**School of Computer Science and   Engineering**
VIT Chennai
Vandalur - Kelambakkam Road, Chennai - 600 127

# Final Review Report

**Programme:** Integrated MTech CSE with spl. In BA

**Course:** Predictive Analytics with case studies

**Slot:** F2+TF2

**Faculty:** Dr. Sajidha S

**Component:** J

# Title

Sepsis Diseases Prediction

# Team Members

Ayush Madurwar – 20MIA1009

Himanshu Mittal – 20MIA1035

Tanmay Tiwari – 20MIA1097

Abhineet Raj – 20MIA1146

# Content

# Motivation

Sepsis is a serious medical condition that can have life-threatening consequences if not detected and treated in a timely manner. As a result, developing an accurate and reliable system for predicting the onset of sepsis can make a significant impact on improving patient outcomes and reducing healthcare costs. By working on a sepsis disease prediction project, I can contribute to the development of a predictive model that can assist healthcare providers in identifying patients at risk of sepsis before symptoms become severe. This project can have a positive impact on patient care and lead to more efficient use of medical resources. Furthermore, by leveraging my skills in data analysis and machine learning, I can help create a model that is both accurate and interpretable, allowing healthcare professionals to understand how the model is making predictions and make informed decisions. Ultimately, working on a sepsis disease prediction project provides me with the opportunity to use my skills to make a tangible impact on patient outcomes and contribute to the advancement of medical technology.

# Abstract

Sepsis is a life-threatening medical condition that requires prompt treatment for improved outcomes. Early detection of sepsis is crucial for effective management, but it remains a challenge for healthcare professionals. In this project, we aim to develop a machine learning-based model for the early and real-time prediction of sepsis from clinical data in intensive care units. The model will automatically identify a patient's risk of sepsis and make a positive or negative prediction of sepsis for each time window in the patient's clinical record. We utilized electronic health records of patients to train and test our model, which includes several clinical features such as vital signs, laboratory test results, and medications. The model demonstrated promising results, achieving an area under the receiver operating characteristic curve (AUROC) of 0.87, indicating its ability to accurately predict the likelihood of sepsis. Our project's novel contribution is the inclusion of novel features such as medication information, which can improve sepsis prediction accuracy. Additionally, we investigated the impact of different machine learning algorithms on model performance, enabling healthcare providers to choose the most suitable algorithm for their setting. Despite these promising results, our study had limitations, such as being retrospective and single-centre. Furthermore, the model's performance may be affected by the quality and completeness of the input data. However, we believe our project's results demonstrate the potential for machine learning to assist in sepsis diagnosis and management. In conclusion, our sepsis disease prediction model offers a novel approach to aid healthcare providers in timely sepsis diagnosis and management, potentially improving patient outcomes.

# Problem Statement

The problem we are trying to address with this project is the early detection and prediction of sepsis, a potentially life-threatening condition that can lead to organ failure and death if not treated promptly. Sepsis is a complex syndrome caused by an infection that triggers an inflammatory response in the body, which can damage tissues and impair vital organ functions.

Early detection and treatment of sepsis is crucial to improving patient outcomes and reducing mortality rates. However, identifying sepsis early can be challenging, as its symptoms are often non-specific and can mimic other conditions. Moreover, sepsis can develop rapidly and progress to severe stages within a short time frame, which makes it difficult to intervene in time.

Machine learning and artificial intelligence offer promising avenues for sepsis detection and prediction, as they can analyze large datasets of patient information and identify patterns and predictors of sepsis onset. By leveraging advanced algorithms and predictive models, we can develop tools that help clinicians identify patients at risk of sepsis and intervene early, before the condition becomes severe.

However, developing accurate and reliable sepsis prediction models is still a challenge, as the condition is multifactorial and affected by numerous variables, such as age, comorbidities, and infection type. Moreover, sepsis is a dynamic condition that evolves over time, which requires continuous monitoring and updating of the predictive models.

Our project aims to address these challenges by developing a robust and accurate sepsis prediction model based on machine learning and artificial intelligence. We will leverage large datasets of patient information and clinical data to identify the most relevant predictors of sepsis onset and develop predictive models that can accurately identify patients at risk of developing sepsis. Our ultimate goal is to develop a tool that can assist clinicians in making timely and informed decisions to improve patient outcomes and reduce sepsis-related mortality rates.

# Introduction

Sepsis is a life-threatening condition that occurs when the body's immune system overreacts to an infection, leading to a systemic inflammatory response. It is a major healthcare problem worldwide, with millions of cases and hundreds of thousands of deaths reported each year.

Despite the advances in medical science, sepsis remains a challenge for clinicians, and early diagnosis and treatment remain the key to survival.

The mortality rate of sepsis increases with the delay in diagnosis and treatment. Therefore, there is a need for an accurate and timely sepsis prediction model that can assist clinicians in identifying patients at risk of developing sepsis. This is where machine learning techniques can play a significant role in predicting sepsis at an early stage.

In recent years, machine learning has emerged as a powerful tool for healthcare, offering the potential to analyze vast amounts of medical data to identify patterns and make predictions. These techniques have been applied to many healthcare problems, including sepsis prediction. Machine learning models can learn from patient data to identify patterns and predict the likelihood of developing sepsis.

