

```
In [1]: 1 # Importing required Libraries -
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from scipy import stats
7 from scipy.stats import ttest_ind # T-test for independent samples
8 from scipy.stats import shapiro # Shapiro-Wilk's test for Normality
9 #from scipy.stats import Levene # Levene's test for Equality of Variance
10 from scipy.stats import f_oneway # One-way ANOVA
11 from scipy.stats import chi2_contingency # Chi-square test of independence
12 #from scipy.stats import kruskal
13
14
15 import warnings
16 warnings.filterwarnings("ignore")
17
18
19 #####
20 # Display Changes
21 #####
22 # cf.go_offline()
23 from matplotlib import rc, rcParams
24 %matplotlib inline
25 rc('axes', linewidth=2)
26 rc('font', weight='bold')
27 sns.set_style('whitegrid')
28 #pd.set_option('max_rows',2000)
29 #pd.set_option('max_columns',255)
30 #rcParams['figure.figsize'] = 24,8
```

```
In [2]: 1 import matplotlib_inline
2 matplotlib_inline.backend_inline.set_matplotlib_formats('svg')
3 from IPython.display import display, HTML
4 display(HTML("<style>.container { width:100% !important; }</style>"))
```

```
In [3]: 1 import os
2 os.getcwd()
```

```
Out[3]: 'C:\\\\Users\\\\gurme\\\\OneDrive\\\\Desktop\\\\scaler'
```

```
In [4]: 1 import urllib.request  
2 url = 'https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/002/492/original/ola_driver_scaler.csv'  
3 urllib.request.urlretrieve(url, 'Ola.csv')
```

```
Out[4]: ('Ola.csv', <http.client.HTTPMessage at 0x1d7d9f18f10>)
```

```
In [5]: 1 df=pd.read_csv('Ola.csv')
```

### Defining Problem Statement and Analysing basic metrics

Using the monthly information for a segment of drivers for 2019 and 2020 and predict whether a driver will be leaving the company or not based on their attributes like

Demographics (city, age, gender etc.) Tenure information (joining date, Last Date) Historical data regarding the performance of the driver (Quarterly rating, Monthly business acquired, grade, Income), etc'

```
In [ ]: 1
```

```
In [ ]: 1
```

Observations on the shape of data, data types of all the attributes, conversion of categorical attributes to 'category' (If required), missing value detection, statistical summary (10 Points)

```
In [6]: 1 df.shape
```

```
Out[6]: (19104, 14)
```

```
In [7]: 1 # none of the variables have null values  
2 # 3 categorical columns and 6 numerical columns  
3 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 19104 entries, 0 to 19103  
Data columns (total 14 columns):  
 #   Column           Non-Null Count  Dtype     
---  --  
 0   Unnamed: 0        19104 non-null   int64    
 1   MMM-YY            19104 non-null   object    
 2   Driver_ID         19104 non-null   int64    
 3   Age               19043 non-null   float64   
 4   Gender             19052 non-null   float64   
 5   City               19104 non-null   object    
 6   Education_Level   19104 non-null   int64    
 7   Income              19104 non-null   int64    
 8   Dateofjoining     19104 non-null   object    
 9   LastWorkingDate    1616  non-null    object    
 10  Joining_Designation 19104 non-null   int64    
 11  Grade              19104 non-null   int64    
 12  Total_Business_Value 19104 non-null   int64    
 13  Quarterly_Rating   19104 non-null   int64    
dtypes: float64(2), int64(8), object(4)  
memory usage: 2.0+ MB
```

```
In [8]: 1 # percentage of null values for different variables  
2 def check_missing(df):  
3     missing = np.round(df.isnull().sum()/len(df)*100 ,2)  
4     return missing[missing>0]  
5 check_missing(df)  
6 #missing values for source_name and destination_name
```

```
Out[8]: Age          0.32  
Gender       0.27  
LastWorkingDate 91.54  
dtype: float64
```

In [9]:

```

1 # 
2 # will explore further to check this
3 df.describe()

```

Out[9]:

	Unnamed: 0	Driver_ID	Age	Gender	Education_Level	Income	Joining Designation	Grade	Total Business Value	Quarterly Rating
count	19104.000000	19104.000000	19043.000000	19052.000000	19104.000000	19104.000000	19104.000000	19104.000000	1.910400e+04	19104.000000
mean	9551.500000	1415.591133	34.668435	0.418749	1.021671	65652.025126	1.690536	2.252670	5.716621e+05	2.008899
std	5514.994107	810.705321	6.257912	0.493367	0.800167	30914.515344	0.836984	1.026512	1.128312e+06	1.009832
min	0.000000	1.000000	21.000000	0.000000	0.000000	10747.000000	1.000000	1.000000	-6.000000e+06	1.000000
25%	4775.750000	710.000000	30.000000	0.000000	0.000000	42383.000000	1.000000	1.000000	0.000000e+00	1.000000
50%	9551.500000	1417.000000	34.000000	0.000000	1.000000	60087.000000	1.000000	2.000000	2.500000e+05	2.000000
75%	14327.250000	2137.000000	39.000000	1.000000	2.000000	83969.000000	2.000000	3.000000	6.997000e+05	3.000000
max	19103.000000	2788.000000	58.000000	1.000000	2.000000	188418.000000	5.000000	5.000000	3.374772e+07	4.000000

In [10]:

```
1 df.head()
```

Out[10]:

	Unnamed: 0	MMM- YY	Driver_ID	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Total Business Value	Quarterly Rating	
0	0	01/01/19	1	28.0	0.0	C23		2	57387	24/12/18	NaN	1	1	2381060	2
1	1	02/01/19	1	28.0	0.0	C23		2	57387	24/12/18	NaN	1	1	-665480	2
2	2	03/01/19	1	28.0	0.0	C23		2	57387	24/12/18	03/11/19	1	1	0	2
3	3	11/01/20	2	31.0	0.0	C7		2	67016	11/06/20	NaN	2	2	0	1
4	4	12/01/20	2	31.0	0.0	C7		2	67016	11/06/20	NaN	2	2	0	1

In [11]:

```

1 df['MMM-YY'] =pd.to_datetime(df['MMM-YY'])
2
3 df['Dateofjoining'] =pd.to_datetime(df['Dateofjoining'])
4

```

In [12]: 1 df.info()

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 19104 entries, 0 to 19103
Data columns (total 14 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Unnamed: 0        19104 non-null   int64  
 1   MMM-YY            19104 non-null   datetime64[ns]
 2   Driver_ID         19104 non-null   int64  
 3   Age               19043 non-null   float64 
 4   Gender             19052 non-null   float64 
 5   City               19104 non-null   object  
 6   Education_Level    19104 non-null   int64  
 7   Income              19104 non-null   int64  
 8   Dateofjoining      19104 non-null   datetime64[ns]
 9   LastWorkingDate     1616  non-null   object  
 10  Joining_Designation 19104 non-null   int64  
 11  Grade              19104 non-null   int64  
 12  Total_Business_Value 19104 non-null   int64  
 13  Quarterly_Rating    19104 non-null   int64  
dtypes: datetime64[ns](2), float64(2), int64(8), object(2)
memory usage: 2.0+ MB

```

## ▼ Missing Values imputation

In [13]: 1 df['LastWorkingDate']=df['LastWorkingDate'].fillna(0)

In [14]: 1 df.head()

Out[14]:

		Unnamed: 0	MMM- YY	Driver_ID	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Total Business Value	Quarterly Rating
0	0	2019-01-01		1	28.0	0.0	C23		2	57387	2018-12-24	0	1	1	2381060
1	1	2019-02-01		1	28.0	0.0	C23		2	57387	2018-12-24	0	1	1	-665480
2	2	2019-03-01		1	28.0	0.0	C23		2	57387	2018-12-24	03/11/19	1	1	0
3	3	2020-11-01		2	31.0	0.0	C7		2	67016	2020-11-06	0	2	2	0
4	4	2020-12-01		2	31.0	0.0	C7		2	67016	2020-11-06	0	2	2	1

In [15]:

```
1 df[df['Driver_ID']==11]
```

Out[15]:

	Unnamed: 0	MMM- YY	Driver_ID	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Total Business Value	Quarterly Rating	
21	21	2020- 12-01	11	28.0	1.0	C19		2	42172	2020-12-07	0	1	1	0	1

In [16]:

```
1 check_missing(df)
```

Out[16]:

Age	0.32
Gender	0.27
	dtype: float64

In [17]:

```
1 def KNNimpute(df):
2     # define imputer
3     from sklearn.impute import KNNImputer
4     imputer = KNNImputer()
5     imputer.fit(df)
6     return imputer.transform(df)
```

In [18]:

```
1 cat_col_list= [col for col in df.columns if df[col].dtype=='object' ]
```

In [19]:

```
1 for col in cat_col_list:
2     print(f'the number of unique items in {col} is {df[col].nunique()}')
```

the number of unique items in City is 29  
the number of unique items in LastWorkingDate is 494

In [20]:

```
1 import warnings
2 warnings.filterwarnings("ignore")
```

In [21]:

```
1 num_col_list =[col for col in df.columns if df[col].dtype=='float64' or df[col].dtype=='int64']
2 num_col_list =[ 
3     'Age',
4     'Income',
5     'Total Business Value']
```

```
In [22]: 1 num_df = df[num_col_list]
          2 num_df
```

Out[22]:

	Age	Income	Total Business Value
0	28.0	57387	2381060
1	28.0	57387	-665480
2	28.0	57387	0
3	31.0	67016	0
4	31.0	67016	0
...	...	...	...
19099	30.0	70254	740280
19100	30.0	70254	448370
19101	30.0	70254	0
19102	30.0	70254	200420
19103	30.0	70254	411480

19104 rows × 3 columns

```
In [23]: 1 num_df = KNNimpute(num_df)
```

```
In [24]: 1 num_col_list
```

Out[24]: ['Age', 'Income', 'Total Business Value']

```
In [25]: 1 num_df[:5]
```

```
Out[25]: array([[ 2.80000e+01,  5.73870e+04,  2.38106e+06],
       [ 2.80000e+01,  5.73870e+04, -6.65480e+05],
       [ 2.80000e+01,  5.73870e+04,  0.00000e+00],
       [ 3.10000e+01,  6.70160e+04,  0.00000e+00],
       [ 3.10000e+01,  6.70160e+04,  0.00000e+00]])
```

```
In [26]: 1 df[num_col_list] = num_df
```

```
In [27]: 1 check_missing(df)
```

```
Out[27]: Gender      0.27
          dtype: float64
```

```
In [28]: 1 df['Gender'].mode()[0]
```

Out[28]: 0.0

```
In [29]: 1 df['Gender'].fillna(df['Gender'].mode()[0], inplace=True)
```

## ▼ Data Aggregation

```
In [30]: 1 df.head()
```

Out[30]:

		Unnamed: 0	MMM- YY	Driver_ID	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Total Business Value	Quarterly Rating
0	0	2019-01-01		1	28.0	0.0	C23		2 57387.0	2018-12-24	0	1	1	2381060.0	2
1	1	2019-02-01		1	28.0	0.0	C23		2 57387.0	2018-12-24	0	1	1	-665480.0	2
2	2	2019-03-01		1	28.0	0.0	C23		2 57387.0	2018-12-24	03/11/19	1	1	0.0	2
3	3	2020-11-01		2	31.0	0.0	C7		2 67016.0	2020-11-06	0	2	2	0.0	1
4	4	2020-12-01		2	31.0	0.0	C7		2 67016.0	2020-11-06	0	2	2	0.0	1

```
In [31]: 1 df_groupby = df.groupby('Driver_ID')
```

```
In [32]: 1 df_clean = pd.DataFrame()
```

```
In [33]: 1 df_clean[['Age', 'Gender', 'City', 'Education_Level', 'Income', 'Dateofjoining', 'LastWorkingDate', 'Joining Designation', 'Grade',  
2 , 'Quarterly Rating']] = df_groupby[['Age', 'Gender', 'City', 'Education_Level', 'Income', 'Dateofjoini  
3 , 'Quarterly Rating']].last()
```

```
In [ ]:
```

## ▼ Feature Engineering

```
In [34]: 1 df_clean['Total Business Value'] = df.groupby(['Total Business Value']).sum()
2 df_clean['Tenure'] = df.groupby(['MMM-YY']).last()-df.groupby(['MMM-YY']).first()
```

```
In [35]: 1 df_clean['rating_change'] = df.groupby(['Quarterly Rating']).last()-df.groupby(['Quarterly Rating']).first()
2 df_clean['income_change'] = df.groupby(['Income']).last()-df.groupby(['Income']).first()
3
```

```
In [36]: 1 df_clean.head()
```

Out[36]:

	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Quarterly Rating	Total Business Value	Tenure	rating_change	income_chang
Driver_ID														
1	28.0	0.0	C23		2 57387.0	2018-12-24	03/11/19	1	1	2 1715580.0	59 days	0	0	0
2	31.0	0.0	C7		2 67016.0	2020-11-06	0	2	2	1 0.0	30 days	0	0	0
4	43.0	0.0	C13		2 65603.0	2019-12-07	27/04/20	2	2	1 350000.0	122 days	0	0	0
5	29.0	0.0	C9		0 46368.0	2019-01-09	03/07/19	1	1	1 120360.0	59 days	0	0	0
6	31.0	1.0	C11		1 78728.0	2020-07-31	0	3	3	2 1265000.0	122 days	1	0	0

```
In [37]: 1 df_clean['rating_change'].unique()
```

Out[37]: array([ 0, 1, -3, -1, -2, 2, 3], dtype=int64)

