

1. Name - Arindam Dev
2. Group - 307 Group Members - ARINDAM DEY - 2020f0425@wlip.bits-pilani.ac.in KAUSHIK DUBEY - 2020f0424@wlip.bits-pilani.ac.in MOHAMMAD ATTAULLAH - 2020F04274@wlip.bits-pilani.ac.in
3. Dataset Name - kyphosis.csv <https://github.com/arindamdey/official/vertebraeSurgeryRiskFactor/blob/main/kyphosis.csv>

RAW-

<https://raw.githubusercontent.com/arindamdey/official/vertebraeSurgeryRiskFactor/7e5ba322ac93a563528e286894521d35f7d4995/kyphosis.csv>

1. Import Libraries/Dataset
- b. Import the required libraries

```
import pandas as pd
#Data normalisation view
import numpy as np
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import plotly
import seaborn as sns
```

```
df = pd.read_csv('https://raw.githubusercontent.com/arindamdey/official/vertebraeSurgeryRiskFactor/7e5ba322ac93a563528e286894521d35f7d4995/kyphosis.csv')
```

1. Data Visualization and Exploration

- a. Print at least 5 rows for sanity check to identify all the features present in the dataset and if the target matches with them.
- In [30]:

```
df.head()
```

	Kyphosis	Age	Number	Start
0	absent	71	3	5
1	absent	71	3	14
2	present	128	4	5
3	absent	2	5	1
4	absent	1	4	15

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 81 entries, 0 to 80
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Kyphosis    81 non-null    object
 1   Age         81 non-null    int64
 2   Number     81 non-null    int64
 3   Start      81 non-null    int64
dtypes: int64(3), object(1)
memory usage: 2.7+ KB
```

```
df.columns
```

```
Index(['Kyphosis', 'Age', 'Number', 'Start'], dtype='object')
```

1. Data Visualization and Exploration

- b. Print the description and shape of the dataset.
- In [6]:

```
df.describe()
```

	Age	Number	Start
count	81.000000	81.000000	81.000000
mean	63.654321	4.048963	11.493827
std	58.104251	1.619423	4.883962
min	1.000000	2.000000	1.000000
25%	26.000000	3.000000	9.000000
50%	67.000000	4.000000	13.000000
75%	130.000000	5.000000	16.000000
max	206.000000	10.000000	18.000000

```
df.shape
```

```
(81, 4)
```

1. Data Visualization and Exploration

- c. Provide appropriate visualization to get an insight about the dataset.
- In [8]:

```
#Number==> number of levels involved
#Start ==> starting vertebrae level of the surgery

n_rows=5
n_cols=2
width=20
height=30

fig,ax=plt.subplots(n_rows,n_cols,sharex=False,sharey=False,figsize=(width,height))

#Normalisation view
wk=df[df['Kyphosis']=='present']

#Age
g=sns.countplot(df["Age"], ax=0,0)
ax[0,0].set_title("Age")

g=sns.distplot(wk("Age"), ax=0,1)
ax[0,1].set_title("Age")

#Number
g=sns.countplot(df["Number"], ax=1,0)
ax[1,0].set_title("Number of Level")

g=sns.distplot(wk("Number"), ax=1,1)
ax[1,1].set_title("Number of Level")

#Start
g=sns.countplot(df["Start"], ax=2,0)
ax[2,0].set_title("Starting Vertebrae")

g=sns.distplot(wk("Start"), ax=2,1)
ax[2,1].set_title("Starting Vertebrae")

#Count Kyphosis
l=len(df[df["Kyphosis"]=="present"])
g=sns.countplot(l,len(df)-1, ax=3,0)
ax[3,0].set_title("With Kyphosis")

#Count Kyphosis for Age vs Number of Level
#Having Kyphosis> Age vs Number Lineplot Age in X axis
wk=df[df["Kyphosis"]=="present"]
g=sns.lineplot(data=wk, x="Age", y="Number", ax=3,1)
ax[3,1].set_title("Kyphosis for Age vs Number of Level")

#Having Kyphosis Number vs Age Countplot Number in X axis
g=sns.lineplot(data=wk, x="Number", y="Age", ax=4,0)
ax[4,0].set_title("Kyphosis for Number of Level vs Age")