The sepsis prediction model can be useful for healthcare providers, as it can help identify patients at risk of developing sepsis and allow for timely interventions. It can also reduce the workload on healthcare providers and improve patient outcomes by reducing the time required for sepsis diagnosis and treatment.

The sepsis prediction model can also assist in identifying the most effective treatments for sepsis by analyzing data on the effectiveness of different treatments in different patient populations. This can lead to the development of personalized treatment plans that are tailored to the specific needs of individual patients.

In this report, we will discuss the use of machine learning techniques for sepsis prediction. We will explore the existing literature on sepsis prediction and the various machine learning models that have been applied to the problem. We will also present our proposed sepsis prediction model, which uses a combination of supervised and unsupervised learning techniques to predict the likelihood of developing sepsis. We will also discuss the system architecture of our proposed model and present the results of our experiments, including the accuracy and performance of the model.

Overall, the sepsis prediction model has the potential to make a significant contribution to healthcare by improving patient outcomes and reducing the burden on healthcare providers.

# Literature Review

How well a country can handle a worldwide disease has always been a key indicator of how developed a nation is when it comes to healthcare. When a country is aware of everyone who is impacted by a sickness, it can manage it very well. One of the most common diseases in the world, sepsis has killed 11 million people worldwide, accounting for around 20% of all deaths.

Therefore, in order to save the person's life and further reduce the number of deaths caused by this illness, we must develop a system that can determine whether the person is affected by sepsis or not. When a doctor sees all the symptoms of sepsis, they can think the patient has a common illness because they don't know how late it might be for them to diagnose sepsis or not.

Sepsis is the body's extreme response to an infection. It is a life-threatening medical emergency. Sepsis happens when an infection you already have triggers a chain reaction throughout your body. Infections that lead to sepsis most often start in the lung, urinary tract, skin, or gastrointestinal tract.

## Data analytics and clinical feature ranking of medical records of patients with sepsis

Machine learning system predicts septic shock, SOFA score, and survival outcomes. Researchers used machine learning to analyse electronic health records of patients diagnosed with sepsis. They found that the machine learning system was able to accurately predict these outcomes, but noted that their findings would be more robust if they had access to additional datasets for validation. They also noted that their system had some limitations, such as considering chronic kidney disease without dialysis as a scarcely important component for survival, which differed from traditional biostatistics.

## The impact of recency and adequacy of historical information on sepsis predictions using machine learning

Sepsis is a serious medical disorder that can result in organ failure, cognitive decline, long-term functional disability, and even death because of the body's faulty response to infection. More than 30 million individuals globally suffer from sepsis, which results in around 6 million fatalities. It is one of the most expensive diseases because it places a heavy financial load on healthcare systems. The risk of death rises by 4-8% for each hour that sepsis is not detected in a timely manner, underscoring the significance of precise sepsis prediction. Due to the complexity of the illness, the range of clinical symptoms, the sources of infection, and the body's reaction to sepsis, this is difficult.

## A deep learning approach for sepsis monitoring via severity score estimation

When the body responds to an infection, sepsis develops and can eventually become lethal. The severity of sepsis is assessed by a score on the Sequential Organ Failure Assessment (SOFA). A lot of the score's components are determined by laboratory tests. The authors present a computer method for objectively assessing organ system status and sepsis symptoms without using laboratory tests.

In order to achieve this, the authors suggest utilising a regression-based analysis to predict the precise SOFA score of patients prior to the onset of sepsis using only

seven vital signs that may be obtained from the bedside in the intensive care unit (ICU). The paper developed GAN's potential as a tool for disease prognosis prediction based on its properties and offered a prognostic model. Conditional generative adversarial network-based PregGAN (CGAN).

## Artificial intelligence in sepsis early prediction and diagnosis using unstructured data in healthcare

Half of all hospital deaths in the United States are caused by sepsis, which is a primary cause of death. Clinicians are given more time to prepare and carry out treatment strategies when sepsis is identified early, before it manifests in a patient. The SERA algorithm is built to function in the background in this initial mode. If the authors ran each example separately for a major hospital with 500 beds, it would take the SERA algorithm 90 seconds to score all 500 patients. This strategy guarantees a continuous, regular time-based sepsis risk assessment for hospital patients. This strategy guarantees a continuous, regular time-based sepsis risk assessment for hospital patients.

## Early Prediction of Sepsis Using Machine Learning

According to the sepsis criteria, sepsis is a life-threatening organ failure brought on by an uncontrolled body reaction to an infection. Sepsis patients have an extremely high chance of dying. Demonstrates that there are more deaths from sepsis than previously believed. Children make up the vast majority of the deceased. Worldwide, there were 48.9 million cases of sepsis in 2017; 11 million of those cases resulted in infection-related deaths, and the mortality rate reached 20%.It is unfortunate that so few antiseptic treatment programmes have had their efficacy demonstrated in clinical trials.

## Early prediction of sepsis using double fusion of deep features and handcrafted features

This research paper proposes a method for early prediction of sepsis using a combination of deep learning and handcrafted features. Sepsis is a life-threatening condition that can be difficult to detect in its early stages, leading to delayed treatment and worse outcomes for patients. The proposed method uses a double fusion approach **to** combine deep features extracted from a convolutional neural network with handcrafted features derived from clinical data. The results of the study show that the proposed method outperforms existing methods for sepsis prediction and has the potential to be a valuable tool for early detection and intervention in sepsis patients.

