

NAME: K S H V SAI HARI KRISHNA

REG NO: 21BCE8069

1.Download the Employee Attrition Dataset

<https://www.kaggle.com/datasets/patelprashant/employee-attrition>

2.Perfrom Data Preprocessing

3.Model Building using Logistic Regression and Decision Tree and Random Forest

4.Calculate Performance metrics

```
#Import the Libraries.
```

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
#Importing the dataset.
```

```
df=pd.read_csv("Employee-Attrition.csv")
```

```
df.head()
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	.
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	.
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	.
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	.
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7	.

5 rows × 35 columns

```
df.shape
```

(1470, 35)

```
df.Age.value_counts()
```

```
35    78
34    77
36    69
31    69
29    68
32    61
30    60
33    58
38    58
40    57
37    50
27    48
28    48
42    46
39    42
45    41
41    40
26    39
44    33
46    33
43    32
50    30
25    26
24    26
49    24
47    24
```

```
55    22
51    19
53    19
48    19
54    18
52    18
22    16
56    14
23    14
58    14
21    13
20    11
59    10
19     9
18     8
60     5
57     4
Name: Age, dtype: int64
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
 #   Column            Non-Null Count  Dtype  
 --- 
 0   Age               1470 non-null    int64  
 1   Attrition         1470 non-null    object 
 2   BusinessTravel    1470 non-null    object 
 3   DailyRate          1470 non-null    int64  
 4   Department         1470 non-null    object 
 5   DistanceFromHome  1470 non-null    int64  
 6   Education          1470 non-null    int64  
 7   EducationField     1470 non-null    object 
 8   EmployeeCount      1470 non-null    int64  
 9   EmployeeNumber     1470 non-null    int64  
 10  EnvironmentSatisfaction  1470 non-null    int64  
 11  Gender             1470 non-null    object 
 12  HourlyRate         1470 non-null    int64  
 13  JobInvolvement    1470 non-null    int64  
 14  JobLevel           1470 non-null    int64  
 15  JobRole            1470 non-null    object 
 16  JobSatisfaction    1470 non-null    int64  
 17  MaritalStatus      1470 non-null    object 
 18  MonthlyIncome      1470 non-null    int64  
 19  MonthlyRate         1470 non-null    int64  
 20  NumCompaniesWorked 1470 non-null    int64  
 21  Over18             1470 non-null    object 
 22  OverTime            1470 non-null    object 
 23  PercentSalaryHike  1470 non-null    int64  
 24  PerformanceRating  1470 non-null    int64  
 25  RelationshipSatisfaction  1470 non-null    int64  
 26  StandardHours       1470 non-null    int64  
 27  StockOptionLevel    1470 non-null    int64  
 28  TotalWorkingYears   1470 non-null    int64  
 29  TrainingTimesLastYear 1470 non-null    int64  
 30  WorkLifeBalance    1470 non-null    int64  
 31  YearsAtCompany     1470 non-null    int64  
 32  YearsInCurrentRole 1470 non-null    int64  
 33  YearsSinceLastPromotion 1470 non-null    int64  
 34  YearsWithCurrManager 1470 non-null    int64  
dtypes: int64(26), object(9)
memory usage: 402.1+ KB
```

```
df.describe()
```

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumber	EnvironmentSatisfaction	HourlyRate	JobInvolvement
count	1470.000000	1470.000000	1470.000000	1470.000000	1470.0	1470.000000	1470.000000	1470.000000	14

```
#Checking for Null Values.
df.isnull().any()

Age          False
Attrition    False
BusinessTravel False
DailyRate     False
Department   False
DistanceFromHome False
Education     False
EducationField False
EmployeeCount False
EmployeeNumber False
EnvironmentSatisfaction False
Gender        False
HourlyRate    False
JobInvolvement False
JobLevel      False
JobRole       False
JobSatisfaction False
MaritalStatus False
MonthlyIncome False
MonthlyRate    False
NumCompaniesWorked False
Over18        False
OverTime      False
PercentSalaryHike False
PerformanceRating False
RelationshipSatisfaction False
StandardHours False
StockOptionLevel False
TotalWorkingYears False
TrainingTimesLastYear False
WorkLifeBalance False
YearsAtCompany False
YearsInCurrentRole False
YearsSinceLastPromotion False
YearsWithCurrManager False
dtype: bool
```