#Having Kyphosis Age vs Starting Vertebrae
g=sns.lineplot(data=wk, x="Age", y="Start", ax=4,1)
ax[4,1].set_title("Kyphosis Age vs Starting Vertebrae")

plt.show()
```

/usr/local/lib/python3.7/dist-packages/seaborn/decorators.py:43: FutureWarning: Pass the following variable as a keyword arg x: From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:255: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

/usr/local/lib/python3.7/dist-packages/seaborn/decorators.py:43: FutureWarning: Pass the following variable as a keyword arg x: From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:255: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

/usr/local/lib/python3.7/dist-packages/seaborn/decorators.py:43: FutureWarning: Pass the following variable as a keyword arg x: From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:255: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

/usr/local/lib/python3.7/dist-packages/seaborn/decorators.py:43: FutureWarning: Pass the following variable as a keyword arg x: From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning



```
sns.pairplot(df,hue='Kyphosis',palette='Set1')
```

```
<seaborn.axisgrid.PairGrid at 0x7fe064be3a10>
```

```
df.bi_variate(df, col1, col2):
    cross_name=pd.crosstab(df[col1], df[col2], normalize='index')
    cross_name.plot.bar()
    plt.xlabel('{}'.format(col1))
    plt.ylabel('{}'.format(col2))
    plt.xticks(rotation=0)
    df.figure(figsize=(16,12))
    plt.show()
    return cross_name*100
```

```
bi_variate(df, 'Start', 'Number')
```

<Figure size 1152x864 with 0 Axes>

```
Start
```

	1	2	3	4	5	6	7	9	10
1	0.000000	0.000000	60.000000	40.000000	0.000000	0.000000	0.000000	0.000000	0.0
2	0.000000	0.000000	0.000000	50.000000	0.000000	50.000000	0.000000	0.000000	0.0
3	33.333333	33.333333	33.333333	0.000000	0.000000	0.000000	33.333333	33.333333	0.0
4	0.000000	25.000000	0.000000	25.000000	0.000000	25.000000	0.000000	0.000000	25.0
5	0.000000	0.000000	0.000000	50.000000	0.000000	0.000000	0.000000	0.000000	0.0
6	0.000000	25.000000	25.000000	50.000000	25.000000	25.000000	0.000000	0.000000	0.0
7	0.000000	25.000000	50.000000	25.000000	0.000000	0.000000	0.000000	0.000000	0.0
9	33.333333	0.000000	33.333333	33.333333	0.000000	0.000000	0.000000	0.000000	0.0
10	0.000000	20.000000	40.000000	20.000000	20.000000	0.000000	0.000000	0.000000	0.0
12	8.333333	33.333333	16.666667	33.333333	0.000000	8.333333	0.000000	0.000000	0.0
14	20.000000	40.000000	0.000000	40.000000	0.000000	0.000000	0.000000	0.000000	0.0
15	0.000000	14.285714	57.142857	14.285714	14.285714	14.285714	0.000000	0.000000	0.0
16	25.529412	52.941176	17.647059	11.764706	0.000000	0.000000	0.000000	0.000000	0.0
17	100.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0
18	0.000000	100.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0

- d. Try exploring the data and see what insights can be drawn from the dataset.
- Exploratory data analysis

1. Data Pre-processing and cleaning

- a. Do the appropriate preprocessing of the data like identifying NULL or Missing Values if any, handling of outliers if present in the dataset, skewed data etc. Apply appropriate feature engineering techniques for them.
- Now first check should be any Null is present or not

```
df.isna().any()
```

```
Kyphosis    False
Age         False
Number      False
Start       False
dtype: bool
```

```
df.isna().sum()
```

```
Kyphosis    0
Age         0
Number      0
Start       0
dtype: int64
```

Now there are no Null values present

