

## HR Analytics Project

### *Machine Learning to Understand & Predict HR Attrition*

Do you want to know what is that one single thing for grand success of any business Enterprise, startup in current competitive era, where razor edge like fine technological upgradation take place continuously along with continuous evolving SOP's, business model? Fundamental answer to this question trace back to core of business which is need & demand of product.

*So next Super Fundamental question comes here : from where product comes ?* Product comes from innovative ideas & solution from talented minds and effort of domain expertise people making idea into reality, a useful product.

The key to success in an organisation is the ability to attract and retain top talents. But where Machine learning comes here? For that we need to go in background of HR analytics.

HR department hire lot of employees every year. The companies invest time and money in training those employees, not just this but there are training programs within the companies for their existing employees as well. HR department often play significant role in designing company compensation programs, creating ambiance in work environment with various team activities, imprinting work culture habits on mindset of peoples and training skill systems that help the organization retain smart minds & talented brain. Most of organisation run different HR analytics vertical to gather data, analyse data to improve process within organisation and to make major decision about human resources.

*The gradual loss of employees' overtime refers as HR attrition. Attrition can happen due to retirement, involuntary (employee is fired or terminated) or voluntary (resigns from an organization).*

A major problem in high employee attrition is its cost to an organization. Job postings, hiring processes, paperwork and new hire training are some of the common expenses of losing employees and replacing them. For Tech companies' loss of talented brain means loss of domain expertise, knowledge base as for these company's Technological capability exists in domain expertise, intellectual property rights of employee. On marketing or sales side customers often prefer to interact with familiar people.

Just like several factors contribute towards building a reliable team and organization, similarly, numerous factors contribute to the attrition rate: Improper work-life balance, better job opportunities, salary hike, Lack of growth or work recognition, unhealthy relations with managers. We can gather data about these factors and utilise Machine learning classification techniques to predict whether employee like to leave organisation or stay in organisation. Here I will show you what leads employee's attrition and Predication of attrition using case study on IBM HR analytics Dataset.

## IBM HR Analytics Employee Attrition & Performance Dataset

In this case study we will use IBM HR Analytics database. This fictional dataset created by IBM employees and available to download from GitHub and Kaggle. You can also download dataset from my [GitHub profile](#). This dataset consists of 1470 rows, 35 features describing each employee's background and characteristics and target variable. Attrition is target variable to be predicted. As target variable is categorial in nature, this case study falls into classification machine learning problem. We have two objectives here:

1. Which key factors result in employee attrition?
2. Building ML Model for predicting attrition.

### Data Preparation : Load, Clean and Format

Let's begin with importing libraries for EDA and dataset itself.

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')

In [2]: df = pd.read_csv("ibm-hr-analytics-employee-attrition-performance.zip")
```

```
In [3]: df
```

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

Checking different datatypes in dataset: -

```
1470 rows × 35 columns

In [4]: df.columns

Out[4]: Index(['Age', 'Attrition', 'BusinessTravel', 'DailyRate', 'Department',
       'DistanceFromHome', 'Education', 'EducationField', 'EmployeeCount',
       'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate',
       'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction',
       'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked',
       'Over18', 'Overtime', 'PercentSalaryHike', 'PerformanceRating',
       'RelationshipSatisfaction', 'StandardHours', 'StockOptionLevel',
       'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance',
       'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion',
       'YearsWithCurrManager'],
      dtype='object')
```

We have 9 features with object datatypes and rest are Numeric feature with int64. Out of all numeric features Education, Environment-Satisfaction, Job-Involvement, Job Satisfaction, Relationship-Satisfaction, Performance Rating, Work Life Balance are ordinal variable. These ordinal features have unique label for each numeric value.

These Ordinal features come with the following label encoding:

- Education: 1 - 'Below College', 2 - 'College', 3 - 'Bachelor', 4 - 'Master', 5 - 'Doctor'
- Environment Satisfaction: 1 - 'Low', 2 - 'Medium', 3 - 'High', 4 - 'Very High'
- Job Involvement: 1 - 'Low', 2 - 'Medium', 3 - 'High', 4 - 'Very High'
- Job Satisfaction: 1 - 'Low', 2 - 'Medium', 3 - 'High', 4 - 'Very High'
- Performance Rating: 1 - 'Low', 2 - 'Average', 3 - 'Good', 4 - 'Excellent', 5 - 'Outstanding'
- Relationship Satisfaction: 1 - 'Low', 2 - 'Medium', 3 - 'High', 4 - 'Very High'
- Work Life Balance: 1 - 'Bad', 2 - 'Good', 3 - 'Better', 4 - 'Best'

*Above nomenclature will help in better understanding of data when we perform EDA in this case study.*

**Data Integrity Check** : Dataset can have missing values, duplicated entries and whitespaces. Now we will perform this integrity check of dataset.

