



Project Title	IBM HR Analytics Employee Attrition & Performance
Tools	Python, ML, SQL, Excel
Technologies	Data Analyst & Data scientist
Project Difficulties level	intermediate

Dataset : Dataset is available in the given link. You can download it at your convenience.

[Click here to download data set](#)

About Dataset

Uncover the factors that lead to employee attrition and explore important questions such as 'show me a breakdown of distance from home by job role and attrition' or 'compare average monthly income by education and attrition'. This is a fictional data set created by IBM data scientists.

Education

1 'Below College'

2 'College'

3 'Bachelor'

4 'Master'

5 'Doctor'

EnvironmentSatisfaction

1 'Low'

2 'Medium'

3 'High'

4 'Very High'

JobInvolvement

1 'Low'

2 'Medium'

3 'High'

4 'Very High'

JobSatisfaction

1 'Low'

2 'Medium'

3 'High'

4 'Very High'

PerformanceRating

1 'Low'

2 'Good'

3 'Excellent'

4 'Outstanding'

RelationshipSatisfaction

1 'Low'

2 'Medium'

3 'High'

4 'Very High'

WorkLifeBalance

1 'Bad'

2 'Good'

3 'Better'

4 'Best'

Example

what steps you should have to follow

Sample code

HR Attrition Analysis¶

In the business world, companies often face the challenge of retaining talented employees. One of the most pressing issues is the increasing rate of employee turnover, commonly known as HR attrition. Turnover can have a significant impact on a company's productivity, stability, and long-term sustainability. High attrition rates can lead to increased recruitment and training costs, disrupt team dynamics, and result in the loss of valuable institutional knowledge. Therefore, understanding the factors contributing to attrition and implementing effective retention strategies is crucial for maintaining a competitive edge and ensuring

Objectives of the Analysis

1. Understand Current Turnover Rates: Gain a comprehensive understanding of the current employee turnover rate and analyze the demographic distribution of attrition by age, gender, education, department,

and job role.

2. **Identify Key Factors Influencing Turnover:** Examine the main factors contributing to employee attrition, including job satisfaction indicators (job involvement and work-life balance), salary factors (monthly income and salary hikes), and benefit factors (stock option levels), to uncover patterns and correlations that drive higher attrition rates.

Data Cleaning

In [1]:

```
# import data manipulation package
import pandas as pd
import numpy as np

# import data visualization package
import matplotlib.pyplot as plt
import seaborn as sns

# importing the warnings library
import warnings
warnings.filterwarnings('ignore')
```

In [2]:

```
# set pandas options
pd.set_option('display.max_columns', 35)

# load dataset
df =
pd.read_csv('/kaggle/input/ibm-hr-analytics-attrition-dataset/WA_Fn-UseC_-HR-Employee-Attrition.csv')
df.head()
```

Out[2]:

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromOffice	Education	Employed	EmployeeCount	EnvironmentSatisfaction	Gender	HourlyRate	JobInvolvement	JobLevel	JobRole	JobSatisfaction	MaritalStatus	MonthlyIncome	MonthlyRate	NumCompaniesWorked	Over18	OverTime	Performance	RelationshipSatisfaction	StandardHours	StockOptions	TotalWorkingHours	TrainingTimesLastYear	WorkLifeBalance	YearsAtCompany
--	-----	-----------	----------------	-----------	------------	--------------------	-----------	----------	---------------	-------------------------	--------	------------	----------------	----------	---------	-----------------	---------------	---------------	-------------	--------------------	--------	----------	-------------	--------------------------	---------------	--------------	-------------------	-----------------------	-----------------	----------------

				avel		t	mHome		ield	unt	mber	sfaction		e	ment			cti	on	tus	ome	e	Worked				ary	Hike	Rating	action	urs	Level	ng	Years	ast	Year	lance	p	any
	0		41	Yes	Travel - Rarely	1102	Sales	1	2				Female	94	3	2		Sales Executive	4		Sing	le	5993	19479	8		Y	Yes	11	3	1		80	0	8	0		1	6
	1		49	No	Travel - Frequently	279	Research & Development	8	1				Male	61	2	2		Research Scientist	2		Married	5130	24907	1		Y	No	23	4	4		80	1	10	3		3	10	
	2		37	Ye	Tra	137	Res	2	2	Oth	1	4	4	Mal	92	2	1	Lab	3	Sin	209	239	6		Y	Ye	15	3	2		80	0	7	3		3	0		

