

```

train =
pd.read_csv('https://raw.githubusercontent.com/sameerCoder/DATA_ANALYST_DATASETS/main/HrAnalytics/HrAnalytics_train.csv')
test =
pd.read_csv('https://raw.githubusercontent.com/sameerCoder/DATA_ANALYST_DATASETS/main/HrAnalytics/HrAnalytics_test.csv')

1.

# getting their shapes
print("Shape of train :", train.shape)
print("Shape of test :", test.shape)

2.

train.head()

3.

test.head()

4.

# describing the training set
train.describe(include = 'all')

5.

train.info()

6.

# checking if there is any NULL value in the dataset
train.isnull().any()

7.

test.isnull().sum()

8.

# looking at the most popular departments

from wordcloud import WordCloud
from wordcloud import STOPWORDS

stopword = set(STOPWORDS)

wordcloud = WordCloud(background_color = 'white', stopwords =
stopword).generate(str(train['department']))

plt.rcParams['figure.figsize'] = (12, 8)
print(wordcloud)

```

```
plt.imshow(wordcloud)
plt.title('Most Popular Departments', fontsize = 30)
plt.axis('off')
plt.show()
```

9.

```
# checkig the no. of Employees Promoted
```

```
train['is_promoted'].value_counts()
```

10.

```
# finding the %age of people promoted
```

```
promoted = (4668/54808)*100
```

```
print("Percentage of Promoted Employees is {:.2f}%".format(promoted))
```

11.

```
#plotting a scatter plot
```

```
plt.hist(train['is_promoted'])
```

```
plt.title('plot to show the gap in Promoted and Non-Promoted Employees',
          fontsize = 30)
```

```
plt.xlabel('0 -No Promotion and 1- Promotion', fontsize = 20)
```

```
plt.ylabel('count')
```

```
plt.show()
```

12.

```
# checking the distribution of the avg_training score of the Employees
```

```
plt.rcParams['figure.figsize'] = (15, 7)
```

```
sns.distplot(train['avg_training_score'], color = 'blue')
```

```
plt.title('Distribution of Training Score among the Employees', fontsize =
          30)
```

```
plt.xlabel('Average Training Score', fontsize = 20)
```

```
plt.ylabel('count')
```

```
plt.show()
```

13.

```
train['awards_won?'].value_counts()
```

14.

```
# plotting a donut chart for visualizing each of the recruitment channel's
share
```

```
size = [53538, 1270]
```

```
colors = ['magenta', 'brown']
```

```

labels = "Awards Won", "NO Awards Won"

my_circle = plt.Circle((0, 0), 0.7, color = 'white')

plt.rcParams['figure.figsize'] = (9, 9)
plt.pie(size, colors = colors, labels = labels, shadow = True, autopct =
'%.2f%%')
plt.title('Showing a Percentage of employees who won awards', fontsize =
30)
p = plt.gcf()
p.gca().add_artist(my_circle)
plt.legend()
plt.show()

```

15.

```
train['KPIs_met >80%'].value_counts()
```

16.

```
# plotting a pie chart
```

```

size = [35517, 19291]
labels = "Not Met KPI > 80%", "Met KPI > 80%"
colors = ['violet', 'grey']
explode = [0, 0.1]

plt.rcParams['figure.figsize'] = (8, 8)
plt.pie(size, labels = labels, colors = colors, explode = explode, shadow =
True, autopct = "%.2f%%")
plt.title('A Pie Chart Representing Gap in Employees in terms of KPI',
fontsize = 30)
plt.axis('off')
plt.legend()
plt.show()

```

17.

```
# checking the distribution of length of service
```

```

sns.distplot(train['length_of_service'], color = 'green')
plt.title('Distribution of length of service among the Employees', fontsize
= 30)
plt.xlabel('Length of Service in years', fontsize = 15)
plt.ylabel('count')
plt.show()

```

18.

```

train['previous_year_rating'].value_counts().sort_values().plot.bar(color =
'violet', figsize = (15, 7))
plt.title('Distribution of Previous year rating of the Employees', fontsize
= 30)

```

```
plt.xlabel('Ratings', fontsize = 15)
plt.ylabel('count')
plt.show()
```

19.

checking the distribution of age of Employees in the company

```
sns.distplot(train['age'], color = 'red')
plt.title('Distribution of Age of Employees', fontsize = 30)
plt.xlabel('Age', fontsize = 15)
plt.ylabel('count')
plt.show()
```

20.

checking the different no. of training done by the employees

```
plt.rcParams['figure.figsize'] = (17, 7)
sns.violinplot(train['no_of_trainings'], color = 'purple')
plt.title('No. of trainings done by the Employees', fontsize = 30)
plt.xlabel('No. of Trainings', fontsize = 15)
plt.ylabel('Frequency')
plt.show()
```

21.

checking the different types of recruitment channels for the company

```
train['recruitment_channel'].value_counts()
```

22.

plotting a donut chart for visualizing each of the recruitment channel's share

```
size = [30446, 23220, 1142]
colors = ['yellow', 'red', 'lightgreen']
labels = "Others", "Sourcing", "Reffered"
```

```
my_circle = plt.Circle((0, 0), 0.7, color = 'white')
```

```
plt.rcParams['figure.figsize'] = (9, 9)
plt.pie(size, colors = colors, labels = labels, shadow = True, autopct =
'%.2f%%')
plt.title('Showing share of different Recruitment Channels', fontsize = 30)
p = plt.gcf()
p.gca().add_artist(my_circle)
plt.legend()
plt.show()
```

23.

```

# checing the most popular education degree among the employees

from wordcloud import WordCloud
from wordcloud import STOPWORDS

stopword = set(STOPWORDS)

wordcloud = WordCloud(background_color = 'white', stopwords = stopword,
max_words = 5).generate(str(train['education']))

plt.rcParams['figure.figsize'] = (12, 8)
print(wordcloud)
plt.imshow(wordcloud)
plt.title('Most Popular Degrees among the Employees', fontsize = 30)
plt.axis('off')
plt.show()

```

24.

```

# checking the gender gap

train['gender'].value_counts()

```

25.

```

# plotting a pie chart

size = [38496, 16312]
labels = "Male", "Female"
colors = ['yellow', 'orange']
explode = [0, 0.1]

plt.rcParams['figure.figsize'] = (8, 8)
plt.pie(size, labels = labels, colors = colors, explode = explode, shadow =
True, autopct = "%.2f%")
plt.title('A Pie Chart Representing GenderGap', fontsize = 30)
plt.axis('off')
plt.legend()
plt.show()

```

26.

```

# checking the different regions of the company

plt.rcParams['figure.figsize'] = (20, 10)
sns.countplot(train['region'], color = 'pink')
plt.title('Different Regions in the company', fontsize = 30)
plt.xticks(rotation = 60)
plt.xlabel('Region Code', fontsize = 15)
plt.ylabel('count', fontsize = 15)
plt.show()

```

27.

```
# scatter plot between average training score and is_promoted

data = pd.crosstab(train['avg_training_score'], train['is_promoted'])
data.div(data.sum(1).astype(float), axis = 0).plot(kind = 'bar', stacked =
True, figsize = (20, 9), color = ['darkred', 'lightgreen'])

plt.title('Looking at the Dependency of Training Score in promotion',
fontsize = 30)
plt.xlabel('Average Training Scores', fontsize = 15)
plt.legend()
plt.show()
```

28.

```
# checking dependency of different regions in promotion

data = pd.crosstab(train['region'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (20, 8), color = ['lightblue', 'purple'])

plt.title('Dependency of Regions in determining Promotion of Employees',
fontsize = 30)
plt.xlabel('Different Regions of the Company', fontsize = 20)
plt.legend()
plt.show()
```

29.

```
# dependency of awards won on promotion

data = pd.crosstab(train['awards_won?'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (10, 8), color = ['magenta', 'purple'])

plt.title('Dependency of Awards in determining Promotion', fontsize = 30)
plt.xlabel('Awards Won or Not', fontsize = 20)
plt.legend()
plt.show()
```

30.

```
#dependency of KPIs with Promotion

data = pd.crosstab(train['KPIs_met >80%'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (10, 8), color = ['pink', 'darkred'])

plt.title('Dependency of KPIs in determining Promotion', fontsize = 30)
plt.xlabel('KPIs Met or Not', fontsize = 20)
plt.legend()
plt.show()
```

31.

checking dependency on previous years' ratings

```
data = pd.crosstab(train['previous_year_rating'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (15, 8), color = ['violet', 'pink'])
```

```
plt.title('Dependency of Previous year Ratings in determining Promotion',
fontsize = 30)
plt.xlabel('Different Ratings', fontsize = 20)
plt.legend()
plt.show()
```

32.

checking how length of service determines the promotion of employees

```
data = pd.crosstab(train['length_of_service'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (20, 8), color = ['pink', 'lightblue'])
```

```
plt.title('Dependency of Length of service in Promotions of Employees',
fontsize = 30)
plt.xlabel('Length of service of employees', fontsize = 20)
plt.legend()
plt.show()
```

33.

checking dependency of age factor in promotion of employees

```
data = pd.crosstab(train['age'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (20, 8), color = ['lightblue', 'green'])
```

```
plt.title('Dependency of Age in determining Promotion of Employees',
fontsize = 30)
plt.xlabel('Age of Employees', fontsize = 20)
plt.legend()
plt.show()
```

34.

checking which department got most number of promotions

```
data = pd.crosstab(train['department'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (20, 8), color = ['orange', 'lightgreen'])
```

```
plt.title('Dependency of Departments in determining Promotion of
Employees', fontsize = 30)
plt.xlabel('Different Departments of the Company', fontsize = 20)
plt.legend()
plt.show()
```

35.

```
# checking dependency of gender over promotion

data = pd.crosstab(train['gender'], train['is_promoted'])
data.div(data.sum(1).astype('float'), axis = 0).plot(kind = 'bar', stacked
= True, figsize = (7, 5), color = ['pink', 'yellow'])

plt.title('Dependency of Genders in determining Promotion of Employees',
fontsize = 30)
plt.xlabel('Gender', fontsize = 20)
plt.legend()
plt.show()
```

36.

```
# filling missing values

train['education'].fillna(train['education'].mode()[0], inplace = True)
train['previous_year_rating'].fillna(1, inplace = True)

# again checking if there is any Null value left in the data
train.isnull().sum().sum()
```

37.

```
# filling missing values

test['education'].fillna(test['education'].mode()[0], inplace = True)
test['previous_year_rating'].fillna(1, inplace = True)

# again checking if there is any Null value left in the data
test.isnull().sum().sum()
```

38.

```
# removing the employee_id column

train = train.drop(['employee_id'], axis = 1)

train.columns
```

39.

```
# saving the employee_id

emp_id = test['employee_id']

# removing the employee_id column

test = test.drop(['employee_id'], axis = 1)
```



```
test.columns
```

```
40.
```

```
# defining the test set
```

```
x_test = test
```

```
x_test.columns
```

```
41.
```

```
# one hot encoding for the test set
```

```
x_test = pd.get_dummies(x_test)
```

```
x_test.columns
```

```
42.
```

```
# splitting the train set into dependent and independent sets
```

```
x = train.iloc[:, :-1]
```

```
y = train.iloc[:, -1]
```

```
print("Shape of x:", x.shape)
```

```
print("Shape of y:", y.shape)
```

```
43.
```

```
# one hot encoding for the train set
```

```
x = pd.get_dummies(x)
```

```
x.columns
```

```
44.
```