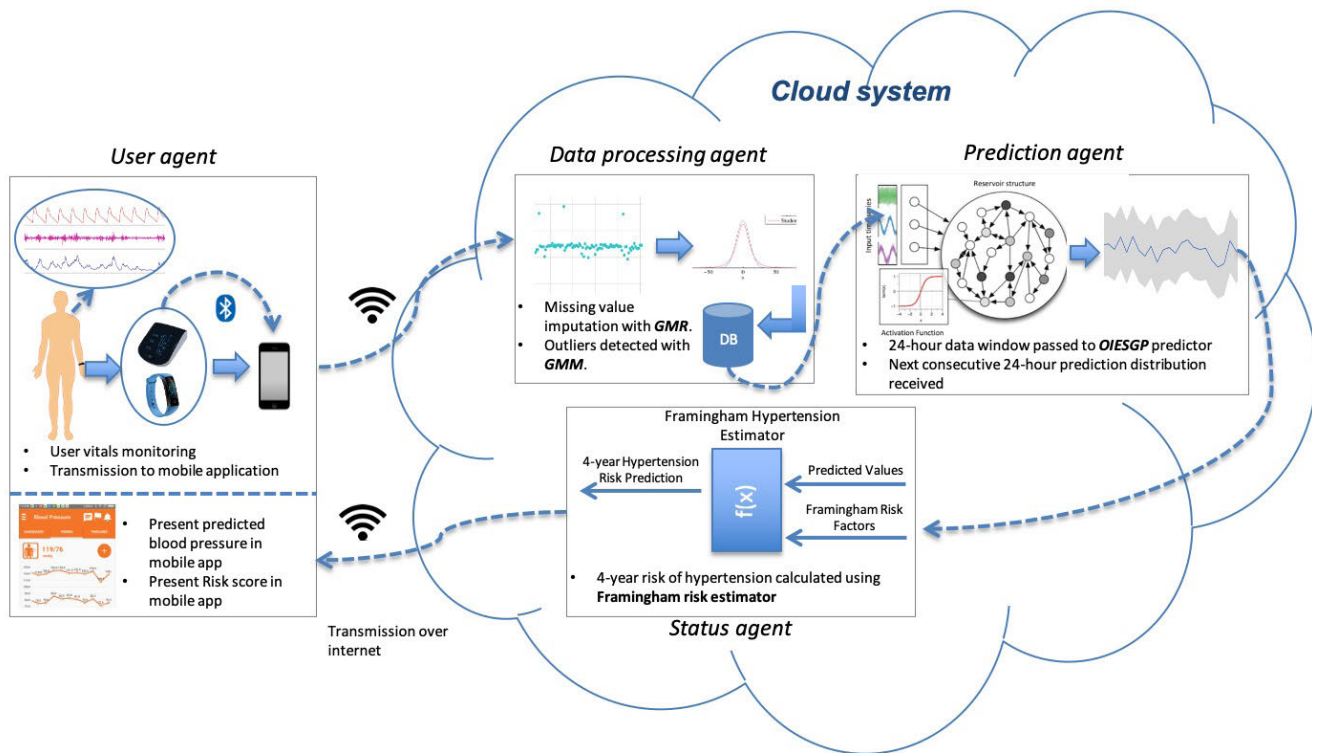


FIGURE 1. Detailed flow chart of proposed framework. Data collection is performed with the help of smart wearables and a smartphone, whereas data preprocessing and predictive analysis is performed on the cloud system. After the hypertension risk estimation, the results are communicated back to the smartphone and presented in a user-friendly interface.

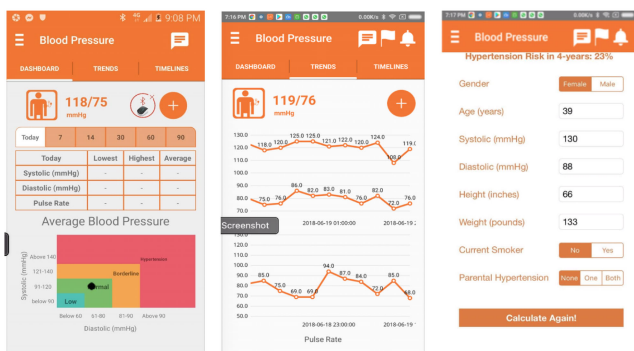


FIGURE 2. Implemented mobile application showing blood pressure prediction and 4-year hypertension risk score.

further processing in an AWS cloud instance. The cloud system hosts a database and an app-server for performing the respective preprocessing tasks.

The developed database contains three main tables: the user profile containing all physiological details is stored in a *User* table, the blood pressure and heart rate records are stored in a *Blood_Pressure* table, and the prediction results are stored in a *Prediction* table. The design of the database is given in Fig. 3. All other agents request the data processing agent for the historical blood pressure readings and user profile values for blood pressure prediction and hypertension risk score calculation.

The physiological data readings received from the wrist band have the potential to be plagued with erroneous values. This can occur by an occasional device malfunction or human error, e.g., it may happen that the user forgot to wear the device after taking it off, or disposition of equipment that might cause faulty measurement value. During the prediction model’s training, this flawed data introduces uncertainty and unreliability in the prediction results. For this reason, it is of utmost importance that the input data is processed and anomalies are removed.

In the general population, blood pressure is known to exhibit a Gaussian (normal) distribution pattern [22] [23] which makes Gaussian mixture models a suitable candidate for modelling blood pressure behaviour. For this study, we compute the missing records in the blood pressure historical data using Gaussian mixture regression (GMR) [24] and identify outliers using Gaussian mixture models (GMM) [25].

1) MISSING DATA ESTIMATION

A significant responsibility of the data processing agent is to ensure that the data stored in the database is a complete-time series without any missing entries. This task is of utmost importance because most of the learning algorithms expect a continuous time series to make accurate predictions, and the presence of missing values reduces the accuracy in the results [26]. Values may be missing due to various reasons depending upon the data source, such as human error,

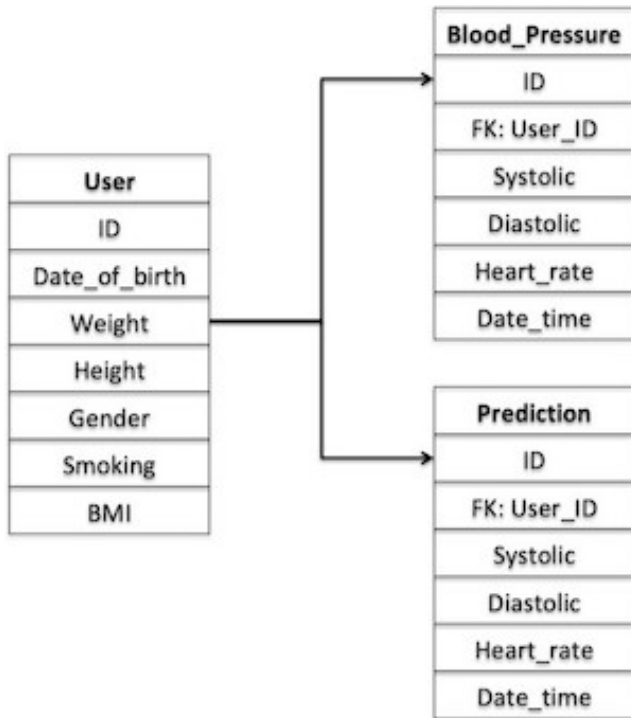


FIGURE 3. Structure of database used in the data processing agent: The User table stores personal profile, whereas hourly readings are stored in the Blood Pressure table, and hourly predictions are stored in the Prediction table.

malfunctioning device, or environmental error. There are several techniques to handle missing values in data. In most situations, simple techniques such as mean/median/mode substitution, complete case analysis and missing-indicator method are used [27]. Most studies treat missing values by complete removal of a specific signal, averaging observed data, replacement or substitution constructed from previous information, normalization, and linear interpolation [28], [29]. These methods may cause a biased model because of the loss of information. Furthermore, they underestimate standard deviation since they do not consider the uncertainty in missing values.

The agent takes a two-step approach for missing data estimation. The first step is the identification of missing values. This is achieved by iterating through the received time series and comparing the timestamps associated with each record. Upon encountering non-consecutive time stamps, missing data and the intermediate data records are inserted into the time series with some initial value. Algorithm 1 presents detailed steps for identifying the missing values.