```
Data Integrity Check

In [6]: df.duplicated().sum()

Out[6]: 0

• There is no duplicate value in dataset.
```

Luckily for us, there is no missing data! this will make it easier to work with the dataset.

Dataset doesn't contain Any duplicate entry, whitespace, 'NA', or '-'.

*So, Yes Good to Go Further!!!*

Statistical parameters like mean, median, quantile can give important details about database. Now is time to look at statistical Matrix of Dataset.

*Few key observations from this statistical matrix are listed below: -*

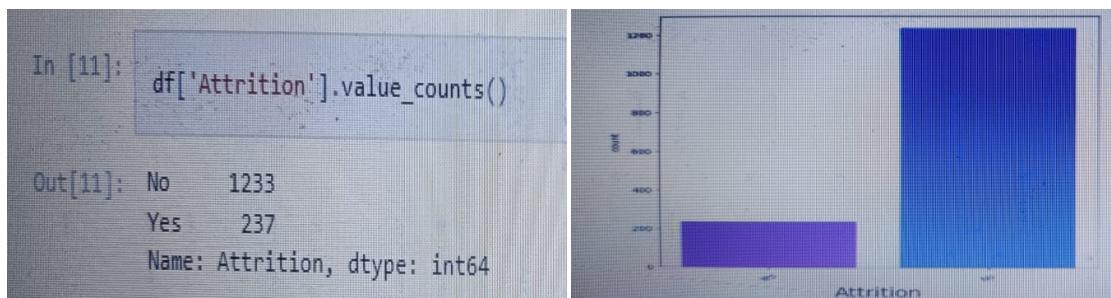
- Minimum Employee Age is 18 and Maximum age of employee 60.
- Average distance from home is 9.1 KM. It means that most of employee travel at least 18 KM in day from home to office.
- Average performance Rating of employees is 3.163 with min value 3.0. This Means that performance of most of employee is 'Good'. This implies that Attrition of Employee with 'Outstanding' or 5 rating need to investigate.
- 50% of Employees has worked at least 2 companies previously.
- For Monthly Income, Monthly Rate by looking at 50% and max column we can say outliers exist in this feature.
- By looking at Mean and Median we see that some of the features are skew in nature.
- For ordinal features statistical terminology like mean, median, std deviation are not applicable.
- Standard Hours and Employee Count contain same value for all statistical parameter. It means they contain one unique value. Exploratory data analysis

Let's begin data exploration of Target variable using count plot.

**Exploratory data analysis :-**

EXPLORATORY DATA ANALYSIS REFERS TO THE CRITICAL PROCESS OF PERFORMING INITIAL INVESTIGATIONS ON DATA SO AS TO DISCOVER PATTERNS, TO SPOT ANOMALIES, TO TEST HYPOTHESIS AND TO CHECK ASSUMPTIONS WITH THE HELP OF SUMMARY STATISTICS AND GRAPHICAL REPRESENTATIONS.

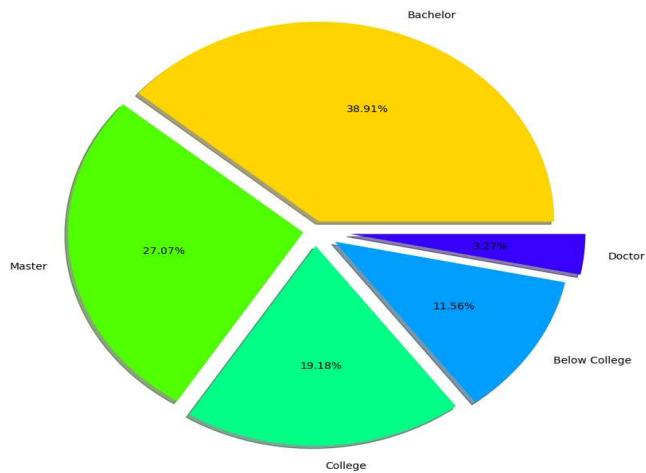
Let's begin data exploration of Target variable using count plot.



**83.87% (1233 employees)** Employees did not leave the organization while **16.13% (237 employees)** did leave the organization making our dataset to be consider as imbalanced since more people stay in the organization than they actually leave.

In this dataset we have features like education, department, education field, job role, job satisfaction which are inter related with each other. Job role & job position not in alignment with educational background can lead attrition. Let investigate this by visualisation of these features one by one to gain more insights.

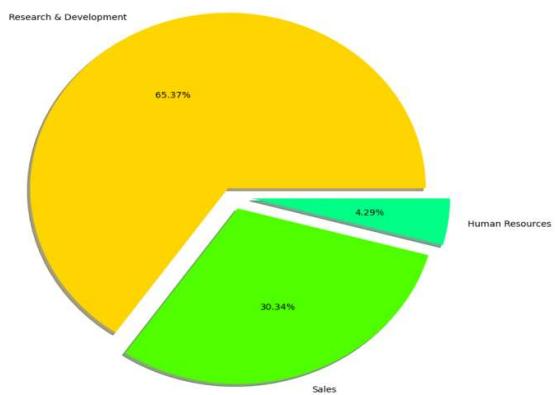
### Education level of Man power available:



#### Key Insights from Pie Plot

1. More than 38 % employees educated at Bachelor level.
2. 30 % of Employees are highly educated which involves master and doctor degree.
3. Almost 19% Employees are educated up to college & 12% are below college.

### Department wise Distribution of Man power :



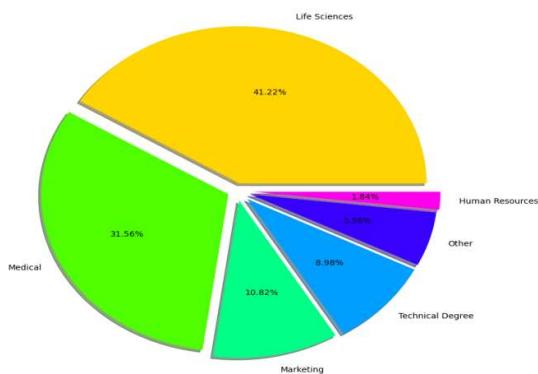
In [28]: pd.crosstab([df.Education],[df.Department], margins=True).style.background_gradient(cmap='summer_r')						
	Department	Human Resources	Research & Development	Sales	All	
Education						
1		5		115	50	170
2		13		182	87	282
3		27		379	166	572
4		15		255	108	398
5		3		30	13	48
All		63		961	446	1470

#### Key Insights on Department and Education Level of employee in each Department

1. 65.37% of Employees work inside Research & Development Department. Out of Total 961 Employee number of employees with education level of Bachelors, Masters, Doctor are 379, 255 and 30 respectively.

2. Only 63 Employee work in HR department.

### Employee distribution as per education field :



#### Key Insights from Pie Plot

- Employees belongs to six different domains.
- 41.22 % Employee comes from Life science background followed by medical profession with 31.56%
- Least number of Employees comes from HR background.

In [25]:	pd.crosstab([df.Department],[df.EducationField], margins=True).style.background_gradient(cmap='summer_r')								
Out[25]:		EducationField	Human Resources	Life Sciences	Marketing	Medical	Other	Technical Degree	All
	Department								
	Human Resources		27	16	0	13	3	4	63
	Research & Development		0	440	0	363	64	94	961
	Sales		0	150	159	88	15	34	446
	All		27	606	159	464	82	132	1470

The probability of Employees Retention is more when there working domain is in alignment with education background. Let check this with crosstab of department against education field.

#### Key Insights from above Cross Tab:

- There are only 27 people with HR background and We know that 63 people work in HR Department from previous result. This implies that at least half employee working in HR department do not have HR background.
- R&D department almost everyone comes from domain expertise or technical background except support staff. These employees usually have high salary, so it will be interesting to investigate attrition in this category.

- There are 159 Employee with Marketing background and all work in Sales Department.
- 50% Employees in sales department have background of Life sciences & Medical. We can clearly see they are working in domain to which their educational background does not belong. So, it will be interesting to see attrition rate in these employees.

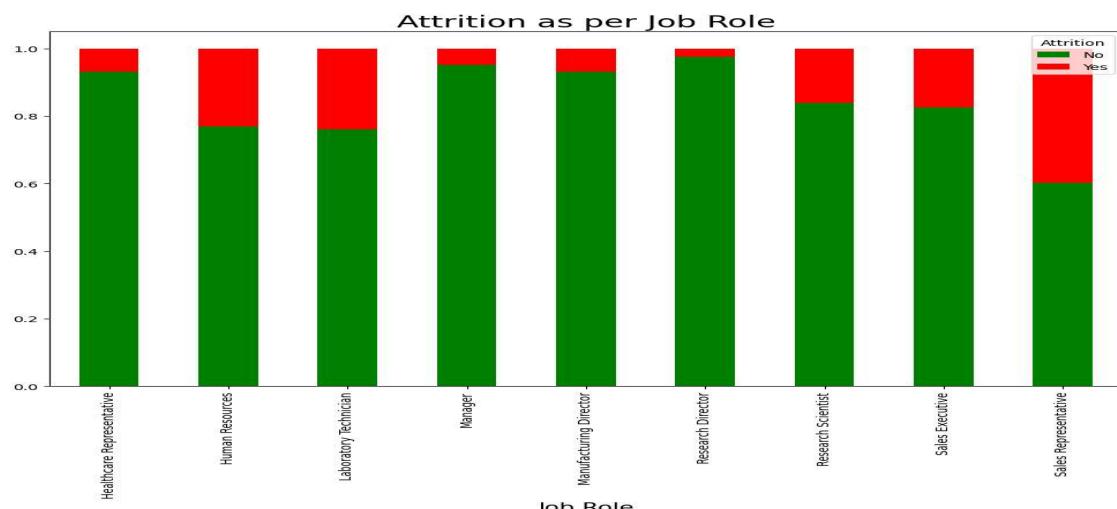
We will Analyse Attrition in department according to education background based on above insight further but before that explore Job role in order to include it in further attrition analyse. First build matrix of department vs job role which will give us idea about number of employees of different job role across department.

