**BITS Pilani**
Pilani | Dubai | Goa | Hyderabad
**Work Integrated Learning Programmes**

## 1: Overview

- **Objective**: We are provided with the data set, this dataset contains 2126 records of features extracted from Cardiotocogram exams, which were then classified by expert obstetrician into 3 classes: "Normal", "Suspect" & "Pathological". Our objective is to perform the following operations on the dataset:
  **Predictive Analytics**
    - Apply Decision Tree Classifier on the dataset. Consider the X variables as baseline value, fetal_movement, uterine_contractions, light_decelerations, severe_decelerations, prolonged_decelerations and abnormal_short_term_variability; and Y variable as fetal_health. Keep the train-test split at 70-30 ratio.
    - Then finding the overall accuracy of the decision tree model, the confusion matrix, the implications of Type-I error and Type-II error in this example. What are the Precision, Recall, F1 score for the three classes, based on the confusion matrix, Plot the ROC curves for the three classes

## 2: Data Acquisition

- **How many features:** 21 (21 independent variables and one dependent target variable)
- **Size of the dataset:** The dataset size: (2126, 22)
- **Multiple files:** No single input file (fetal_health.csv downloaded from https://www.kaggle.com/datasets/andrewmvd/fetal-health-classification)
- **What kind of data – numerical or character**: Numerical Data
- **Balanced or imbalanced – what is the distribution:** imbalanced as normal is dominant
- **Distribution of Training set, validation set, testing set:** Training size:70% and test size: 30%

## 3: Data Wrangling

- There is no null value observed in the given dataset

Highly imbalanced dataset

Now all the classes have same no of samples after oversampling of minority classes



```
data.isnull().sum()

baseline value                                          0
accelerations                                           0
fetal_movement                                          0
uterine_contractions                                    0
light_decelerations                                     0
severe_decelerations                                    0
prolongued_decelerations                                0
abnormal_short_term_variability                         0
mean_value_of_short_term_variability                    0
percentage_of_time_with_abnormal_long_term_variability  0
mean_value_of_long_term_variability                     0
histogram_width                                         0
histogram_min                                           0
histogram_max                                           0
histogram_number_of_peaks                               0
histogram_number_of_zeroes                              0
histogram_mode                                          0
histogram_mean                                          0
histogram_median                                        0
histogram_variance                                      0
histogram_tendency                                      0
fetal_health                                            0
dtype: int64
```
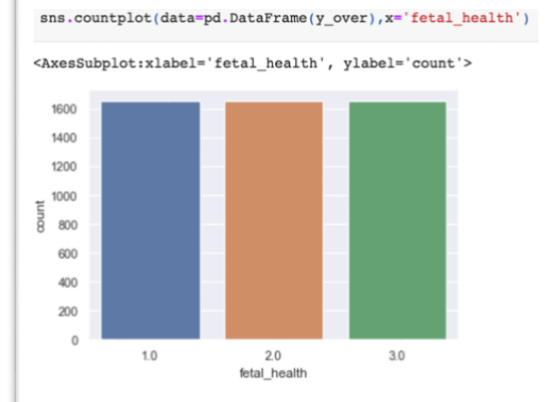There is no null value let's proceed further to check class blancing



Text(0.5, 1.0, 'Number of samples of each class:')

Observation: An highly Imbalanced Dataset. Which is obvious as Normal would be dominant.

Solution: Oversampling of the minority classes. To make better predictions.(Will perform later in the notebook)



sns.countplot(data=pd.DataFrame(y_over),x='fetal_health')
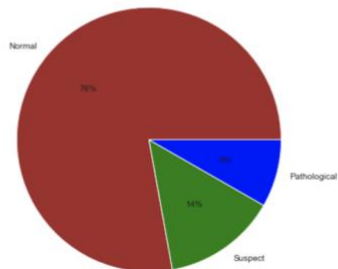
<AxesSubplot:xlabel='fetal_health', ylabel='count'>

Observation: Attributes such as severe_decelerations, prolongues_decelerations are more or less constant with few variations. We removed them using PCA by getting 3 components with highest variance and 95% total variance in Dimensionality Reduction.
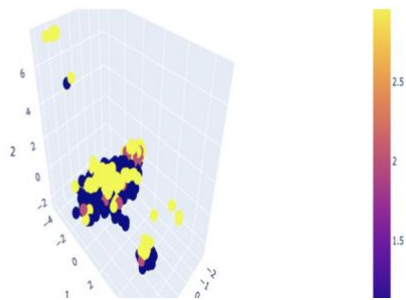
## 4: Feature Engineering Techniques and Results

Feature engineering plays an important role as everything from the data to the output of the same is dependent on the feature engineering which is being performed. Firstly a pie chart shown is being visualised which shows about the different health types of the fetal and get to know if it is a significant feature or not. While we are getting a correlation matrix of all the features to their significance. The correlation matrix is used in getting the relation of each feature or attribute with itself and also with the other features present out in the dataset.
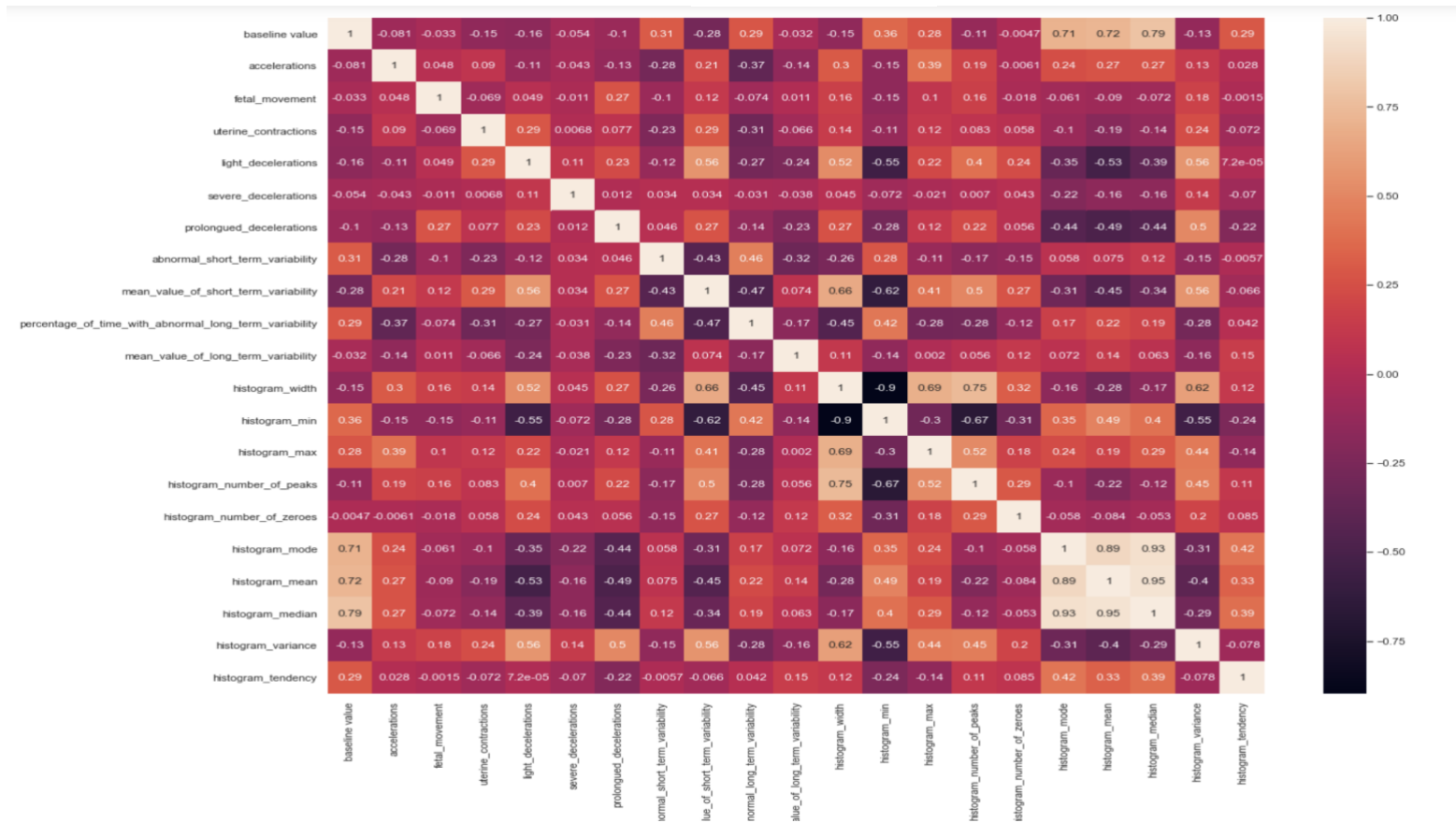
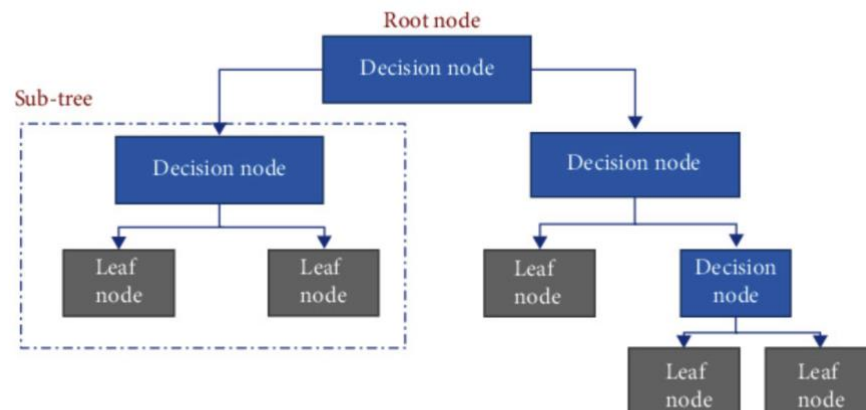Correlation Matrix



