

1.Download the Employee Attrition Dataset <https://www.kaggle.com/datasets/pareshprashant/employee-attrition>

2.Perform Data Preprocessing

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

4.Calculate Performance metrics

```
In [1486]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

In [1487]: df=pd.read_csv("../content/HR_An_Fn-Use-C-1R-Employee-Attrition.csv")

Out[1487]:
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	RelationshipSatisfaction	StandardHours	StockOptionLevel	TotalWorkingYears
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	...	1	80	0	8
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	...	4	80	1	10
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	...	2	80	0	7
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	...	3	80	0	8
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7	...	4	80	1	6
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
1465	36	No	Travel_Frequently	884	Research & Development	23	2	Medical	1	2061	...	3	80	1	17
1466	39	No	Travel_Rarely	613	Research & Development	6	1	Medical	1	2062	...	1	80	1	9
1467	27	No	Travel_Rarely	158	Research & Development	4	3	Life Sciences	1	2064	...	2	80	1	6
1468	49	No	Travel_Frequently	1023	Sales	2	3	Life Sciences	1	2065	...	4	80	0	17
1469	34	No	Travel_Rarely	628	Research & Development	8	3	Medical	1	2068	...	1	80	0	6

1470 rows x 35 columns

```
In [1488]: df.head()
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	RelationshipSatisfaction	StandardHours	StockOptionLevel	TotalWorkingYears	Tr
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	...	1	80	0	8	
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	...	4	80	1	10	
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	...	2	80	0	7	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	...	3	80	0	8	
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7	...	4	80	1	6	

5 rows x 35 columns

```
In [1489]: <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   bool
22  OverTime               1470 non-null   bool
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  YearWithCurrManager    1470 non-null   int64
dtypes: int64(29), object(6)
memory usage: 402.1+ KB
```

```
In [1490]: df.describe()
```

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumber	EnvironmentSatisfaction	HourlyRate	JobInvolvement	JobLevel	...	RelationshipSatisfaction	StandardHours	StockOptionLevel
count	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	...	1470.000000	1470.0	1470.00
mean	36.923810	802.488714	9.132517	2.912925	1.0	1024.865306	2.721769	65.891156	2.729932	2.063946	...	2.712245	80.0	0.79
std	13.53373	403.509100	8.106864	1.024165	0.0	602.024335	1.093002	20.329428	0.711661	1.109440	...	1.081209	0.0	0.85
min	18.000000	102.000000	1.000000	1.000000	1.000000	30.000000	1.000000	30.000000	1.000000	1.000000	...	1.000000	80.0	0.00
25th percentile	45.000000	65.000000	2.000000	2.000000	1.0	40.000000	2.000000	40.000000	2.000000	1.000000	...	2.000000	80.0	0.00
50th percentile	36.000000	802.000000	7.000000	3.000000	1.0	1025.000000	3.000000	66.000000	3.000000	2.000000	...	3.000000	80.0	1.00
75th percentile	45.000000	1157.000000	14.000000	4.000000	1.0	1555.750000	4.000000	83.750000	3.000000	3.000000	...	4.000000	80.0	1.00
max	60.000000	1459.000000	29.000000	5.000000	1.0	2068.000000	4.000000	100.000000	4.000000	5.000000	...	4.000000	80.0	3.00

8 rows x 26 columns

```
In [1491]: df.shape
(1470, 35)
```

```
In [1492]: #Checking for Null values
df.isnull().any()
```

```
Out[1492]: 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
YearWithCurrManager    False
dtype: bool
```

```
In [1493]: df.isnull().sum()
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
YearWithCurrManager 0
dtype: int64
```

```
In [1494]: sns.distplot(df['Age'])
```

`<matplotlib.figure.Figure at 0x7afe04ea375b>: UserWarning: 'distplot' is a deprecated function and will be removed in seaborn v0.14.0. Please adapt your code to use either 'displot' (an 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/maskor/d644147ed2974457a063727580be5751>

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

```
In [1495]: df.corr()
```

`<matplotlib.figure.Figure at 0x7f6688a2431c>: 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.`

```
Out[1495]:
```