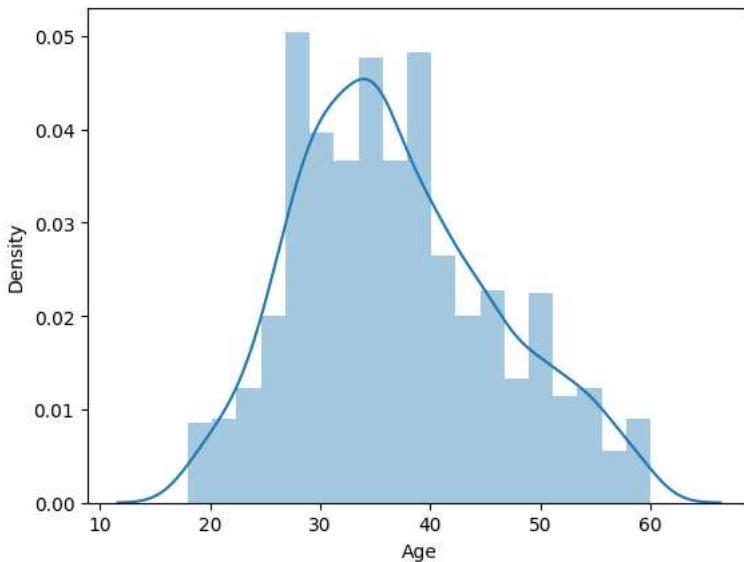


```
df.isnull().sum()

Age          0
Attrition    0
BusinessTravel 0
DailyRate     0
Department   0
DistanceFromHome 0
Education     0
EducationField 0
EmployeeCount 0
EmployeeNumber 0
EnvironmentSatisfaction 0
Gender        0
HourlyRate    0
JobInvolvement 0
JobLevel      0
JobRole       0
JobSatisfaction 0
MaritalStatus 0
MonthlyIncome 0
MonthlyRate    0
NumCompaniesWorked 0
Over18        0
OverTime      0
PercentSalaryHike 0
PerformanceRating 0
RelationshipSatisfaction 0
StandardHours 0
StockOptionLevel 0
TotalWorkingYears 0
TrainingTimesLastYear 0
WorkLifeBalance 0
YearsAtCompany 0
YearsInCurrentRole 0
YearsSinceLastPromotion 0
YearsWithCurrManager 0
dtype: int64
```

```
#Data Visualization.  
sns.distplot(df["Age"])  
  
<ipython-input-84-25fc8198007f>:2: UserWarning:  
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.  
Please adapt your code to use either `displot` (a figure-level function with  
similar flexibility) or `histplot` (an axes-level function for histograms).  
For a guide to updating your code to use the new functions, please see  
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
```

```
sns.distplot(df["Age"])  
<Axes: xlabel='Age', ylabel='Density'>
```



```
df.corr()
```

```
<ipython-input-85-2f6f6606aa2c>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version
df.corr()
```

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumber	EnvironmentSatisfaction	Hourly
Age	1.000000	0.010661	-0.001686	0.208034	NaN	-0.010145	0.010146	0.02
DailyRate	0.010661	1.000000	-0.004985	-0.016806	NaN	-0.050990	0.018355	0.02
DistanceFromHome	-0.001686	-0.004985	1.000000	0.021042	NaN	0.032916	-0.016075	0.03
Education	0.208034	-0.016806	0.021042	1.000000	NaN	0.042070	-0.027128	0.01
EmployeeCount	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
EmployeeNumber	-0.010145	-0.050990	0.032916	0.042070	NaN	1.000000	0.017621	0.03

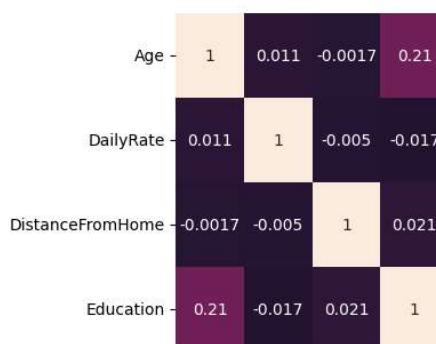
df.head()

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	..
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	..
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	..
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	..
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7	..