The next step is to estimate values for the missing records. We use Gaussian mixture regression [24] to estimate the missing heart rate, systolic, and diastolic blood pressure values. A Gaussian distribution is modelled using the input data. A Gaussian mixture comprises several Gaussians, with each containing: a mean that defines the centre of the distribution and covariance that represents its width and weight that determines the size of the Gaussian function. The data is

Algorithm 1 Missing Value Detection

Global variable

D , Time series with missing values
 X , Single record
 x , feature in a record
 d_c , Current record
 d_p , Previous record
 d_m , Value for Missing record
 T_d , Time interval between two records

end Global variable

$X = \{x_1, x_2, \dots, x_n\}$

$D = \{X_0, \dots, X_t\}$

$X \geq 0$

$i \leftarrow 1$

while $i \leq \text{Length}(D)$ **do**

$d_c \leftarrow D[i]$, $d_p \leftarrow D[i-1]$

$T_d \leftarrow \text{CalculateTimeDifference}(d, d_p)$

if $T_d > 1$ **then:**

for $j \leftarrow 0$ to T_d **do:**

$X'' \leftarrow \text{NewRecord}(d_m)$

$D.\text{InsertAtIndex}(j, X'')$

$d_p \leftarrow X''$

end for

end if

$i \leftarrow i + 1$

end while

sequentially scanned for missing values, and a buffer with the continuous values is populated. A GMR then constructs a sequence of Gaussian mixture models (GMM) from the time series data to compute the joint probability density of the real data. Then the conditional density and regression functions for each model are constructed. The underlying joint density $f_{X,Y}$ is given by Equation (1):

$$f_{X,Y}(x, y) = \sum_{j=1}^K \pi_j \varnothing(x, y; \mu_j \varepsilon_j) \quad (1)$$

where π_j represents the weight of the j -th mixture component and $\varnothing(x, y; \mu, \varepsilon)$ is the probability density function (PDF) of a multivariate Gaussian distribution with the mean μ and covariance ε . Partitioning each component \varnothing_j as in [24], the joint density can be expressed as in Equation (2):

$$f_{X,Y}(x, y) = \sum_{j=1}^K \pi_j \varnothing(y|x; m_j(x), \sigma_j^2) \varnothing(x; \mu_{jX}, \varepsilon_{jX}) \quad (2)$$

The mean vector $m_j(x)$ and covariance matrix σ_j^2 are calculated in Equation (3) and Equation (4):

$$m_j(x) = \mu_{jY} + \varepsilon_{jYX} \varepsilon_{jX}^{-1} (x - \mu_{jX}) \quad (3)$$

$$\sigma_j^2 = \varepsilon_{jYY} - \varepsilon_{jYX} \varepsilon_{jX}^{-1} \varepsilon_{jXY} \quad (4)$$

Then the conditional PDF $Y|X$ from the GMMs is generated, which is given in Equation (5):

$$f_{Y|X}(y|x) = \sum_{j=1}^K w_j(x) \varnothing(y; m_j(x), \sigma_j^2) \quad (5)$$

where $w_j(x)$ is the mixing weight and is given in Equation (6):

$$w_j(x) = \frac{\pi_j \mathcal{O}(x; \mu_{jX}, \varepsilon_{jX})}{\sum_{j=1}^K \pi_j \mathcal{O}(x; \mu_{jX}, \varepsilon_{jX})} \quad (6)$$

Finally, the missing value is generated using the regression function $m(x)$ given in Equation (7):

$$m(x) = E[Y|X = x] = \sum_{j=1}^K w_j(x) m_j(x) \quad (7)$$

Algorithm 2 presents the process of imputing the identified missing values using GMR.

Algorithm 2 Missing Value Imputation Using GMR

Global variable

K , Complete time series

G , Array of Gaussian mixture models

d_m , Missing value

W_p , Train data buffer

S_p , Size of prediction buffer

for $j \leftarrow 0$ to $\text{Length}(K)$ **do**:

$V \leftarrow K[j]$

if $V = d_m$ **then**:

G .ComputeDensity(W_p) {Equation. 5}

$V'' \leftarrow G$.PredictNextValue() {Equation. 7}

$K[j] \leftarrow V''$

else

W_p .insert X''

end if

if $\text{Length}(W_p) > S_p$ **then**:

W_p .pop()

end if

end for

When a missing value is encountered in the time series, the probability density is recomputed using the previous available data and the next value is estimated. The missing value is then replaced with this predicted value. To control the influence of past values on the prediction, a buffer size is maintained and the past data is flushed to make space of only recent data which is used for prediction.

2) OUTLIER DETECTION

An outlier is a value of a variable different from the generally observed pattern for that variable. It arises when an unexpected or inconsistent data point is observed where a design or distribution was found previously. The leading cause of outliers is the presence of anomalous data. There may be an error in data transmission and its reception, human error, system error, or instrumental error. In biomedical data, a sudden spike can also be observed during a workout or a physical tasking activity. These abnormally high values are found for a short period before they reach back to their original values. However, if recorded, these values act as outliers in the actual trend of the biomedical series.

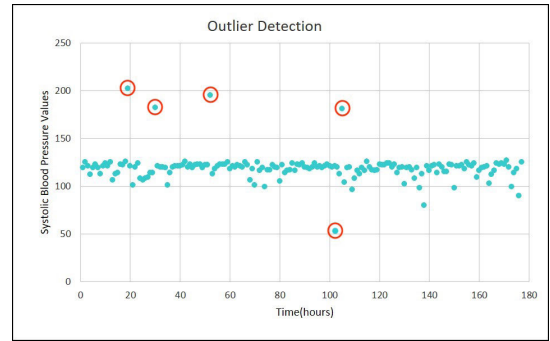


FIGURE 4. Existence of outliers in time series data.

Fig. 4 shows the existence of anomalous values in blood pressure values.

The presence of outliers or anomalies in the blood pressure time series weakens the results in the prediction agent. Therefore, the last step in the data processing is the removal of any outliers. Healthy adults exhibit a rhythm in their blood pressure in a period of 24-hours, with the highest blood pressure during the day and the lowest during the night [30]. This pattern is generally believed to be caused by the behavioural triggers and daily variations in postures and physical activity throughout the day [31], [32], and GMM detects seasonal variations in time series data [25]. Therefore, we utilize the same approach to model this seasonal variation in the blood pressure data and detect any outliers in the time series.

The time series is divided into time bins corresponding to the total number of seasonal variations that can be expected in the time series. Let $B = [B_1, B_2, \dots, B_n]$ represent the total bins that the data is divided into, and each $B_n = \{T_1, T_2, \dots, T_n\}$. For blood pressure data, we divided the data into 24 time bins, with each bin corresponding to a single hour of the day. Each T_n contains all the physiological data of the user for the hour value n . The next step is to compute the density of each time bin using GMM.

Let the time series is represented by $X = \{x_1, \dots, x_n\}$. We can assume that X is generated by a GMM with K number of components. The function $f_k(x_i)$ is then the probability density function of the k -th component, representing the probability of x_i generated by the k -th component. So, the PDF of the GMM is given by Equation (8) as below:

$$P(x_i) = \sum_{k=1}^K \pi_k f_k(x_i | \mu_k, \varepsilon_k) \quad (8)$$

In the above Equation (8), π_k is the weight of the k -th component; μ_k and ε_k are the mean vector and covariance matrix of the k -th component and $P(x_i)$ represents the probability of x_i generated by the GMM. Moreover, $\sum_{k=1}^K \pi_k = 1$. The PDF of the k -th component is calculated using Equation (9):

$$f_k(x_i | \mu_k, \varepsilon_k) = \frac{\exp\{-\frac{1}{2}(x_i - \mu_k)^T \sum_k^{-1}(x_i - \mu_k)\}}{(2\pi)^{\frac{p}{2}} |\varepsilon_k|^{\frac{1}{2}}} \quad (9)$$

For every data point in the time bin, its density is calculated using Equation (9), and a score is assigned to it by examining its distance from the other data points on the overall probability density scale. The scores are calculated by taking the log of probabilities as given in Equation (10):

$$OS = \log(p(x))^{2f} \quad (10)$$

Here $p(x)$ represents the density function of x , and f is a value used to scale the log values. The value of f is directly proportional to the calculated score, and a larger value of f would create a larger difference between the score values of the outliers.

Algorithm 3 Outlier Detection Using GMM

Global variable

T_x , Time bin containing records for an hour value x
 P , Probability densities
 OS , Outlier score
 G , Array of Gaussian mixture models

end Global variable

initialise;

$P \leftarrow G.ComputeDensity(T_x)$

for $j \leftarrow 0$ to $Length(P)$ **do**:

$p = P[j]$

$OS \leftarrow CalculateScoreFor(p)$ {Equation 10 is used here}

if $OS > 5$:

$T_x[j].IsOutlier \leftarrow true$

end if

end for

Algorithm 3 gives an overview of the outlier detection process. A scaled value of 0 or 1 indicates that the data point in consideration is normal, whereas a scaled score of 9 or 10 indicates extreme outlier points. For our experiments, all scaled score values greater than 5 are classified as outliers. The OS is further scaled to a value in a range of 0 to 10 using Equation (11). Presenting outlier scores in such a way helps in better visualization of the scores, and comparison with other techniques becomes easy.

$$ScaledScore = \frac{OS - \min(OS)}{\max(OS) - \min(OS)} \times 10 \quad (11)$$

C. PREDICTION AGENT

The main task of the prediction agent is to perform blood pressure prediction for the next 24-hours. It requests the data processing agent for the past 24-hour readings and the updated user profile and passes them to the machine learning algorithm, which performs the prediction. The prediction algorithm receives the time-series data in the form of 24-hour windows, and the predicted output is, in turn, fed back to the system to update the model. The prediction agent receives new information about the blood pressure and user profile continuously, updating the model for future predictions. The constant updating of new information poses two potential problems for the hypertension risk prediction system:

- 1) As continuous biological data is infinite in nature, this would lead to an infinite model size.
- 2) As the size of the model grows indefinitely, the computation time for updates and predictions would also increase.

