

Data Collection and Preprocessing Phase

Date	9 July 2024
Team ID	SWTID1720104839
Project Title	Human Resource Management: Predicting Employee Promotions Using Machine Learning
Maximum Marks	6 Marks

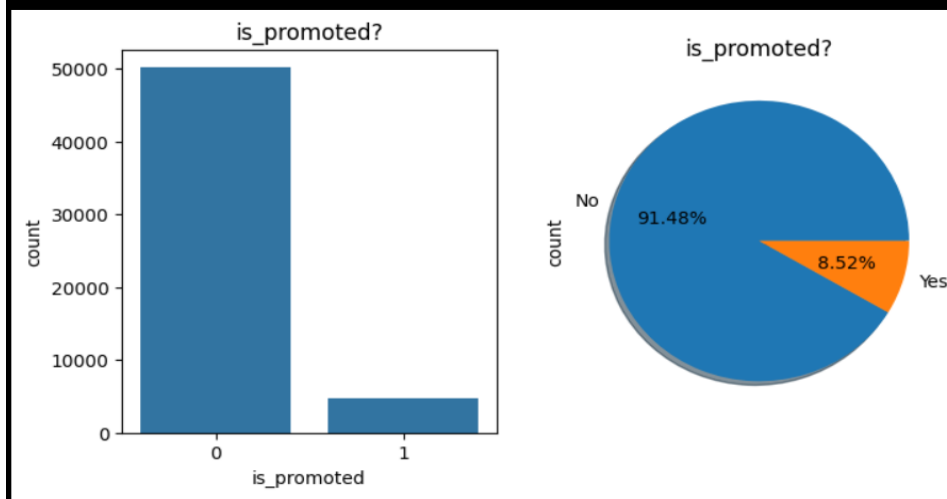
Data Exploration and Preprocessing Template

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

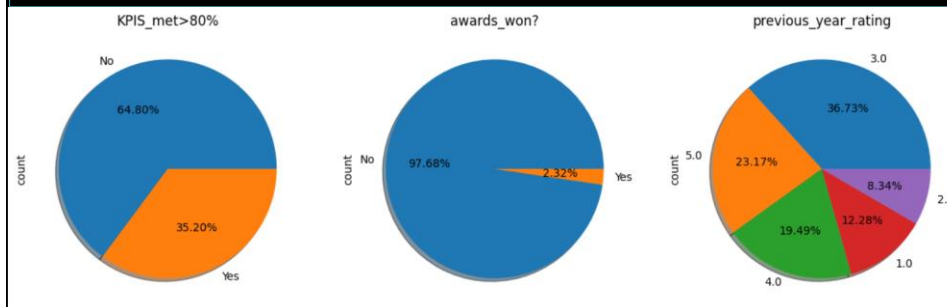
Section	Description
Data Overview	<pre>df.describe(include="all")</pre>

Univariate Analysis

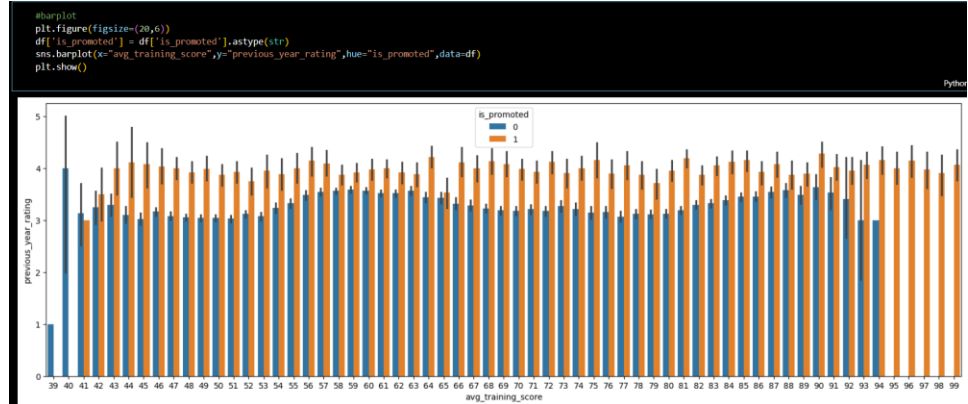
```
plt.figure(figsize=(8,4))
plt.subplot(121)
#countplot
sns.countplot(x="is_promoted",data=df)
plt.title("is_promoted?")
plt.subplot(122)
#pieplot
plt.title("is_promoted?")
val_count=df["is_promoted"].value_counts()
val_count.plot(kind="pie",autopct='%.2f%%',shadow=True,labels=["No","Yes"])
plt.show()
```



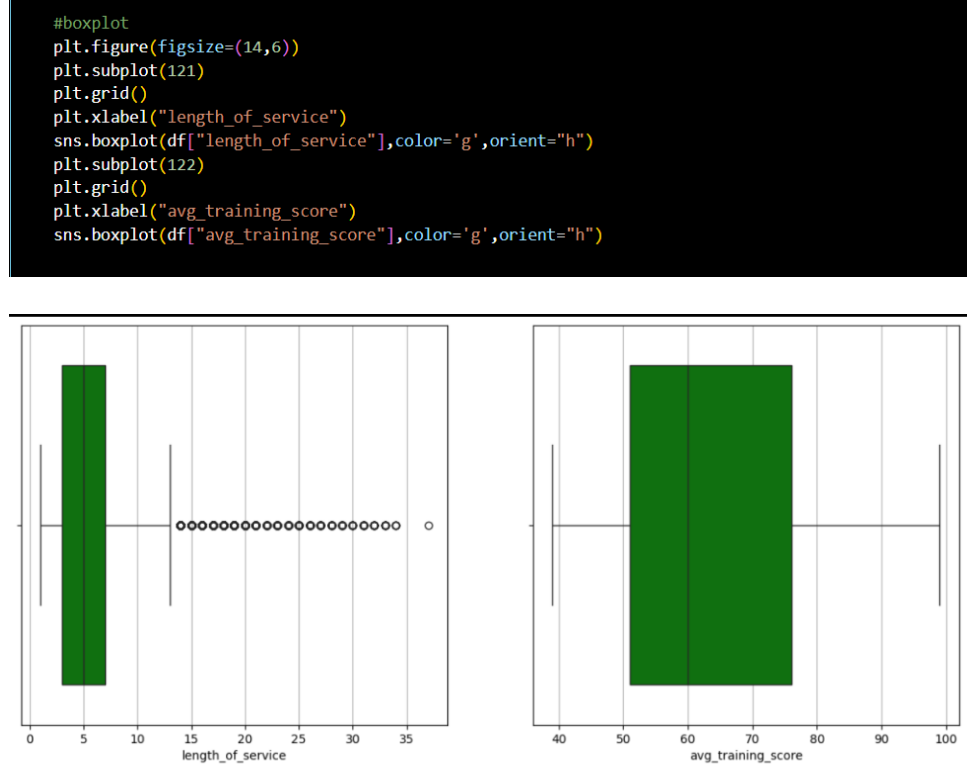
```
plt.figure(figsize=(15,15))
plt.subplot(131)
plt.title("KPIS_met>80%")
df["KPIS_met >80%"].value_counts().plot(kind="pie",shadow=True,autopct='%.2f%%',labels=["No","Yes"])
plt.subplot(132)
plt.title("awards_won?")
df["awards_won?"].value_counts().plot(kind="pie",shadow=True,autopct='%.2f%%',labels=["No","Yes"])
plt.subplot(133)
plt.title("previous_year_rating")
df["previous_year_rating"].value_counts().plot(kind="pie",shadow=True,autopct='%.2f%%')
plt.show()
```



Multivariate Analysis



Outliers and Anomalies



Data Preprocessing Code Screenshots

Loading Data

```
df=pd.read_csv("emp_promotion.csv")
```

Handling Missing Data

Drop unwanted features

```
df=df.drop(["employee_id","gender","region","recruitment_channel"],axis=1)
```

```
df.head()
```

	department	education	no_of_trainings	age	previous_year_rating	length_of_service	KPIs_met >80%	awards_won?	avg_training_score	is_promoted
0	Sales & Marketing	Master's & above	1	35	5.0	8	1	0	49	0
1	Operations	Bachelor's	1	30	5.0	4	0	0	60	0
2	Sales & Marketing	Bachelor's	1	34	3.0	7	0	0	50	0
3	Sales & Marketing	Bachelor's	2	39	1.0	10	0	0	50	0
4	Technology	Bachelor's	1	45	3.0	2	0	0	73	0

Checking for null values

```
df.isnull().sum()
```

```
department      0
education      2409
no_of_trainings 0
age             0
previous_year_rating 4124
length_of_service 0
KPIs_met >80%   0
awards_won?     0
avg_training_score 0
is_promoted     0
dtype: int64
```

```
#handling null values
df["education"]=df["education"].fillna(df["education"].mode()[0])
df["previous_year_rating"]=df["previous_year_rating"].fillna(df["previous_year_rating"].mode()[0])
```

```
df.isnull().sum()
```

```
department      0
education        0
no_of_trainings 0
age             0
previous_year_rating 0
length_of_service 0
KPIs_met >80%   0
awards_won?     0
avg_training_score 0
is_promoted     0
dtype: int64
```

Data Transformation

```
le=LabelEncoder()
df["department"]=le.fit_transform(df["department"])
df["education"]=le.fit_transform(df["education"])
df.head()
```

	department	education	no_of_trainings	age	previous_year_rating	length_of_service	KPIs_met >80%	awards_won?	avg_training_score	is_promoted
0	7	2	1	35	5.0	8.0	1	0	49	0
1	4	0	1	30	5.0	4.0	0	0	60	0
2	7	0	1	34	3.0	7.0	0	0	50	0
3	7	0	2	39	1.0	10.0	0	0	50	0
4	8	0	1	45	3.0	2.0	0	0	73	0

<p>Handling imbalanced data</p>	<h3>Handling Imbalanced data</h3> <pre> #splitting data and resampling it x=df.drop("is_promoted",axis=1) y=df["is_promoted"] print(x.shape) print(y.shape) </pre> <p>[26]</p> <pre> ... (54808, 9) (54808,) </pre> <pre> from imblearn.over_sampling import SMOTE sm=SMOTE() x_resample,y_resample=sm.fit_resample(x,y) </pre> <p>[27]</p> <pre> #visualization before and after smote technique plt.figure(figsize=(10,6)) plt.subplot(121) sns.countplot(x=y,data=df) plt.title("before smote technique") plt.subplot(122) sns.countplot(x=y_resample,data=df) plt.title("after smote technique") plt.show() </pre> <p>[28]</p>
<p>Spilting data</p>	<h3>Splitting data into train and test</h3> <pre> x_train,x_test,y_train,y_test=train_test_split(x_resample,y_resample,test_size=0.2,random_state=0) </pre> <p>[29]</p>