	Department	Human Resources	Research & Development	Sales	All
	JobRole				
Healthcare Representative		0	131	0	131
Human Resources		52	0	0	52
Laboratory Technician		0	259	0	259
Manager		11	54	37	102
Manufacturing Director		0	145	0	145
Research Director		0	80	0	80
Research Scientist		0	292	0	292
Sales Executive		0	0	326	326
Sales Representative		0	0	83	83
All		63	961	446	1470

#### Key Insights from above Cross Tab:

- There are 3 job roles in HR Department, maximum of which are sales Executive with 446 Total Employees.
- Human Resources department has 2 Job role i.e., HR & Manager.
- There are 6 different Job role in R&D department with total 961 employees and until now we know that all of them belong to their respective domain background.

#### Attrition by Job role :



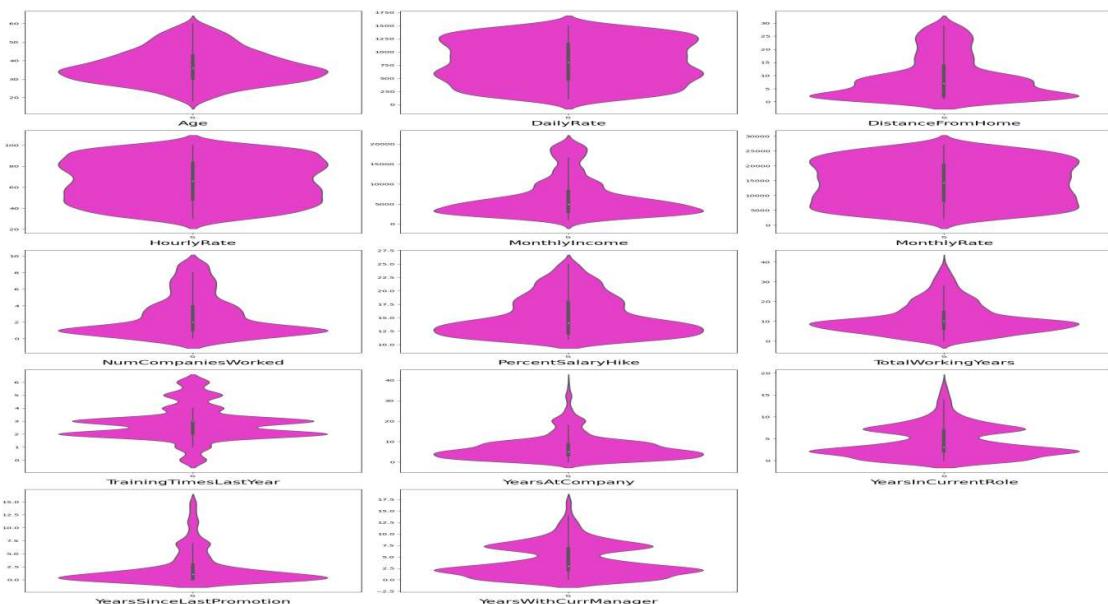
Let's check absolute number matrix of attrition according job role, again this time using crosstab.

		Attrition	No	Yes	All
		JobRole	Department		
Healthcare Representative	Research & Development	122	9	131	
Human Resources	Human Resources	40	12	52	
Laboratory Technician	Research & Development	197	62	259	
	Human Resources	11	0	11	
Manager	Research & Development	51	3	54	
	Sales	35	2	37	
Manufacturing Director	Research & Development	135	10	145	
Research Director	Research & Development	78	2	80	
Research Scientist	Research & Development	245	47	292	
	Sales	269	57	326	
Sales Executive	Sales	50	33	83	
Sales Representative	Sales				
All		1233	237	1470	

#### Key Insights from above Cross Tab:

- Percentage of attrition is high in Sales Representative, Laboratory Technician, Human Resources.
- At the Top chart 62 Laboratory Technician has resign from job, followed by 57 sales executive and 47 Research Scientist.
- 16 % attrition rate for Research Scientist, which involve huge investment from company. Company not only loses employee but its knowledge base, expertise & Intellectual property rights in some cases.

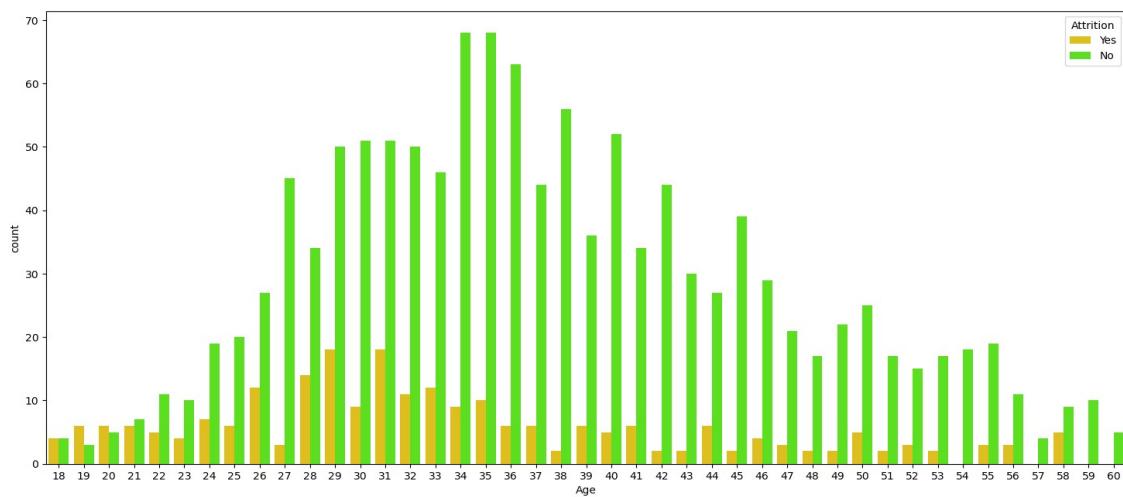
Let's check violin plot of some numerical features to gain more insight.



#### **Key Insights from above Cross Tab:**

- For Majority of people have spent 3 to 10 years at company.
  - Most of people staying company up to 2 years after promotion.
  - Majority of people are train 2-3 times in last year
  - Majority of people stay in same role for maximum 4 yrs.
  - Majority of Employees have salary hike of 10 to 15%.

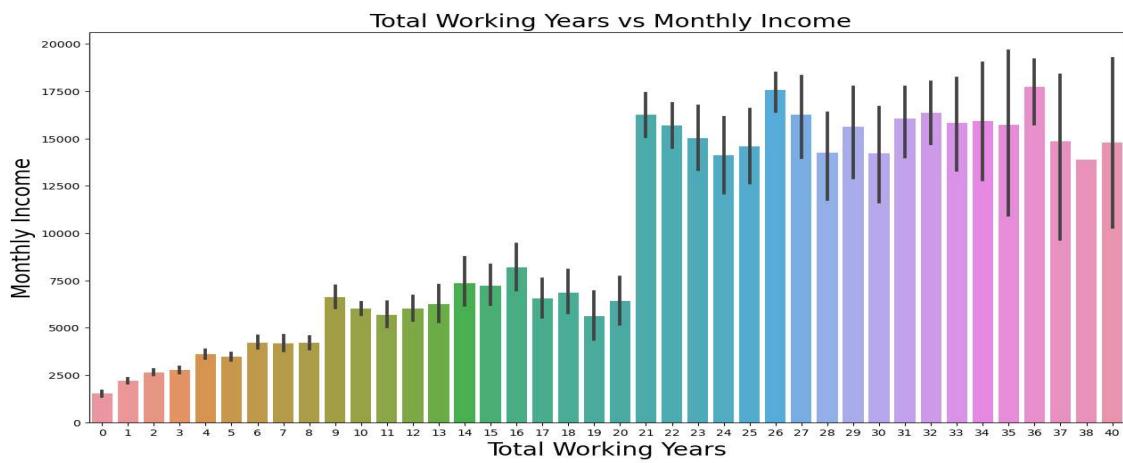
***Q. In which age group attrition rate is high ?***



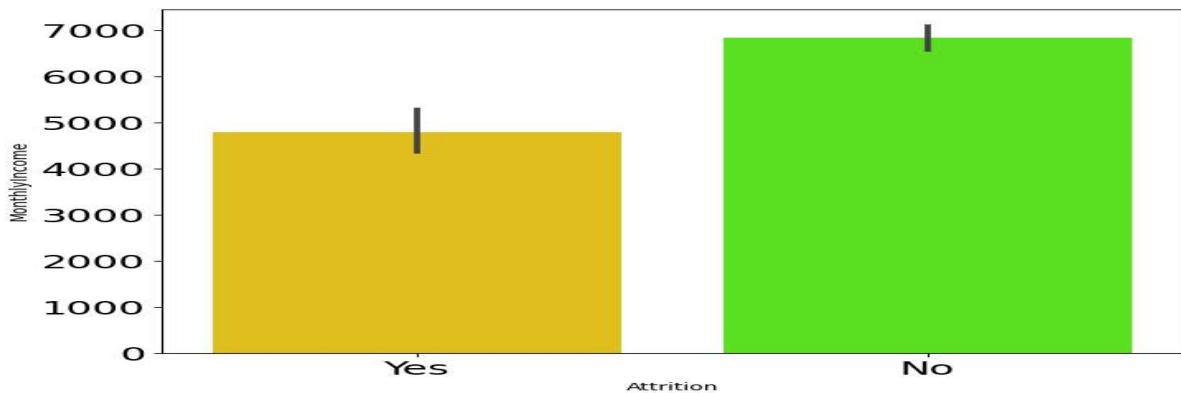
## Key Insights from count plot of Age Vs Attrition:

1. The Attrition rate is minimum between the Age years of 34 and 45.