```
In [38]: 1 df_clean['rating_change']= df_clean['rating_change'].apply(lambda x: 'Increased' if x>0 else 'Same' if x==0 else "Decreased")
2 df_clean['income_change']= df_clean['income_change'].apply(lambda x: 'Increased' if x>0 else 'Same' if x==0 else "Decreased")
```

In [39]: 1 df\_clean.head()

Out[39]:

	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Quarterly Rating	Total Business Value	Tenure	rating_change	income_change
Driver_ID														
1	28.0	0.0	C23		2 57387.0	2018-12-24	03/11/19	1	1	2	1715580.0	59 days	Same	San
2	31.0	0.0	C7		2 67016.0	2020-11-06	0	2	2	1	0.0	30 days	Same	San
4	43.0	0.0	C13		2 65603.0	2019-12-07	27/04/20	2	2	1	350000.0	122 days	Same	San
5	29.0	0.0	C9		0 46368.0	2019-01-09	03/07/19	1	1	1	120360.0	59 days	Same	San
6	31.0	1.0	C11		1 78728.0	2020-07-31	0	3	3	2	1265000.0	122 days	Increased	San

In [40]: 1 df\_clean['Target']= np.where(df\_clean['LastWorkingDate'] !=0,1,0)  
2 df\_clean.reset\_index(inplace=True)

In [41]: 1 # Creating some more features  
2 df\_clean['Joining Month'] = df\_clean['Dateofjoining'].dt.month  
3 df\_clean['Joining year'] = df\_clean['Dateofjoining'].dt.year  
4  
5 df\_clean['Tenure'] = df\_clean['Tenure'].dt.days

In [42]: 1 df\_clean.head()

Out[42]:

	Driver_ID	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Quarterly Rating	Total Business Value	Tenure	rating_change	income_ch
Driver_ID															
0	1	28.0	0.0	C23		2 57387.0	2018-12-24	03/11/19	1	1	2	1715580.0	59	Same	\$
1	2	31.0	0.0	C7		2 67016.0	2020-11-06	0	2	2	1	0.0	30	Same	\$
2	4	43.0	0.0	C13		2 65603.0	2019-12-07	27/04/20	2	2	1	350000.0	122	Same	\$
3	5	29.0	0.0	C9		0 46368.0	2019-01-09	03/07/19	1	1	1	120360.0	59	Same	\$
4	6	31.0	1.0	C11		1 78728.0	2020-07-31	0	3	3	2	1265000.0	122	Increased	\$

In [43]: 1 df\_clean.head()

Out[43]:

	Driver_ID	Age	Gender	City	Education_Level	Income	Dateofjoining	LastWorkingDate	Joining Designation	Grade	Quarterly Rating	Total Business Value	Tenure	rating_change	income_ch
0	1	28.0	0.0	C23	2	57387.0	2018-12-24	03/11/19	1	1	2	1715580.0	59	Same	\$
1	2	31.0	0.0	C7	2	67016.0	2020-11-06	0	2	2	1	0.0	30	Same	\$
2	4	43.0	0.0	C13	2	65603.0	2019-12-07	27/04/20	2	2	1	350000.0	122	Same	\$
3	5	29.0	0.0	C9	0	46368.0	2019-01-09	03/07/19	1	1	1	120360.0	59	Same	\$
4	6	31.0	1.0	C11	1	78728.0	2020-07-31	0	3	3	2	1265000.0	122	Increased	\$

In [44]: 1 #Dropping unnecessary columns  
2 df\_clean.drop(['Driver\_ID','LastWorkingDate'],inplace=True,axis=1)

In [45]: 1 df\_clean.drop(['Dateofjoining'],inplace=True,axis=1)  
2

## ▼ Data Exploration -

In [46]: 1 df\_clean.head()

Out[46]:

	Age	Gender	City	Education_Level	Income	Joining Designation	Grade	Quarterly Rating	Total Business Value	Tenure	rating_change	income_change	Target	Joining Month	Joining year
0	28.0	0.0	C23	2	57387.0	1	1	2	1715580.0	59	Same	Same	1	12	2018
1	31.0	0.0	C7	2	67016.0	2	2	1	0.0	30	Same	Same	0	11	2020
2	43.0	0.0	C13	2	65603.0	2	2	1	350000.0	122	Same	Same	1	12	2019
3	29.0	0.0	C9	0	46368.0	1	1	1	120360.0	59	Same	Same	1	1	2019
4	31.0	1.0	C11	1	78728.0	3	3	2	1265000.0	122	Increased	Same	0	7	2020

In [ ]:

1

In [47]: 1 df\_clean.groupby(by='Target')['Income'].describe()

Out[47]:

	count	mean	std	min	25%	50%	75%	max
<b>Target</b>								
0	765.0	67662.90719	29578.579157	12938.0	46131.0	64154.0	84573.00	188418.0
1	1616.0	55391.40099	26924.959518	10747.0	36117.5	51630.0	69816.75	167758.0

In [48]: 1 df\_clean['Education\_Level'].value\_counts()

Out[48]: 2 802  
1 795  
0 784  
Name: Education\_Level, dtype: int64

In [49]: 1 df\_clean.groupby(by='Target')['Total Business Value'].describe()

Out[49]:

	count	mean	std	min	25%	50%	75%	max
<b>Target</b>								
0	765.0	9.620626e+06	1.323157e+07	0.0	240920.0	2636210.0	16290800.0	95331060.0
1	1616.0	2.203746e+06	4.717754e+06	-1385530.0	0.0	465025.0	2345702.5	60153830.0

In [50]: 1 df\_clean.groupby(by='Target')['Tenure'].describe()

Out[50]:

	count	mean	std	min	25%	50%	75%	max
<b>Target</b>								
0	765.0	317.666667	274.769928	0.0	61.0	183.0	700.0	700.0
1	1616.0	164.916460	140.248902	0.0	61.0	122.0	214.0	700.0

In [ ]:

1

In [ ]:

1

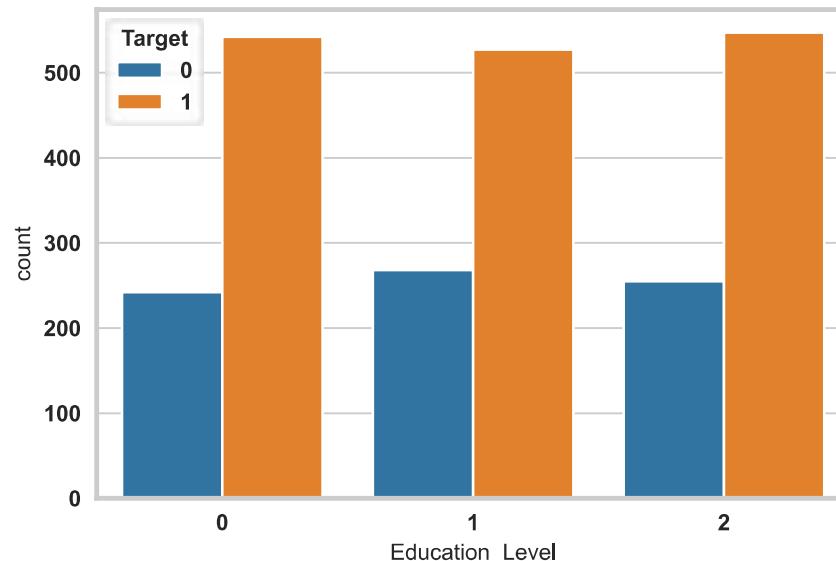
In [ ]:

1

## ▼ Data Visualization

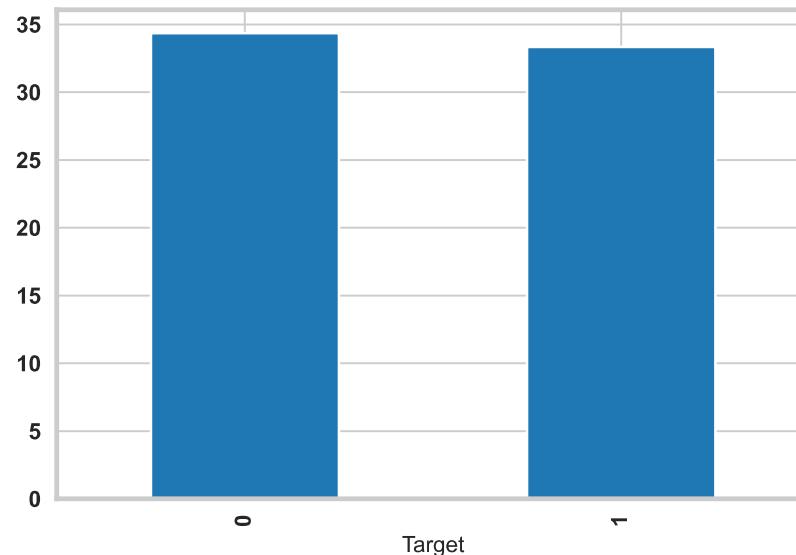
```
In [51]: 1 sns.countplot(data=df_clean,x='Education_Level',hue='Target')
2 plt.xticks()
```

```
Out[51]: (array([0, 1, 2]), [Text(0, 0, '0'), Text(1, 0, '1'), Text(2, 0, '2')])
```



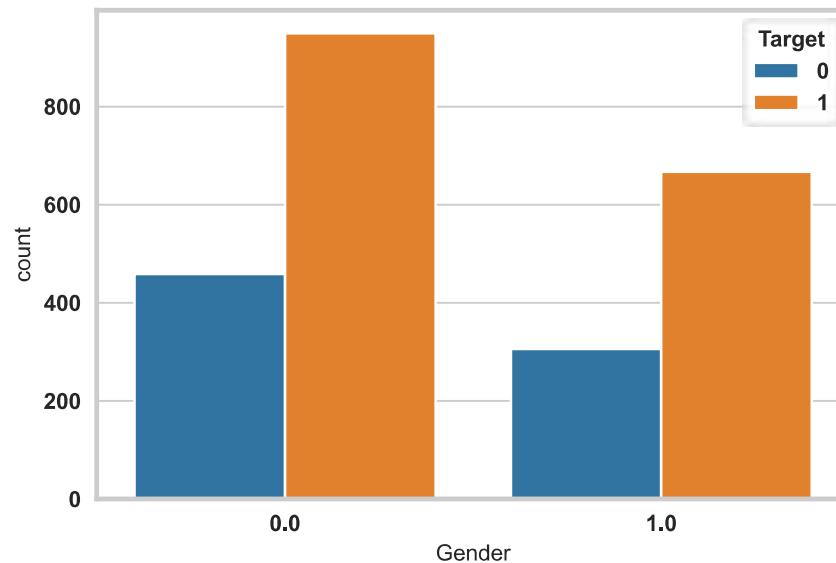
```
In [52]:  
1  
2 df_clean.groupby('Target')['Age'].mean().plot(kind='bar')
```

Out[52]: <AxesSubplot:xlabel='Target'>



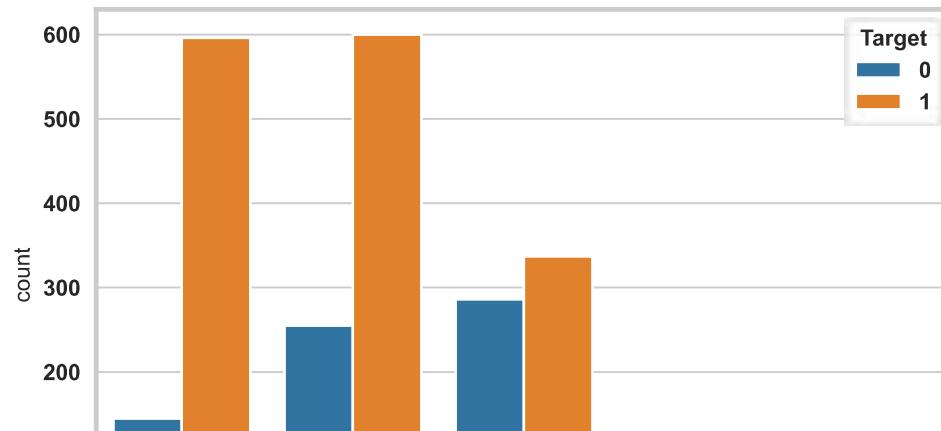
```
In [53]: 1 sns.countplot(data=df_clean,x='Gender',hue='Target')
```

```
Out[53]: <AxesSubplot:xlabel='Gender', ylabel='count'>
```



```
In [54]: 1 plt.figure(figsize=(15, 20))
2
3 plt.subplot(4, 2, 1)
4 sns.countplot(x='Grade', data=df_clean, hue='Target')
5
6 plt.subplot(4, 2, 2)
7 sns.countplot(x='Quarterly Rating', data=df_clean, hue='Target')
8
9 plt.subplot(4, 2, 3)
10 sns.countplot(x='Joining Designation', data=df_clean, hue='Target')
11
12 plt.subplot(4, 2, 4)
13 sns.countplot(x='rating_change', data=df_clean, hue='Target')
14
15 plt.subplot(4, 2, 5)
16 sns.countplot(x='income_change', data=df_clean, hue='Target')
```

```
Out[54]: <AxesSubplot:xlabel='income_change', ylabel='count'>
```



as expected , People who had increased rating change and income change are less likely to churn

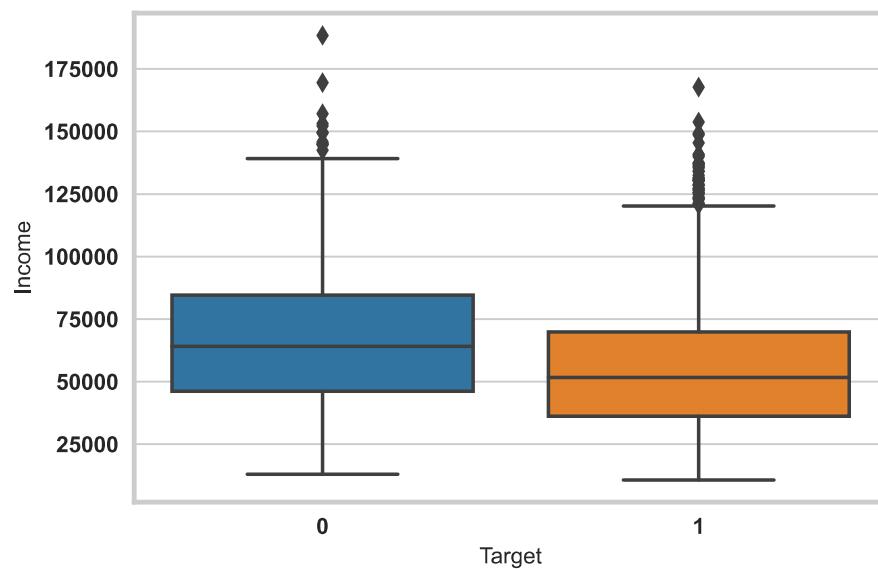
In [ ]:

1

In [55]:

1 sns.boxplot(x='Target',y='Income',data=df\_clean)

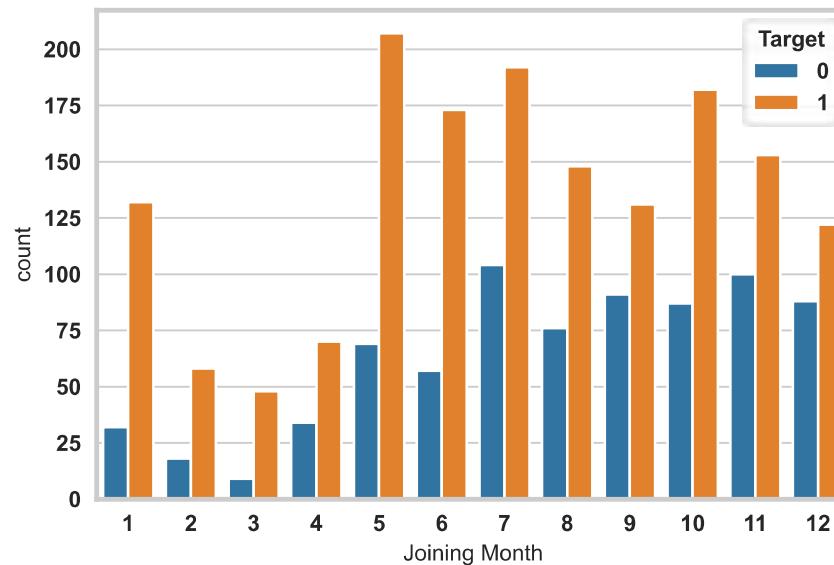
Out[55]: <AxesSubplot:xlabel='Target', ylabel='Income'>



Clealy People who churn have lower median income

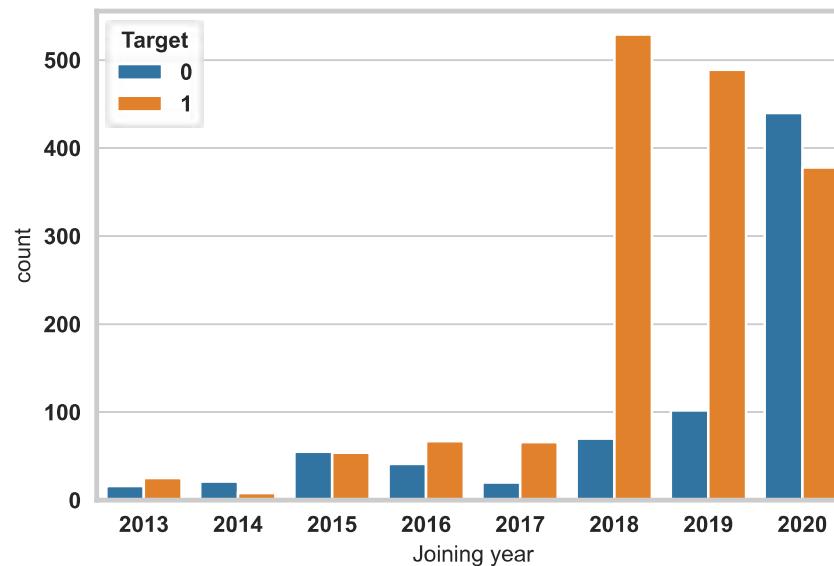
```
In [56]: 1 sns.countplot(data=df_clean,x='Joining Month',hue='Target')
```

Out[56]: <AxesSubplot:xlabel='Joining Month', ylabel='count'>



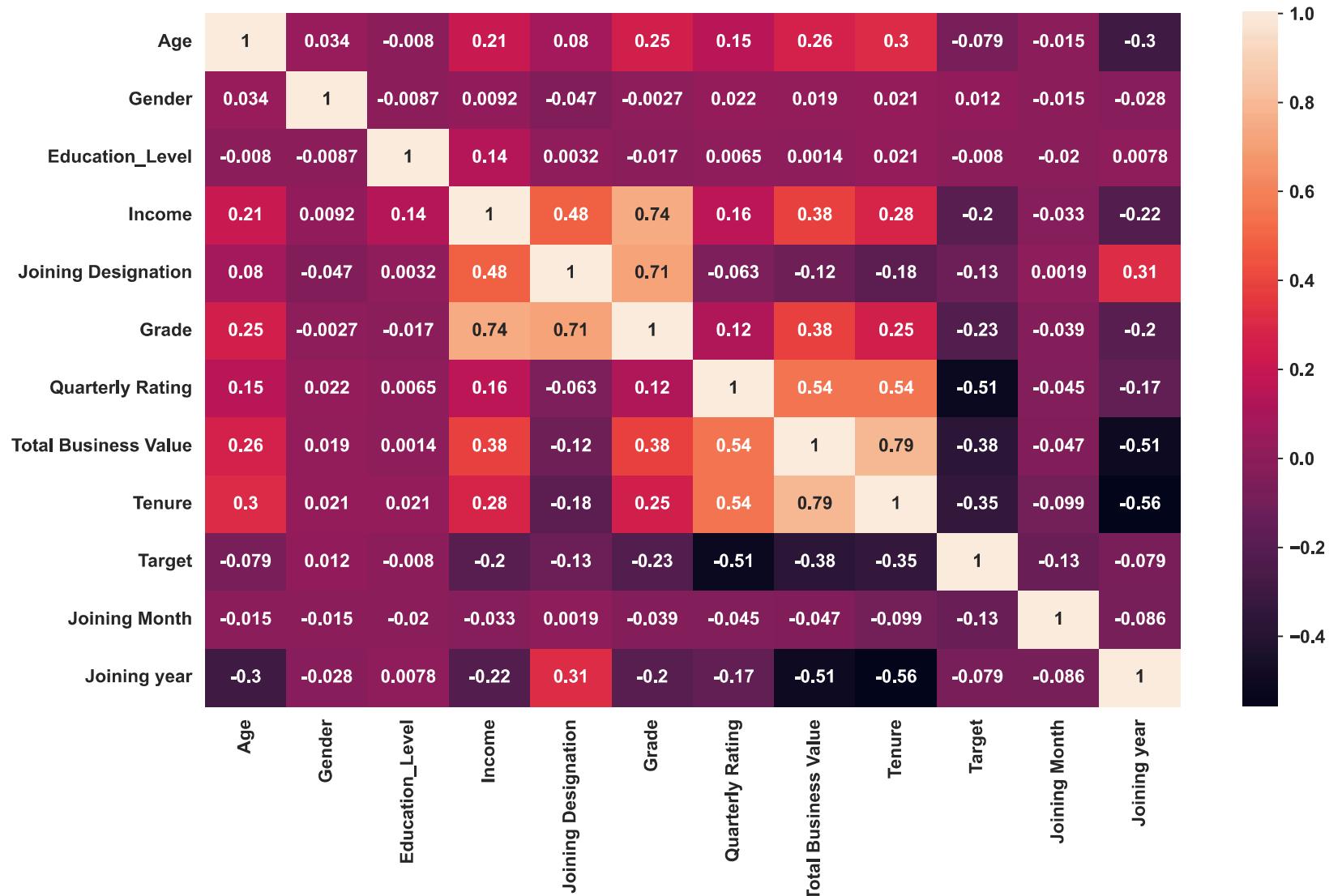
```
In [57]: 1 sns.countplot(data=df_clean,x='Joining year',hue='Target')
```

Out[57]: <AxesSubplot:xlabel='Joining year', ylabel='count'>



```
In [58]: 1 plt.figure(figsize=(12,7))
2 sns.heatmap(df_clean.corr(), annot=True)
3
```

Out[58]: <AxesSubplot:>



No one of the variables seems highly correlated, let's continue for now

In [ ]:

1

## ▼ Encoding & Standardization

In [59]:

```
1 from imblearn.over_sampling import SMOTE
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler, OneHotEncoder
4 from sklearn.pipeline import make_pipeline
5 from sklearn.compose import ColumnTransformer
```

In [60]:

```
1 df_clean.columns
```

Out[60]: Index(['Age', 'Gender', 'City', 'Education\_Level', 'Income',  
 'Joining Designation', 'Grade', 'Quarterly Rating',  
 'Total Business Value', 'Tenure', 'rating\_change', 'income\_change',  
 'Target', 'Joining Month', 'Joining year'],  
 dtype='object')

In [61]:

```
1 categorical_columns = ['City', 'rating_change', 'income_change']
2 non_categorical_columns = ['Age', 'Gender', 'Education_Level', 'Income',
   'Joining Designation', 'Grade', 'Quarterly Rating',
   'Total Business Value', 'Tenure',
   'Joining Month', 'Joining year']
```

In [62]:

```
1 preprocessor = ColumnTransformer(
2     transformers=[
3         ('cat', OneHotEncoder(), categorical_columns), # One-hot encode categorical columns
4         ('num', 'passthrough', non_categorical_columns) # Leave non-categorical columns as-is
5     ])
6 scaler = StandardScaler()
7 pipeline = make_pipeline(preprocessor, scaler)
```

In [63]:

```
1 X = df_clean.drop(['Target'], axis=1)
2 y = df_clean['Target']
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

In [64]:

```
1 X_train_scaled = pipeline.fit_transform(X_train)
2 X_test_scaled = pipeline.transform(X_test)
3
```

## ▼ Class Imbalance

we have class imbalance, lets use SMOTE

```
In [65]: 1 y_train.value_counts()
```

```
Out[65]: 1    1289  
0     615  
Name: Target, dtype: int64
```

```
In [66]: 1 ros = SMOTE()
```

```
In [67]: 1 X_resampled, y_resampled = ros.fit_resample(X_train_scaled, y_train)
```

```
In [ ]:
```

## ▼ Model Building

### ▼ Ensemble - Bagging Algorithm Random forest

```
In [68]: 1 from sklearn.ensemble import RandomForestClassifier  
2 from sklearn.model_selection import KFold, cross_validate
```

```
In [69]: 1 rfc = RandomForestClassifier(random_state=7, n_estimators=100) ## i am letting each DT overfit  
2  
3 kfold = KFold(n_splits=10)  
4 cv_acc_results = cross_validate(rfc, X_resampled, y_resampled, cv = kfold,  
5                                 scoring = 'accuracy',  
6                                 return_train_score = True)  
7  
8 print(f"K-Fold Accuracy Mean: Train: {cv_acc_results['train_score'].mean().round(3)*100} Validation: {cv_acc_results['test_score'].mean().round(3)*100}")
```

K-Fold Accuracy Mean: Train: 100.0 Validation: 93.60000000000001

```
In [70]:  
1 # lets try Grid search to tune hyperparameters  
2  
3 from sklearn.model_selection import GridSearchCV , RandomizedSearchCV  
4 params = {  
5     'max_depth' : [7,10, 15, 20],  
6     'n_estimators' : [100,200,300,400],  
7     'max_features' : [5, 10, 20, 30]  
8 }  
9
```

```
In [71]: 1 model = RandomizedSearchCV(rfc, param_distributions = params, scoring="accuracy", cv = 3, verbose=5)
2
3 # model
4 model.fit(X_resampled, np.ravel(y_resampled))
```

Fitting 3 folds for each of 10 candidates, totalling 30 fits

```
[CV 1/3] END max_depth=10, max_features=5, n_estimators=400;, score=0.908 total time= 2.2s
[CV 2/3] END max_depth=10, max_features=5, n_estimators=400;, score=0.929 total time= 2.2s
[CV 3/3] END max_depth=10, max_features=5, n_estimators=400;, score=0.939 total time= 2.2s
[CV 1/3] END max_depth=10, max_features=10, n_estimators=300;, score=0.909 total time= 2.2s
[CV 2/3] END max_depth=10, max_features=10, n_estimators=300;, score=0.932 total time= 2.1s
[CV 3/3] END max_depth=10, max_features=10, n_estimators=300;, score=0.942 total time= 2.0s
[CV 1/3] END max_depth=20, max_features=5, n_estimators=400;, score=0.903 total time= 2.5s
[CV 2/3] END max_depth=20, max_features=5, n_estimators=400;, score=0.928 total time= 2.4s
[CV 3/3] END max_depth=20, max_features=5, n_estimators=400;, score=0.951 total time= 2.5s
[CV 1/3] END max_depth=20, max_features=30, n_estimators=400;, score=0.902 total time= 5.3s
[CV 2/3] END max_depth=20, max_features=30, n_estimators=400;, score=0.929 total time= 5.1s
[CV 3/3] END max_depth=20, max_features=30, n_estimators=400;, score=0.959 total time= 5.1s
[CV 1/3] END max_depth=15, max_features=20, n_estimators=100;, score=0.910 total time= 1.0s
[CV 2/3] END max_depth=15, max_features=20, n_estimators=100;, score=0.931 total time= 0.9s
[CV 3/3] END max_depth=15, max_features=20, n_estimators=100;, score=0.956 total time= 0.9s
[CV 1/3] END max_depth=15, max_features=30, n_estimators=100;, score=0.899 total time= 1.4s
[CV 2/3] END max_depth=15, max_features=30, n_estimators=100;, score=0.929 total time= 1.2s
[CV 3/3] END max_depth=15, max_features=30, n_estimators=100;, score=0.958 total time= 1.2s
[CV 1/3] END max_depth=10, max_features=10, n_estimators=200;, score=0.912 total time= 1.5s
[CV 2/3] END max_depth=10, max_features=10, n_estimators=200;, score=0.930 total time= 1.4s
[CV 3/3] END max_depth=10, max_features=10, n_estimators=200;, score=0.943 total time= 1.3s
[CV 1/3] END max_depth=20, max_features=20, n_estimators=300;, score=0.908 total time= 3.1s
[CV 2/3] END max_depth=20, max_features=20, n_estimators=300;, score=0.931 total time= 3.0s
[CV 3/3] END max_depth=20, max_features=20, n_estimators=300;, score=0.959 total time= 3.0s
[CV 1/3] END max_depth=10, max_features=10, n_estimators=400;, score=0.907 total time= 2.9s
[CV 2/3] END max_depth=10, max_features=10, n_estimators=400;, score=0.934 total time= 2.7s
[CV 3/3] END max_depth=10, max_features=10, n_estimators=400;, score=0.945 total time= 2.8s
[CV 1/3] END max_depth=20, max_features=30, n_estimators=300;, score=0.899 total time= 4.0s
[CV 2/3] END max_depth=20, max_features=30, n_estimators=300;, score=0.927 total time= 3.9s
[CV 3/3] END max_depth=20, max_features=30, n_estimators=300;, score=0.959 total time= 3.8s
```

```
Out[71]: RandomizedSearchCV(cv=3, estimator=RandomForestClassifier(random_state=7),
                           param_distributions={'max_depth': [7, 10, 15, 20],
                                                'max_features': [5, 10, 20, 30],
                                                'n_estimators': [100, 200, 300, 400]},
                           scoring='accuracy', verbose=5)
```

```
In [72]: 1 model.best_params_
```

```
Out[72]: {'n_estimators': 300, 'max_features': 20, 'max_depth': 20}
```

```
In [73]: 1 print(model.best_score_)
```

0.9329033218723772

```
In [74]: 1 best_model =model.best_estimator_
```

```
In [79]: 1 from sklearn.metrics import accuracy_score,confusion_matrix, classification_report, roc_curve, roc_auc_score, precision_recall_curve
2
3 y_pred = best_model.predict(X_test_scaled)
4 accuracy_score(y_test,y_pred)
5
```

Out[79]: 0.9329140461215933

```
In [80]: 1 # we have the accuracy of 93 % on test data which is pretty decent ,
```

```
In [81]: 1 y_prob = best_model.predict_proba(X_test_scaled)
2
```

```
In [82]: 1 # Evaluate performance using accuracy score and confusion matrix
2 print(f"Accuracy score: {round(accuracy_score(y_test, y_pred), 3)}\n")
3 print(f"Confusion matrix\n {confusion_matrix(y_test, y_pred)}\n")
4 print(f"Classification report\n {classification_report(y_test, y_pred)}\n")
5 print(f"ROC-AUC Score\n {round(roc_auc_score(y_test, y_prob[:,1]), 3)}")
```

Accuracy score: 0.933

Confusion matrix

```
[[133 17]
 [ 15 312]]
```

Classification report

	precision	recall	f1-score	support
0	0.90	0.89	0.89	150
1	0.95	0.95	0.95	327
accuracy			0.93	477
macro avg	0.92	0.92	0.92	477
weighted avg	0.93	0.93	0.93	477

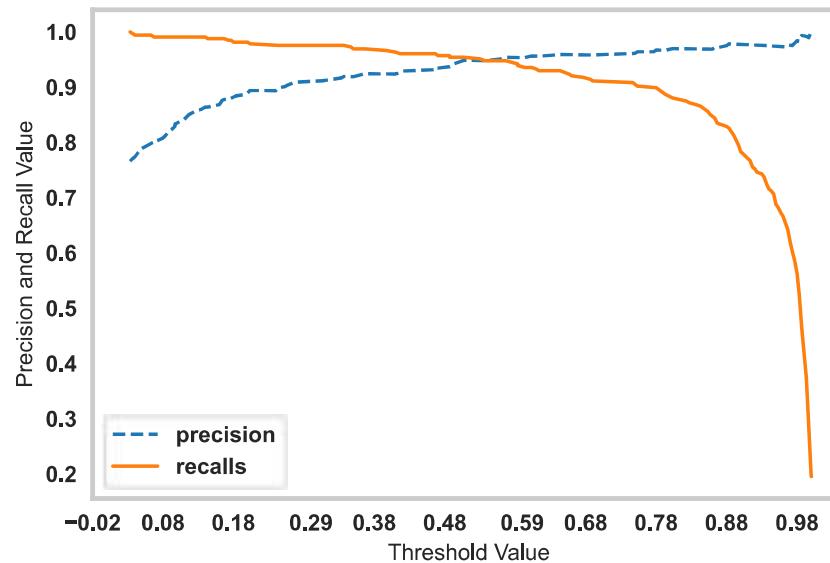
ROC-AUC Score

0.964

In [ ]:

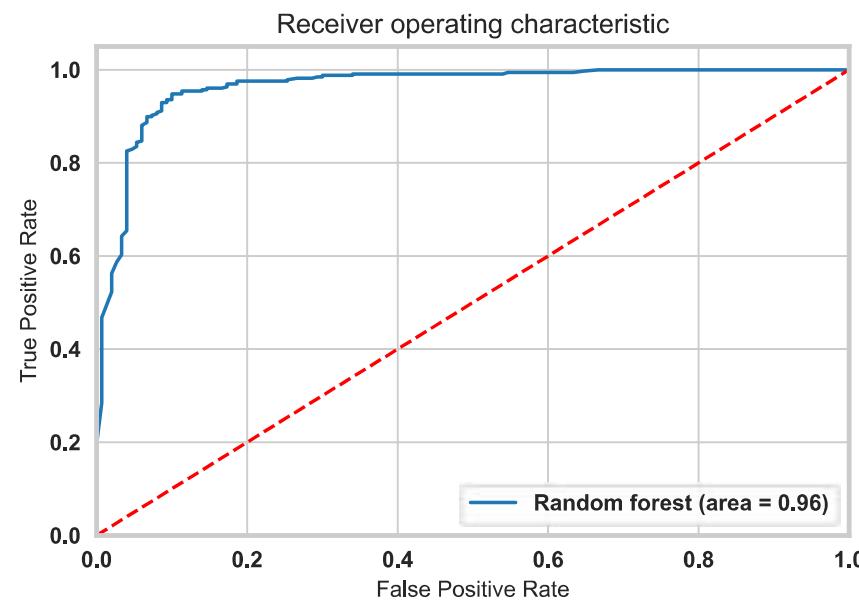
1

```
In [83]: 1 def precision_recall_curve_plot(y_test, pred_proba_c1):
2     precisions, recalls, thresholds = precision_recall_curve(y_test, pred_proba_c1)
3
4     threshold_boundary = thresholds.shape[0]
5     # plot precision
6     plt.plot(thresholds, precisions[0:threshold_boundary], linestyle='--', label='precision')
7     # plot recall
8     plt.plot(thresholds, recalls[0:threshold_boundary], label='recalls')
9
10    start, end = plt.xlim()
11    plt.xticks(np.round(np.arange(start, end, 0.1), 2))
12
13    plt.xlabel('Threshold Value'); plt.ylabel('Precision and Recall Value')
14    plt.legend(); plt.grid()
15    plt.show()
16
17 precision_recall_curve_plot(y_test, y_prob[:,1])
```



In [84]:

```
1 roc_auc = roc_auc_score(y_test, y_prob[:,1])
2 fpr, tpr, thresholds = roc_curve(y_test, y_prob[:,1])
3 plt.figure()
4 plt.plot(fpr, tpr, label='Random forest (area = %0.2f)' % roc_auc)
5 plt.plot([0, 1], [0, 1], 'r--')
6 plt.xlim([0.0, 1.0])
7 plt.ylim([0.0, 1.05])
8 plt.xlabel('False Positive Rate')
9 plt.ylabel('True Positive Rate')
10 plt.title('Receiver operating characteristic')
11 plt.legend(loc="lower right")
12 plt.savefig('Log_ROC')
13 plt.show()
```



In [ ]:

1

In [ ]:

1

In [ ]:

1

In [ ]:

1

## ▼ Boosting - LightGBM & XGboost

### ▼ XGBoost

```
In [85]: 1 !pip install xgboost
```

```
Requirement already satisfied: xgboost in c:\users\gurme\anaconda3\lib\site-packages (2.0.0)
Requirement already satisfied: numpy in c:\users\gurme\anaconda3\lib\site-packages (from xgboost) (1.22.4)
Requirement already satisfied: scipy in c:\users\gurme\anaconda3\lib\site-packages (from xgboost) (1.7.3)
```

```
In [86]: 1 from xgboost import XGBClassifier
2 from sklearn.metrics import confusion_matrix, classification_report, roc_curve, roc_auc_score, precision_recall_curve
3
```

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [87]: 1 params = {
2     'min_child_weight': [1, 5, 10],
3     'gamma': [0.5, 1, 1.5, 2, 5],
4     'subsample': [0.6, 0.8, 1.0],
5     'colsample_bytree': [0.6, 0.8, 1.0],
6     'max_depth': [3, 4, 5]
7 }
```

```
In [88]: 1 xgb = XGBClassifier(learning_rate=0.02, n_estimators=600, objective='binary:logistic',
2                         silent=True, nthread=1)
```

```
In [89]: 1 folds = 3
2 param_comb = 5
3 from sklearn.model_selection import StratifiedKFold
4 skf = StratifiedKFold(n_splits=folds, shuffle = True, random_state = 1001)
5
6 random_search = RandomizedSearchCV(xgb, param_distributions=params, n_iter=param_comb, scoring='accuracy',
7                                     n_jobs=4, cv=3, verbose=3, random_state=1001 )
```

```
In [90]:  
1 random_search.fit(X_resampled, y_resampled)
```

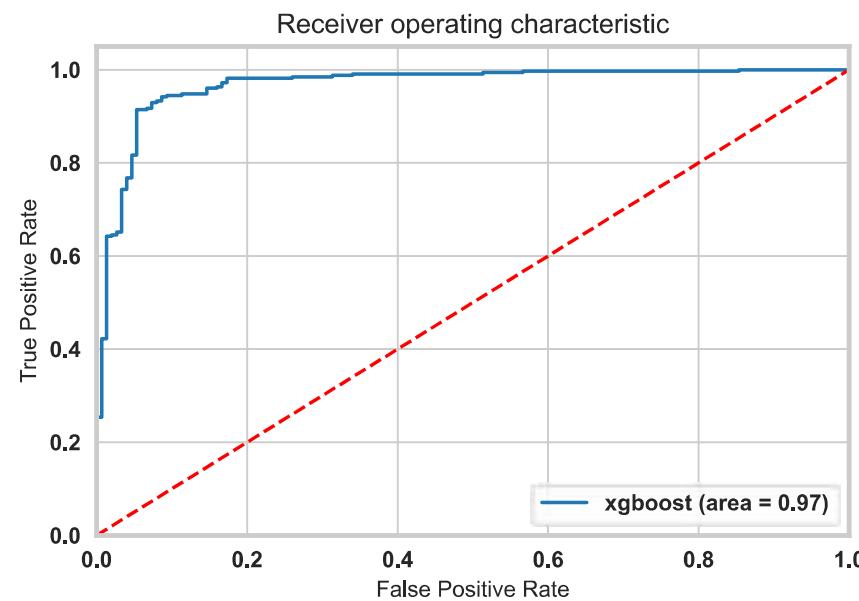
Fitting 3 folds for each of 5 candidates, totalling 15 fits

```
Out[90]: RandomizedSearchCV(cv=3,  
                           estimator=XGBClassifier(base_score=None, booster=None,  
                           callbacks=None,  
                           colsample_bylevel=None,  
                           colsample_bynode=None,  
                           colsample_bytree=None, device=None,  
                           early_stopping_rounds=None,  
                           enable_categorical=False,  
                           eval_metric=None, feature_types=None,  
                           gamma=None, grow_policy=None,  
                           importance_type=None,  
                           interaction_constraints=None,  
                           learning_rate...  
                           min_child_weight=None, missing=nan,  
                           monotone_constraints=None,  
                           multi_strategy=None,  
                           n_estimators=600, n_jobs=None,  
                           nthread=1, num_parallel_tree=None, ...),  
                           n_iter=5, n_jobs=4,  
                           param_distributions={'colsample_bytree': [0.6, 0.8, 1.0],  
                           'gamma': [0.5, 1, 1.5, 2, 5],  
                           'max_depth': [3, 4, 5],  
                           'min_child_weight': [1, 5, 10],  
                           'subsample': [0.6, 0.8, 1.0]},  
                           random_state=1001, scoring='accuracy', verbose=3)
```

```
In [91]:  
1 print('\n All results:')  
2 print(random_search.cv_results_)  
3 print('\n Best estimator:')  
4 print(random_search.best_estimator_)  
5 print('\n Best normalized gini score for %d-fold search with %d parameter combinations:' % (folds, param_comb))  
6 print(random_search.best_score_ * 2 - 1)  
7 print('\n Best hyperparameters:')  
8 print(random_search.best_params_)  
9 results = pd.DataFrame(random_search.cv_results_)  
10 results.to_csv('xgb-random-grid-search-results-01.csv', index=False)
```

In [97]:

```
1 roc_auc = roc_auc_score(y_test, y_prob[:,1])
2 fpr, tpr, thresholds = roc_curve(y_test, y_prob[:,1])
3 plt.figure()
4 plt.plot(fpr, tpr, label='xgboost (area = %0.2f)' % roc_auc)
5 plt.plot([0, 1], [0, 1], 'r--')
6 plt.xlim([0.0, 1.0])
7 plt.ylim([0.0, 1.05])
8 plt.xlabel('False Positive Rate')
9 plt.ylabel('True Positive Rate')
10 plt.title('Receiver operating characteristic')
11 plt.legend(loc="lower right")
12 plt.savefig('Log_ROC')
13 plt.show()
```



▼ LightBGM

In [98]:

```
1 param_grid = {
2     "num_leaves": [31, 63, 127],
3     "max_depth": [-1, 3, 5],
4     "subsample": [0.8, 1.0],
5     "colsample_bytree": [0.8, 1.0]
6 }
```

```
In [99]: 1 !pip install lightgbm
```

```
Requirement already satisfied: lightgbm in c:\users\gurme\anaconda3\lib\site-packages (4.1.0)
Requirement already satisfied: numpy in c:\users\gurme\anaconda3\lib\site-packages (from lightgbm) (1.22.4)
Requirement already satisfied: scipy in c:\users\gurme\anaconda3\lib\site-packages (from lightgbm) (1.7.3)
```

```
In [100]: 1 import lightgbm as lgb
```

```
In [101]: 1 # Create an instance of the LightGBM classifier
2 lgbm = lgb.LGBMClassifier(objective="binary", metric="accuracy", random_state=42)
```

```
In [102]: 1 random_search = RandomizedSearchCV(lgbm, param_distributions=param_grid, n_iter=5, scoring='accuracy',
2                                         n_jobs=4, cv=3, verbose=3, random_state=1001 )
```

```
In [103]: 1 random_search.fit(X_resampled,y_resampled)
```

```
Out[103]: RandomizedSearchCV(cv=3,
                           estimator=LGBMClassifier(metric='accuracy',
                                                     objective='binary',
                                                     random_state=42),
                           n_iter=5, n_jobs=4,
                           param_distributions={'colsample_bytree': [0.8, 1.0],
                                                'max_depth': [-1, 3, 5],
                                                'num_leaves': [31, 63, 127],
                                                'subsample': [0.8, 1.0]},
                           random_state=1001, scoring='accuracy', verbose=3)
```

In [104]:

```

1 print('\n All results:')
2 print(random_search.cv_results_)
3 print('\n Best estimator:')
4 print(random_search.best_estimator_)
5 print('\n Best normalized gini score for %d-fold search with %d parameter combinations:' % (folds, param_comb))
6 print(random_search.best_score_* 2 - 1)
7 print('\n Best hyperparameters:')
8 print(random_search.best_params_)
9 results = pd.DataFrame(random_search.cv_results_)
10 results.to_csv('xgb-random-grid-search-results-01.csv', index=False)

```

All results:

```

{'mean_fit_time': array([0.71018442, 0.27947617, 0.2317001 , 0.22844442, 0.17791533]), 'std_fit_time': array([0.02739269, 0.09187242,
0.02029487, 0.01629035, 0.04128941]), 'mean_score_time': array([0.01041738, 0.01041834, 0.00838017, 0.          , 0.01041722]), 'std_sco
re_time': array([0.00736662 , 0.00736688, 0.00642829, 0.          , 0.00736609]), 'param_subsample': masked_array(data=[1.0, 0.8, 0.8, 1.
0, 1.0],
mask=[False, False, False, False, False],
fill_value='?'),
dtype=object), 'param_num_leaves': masked_array(data=[63, 127, 127, 127, 127],
mask=[False, False, False, False, False],
fill_value='?'),
dtype=object), 'param_max_depth': masked_array(data=[5, 3, 3, 3, 3],
mask=[False, False, False, False, False],
fill_value='?'),
dtype=object), 'param_colsample_bytree': masked_array(data=[0.8, 1.0, 0.8, 0.8, 1.0],
mask=[False, False, False, False, False],
fill_value='?'),
dtype=object), 'params': [{{'subsample': 1.0, 'num_leaves': 63, 'max_depth': 5, 'colsample_bytree': 0.8}, {'subsample': 0.
8, 'num_leaves': 127, 'max_depth': 3, 'colsample_bytree': 1.0}, {'subsample': 0.8, 'num_leaves': 127, 'max_depth': 3, 'colsample_bytre
e': 0.8}, {'subsample': 1.0, 'num_leaves': 127, 'max_depth': 3, 'colsample_bytree': 0.8}, {'subsample': 1.0, 'num_leaves': 127, 'max_d
epth': 3, 'colsample_bytree': 1.0}], 'split0_test_score': array([0.90116279, 0.90813953, 0.91744186, 0.91744186, 0.90813953]), 'split1
_test_score': array([0.93946449, 0.93597206, 0.94062864, 0.94062864, 0.93597206]), 'split2_test_score': array([0.96274738, 0.95343423,
0.95925495, 0.95925495, 0.95343423]), 'mean_test_score': array([0.93445822, 0.93251527, 0.93910848, 0.93910848, 0.93251527]), 'std_tes
t_score': array([0.02538979, 0.01865233, 0.01710393, 0.01710393, 0.01865233]), 'rank_test_score': array([3, 4, 1, 1, 4])}

```

Best estimator:

```

LGBMClassifier(colsample_bytree=0.8, max_depth=3, metric='accuracy',
               num_leaves=127, objective='binary', random_state=42,
               subsample=0.8)

```

Best normalized gini score for 3-fold search with 5 parameter combinations:  
0.8782169640198176

Best hyperparameters:

```

{'subsample': 0.8, 'num_leaves': 127, 'max_depth': 3, 'colsample_bytree': 0.8}

```

```
In [105]: 1 best_model = random_search.best_estimator_
```

```
In [106]: 1 y_pred = best_model.predict(X_test_scaled)
2 accuracy_score(y_test,y_pred)
```

Out[106]: 0.9203354297693921

```
In [107]: 1 y_prob = best_model.predict_proba(X_test_scaled)
2
```

```
In [108]: 1 # Evaluate performance using accuracy score and confusion matrix
2 print(f"Accuracy score: {round(accuracy_score(y_test, y_pred), 3)}\n")
3 print(f"Confusion matrix\n {confusion_matrix(y_test, y_pred)}\n")
4 print(f"Classification report\n {classification_report(y_test, y_pred)}\n")
5 print(f"ROC-AUC Score\n {round(roc_auc_score(y_test, y_prob[:,1]), 3)}")
```

Accuracy score: 0.92

Confusion matrix  
[[131 19]  
 [ 19 308]]

Classification report

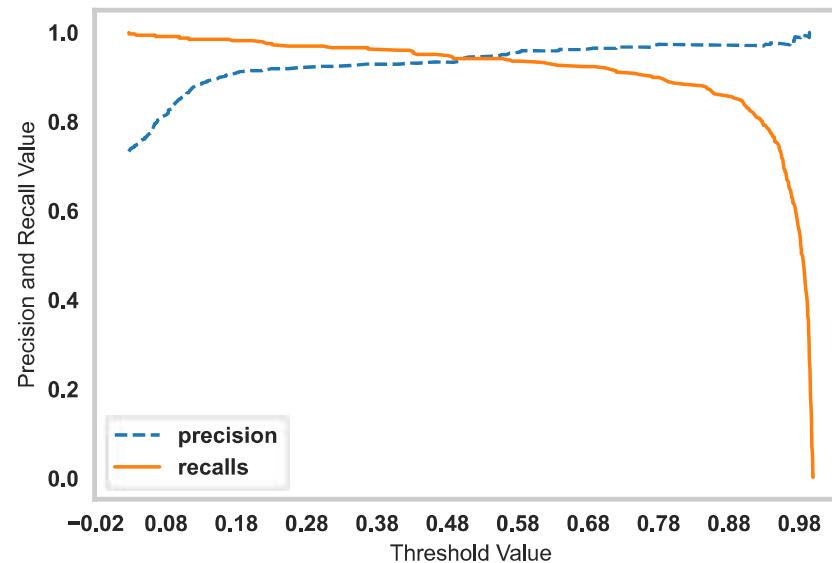
	precision	recall	f1-score	support
0	0.87	0.87	0.87	150
1	0.94	0.94	0.94	327
accuracy			0.92	477
macro avg	0.91	0.91	0.91	477
weighted avg	0.92	0.92	0.92	477

ROC-AUC Score  
0.965

```
In [ ]:
```

1

```
In [109]: 1 def precision_recall_curve_plot(y_test, pred_proba_c1):
2     precisions, recalls, thresholds = precision_recall_curve(y_test, pred_proba_c1)
3
4     threshold_boundary = thresholds.shape[0]
5     # plot precision
6     plt.plot(thresholds, precisions[0:threshold_boundary], linestyle='--', label='precision')
7     # plot recall
8     plt.plot(thresholds, recalls[0:threshold_boundary], label='recalls')
9
10    start, end = plt.xlim()
11    plt.xticks(np.round(np.arange(start, end, 0.1), 2))
12
13    plt.xlabel('Threshold Value'); plt.ylabel('Precision and Recall Value')
14    plt.legend(); plt.grid()
15    plt.show()
16
17 precision_recall_curve_plot(y_test, y_prob[:,1])
```

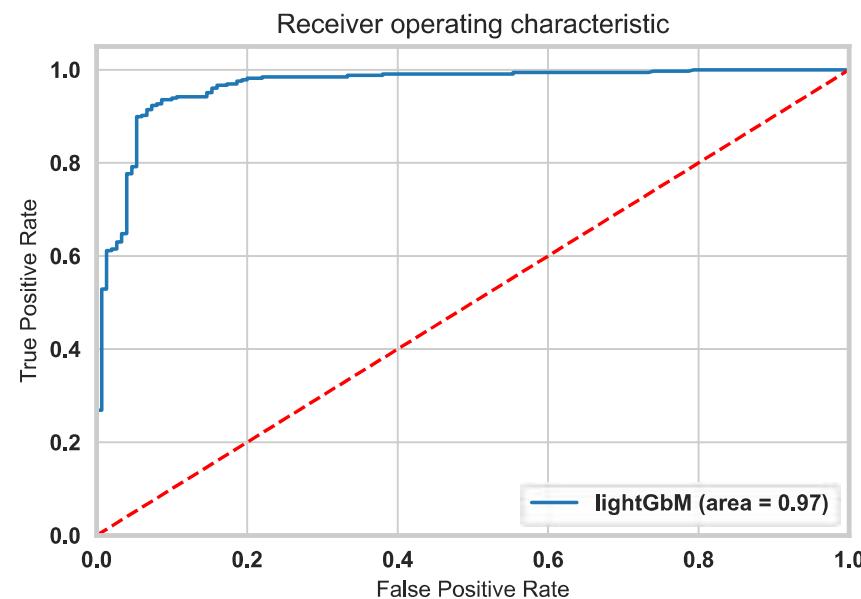


In [110]:

```

1 roc_auc = roc_auc_score(y_test, y_prob[:,1])
2 fpr, tpr, thresholds = roc_curve(y_test, y_prob[:,1])
3 plt.figure()
4 plt.plot(fpr, tpr, label='lightGbM (area = %0.2f)' % roc_auc)
5 plt.plot([0, 1], [0, 1], 'r--')
6 plt.xlim([0.0, 1.0])
7 plt.ylim([0.0, 1.05])
8 plt.xlabel('False Positive Rate')
9 plt.ylabel('True Positive Rate')
10 plt.title('Receiver operating characteristic')
11 plt.legend(loc="lower right")
12 plt.savefig('Log_ROC')
13 plt.show()

```



In [ ]:

1

So XGBoost gives the best performance out of all the three models with ROC-AUC score of 97%

In [ ]:

1

## ## Observation

- 1 People who had increased rating change and income change are less likely to churn
- 2 People who have less median income are less likely to churn
- 3
- 4
- 5

```
6 people with good grade are less likely to churn  
7  
8 People getting good quarterly ratings are less likely to churn and vice versa  
9  
10 People with better designation are less likely to shurn and vice versa  
11  
12 We are gettting a good precision, recall and F1 score for the all the models we tried.  
13  
14 Best model that was selected was LightGBM although there was only marginal increaese in the accuracy and other metrics.  
15  
16  
17  
18
```

In [ ]:

1

## ## Recommendations

```
1 Focus on Happy Customers: Pay extra attention to customers who are reporting increased satisfaction (higher ratings) and income  
2 growth. Engage with them, offer incentives, and build loyalty programs.  
3  
4 Support Lower-Income Customers: Customers with lower median incomes tend to stay longer. Create special offers or tailored  
5 services to retain their loyalty.  
6 Reward Good Grades: Encourage and reward customers who consistently achieve good grades or ratings. This can reinforce their  
7 loyalty.  
8 Quarterly Ratings: Keep an eye on quarterly ratings. Reach out to customers with high ratings to nurture their loyalty and  
9 address issues with lower-rated customers promptly.  
10 Value Customers with Better Designations: Offer exclusive benefits to customers with prestigious designations to make them feel  
11 valued.  
12 Opt for Efficient Models: Stick with LightGBM since it's efficient, even if it only marginally improved accuracy. This will  
13 enable faster predictions and quicker action to prevent churn.  
14 Segment and Personalize: Create customer segments based on these factors and tailor retention strategies to each group.  
15  
16 Collect Feedback: Gather feedback from churned customers to understand why they left and make improvements based on their input.  
17  
18 Continuous Monitoring: Keep an eye on customer behavior and model performance. Adapt strategies as needed to minimize churn  
19
```

In [ ]:

1