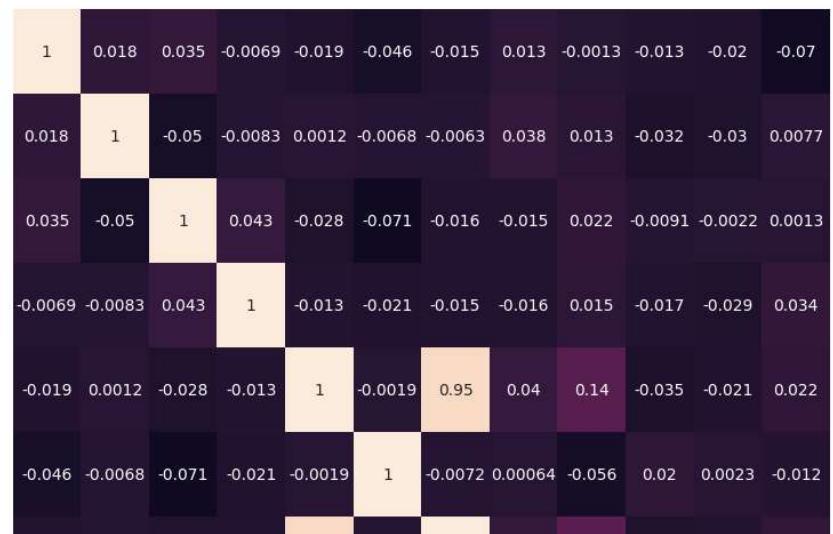
5 rows × 35 columns

```
plt.subplots(figsize = (25,25))
sns.heatmap(df.corr(), annot=True)
```

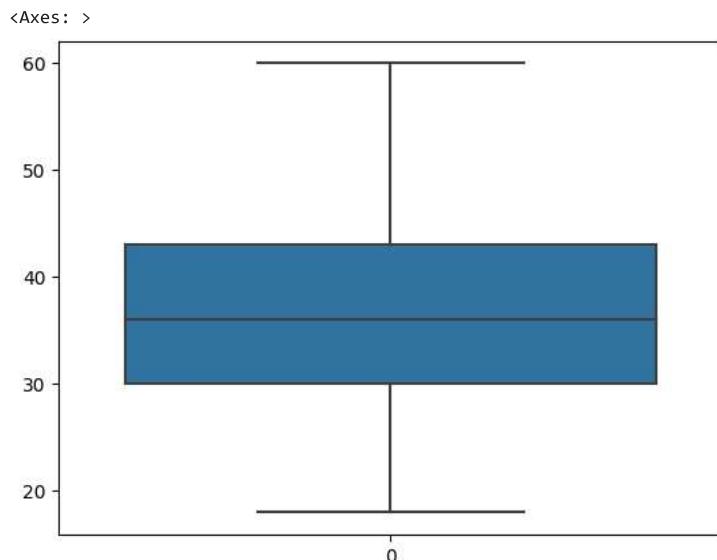
```
<ipython-input-87-9329d5e70af4>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version  
  sns.heatmap(df.corr(), annot=True)  
<Axes: >
```



EmployeeCount -



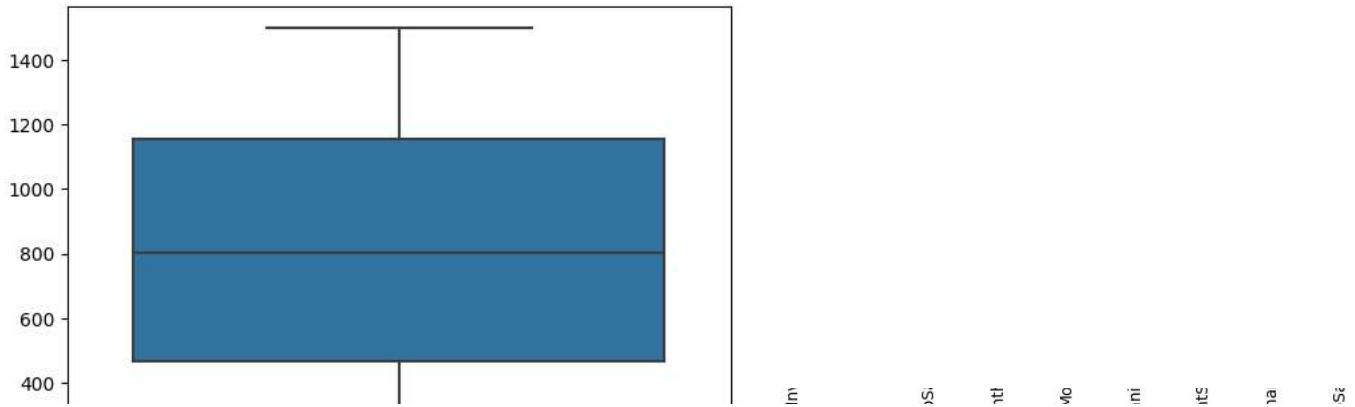
```
sns.boxplot(df['Age'])
```



```
sns.boxplot(df['DailyRate'])
```



