# An Explainable AI based Clinical Assistance Model for Identifying Patients with the Onset of Sepsis

Snehashis Chakraborty[1], Komal Kumar[1], Balakrishna Pailla Reddy[2], Tanushree Meena[1], Sudipta Roy[1,*]

[1]*Artificial Intelligence & Data Science, Jio Institute, Navi Mumbai-410206, India*
[2]*Artificial Intelligence Centre of Excellence (AICoE), Reliance Jio, Hyderabad, India*
snehashis1.C@jioinstitute.edu.in; komal2.kumar@jioinstitute.edu.in; balakrishna.pailla@ril.com; tanushree.meena@jioinstitute.edu.in; sudipta1.roy@jioinstitute.edu.in;

*Abstract*— **The high mortality rate of sepsis, especially in Intensive Care Unit (ICU) makes it third-highest mortality disease globally. The treatment of sepsis is also time consuming and depends on multi-parametric tests, hence early identification of patients with sepsis becomes crucial. The recent rise in the development of Artificial Intelligence (AI) based models, especially in early prediction of sepsis, have improved the patient outcome. However, drawbacks like low sensitivity, use of excess features that leads to overfitting, and lack of interpretability limit their ability to be used in a clinical setting. So, in this research we have developed a smart, explainable and a highly accurate AI based model (called XAutoNet) that provides quick and early prediction of sepsis with a minimal number of features as input. An application based novel convolutional neural network (CNN) based autoencoder is also implemented that improves the performance of XAutoNet by dimensional reduction. Finally, to unbox the "Black Box" nature of these models, Gradient based Class Activation Map (GradCAM) and SHapley Additive exPlanations (SHAP) are implemented to provide interpretability of autoencoder and XAutoNet in the form of visualization graphs to assist clinicians in diagnosis and treatment.**

*Keywords—Healthcare, XAI, Sepsis Prediction, Autoencoders.*

## I. INTRODUCTION

Sepsis is a serious medical condition caused by an overwhelming immune response to an infection. Nearly 50 million people get affected every year with a mortality rate of at least 11 million. In 2017, nearly half of global sepsis cases were accounted by children under five years of age with 2.9 million deaths [1]. Immediate treatment like Broad-spectrum antibiotics shows promising results in improving patient's outcome. For these treatments to be effective, they must be initiated as early as possible. To address this issue, hospitals came up with different ICU related scoring systems like Glasgow Coma Scale, Systemic Inflammatory Response Syndrome (SIRS) score, Quick Sequential Organ Failure Assessment and many more. Though these indicators provide early identification of sepsis with a high sensitivity rate, their specificity is low due to similarity in the symptoms with other diseases like cancer, pneumonia and so on, hence can give false alarms.

Over the past few years, AI has made significant progress in the medical field by assisting healthcare professionals in analysing and identifying illness due to high availability of medical data [2-6]. There are various frameworks based on Machine Learning (ML) and Deep Learning (DL) that not only utilizes structured data such as Electronic Medical Record (EMR) [7] but also unstructured data in the form of narrative notes and prescriptions written by doctors [8] to predict sepsis early. In [9], the authors proposed an ensemble ML model comprising of Gradient Boosting Machine (light GBM), Random Forest (RF) and Extreme Gradient Boosting (XGBoost) to predict sepsis 6 hours in advance. Though the performance of their model was better compared to the individual models, its overall performance was not very high, also it lacked interpretability. In [10], the author used light GBM along with a new imputation method named mixed filling to predict sepsis 6 hours in advance and used interpretable tools like Local Interpretable Model-Agnostic Explanations to provide local feature importance. However, drawbacks like lack of global explanations for the model and strong theoretical foundation of LIME limits their work. In [11], the author proposed two methods named mean processing and feature generation method along with use of Light GBM and XGBoost to predict sepsis 6 hours in advance. The authors chose light GBM with feature generation over XGBoost due to its good performance, speed, and better generalization ability. However, use of many generated features leads to "Curse of Dimensionality," which results in poor performance of the model [12]. Similarly, in [13] the author used 65 high-resolutions vital sign features from EMR data to develop Artificial Intelligence Sepsis Expert (AISE). The AISE was able to predict sepsis 4 to 12 hours before and can provide features that has a significant impact on its predictions. Though it achieved a good Area Under the Receiver Operating Characteristics (AUROC), the high false positive rate was over two-fold higher than false negative rate, making it more sensitive. In [14], the authors proposed a sepsis prediction model by combining Kernel Extreme Learning Machine, which was trained on data extracted from blood samples, Chaotic Fruit fly Optimization which was used to enhance the predictive power of the model and RF for feature selection. Their proposed model outperformed all other ML models with a reasonably good sensitivity rate. But the low specificity of the model was the drawback of their work.