  2. The Attrition rate is maximum between the Age years of 29 and 33.

**Q. What is variation in monthly income as Total working year increases ?**



**Monthly Income is higher for the employees with 21 or more number of Total working years. For first 9 years monthly income is less than 5000\$. But what about attrition, let's bar chart of Monthly income so we can come across some benchmark of average monthly income in both attrition categories.**



#### Key Insights on Average Monthly Income as per attrition

- We can see that Average monthly income is less in employees who choose to resign compare to rest. Less Monthly Income is major reason behind attrition.
- To prevent attrition average monthly to be greater than 6900\$ is recommended.

## Feature Engineering: Data Pre-processing

*Feature engineering is the process of transforming raw data into features that better represent the underlying problem to the predictive models, resulting in improved model accuracy on unseen data.*

Feature Engineering is very important step in building Machine Learning model. Some machine learning projects succeed and some fail. What makes the difference? Easily the most important factor is the features used. In Feature engineering can be done for various reason. Some of them are mention below:

1. Feature Importance: An estimate of the usefulness of a feature
2. Feature Extraction: The automatic construction of new features from raw data (Dimensionality reduction Technique like PCA)
3. Feature Selection: From many features to a few that are useful
4. Feature Construction: The manual construction of new features from raw data (For example, construction of new column for month out date - mm/dd/yy)

There are Varity of techniques use to achieve above mention means as per need of dataset. Some of Techniques important are as below:

- Handling missing values
- Handling imbalanced data using SMOTE

- Outliers' detection and removal using Z-score, IQR
- Scaling of data using Standard Scalar or Minmax Scalar
- Binning whenever needed
- Encoding categorical data using one hot encoding, label / ordinal encoding ➤ Skewness correction using Boxcox or yeo-Johnson method
- Handling Multicollinearity among feature using variance inflation factor
- Feature selection Techniques:
  - ✓ Correlation Matrix with Heatmap
  - ✓ Univariate Selection – SelectKBest
  - ✓ ExtraTreesClassifier method

In this case study we will use some of the mention feature engineering Techniques one by one.

## 1. Dropping unnecessary features

Feature like 'Over18', 'StandardHours' contain single unique value. Features like EmployeeCount, EmployeeNumber are irrelevant from ML model building perspective. We will drop these features.e will drop these features.

```
In [40]: df.drop(["EmployeeCount", "EmployeeNumber", "Over18", "StandardHours"], axis=1, inplace=True)

In [41]: df.shape

Out[41]: (1470, 31)
```

## 2. Encoding Categorical & Ordinal Features

Label Encoding is employed over target variable 'Attrition' while Ordinal encoding employ for rest categorical features.

# Encoding Categorical Data

```
In [39]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df["Attrition"] = le.fit_transform(df["Attrition"])
df.head()
```

```
In [42]: from sklearn.preprocessing import OrdinalEncoder
oe = OrdinalEncoder()
def ordinal_encode(df, column):
    df[column] = oe.fit_transform(df[column])
    return df

oe_col = ['BusinessTravel', 'Department', 'EducationField', 'Gender', 'JobRole', 'MaritalStatus', 'OverTime']
df=ordinal_encode(df, oe_col)
df.head()
```

Since now encoding is done we will move towards outliers' detection and removal.

### 3. Outliers' detection and removal

*Identifying outliers and bad data in your dataset is probably one of the most difficult parts of data clean-up, and it takes time to get right. Even if you have a deep understanding of statistics and how outliers might affect your data, it's always a topic to explore cautiously.*

Machine learning algorithms are sensitive to the range and distribution of attribute values. Data outliers can spoil and mislead the training process resulting in longer training times, less accurate models and ultimately poorer results. Outliers can be seen in boxplot of numerical feature. We did not added boxplot here as it will make this article length, I left it to reader to further investigate. Now we will use Zscore method for outliers' detection.

```
In [44]: from scipy.stats import zscore
z = np.abs(zscore(df))
threshold = 3
df1 = df[(z<3).all(axis = 1)]

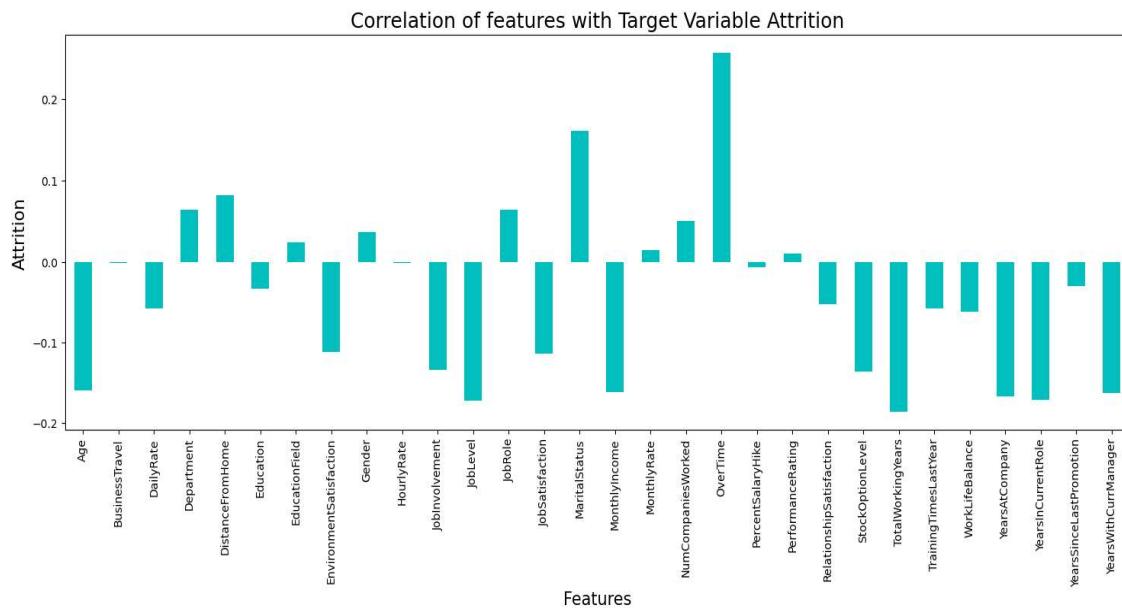
print ("Shape of the dataframe before removing outliers: ", df.shape)
print ("Shape of the dataframe after removing outliers: ", df1.shape)
print ("Percentage of data loss post outlier removal: ", (df.shape[0]-df1.shape[0])/df.shape[0]*100)

df=df1.copy()
```

Shape of the dataframe before removing outliers: (1470, 31)  
 Shape of the dataframe after removing outliers: (1387, 31)  
 Percentage of data loss post outlier removal: 5.646258503401361

### 4. Correlation Heatmap

Correlation Heatmap show in a glance which variables are correlated, to what degree, in which direction, and alerts us to potential multicollinearity problems. The bar plot of correlation coefficient of target variable with independent features shown below



## 5. Multicollinearity between features

Variance Inflation factor imported from statsmodels.stats.outliers\_influence to check multicollinearity between features.

**Check Multicollinearity between features using Variance\_inflation\_factor**