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumber	EnvironmentSatisfaction	HourlyRate	JobInvolvement	JobLevel	...	RelationshipSatisfaction	StandardHours	StockOptionLevel
Age	1.000000	0.010661	-0.001696	0.209034	NAN	NAN	0.010146	0.024287	0.029820	0.009004	...	0.030375	NAN	NAN
DailyRate	0.010661	1.000000	-0.004985	0.018066	NAN	NAN	-0.050900	0.018355	0.023361	0.046135	...	0.002966	...	0.007946
DistanceFromHome	-0.001696	-0.004985	1.000000	0.021042	NAN	NAN	0.032916	-0.016075	0.013111	0.008783	...	0.005031	...	0.000557
Education	0.209034	0.018066	0.021042	1.000000	NAN	NAN	0.042070	-0.027128	0.016775	0.042438	...	0.011589	...	-0.009118
EmployeeCount	NAN	NAN	NAN	NAN	1.000000	NAN	0.017621	0.035119	-0.006888	-0.018519	...	-0.006961	...	0.007495
EmployeeNumber	0.010146	0.050900	0.032916	0.042070	NAN	1.000000	0.017621	0.035119	-0.006888	-0.018519	...	-0.006961	...	0.007495
EnvironmentSatisfaction	0.010146	0.018355	-0.001697	0.021128	NAN	0.017621	1.000000	-0.049857	-0.008278	0.001212	...	0.001212	...	0.007495
HourlyRate	0.024287	0.023361	0.021123	0.018775	NAN	0.035119	-0.049857	1.000000	0.045861	-0.027883	...	0.001360	...	0.007495
JobInvolvement	0.029820	0.046135	0.005031	0.042438	NAN	-0.006888	-0.008278	0.045861	1.000000	-0.013630	...	0.024297	...	0.004297
JobLevel	0.009004	0.002966	0.000557	0.002122	NAN	-0.018519	0.001212	-0.027883	-0.013630	1.000000	...	0.021642	...	NAN
JobSatisfaction	0.046135	0.002966	0.000557	-0.002122	NAN	-0.006888	-0.008278	-0.013630	-0.027883	-0.013630	...	-0.012654	...	NAN
MonthlyIncome	0.497855	0.007707	-0.017014	0.094961	NAN	-0.014829	-0.006259	-0.015794	-0.012371	0.963030	...	0.025873	...	NAN
MonthlyRate	0.028051	0.032182	0.027473	0.026294	NAN	0.012648	-0.007600	-0.015097	-0.016332	0.395643	...	-0.004085	...	NAN
NumCompaniesWorked	0.299635	0.028153	-0.028251	0.126317	NAN	-0.001251	0.012594	0.023157	0.015012	0.142501	...	0.002733	...	NAN
PercentSalaryHike	0.003634	0.022704	0.040235	0.011111	NAN	-0.012944	-0.031701	-0.009062	-0.027970	-0.034740	...	-0.040490	...	NAN
PerformanceRating	0.001934	0.000473	0.027110	0.024539	NAN	-0.010359	-0.029548	-0.001372	-0.021222	-0.013151	...	-0.013151	...	NAN
RelationshipSatisfaction	0.053535	0.007846	0.000557	-0.009118	NAN	-0.026861	0.007655	0.001300	0.034249	0.021642	...	1.000000	...	NAN
StandardHours	NAN	NAN	NAN	NAN	NAN	NAN	NAN	NAN	NAN	NAN	...	NAN	...	NAN
StockOptionLevel	0.079710	0.042143	0.044872	0.018422	NAN	0.062227	0.003432	0.050283	0.021523	0.013984	...	-0.045952	...	NAN
TotalWorkingYears	0.020019	0.000015	0.000000	0.000000	NAN	-0.014865	-0.000993	-0.002024	-0.005529	0.707026	...	0.004854	...	NAN
TrainingTimesLastYear	0.019621	0.002463	-0.008462	0.025100	NAN	0.012602	-0.015359	-0.008548	-0.015338	-0.018191	...	0.024297	...	NAN
WorkLifeBalance	-0.021490	-0.023748	-0.026556	0.009819	NAN	0.010309	0.027627	-0.004907	-0.034617	0.037818	...	-0.019604	...	NAN
YearsAtCompany	0.113109	0.034055	0.005008	0.069114	NAN	-0.011240	0.001458	-0.015962	-0.021355	0.934739	...	0.019367	...	NAN
YearsInCurrentRole	0.212901	0.005932	0.018845	0.060236	NAN	-0.008416	0.018007	-0.024106	0.008717	0.389447	...	-0.015123	...	NAN
YearsSinceLastPromotion	0.216513	0.033229	0.010209	0.054254	NAN	-0.009019	0.016194	-0.028716	-0.024184	0.353885	...	0.033493	...	NAN
YearWithCurrManager	0.202089	-0.028163	0.014406	0.069065	NAN	-0.009197	-0.004999	-0.020123	0.025976	0.375281	...	-0.008667	...	NAN

26 rows x 26 columns

```
In [1496]: sns.heatmap(df.corr())
```

`<matplotlib.figure.Figure at 0x7f6688a2431c>: 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.`

```
Out[1496]:
```

`<Axes: >`

```
In [1497]: df.head()
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	RelationshipSatisfaction	StandardHours	StockOptionLevel	TotalWorkingYears	Tr
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	...	1	80	0	8	
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	...	4	80	1	10	
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	...	2	80	0	7	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	...	3	80	0	8	
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7	...	4	80	1	6	

5 rows x 35 columns

```
In [1498]: sns.boxplot(df['Age'])
```

`<Axes: >`

```
In [1499]: sns.boxplot(df['DailyRate'],color='cyan')
```

`<Axes: >`

```
In [1500]: df.head()
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	RelationshipSatisfaction	StandardHours	StockOptionLevel	TotalWorkingYears	Tr
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	...	1	80	0	8	
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	...	4	80	1	10	
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	...	2	80	0	7	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	...	3	80	0	8	
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical								