**Harshvardhan Thakur, 24231450, Computer Science Data Analytics**

## 1.  <u>Understanding the Dataset</u>

Given dataset is Perinatal Risk Information with some health parameters to classify women's maternal health using machine learning. Given data has Range Index of 1014 (0 to 1013) and total of 7 data columns with different datatypes as shown in Fig 1.1. Some columns such as Age, SystolicBP, DiastolicBP are independent variable and Type column is dependent variable and it is of categorical data, hence Classification model is preferable for this problem.





**Fig 1.1** Dataframe information.　　　　　　　**Fig. 2.1** Dependent Variable

## 2.  <u>Data Exploration and Pre-processing</u>

- Dataset is balanced since the dependent variable 'Type' has roughly same number of instances. The categorical values ranging between 26.82% to 40.03% as shown in Fig. 2.1.
- There are no missing values in the data set as shown in Fig. 2.2, thus imputation is not required for this dataset.
- On further exploring, total of 562 duplicated values are present in the dataset as shown in Fig. 2.3.





**Fig. 2.3** Duplicated Values

**Fig. 2.2** No Missing Values

- **Histplot:** sns.histplot() function to plot each attributes against 'Type' dependent variable to categorize the datapoints as shown in Fig. 2.4& 2.5 below, but from the graphs it is clearly visible datapoints are not easily separable into high, low, mid risk categories hence we use machine learning model that captures complexity such dataset.
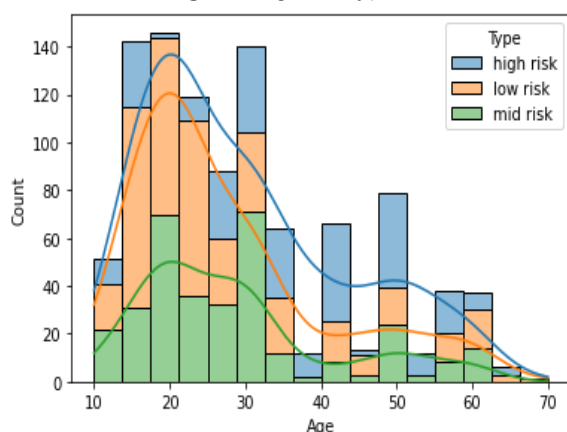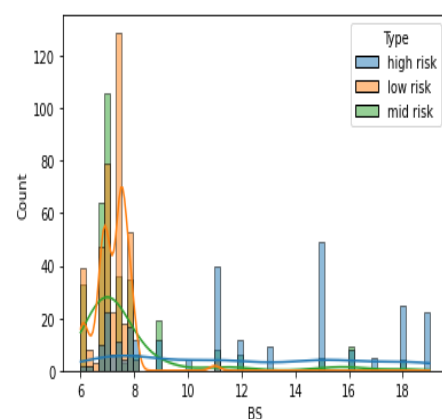
**Fig. 2.4** Age vs Type　　　　　　　　　　　**Fig. 2.5** BS vs Type

- Before Model Selection, we need to check for outliers as they can overfit the model performance. Many methods can be used to check for outliers such as IQR, line plots, z-score method etc. We made a custom function that takes each attribute and plots a line graph to check for outliers.

```python
#Checking outliers using line graphs

def custom_graph(df):
    for column in df.columns:
        plt.figure(figsize=(10,6))
        plt.plot(df.index ,df[column], label=f'{column} values', linewidth=1)


        plt.title(f'{column} Line Plot')
        plt.xlabel('Values')
        plt.ylabel(column)
        plt.legend()

        plt.show()
custom_graph(df)
```
[128]  ✓ 1.3s

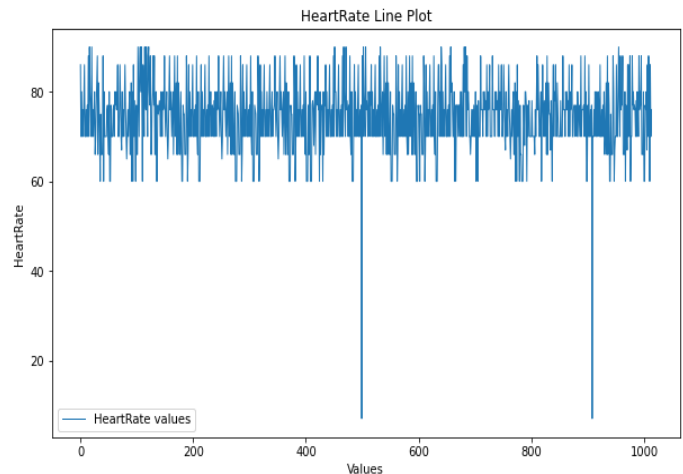**Fig. 2.6** Custom Outlier Function



**Fig. 2.7** Outlier Detection in HeartRate

- Although there's no significant change in model performance upon removing the outliers.
- **Correlation:** Correlation matrix of attributes to check which attributes are strongly correlated. SystolicBP and DiastolicBP have strong correlation coefficient. Although, matrix suggests that SystolicBP and DiastolicBP are strongly correlated and can be considered same in machine learning domain, however in medical domain both type of blood pressure is different and have some factors to them. Hence, we keep both the attributes.

```python
cor_matrix=df.corr()
print(cor_matrix)
#SystolicBP and DiastolicBP have strong correlation coefficient
```
[47]  ✓ 0.0s

```
                 Age  SystolicBP  DiastolicBP        BS  BodyTemp  HeartRate
Age         1.000000    0.416045     0.398026  0.473284 -0.255323   0.079798
SystolicBP  0.416045    1.000000     0.787006  0.425172 -0.286616  -0.023108
DiastolicBP 0.398026    0.787006     1.000000  0.423824 -0.257538  -0.046151
BS          0.473284    0.425172     0.423824  1.000000 -0.103493   0.142867
BodyTemp   -0.255323   -0.286616    -0.257538 -0.103493  1.000000   0.098771
HeartRate   0.079798   -0.023108    -0.046151  0.142867  0.098771   1.000000
```

**Fig. 2.8** Correlation Matrix.

## 3. <u>Algorithm Selection and Application</u>

- **Splitting the Data:** Before we feed dataset to the model, we split dataset into train, validation and test sets using train_test_split. First, data is split into test and temp and then further we split temp into train and validation data as shown in Fig. 3.1. Validation set prevents model from memorising the training set.

```python
from sklearn.model_selection import train_test_split
```
[53]  ✓ 0.0s

```python
X_temp, X_test, y_temp, y_test = train_test_split(X, y, test_size=.20, random_state=1)
```
[54]  ✓ 0.0s

```python
X_train, X_val, y_train, y_val= train_test_split(X_temp, y_temp, test_size=.20, random_state=1)
```
[55]  ✓ 0.0s

**Fig. 3.1** Train Test Split.

- The problem is categorical, classification model such as Decision Tree and Random Forest perform better in such problems. Decision Tree handles binary classification as well as multi-class classification since the dependent variable is multi-class, and Decision Tree does not require data transformation and it handles imbalance data as well. However, Decision Tree can sometimes overfit to overcome that issue Random Forest model is used

**Fig. 3.2** Decision Tree Implementation



**Fig. 3.3** Random Forest Implementation

- **Hyperparameter Tuning:** RandomizedSearchCV is one of the hyperparameter which selects samples from the combinations of hyperparameter compared to which GridSearchCV tests all the possible combinations of hyperparameter. An n_estimator defines values of decision trees in random forest ensemble; for moderate dataset size values like 100,200,300 can help generalization of model.
  Upon fitting the random_search model, it fits data 3 times (as cv=3) and iterates over 100; making 300 fits.



**Fig. 3.4** Hyperparameter Space



**Fig. 3.5** Best Parameters and Score for RandomSearch

- Training, Validation and Testing Accuracy comes out to be 0.924 and 0.834 and 0.835 respectively indicating a good fit and model performance for Random Forest.



**Fig. 3.6** Training and Validation Accuracy
Of Random Forest after Tuning.



**Fig. 4.1** Decision Tree Classification Report

```
best_rf = random_search.best_estimator_
y_pred_h3 = best_rf.predict(X_test)

# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred_h3)
print("Testing Accuracy: ", accuracy)
```
320]  ✓ 0.0s

··  Testing Accuracy:  0.8325123152709359

**Fig. 3.7** Testing Accuracy

## 4. Model Evaluation and Comparative Analysis

- **Decision Tree and Random Forest (Tuned) comparison;** Decision Tree Classification report can be seen in Fig. 4.1, it's performance on training dataset and validation dataset comes out be 0.92 and 0.80 respectively which can be seen in the submitted code file. Fig. 4.2 shows classification report of Random Forest after Hyperparameter Tuning. Thus, on comparing the scores of both model their performance scores come out to be almost similar as shown in bar plot Fig. 4.3

```
print(f"Random Forest Tuning Report:\n{classification_report(y_test, y_pred_h3)}")
✓ 0.0s

Random Forest Tuning Report:
              precision    recall  f1-score   support

   high risk       0.88      0.92      0.90        63
    low risk       0.84      0.83      0.84        77
    mid risk       0.77      0.75      0.76        63

    accuracy                           0.83       203
   macro avg       0.83      0.83      0.83       203
weighted avg       0.83      0.83      0.83       203
```



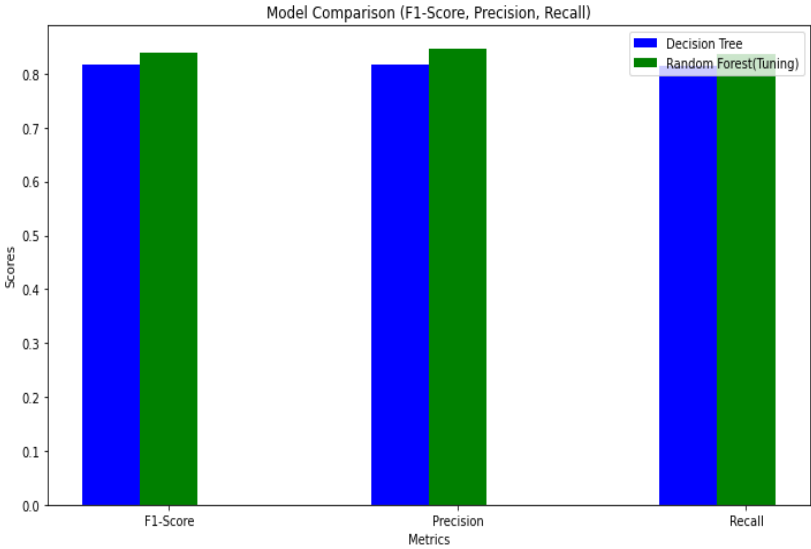**Fig. 4.2** Random Forest Hyperparameter Tuning

Classification Report

**Fig. 4.3** Model Comparison

## 5. Ethical Considerations

- For ML model to perform better, good dataset must be used with different attributes. We can't be sure whether this model will work for wider population as this model is trained on only 7 attributes. If such Machine Learning starts predicting wrong results; as consequences people might suffer major health issues.