

```
[1]: import numpy as np
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline

In [2]: df=pd.read_csv("fetal_health.csv")

In [3]: df.head()

Out[3]:
baseline      accelerations      fetal_movement      uterine_contractions      light_decelerations      severe_decelerations      prolonged_decelerations      abnormal_short_term_variability      mean_value_of_short_term_variability      percentage_of_time_with_abnormal_long_term_variability
0      125.0      0.000      0.0      0.000      0.000      0.0      0.0      32.0      0.5
1      132.0      0.006      0.0      0.006      0.003      0.0      0.0      17.0      2.1
2      133.0      0.003      0.0      0.006      0.003      0.0      0.0      16.0      2.1
3      134.0      0.003      0.0      0.006      0.003      0.0      0.0      36.0      2.4
4      132.0      0.007      0.0      0.006      0.000      0.0      0.0      16.0      2.4

5 rows x 22 columns

In [4]: df.isnull().sum()

Out[4]:
baseline value      0
accelerations      0
fetal_movement     0
uterine_contractions      0
light_decelerations      0
severe_decelerations      0
prolonged_decelerations      0
abnormal_short_term_variability      0
mean_value_of_short_term_variability      0
percentage_of_time_with_abnormal_long_term_variability      0
histogram_width      0
histogram_min      0
histogram_max      0
histogram_number_of_peaks      0
histogram_number_of_zeros      0
histogram_mode      0
histogram_mean      0
histogram_median      0
histogram_variance      0
histogram_tendency      0
fetal_health      0
dtype: int64

In [5]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2126 entries, 0 to 2125
Data columns (total 22 columns):
#   Column      Non-Null Count  Dtype
---  --
0   baseline value      2126 non-null    float64
1   accelerations      2126 non-null    float64
2   fetal_movement     2126 non-null    float64
3   uterine_contractions      2126 non-null    float64
4   light_decelerations      2126 non-null    float64
5   severe_decelerations      2126 non-null    float64
6   prolonged_decelerations      2126 non-null    float64
7   abnormal_short_term_variability      2126 non-null    float64
8   mean_value_of_short_term_variability      2126 non-null    float64
9   percentage_of_time_with_abnormal_long_term_variability      2126 non-null    float64
10  mean_value_of_long_term_variability      2126 non-null    float64
11  histogram_width      2126 non-null    float64
12  histogram_min      2126 non-null    float64
13  histogram_max      2126 non-null    float64
14  histogram_number_of_peaks      2126 non-null    float64
15  histogram_number_of_zeros      2126 non-null    float64
16  histogram_mode      2126 non-null    float64
17  histogram_mean      2126 non-null    float64
18  histogram_median      2126 non-null    float64
