# What is Employee Turnover?

Employee Turnover or Employee Turnover ratio is the measurement of the total number of employees who leave an organization in a particular year. Employee Turnover Prediction means to predict whether an employee is going to leave the organization in the coming period.

A Company uses this predictive analysis to measure how many employees they will need if the potential employees will leave their organization. A company also uses this predictive analysis to make the workplace better for employees by understanding the core reasons for the high turnover ratio.

## Data Preprocessing

Now let's dive into the data to move further with this project on Employee Turnover Prediction. You can download the dataset I have used in this article below.

In [26]:
```python
import pandas as pd
hr = pd.read_csv('HR.csv')
col_names = hr.columns.tolist()
print("Column names:")
print(col_names)
print("\nSample data:")
hr.head()
```

```
Column names:
['satisfaction_level', 'last_evaluation', 'number_project', 'average_montly_hours', 'time_spend_company', 'Work_a
ccident', 'left', 'promotion_last_5years', 'sales', 'salary']

Sample data:
```

Out[26]:

| | satisfaction_level | last_evaluation | number_project | average_montly_hours | time_spend_company | Work_accident | left | promotion_last_5years | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.38 | 0.53 | 2 | 157 | 3 | 0 | 1 | 0 | |
| 1 | 0.80 | 0.86 | 5 | 262 | 6 | 0 | 1 | 0 | |
| 2 | 0.11 | 0.88 | 7 | 272 | 4 | 0 | 1 | 0 | |
| 3 | 0.72 | 0.87 | 5 | 223 | 5 | 0 | 1 | 0 | |
| 4 | 0.37 | 0.52 | 2 | 159 | 3 | 0 | 1 | 0 | |

# Rename column name from "sales" to "department":

In [4]:
```python
hr=hr.rename(columns = {'sales':'department'})
```

In [5]:
```python
hr.shape
```

Out[5]:
```
(14999, 10)
```

In [6]:
```python
hr['department'].unique()
```

Out[6]:
```
array(['sales', 'accounting', 'hr', 'technical', 'support', 'management',
       'IT', 'product_mng', 'marketing', 'RandD'], dtype=object)
```

In [7]:
```python
import numpy as np
hr['department']=np.where(hr['department'] =='support', 'technical', hr['department'])
hr['department']=np.where(hr['department'] =='IT', 'technical', hr['department'])
```

## Creating Variables for Categorical Variables

As there are two categorical variables (department, salary) in the dataset and they need to be

converted to dummy variables before they can be used for modelling

In [8]:
```python
cat_vars=['department','salary']
for var in cat_vars:
    cat_list='var'+'_'+var
    cat_list = pd.get_dummies(hr[var], prefix=var)
    hr1=hr.join(cat_list)
    hr=hr1
```

In [9]:
```python
hr.drop(hr.columns[[8, 9]], axis=1, inplace=True)
hr.columns.values
```

Out[9]:
```
array(['satisfaction_level', 'last_evaluation', 'number_project',
       'average_montly_hours', 'time_spend_company', 'Work_accident',
       'left', 'promotion_last_5years', 'department_RandD',
       'department_accounting', 'department_hr', 'department_management',
       'department_marketing', 'department_product_mng',
       'department_sales', 'department_technical', 'salary_high',
       'salary_low', 'salary_medium'], dtype=object)
```

In [10]:
```python
hr_vars=hr.columns.values.tolist()
y=['left']
X=[i for i in hr_vars if i not in y]
```

# Feature Selection for Employee Turnover Prediction

Let's use the feature selection method to decide which variables are the best option that can predict employee turnover with great accuracy. There are a total of 18 columns in X, and now let's see how we can select about 10 from them

In [11]:
```python
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
rfe = RFE(model, 10)
rfe = rfe.fit(hr[X], hr[y])
print(rfe.support_)
print(rfe.ranking_)
```

```
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:70: FutureWarning: Pass n_features_to_sele
ct=10 as keyword args. From version 1.0 (renaming of 0.25) passing these as positional arguments will result in a
n error
  warnings.warn(f"Pass {args_msg} as keyword args. From version "
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector
y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel
().
  return f(*args, **kwargs)
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:763: ConvergenceWarning: lbfgs faile
d to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector
y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel
().
  return f(*args, **kwargs)
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector
y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel
().
  return f(*args, **kwargs)
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector
y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel
().
  return f(*args, **kwargs)
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector
y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel
().
  return f(*args, **kwargs)
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector
y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel
().
  return f(*args, **kwargs)
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector
y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel
```

```
().
  return f(*args, **kwargs)
```

```
[ True  True False False  True  True  True  True False  True  True False
 False False False  True  True False]
[1 1 3 9 1 1 1 1 5 1 1 6 8 7 4 1 1 2]
```

In [12]:
```python
cols=['satisfaction_level', 'last_evaluation', 'time_spend_company', 'Work_accident', 'promotion_last_5years',
      'department_RandD', 'department_hr', 'department_management', 'salary_high', 'salary_low']
X=hr[cols]
y=hr['left']
```

# Logistic Regression Model to Predict Employee Turnover

In [15]:
```python
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
from sklearn.linear_model import LogisticRegression
from sklearn import metrics
logreg = LogisticRegression()
logreg.fit(X_train, y_train)
```

Out[15]: 
```
LogisticRegression()
```

In [16]:
```python
from sklearn.metrics import accuracy_score
print('Logistic regression accuracy: {:.3f}'.format(accuracy_score(y_test, logreg.predict(X_test))))
```

```
Logistic regression accuracy: 0.771
```

# Random Forest Classification Model

In [17]:
```python
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
```

Out[17]: 
```
RandomForestClassifier()
```

# Now let's check the accuracy of our Random Forest Classification Model:

In [18]:
```python
print('Random Forest Accuracy: {:.3f}'.format(accuracy_score(y_test, rf.predict(X_test))))
```

```
Random Forest Accuracy: 0.978
```

# Confusion Matrix for our Machine Learning Models

Now I will construct a confusion matrix to visualize predictions made by our classifier and evaluate the accuracy of our machine learning classification

In [19]:
```python
from sklearn.metrics import classification_report
```

```python
print(classification_report(y_test, rf.predict(X_test)))
```

```
              precision    recall  f1-score   support

           0       0.99      0.98      0.99      3462
           1       0.94      0.96      0.95      1038

    accuracy                           0.98      4500
   macro avg       0.97      0.97      0.97      4500
weighted avg       0.98      0.98      0.98      4500
```

In [20]:
```python
y_pred = rf.predict(X_test)
from sklearn.metrics import confusion_matrix
import seaborn as sns
forest_cm = metrics.confusion_matrix(y_pred, y_test, [1,0])
sns.heatmap(forest_cm, annot=True, fmt='.2f',xticklabels = ["Left", "Stayed"] , yticklabels = ["Left", "Stayed"]
plt.ylabel('True class')
plt.xlabel('Predicted class')
plt.title('Random Forest')
```

```
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:70: FutureWarning: Pass labels=[1, 0] as k
eyword args. From version 1.0 (renaming of 0.25) passing these as positional arguments will result in an error
  warnings.warn(f"Pass {args_msg} as keyword args. From version "
```

```
---------------------------------------------------------------------------
NameError                                 Traceback (most recent call last)
~\AppData\Local\Temp/ipykernel_23288/4074542347.py in <module>
      4 forest_cm = metrics.confusion_matrix(y_pred, y_test, [1,0])
      5 sns.heatmap(forest_cm, annot=True, fmt='.2f',xticklabels = ["Left", "Stayed"] , yticklabels = ["Left", "S
tayed"] )
----> 6 plt.ylabel('True class')
      7 plt.xlabel('Predicted class')
      8 plt.title('Random Forest')

NameError: name 'plt' is not defined
```



# Logistic Regression

In [21]:
```python
print(classification_report(y_test, logreg.predict(X_test)))
```

```
              precision    recall  f1-score   support

           0       0.81      0.92      0.86      3462
           1       0.51      0.26      0.35      1038

    accuracy                           0.77      4500
   macro avg       0.66      0.59      0.60      4500
weighted avg       0.74      0.77      0.74      4500
```