To avoid an infinitely growing model and maintain a shorter computation time, the proposed solution should have the ability to perform online learning and somehow prevent the model size from growing indefinitely. For this purpose, we employ an online learning algorithm built on the echo state networks (ESN) [33] called the Online Infinite Echo State Gaussian Process (OIESGP) [34]. The following sections explain the working of ESN and OIESGP algorithms in detail.

1) ECHO STATE NETWORKS

The echo state network is a recurrent neural network proposed by Jaeger [33] that belongs to the Reservoir computing framework. In this approach, the weights of the hidden layer neurons and the reservoir are not trainable, and only the output weights are trained. A randomly generated reservoir (which is the RNN) is driven using input signals, and the output is received by use of a combination of the reservoir units. The inputs are connected to the reservoir with an activation function, and the outputs have weights that are learned with Linear Regression and are connected to the reservoir as well. The state of the reservoir is updated during training using Equation (12),

$$x_{t+1} = (1 - \gamma)h(Wx_t + W_i u_{t+1} + W_b d_t) + \gamma x_t \quad (12)$$

where x_t is the state of the reservoir at a given time t , γ is the leak rate, $h(\cdot)$ is the activation function, W represents the reservoir weight matrix, W_i is the input weight matrix, u_t is the input, W_b is the output feedback weight matrix and d_t is the desired output. After training, the update equation is represented as in Equation (13):

$$x_{t+1} = (1 - \gamma)h(Wx_t + W_i u_{t+1} + W_b y_t) + \gamma x_t \quad (13)$$

The predicted outputs y_t are obtained using Equation (14),

$$y_{t+1} = W_o \psi_{y+1} \quad (14)$$

In Equation (14), W_o is the linear output weight matrix and ψ_{t+1} represents the augmented reservoir state and input vector, given in Equation (15) as,

$$\psi_{t+1} \triangleq [x_{t+1}; u_{t+1}] \quad (15)$$

2) ONLINE INFINITE ECHO STATE GAUSSIAN PROCESS

The OIESGP [34] is an online variant of echo state networks combined with Bayesian learning for Gaussian processes. The reservoir is assumed to have an infinite size, and the recurrent kernel incorporates automatic relevance determination. The structure of the network is given in Fig. 5. In this technique, the model performs fast successive updates just as new incoming data arrives. The model size is maintained by only storing unique neural states; non-unique states do not affect the model size and do not affect the computation

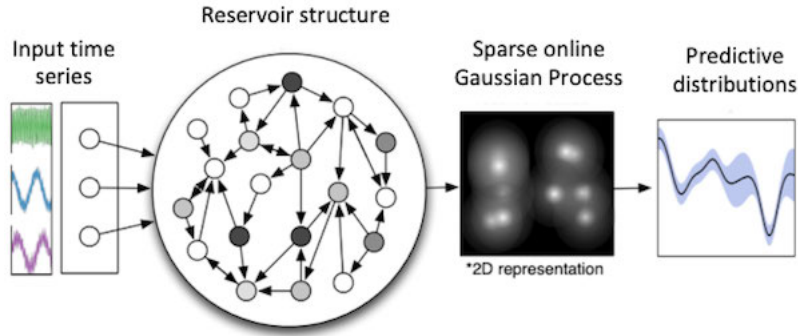


FIGURE 5. Online infinite echo state Gaussian process: learns from temporal sequences and produces predictive distributions (Soh & Demiris, 2012).

time. The output is a predictive distribution with estimated variances, which helps specify a possible range of predicted values instead of a hard predicted value.

To model the OIESGP network, the ESN is updated using Equation (12) and a new composite state ψ_{t+1} is calculated. Then Bayesian online learning for regular Gaussian processes is applied to yield a posterior as shown in Equation (16),

$$\hat{p}(f|\tilde{y}_{t+1}) = \frac{P(\tilde{y}_{t+1}|f(\psi_{t+1}))p_t(f)}{(P(\tilde{y}_{t+1}|f(\psi_{t+1}))p_t(f))_t} \quad (16)$$

This posterior is then projected to the closest GP measured via the Kullback-Leibler divergence [34], $KL(\hat{p} \parallel q)$ and q is the required approximation. To maintain the model size, the number of reservoir states to be retained (also termed as the *basis vectors*, $b \in \beta$) is reduced. This is done so by calculating a score presented by [35] for each state as in Equation (17):

$$\gamma(\psi_{t+1}) = k_r(\psi_{t+1}, \psi_{t+1}) - k_{\beta,t+1}^T K_{\beta,t}^{-1} k_{\beta,t+1} \quad (17)$$

where $k_r(\psi_{t+1}, \psi_{t+1})$ is the reservoir kernel function such as the radial basis function, and $k_{\beta,t+1} = [k_r(b_i, \psi_{t+1})]_{b_i \in \beta}$ and $K_{\beta,t}^{-1} = [k_r(b_i, b_j)]_{b_i, b_j \in \beta}$. For a certain state ψ_{t+1} , if the computed score $\gamma(\psi_{t+1})$ is higher than some defined threshold, then the reservoir is updated. The mean and variance of the predictive distribution are given in Equation (18) and Equation (19) which are used to estimate the prediction distribution:

$$\mu_* = k_{\beta,t}(\psi_{t*})^T \alpha_t \quad (18)$$

$$\sigma_*^2 = k_r(\psi_{t*}, \psi_{t*}) + k_{\beta,t}(\psi_{t*})^T C_t k_{\beta,t}(\psi_{t*}) \quad (19)$$

D. BLOOD PRESSURE STATUS AGENT

The blood pressure status agent is responsible for interpreting the predicted results and generate user-readable responses. After performing a prediction, the prediction agent passes the prediction results to the data access agent and informs the blood pressure status agent. After getting notified, the blood pressure status agent can safely request all the predicted results from the data access agent and interpret those results. A popular means of explaining the risk of hypertension is the Framingham risk calculator [36].

1) FRAMINGHAM HEART STUDY

The Framingham Heart Study (FHS) is a long-term ongoing study on the people of the town Framingham, Massachusetts, USA [37]. The study aimed to identify the risk factors that influence the development of cardiovascular diseases. The investigation began in 1948 with 5,209 adults, which is now on its third generation of participants. The participant's ages ranged from 30 to 62 years old, and none of them had any history of cardiovascular diseases or accidents. With time, the descendants of the original group also became a part of the study.

2) FRAMINGHAM RISK SCORE FOR HYPERTENSION

The Framingham risk-score-calculator estimates the risk of coronary heart diseases for 10 years based on several risk factors that include age, gender, smoking history, previous treatment of hypertension, BMI, and last blood pressure values. These risk factors are termed cardiovascular risk factors in the Framingham Heart Study. The study uses the regression model [38] for the risk calculation, and the result is called the Framingham risk score.

The score represents the risk of developing any cardiovascular disease in a period of 4-years at any given time. The input parameters for the risk calculator would be age, sex, systolic and diastolic blood pressure, smoking habit, if parents had hypertension, and body mass index. The complete scoring system is explained in Fig. 6. In order to determine a user's risk of developing hypertension, the BP status agent makes use of the Framingham risk calculator, which is in Equation (20),

$$FHS = 1 - \exp \left[-\exp \left(\frac{\ln(4) - 22.94954 + \sum X\beta}{0.8769} \right) \right] \quad (20)$$

where,

β = the coefficient of regression,

X = the level of each variable

If the gender is male, the variable is assigned a 0, if it is female then it is assigned a value of 1.

Variable	Beta**	p-value	Hazard Ratio	95% CI
Age	-0.15641	< 0.001	1.195	(1.089, 1.312)
Sex	-0.20293	0.004	1.260	(1.091, 1.456)
SBP	-0.05933	<0.001	1.070	(1.060, 1.080)
DBP	-0.12847	<0.001	1.158	(1.087, 1.234)
Smoking	-0.19073	0.013	1.243	(1.058, 1.460)
Parental Hypertension*	-0.16612	0.014	1.209	(1.047, 1.395)
BMI	-0.03388	<0.001	1.039	(1.025, 1.054)
Age times DBP interaction	0.00162	0.005	0.998	(0.997, 0.999)

(Scale = ± 0.87692 , Intercept = ± 22.94954)

FIGURE 6. Standard Framingham risk score parameters [38].

The risk calculator evaluates the input variables and presents a 1, 2- and 4-year risk of hypertension for the said user. These results are passed on to the user agent which displays them in user readable format.

E. PERSONALISED PREDICTION MODEL

The global burden of hypertension increased drastically over the past few years [4]. The main factors leading to this increase are population growth, unhealthy lifestyles, obesity and ageing. It also does not come as a surprise that hypertension is controlled by less than 20% of the population that is suffering from the disease [39]. Early diagnosis is key to controlling hypertension and, if treated correctly, can act as a modifiable risk factor for other cardiovascular diseases.

Evidence suggests that the causes and, in turn, the effect of hypertension are unique for every individual. The blood pressure level and extent of organ damage in individuals depend not only on the environmental factors (such as dietary intake, physical inactivity, mental health) but also on the individual genetic structure [40]. Hence any single method dedicated to treating hypertension or the relevant organ damage may not be effective for all hypertensive patients. To date, most solutions for managing hypertension utilise standard universal guidelines that do not take into account personal and environmental factors in the control of the disease [41]. There is a need for personalised data modelling and a more personalised approach to study patients with hypertension to develop an effective solution.