19  histogram_variance      2126 non-null    float64
20  histogram_tendency      2126 non-null    float64
21  fetal_health      2126 non-null    float64
dtypes: float64(22)
memory usage: 365.5 KB

In [6]: df.describe()

Out[6]:
baseline      accelerations      fetal_movement      uterine_contractions      light_decelerations      severe_decelerations      prolonged_decelerations      abnormal_short_term_variability      mean_value_of_short_term_variability      percentage_of_time_with_abnormal_long_term_variability
count      2126.000000      2126.000000      2126.000000      2126.000000      2126.000000      2126.000000      2126.000000      2126.000000      2126.000000
mean      133.303857      0.003178      0.009481      0.004366      0.001889      0.000003      0.000159      46.890122      1.332785
std      9.840844      0.003866      0.046666      0.002946      0.002960      0.000057      0.000590      17.192814      0.883241
min      106.003050      0.000000      0.000000      0.000000      0.000000      0.000000      0.000000      12.000000      0.200000
25%      126.000000      0.000000      0.000000      0.000000      0.000000      0.000000      0.000000      32.000000      0.700000
50%      133.000000      0.002000      0.000000      0.004000      0.000000      0.000000      0.000000      49.000000      1.200000
75%      140.000000      0.008000      0.003000      0.007000      0.003000      0.000000      0.000000      61.000000      1.700000
max      160.000000      0.019000      0.481000      0.015000      0.010000      0.005000      0.005000      87.000000      7.000000

8 rows x 22 columns

In [7]: X=df.loc[:,df.columns != 'fetal_health']
y=df["fetal_health"]

In [8]: from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X,y, test_size=0.20, random_state=42)

In [9]: X_train.head()

Out[9]:
baseline      accelerations      fetal_movement      uterine_contractions      light_decelerations      severe_decelerations      prolonged_decelerations      abnormal_short_term_variability      mean_value_of_short_term_variability      percentage_of_time_with_abnormal_long_term_variability
1223      125.0      0.000      0.0      0.008      0.000      0.0      0.0      32.0      1.1
480      140.0      0.000      0.0      0.001      0.000      0.0      0.0      60.0      0.8
1111      122.0      0.000      0.0      0.000      0.001      0.0      0.0      20.0      1.8
1303      137.0      0.005      0.0      0.005      0.002      0.0      0.0      36.0      0.9
861      142.0      0.003      0.0      0.004      0.000      0.0      0.0      46.0      0.7

5 rows x 21 columns

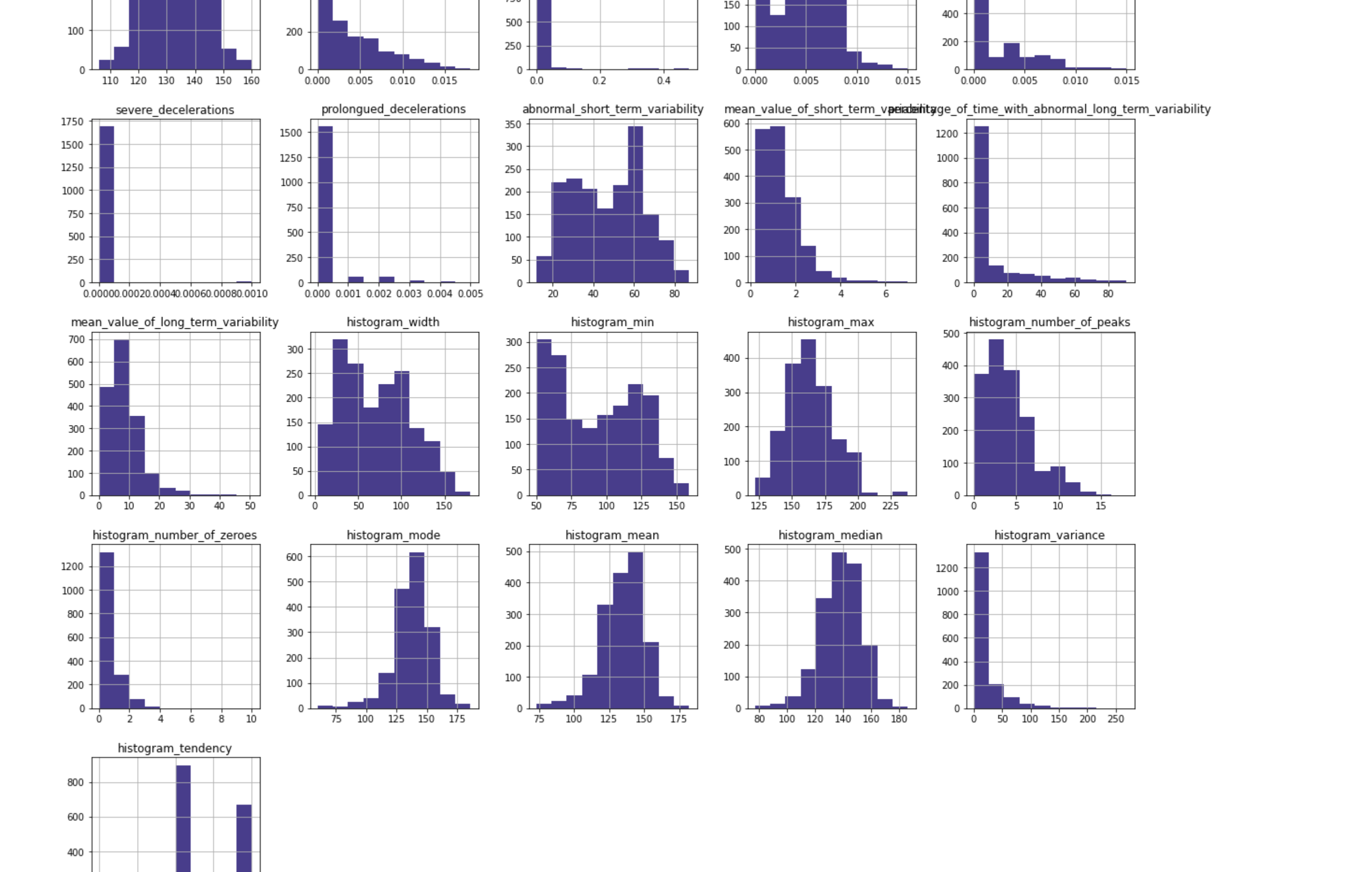
In [10]: sns.countplot(x=y_train)

Out[10]:
<matplotlib.axes._subplots.AxesSubplot at 0x7f6ad61915b9>