However, due to missingness of the data and the class imbalance problem previous research suffered from information loss due to reduction in the size of datasets. Other drawbacks like high false alarm rate, risk of automation bias and lack of interpretability to the clinicians question their adoption as sepsis prediction system. To overcome these, development of a smart AI based model is required that should not only reduce the false alarms but also should provide reasonable explanations behind its decision making, thus winning back the trust of the medical practitioners that it lagged. Therefore, the purpose of this paper is to develop a

smart model that can predict sepsis 6 hours in advance considering the false positives it encounters, as well as providing practitioners with diagnostic maps explaining how it predicts sepsis. This paper has the following contributions:

- Development of an early and highly accurate predictive model that identifies patients with sepsis by considering a minimal number of features as input.

- Implementation of a CNN-based autoencoder for extracting useful information from complex EMR data that boosts the performance of the classifier.

- Implementation of proper data pre-processing to retain important information.

- Use of GradCAM to explain the autoencoder's feature extraction process and SHAP for uncovering the hidden reasons behind the model's prediction.

## II. METHODOLOGY

In this section, we discuss the complete workflow of our research which is depicted in Fig. 1. It consists of three main steps namely, pre-processing, model building and prediction along with explanation. The raw data was passed through different stages of pre-processing to get a clean data on which model was trained.
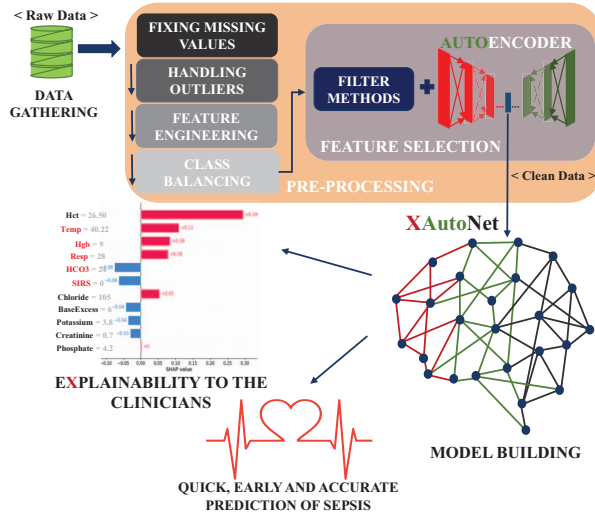


Fig. 1. The overall overview of proposed method.

### A. Data Gathering

The data which was used to train XAutoNet was downloaded from Kaggle and was a part of PhysioNet Challenge [15]. The data represents 20,336 patients' files, where in each file, each row represents a single hour worth of data of that patient within their ICU stay and were categorised between no sepsis and onset of sepsis 6 hours in advance.

### B. Pre-processing

The two main problems related to this dataset were missing values and class imbalance which was handled properly in the following steps.

*1) Fixing missing values:* Due to high percentage of missing values, Multiple Imputation Using Chained Equations (MICE) algorithm was used to impute them. It is a technique that imputes missing values in a dataset by leveraging information from other columns to estimate the best predictions for each missing value. MICE imputation calculates missing values, considers the data generation process for missingness, and maintains relationships while incorporating uncertainty [16].

*2) Handling outliers:* After imputing missing data, outliers were identified using Z-Score and Interquartile Range (IQR) methods for both normally distributed and skewed data. Invalid outliers were fixed by clipping them within their permissible range based on domain expert advice, while valid outliers were retained.

*3) Feature engineering:* Generation of additional feature named SIRS was done by combining existing features as it is considered as one of the helping tool in the identification of sepsis and hence can provides useful information while model building.

*4) Class balancing:* Due to high class imbalance problem in the data (2% of positive samples), a joint approach of oversampling of the minority class by Synthetic Minority Over-Sampling Technique [17] and then undersampling of the majority class based on clusters was incorporated to retain information as much as possible. The final dataset consisted of 32,000 normal records and 25,000 sepsis records.

*5) Feature selection:* For building a highly accurate model with least number of input features, the below mentioned methods were used.

*Filter Method:* To select the best 19 features from the total feature space, methods like Mutual Information (MI) score and Chi-square ($\chi2$) were used. MI helped us to get the best 18 features from the continuous feature space with the highest MI scores. It quantifies the feature-label dependency on a scale of 0 to 1, with a high score indicating strong dependency and 0 denoting independence. On the other hand, $\chi2$ test was used to select the best feature from the categorical feature space (SIRS and Gender). $\chi2$ assesses the association between categorical features, providing $\chi2$ statistics and p-value. Notable features for label determination have high $\chi2$ statistics and p-values below 0.05.
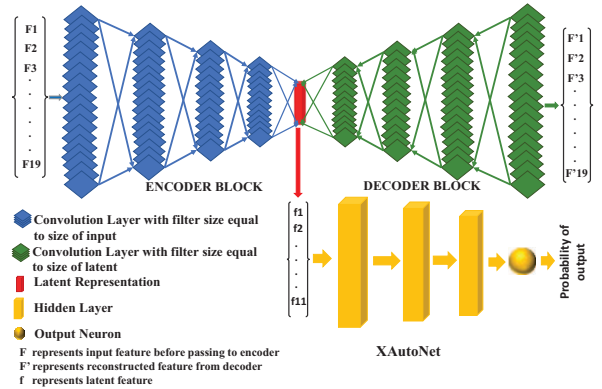


Fig. 2. The complete architecture of the CNN-CNN based Autoencoder with a classifier head for final classification.

*Autoencoder:* Further reduction of 19 features to an optimal size of 11 with the help of a CNN based autoencoder [18], [19]. The working and the description of the autoencoder architecture is mentioned below as well as shown in Fig. 2.

298

*a)* Encoder block: This block encodes or compresses the input dimension into a smaller latent dimension. It comprises of four 1D convolutional layers with decreasing number of filters in each layer. The role of each layer is to extract useful information from their input and compress them into a lower dimension. The output from the last convolution layer which is known as the latent dimension contains most of the information of the input layer in a compact form.

*b)* Decoder block: The role of this block is to decompress the latent dimension of the encoder back to the original input dimension thereby reconstructing them. The convolution layers used in this block are also connected sequentially with increasing number of filters, making it symmetric with the encoder. The output of the last convolution layer was flattened and passed to a final dense layer with nineteen neurons representing the corresponding input dimensions to be reconstructed with the help of linear transformation. Our objective is to minimize the total reconstruction error with the help of the expected mean squared error mentioned in (1).

$$L(\theta) = E\left[\left(I - Output\ (\theta)\right)^2\right] \quad (1)$$

Where $Output\ (\theta)$ is value predicted by the autoencoder which is the function of all the parameters $(\theta)$ and $I$ is the input.

### C. Model Building

A deep neural network model named XAutoNet was trained on the latent dimension outputted from the encoder block of the autoencoder (once it was trained) to classify them between normal or sepsis 6 hour before, shown in Fig. 2. The objective is to minimize the binary cross entropy between its output $(P'(\theta'))$ and ground truth $(Q)$ with respect to all parameters $(\theta')$.

$$Loss(\theta) = -\frac{1}{M}\sum_{t=0}^{M} Q_t\ log\left(P'_t(\theta')\right) + (1 - Q_t)\ log\left(1 - \left(P'_t(\theta')\right)\right) \quad (2)$$

Where M is the length of the dataset.

### D. Explainability

Interpretable methods like GradCAM [20] and SHAP [21] were implemented with a goal to increase the transparency and trustworthiness of both autoencoder and XAutoNet.

*1) GradCAM for autoencoder:* Participation of features in the formation of latent dimensions is shown through heatmaps by 1D GradCAM. To get the heatmap $hm \in \mathbb{R}^N$ for N feature, we first compute the gradient of class cls, $y^{cls}\ where\ y \in R^{cls}$ is the output from bottleneck, with respect to the feature map vector ( $F \in \mathbb{R}^{N \times C}$ ) from convolutional operation shown in (3):

$$\nabla = \frac{\partial y^{cls}}{\partial F} \quad (3)$$

Where $y^{cls}$ is calculated using gradient $(\nabla \in R^{N \times C})$ via backpropagation with respect to feature maps $(F \in \mathbb{R}^{N \times C})$ where $C$ is number of channels in feature map. Then we calculate the heatmap H = {$hm_i : hm_i \in \mathbb{R}^N, \forall i \in [1, N]$} by (4):

$$H = \frac{1}{C}\sum_{\forall i \in [1,C]} hm_i \odot \nabla_I \quad (4)$$

Where {$\nabla_i : \nabla_i \in \mathbb{R}^N, \forall i \in [1, N]$} = $\nabla$, C is the number of channels in feature map, and $\odot$ denotes the Hadamard product.

*2) SHAP for XAutoNet:* To check the contribution of features both globally and locally in model's prediction, we implemented deep SHAP. Deep SHAP makes use of Deep Learning Important FeaTures (DeepLIFT) [22] and Shapley value [23] and assigns each feature an importance value for a particular prediction. For each input, it assigns an attribute which represents the effect to the corresponding output.

On the other hand, Shapley assigns a value for each feature which represents the effect of the model predictions. For the set of features F, we use a pretrained network f on feature subset (S ⊆ F), both including it (S ∪ I) as well as excluding (S). Then, these two outputs are compared on current input i, shown in (5),

$$OD_S = f_{S \cup i}(x_{S \cup i}) - f_S(x_S), \quad (5)$$

Weights for corresponding $OD_S$ is calculated by the cardinality of the subsets (|S|) and set (|F|) as shown (6),

$$W_S = \frac{|S|!(|F|! - |S|! - 1)}{|F|!} \quad (6)$$

Shapley values ($SV_i$) are calculated using the weighted average of $OD_S$ with weights $W_S$ by (7),

$$SV_i = \sum_{S \subseteq F \setminus \{i\}} OD_S W_S\ (A) \quad (7)$$

Property of local accuracy explain the model g(x′) matches the original model f(x) when x = $h_x$(x′) shown in (8) below,

$$f(x) = g(x') = \phi_0 + \sum_{i=0}^{M} \phi_i x'_i \quad (8)$$

Where g(x′) represents the explanation model which consists of Shapley values.

## III. RESULTS AND DISCUSSION

*A. Comparative study:* XAutoNet was evaluated using 5-fold cross-validation on the entire dataset to prevent overfitting. Table 1 displays its good performance on each fold for various metrics, with mean scores above 0.90 and a low standard deviation (<0.015), indicating strong generalization ability.

TABLE I.  PERFORMACE OF XAUTONET IN 5-FOLD CROSS VALIDATION

| Fold | F1 Score | Precision | Recall | Accuracy |
|---|---|---|---|---|
| 1 | 0.92 | 0.92 | 0.91 | 0.93 |
| 2 | 0.92 | 0.92 | 0.92 | 0.93 |
| 3 | 0.93 | 0.95 | 0.91 | 0.93 |
| 4 | 0.93 | 0.94 | 0.93 | 0.94 |
| 5 | 0.94 | 0.94 | 0.94 | 0.94 |
| Mean ± SD | 0.93± 0.008 | 0.93± 0.012 | 0.92±0.012 | 0.94± 0.007 |

Furthermore, XAutoNet's performance was compared to traditional ML algorithms in Table 2, demonstrating its superior predictive power. XAutoNet outperformed all traditional models in Accuracy, Precision, Recall, and F1 Score.

It also outperformed most of the existing models not only with respect to accuracy but also with respect to the minimum number of input features. Moreover, its high recall value demonstrates its efficiency in identifying positive cases better than the other models.

299

| Model | Accuracy | F1 Score | Precision | Recall |
|---|---|---|---|---|
| KNN | 0.88 | 0.85 | 0.89 | 0.78 |
| Gradient Boost | 0.89 | 0.87 | 0.88 | 0.86 |
| Random Forest | 0.89 | 0.88 | 0.89 | 0.88 |
| Naïve Bayes | 0.62 | 0.49 | 0.59 | 0.43 |
| XG Boost | 0.90 | 0.89 | 0.90 | 0.89 |
| Decision Tree | 0.86 | 0.84 | 0.85 | 0.84 |
| SVM | 0.87 | 0.86 | 0.87 | 0.86 |
| Logistic Regression | 0.64 | 0.52 | 0.62 | 0.45 |
| ADA Boost | 0.88 | 0.85 | 0.88 | 0.82 |
| XAutoNet | **0.93** | **0.92** | **0.90** | **0.94** |

*B. Visualization of encoder using GradCAM:* After training the autoencoder, we calculated the heatmap (EB) for each encoder layer (E) in Fig. 3(A), with the help of GradCAM as discussed in explainability section. Based on the result, we can see that in the first encoder layer (E1), Hct is the most active feature while FiO2 is the least among all the features. Similarly, explainable block (EB2) that corresponds to second encoder layer (E2) again has Hct as most active feature while Lactate being the least. EB3 that corresponds to third encoder block (E3) shows that Phosphate is most active and again Lactate as least, and for the final encoder layer (E4), Potassium is the most active while pH and Calcium are the least active features according to EB4.
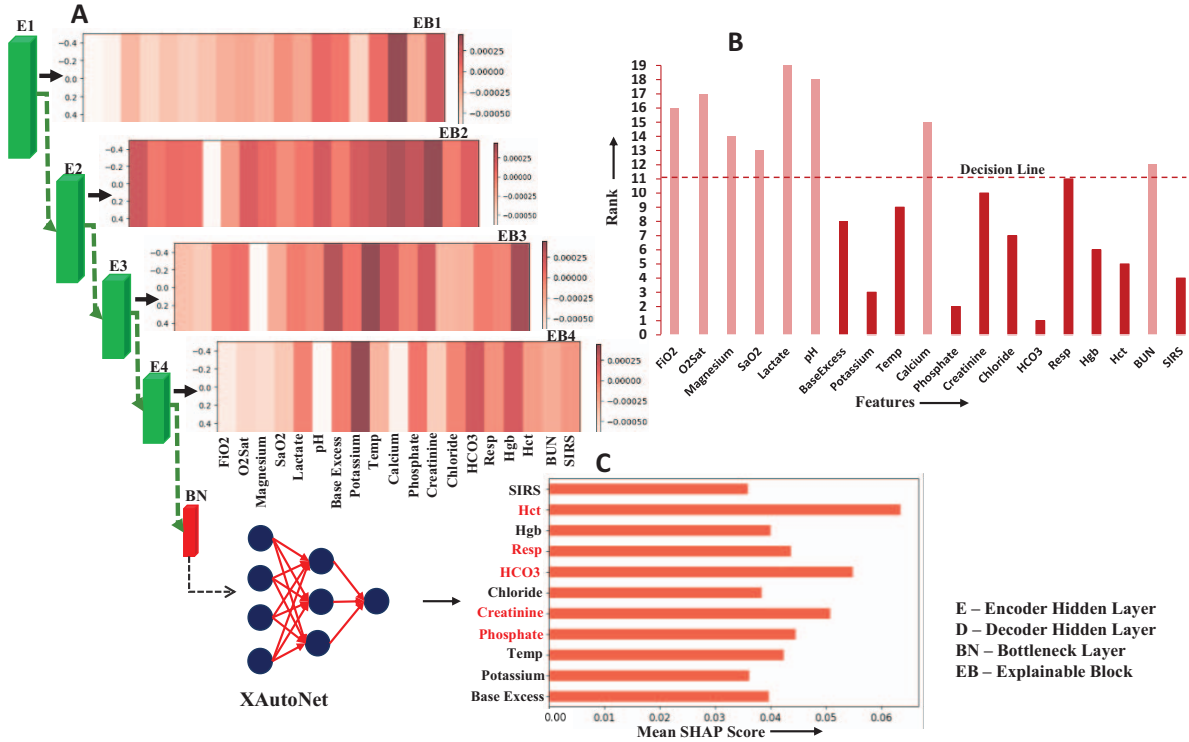