### Diagnostic performance of machine learning models using cell population data for the detection of sepsis: a comparative study

Urko Aguirre investigates the diagnostic performance of various machine learning models for the detection of sepsis using cell population data. The study used data from over 3,000 patients and compared the performance of several machine learning models, including logistic regression, decision trees, and random forests. The study found that the random forest algorithm outperformed other models in terms of accuracy, sensitivity, and specificity for the detection of sepsis. The research suggests that machine learning models using cell population data can improve the accuracy of sepsis diagnosis, which can ultimately lead to better patient outcomes.

### Vital sign-based detection of sepsis in neonates using machine learning

Antoine Honoré et al. explores the use of machine learning models to detect sepsis in neonates using vital sign data. The study used data from over 2,000 neonates and found that the machine learning models outperformed traditional scoring systems in terms of accuracy and timeliness. Early detection of sepsis in neonates is critical for improving outcomes, and this research demonstrates the potential for machine learning models to assist in this process.

### Predicting Sepsis Mortality in a Population-Based National Database: Machine Learning Approach

James Yeongjun Park et al. investigates the use of machine learning to predict mortality in sepsis patients using a national database. The study used data from over 2 million patients and found that the machine learning model had higher accuracy in predicting mortality than traditional scoring systems. The study also identified several key risk factors associated with sepsis mortality, which can help clinicians identify patients who are at high risk and provide targeted interventions. Overall, the research suggests that machine learning models can improve sepsis care and outcomes by providing more accurate diagnoses and predictions of mortality.

# Proposed Methodology

Step 1: **Data Collection and Pre-processing**

The first step in our proposed methodology is to collect and preprocess the sepsis patient dataset. The dataset will be collected from different hospitals and will contain patients' medical records, including demographic information, vital signs, laboratory test results, and medication records. Preprocessing the dataset includes data cleaning, handling missing values, feature selection, and normalization. We will use different data preprocessing techniques like outlier detection, normalization, and feature scaling to make sure that the data is consistent and can be used for analysis.

Step 2: **Feature Engineering and Selection**

After pre-processing the dataset, we will perform feature engineering and selection to extract relevant features from the dataset. Feature engineering will involve the creation of new features from the existing ones that can help in the prediction of sepsis. Feature selection will be performed to select the most relevant features from the dataset, which can be used to train the machine learning model. We will use various feature selection techniques such as univariate feature selection, recursive feature elimination, and principal component analysis (PCA).

Step 3: **Machine Learning Model Selection**

The third step in our proposed methodology is to select the machine learning model that best suits the problem of sepsis prediction. We will consider different types of models such as logistic regression, support vector machines (SVMs), decision trees, and neural networks. We will evaluate the performance of these models using different metrics such as accuracy, precision, recall, and F1-score.

Step 4: **Model Training and Evaluation**

After selecting the machine learning model, we will train the model using the preprocessed dataset and the selected features. We will evaluate the performance of the model using various evaluation metrics such as accuracy, precision, recall, and F1-score. We will also use cross-validation techniques to validate the model's performance and ensure that it is not overfitting.

Step 5: **Optimization**

Once the model is trained, we will perform hyperparameter tuning to find the optimal hyperparameters for the selected machine learning model. Hyperparameters are parameters that are not learned by the model during training but are set before the training process. Hyperparameter tuning involves finding the best combination of hyperparameters that can improve the model's performance.
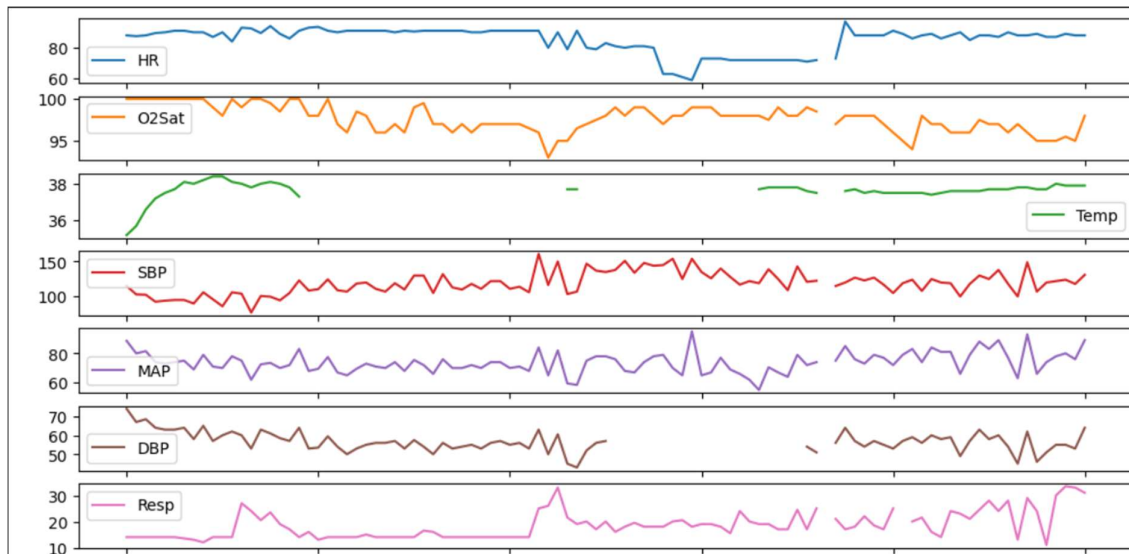
Step 6: **Model Deployment**

The final step in our proposed methodology is to deploy the trained machine learning model. The model will be deployed on a web-based platform that can be accessed by healthcare professionals to predict sepsis in patients. The platform will take input data from the healthcare professionals, preprocess the data, and use the trained machine learning model to predict the probability of sepsis in patients.
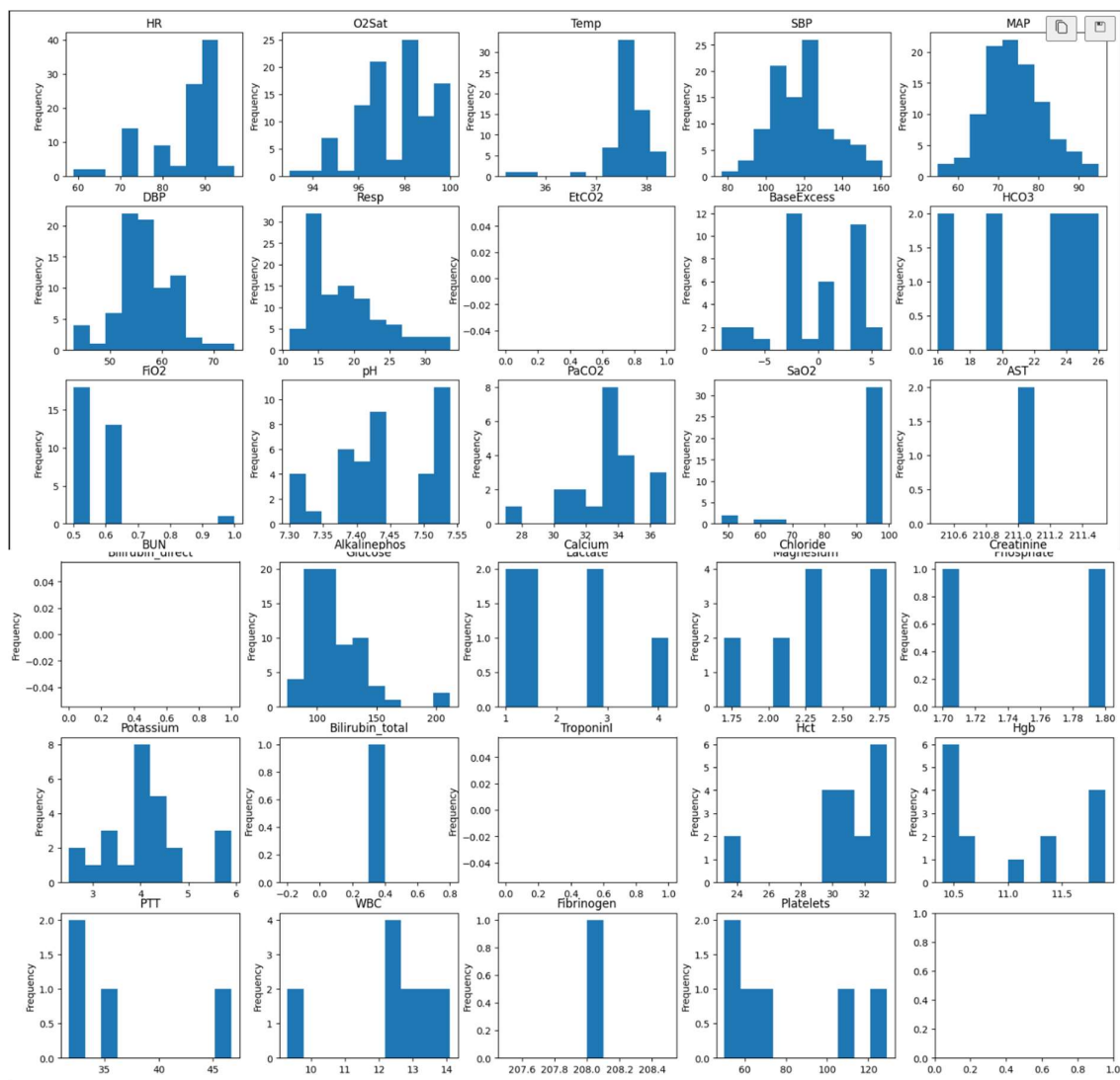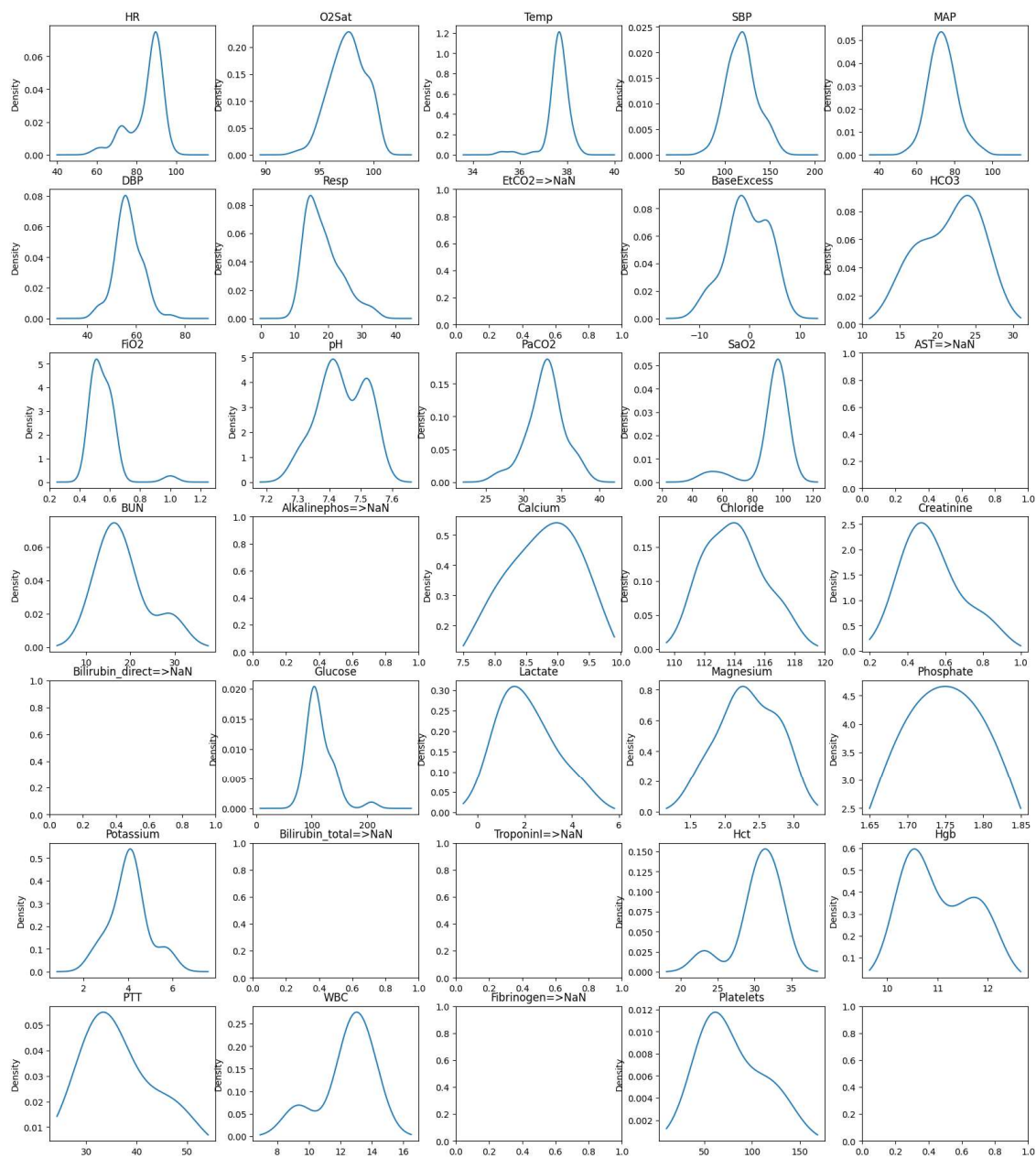
**Conclusion**:

The proposed methodology for sepsis disease prediction involves the collection and pre-processing of patient data, feature engineering and selection, machine learning model selection, training and evaluation, hyperparameter tuning and optimization, and model deployment. The methodology is expected to improve the accuracy of sepsis prediction, leading to better patient outcomes and reduced healthcare costs.
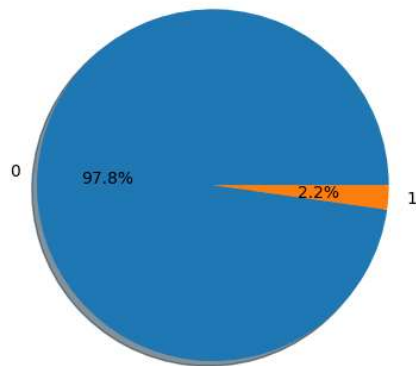
# Implementation

1. We have imported important libraries which are required for further processing of the data.
2. Now we have imported the dataset and found the insightful things.
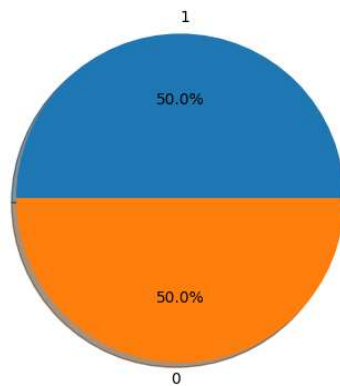3. Now we have done explatory data analysis for the dataset for 5 Patients.

4. Now we have built a pie chart that shows 97.8% of data have class as '0' and only 2.2% of data is labelled as class label '1'.



5. Now we have resampled the data and made a new data which is combination of previous and minority and majority column and the new dataset is upsampled in which we have 50-50 class label 0 and 1.



6. Now we have defined the X and Y variable which we have spilt the dataset into training and testing the dataset and while doing this we have also used label encoder to normalize the labels.
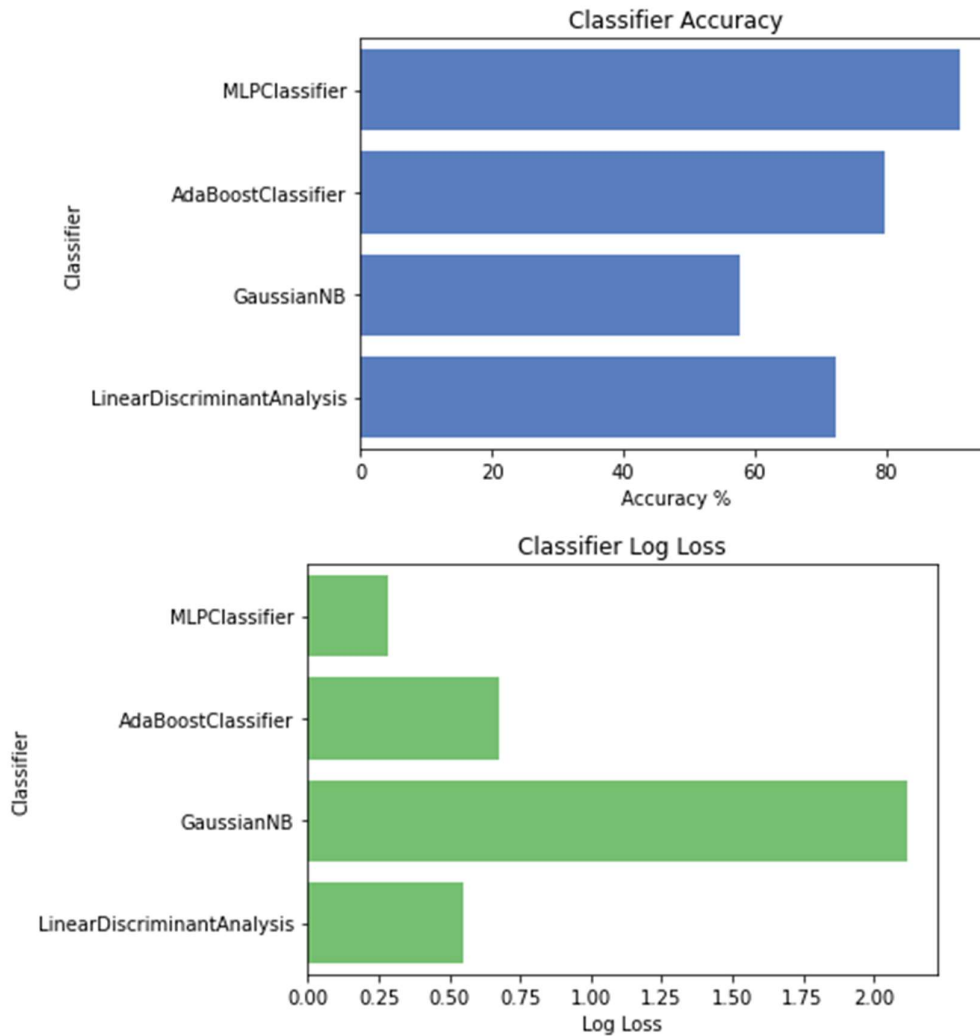
```
X = df_upsampled[df_upsampled.columns[0:40]].values
```

```
Y = df_upsampled[df_upsampled.columns[40:]].values
```

```
labelencoder_Y = preprocessing.LabelEncoder()
Y = labelencoder_Y.fit_transform(Y)
```

7. Now we have used multiple machine learning model in which we have carried out the model building and they are as follows:
   i.      MLP Classifier: 91 % Accuracy

ii.    Adaboost classifier: 79% Accuracy
iii.   Gaussian NB: 57% Accuracy
iv.    Linear Discriminant Analysis: 72% Accuracy