```

EDA(Exploratory data analysis)

```
[12]: # Evaluating distributions of the features
hist_plot = X_train.hist(figsize = (20,28), color = "#4683B8")
```



It is a unbalanced data set

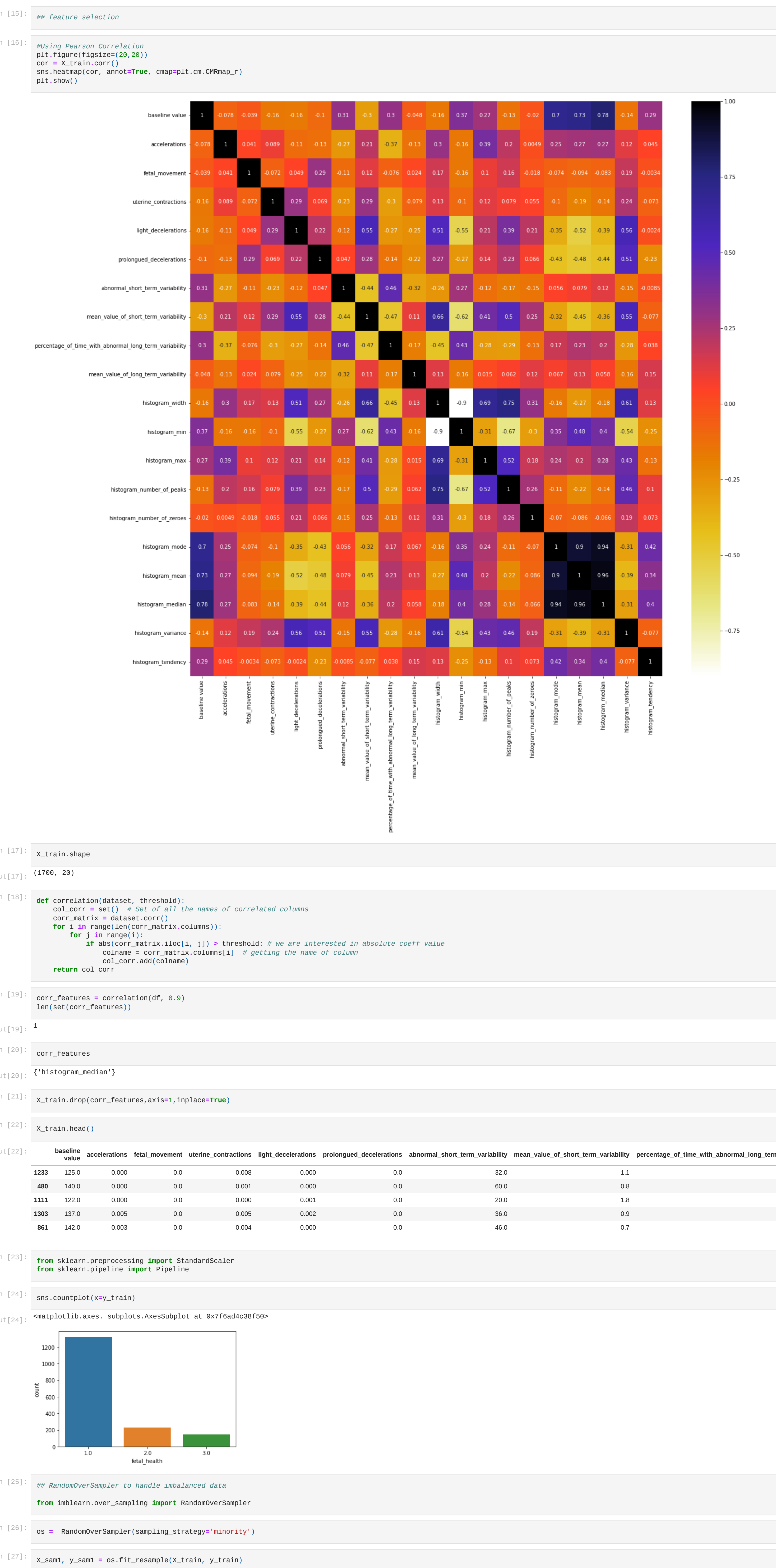
```
In [13]: X_train.drop("severe_decelerations", axis=1, inplace=True)
```

```
In [14]: X_train.head()
```