Fig. 3. Interpretability of Autoencoder and XAutoNet. (A) heatmaps (EB1 to EB4) showing features activation in each encoder layer in the form of heatmap using Grad-CAM. The darker the shade of a feature in the heatmap, more active that feature is and vice versa. (B) Represents feature ranking with respect to the bottleneck layer (BN) for the reconstruction using Autoencoder. (C) Global feature importance of XAutoNet using Mean SHAP Score.

Now to check the overall impact of these features in the formation of bottleneck, we combine the heatmaps of each encoder layer by discounting them from last encoder layer (E4) by (9).

$$DisHM = \sum_{i=1}^{4} \beta^i HM_i(x_1, x_2, \dots, x_{19}) \qquad (9)$$

Based on Discounting HeatMap (DHM), we showed the rank of the input features for discounting factor ($\beta$) of 0.9 in the form of bar plot in Fig. 3(B). This bar plot helps us to identify the top 11 features (lying on or below the decision line) that the BN layer comprises of.

*C. Interpretation of XAutoNet by DeepSHAP:* We assume a hypothesis that features from the bottleneck are the top 11

features based on DHM feature ranking method shown in Fig. 3(B). We calculated the mean SHAP value to provide the global feature importance by XAutoNet in determining sepsis, shown in Fig. 3(C). Out of all the features, Creatinine, Hct, Phosphate, HCO3 and Resp (marked in red color) are the five most impactful features due to their high magnitude as shown in Fig. 3(C), based on our hypothesis. With respect to local feature importance, waterfall plot becomes very insightful in providing both qualitative and quantitative explanations of feature's impact in XAutoNet's prediction. The waterfall plots in Fig. 4 represent the comparison of the impact of each feature, whether contributing (red bar) or offsetting (blue bar), in XAutoNet's prediction for a normal patient (A) and for a patient going to have sepsis in the next

300

6 hours (B). For patient A, correctly identified as normal by the model, 8 out of 11 features reduced the prediction probability of this person being prone to sepsis, with Hgb (hemoglobin) value of 9.99 g/dL contributing the most and Creatinine and Resp contributing the least. Whereas for patient B, correctly identified as septic 6 hours before, 6 out of 11 features increased the prediction probability for being prone to sepsis, with the most contributing feature as Temp (Temperature) with a value of 40.22 ℃ and Phosphate being the least.
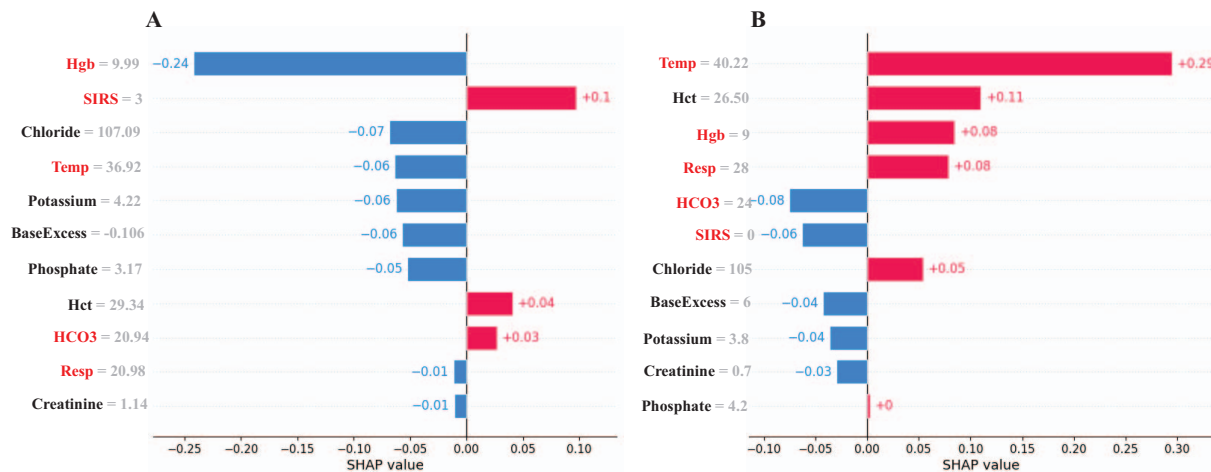


Fig. 4. Waterfall plots showing local explanations of the *XAutoNet* prediction on a (A)true negative instance and (B) true positive instance.

According to Fig. 4, features like Temp, SIRS, Hgb, Resp, and HCO3 of the patients (marked in red) played a crucial role in differentiating them. Other interesting facts with respect to the mentioned figure are that a SIRS score higher than 2 indicates a higher chance of sepsis as per its definition, which is also correctly depicted by the model. The mentioned fact is validated by a positive SHAP value of 0.1 as reported in Fig. 4(A) and a SHAP value of -0.06 for SIRS score of 0, in Fig. 4(B). Regarding normal body temperature (Temp) that ranges between 36.1℃ and 37.2℃, a high temperature of 40.22℃ observed in patient B contributes highly for being sepsis positive, with a SHAP score of +0.29. In contrast, the temperature for patient A that falls within the normal range, the SHAP score of -0.06 offered by XAutoNet implies that the temperature of this patient is not a point of concern for becoming positive. Similarly, for respiration (Resp) whose normal range is between 12 and 20, a high respiration rate of 28 in patient B indicates its high impact in becoming sepsis positive with a SHAP score of +0.08 but for patient A whose respiration rate is considered normal, SHAP score of -0.01 infers low risk with respect to respiration rate for this patient.

## IV. CONCLUSION

Our proposed model predicts sepsis 6 hours in advance with only 11 inputs, making it less complex than existing models. Our model outperforms other methods in the literature and traditional ML models, with good generalizability and improvement in recall. The autoencoder architecture effectively captures non-linearity in the data through pre-processing, enabling efficient classification. Future work will be inclusion of more relevant biomarkers, that are proven to be efficient in sepsis prediction, which can further improve the model's generalizability. Exploration of other pre-processing techniques including extraction of more features, and other ways of handling missing data.

SOURCE CODE

The complete source code of this work is available in GitHub repository at: https://github.com/Snehashis100/XAutoNet.

REFERENCES

[1] Rudd, K. E., Johnson, S. C., Agesa, K. M., Shackelford, K. A., Tsoi, D., Kievlan, D. R., ... & Naghavi, M. (2020). Global, regional, and national sepsis incidence and mortality, 1990–2017: analysis for the Global Burden of Disease Study. *The Lancet, 395(10219), 200-211.*

[2] Roy, Sudipta, Tanushree Meena, and Se-Jung Lim. "Demystifying supervised learning in healthcare 4.0: A new reality of transforming diagnostic medicine." *Diagnostics* 12, no. 10 (2022): 2549.

[3] Halder, S. et al. (2023). Fetal Brain Component Segmentation Using 2-Way Ensemble U-Net. In: Sharma, N., Goje, A., Chakrabarti, A., Bruckstein, A.M. (eds) Data Management, Analytics and Innovation. ICDMAI 2023. *Lecture Notes in Networks and Systems, vol 662. Springer, Singapore. https://doi.org/10.1007/978-981-99-1414-2_28*

[4] Kabiraj, Anwesh, Tanushree Meena, Pailla Balakrishna Reddy, and Sudipta Roy. "Detection and Classification of Lung Disease Using Deep Learning Architecture from X-ray Images." In Advances in Visual Computing: 17th International Symposium, ISVC 2022, San Diego, CA, USA, October 3–5, 2022, Proceedings, Part I, pp. 444-455. Cham: Springer International Publishing, 2022.

[5] Khan, M.A., Mittal, M., Goyal, L.M. et al. A deep survey on supervised learning based human detection and activity classification methods. Multimed Tools Appl 80, 27867–27923 (2021). https://doi.org/10.1007/s11042-021-10811-5

[6] Zargoush, M., Sameh, A., Javadi, M., Shabani, S., Ghazalbash, S., & Perri, D. (2021). The impact of recency and adequacy of historical information on sepsis predictions using machine learning. *Scientific reports*, *11*(1), 1-12.

[7] Awad, A., Bader-El-Den, M., McNicholas, J., & Briggs, J. (2017). Early hospital mortality prediction of intensive care unit patients using an ensemble learning approach. *International journal of medical informatics*, *108*, 185-195.

[8] Liu, R., Greenstein, J. L., Sarma, S. V., & Winslow, R. L. (2019, July). Natural language processing of clinical notes for improved early prediction of septic shock in the ICU. In *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)* (pp. 6103-6108). IEEE.

[9] R Fu, M., Yuan, J., Lu, M., Hong, P., & Zeng, M. (2019, September). An ensemble machine learning model for the early detection of sepsis from clinical data. In *2019 Computing in Cardiology (CinC)* (pp. Page-1). IEEE.

[10] Shankar, A., Diwan, M., Singh, S., Nahrpurawala, H., & Bhowmick, T. (2021, January). Early Prediction of Sepsis using Machine Learning. In *2021 11th International Conference on Cloud Computing, Data Science & Engineering (Confluence)* (pp. 837-842). IEEE.

[11] Zhao, X., Shen, W., & Wang, G. (2021). Early prediction of sepsis based on machine learning algorithm. *Computational Intelligence and Neuroscience*, *2021*.

[12] Debie, E., & Shafi, K. (2019). Implications of the curse of dimensionality for supervised learning classifier systems: theoretical and empirical analyses. *Pattern Analysis and Applications*, *22*(2), 519-536.

[13] Nemati, S., Holder, A., Razmi, F., Stanley, M. D., Clifford, G. D., & Buchman, T. G. (2018). An interpretable machine learning model for accurate prediction of sepsis in the ICU. *Critical care medicine*, *46*(4), 547.

[14] Wang, X., Wang, Z., Weng, J., Wen, C., Chen, H., & Wang, X. (2018). A new effective machine learning framework for sepsis diagnosis. *IEEE access*, *6*, 48300-48310.L'L

[15] Reyna, Matthew A., Josef, Christopher S., Jeter, Russell , Shashikumar, Supreeth P., Westover, M. Brandon, Nemati, Shamim, Clifford, Gari D, Sharma, Ashish. Early Prediction of Sepsis From Clinical Data: The PhysioNet/Computing in Cardiology Challenge 2019. Critical Care Medicine: February 2020 - Volume 48 - Issue 2 - p 210-217doi: 10.1097/CCM.0000000000004145.

[16] Van Buuren, S., & Groothuis-Oudshoorn, K. (2011). mice: Multivariate imputation by chained equations in R. *Journal of statistical software*, *45*, 1-67.

[17] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: synthetic minority over-sampling technique. Journal of artificial intelligence research, 16, 321-357.

[18] Wang, Y., Yao, H., & Zhao, S. (2016). Auto-encoder based dimensionality reduction. Neurocomputing, 184, 232-242.

[19] Kumar, K., Kumar, H., & Wadhwa, P. (2023). Encoder–Decoder Network-Based Prediction Model for Trend Forecasting in Currency Market. In Soft Computing for Problem Solving: Proceedings of the SocProS 2022 (pp. 211-223).

[20] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision* (pp. 618-626).

[21] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in neural information processing systems*, *30*.

[22] Shrikumar, A., Greenside, P., & Kundaje, A. (2017, July). Learning important features through propagating activation differences. In *International conference on machine learning* (pp. 3145-3153). PMLR.

[23] Fryer, D., Strümke, I., & Nguyen, H. (2021). Shapley values for feature selection: the good, the bad, and the axioms. *IEEE Access*, *9*, 144352-144360.