

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

```
In [2]: data = pd.read_csv('WA_Fn-UseC_-HR-Employee-Attrition.csv')
data.head()
```

```
Out[2]:
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Educa
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Lif
1	49	No	Travel_Frequently	279	Research & Development	8	1	Lif
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Lif
4	27	No	Travel_Rarely	591	Research & Development	2	1	

5 rows × 35 columns

```
In [3]: data.tail()
```

```
Out[3]:
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Educa
1465	36	No	Travel_Frequently	884	Research & Development	23	2	
1466	39	No	Travel_Rarely	613	Research & Development	6	1	
1467	27	No	Travel_Rarely	155	Research & Development	4	3	
1468	49	No	Travel_Frequently	1023	Sales	2	3	
1469	34	No	Travel_Rarely	628	Research & Development	8	3	

5 rows × 35 columns

```
In [4]: data.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
```

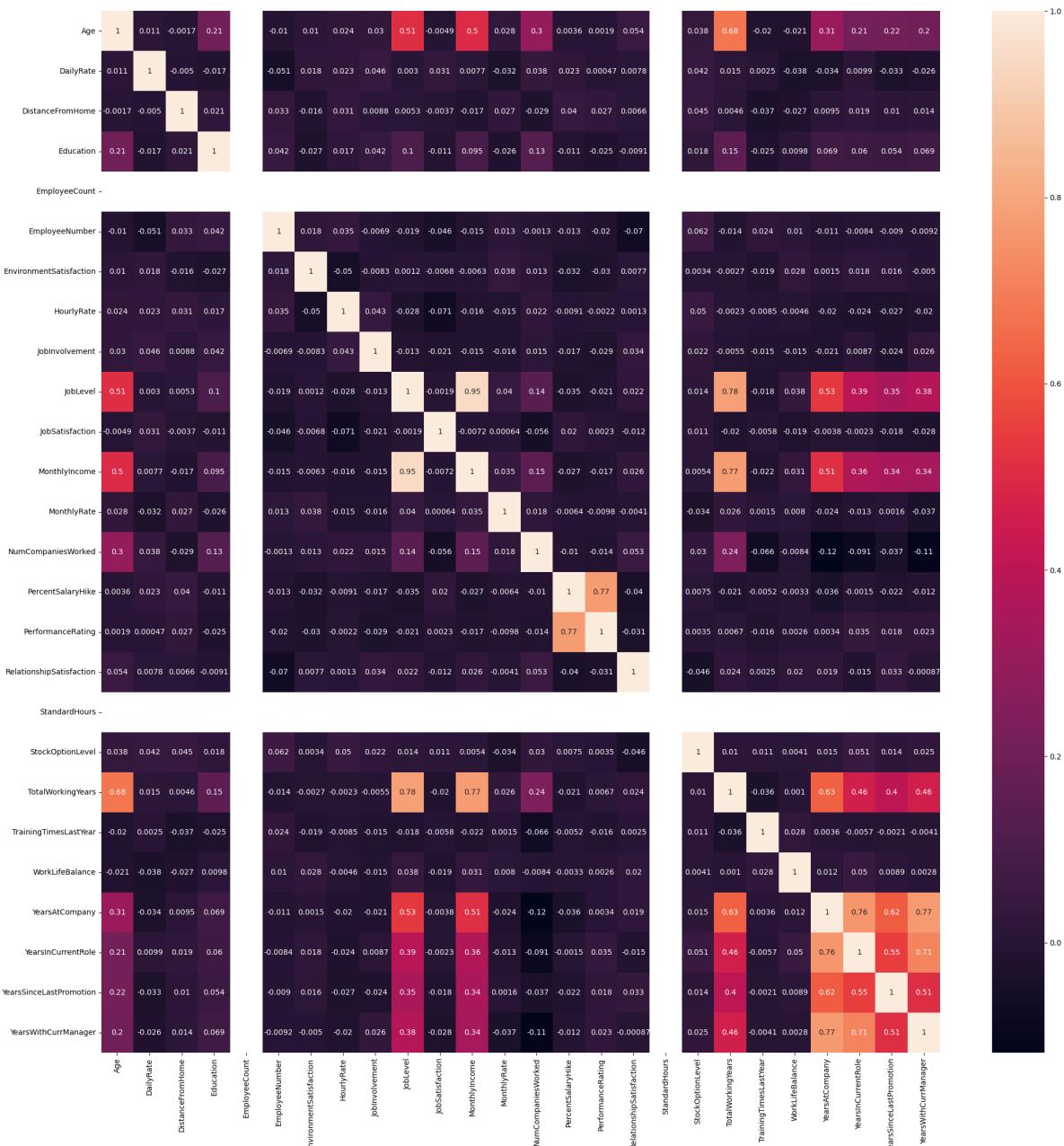
In [5]: `data.describe()`


```
C:\Users\karth\AppData\Local\Temp\ipykernel_37264\2763113851.py:1: FutureWarning:  
The default value of numeric_only in DataFrame.corr is deprecated. In a future ver  
sion, it will default to False. Select only valid columns or specify the value of  
numeric_only to silence this warning.
```

```
cor = data.corr()
```

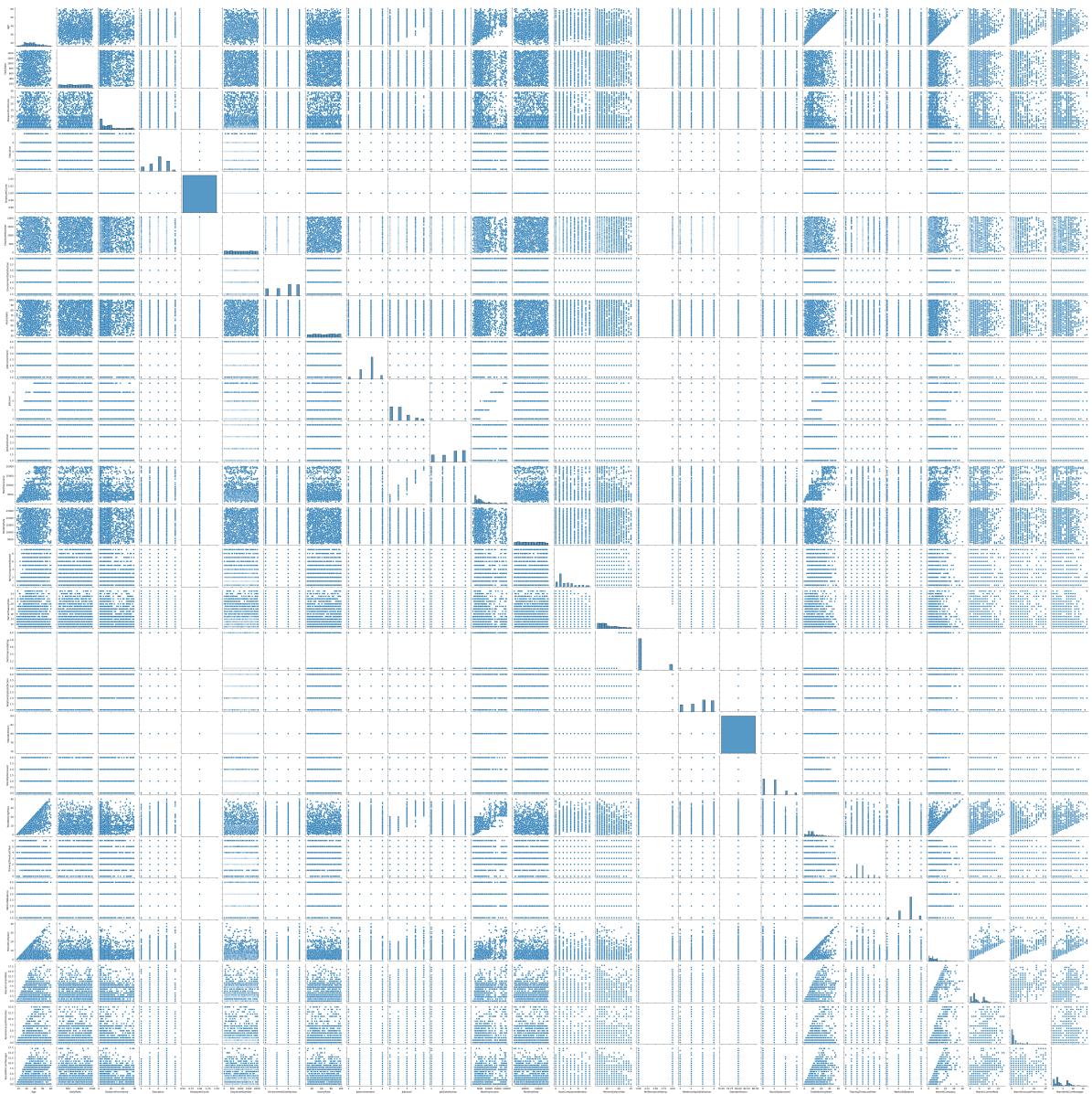
```
In [8]: fig, ax = plt.subplots(figsize=(25,25))  
sns.heatmap(cor, annot=True)
```

```
Out[8]: <Axes: >
```



```
In [9]: sns.pairplot(data)
```

```
Out[9]: <seaborn.axisgrid.PairGrid at 0x2402822bb10>
```



```
In [10]: from sklearn.preprocessing import LabelEncoder
```

```
In [11]: le=LabelEncoder()
```

```
In [12]: data["BusinessTravel"]=le.fit_transform(data["BusinessTravel"])
```

```
In [13]: data["Department"]=le.fit_transform(data["Department"])
```

```
In [14]: data["EducationField"]=le.fit_transform(data["EducationField"])
```

```
In [15]: data["Gender"]=le.fit_transform(data["Gender"])
```

```
In [16]: data["JobRole"]=le.fit_transform(data["JobRole"])
```

```
In [17]: data["MaritalStatus"]=le.fit_transform(data["MaritalStatus"])
```

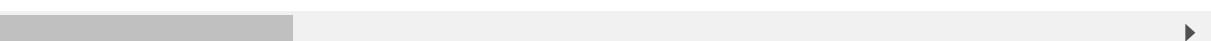
```
In [18]: data["Over18"]=le.fit_transform(data["Over18"])
```

```
In [19]: data["OverTime"]=le.fit_transform(data["OverTime"])
```

```
In [20]: data.head()
```

```
Out[20]:    Age Attrition BusinessTravel DailyRate Department DistanceFromHome Education Education
      0  41       Yes           2      1102           2              1        2
      1  49        No           1      279            1              8        1
      2  37       Yes           2     1373            1              2        2
      3  33        No           1     1392            1              3        4
      4  27        No           2      591            1              2        1
```

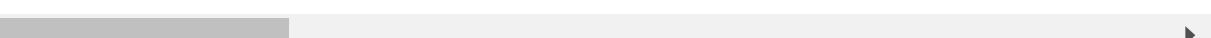
5 rows × 35 columns



```
In [21]: data.tail()
```

```
Out[21]:    Age Attrition BusinessTravel DailyRate Department DistanceFromHome Education Education
      1465   36        No           1      884           1             23        2
      1466   39        No           2      613           1              6        1
      1467   27        No           2      155           1              4        3
      1468   49        No           1     1023           2              2        3
      1469   34        No           2      628           1              8        3
```

5 rows × 35 columns



```
In [22]: X = data.drop(columns=["EmployeeNumber", "EmployeeCount", "StandardHours", "Attrition", "DistanceFromHome"], axis=1)
```

```
In [23]: y = data["Attrition"]
```

```
In [24]: from sklearn.preprocessing import MinMaxScaler
ms = MinMaxScaler()
```

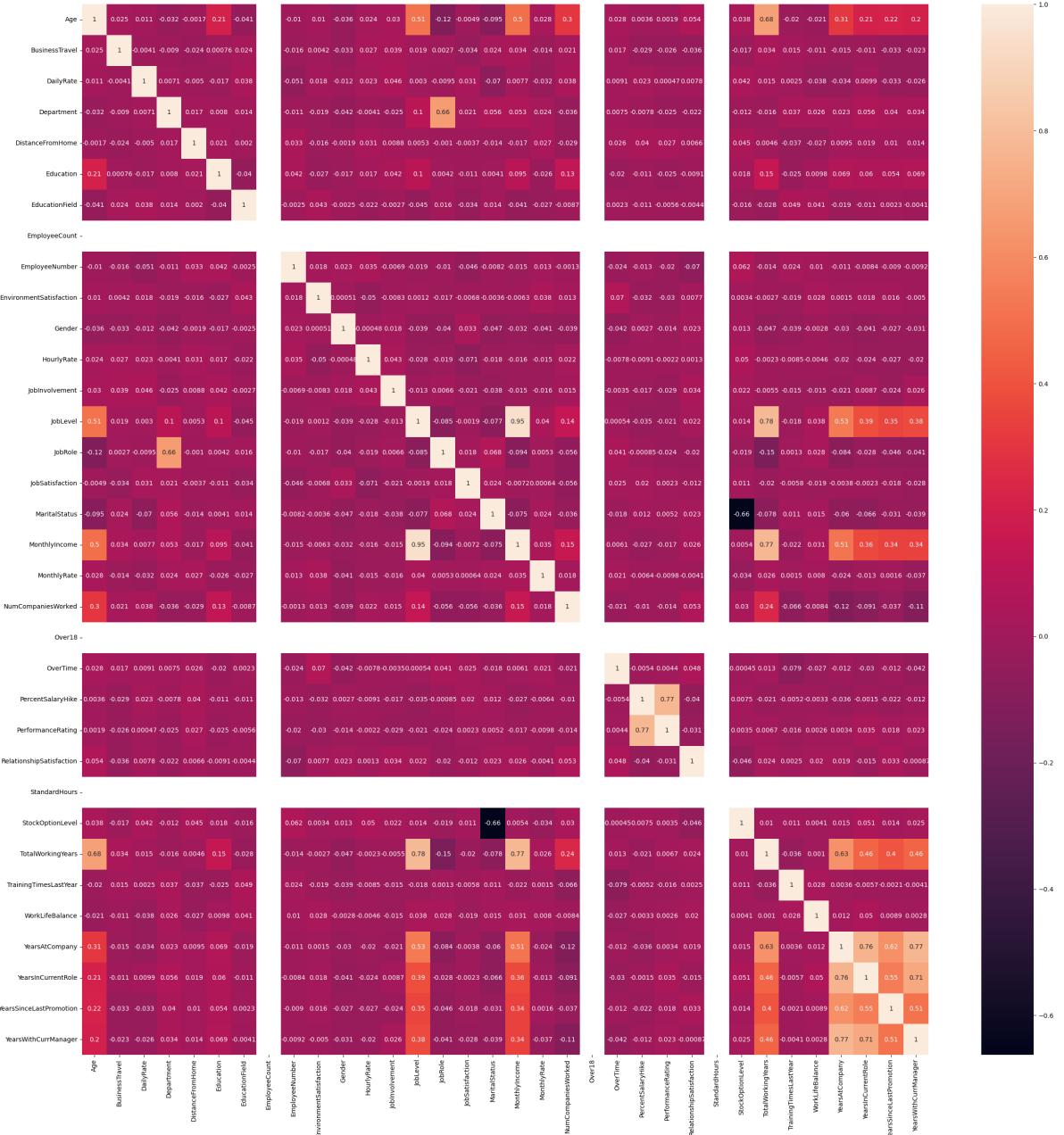
```
In [25]: X_Scaled = ms.fit_transform(X)
```

```
In [26]: cor = data.corr()
```

```
C:\Users\karth\AppData\Local\Temp\ipykernel_37264\1426905697.py:1: FutureWarning:
The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
cor = data.corr()
```

```
In [27]: fig, ax = plt.subplots(figsize=(30,30))
sns.heatmap(cor, annot=True)
```

```
Out[27]: <Axes: >
```



```
In [28]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(X_Scaled,y,test_size =0.2,random_
```

```
In [29]: from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state=0)
classifier.fit(x_train,y_train)
```

```
Out[29]: LogisticRegression
LogisticRegression(random_state=0)
```

```
In [30]: from sklearn.metrics import accuracy_score,confusion_matrix
y_pred = classifier.predict(x_test)
cm = confusion_matrix(y_test, y_pred)
print(cm)
accuracy_score(y_test, y_pred)*100
```

```
[[242  3]
 [ 32 17]]
88.09523809523809
```

Out[30]:

```
In [31]: from sklearn.metrics import accuracy_score,confusion_matrix,classification_report,
```

```
In [32]: print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
No	0.88	0.99	0.93	245
Yes	0.85	0.35	0.49	49
accuracy			0.88	294
macro avg	0.87	0.67	0.71	294
weighted avg	0.88	0.88	0.86	294

```
In [33]: from sklearn.tree import DecisionTreeClassifier  
dtc=DecisionTreeClassifier()
```

```
In [34]: dtc.fit(x_train,y_train)
```

```
Out[34]: ▾ DecisionTreeClassifier  
DecisionTreeClassifier()
```

```
In [35]: from sklearn.metrics import accuracy_score,confusion_matrix  
y_pred = dtc.predict(x_test)  
cm = confusion_matrix(y_test, y_pred)  
print(cm)  
accuracy_score(y_test, y_pred)*100
```

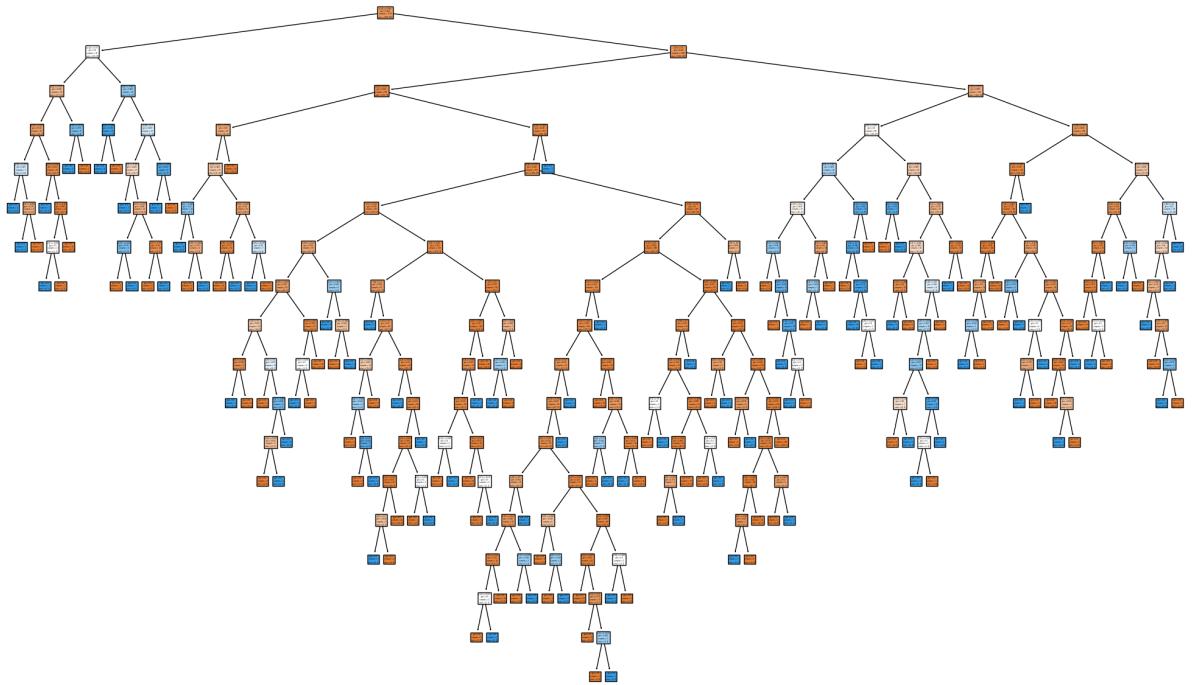
```
[[202  43]  
 [ 35  14]]
```

```
Out[35]: 73.46938775510205
```

```
In [36]: from sklearn import tree  
plt.figure(figsize=(25,15))  
tree.plot_tree(dtc,filled=True)
```



```
Text(0.9530201342281879, 0.6388888888888888, 'x[25] <= 0.833\ngini = 0.444\nsamples = 3\nvalue = [1, 2]'),
Text(0.9463087248322147, 0.5833333333333334, 'gini = 0.0\nsamples = 2\nvalue = [0, 2]'),
Text(0.959731543624161, 0.5833333333333334, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.9865771812080537, 0.6944444444444444, 'x[28] <= 0.1\ngini = 0.49\nsamples = 28\nvalue = [12, 16]'),
Text(0.9798657718120806, 0.6388888888888888, 'x[10] <= 0.833\ngini = 0.48\nsamples = 20\nvalue = [12, 8]'),
Text(0.9731543624161074, 0.5833333333333334, 'x[4] <= 0.018\ngini = 0.415\nsamples = 17\nvalue = [12, 5]'),
Text(0.9664429530201343, 0.5277777777777778, 'gini = 0.0\nsamples = 2\nvalue = [0, 2]'),
Text(0.9798657718120806, 0.5277777777777778, 'x[15] <= 0.365\ngini = 0.32\nsamples = 15\nvalue = [12, 3]'),
Text(0.9731543624161074, 0.4722222222222222, 'gini = 0.0\nsamples = 11\nvalue = [11, 0]'),
Text(0.9865771812080537, 0.4722222222222222, 'x[23] <= 0.325\ngini = 0.375\nsamples = 4\nvalue = [1, 3]'),
Text(0.9798657718120806, 0.4166666666666667, 'gini = 0.0\nsamples = 3\nvalue = [0, 3]'),
Text(0.9932885906040269, 0.4166666666666667, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.9865771812080537, 0.5833333333333334, 'gini = 0.0\nsamples = 3\nvalue = [0, 3]'),
Text(0.9932885906040269, 0.6388888888888888, 'gini = 0.0\nsamples = 8\nvalue = [0, 8]])
```



```
In [37]: from sklearn.model_selection import GridSearchCV
parameter={

  'criterion':['gini', 'entropy'],
  'splitter':['best', 'random'],
  'max_depth':[1,2,3,4,5,6,7,8,9,10],
  'max_features':['auto', 'sqrt', 'log2']

}
```

```
In [38]: grid_search=GridSearchCV(estimator=dtc,param_grid=parameter,cv=5,scoring="accuracy")
```



```
Out[40]: {'criterion': 'gini',
          'max_depth': 3,
          'max_features': 'sqrt',
          'splitter': 'best'}
```

```
In [41]: dtc_cv=DecisionTreeClassifier(criterion= 'entropy',
                                       max_depth= 4,
                                       max_features= 'sqrt',
                                       splitter= 'best')
dtc_cv.fit(x_train,y_train)
```

```
Out[41]: ▾ DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy', max_depth=4, max_features='sqrt')
```

```
In [42]: print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
No	0.85	0.82	0.84	245
Yes	0.25	0.29	0.26	49
accuracy			0.73	294
macro avg	0.55	0.56	0.55	294
weighted avg	0.75	0.73	0.74	294

```
In [43]: from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 1000, criterion = 'entropy', random_state=0)
classifier.fit(x_train, y_train)
```

```
Out[43]: ▾ RandomForestClassifier
RandomForestClassifier(criterion='entropy', n_estimators=1000, random_state=0)
```

```
In [44]: from sklearn.metrics import confusion_matrix, accuracy_score
y_pred = classifier.predict(x_test)
cm = confusion_matrix(y_test, y_pred)
print(cm)
accuracy_score(y_test, y_pred)
```

```
[[243  2]
 [ 41  8]]
```

```
Out[44]: 0.8537414965986394
```

```
In [45]: from sklearn.ensemble import RandomForestClassifier
```

```
In [46]: rfc=RandomForestClassifier()
```

```
In [47]: forest_params = [{ 'max_depth': list(range(10, 15)), 'max_features': list(range(0,1))}]
```

```
In [48]: rfc_cv=GridSearchCV(rfc,param_grid=forest_params,cv=10,scoring="accuracy")
```

```
In [49]: rfc_cv.fit(x_train,y_train)
```