	baseline value	accelerations	fetal_movement	uterine_contractions	light_decelerations	prolonged_decelerations	abnormal_short_term_variability	mean_value_of_short_term_variability	percentage_of_time_with_abnormal_long_term_variability
1223	125.0	0.000	0.0	0.008	0.000	0.0	0.0	32.0	1.1
480	140.0	0.000	0.0	0.001	0.000	0.0	0.0	60.0	0.8
1111	122.0	0.000	0.0	0.000	0.001	0.0	0.0	20.0	1.8
1303	137.0	0.005	0.0	0.005	0.002	0.0	0.0	36.0	0.9
861	142.0	0.003	0.0	0.004	0.000	0.0	0.0	46.0	0.7

```
In [15]: ## Feature Selection
```

```
In [16]: #Using Pearson Correlation
plt.figure(figsize=(20,28))
cor = X_train.corr()
sns.heatmap(corr, annot=True, cmap=plt.cm.CmRmap_r)
plt.show()
```



```
In [17]: X_train.shape
```

(1780, 28)

```
In [18]: def correlation_matrix(threshold):
col_corr = set() # set of all the names of correlated columns
corr_matrix = dataset.corr()
for i in range(len(corr_matrix.columns)):
    for j in range(i):
        if abs(corr_matrix.loc[i, j]) > threshold: # we are interested in absolute coeff value
            colname = corr_matrix.columns[i] # getting the name of column
            col_corr.add(colname)
return col_corr
```

```
In [19]: corr_features = correlation(df, 0.9)
len(set(corr_features))
```

1

```
In [20]: corr_features
```

['histogram\_median']