Data Distribution



Dimensionality Reduction and Visualization

## 6: Modelling

It depicts the whole decision tree design flowchart. This assignment makes use of a decision tree classifier. This classifier seems to recursively divide the example space. It is a predictive paradigm that acts as a mapping between the characteristics of an item and their values. It regularly splits each potential data result into parts. Each non leaf node corresponds to a feature experiment, each branch to the outcome of the experiment, and each leaf node to a judgment or classification. The root node of the tree, which is at the very top, reflects the most often used prediction model. The decision node and the leaf node are the two nodes in a decision tree. The choice nodes are used to make those selections and have numerous branches, whereas the leaf nodes are the result of those choices and contain no additional branches. The outcomes of the tests or judgments are contingent on the dataset's properties.

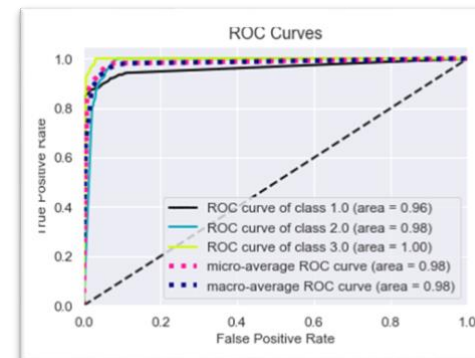Flowchart of decision tree classifier



## 7: Modeling Results and Conclusion

It shows the DT model's classification report. Here, the overall achieved F1-score is 93%. The individual F1-score is 96% for normal, 81% for suspected, and 93% for pathological. It displays the prediction of the DT model. The projected result is displayed in the confusion matrix, as well as the model's computed performance. The total number of correct predictions is 595, with 43 incorrect forecasts.

Classification Report



```
print(classification_report(y_test, dt_best.predict(X_test)))

              precision    recall  f1-score   support

         1.0       0.94      0.98      0.96       472
         2.0       0.92      0.72      0.81       105
         3.0       0.93      0.92      0.93        61

    accuracy                           0.93       638
   macro avg       0.93      0.87      0.90       638
weighted avg       0.93      0.93      0.93       638
```

ROC Curves



Confusion Matrix

```
evaluate_model(dt_best)

Train Accuracy : 0.9599785119527263
Train Confusion Matrix:
[[1169   61   11]
 [  34 1189   18]
 [   6   19 1216]]
------------------------------------------
Test Accuracy : 0.9227053140096618
Test Confusion Matrix:
[[361  46    7]
 [ 17 387   10]
 [  9   7 398]]
```

Type I Error: Null hypothesis mistakenly rejected or False Positive Error = 0.054
Type II Error: Null hypothesis mistakenly accepted or False Negative Error = 0.082