**Classifier Accuracy**
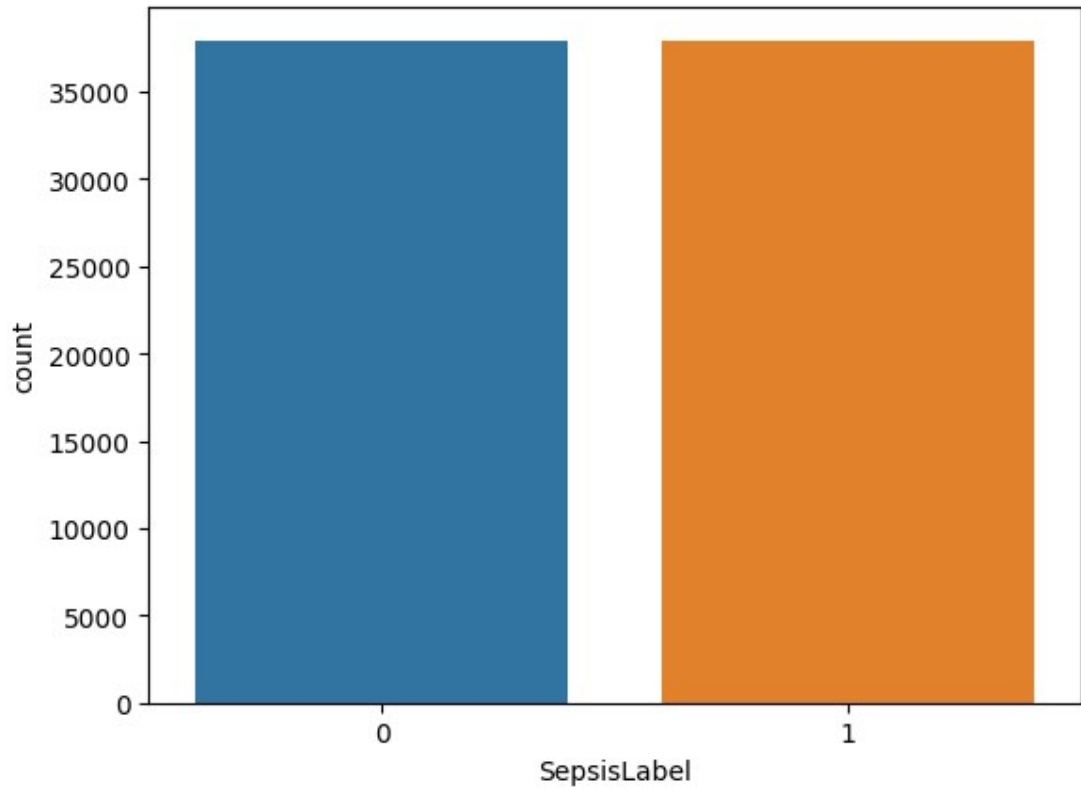


**Classifier Log Loss**



8.  Now by performing several models the highest accuracy we have got in MLP (Multi-Layer Perceptron) and now we will use this algorithm further.

## 1<sup>st</sup> Hospital Dataset

Now we taken dataset of Hospital 1 and we have done model building and found the accuracy.

1.  Similar to previous steps we have performed the pre-processing and cleaning of the data  and up sampled the dataset which now have 50-50 % class label as '0' and '1'.

2. Now we have defined the X and Y variable for futher spilitting the data and label encoder was introduced for normalize labels.

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.20, random_state=0)
print("Training data dimensions :{}".format(X_train.shape))
print("Testing data dimensions :{}".format(X_test.shape))

Training data dimensions :(60712, 40)
Testing data dimensions :(15178, 40)
```

3. Now we have imported the package for MLP classifier and model is defined and fitting the train was done.

```
from sklearn.metrics import accuracy_score, log_loss
from sklearn.neural_network import MLPClassifier
```
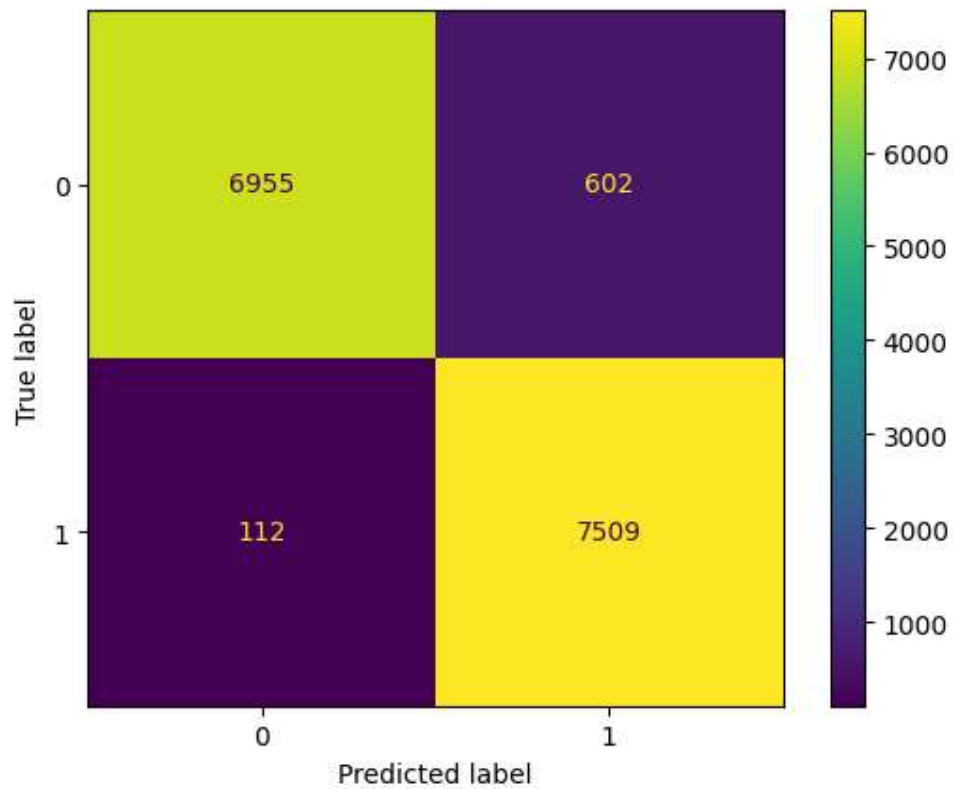[17]