		s	v e l _ R a r e l y	3	e a r c h & D e v e l o p m e n t				r				e				o r a t o r y T e c h n i c i a n		g l e	0	6				s												
	3	3 3	N o	T r a v e l _ F r e q u e n t l y	1 3 9 2	R e s e a r c h & D e v e l o p m e n t	3	4	L i f e S c i e n c e s	1	5	4	F e m a l e	5 6	3	1	R e s e a r c h S c i e n t i s t	3	M a r r i e d	2 9 0 9	2 3 1 5 9	1	Y	Y e s	1 1	3	3	8 0	0	8	3	3	8				
	4	2 7	N o	T r a v e l _ R a r e	5 9 1	R e s e a r c h & D e v e l o p m e n t	2	1	M e d i c a l	1	7	1	M a l e	4 0	3	1	L a b o r a t o r y T	2	M a r r i e d	3 4 6 8	1 6 6 3 2	9	Y	N o	1 2	3	4	8 0	1	6	3	3	2				

DistanceFromHome	0.0
Education	0.0
EducationField	0.0
EmployeeCount	0.0
EmployeeNumber	0.0
EnvironmentSatisfaction	0.0
Gender	0.0
HourlyRate	0.0
JobInvolvement	0.0
JobLevel	0.0
JobRole	0.0
JobSatisfaction	0.0
MaritalStatus	0.0
MonthlyIncome	0.0
MonthlyRate	0.0
NumCompaniesWorked	0.0
Over18	0.0
Overtime	0.0
PercentSalaryHike	0.0
PerformanceRating	0.0
RelationshipSatisfaction	0.0
StandardHours	0.0
StockOptionLevel	0.0
TotalWorkingYears	0.0
TrainingTimesLastYear	0.0
WorkLifeBalance	0.0
YearsAtCompany	0.0
YearsInCurrentRole	0.0
YearsSinceLastPromotion	0.0
YearsWithCurrManager	0.0

dtype: float64

There are any missing values in the dataset.

In [6]:

```
# check data types
df.dtypes
```

Out[6]:

Age	int64
Attrition	object
BusinessTravel	object
DailyRate	int64
Department	object


```
DistanceFromHome      int64
Education              int64
EducationField         object
EmployeeCount          int64
EmployeeNumber         int64
EnvironmentSatisfaction int64
Gender                object
HourlyRate             int64
JobInvolvement         int64
JobLevel              int64
JobRole               object
JobSatisfaction        int64
MaritalStatus         object
MonthlyIncome         int64
MonthlyRate           int64
NumCompaniesWorked    int64
Over18               object
OverTime              object
PercentSalaryHike     int64
PerformanceRating     int64
RelationshipSatisfaction int64
StandardHours         int64
StockOptionLevel      int64
TotalWorkingYears     int64
TrainingTimesLastYear int64
WorkLifeBalance       int64
YearsAtCompany        int64
YearsInCurrentRole    int64
YearsSinceLastPromotion int64
YearsWithCurrManager  int64
```

```
dtype: object
```

All columns have appropriate data types, ensuring that the data is correctly formatted for analysis.

In [7]:

```
# check data decribe
df.describe()
```

Out[7]:

	A	D	Di	E	E	E	En	H	J	J	J	M	M	Nu	P	P	Re	S	S	T	Tr	W	Y	Y	Ye	Ye
	g	a	st	d	m	m	vir	o	o	o	o	o	o	m	er	er	lati	t	t	ot	ain	o	e	e	ars	ars
		i	a	u	p	pl	on	u	b	b	b	n	n	Co	c	fo	on	a	o	al	ing	r	a	ar	Sin	Wit

[illegible]

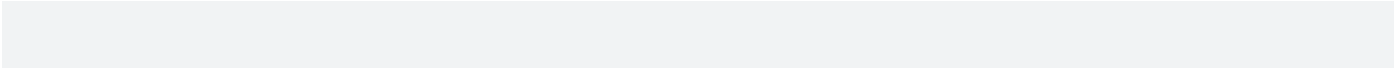
	00000	0000000	0000	000		500000	0	00000	0000	0000	0000	0000	0000000	.5000000	0	0000	00	0	0000	0000	0	0000	0000	0	0	0
max	60000000	14990000000	290000000	50000000	10	20680000000	400000	1000000000	40000000	50000000	40000000	19990000000	26999000000	900000	250000000	40000000	40000000	80000	30000000	40000000	600000	40000000	40000000	18000000	15000000	17000000

Based on this summary, there are no apparent outliers in the dataset, as the values fall within expected ranges.

Exploratory Data Analysis

In [8]:

```
df.head()
```



Out[8]:

		Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EmployeeCount	EmployeeNumber	EnvironmentSatisfaction	Gender	HourlyRate	JobInvolvement	JobLevel	JobRole	JobSatisfaction	MaritalStatus	MonthlyIncome	MonthlyRate	NumCompaniesWorked	Overtime	Over18	PerformanceRating	RelationshipSatisfaction	StandardHours	StockOptionsLevel	TotalWorkingYears	TrainingTimesLastYear	WorkLifeBalance	YearsAtCompany
--	--	-----	-----------	----------------	-----------	------------	------------------	-----------	---------------	----------------	-------------------------	--------	------------	----------------	----------	---------	-----------------	---------------	---------------	-------------	--------------------	----------	--------	-------------------	--------------------------	---------------	-------------------	-------------------	-----------------------	-----------------	----------------

			l		e	d	t	r	n		t		n	s	e	d		i	g	n	s	l	a	ar	e	y						
0	41	Yes	Travel - Rarely	1102	Sales	1	2	Life Sciences	1	1	2	Female	94	3	2	Sales Executive	4	Singl	5993	19479	8	Y	Yes	11	3	1	80	0	8	0	1	6
1	49	No	Travel - Frequently	279	Research & Development	8	1	Life Sciences	1	2	3	Male	61	2	2	Research Scientist	2	Married	5130	24907	1	Y	No	23	4	4	80	1	10	3	3	10
2	37	Yes	Travel -	1373	Research	2	2	Other	1	4	4	Male	92	2	1	Laborat	3	Singl	2090	2396	6	Y	Yes	15	3	2	80	0	7	3	3	0

3	33	No	Travel - Frequently	1392	Research & Development	3	4	Lifesciences	1	5	4	Female	56	3	1	Research Scientist	3	Married	2909	23159	1	Y	Yes	11	3	3	80	0	8	3	3	8
4	27	No	Travel - Rarely	591	Research & Development	2	1	Medical	1	7	1	Male	40	3	1	Laboratory Technician	2	Married	3468	16632	9	Y	No	12	3	4	80	1	6	3	3	2

[illegible]

Attrition Rate

Attrition rate: The attrition rate measures the percentage of employees who leave the company in a given period of time. It is usually calculated within a year and is expressed as a percentage of the total number of employees.

In [9]:

```
df['Attrition'].value_counts(normalize=True)
```

Out[9]:

Attrition	
No	0.838776
Yes	0.161224

Name: proportion, dtype: float64

The output displays the proportion of employees with regard to attrition status in the dataset. Let's visualize it!

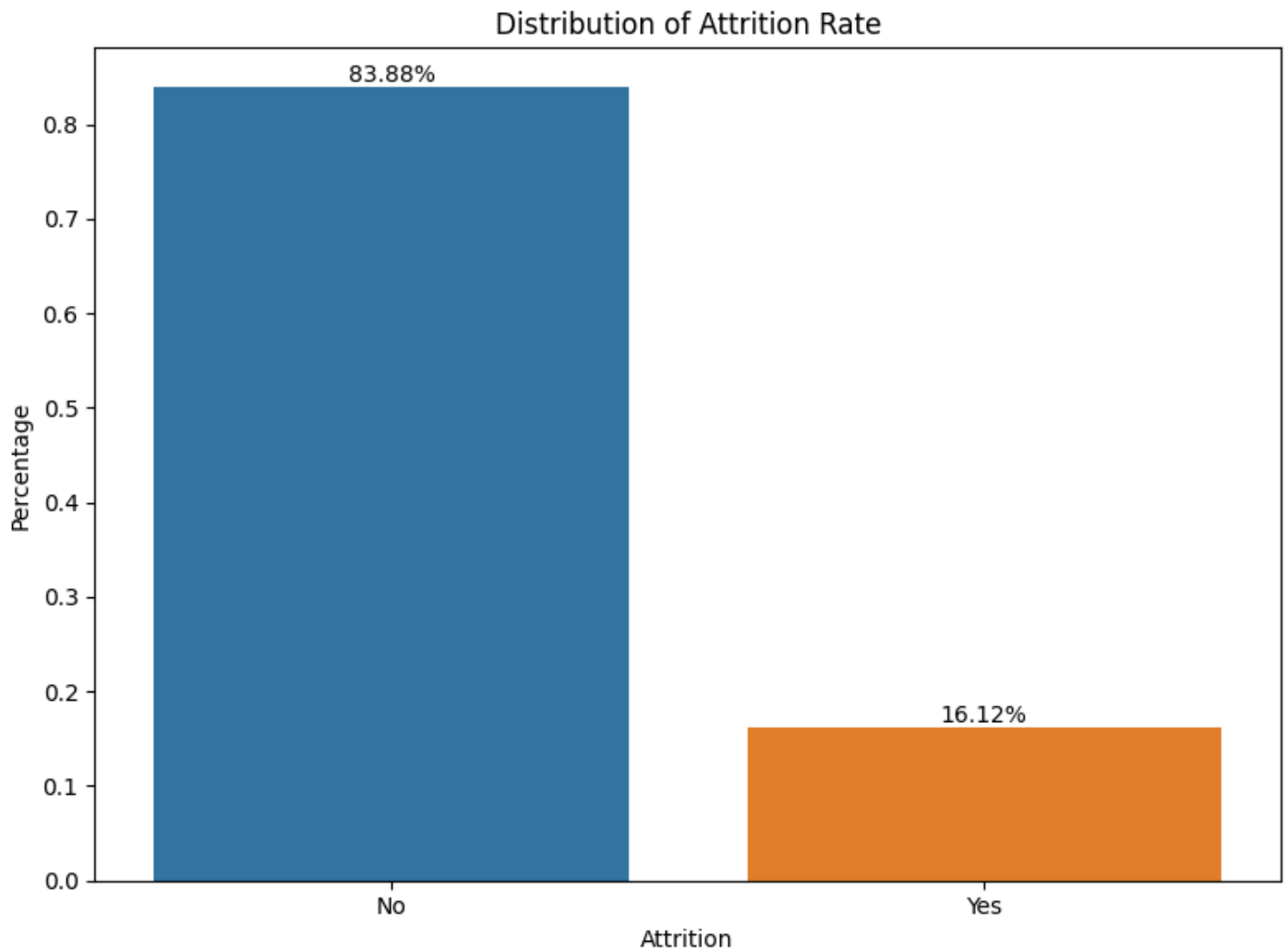
```
In [10]:
```

```
attrition = df['Attrition'].value_counts(normalize=True)

plt.figure(figsize=(8,6))
ax = sns.barplot(x=attrition.index, y=attrition)

for p in ax.patches:
    ax.annotate(f'{p.get_height() * 100:.2f}%',
                (p.get_x() + p.get_width() / 2., p.get_height()),
                ha='center', va='bottom')

plt.title('Distribution of Attrition Rate')
plt.xlabel('Attrition')
plt.ylabel('Percentage')
plt.tight_layout()
plt.show()
```



Based on the analysis, the company's attrition rate is 16.12%. This means that about 16.12% of the employees decided to leave the company during the analyzed period.

Average of Tenure

Average tenure: The average tenure measures the average number of years an employee stays with the company before leaving. It can provide insight into workforce stability and employee satisfaction within the organization.

In [11]:

```
avg_tenure = df['YearsAtCompany'].mean()  
print(f'Average years of employee to leave the company is {avg_tenure} years')
```

```
Average years of employee to leave the company is 7.0081632653061225 years
```


The average tenure of employees before they decided to leave was 7.01 years. With this average tenure, it can be concluded that many employees feel comfortable and have been with the company for a long time.

Employee's Demographics

In [12]:

```
fig, axes = plt.subplots(nrows=1, ncols=3, figsize=(15,5))

sns.histplot(data=df, x='Age', kde=True, ax=axes[0])
axes[0].set_title('Distribution Employee by Age')
axes[0].set_xlabel('Age')
axes[0].set_ylabel('Count')

sns.countplot(data=df, x='Gender', ax=axes[1])
axes[1].set_title('Distribution Employee by Gender')
axes[1].set_xlabel('Gender')
axes[1].set_ylabel('Count')

sns.countplot(data=df, x='Department', ax=axes[2])
axes[2].set_title('Distribution Employee by Department')
axes[2].set_xlabel('Department')
axes[2].set_ylabel('Count')

plt.tight_layout()
plt.show()
```



1. Age: Most of the company's employees are in the 30-35 age group. This indicates that the company has

many employees who are at a productive and experienced age.

- 2. Gender: The majority of employees at this company are male. There are significantly more male employees than female employees.
- 3. Department: Most of the company's employees are concentrated in the research and development department. This indicates that the company is heavily focused on product or service research and development activities.