```
!pip install missingno
import missingno
missingno.matrix(df,figsize=(12,8))
```

Requirement already satisfied: missingno in /usr/local/lib/python3.7/dist-packages (0.5.0)

Requirement already satisfied: scipy in /usr/local/lib/python3.7/dist-packages (from missingno) (1.4.1)

Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from missingno) (1.19.5)

Requirement already satisfied: matplotlib in /usr/local/lib/python3.7/dist-packages (from missingno) (3.2.2)

Requirement already satisfied: seaborn in /usr/local/lib/python3.7/dist-packages (from missingno) (0.11.1)

Requirement already satisfied: python-dateutil<=2.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=missingno) (2.8.2)

Requirement already satisfied: pytz in /usr/local/lib/python3.7/dist-packages (from python-dateutil<=2.1) (2.0.1)

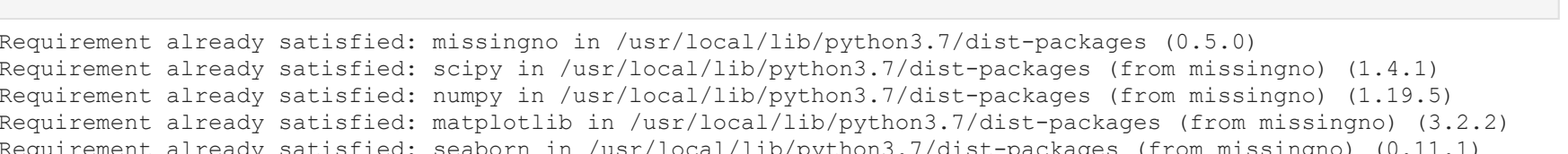
Requirement already satisfied: kiwisolver<=1.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=missingno) (1.3.1)

Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from cycler>=0.10 to matplotlib>=missingno) (1.15.0)

Requirement already satisfied: pandas>=2017.2 in /usr/local/lib/python3.7/dist-packages (from pandas>=0.23 to seaborn>=2018.9) (2018.9)

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe061dcaad0>
```

1



81

4 4

Outlier identification and removal

Box plot use the Inter Quantile Range(IQR) method. This is easier so will be using

```
df.boxplot()
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe057ae50d0>
```

```
df.boxplot(column='Number')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe0529716d0>
```

```
df.plot.scatter('Age', 'Number')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe05295ad50>
```

Z-score is the signed number of standard deviations by which the value of an observation or data point is above the mean value of what is being observed or measured Through this method also we can identify Outlier

```
mean = np.mean(df["Number"])
mean
```

```
4.049382716049383
```

```
std = np.std(df["Number"])
std
```

```
1.6093955200472556
```

```
threshold = 2
outlierCombo = []
outlier = []
score = []
for i in df["Number"]:
```

```
    z = (i-mean)/std
    if z > threshold:
```

```
        outlierCombo.append(i)
        outlier.append(i)
        score.append(i)
    outlierCombo
```

```
Out[118]:
```

```
[9, 3.076072489505411], [10, 3.697423790378075]]
```

```
df[df["Number"]>=9]
```

```
Kyphosis Age Number Start
42 absent 143 9 3
52 present 139 10 6
```

This is really unrealistic. Age people live according to age also and as per scatter plot and ZScore. But data is skewed toward 100+ years of age and data quantity also less so preferable not drop but replace in 3 sigma method the 2 outliers.

```
leftVal = mean - 3*std
rightVal = mean + 3*std
print(mean, std, leftVal, rightVal)
outli = df[df["Number"]<leftVal]
outliR = df[df["Number"]>rightVal]
print(outli,outliR)
print("threshold value counts", df[df["Number"] == rightVal]["Number"].count())
print(outliR)
```

```
4.049382716049383 1.6093955200472556 -0.7788038440923835 8.87756927619115
0 2
threshold value counts: 0
Kyphosis Age Number Start
42 absent 143 9 3
52 present 139 10 6
```

```
df.loc[df["Number"] > rightVal, "Number"] = rightVal
```

```
df[df["Number"] > rightVal]
```

```
Kyphosis Age Number Start
0 0 71 3.0 5
1 0 158 3.0 14
2 1 128 4.0 5
3 0 2 5.0 1
4 0 1 4.0 15
5 0 1 2.0 16
6 0 61 2.0 17
7 0 37 3.0 16
8 0 113 2.0 16
9 1 59 6.0 12
```

1. Data Pre-processing and cleaning

- c. Do the correlational analysis on the dataset. Provide a visualization for the same.
- Feature Selection

We can do by correlation, mostly continuous data as input and categorical data as output. SO we can use below methods.

