



Correlation between columns are not between -1 to + 1 i.e. the best but not greater than 0.5 also and between kyphosis as target and other column is generally greater. so we are good.

ANOVA F measure

```
In [38]: import scipy.stats as stats
        f_stat, f_oneway, functions takes the groups as input and returns ANOVA F and p value
        Fvalue, pvalue = stats.f_oneway(df['Age'], df['Number'], df['Start'])
        print(fvalue, pvalue)

138.10440440436454 1.2198309671987674e-40

p value is e to the power 40 i.e. very small. So null hypothesis rejected. Variance of these columns are not same. So they are less correlated.
We can take all 3 to build model.
```

1. Data Preparation

a. Do the final feature selection and extract them into Column X and the class label into Column into Y.

```
In [39]: x=df.drop(['Kyphosis'],axis=1)
        y=df['Kyphosis']

1. Data Preparation
b. Split the dataset into training and test sets.
```

```
In [40]: from sklearn.model_selection import train_test_split
        x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.33,random_state=101)

Part B

1. Model Building

a. Perform Model Development using at least three models, separately. You are free to apply any Machine Learning Models on the dataset.
Deep Learning Models are strictly not allowed.
We will try Decision Tree as data is too small and number of column is also less preferably. And also try Random forest and Logistics
regression as all of them are Classification models.
We have general common classification techniques

1. Logistic Regression
2. Naive Bayes
3. Stochastic Gradient Descent
4. K-Nearest Neighbours
5. Decision Tree
6. Random Forest
7. Support Vector Machine

1. Model Building

Logistics Regression
```

```
In [41]: from sklearn.metrics import classification_report,confusion_matrix
        def viewConfusionMatrix(y_test, y_pred):
            print(confusion_matrix(y_test,y_pred))
            print(classification_report(y_test,y_pred))

In [42]: from sklearn.linear_model import LogisticRegression
        def logicalRegres(x,y,x_train,y_train,x_test,y_test):
            clf = LogisticRegression(random_state=0).fit(x, y)
            logreg = LogisticRegression()
            logreg.fit(x_train, y_train)
            y_pred = logreg.predict(x_test)
            print('Accuracy of logistic regression classifier on test set: {:.2f}'.format(logreg.score(x_test, y_test)))
            return y_pred

Now logical regression using all column
```

```
In [43]: y_pred = logicalRegres(x,y,x_train,y_train,x_test,y_test)

Accuracy of logistic regression classifier on test set: 0.74

viewConfusionMatrix(y_test, y_pred)

[[19  0]
 [ 7 11]]

      precision    recall  f1-score   support

   0       0.73       1.00       0.84        19
   1       1.00       0.12       0.22         8

 accuracy          macro avg          weighted avg
   0.87          0.87          0.56          0.53        27
   0.81          0.81          0.74          0.66        27
```

Logical regression not using "Start" column

Need to prepare other sets of x i.e. input as Start correlation is not very good and will compare them in same Logical Regression model. Their accuracy and performance.

```
In [45]: x_no_start=df.drop(['Kyphosis', 'Start'],axis=1)

In [46]: from sklearn.model_selection import train_test_split
        x_no_start_train,x_no_start_test,y_no_start_train,y_no_start_test=train_test_split(x_no_start,y,test_size=0.33,random_state=101)

In [47]: y_no_start_pred = logicalRegres(x_no_start,y,x_no_start_train,y_no_start_train,x_no_start_test,y_no_start_test)

Accuracy of logistic regression classifier on test set: 0.74

viewConfusionMatrix(y_no_start_test, y_no_start_pred)

[[19  0]
 [ 7 11]]

      precision    recall  f1-score   support

   0       0.73       1.00       0.84        19
   1       1.00       0.12       0.22         8

 accuracy          macro avg          weighted avg
   0.87          0.87          0.56          0.53        27
   0.81          0.81          0.74          0.66        27
```