<Axes: >



df.head()

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	.
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	.
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	.
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	.
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7	.

5 rows × 35 columns

```
x=df.iloc[:,1:4]
x.head()
```

	Attrition	BusinessTravel	DailyRate	grid
0	Yes	Travel_Rarely	1102	blue
1	No	Travel_Frequently	279	
2	Yes	Travel_Rarely	1373	
3	No	Travel_Frequently	1392	
4	No	Travel_Rarely	591	

```
y=df.Attrition
y.head()
```

```
0    Yes
1    No
2    Yes
3    No
4    No
Name: Attrition, dtype: object
```

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
y=le.fit_transform(y)
```

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
y_test=le.fit_transform(y_test)
```

y

```
array([1, 0, 1, ..., 0, 0, 0])
```

y_test

```
array([0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0,
       0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
       1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1,
       0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
```

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
x.Attrition=le.fit_transform(x.Attrition)
x.head()
```

	Attrition	BusinessTravel	DailyRate	
0	1	Travel_Rarely	1102	
1	0	Travel_Frequently	279	
2	1	Travel_Rarely	1373	
3	0	Travel_Frequently	1392	
4	0	Travel_Rarely	591	

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
x.BusinessTravel =le.fit_transform(x.BusinessTravel )
x.head()
```

	Attrition	BusinessTravel	DailyRate	
0	1	2	1102	
1	0	1	279	
2	1	2	1373	
3	0	1	1392	
4	0	2	591	

```
#feature scaling
from sklearn.preprocessing import MinMaxScaler
ms=MinMaxScaler()
x_scaled=pd.DataFrame(ms.fit_transform(x),columns=x.columns)
```

x_scaled

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7
...
1465	36	No	Travel_Frequently	884	Research & Development	23	2	Medical	1	2061
1466	39	No	Travel_Rarely	613	Research & Development	6	1	Medical	1	2062
1467	27	No	Travel_Rarely	155	Research & Development	4	3	Life Sciences	1	2064

```
dtc.predict(ms.transform([[1,19,19000]]))
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but MinMaxScaler was fitted
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but DecisionTreeClassifier was fitted
  warnings.warn(
array([1])
```

▼ Evaluation of classification model

```
#Accuracy score
from sklearn.metrics import accuracy_score,confusion_matrix,classification_report,roc_auc_score,roc_curve
```

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
y=le.fit_transform(y)
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
y_test=le.fit_transform(y_test)

y_test

array([0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
       0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1,
       1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1,
       0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,
       0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0,
       0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,
```

```
accuracy_score(y_test,pred)
```

1.0

```
confusion_matrix(y_test,pred)
```

```
array([[245,   0],
       [  0, 49]])
```


ROC CURVE



```
from sklearn import tree
plt.figure(figsize=(25,15))
tree.plot_tree(dtc,filled=True)

[Text(0.5, 0.75, 'x[0] <= 0.5\n gini = 0.269\n samples = 1176\n value = [988, 188]'),
 Text(0.25, 0.25, 'gini = 0.0\n samples = 988\n value = [988, 0]'),
 Text(0.75, 0.25, 'gini = 0.0\n samples = 188\n value = [0, 188]')]
```

x[0] <= 0.
gini = 0.26
samples = 1176
value = [988, 188]

gini = 0.0
samples = 988
value = [988, 0]

```
from sklearn.model_selection import GridSearchCV
parameter={
    'criterion':['gini','entropy'],
    'splitter':['best','random'],
    'max_depth':[1,2,3,4,5],
    'max_features':['auto', 'sqrt', 'log2']
}

grid_search=GridSearchCV(estimator=dtc,param_grid=parameter,cv=5,scoring="accuracy")

grid_search.fit(x_train,y_train)
```



```
warnings.warn(  
    /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
    warnings.warn(  
        /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
        warnings.warn(  
            /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
            warnings.warn(  
                /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                warnings.warn(  
                    /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                    warnings.warn(  
                        /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                        warnings.warn(  
                            /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                            warnings.warn(  
                                /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                                warnings.warn(  
                                    /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                                    warnings.warn(  
                                        /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                                        warnings.warn(  
                                            /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                                            warnings.warn(  
                                                /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                                                warnings.warn(  
                                                    /usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and  
                                                    warnings.warn(  
                                                       ▾ GridSearchCV  
  
grid_search.best_params_  
  
{'criterion': 'entropy',  
 'max_depth': 5,  
 'max_features': 'log2',  
 'splitter': 'best'}  
  
dtc_cv=DecisionTreeClassifier(criterion= 'entropy',  
                             max_depth=3,  
                             max_features='sqrt',  
                             splitter='best')  
dtc_cv.fit(x_train,y_train)  
  
              ▾ DecisionTreeClassifier  
DecisionTreeClassifier(criterion='entropy', max_depth=3, max_features='sqrt')  
  
pred=dtc_cv.predict(x_test)  
  
print(classification_report(y_test,pred))  
  


|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 1.00      | 1.00   | 1.00     | 245     |
| 1            | 1.00      | 1.00   | 1.00     | 49      |
| accuracy     |           |        | 1.00     | 294     |
| macro avg    | 1.00      | 1.00   | 1.00     | 294     |
| weighted avg | 1.00      | 1.00   | 1.00     | 294     |


```

RandomForestClassifier

```
from sklearn.ensemble import RandomForestClassifier  
rfc=RandomForestClassifier()  
  
forest_params = [{"max_depth": list(range(10, 15)), "max_features": list(range(0,14))}]  
  
rfc_cv= GridSearchCV(rfc,param_grid=forest_params,cv=10,scoring="accuracy")  
  
rfc_cv.fit(x_train,y_train)
```

```
sr/local/lib/python3.10/dist-packages/sklearn/model_selection/_validation.py:378: FitFailedWarning: Some fits failed out of a total of 700.
  score on these train-test partitions for these parameters will be set to nan.
  these failures are not expected, you can try to debug them by setting error_score='raise'.
  low are more details about the failures:
-----
  fits failed with the following error:
  traceback (most recent call last):
  file "/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_validation.py",
    estimator.fit(X_train, y_train, **fit_params)
  file "/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_forest.py", line 340,
    self._validate_params()
  file "/usr/local/lib/python3.10/dist-packages/sklearn/base.py", line 600, in _validate_parameters
    validate_parameter_constraints(
  file "/usr/local/lib/python3.10/dist-packages/sklearn/utils/_param_validation.py", line 144, in validate_parameters
    raise InvalidParameterError(
sklearn.utils._param_validation.InvalidParameterError: The 'max_features' parameter of RandomForestClassifier is not valid for this estimator. It must be a float between 0.0 and 1.0, or an integer between 0 and n_features inclusive. Got 'None'.
```

`warnings.warn(some_fits_failed_message, FitFailedWarning)`

sr/local/lib/python3.10/dist-packages/sklearn/model_selection/_search.py:952: UserWarning

L 1. 1. 1. 1. 1. 1. 1. nan 1. 1. 1. 1. 1. 1. 1.
L 1. 1. 1. 1. 1. nan 1. 1. 1. 1. 1. 1. 1. 1. 1.
L 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]

`warnings.warn(`

`GridSearchCV`

`estimator: RandomForestClassifier`

`RandomForestClassifier`

+ Code + Text

```
pred=rfc_cv.predict(x_test)
```

```
print(classification_report(y_test,pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	245
1	1.00	1.00	1.00	49
accuracy			1.00	294
macro avg	1.00	1.00	1.00	294
weighted avg	1.00	1.00	1.00	294

```
rfc_cv.best_params_
```

```
{'max_depth': 10, 'max_features': 1}
```