1. ANOVA correlation coefficient (linear)
2. Kendall's rank coefficient (nonlinear)

kendall

```
#Correlation
#Correlation analysis shows us how to determine both the nature and strength of relationship between two variab
#Correlation lies between -1 to 1 (0: No correlation; -1: perfect negative correlation; +1: positive core)
#Correlation= df.corr(method='kendall')
correlation
```

```
#Plotting correlation Matrix
plt.figure(figsize=(20,15))
cm = sns.light_palette("green", as_cmap=True)
sns.heatmap(correlation, xticklabels=correlation.columns.values, yticklabels=correlation.columns.values, cmap=cm,
            plt.xticks(rotation=90))
```

```
(array([0.5, 1.5, 2.5, 3.5]), <a list of 4 Text major ticklabel objects>)
```

```
array([[0.5, 0.1, 0.3, 0.39],
       [0.1, 1, 0.025, 0.013],
       [0.3, 0.025, 1, 0.38],
       [0.39, 0.013, 0.38, 1]])
```

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

ANOVA F measure

```
# ANOVA feature selection for numeric input and categorical output
from sklearn.datasets import load_iris
from sklearn.feature_selection import SelectFromModel
from sklearn.feature_selection import f_classif
# generate dataset
X, y = load_iris(return_X_y=True)
# define feature selection
fs = SelectFromModel(f_classif, k=2)
# apply feature selection
X_selected = fs.fit_transform(X, y)
print(X_selected.shape)
```

```
(100, 2)
```

1. Data Preparation

- a. Do the final feature selection and extract them into Column X and the class label into Column into Y.
- In [1]:

```
X=df.drop(['Kyphosis','Start'],axis=1)
Y=df['Kyphosis']
```

1. Data Preparation

- b. Split the dataset into training and test sets.
- In [1]:

```
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 Logistic regression as all of them are Classification models.

We have general common classification techniques

1. Logistic Regression

2. Naive Bayes

3. Gradient Descent

4. K-Nearest Neighbors

5. Decision Tree

6. Random Forest

7. Support Vector Machine

1. Model Building

DecisionTree

```
from sklearn.tree import DecisionTreeClassifier
dtree=DecisionTreeClassifier()
```

- b. Train the model and print the training accuracy and loss values
- In [1]:

```
dtree.fit(X_train,Y_train)
```

```
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                        max_depth=None, max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort='deprecated',
                        random_state=None, splitter='best')
```

1. Performance Evaluation

- b. Do the prediction for the test data and display the results for the inference
- In [1]:

```
prediction=dtree.predict(X_test)
```



```
In [ ]:
from sklearn.metrics import classification_report, confusion_matrix

In [ ]:
print(confusion_matrix(y_test, prediction))

print(classification_report(y_test, prediction))

[[14  5]
 [ 6 21]]

      precision    recall  f1-score   support

    0       0.70      0.74      0.72         19
    1       0.29      0.25      0.27         8

 accuracy          0.49      0.49      0.49         27
 macro avg          0.58      0.59      0.58         27
 weighted avg          0.58      0.59      0.58         27

RandomForest

In [ ]:
from sklearn.ensemble import RandomForestClassifier

In [ ]:
rfc=RandomForestClassifier(n_estimators=20000)

In [ ]:
rfc.fit(x_train, y_train)

Out [ ]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
criterion='gini', max_depth=None, max_features='auto',
max_leaf_nodes=None, max_samples=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=20000,
n_jobs=None, oob_score=False, random_state=None,
verbose=0, warm_start=False)

In [ ]:
rfc_pred=rfc.predict(x_test)

In [ ]:
print(confusion_matrix(y_test, rfc_pred))

print(classification_report(y_test, rfc_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          0.88      0.62      0.63         27
 macro avg          0.83      0.78      0.73         27
 weighted avg          0.83      0.78      0.73         27

Logistics Regression

In [ ]:
from sklearn.linear_model import LogisticRegression

In [ ]:
clf = LogisticRegression(random_state=0).fit(x, y)

In [ ]:
logreg = LogisticRegression()

In [ ]:
logreg.fit(x_train, y_train)

Out [ ]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, l1_ratio=None, max_iter=100,
multi_class='auto', n_jobs=None, penalty='l2',
random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
warm_start=False)

In [ ]:
y_pred = logreg.predict(x_test)

In [ ]:
print('Accuracy of logistic regression classifier on test set: {:.2f}'.format(logreg.score(x_test, y_test)))

Accuracy of logistic regression classifier on test set: 0.85

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

print(classification_report(y_test, y_pred))

[[22  0]
 [ 4 21]]

      precision    recall  f1-score   support

    0       0.85      1.00      0.92         22
    1       1.00      0.20      0.33         5

 accuracy          0.92      0.60      0.62         27
 macro avg          0.92      0.60      0.62         27
 weighted avg          0.87      0.85      0.81         27
```