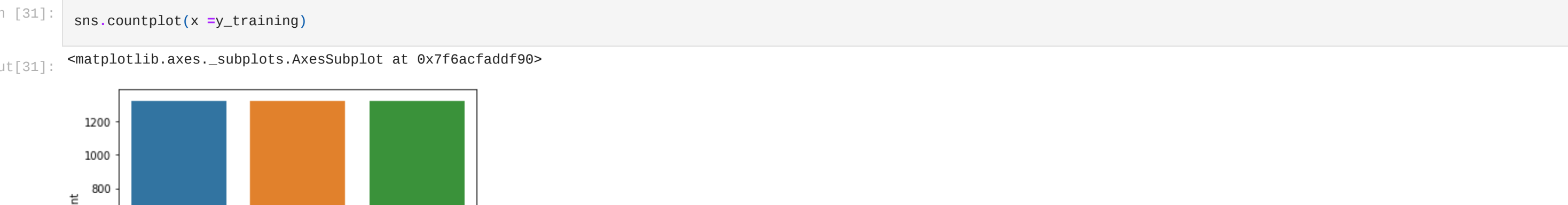
```
In [21]: X_train.drop(corr_features,axis=1,inplace=True)
```

```
In [22]: X_train.head()
```

	baseline value	accelerations	fetal_movement	uterine_contractions	light_decelerations	prolonged_decelerations	abnormal_short_term_variability	mean_value_of_short_term_variability	percentage_of_time_with_abnormal_long_term_variability
1223	125.0	0.000	0.0	0.008	0.000	0.0	0.0	32.0	1.1
480	140.0	0.000	0.0	0.001	0.000	0.0	0.0	60.0	0.8
1111	122.0	0.000	0.0	0.000	0.001	0.0	0.0	20.0	1.8
1303	137.0	0.005	0.0	0.005	0.002	0.0	0.0	36.0	0.9
861	142.0	0.003	0.0	0.004	0.000	0.0	0.0	46.0	0.7

```
In [23]: from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
```

```
In [24]: sns.countplot(x=y_train)
```



```
In [25]: ## RandomOverSampler to handle imbalanced data
from imblearn.over_sampling import RandomOverSampler
```

```
In [26]: os = RandomOverSampler(sampling_strategy='minority')
```

```
In [27]: X_sam1, y_sam1 = os.fit_resample(X_train, y_train)
```

```
In [28]: from collections import Counter
print('Original dataset shape {}'.format(Counter(y_train)))
print('Resampled dataset shape {}'.format(Counter(y_sam1)))
```

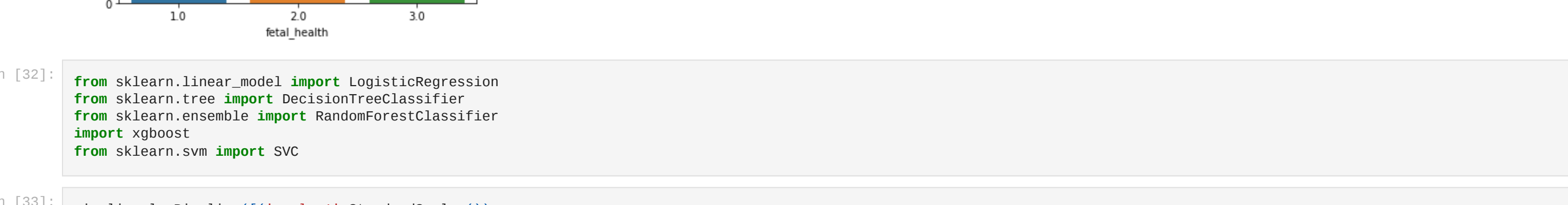
Original dataset shape Counter{(1.0: 1322, 2.0: 231, 3.0: 147)}  
Resampled dataset shape Counter{(1.0: 1322, 2.0: 1322, 3.0: 1322)}

```
In [29]: X_training, y_training = os.fit_resample(X_sam1, y_sam1)
```

```
In [30]: from collections import Counter
print('Original dataset shape {}'.format(Counter(y_train)))
print('Resampled dataset shape {}'.format(Counter(y_training)))
```

Original dataset shape Counter{(1.0: 1322, 2.0: 231, 3.0: 147)}  
Resampled dataset shape Counter{(1.0: 1322, 2.0: 1322, 3.0: 1322)}

```
In [31]: sns.countplot(x=y_training)
```



```
In [32]: from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
import xgboost
from sklearn.svm import SVC
```

```
In [33]: pipeline_lr=Pipeline([('scalar1',StandardScaler()),
                           ('lr_classifier',LogisticRegression())])
```

```
In [34]: pipeline_dt=Pipeline([('scalar2',StandardScaler()),
                           ('dt_classifier',DecisionTreeClassifier())])
```

```
In [35]: pipeline_randomforest=Pipeline([('scalar3',StandardScaler()),
                                      ('rf_classifier',RandomForestClassifier())])
```

```
In [36]: pipeline_xg=Pipeline([('scalar4',StandardScaler()),
                            ('xg_classifier',xgboost.XGBClassifier())
                            ])
```

```
In [37]: pipeline_svc=Pipeline([('scalar5',StandardScaler()),
                             ('svc_classifier',SVC())])
```

```
In [38]: pipelines = [pipeline_lr, pipeline_dt, pipeline_randomforest, pipeline_xg, pipeline_svc]
```

```
In [39]: best_accuracy=0
best_classifier=""
best_pipeline=""
```

```
In [40]: # Dictionary of pipelines and classifier types for ease of reference
pipe_dict = {0: 'Logistic Regression', 1: 'Decision Tree', 2: 'RandomForest', 3:'XG boost',4:'SVC'}
```

```
In [41]: # Fit the pipelines
for pipe in pipelines:
    pipe.fit(X_training, y_training)
```

```
In [42]: X_test.drop(["histogram_median","severe_decelerations"],axis=1,inplace=True)
```

```
In [43]: # For i model in enumerate(pipelines):
print("Test Accuracy: {}".format(pipe_dict[i],model.score(X_test,y_test)))
```

Logistic Regression Test Accuracy: 0.828638497651921  
Decision Tree Test Accuracy: 0.9201877934723  
RandomForest Test Accuracy: 0.940893899713615  
XG boost Test Accuracy: 0.9201877934723  
SVC Test Accuracy: 0.87586854468939

```
In [43]: for i, model in enumerate(pipelines):
if model.score(X_test,y_test)>best_accuracy:
    best_accuracy=model.score(X_test,y_test)
    best_pipeline=model
    best_classifier=i
print('Classifier with best accuracy:{}'.format(pipe_dict[best_classifier]))
```

Classifier with best accuracy:RandomForest

```
In [44]: from sklearn.metrics import Classification_report
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
```

```
In [45]: for i, model in enumerate(pipelines):
y_pred=model.predict(X_test)
print("{} Accuracy Score {}".format(pipe_dict[i]), format(accuracy_score(y_test,y_pred)))
print("{} Classification report: {}".format(pipe_dict[i]), format(classification_report(y_test,y_pred)))
print("{} Confusion matrix: {}".format(pipe_dict[i]), format(confusion_matrix(y_test,y_pred)))
print("{} Confusion matrix: {}".format(pipe_dict[i]), format(confusion_matrix(y_test,y_pred)))
```

Logistic Regression:Accuracy Score  
0.828638497651921  
Logistic Regression:Classification report:  
precision recall f1-score support  
1.0 0.98 0.85 0.91 333  
2.0 0.62 0.88 0.72 64  
3.0 0.50 0.79 0.61 29  
accuracy 0.83 426  
macro avg 0.67 0.88 0.71 426  
weighted avg 0.88 0.83 0.84 426

Decision Tree:Confusion matrix:  
[[282 40 11]  
 [ 57 0 1]  
 [ 1 9 21]  
 Decision Tree:Accuracy Score  
0.9201877934723  
Decision Tree:Classification report:  
precision recall f1-score support  
1.0 0.95 0.95 0.95 333  
2.0 0.74 0.77 0.75 64  
3.0 0.93 0.97 0.95 29  
accuracy 0.92 426  
macro avg 0.88 0.88 0.88 426  
weighted avg 0.92 0.92 0.92 426

RandomForest:Confusion matrix:  
[[321 21 1]  
 [ 54 21 1]  
 [ 1 0 81]  
 RandomForest:Accuracy Score  
0.940893899713615  
RandomForest:Classification report:  
precision recall f1-score support  
1.0 0.97 0.96 0.97 333  
2.0 0.83 0.84 0.84 64  
3.0 0.80 0.87 0.83 29  
accuracy 0.95 426  
macro avg 0.90 0.92 0.91 426  
weighted avg 0.95 0.95 0.95 426

XG boost:Confusion matrix:  
[[289 23 8]  
 [ 56 31 1]  
 [ 0 2 77]  
 XG boost:Accuracy Score  
0.87586854468939  
XG boost:Classification report:  
precision recall f1-score support  
1.0 0.99 0.88 0.93 333  
2.0 0.57 0.84 0.68 64  
3.0 0.75 0.93 0.83 29  
accuracy 0.77 426  
macro avg 0.77 0.88 0.81 426  
weighted avg 0.91 0.88 0.89 426

SVC:Confusion matrix:  
[[282 39 2]  
 [ 5 54 7]  
 [ 1 1 277]  
 SVC:Accuracy Score  
0.87586854468939  
SVC:Classification report:  
precision recall f1-score support  
1.0 0.99 0.88 0.93 333  
2.0 0.57 0.84 0.68 64  
3.0 0.75 0.93 0.83 29  
accuracy 0.77 426  
macro avg 0.77 0.88 0.81 426  
weighted avg 0.91 0.88 0.89 426

```
In [48]: from sklearn.model_selection import GridSearchCV
```

```
In [49]: bestpipe = Pipeline([('scalar6',StandardScaler()),
                          ('classifier',RandomForestClassifier())])
```

```
In [50]: grid_param = {
"classifier__n_estimators": [10, 100, 200, 300, 1000],
"classifier__max_depth": [5, 8, 15, 25, 38, None],
"classifier__min_samples_leaf": [1, 2, 5, 20, 15, 100],
"classifier__max_leaf_nodes": [2, 5, 10]}
# create a gridsearch = GridSearchCV(bestpipe, grid_param, cv=5, verbose=0,n_jobs=-1) # Fit Grid search
best_model = gridsearch.fit(X_train,y_train)
```

```
In [51]: best_model.best_params_
```

{'classifier\_\_max\_depth': 30,  
'classifier\_\_max\_leaf\_nodes': 10,  
'classifier\_\_min\_samples\_leaf': 2,  
'classifier\_\_n\_estimators': 200}

```
In [52]: pipelinef = Pipeline([('scalar7',StandardScaler()),
                           ('classifier',RandomForestClassifier())])
```

```
In [53]: pipelinef.fit(X_training,y_training)
```

```
In [54]: Pipeline(steps=[('scalar7', StandardScaler()),
                    ('classifier', RandomForestClassifier(max_depth=30, max_leaf_nodes=10, min_samples_leaf=2, n_estimators=200))])
```

```
In [55]: pipelinef.fit(X_training,y_training)
```

```
In [56]: y_pred=pipefinal.predict(X_test)
```

```
In [57]: print("Accuracy Score ")
print("{} {}".format(accuracy_score(y_test,y_pred)))
print("{} Classification report: {}".format(pipe_dict[i]), format(classification_report(y_test,y_pred)))
print("{} Confusion matrix: {}".format(pipe_dict[i]), format(confusion_matrix(y_test,y_pred)))
print("{} Confusion matrix: {}".format(pipe_dict[i]), format(confusion_matrix(y_test,y_pred)))
```

Accuracy Score  
0.87586854468939  
Classification report:  
precision recall f1-score support  
1.0 0.98 0.87 0.92 333  
2.0 0.62 0.88 0.72 64  
3.0 0.50 0.79 0.61 29  
accuracy 0.83 426  
macro avg 0.76 0.88 0.81 426  
weighted avg 0.91 0.88 0.89 426

Confusion matrix:  
[[289 23 8]  
 [ 5 56 3]  
 [ 0 2 77]]

```
In [ ]:
```