In [13]:

```
df_attrition = df[df['Attrition'] == 'Yes']
df_attrition.head()
```

Out[13]:

Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Field	EmployeeCount	EnvironmentSatisfaction	Gender	HourlyRate	JobInvolvement	JobLevel	JobRole	JobSatisfaction	MonthlyIncome	MonthlyRate	NumCompaniesWorked	Overtime	PercentSalaryHike	PerformanceRating	RelationshipSatisfaction	StandardHours	StockOptionsLevel	TotalWorkingYears	TrainingTimesLastYear	WorkLifeBalance	YearsAtCompany			
0	41	Yes	1102	Sales	1	2	LifeSciences	1	1	2	Female	94	3	2	SalesExecutive	4	Singl	5993	19479	8	Y	Yes	11	3	1	80	0	8	0	1	6

2	3	Yes	Travel - Rarely	1373	Research & Development	2	2	Other	1	4	4	Male	92	2	1	Laboratory Technician	3	Single	2090	2396	6	Y	Yes	15	3	2		80	07	3	3	0	
14	28	Yes	Travel - Rarely	103	Research & Development	24	3	Life Sciences	1	19	3	Male	50	2	1	Laboratory Technician	3	Single	2028	12947	5	Y	Yes	14	3	2		80	06	4	3	4	
21	36	Yes	Travel	1218	Sales	9	4	Life Sc	1	27	3	Male	82	2	1	Sales R	1	Single	3407	6986	7	Y	No	23	4	2		80	01	4	3	5	

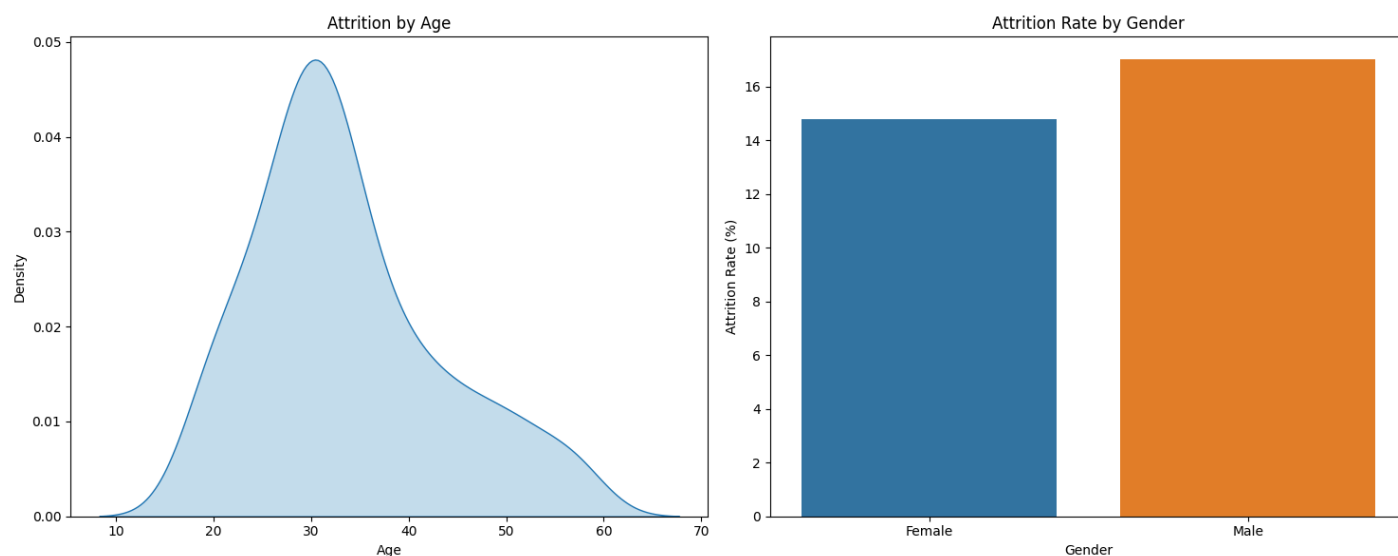
In [15]:

```
fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(15,6))

# Plot 1: KDE plot of Age with Attrition hue
sns.kdeplot(data=df_attrition, x='Age', fill=True, ax=axes[0])
axes[0].set_title('Attrition by Age')
axes[0].set_xlabel('Age')
axes[0].set_ylabel('Density')

# Plot 2: Bar plot of Gender count with Attrition hue
attrition_rate_df = calculate_attrition_rate(df, 'Gender')
sns.barplot(data=attrition_rate_df, x='Gender', y='AttritionRate', ax=axes[1])
axes[1].set_title('Attrition Rate by Gender')
axes[1].set_xlabel('Gender')
axes[1].set_ylabel('Attrition Rate (%)')

plt.tight_layout()
plt.show()
```



1 Reference link