So no change in accuracy though. SO we are good

DecisionTree

```
In [49]: from sklearn.tree import DecisionTreeClassifier
        dtree=DecisionTreeClassifier()
        #B. Train the model and print the training accuracy and loss values
        dtree.fit(x_train,y_train)
        #2. Performance Evaluation
        #B. Do the prediction for the test data and display the results for the inference
        y_pred=dtree.predict(x_test)
        #A. Print the confusion matrix. Provide appropriate analysis for the same.
```

a. Print the confusion matrix. Provide appropriate analysis for the same.

```
In [50]: viewConfusionMatrix(y_test, y_pred)

[[15  4]
 [ 6 21]]

      precision    recall  f1-score   support

   0       0.71       0.79       0.75        19
   1       0.33       0.25       0.29         8

 accuracy          macro avg          weighted avg
   0.52          0.52          0.52          0.52        27
   0.60          0.60          0.63          0.61        27
```

Random Forest

```
In [51]: from sklearn.ensemble import RandomForestClassifier
        rfc=RandomForestClassifier(n_estimators=20000)
        rfc.fit(x_train,y_train)
        y_pred=rfc.predict(x_test)

In [52]: viewConfusionMatrix(y_test, y_pred)

[[19  0]
 [ 6 21]]

      precision    recall  f1-score   support

   0       0.76       1.00       0.86        19
   1       1.00       0.25       0.40         8

 accuracy          macro avg          weighted avg
   0.88          0.88          0.62          0.63        27
   0.83          0.83          0.78          0.73        27
```

data is not very much imbalanced though but we can try SMOTE approach

```
In [53]: !pip install imbalanced-learn
import imblearn
print(imblearn.__version__)
from collections import Counter
from imblearn.over_sampling import SMOTE
from imblearn.under_sampling import RandomUnderSampler
from imblearn.pipeline import Pipeline
from sklearn.model_selection import import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
```

Requirement already satisfied: imbalanced-learn in /usr/local/lib/python3.7/dist-packages (0.4.3)
Requirement already satisfied: scikit-learn>=0.20 in /usr/local/lib/python3.7/dist-packages (from imbalanced-learn) (0.22.2.post1)
Requirement already satisfied: scikit-learn>=0.13.3 in /usr/local/lib/python3.7/dist-packages (from imbalanced-learn) (1.4.1)
Requirement already satisfied: numpy>=1.8.2 in /usr/local/lib/python3.7/dist-packages (from imbalanced-learn) (1.19.2)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.20-0.4.3) (1.0.1)

/usr/local/lib/python3.7/dist-packages/sklearn/externals/six.py:31: FutureWarning: The module is deprecated in version 0.21 and will be removed in version 0.23 since we've dropped support for Python 2.7. Please rely on the official version of six (https://pypi.org/project/six/).
"https://pypi.org/project/six/"). FutureWarning)
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:144: FutureWarning: The sklearn.neighbors.basemodule is deprecated in version 0.22 and will be removed in version 0.24. The corresponding classes / functions should instead be imported from sklearn.neighbors. Anything that cannot be imported from sklearn.neighbors is now part of the private API.
warnings.warn(message, FutureWarning)

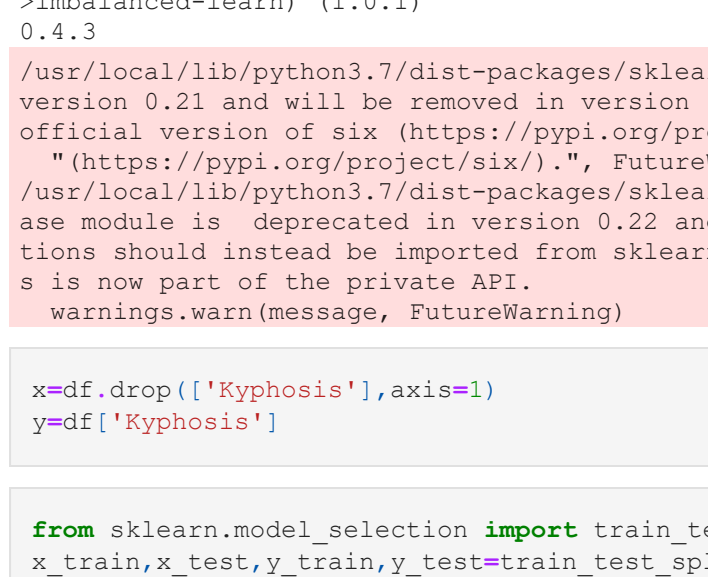
```
In [54]: x=df.drop(['Kyphosis'],axis=1)
        y=df['Kyphosis']

In [55]: from sklearn.model_selection import train_test_split
        x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.33,random_state=101)

In [56]: oversample = SMOTE()
        counter = Counter(y)
        X, y = oversample.fit_resample(x, y)
        # summarise the new class distribution
        counter = Counter(y)
        print(counter)
        # scatter plot of examples by class label
        for label, _ in counter.items():
            row_ix = np.where(y == label)[0]
            plt.scatter(X[row_ix, 0], X[row_ix, 1], label=str(label))
        plt.legend()
        plt.show()

Counter({0: 64, 1: 17})
Counter({0: 64, 1: 64})
```

/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function safe_indexing is deprecated; safe_indexing is deprecated in version 0.22 and will be removed in version 0.24.
warnings.warn(msg, category=FutureWarning)



Generalized Linear Model Table will help for linear regression and predict drop column. It is needed as between target i.e. kyphosis and one of the X i.e column "Start" correlation is not so good. And same between column "Start" and "Number".

Generalized Linear Models (GLM) estimate regression models for outcomes following exponential distributions. In addition to the Gaussian (i.e. normal) distribution, these include Poisson, binomial, and gamma distributions. Each serves a different purpose, and depending on distribution and link function choice, can be used either for prediction or classification.

```
In [37]: from google.colab import files
        df.to_csv("GlnData.csv", sep=',', index=False)
        files.download("GlnData.csv")
```

```
In [59]: !pip install h2o
import h2o
h2o.init()
from h2o.estimators.glm import H2OGeneralizedLinearEstimator
glm = h2o.import_file("https://raw.githubusercontent.com/ArindonyOfficial/vertebraeSurgery81akFactor/main/GlnData.csv")
print(glm.as_data_frame())
glm["Kyphosis"] = glm["Age"].asfactor()
glm["Age"] = glm["Age"].asfactor()
glm["Number"] = glm["Age"].asfactor()
glm["Start"] = glm["Age"].asfactor()
predictors = ["Age", "Number", "Start"]
response_col = "Kyphosis"

glm_model = H2OGeneralizedLinearEstimator(family="AUTO")
glm_model.train(predictors, response_col, training_frame= glm)

# Coefficients that can be applied to the non-standardized data.
print(glm_model.coef())

# Coefficients fitted on the standardized data (requires standardize = True, which is on by default)
print(glm_model.coef_norm())

# Print the Coefficients table
print(glm_model.model_json['output']['coefficients_table'])
#print(glm_model.model_json['output']['std_error'])

# Print the Standard error
#print(glm_model.model_json['output']['coefficients_table']['std_error'])

# Print the p values
#print(glm_model.model_json['output']['coefficients_table']['p_value'])

# Print the z values
#print(glm_model.model_json['output']['coefficients_table']['z_value'])

# Retrieve a graphical plot of the standardized coefficient magnitudes
glm_model.std_coef_plot()
```


[illegible]

[illegible]

[illegible]

```
# evalua
```

```
scores = cross_val_score(pipeline, X, y, scoring='roc_auc', cv=cv, n_jobs=-1)
print('Mean ROC AUC: ', scores.mean())
```

Positive Position 0535