The advanced learning and prediction abilities of AI can be leveraged to achieve the goal of personalised prediction. As detailed in Section II, AI is mainly used to investigate risk factors instead of managing the disease. There is no study on estimating the personalised risk of hypertension in individuals to the best of our knowledge. The proposed model provides personalised blood pressure predictions for the next 24-hours and a 4-year hypertension risk score. All this is

presented in an interactive mobile application to promote self-awareness, empowerment, a healthy lifestyle and medication adherence.

IV. EXPERIMENTAL SETUP

To evaluate the performance of the proposed system, we performed three experiments; the first was to assess the missing data estimation in the blood pressure time series. The second experiment removes the outliers in the data set. The third experiment evaluates the predicted systolic, diastolic, and heart rate values for the next 24-hours and calculates the 4-year hypertension risk score. The input data is divided into training and testing data and further divided into 24-hour windows representing the blood pressure recordings of a single day. The prediction algorithm accepts a window of data and predicts the successive 24-hour blood pressure window. The algorithms were coded in Python programming language and ran in a Linux environment.

A. DATA DESCRIPTION AND FORMULATION

The dataset for this study was gathered from the Malaysian population through clinical pilot programs held at the University Malaya Medical Centre (UMMC). Their basic physiological profile and their systolic, diastolic and heart rate readings for six months were collected. Their initial physiological readings were measured by a medical health professional who served as a baseline for future readings. The recorded data included systolic pressure, diastolic pressure, heart rate, age, BMI and smoking status. Table 3 and Table 4 show the data structure and the detail data types of each item in the data set.

To ensure continuous blood pressure monitoring, the patients used a mobile application and wrist bands that measure blood pressure and pulse rate every hour. The wrist band used photoplethysmography [42], which is a process of using light waves to measure blood flow. The wrist band emitted a burst of green light on the subject's wrist for a short interval and captured the refracted response. This information, along with the user's position and motion information gathered from the device's accelerometer, determined the heart rate and blood pressure reading. After successfully recording a reading, it was transmitted to the mobile device for further processing.

In case the user took off the wrist band at the time of measurement of the reading, an option to manually store their BP reading in the mobile application is also available. The mobile app passed this data to the data processing agent for further processing. The system handled each patient data set separately to obtain reliable personalized predictions. For this research, a patient was selected at random, and their blood pressure data set was used for further experiments. If the total duration of data collection is denoted by t , and the number of daily readings is denoted as d , then the total number of records for each patient, denoted by L can be determined as in Equation (21):

$$L = d \times t \quad (21)$$

TABLE 3. Data values collected from the user.

Variable	Description
Systolic Blood Pressure	Continuous Variable (mmHg)
Diastolic Blood Pressure (DBP)	Continuous Variable (mmHg)
Age	Continuous Variable (year)
Gender	Binary Variable (Male:1 / Female:0)
Body Weight	Continuous Variable (kg)
Body Mass Index	Continuous Variable (kg/m ²)
Smoking Status	Binary Variable (Yes:1 / No:0)
Pulse Rate	Continuous Variable (bpm)
Date Time	Continuous Variable (dd/MM/yy hh:mm:ss)

TABLE 4. Detail of data types.

Data Item	Detail
Subject Information	Age Birthday Gender (Male/Female) Smoking (Yes / No)
Blood Pressure	Measurement date and Time Systolic Blood Pressure (SBP) Diastolic Blood Pressure (DBP) Heart Rate(HR)
Body Composition	Measurement date and time Body Weight Body Mass Index (BMI)

For our experiments, the data was collected for 180 days. Data was stored in the form of hourly records for a day, then the total number of data records, $L = 4320$. Each record was a collection of systolic pressure, diastolic pressure and heart rate. These values were further combined with patient profile data, i.e. age, gender, BMI and smoking status, to create an input matrix X with the above features. If the total features are denoted by $p = 8$, then X can be expressed as in Equation (22):

$$X = \begin{bmatrix} x_{1,1} & \cdots & x_{1,p} \\ \vdots & \ddots & \vdots \\ x_{L,1} & \cdots & x_{L,p} \end{bmatrix} \quad (22)$$

B. EVALUATION CRITERION

Among the most popular metrics available for prediction accuracy comparison in time series analysis, we chose Mean Absolute Error (MAE), Root-Mean-Square Error (RMSE) and Mean Square Error (MSE) to evaluate the prediction performance. These are the most common metrics for measuring the errors in continuous variables.

The mean absolute error [43] for the prediction model is given by the following formula,

$$MAE = \frac{1}{n} \sum_{t=1}^n |A_t - \hat{A}_t| \quad (23)$$

A lower MAE value indicates better fit for the prediction model thus indicating superior prediction accuracy. Root mean square error or root-mean-square deviation [44] is

represented by the following equation,

$$RMSE = \sqrt{\frac{\sum_{t=1}^n (y_t - \hat{y}_t)^2}{n}} \quad (24)$$

Mean square error [45] is the average of the square error which is denoted by the following formula,

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (25)$$

V. RESULTS

In this section, we present the experimental results to test the personalized blood pressure prediction system.

A. MISSING VALUE IMPUTATION

To maintain data integrity and ensure no bias in data is introduced, we used Gaussian mixture models to estimate missing values in the original data set. As shown in Fig. 7, missing values of pulse rate, systolic and diastolic blood pressure record are predicted by GMR. Further experiments were conducted to verify the performance of the algorithm. The number of missing values was gradually increased in the test data set, and the MAE, MSE, and RMSE were recorded. The results are presented in the Fig. 8. The increase in error rate is quite low, which indicates that the algorithm stability even in the presence of missing values.

B. OUTLIER DETECTION

The outlier score for each data point in the systolic, diastolic and heart rate time series is calculated using Equation (10). Fig. 9 shows the outlier scores calculated for each value in the sample data set. Outliers are identified based on these scores. A higher score indicates an outlier, whereas a lower score indicates a correct reading. The results show that the system assigns extreme values a higher score, and normal values are assigned lower scores. Thus the outliers are correctly identified and later removed.

As part of further experimentation, we introduced varying outliers in systolic, diastolic and heart rate data and investigated how accurately our system detected those outliers. Table 5 shows results of detected outliers.

C. BLOOD PRESSURE PREDICTION

The continuous blood pressure and heart rate for 24 hours of a random day were selected from the completed time series to predict the next 24 hours. Fig. 10 depicts the prediction distribution of systolic, diastolic and heart rate values predicted using the OIESGP algorithm. The shaded region in the graph represents the prediction probability distribution. It shows that the actual and predicted values lie within the prediction probability distribution, hence proving the proposed algorithm to be a reasonable estimate for blood pressure prediction.

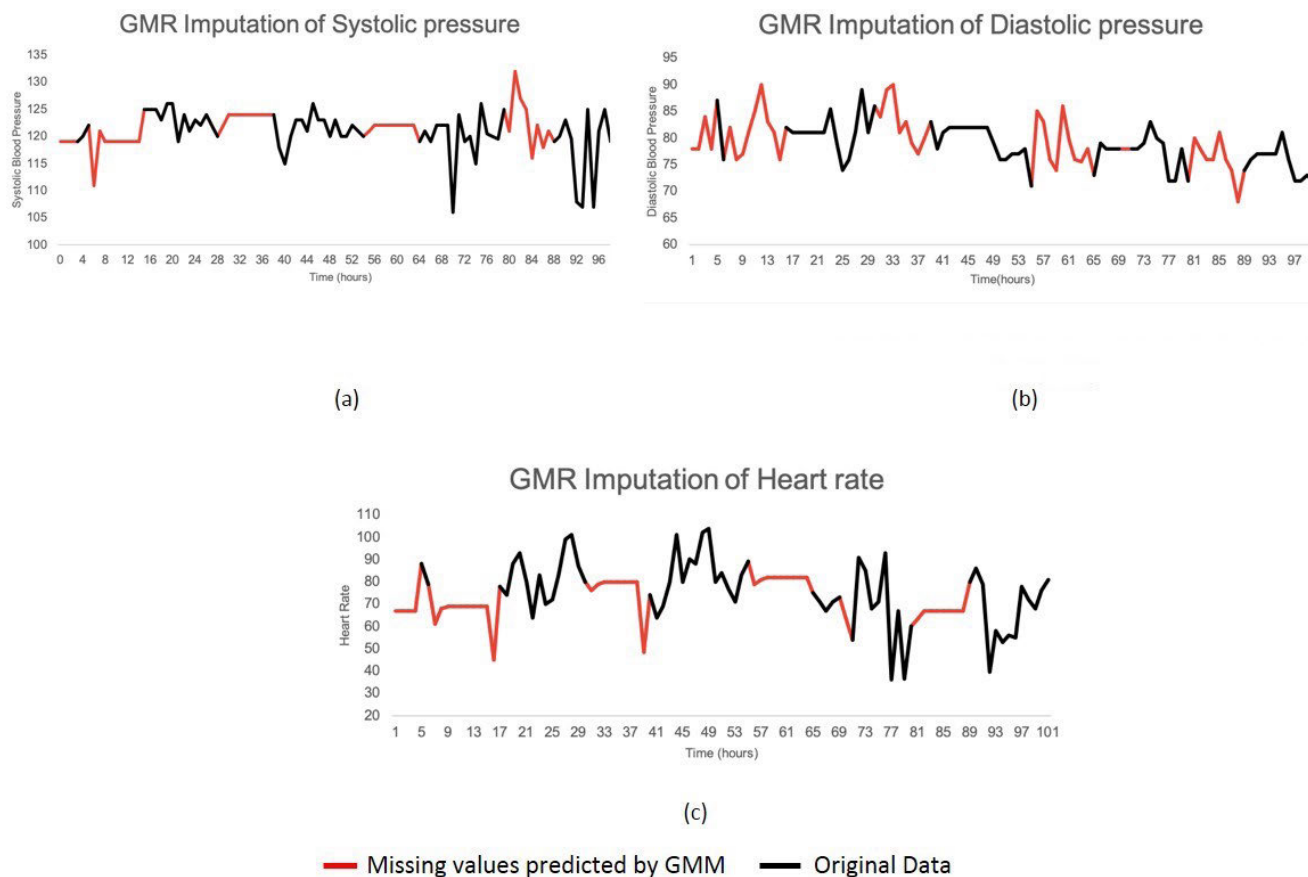


FIGURE 7. Incomplete time series is estimated using Gaussian mixture regression. The figures above show the missing data estimation in (a) Systolic, (b) Diastolic, (c) Heart rate data.

TABLE 5. Introduced outliers in data and detection with GMM.

Outliers Introduced	Outliers Detected	Outliers Removed
Systolic		
6	5	5
8	8	8
10	9	9
Diastolic		
5	5	5
8	8	8
9	9	9
Heart Rate		
4	4	4
6	6	6
7	7	7

For further verification, prediction for consecutive days was also performed. Blood pressure data for five days was selected at random, and the prediction results of the next consecutive five days were recorded. Each day is considered a window, and the prediction error in each window was recorded. Fig. 11 shows the prediction error in each successive window. It is observed that with each successive window, the error value decreases, indicating an increase

in the prediction accuracy of the algorithm. The following section presents a comparison of the prediction results with state of the art offline and online prediction algorithms.

1) COMPARISON OF RESULTS

We conducted three experiments to further verify the superiority of the proposed technique. The prediction performance was compared against three existing state of the art online methods namely: passive-aggressive regressor [46], online recurrent extreme learning machine (OR-ELM) [47] and fully online sequential extreme learning machine (FOS-ELM) [48]. These algorithms were selected based on their recurrent nature and online learning ability. Systolic, diastolic and heart rate values of each user were predicted separately, and the average prediction errors were recorded as given in Table 6. The results show higher accuracy for OIESGP predictions as compared to the other online methods.

To demonstrate the superiority of the proposed online prediction scheme, we compared the prediction results with existing batch learning algorithms currently used in medical time series and hypertension risk prediction. We compare the results with Artificial neural networks [15], [19] namely the Long Short Term Memory (LSTM) [49], Bayesian based Gaussian regression [50], [51] and Support Vector Machine

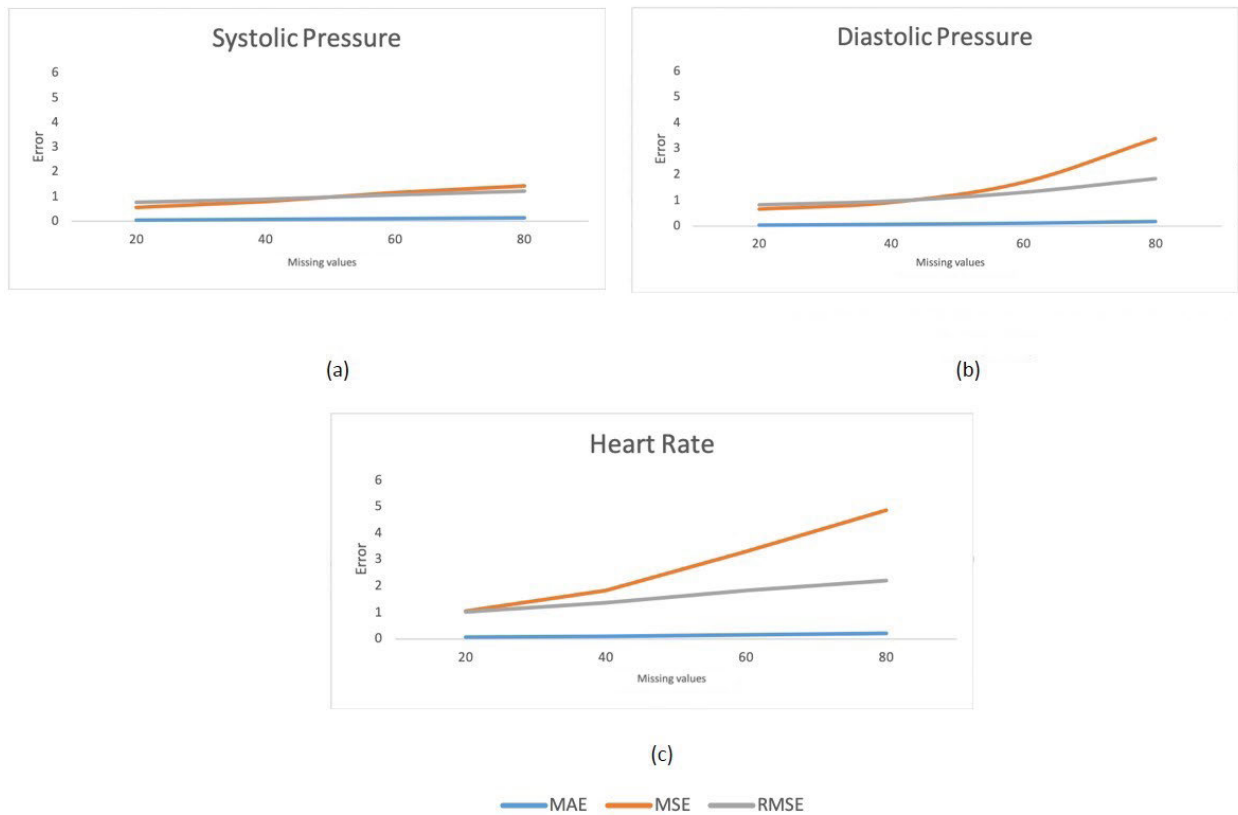


FIGURE 8. Analysis result of prediction error when missing values are introduced in (a) Systolic, (b) Diastolic, (c) Heart rate data. The model shows minor deviation when predicting an increasing number of missing values and still observes the trend in the data.

(SVM) [16]. The data set was divided into training and testing sets. The algorithms were trained with data of 24-hours, and the prediction results for the next five consecutive days were recorded. Table 7 shows a comparison of the prediction results. The batch learning methods cannot update their models, and it was observed that the prediction accuracy does not improve over time. The proposed method is online and capable of learning new patterns from the new input, thus showing higher accuracy.

VI. DISCUSSION

Gaussian mixture models provide an adequate representation of the blood pressure time series [50]. Fig. 8 shows that the impact of missing values on the accuracy of data is low as the Gaussian mixture model can infer the missing values. Moreover, experimental results in Table 5 demonstrate that the same modelling technique performs very well in detecting outliers in the blood pressure time series. However, these techniques are used only for data processing, and further experimentation is required to validate the accuracy of these methods, which is out of the scope of this research. From Table 6 and Table 7 it is evident that the prediction accuracy of the proposed model is higher than existing online as well as batch learning (offline) prediction techniques. Commonly used methods such as ANNs require a large amount of

TABLE 6. Prediction Error Comparison with state of the art online methods .

	FOS-ELM	OR-ELM	Passive Aggressive Regressor	OIESGP
Systolic				
MAE	0.73386	0.41339	0.45065	0.00671
MSE	2.19977	1.54252	0.203621	0.00011
RMSE	1.48316	1.24198	0.45124	0.0109
Diastolic				
MAE	0.59165	0.27813	0.79348	0.01071
MSE	0.78848	0.24843	0.63258	0.00043
RMSE	0.88797	0.49843	0.79535	0.02091
Heart Rate				
MAE	0.49871	0.20071	0.57029	0.00927
MSE	0.64315	0.12736	0.32765	0.00018
RMSE	0.80197	0.35687	0.57240	0.01350

training data for the most accurate results and are often prone to overfitting. The accuracy is also dependant on hyperparameter tuning, which can be time-consuming. SVM can learn from smaller datasets, but upon updating them with larger datasets, the model size increases many folds. This is because the kernel matrix requires memory that scales with the number of data points, and the training time increases

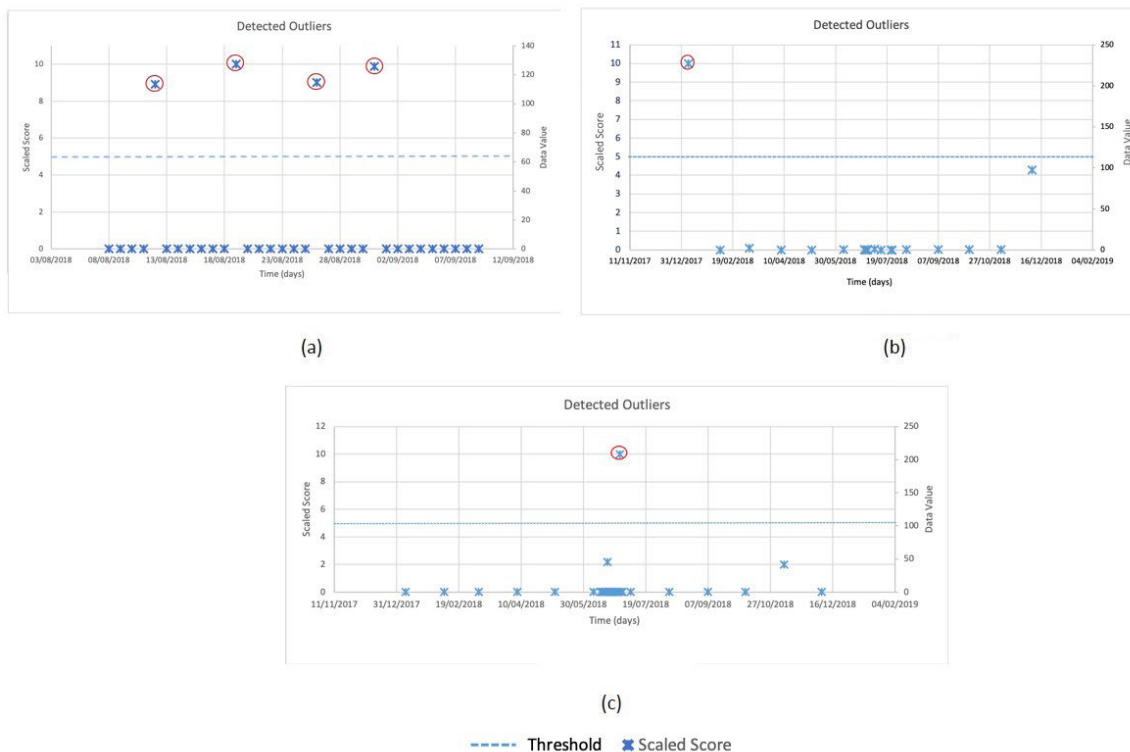


FIGURE 9. Outliers detected in (a) Systolic, (b) Diastolic, (c) Heart Rate data. Data points with higher outlier scores are identified as outliers and highlighted with red circles along their boundaries.

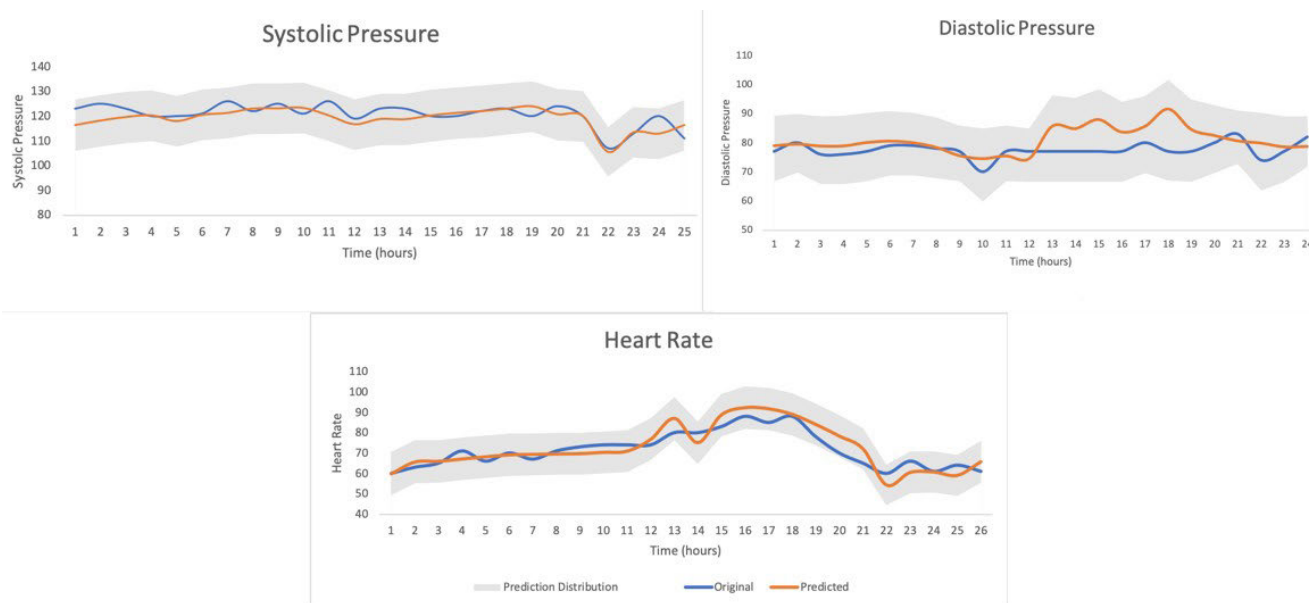


FIGURE 10. Prediction distributions using OIESGP model in (a) Systolic, (b) Diastolic (c) Heart rate. The model provides a maximum and minimum bound within which the predicted value lies.

linearly with model size. Online methods such as FOS-ELM, OR-ELM and passive-aggressive regressor update the model rather than entirely retraining; however, the prediction performance is still not par with OIESGP. However, the low prediction performance of FOS-ELM is attributed to the fact

that it cannot be used to train recurrent neural networks. Upon updating the input weights in FOS-ELM, the entire weight distribution in the hidden layer also changes, which results in lower performances and instability in results. OR-ELM shows some improvement due to its recurrent structure,

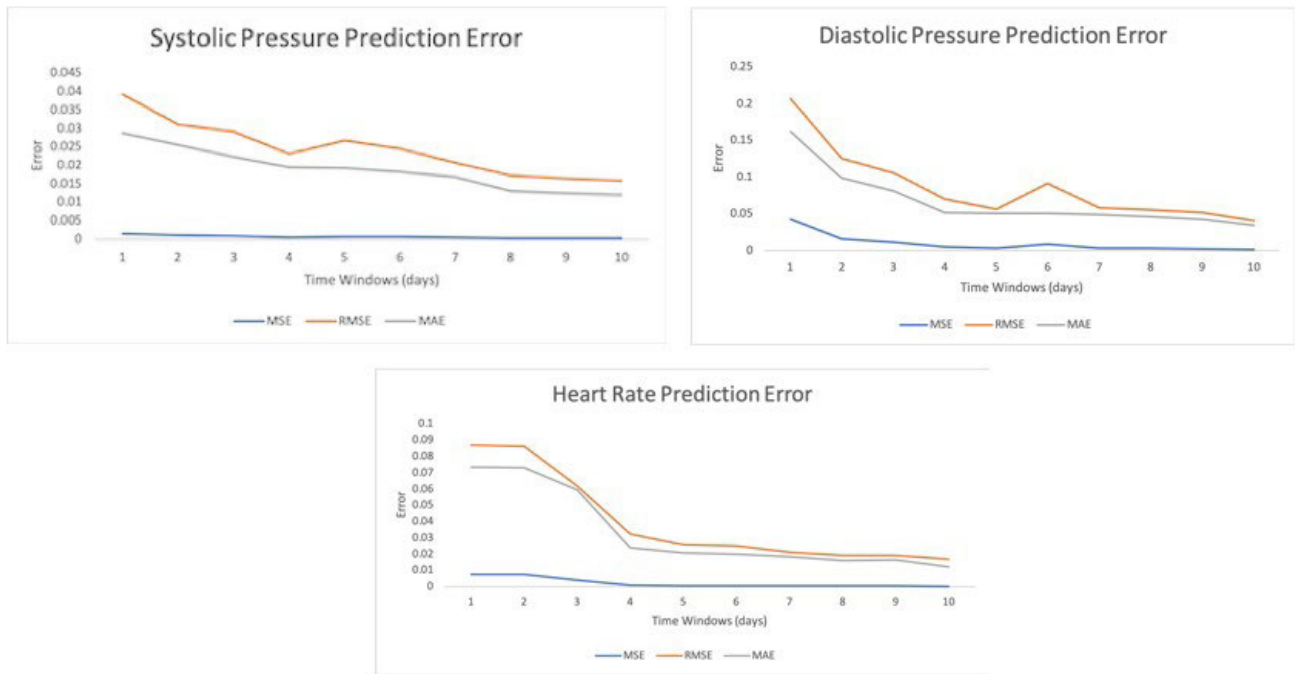


FIGURE 11. Prediction error in successive windows in (a) Systolic, (b) Diastolic, (c) Heart rate. With each successive window the model updates itself and after successive iterations, reduced prediction error is observed.

TABLE 7. Prediction Error Comparison with state of the art offline methods .

	Deep Learning (LSTM)	Gaussian Regression	SVM	OIESGP
Systolic				
MAE	0.014279	0.03364	0.01841	0.00671
MSE	0.000855	0.890052	0.00083	0.0011
RMSE	0.029239	0.943425	0.02887	0.0109
Diastolic				
MAE	0.033859	0.03313	0.04552	0.01071
MSE	0.003157	0.00400	0.00367	0.00043
RMSE	0.056188	0.06328	0.06054	0.02091
Heart Rate				
MAE	0.019802	0.01890	0.01886	0.00927
MSE	0.001279	0.00417	0.00140	0.00018
RMSE	0.035761	0.06462	0.03746	0.01350

and the model enhances as generated output is fed back to the network for computation of the subsequent output. Passive-aggressive regressor is suitable for large scale data, and the model updates only when a significant change in the input is observed, thus conserving the model size and ensuring stability. However, the most stable results are received using the proposed approach based on OIESGP. The proposed scheme enhances the data before performing any predictions, which decreases the chances of any bias or loss of information in the input data. Secondly, as the model updates with new incoming data, the online learning ability of OIESGP further improves the prediction accuracy. The reservoir structure maintains the intermediate states, which are also connected.

Thus, the generated results are much more stable. Furthermore, the reservoir only takes the most relevant states while training the model, which, given the infinite nature of biomedical time series data, ensures the model does not grow uncontrollably. The prediction algorithm takes the input data of a specific user, and the system generates a unique model for that user. In this way, each user of the system has a dedicated prediction model for their personalized blood pressure prediction.

VII. CONCLUSION

Hypertension is a complex disease with various risk factors, including individual lifestyle, genetic structure, and blood pressure history. Personalized blood pressure prediction provides an efficient means to diagnose hypertension in an individual in its early stages. However, the prediction accuracy is highly dependent on the input data, and it is of paramount importance that this data is in its purest form possible. This research work contributes to hypertension risk prediction in the Malaysian population:

- 1) Integrated multi-agent architecture for personalized hypertension risk prediction
- 2) Estimation of missing values and outliers in the input data using Gaussian mixture models,
- 3) Personalized prediction of blood pressure using online echo state Gaussian process
- 4) Calculation of a 4-year risk of hypertension using the Framingham risk calculator

In this research work, we proposed an integrated multi-agent-based system for personalized hypertension risk

prediction. A user agent is designed, which collects biomedical data from users for training and presents the predicted blood pressure results. A data processing agent receives the data from the user agent and uses Gaussian mixture regression to evaluate missing values in the historical time series. To further enhance the input data, an outlier detection mechanism is applied that employs Gaussian mixture models to identify and remove outliers in the data. The prediction agent receives the complete time series and feeds it to the OIEGSP prediction model, which produces blood pressure predictive distribution. The algorithm processes the time series by dividing it into 24-hour windows and takes a single window to predict the next window. The predictive distributions determine the range in which the blood pressure lies, which helps generate alerts for the user in dangerously high blood pressure predictions. The model is updated as new incoming data is recorded, which further improves the prediction accuracy. The Framingham risk calculator uses the average predicted blood pressure to estimate the 4-year risk of Hypertension, and the results are presented to the user in a mobile application.

The proposed system presents itself as a unique modular system, with each module responsible for a single task and operating independently. However, some research challenges remain to be addressed. The data set used in this research has been collected from the Malaysian population, and the prediction algorithm requires at least one window of data, i.e., 24 values, to be present before performing any prediction. This limitation can be overcome by generating a pool-data prediction scheme that identifies users based on the physical profile and uses an initial blood pressure profile from the collected data. The prediction results from this initial step can serve as a baseline, and the algorithm will further improve itself as the system receives actual values. The research can be further extended to include other chronic diseases such as diabetes, which can further improve the long term hypertension risk prediction.

REFERENCES

- [1] M. A. Omar, N. I. Irfan, K. Y. Yi, N. Muksan, N. L. A. Majid, and M. F. M. Yusoff, "Prevalence of young adult hypertension in Malaysia and its associated factors: Findings from national health and morbidity survey 2011," *Malaysian J. Public Health Med.*, vol. 16, no. 3, pp. 274–283, 2016.
- [2] G. Ehret, P. Munroe, K. Rice, M. Bochud, A. Johnson, D. Chasman, A. Smith, M. Tobin, G. Verwoert, S.-J. Hwang, V. Pihur, P. Vollenweider, P. O'Reilly, N. Amin, J. Bragg-Gresham, A. Teumer, N. Glazer, L. Launer, J. H. Zhao, and T. Johnson, "Genetic variants in novel pathways influence blood pressure and cardiovascular disease risk," *Nature*, vol. 478, pp. 9–103, Jan. 2011.
- [3] P. E. Marik and J. Varon, "Hypertensive crises: Challenges and management," *Chest*, vol. 131, no. 6, pp. 1949–1962, 2007.
- [4] P. M. Kearney, M. Whelton, K. Reynolds, P. Muntner, P. K. Whelton, and J. He, "Global burden of hypertension: Analysis of worldwide data," *Lancet*, vol. 365, no. 9455, pp. 217–223, Jan. 2005.
- [5] *National Health and Morbidity Survey 2011 (NHMS 2011)*, Inst. Public Health, Kuala Lumpur, Malaysia, 2015.
- [6] M. Nor, N. Safiza, G. L. Khor, S. Shahar, K. Cheong, J. Haniff, G. Appannah, R. Rasat, N. Wong, and A. A. Zainuddin, "The third national health and morbidity survey (NHMS III) 2006: Nutritional status of adults aged 18 years and above," *Malaysian J. Nutrition*, vol. 14, pp. 1–87, Jan. 2008.
- [7] D. Mozaffarian, D. Mozaffarian, E. J. Benjamin, A. S. Go, D. K. Arnett, M. J. Blaha, M. Cushman, S. R. Das, S. de Ferranti, J. P. Després, and H. J. Fullerton, "Executive summary: Heart disease and stroke statistics—2016 update: A report from the American heart association," *Circulation*, vol. 133, no. 4, pp. 447–454, Jan. 2016.
- [8] D. C. Cireşan, A. Giusti, L. M. Gambardella, and J. Schmidhuber, *Mitosis Detection in Breast Cancer Histology Images With Deep Neural Networks*. Berlin, Germany: Springer, 2013.
- [9] F. Amato, A. López, E. M. Pe na-Méndez, P. Vañhara, A. Hampl, and J. Havel, "Artificial neural networks in medical diagnosis," *J. Appl. Biomed.*, vol. 11, no. 2, pp. 47–58, 2013.
- [10] A. Narin, Y. Isler, M. Ozer, and M. Perc, "Early prediction of paroxysmal atrial fibrillation based on short-term heart rate variability," *Phys. A, Stat. Mech. Appl.*, vol. 509, pp. 56–65, Nov. 2018.
- [11] O. Erkaymaz, M. Ozer, and M. Perc, "Performance of small-world feedforward neural networks for the diagnosis of diabetes," *Appl. Math. Comput.*, vol. 311, pp. 22–28, Oct. 2017.
- [12] Y. Isler, A. Narin, M. Ozer, and M. Perc, "Multi-stage classification of congestive heart failure based on short-term heart rate variability," *Chaos, Solitons Fractals*, vol. 118, pp. 145–151, Jan. 2019.
- [13] T. Takeda, H. Nakajima, N. Tsuchiya, and Y. Hata, "A fuzzy human model for blood pressure estimation," in *Advanced Intelligent Systems*. Cham, Switzerland: Springer, 2014, pp. 109–124.
- [14] X. Li, S. Wu, and L. Wang, *Blood Pressure Prediction Via Recurrent Models With Contextual Layer*. Perth, WA, Australia: ACM, 2017.
- [15] D. LaFreniere, F. Zulkernine, D. Barber, and K. Martin, *Using Machine Learning to Predict Hypertension From a Clinical Dataset*. Athens, Greece: IEEE, 2016.
- [16] M. Ambika, G. Raghuraman, and L. SaiRamesh, "Enhanced decision support system to predict and prevent hypertension using computational intelligence techniques," *Soft Comput.*, vol. 24, pp. 13293–13304, Feb. 2020.
- [17] R. Mohammadi, S. Jain, R. Palacholla, S. Kamarthi, and B. Wallace, "Learning to identify patients at risk of uncontrolled hypertension using electronic health records data," *AMIA Joint Summits Transl. Sci. Proc. AMIA Joint Summits Transl. Sci.*, vol. 2019, pp. 533–542, May 2019.
- [18] H. Kanegae, K. Suzuki, K. Fukatani, T. Ito, N. Harada, and K. Kario, "Highly precise risk prediction model for new-onset hypertension using artificial intelligence techniques," *J. Clin. Hypertension*, vol. 22, no. 3, pp. 445–450, Mar. 2020.
- [19] P. Melin, I. Miramontes, and G. Prado-Arechiga, "A hybrid model based on modular neural networks and fuzzy systems for classification of blood pressure and hypertension risk diagnosis," *Expert Syst. Appl.*, vol. 107, pp. 146–164, Oct. 2018, doi: [10.1016/j.eswa.2018.04.023](https://doi.org/10.1016/j.eswa.2018.04.023).
- [20] C. Ye, T. Fu, S. Hao, Y. Zhang, O. Wang, B. Jin, M. Xia, M. Liu, X. Zhou, Q. Wu, Y. Guo, C. Zhu, Y.-M. Li, D. S. Culver, S. T. Alfrede, F. Stearns, K. G. Sylvester, E. Widen, D. McElhinney, and X. Ling, "Prediction of incident hypertension within the next year: Prospective study using statewide electronic health records and machine learning," *J. Med. Internet Res.*, vol. 20, no. 1, p. e22, Jan. 2018, doi: [10.2196/jmir.9268](https://doi.org/10.2196/jmir.9268).
- [21] A. Wang, N. An, G. Chen, L. Li, and G. Alterovitz, "Predicting hypertension without measurement: A non-invasive, questionnaire-based approach," *Expert Syst. Appl.*, vol. 42, no. 21, pp. 7601–7609, Nov. 2015.
- [22] C. Pater, "The blood pressure 'uncertainty range'—A pragmatic approach to overcome current diagnostic uncertainties (II)," *Current Controlled Trials Cardiovascular Med.*, vol. 6, no. 1, pp. 1–16, Apr. 2005, doi: [10.1186/1468-6708-6-5](https://doi.org/10.1186/1468-6708-6-5).
- [23] (Jan. 2019). *Summarising and Presenting Data*. [Online]. Available: <https://surfstat.anu.edu.au/surfstat-home/1-3-1.html>
- [24] H. G. Sung, "Gaussian mixture regression and classification," Rice Univ., Houston, TX, USA, Tech. Rep., 2004.
- [25] A. Reddy, M. Ordway-West, M. Lee, M. Dugan, J. Whitney, R. Kahana, B. Ford, J. Muedsam, A. Henslee, and M. Rao, "Using Gaussian mixture models to detect outliers in seasonal univariate network traffic," in *Proc. IEEE Secur. Privacy Workshops (SPW)*, May 2017, pp. 229–234, doi: [10.1109/SPW.2017.9](https://doi.org/10.1109/SPW.2017.9).
- [26] J. Chorowski, J. Wang, and J. M. Zurada, "Review and performance comparison of SVM- and ELM-based classifiers," *Neurocomputing*, vol. 128, pp. 507–516, Mar. 2014.
- [27] A. R. T. Donders, G. J. Van Der Heijden, T. Stijnen, and K. G. Moons, "A gentle introduction to imputation of missing values," *J. Clin. Epidemiology*, vol. 59, no. 10, pp. 1087–1091, 2006.
- [28] I. Pratama, A. E. Permanasari, I. Ardiyanto, and R. Indrayani, "A review of missing values handling methods on time-series data," in *Proc. Int. Conf. Inf. Technol. Syst. Innov. (ICITSI)*, Oct. 2016, pp. 1–6, doi: [10.1109/ICITSI.2016.7858189](https://doi.org/10.1109/ICITSI.2016.7858189).