In [22]:
```python
logreg_y_pred = logreg.predict(X_test)
logreg_cm = metrics.confusion_matrix(logreg_y_pred, y_test, [1,0])
sns.heatmap(logreg_cm, annot=True, fmt='.2f',xticklabels = ["Left", "Stayed"] , yticklabels = ["Left", "Stayed"]
plt.ylabel('True class')
plt.xlabel('Predicted class')
plt.title('Logistic Regression')
```

```
C:\Users\rosha\anaconda3\lib\site-packages\sklearn\utils\validation.py:70: FutureWarning: Pass labels=[1, 0] as k
eyword args. From version 1.0 (renaming of 0.25) passing these as positional arguments will result in an error
  warnings.warn(f"Pass {args_msg} as keyword args. From version "
---------------------------------------------------------------------------
NameError                                 Traceback (most recent call last)
~\AppData\Local\Temp/ipykernel_23288/2388471023.py in <module>
      2 logreg_cm = metrics.confusion_matrix(logreg_y_pred, y_test, [1,0])
      3 sns.heatmap(logreg_cm, annot=True, fmt='.2f',xticklabels = ["Left", "Stayed"] , yticklabels = ["Left", "S
tayed"] )
----> 4 plt.ylabel('True class')
      5 plt.xlabel('Predicted class')
      6 plt.title('Logistic Regression')

NameError: name 'plt' is not defined
```
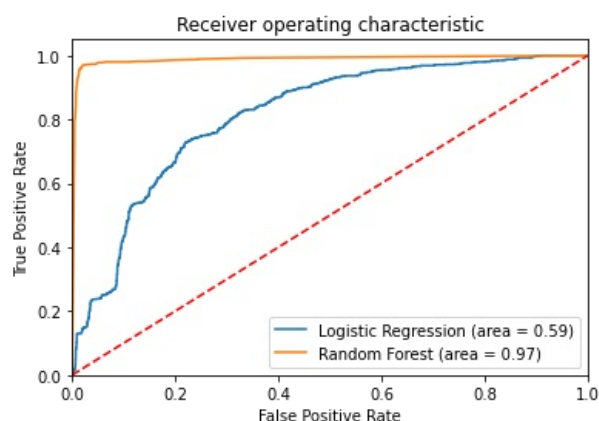


# Employee Turnover Prediction Curve

In [24]:
```python
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve

from matplotlib import pyplot as plt

logit_roc_auc = roc_auc_score(y_test, logreg.predict(X_test))
fpr, tpr, thresholds = roc_curve(y_test, logreg.predict_proba(X_test)[:,1])
rf_roc_auc = roc_auc_score(y_test, rf.predict(X_test))
rf_fpr, rf_tpr, rf_thresholds = roc_curve(y_test, rf.predict_proba(X_test)[:,1])
plt.figure()
plt.plot(fpr, tpr, label='Logistic Regression (area = %0.2f)' % logit_roc_auc)
plt.plot(rf_fpr, rf_tpr, label='Random Forest (area = %0.2f)' % rf_roc_auc)
plt.plot([0, 1], [0, 1],'r--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic')
plt.legend(loc="lower right")
plt.show()
```



In [25]:
```python
feature_labels = np.array(['satisfaction_level', 'last_evaluation', 'time_spend_company', 'Work_accident', 'promo
        'department_RandD', 'department_hr', 'department_management', 'salary_high', 'salary_low'])
importance = rf.feature_importances_
feature_indexes_by_importance = importance.argsort()
for index in feature_indexes_by_importance:
```

```
    print('{}-{:.2f}%'.format(feature_labels[index], (importance[index] *100.0)))
```

department_management-0.23%
promotion_last_5years-0.24%
department_hr-0.30%
department_RandD-0.34%
salary_high-0.69%
salary_low-1.21%
Work_accident-1.45%
last_evaluation-18.97%
time_spend_company-25.67%
satisfaction_level-50.90%

According to our Random Forest classification model, the above aspects show the most important features which will influence whether an employee will leave the company, in ascending order.

In [ ]: