# Improving Accuracy of Patient Synthetic Data for Testing Medical Cyber-Physical Systems

Leonardo C. Santos, Lenardo C. Silva, Ana Luisa Medeiros, Hyggo Almeida and Angelo Perkusich

Federal University of Campina Grande, Campina Grande, Brazil
{leonardo, lenardo.silva, ana.medeiros, hyggo, perkusic}@embedded.ufcg.edu.br

## Abstract

*Medical Cyber-Physical Systems (MCPS) integrate the cyber space and physical world elements for promoting support for health assurance activities. MCPS are life-critical systems, demanding a strong engineering effort to guarantee safety, what directly impacts on testing process. Testing MCPS using real patients is very expensive and complex, since their lives are involved. Thus, the use of patient synthetic data becomes a promising approach. In this paper we propose a model for improving accuracy of patient synthetic data for testing MCPS based on regression models. We use an existing Patient Baseline Model to generate vital signs of patients, but improving the statistical analysis. Using our approach we increased in about 73.9% the quality of the regression models and, consequently, their accuracies.*

***Medical Cyber-Physical Systems; Statistical Analysis; Simulation; Testing; Patient Baseline Model (key words)***

## 1. INTRODUCTION

The use of computing resources is increasing daily in personal and corporate environments. Since virtual entities directly react to stimuli produced by physical entities, these elements end up becoming a huge source of information for its users. This scenario, in which embedded computing units are in constant interaction with real world elements to monitor and control physical processes, forms the so-called *Cyber-Physical Systems* (CPS) [1].

CPS applied to health are commonly called *Medical Cyber-Physical Systems* (MCPS) [2]. In this sense, designing such systems has become an increasingly complex task due to the need to ensure the patient safety at runtime. This guarantee can be achieved through system verification and validation, what requires high abstraction level, realistic simulations and relevant tests.

Several models have been adopted as a way of representing the physical and cyber elements in the health field. Hotehama et al. [3] presented a cardiovascular model to predict blood pressure and heart rate during physical exercises. Van Heusden et al. [4] proposed an artificial pancreas model for patients with type 1 diabetes mellitus, with the goal of improving glucose control in such patients. Wu et al. [5] and Bhaduri et al. [6] used artificial neural networks to investigate the correlation between blood pressure and some variables such as alcohol consumption, body mass index (BMI), age, and exercise, thereby building patient models to represent specific behaviors of the human body. Finally, Khan et al. [7] provided a glucose control system to be used to prevent hypoglycemic episodes, in which the patient model (i.e., the artificial pancreas) establishes a relationship between heart rate and blood glucose level.

Although there are several related works to develop MCPS, testing MCPS is still a challenge. Using real patients is very expensive and complex, since their lives are involved. Thus, the use of patient synthetic data becomes a promising approach. In this context, Silva et al.[8] presented a model-based architecture to support testing of MCPS. The authors introduced a Patient Baseline Model that uses regression models to generate synthetic data for the heart rate (HR), respiratory rate (RR), blood pressure (BP), and body temperature (BT) vital signs. In addition to proposing a new model, Silva et al. discussed the patient models proposed in other works, thus proving that the use of statistical data in order to obtain knowledge on human behavior is a common - but nontrivial - method. However, the potential predictor variables selected for the statistical analysis, as well as the use of the regression models for the vital variables, generated insignificant statistical results. In this case, the regression models inherent to the heart and respiratory rates, systolic blood pressure and body temperature represented only 48.9%, 31%, 51.1% and 48%, respectively, of the variability of the data contained in the samples selected for analysis.

Two main issues must be considered in the Patient Baseline Model: (i) predictor variables are not sufficient to explain the vital variables because the linear correlation among them is weak; (ii) the sample selected to perform the statistical analysis is heterogeneous, since the records

contained in this sample were collected in a time in which the individuals were admitted to intensive care units. This means that the patients were presenting the most varied critical health conditions.

In this research, we investigate the above mentioned issues to improve the quality of the prediction and accuracy of the regression models that compose the Patient Baseline Model. Therefore, in order to answer the following research questions, we declare their respective null hypothesis:

**Q1:** Can the regression models proposed by Silva et al. [8] be improved by modifying the predictor variables?

> **H1-0:** There is no way to improve the regression models proposed by Silva et al. [8] by modifying the predictor variables.

**Q2:** Can the regression models proposed by Silva et al. [8] be improved by selecting a homogeneous sample?

> **H2-0:** There is no way to improve the regression models proposed by Silva et al. [8] with the selection of another homogeneous sample from the database.

To investigate such issues, we present an experiment that was divided into the following three steps:

**Step 1:** Investigate the literature in order to identify the potential variables that are strongly correlated with each vital variable of interest (i.e., HR, RR, BP and BT);

**Step 2:** Select a sample from a database containing patient records in intermediate treatment periods in order to obtain better quality indicators for the regression models;

**Step 3:** Perform statistical analysis in order to obtain regression models that can better explain the vital variables of interest.

The remainder of this paper is organized as follows. Section 2 describes the characterization of the population of interest, as well as the statistical analysis performed to obtain the regression models for the vital variables. In Section 4, we discuss the threats to the experiment's validity. Finally, in Section 5, we expose the final considerations.

## 2. MATERIALS AND METHODS

In the literature review, we identified some works that define a set of potential predictor variables for each vital sign. These variables are presented in Table 1. With the possibility of a multicollinearity problem in the construction of the regression models, only the systolic blood pressure was used, as safeguarded by Gavish [9], ignoring the diastolic blood pressure, as they have a strong correlation.

We used the same database used by Silva et al. [8] to collect patient's records containing the larger number of potential predictor variables identified for obtaining the new regression models for vital signs. This database, so-called MIMIC II Clinical Database [13], is made available freely by the American service PhysioNet. The information in this database refers to patients admitted to Intensive Care Units (ICU) whose data were collected from bedside monitors and

Table 1: Potential predictor variables for each vital sign of interest, grouped by related work.

| Vital sign | Predictor variables |
|---|---|
| BP [6] | Environment temperature, age, gender, body mass index (BMI), alcohol consumption, smoking, cholesterol and blood glucose. |
| BP [5] | Alcohol consumption, age and exercises. |
| BP and HR [3] | Weight, age, blood pressure at rest, heart rate at rest, exercise intensity, type of exercise and oxygen consumption. |
| All [10] | Age, gender, exercises, pregnancy, emotional state, hormones, medications, fever and hemorrhage. |
| BP | Age, exercises, stress, medications and diseases. |
| HR, RR and BT [11] | Age, exercises, stress, environment temperature, medications and diseases. |
| BP | Age, gender, environment temperature, emotional state, exercises, body position, medications, pain, recent meal, caffeine, smoking and bladder distention. |
| HR | Age, gender, exercises, emotional state, metabolism, fever and medications. |
| RR | Age, exercises, emotional state, fever and medications. |
| BT [12] | Age, environment temperature, emotional state, environment, exercises, patient's normal body temperature and pregnancy. |

hospital files. In addition to general patient data, other information can be found such as patient's conditions at time of admission, vital signs and physiological parameters, drug administration, laboratory tests, and other information described in the doctor's report.

From the analysis of the patients' records found in the MIMIC II database, we identified and collected some of the variables presented in Table 1, such as age, gender, alcohol consumption, cholesterol, blood glucose, weight, height, oxygen consumption, medications and diseases, and the vital variables of respiratory and heart rate, blood pressure, and body temperature. These variables served as basis for the establishment of the linear regression models that compose the Patient Baseline Model.

### 2.1. Definition of the Population of Interest

Due to the large amount of records found in the MIMIC II database, as well as the possibility that some of the records contained errors or were duplicated, we identified the need of defining a population of patients for the study and a second sample for the validation of the proposed models. In the process of defining the population of interest, shown in Figure 1, we determined a set of rules to be ap-

plied, wherein the first rule was the "Specification of the population". In this first rule, we selected the records of patients over the age of 15 years (because they have more stable vital signs), whose respiratory and heart rates, blood pressure, body temperature, glucose and $CO_2$ consumption were measured simultaneously.
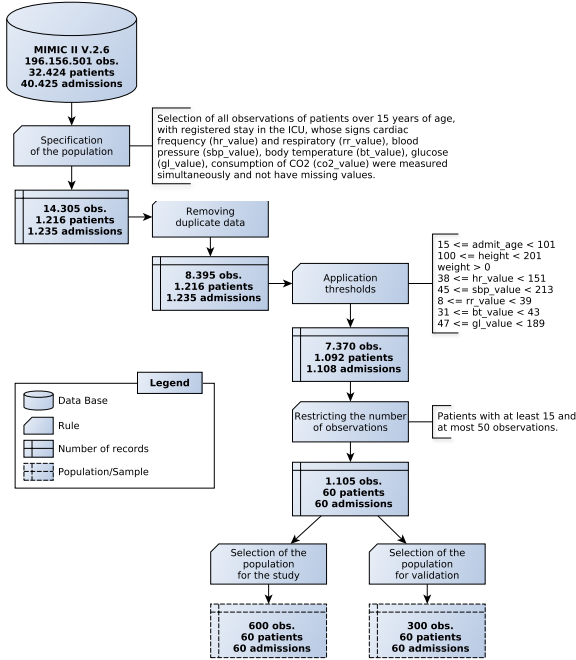


Figure 1: Defining the population of interest.

After the "Specification of the population", we removed all of the duplicate data ("Removing duplicate data") and then, in the "Application thresholds" rule, we defined the minimum and maximum values for the non-categorical variables with the purpose of removing data inconsistencies. As a reference for thresholds definition, we used a set of Clinical Guidelines described in [14, 15, 16, 17, 18]. With the "Restricting the number of observations" rule, we removed patient's records who remained for a short (15 observations or less) or long time (50 observations or more) in the ICU. This allow us to select more stable patients.

The latest rules ("Selection of the population for the study" and "For validation") are specifically related to the selection of two populations of interest for the study. At first, we selected randomly two records for each patient, resulting in 600 observations relating to 60 patients. The second population was defined to validate the regression models that were obtained from the first population. Thus, selecting randomly only one record for each patient results in 300 observations.

### 2.2. Statistical Analysis

This section describes the process of obtaining the regression models for the vital variables of interest which,

when interrelated, will provide the basic behavior of the Patient Baseline Model. The main statistical metric adopted in the evaluation process of the regression models was the square of the linear correlation coefficient between the answer variable ($\widehat{y}$) and the adjusted values ($\widehat{\mu}$), given by (1).

$$R^{2*} = cor(\widehat{y}, \widehat{\mu})^2 \qquad (1)$$

In order to obtain the regression models, we used the generalized linear models (GLM) class. The method used to adjust the regression models was the Backwards Elimination [19] method. The reason for such choices is related to the large number of variables to be analyzed. In this method, the first linear regression is obtained with all of the potential predictor variables for each regression model, and then the variables are disregarded one by one according to the *p-value*[1] calculated by the t-test for significance testing.

In order to validate the obtained regression models, we used the following methods: (i) verifying the normality of the errors through the Shapiro-Wilk [20] normality test, which allows us to check if the model used is suitable for the data; (ii) comparison between the data obtained from the test sample and data generated by each vital sign regression model for this sample. In this comparison we performed the visual analysis of line graphs and the t-test for significance testing.

## 3. RESULTS

### 3.1. Regression Model for Respiratory Rate

As a starting point, we considered all of the variables present in the population of interest and the interactions between the variables $admit\_age$, $height$, $weight$, $rr\_value$, $hr\_value$, $sbp\_value$, $bt\_value$, $gl\_value$ e $co2\_value$, $sex$, due to the possibility that the interactions between them are significant.

The linear regression model used was the Normal Inverse given by (2) with canonical link function defined in (3). This is a particular case of the MLGs class [19] and was chosen since it best represents the sample. It is noteworthy that other models were tested.

$$\eta = \frac{1}{\mu^2} \Leftrightarrow \mu = \eta^{-\frac{1}{2}} \qquad (2)$$

Where $\mu$ is the average of the respiratory rate ($rr\_value$) variable, which we wish to model and

$$\eta = \beta_0 + \beta_1 X_1 + ... + \beta_n X_n \qquad (3)$$

is a systematic component, or linear predictor, in which $\beta_0$ represents the intercept coefficient, that is, when the value of all of the predictor variables of the model take the value 0, $\eta = \beta_0$.

---

[1]According to Diez et al. [19], *p-value* is the variable used to obtain the test statistic that is equal or even more extreme than the one observed in a sample.

After the regression model for RR ($MLG\_RR$) adjusted, the $R^{2*}$ obtained for this model was 0.711. This means that $MLG\_RR$ explains 71.1% of the data variability contained in the vital variable $rr\_value$. In practical terms, the $MLG\_RR$ representation, which represents the estimated average respiratory rate is given by (4)

$$\hat{\eta} = \beta_0 + \sum_{i=1}^{22} \beta_i X_i - \beta_{23} X_{11} X_{12} - \beta_{24} X_{12} X_{18}$$
$$-\beta_{25} X_{19} X_{20} - \beta_{26} X_{11} X_{19} - \beta_{27} X_{15} X_{21} + \beta_{28} X_{15} X_{21}$$
$$-\beta_{29} X_{15} X_{19} + \beta_{30} X_{20} X_{22} + \beta_{31} X_{12} X_{21} - \beta_{32} X_{21} X_{22}$$
$$+\beta_{33} X_{19} X_{22} + \beta_{34} X_{15} X_{12} - \beta_{35} X_{19} X_{21} \qquad (4)$$

where $\beta_0$ is the intercept, $\beta_{1-35}$ are the coefficients inherent to each predictor variable, and $X_{1-22}$, a subset of the variables present in the population of interest. For more details see https://github.com/leonardocsantoss/patient-baseline-model.

Regarding the Shapiro-Wilk normality test used for determining the residuals present in the model, the *p-value* calculated was 0.995. This way, we can say with 95% confidence that the residuals present in $MLG\_RR$ have normal distribution and the model fits the data contained in the sample.
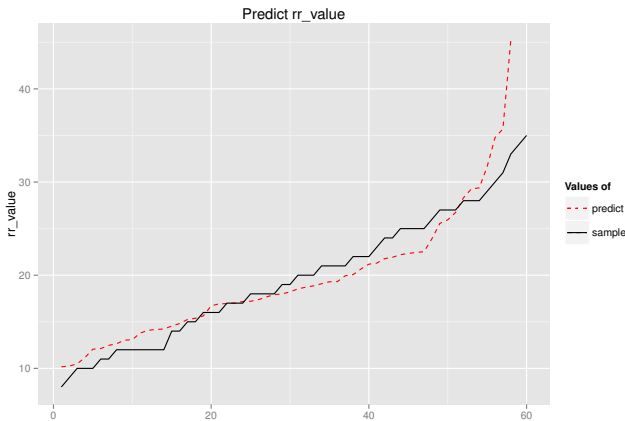


Figure 2: Comparison between the real data and the synthetic data calculated by $MLG\_RR$ for the $rr\_value$ variable.

In addition, to measure the accuracy of the $MLG\_RR$, we compared the values of the $rr\_value$ variable in the test sample (real data), with the values calculated by the regression model for the same sample (synthetic data). In the comparative graph shown in Figure 2, it is possible to verify that the values calculated by the model (dashed line) are close to the real data (solid line).

The statistical evidence that these two data sets are equal is shown by the result of the t-test, in which the hypothesis are $H0 : \beta_i = \beta_j$ and $H1 : \beta_i \neq \beta_j$. Thus, to refute

$H0$ (null hypothesis) implies that the two data sets are different. The result of the t-test for $MLG\_RR$ calculated a $p - value = 0.5012$. Thus, with 95% confidence, the null hypothesis was refuted, which leads us to conclude that statistically the two data sets are equal.

Once the process to obtain and fitting the regression models for vital signs was presented, we present only the results of the regression models for *hr*, *sbp*, and *bt* vital variables. Thus, we omitted the equations related to the linear predictor of these regression models.

### 3.2. Regression Model for Heart Rate

The regression model chosen for the heart rate variable ($hr\_value$) was the Gama Linear Model with canonical link function given by (5). This regression model also belongs to the MLGs class and, when related to heart rate, is what best represents the variability of the data found in the sample. It is noteworthy that other models were tested.

$$\eta = \frac{1}{\mu} \Leftrightarrow \mu = \eta^{-1} \qquad (5)$$

After adjustment of the regression model for HR ($MLG\_HR$), we obtained the $R^{2*} = 0.822$.
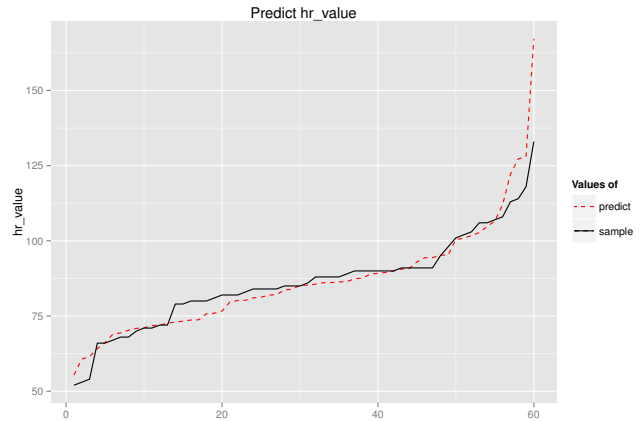


Figure 3: Comparison between the real data and the synthetic data calculated by $MLG\_HR$ for the $hr\_value$ variable.

The Shapiro-Wilk normality test used to verify the residuals in $MLG\_HR$ calculated a $p - value = 0.9918$. Thus, with 95% confidence, the residuals in these models also have normal distribution. Therefore, the $MLG\_HR$ is appropriate to the real sample data. Comparing the values calculated for this regression model with the variable $hr\_value$ values of the test sample (see Figure 3) using t-test for significance testing, we obtained a $p - value = 0.7036$. Statistically speaking, with 95% confidence, both data sets are equal.

### 3.3. Regression model for Systolic Blood Pressure

For the systolic blood pressure ($sbp\_value$) variable, the Gama regression model was also used, whose canonical link

function was previously shown in (5). Other models have been tested, however, this was best represented the variability of the data. After adjusting the regression model for SBP ($MLG\_SBP$), we obtained the $R^{2*} = 0.825$.

The Shapiro-Wilk normality test used to verify the residuals in $MLG\_SBP$, calculated a $p - value = 0.06131$. Thus, with 95% confidence, the residuals in these models also have normal distribution. Therefore, the $MLG\_SBP$ fits to the real sample data. Comparing the calculated values for this regression model with the values of the $sbp\_value$ variable of the test sample (see Figure 4) using the t-test for significance testing, we obtained a $p - value = 0.3837$. Thus, with 95% confidence, the two data sets are statistically equal.
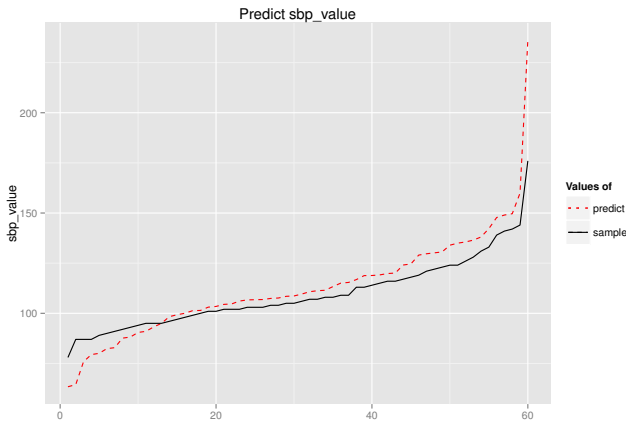


Figure 4: Comparison between the real data and the synthetic data calculated by $MLG\_SBP$ for the $sbp\_value$ variable.

### 3.4. Regression Model for Body Temperature

Finally, for the body temperature ($bt\_value$) variable, we used the Gama regression model, whose canonical link function was presented previously in (5). After adjusting the regression model for BT ($MLG\_BT$), we obtained the $R^{2*} = 0.755$.

The Shapiro-Wilk normality test used to verify the residuals in $MLG\_BT$ calculated a $p - value = 0.4675$. Thus, with 95% confidence, the residuals in these models also have a normal distribution. Therefore, the $MLG\_BT$ is suitable to the real sample data. Comparing the calculated values for this regression model with the values of the $bt\_value$ variable of the test sample (see Figure 5) using the t-test for significance testing, we obtained a $p - value = 0.3914$. Statistically speaking, with 95% confidence, the two data sets are equal.