- [29] J. Paalasmaa, D. Murphy, and O. Holmqvist, "Analysis of noisy biosignals for musical performance," in *Advances in Intelligent Data Analysis XI*, vol. 7619. Berlin, Germany: Springer, Dec. 2012, pp. 241–252, doi: 10.1007/978-3-642-34156-4_23.
- [30] A. T. Bevan, "Direct arterial pressure recording in unrestricted man," *Clin Sci*, vol. 36, pp. 329–344, Apr. 1969.
- [31] D. Athanassiadis, "Variability of automatic blood pressure measurements over 24-hour periods," *Clin. Sci.*, vol. 36, pp. 147–156, Feb. 1969.
- [32] S. Mann, M. M. Craig, D. I. Melville, V. Balasubramanian, and E. B. Raftery, "Physical activity and the circadian rhythm of blood pressure," *Clin. Sci.*, London, U.K., Tech. Rep., 1979.
- [33] H. Jaeger, "'The 'echo state' approach to analysing and training recurrent neural networks-with an erratum note,'" German Nat. Res. Center Inf. Technol. GMD, Tech. Rep., Jan. 2001, vol. 148, no. 34, p.13
- [34] H. Soh and Y. Demiris, "Iterative temporal learning and prediction with the sparse online echo state Gaussian process," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, Jun. 2012, pp. 1–8, doi: 10.1109/IJCNN.2012.6252504.
- [35] L. Csató and M. Opper, "Sparse on-line Gaussian processes," *Neural Comput.*, vol. 14, pp. 641–668, Apr. 2002, doi: 10.1162/089976602317250933.
- [36] *Framingham Risk Calculator*. Accessed: May 1, 2020. [Online]. Available: <https://www.framinghamheartstudy.org/fhs-risk-functions/hypertension/>
- [37] W. Kannel, "Risk stratification in hypertension: New insights from the framingham study*1," *Amer. J. Hypertension*, vol. 13, no. 1, pp. S3–S10, Jan. 2000.
- [38] D. R. Cox and D. Oakes, *Analysis of Survival Data*, vol. 21. Boca Raton, FL, USA: CRC Press, 1984.
- [39] B. M. Egan, S. E. Kjeldsen, G. Grassi, M. Esler, and G. Mancia, "The global burden of hypertension exceeds 1.4 billion people: Should a systolic blood pressure target below 130 become the universal standard?" *J. Hypertension*, vol. 37, no. 6, pp. 1148–1153, 2019.
- [40] S. T. Turner, G. L. Schwartz, and E. Boerwinkle, "Personalized medicine for high blood pressure," *J. Amer. Heart Assoc.*, vol. 50, no. 1, 2007.
- [41] C. Krittanawong, A. S. Bomback, U. Baber, S. Bangalore, F. H. Messerli, and W. H. W. Tang, "Future direction for using artificial intelligence to predict and manage hypertension," *Current Hypertension Rep.*, vol. 20, no. 9, p. 75, Sep. 2018.
- [42] J. Allen, "Photoplethysmography and its application in clinical physiological measurement," *Physiol. Meas.*, vol. 28, no. 3, pp. R1–R39, Mar. 2007.
- [43] C. Tofallis, "A better measure of relative prediction accuracy for model selection and model estimation," *J. Oper. Res. Soc.*, vol. 66, no. 8, pp. 1352–1362, Aug. 2015.
- [44] R. G. Pontius, O. Thontteh, and H. Chen, "Components of information for multiple resolution comparison between maps that share a real variable," *Environ. Ecol. Statist.*, vol. 15, no. 2, pp. 111–142, Jun. 2008.
- [45] R. J. Hyndman and A. B. Koehler, "Another look at measures of forecast accuracy," *Int. J. Forecasting*, vol. 22, no. 4, pp. 679–688, Oct. 2006.
- [46] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer, "Online passive-aggressive algorithms," *J. Mach. Learn. Res.*, vol. 7, pp. 551–585, Dec. 2006.
- [47] J.-M. Park and J.-H. Kim, "Online recurrent extreme learning machine and its application to time-series prediction," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, May 2017, pp. 1983–1990.
- [48] G.-B. Huang, N.-Y. Liang, H.-J. Rong, P. Saratchandran, and N. Sundararajan, "On-line sequential extreme learning machine," *Comput. Intell.*, vol. 2005, pp. 232–237, Jul. 2005.
- [49] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [50] L.-F. Cheng, G. Darnell, B. Dumitrascu, C. Chivers, M. E. Draugelis, K. Li, and B. E. Engelhardt, "Sparse multi-output Gaussian processes for medical time series prediction," 2017, *arXiv:1703.09112*. [Online]. Available: <http://arxiv.org/abs/1703.09112>
- [51] C. Rasmussen and C. Williams, "Gaussian processes for machine learning," in *Adaptive Computation and Machine Learning*, vol. 38. Cambridge, MA, USA: MIT Press, 2006, pp. 715–719.



SUNDUS ABRAR received the B.S. degree (Hons.) in computer systems engineering from the National University of Science and Technology, in 2011. She is currently pursuing the M.S. degree in computer science with specialisation in applied computing with the University of Malaya. She is working at the Smart Robotics Lab, University of Malaya. Her research interests include time series prediction and health data analytics.



CHU KIONG LOO (Senior Member, IEEE) received the B.Eng. degree (Hons.) in mechanical engineering from the University of Malaya, in 1996, and the Ph.D. degree specializing in neurorobotics from Universiti Sains Malaysia, in 2004. He is currently a Full Professor with the Department of Artificial Intelligence, Faculty of Computer Science and Information Technology, University of Malaya. His main research interest includes neuroscience inspired machine intelligence. He was the IEEE Systems, Man, and Cybernetics SMC Society Vice-Chairman for Malaysia Chapter, from 2013 to 2014. He was the President of the Asia Pacific Neural Network Assembly (APNNA), in 2014.



NAOYUKI KUBOTA (Member, IEEE) received the B.Sc. degree from Osaka Kyoiku University, Kashiwara, Japan, in 1992, the M.Eng. degree from Hokkaido University, Japan, in 1994, and the D.E. degree from Nagoya University, Nagoya, Japan, in 1997. He joined the Osaka Institute of Technology, Osaka, Japan, in 1997. In 2000, he joined the Department of Human and Artificial Intelligence Systems, University of Fukui, Fukui, Japan, as an Associate Professor. He joined the Department of Mechanical Engineering, Tokyo Metropolitan University, Tokyo, in 2004. From 2005 to 2012, he was an Associate Professor with the Department of System Design, Tokyo Metropolitan University, where he has been a Professor, since 2012.

• • •