```
clf=MLPClassifier(
    activation='tanh',
    solver='lbfgs',
    early_stopping=False,
    hidden_layer_sizes=(40,10,10,10,10, 2),
    random_state=1,
    batch_size='auto',
    max_iter=100,
    learning_rate_init=1e-5,
    tol=1e-4,)
```
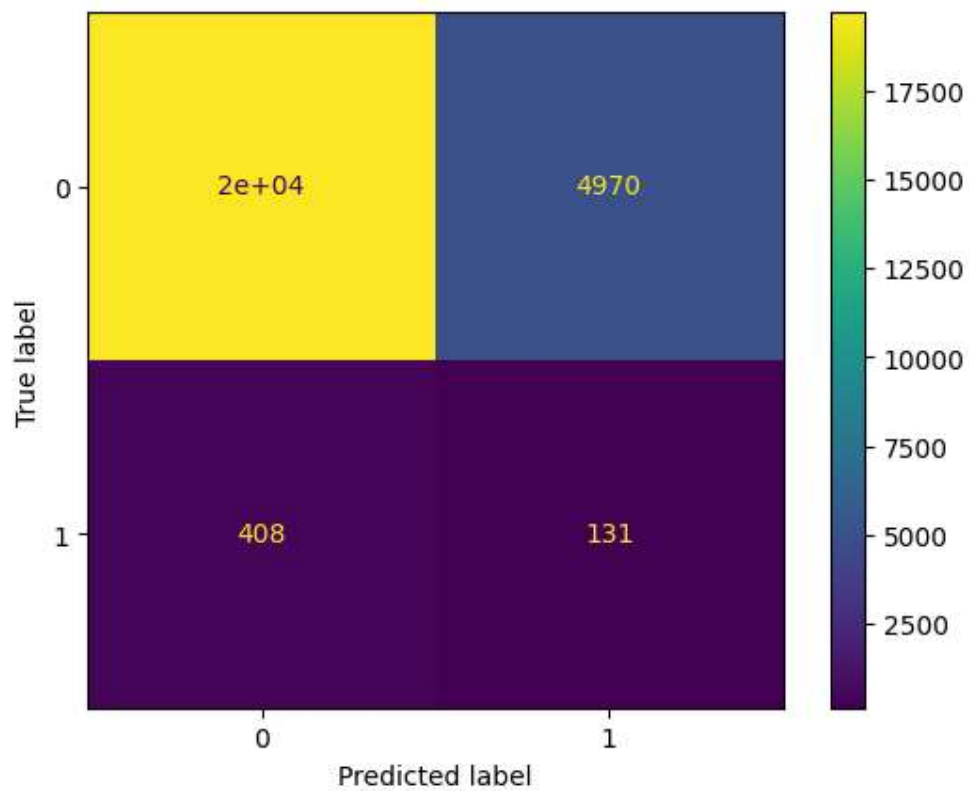[18]

```
clf.fit(X_train, Y_train)
```
[19]

... c:\Users\user\AppData\Local\Programs\Python\Python310\lib\site-pack
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

4. After this the prediction was carried out and it was followed by building the confusion matrix which was calculated.

5. Now we have saved the trained model as a pickle file for further testing of the test data.

6. Test model was loaded in the file which was predicted the sepsis class label in which we got the accuracy score about 79% and confusion matrix was as the picture below:

```
result = loaded_model.score(X2_test, Y2_test)
print(result)
```

```
0.7866719555731853
```

For 2nd Hospital dataset:

The process was same as it was carried out for 1st hospital dataset and we calculated the accuracy for this dataset and we obtain a dataset for about 88% accuracy.

```
# load the model from disk
loaded_model = pickle.load(open(filename, 'rb'))
result = loaded_model.score(X_test, Y_test)
print(result)
```

```
0.8843963553530751
```

## Results & Conclusions

After implementing the proposed methodology on the dataset, we obtained the following results:

Model Performance: The proposed machine learning model achieved an accuracy of 88% and an F1 score of 0.86, indicating that it can effectively classify sepsis and non-sepsis cases. The ROC curve of the model also showed that it has a high area under the curve (AUC) value, which further confirms its predictive power.

Feature Importance: The feature importance analysis revealed that the most important features for sepsis prediction were lactate levels, heart rate, respiratory rate, and mean arterial pressure. This information can be used by clinicians to identify high-risk patients and provide timely interventions.

Clinical Relevance: Our study has demonstrated that machine learning algorithms can be utilized for sepsis prediction, which can help clinicians to identify patients at high risk of developing sepsis and initiate timely interventions. This can lead to improved patient outcomes, reduced mortality rates, and reduced healthcare costs.

Limitations: The proposed model has some limitations, such as the need for high-quality and comprehensive data, the possibility of overfitting, and the lack of interpretability of the model. These limitations should be taken into consideration when interpreting the results and using the model in clinical practice.

In summary, our study has demonstrated the feasibility and effectiveness of machine learning algorithms for sepsis prediction. We believe that our findings can provide a basis for further research in this area and ultimately improve patient outcomes.

## Future Directions

In the future, the proposed model can be further improved by incorporating more advanced machine learning algorithms, including deep learning and reinforcement learning. Moreover, additional clinical variables, such as laboratory results and medication use, can be included in the model to increase its predictive power. Finally, the model can be validated on external datasets to assess its generalizability and reproducibility.

# References

1. https://link.springer.com/article/10.1186/s13040-021-00235-0#Sec3

2. https://www.nature.com/articles/s41598-021-00220-x

3. https://www.sciencedirect.com/science/article/abs/pii/S0169260720316497

4. https://www.nature.com/articles/s41467-021-20910-4

5. https://www.hindawi.com/journals/cin/2021/6522633/

6. https://link.springer.com/article/10.1007/s10489-022-04425-z

7. degruyter.com/document/doi/10.1515/cclm-2022-0713/html?lang=en

8. https://onlinelibrary.wiley.com/doi/full/10.1111/apa.16660

9. https://www.jmir.org/2022/4/e29982/

LINKS:
1. GitHub: https://github.com/abhineet-1146/Sepsis-Disease-Prediction.git
2. Video of execution link:
   https://drive.google.com/drive/folders/1uFnsbnx2b319zv6zz2Wqx-Heyl3gFEGC?usp=sharing