## 4. DISCUSSION

In order to discuss the results obtained in this work, it is necessary to resume the two research questions defined in Section 1. According to the results obtained in the hypoth-
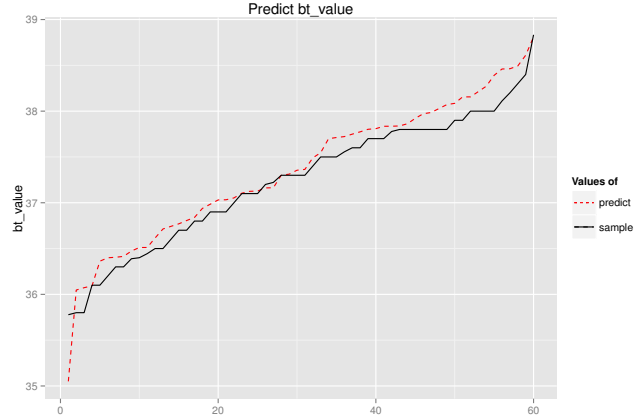


Figure 5: Comparison between the real data and the synthetic data calculated by $MLG\_BT$ for the $bt\_value$ variable.

esis tests related to their research questions, it was possible to refute the null hypothesis **H1-0** and **H2-0**, as shown in Table 2. This means that the use of a homogeneous sample from the MIMIC II clinical database, along with a set of predictor variables different from those proposed by Silva et al. [8], allowed us to obtain better regression models for the vital variables that make up the Patient Baseline Model. This was evidenced statistically through the values of the calculated $R^{2*}$ metric for each regression model, as well as through visual analysis of real data compared to synthetic data generated by these models.

Table 2: Summary of the Results for the MLGs.

| MLG | $R^{2*}$ | Shapiro-Wilk test (*p-value*) | t-test (*p-value*) |
|---|---|---|---|
| MLG_RR | 0.711 | 0.9950 | 0.5012 |
| MLG_HR | 0.822 | 0.9918 | 0.7036 |
| MLG_SBP | 0.825 | 0.0613 | 0.3837 |
| MLG_BT | 0.755 | 0.4675 | 0.3914 |

Regarding to threats to the validity of this study, we can take into account the following issue:

**Restricted application domain:** the used clinical database contains only data from intensive care units related to patients in critical health condition. This can interfere in the adjustment of regression models and the accuracy of the prediction of their vital signs. Therefore, this feature of the data sample restricts the application scope of the Patient Baseline Model;

**Effectiveness of the Accuracy:** some of the predictor variables shown in Table 1, such as exercise, environment temperature, stress, emotional state, pain, and so on, were not found in the clinical database used for this statistical analysis. Thus, these variables were excluded from the analysis, which can negatively impact in the accuracy of the regression models extracted for the vital signs considered in

the Patient Baseline Model. Furthermore, the number of variables used to generate the models was too high, which resulted in the increase of the model's complexity and, at the same time, raised its accuracy. Finally, it was not possible to determine the level of difficulty to use the model.

## 5. CONCLUSION

In this paper, we presented the Patient Model based model proposed by Silva et al. [8], which was built to serve as a basis to used in the MCPS validation and verification processes that requires interaction with real patients.

Initially, we discussed some related work that list the potential predictor variables for the respiratory and heart rate, blood pressure and body temperature vital signs. In addition, we presented a process to define the populations of interest for the study from a clinical database, including the data set used for statistical analysis and validation. Finally, we described the whole process to obtain and fit the regression models for the vital variables considered in the analysis. In the validation process of these models, we verified the normality of errors and compared the data generated by the regression models with the real values extracted from the test sample.

With the Patient Baseline Model representing the variability found in their vital variables (i.e., 71.1% of the respiratory rate, 82.2% of the heart rate, 82.5% of the systolic blood pressure, and 75.5% of the body temperature), we obtained regression models counting a gain of approximately 73%, when compared to the original model proposed by Silva et al. [8]. The results achieved in this study allow us to make more realistic simulations and generating relevant tests, increasing the confidence of such tests and minimizing the need to do clinical trials during the initial tests of a MCPS.

For future work, we intend to build a model representing the interaction between these regression models, allowing us to simulate different conditions regarding the basic health condition of a patient, characterized by the four main vital signs considered in this work.

## References

[1] M. Conti, S. K. Das, C. Bisdikian, M. Kumar, L. M. Ni, A. Passarella, G. Roussos, G. Tröster, G. Tsudik, and F. Zambonelli, "Looking ahead in pervasive computing: Challenges and opportunities in the era of cyber–physical convergence," *Pervasive and Mobile Computing*, vol. 8, no. 1, pp. 2–21, 2012.

[2] I. Lee, O. Sokolsky, S. Chen, J. Hatcliff, E. Jee, B. Kim, A. King, M. Mullen-Fortino, S. Park, A. Roederer *et al.*, "Challenges and research directions in medical cyber–physical systems," *Proceedings of the IEEE*, vol. 100, no. 1, pp. 75–90, 2012.

[3] M. Hotehama, S. T. Lindsay, and T. Takemori, "Simulation of living behavior: a cardiovascular model for predicting blood pressure and heart rate," in *SICE 2003 Annual Conference (IEEE Cat. No. 03TH8734)*, 2003.

[4] K. van Heusden, E. Dassau, H. C. Zisser, D. E. Seborg, and F. J. Doyle, "Control-relevant models for glucose control using a priori patient characteristics," *Biomedical Engineering, IEEE Transactions on*, vol. 59, no. 7, pp. 1839–1849, 2012.

[5] T. H. Wu, G. K.-H. Pang, and E. W.-Y. Kwong, "Predicting systolic blood pressure using machine learning," in *Information and Automation for Sustainability (ICIAfS), 2014 7th International Conference on*. IEEE, 2014, pp. 1–6.

[6] A. Bhaduri, A. Bhaduri, A. Bhaduri, and P. Mohapatra, "Blood pressure modeling using statistical and computational intelligence approaches," in *Advance Computing Conference, 2009. IACC 2009. IEEE International*. IEEE, 2009, pp. 1026–1030.

[7] S. H. Khan, A. H. Khan, and Z. H. Khan, "Artificial pancreas coupled vital signs monitoring for improved patient safety," *Arabian Journal for Science and Engineering*, vol. 38, no. 11, pp. 3093–3102, 2013.

[8] L. Silva, M. Perkusich, H. Almeida, A. Perkusich, M. Lima, and K. Gorgônio, "A baseline patient model to support testing of medical cyber-physical systems." *Studies in health technology and informatics*, vol. 216, pp. 549–553, 2015.

[9] B. Gavish, I. Z. Ben-Dov, and M. Bursztyn, "Linear relationship between systolic and diastolic blood pressure monitored over 24 h: assessment and correlates," *Journal of hypertension*, vol. 26, no. 2, pp. 199–209, 2008.

[10] J. V. González, O. A. V. Arenas, and V. V. González, "Vitals sign semiology: the new look to an actual problem," *Archivos de Medicina (Manizales)*, vol. 12, no. 2, pp. 221–240, 2012.

[11] B. Vaughans, *Nursing Fundamentals DeMYSTiFieD: A Self-Teaching Guide*. McGraw Hill Professional, 2010.

[12] K. Bonewit-West, S. Hunt, and E. Applegate, *Today's Medical Assistant: Clinical & Administrative Procedures*. Elsevier Health Sciences, 2014.

[13] G. D. Clifford, D. J. Scott, and M. Villarroel, "User guide and documentation for the mimic ii database," *MIMIC-II database version*, vol. 2, 2009.

[14] A. D. Association *et al.*, "Diagnosis and classification of diabetes mellitus," *Diabetes care*, vol. 33, no. Supplement 1, pp. S62–S69, 2010.

[15] A. V. Chobanian, G. L. Bakris, H. R. Black, W. C. Cushman, L. A. Green, J. L. Izzo, D. W. Jones, B. J. Materson, S. Oparil, J. T. Wright *et al.*, "Seventh report of the joint national committee on prevention, detection, evaluation, and treatment of high blood pressure," *Hypertension*, vol. 42, no. 6, pp. 1206–1252, 2003.

[16] Y. Handelsman, J. Mechanick, L. Blonde, G. Grunberger, Z. Bloomgarden, G. Bray, S. Dagogo-Jack, J. Davidson, D. Einhorn, O. Ganda *et al.*, "American association of clinical endocrinologists medical guidelines for clinical practice for developing a diabetes mellitus comprehensive care plan," *Endocrine Practice*, vol. 17, no. Supplement 2, pp. 1–53, 2011.

[17] S. McGee, *Evidence-based physical diagnosis*. Elsevier Health Sciences, 2012.

[18] U. D. of Health, H. Services *et al.*, "The fourth report on the diagnosis, evaluation, and treatment of high blood pressure in children and adolescents. 2005," 2012.

[19] D. M. Diez, C. D. Barr, and M. Cetinkaya-Rundel, *OpenIntro statistics*. CreateSpace independent publishing platform, 2012.

[20] J. Royston, "An extension of shapiro and wilk's w test for normality to large samples," *Applied Statistics*, pp. 115–124, 1982.