```
In [52]: from statsmodels.stats.outliers_influence import variance_inflation_factor
vif= pd.DataFrame()
vif['VIF']= [variance_inflation_factor(df.values,i) for i in range(df.shape[1])]
vif['Features']= df.columns
vif
```

Out[52]:	VIF	Features		
	0 1.930457	Age	15 5.842828	MonthlyIncome
	1 1.014314	BusinessTravel	16 1.022108	MonthlyRate
	2 1.025841	DailyRate	17 1.426763	NumCompaniesWorked
	3 2.172093	Department	18 1.028400	Overtime
	4 1.017385	DistanceFromHome	19 1.016867	PercentSalaryHike
	5 1.065266	Education	20 1.747581	PerformanceRating
	6 1.030480	EducationField	21 1.022260	RelationshipSatisfaction
	7 1.024396	EnvironmentSatisfaction	22 2.279101	StockOptionLevel
	8 1.024366	Gender	23 4.093506	TotalWorkingYears
	9 1.024189	HourlyRate	24 1.025519	TrainingTimesLastYear
	10 1.020167	JobInvolvement	25 1.017093	WorkLifeBalance
	11 5.976707	JobLevel	26 6.296064	YearsAtCompany
	12 2.023213	JobRole	27 3.513852	YearsInCurrentRole
	13 1.023909	JobSatisfaction	28 1.373189	YearsSinceLastPromotion
	14 2.298943	MaritalStatus	29 3.433437	YearsWithCurrManager

We can see that for all features Variance inflation factor is within permissible limit of 10. Multicollinearity do not pose any threat here.

## 6. Handling imbalanced data using SMOTE

This two-class dataset is imbalanced (84% vs 16%). As a result, there is a possibility that the model built might be biased towards the majority and overrepresented class. We can resolve this by Synthetic Minority Oversampling Technique (SMOTE) to over-sample the minority class.

```
In [53]: from imblearn.over_sampling import SMOTE  
  
In [54]: oversample = SMOTE()  
X, Y = oversample.fit_resample(X, Y)
```

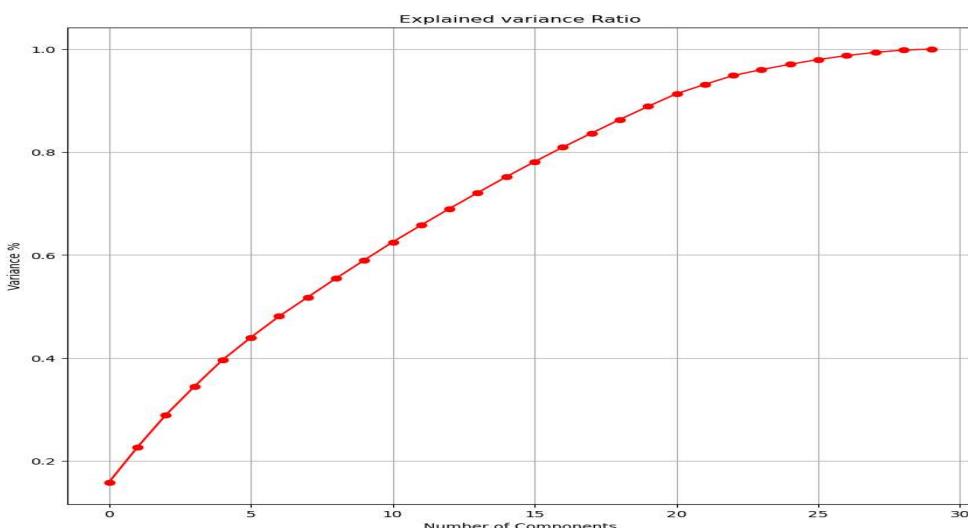
## 7. Scaling of data using Standard Scalar

### Standard Scaling

```
In [56]: from sklearn.preprocessing import StandardScaler  
scaler= StandardScaler()  
X_scale = scaler.fit_transform(X)
```

## 8. Dimensionality Reduction Using PCA

PCA used find patterns and extract the latent features from our dataset.



We can see that 21 principal components attribute for 90% of variation in the data. PCA applied for 21 components.

```
In [58]: pca_new = PCA(n_components=21)
x_new = pca_new.fit_transform(X_scale)

In [59]: principle_x=pd.DataFrame(x_new,columns=np.arange(21))
```

## Machine Learning Model Building:

In this section we will build Supervised learning ML model-based classification algorithm. As objective is to predict attrition in 'Yes' or 'No' leads to fall problem in domain of classification algorithm. train\_test\_split used to split data with size of 0.33

```
In [60]: from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, confusion_matrix,classification_report,f1_score

In [61]: from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import BaggingClassifier

In [62]: X_train, X_test, Y_train, Y_test = train_test_split(principle_x, Y, random_state=42, test_size=.33)
print('Training feature matrix size:',X_train.shape)
print('Training target vector size:',Y_train.shape)
print('Test feature matrix size:',X_test.shape)
print('Test target vector size:',Y_test.shape)

Training feature matrix size: (1551, 21)
Training target vector size: (1551,)
Test feature matrix size: (765, 21)
Test target vector size: (765,)
```

Activate Wind  
Go to Settings to a

First we will build base model using logistic regression algorithm. Best random state is investigated using for loop for random state in range of (0,250).

```
In [63]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix,classification_report,f1_score
maxAccu=0
maxRS=0
for i in range(1,250):
    X_train,X_test,Y_train,Y_test = train_test_split(principle_x,Y,test_size = 0.33, random_state=i)
    log_reg=LogisticRegression()
    log_reg.fit(X_train,Y_train)
    y_pred=log_reg.predict(X_test)
    acc=accuracy_score(Y_test,y_pred)
    if acc>maxAccu:
        maxAccu=acc
        maxRS=i
print('Best accuracy is', maxAccu , 'on Random_state', maxRS)

Best accuracy is 0.8666666666666667 on Random_state 53
```

Logistics regression model is train with random state 53. The evalution matrix along with classification report is as below :

```
Logistics Regression Evaluation

Accuracy Score of Logistics Regression : 0.8666666666666667

Confusion matrix of Logistics Regression :
[[327 43]
 [ 59 336]]

classification Report of Logistics Regression
precision    recall   f1-score   support
          0       0.85      0.88      0.87      370
          1       0.89      0.85      0.87      395

accuracy                           0.87      765
macro avg       0.87      0.87      0.87      765
weighted avg    0.87      0.87      0.87      765
```

As Now base model is ready with f1-score of 0.87, we will train model with different classification algorithm along with k-5 fold cross validation. The final evaluation matrix different classification algorithm is as shown table below:

ML Algorithm	Accuracy Score	CV Mean Score	f-1 Score
Logistics Regression	0.8666	0.6856	0.87
SVC	0.9111	0.6088	0.91
Gaussian NB	0.8496	0.7331	0.85
DecisionTree Classifier	0.7908	0.8342	0.79
KNeighbors Classifier	0.8444	0.7236	0.84
Random Forest Classifier	0.8993	0.9111	0.90
AdaBoost Classifier	0.8575	0.8740	0.86
Gradient Boosting Classifier	0.8836	0.8873	0.88
Bagging Classifier	0.8601	0.8887	0.86

(Min Value in column -Green, Max Value in column - Yellow Colour )

We can see that Random Forest Classifier gives us maximum f1-score & mean cross validation score. We will perform hyper parameter tuning on random forest classifier to build final ML Model.

```
GridSearchCV

In [69]: from sklearn.model_selection import GridSearchCV

In [70]: parameter = { 'bootstrap': [True], 'max_depth': [5, 10, 20, 40, 50, None],
   'max_features': ['auto', 'log2'],
   'criterion': ['gini', 'entropy'],
   'n_estimators': [5, 10, 15, 25, 50, 100]}

In [71]: GCV = GridSearchCV(RandomForestClassifier(), parameter, cv=5, n_jobs = -1, verbose=3)
GCV.fit(X_train, Y_train)

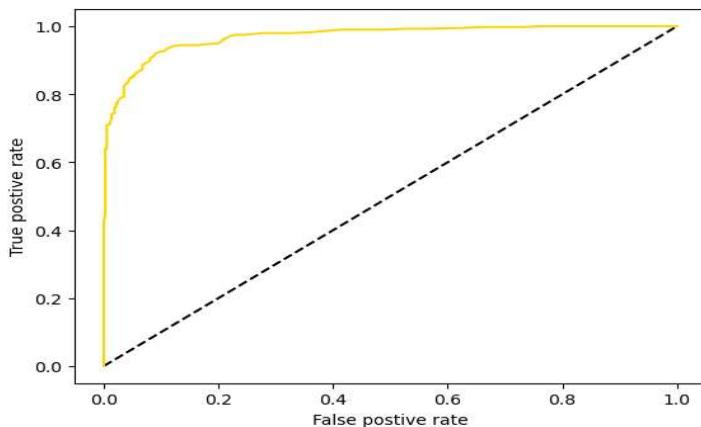
Fitting 5 folds for each of 144 candidates, totalling 720 fits
Out[71]: GridSearchCV(cv=5, estimator=RandomForestClassifier(), n_jobs=-1,
param_grid=[{'bootstrap': [True], 'criterion': ['gini', 'entropy'],
'max_depth': [5, 10, 20, 40, 50, None],
'max_features': ['auto', 'log2'],
'n_estimators': [5, 10, 15, 25, 50, 100]}],
verbose=3)
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.
```

Next step is to build final machine learning model over best params in Hyper parameter tuning.

```
Final Model

In [73]: Final_mod = RandomForestClassifier(bootstrap=True, criterion='entropy', n_estimators= 100, max_depth=20, max_features='auto')
Final_mod.fit(X_train, Y_train)
y_pred=Final_mod.predict(X_test)
print('Accuracy Score : '+str(accuracy_score(Y_test, y_pred)))
4
Accuracy Score :
0.9084967320261438
```

We can see that Final model with hyper parameter tuning leads to slight decrease in accuracy score from 0.8980 in original model to 0.9089. This complete possible We will use model with default values as our final model. AOC-ROC score of final random forest classifier model is shown below:



At last, we will save final model with joblib library, so it can be deploy on cloud platform.

## Save the Model

```
In [75]: import joblib
joblib.dump(Final_mod,'IBM_HR_Analytics_Employee_Attrition_Performance.pkl')

Out[75]: ['IBM_HR_Analytics_Employee_Attrition_Performance.pkl']
```

### Concluding Remarks on EDA and ML Model

- Bench mark of 6900\$ monthly income is recommended to Prevent attrition.
- Attrition rate is high in age group of 29 to 33. HR need to keep eye over need & expectation of this age group from company.
- Percentage of attrition is high in Sales Representative, Laboratory Technician
- 16 % attrition rate among Research Scientist and no company afford to lose them.
- Almost 50% employs in sales department from different education background. There is possibility of dissatisfaction among them as attrition high among them.
- Different feature engineering techniques like balancing data, outliers' removal, label encoding, feature selection & PCA are performed on data.
- Random Forest Classifier model gives maximum Accuracy.

You can get code of this case study from my [GitHub Profile](#).