

Data Collection and Preprocessing Phase

Date	13 JULY 2024
Team ID	SWTID1720174957
Project Title	Human resource management employee promotion prediction 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	Basic statistics, dimensions, and structure of the data.
Univariate Analysis	Exploration of individual variables (mean, median, mode, etc.).
Bivariate Analysis	Relationships between two variables (correlation, scatter plots).
Multivariate Analysis	Patterns and relationships involving multiple variables.
Outliers and Anomalies	Identification and treatment of outliers.
Data Preprocessing Code Screenshots	
Loading Data	<pre>df = pd.read_csv('C:\Dataset\emp_promotion.csv') print('Shape of train data {}'.format(df.shape))</pre> <p>Shape of train data (54808, 14)</p>

Handling Missing Data

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

```
department      0
education      2489
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
```

```
print(df['education'].value_counts())
df['education']=df['education'].fillna(df['education'].mode()[0])
```

```
education
Bachelor's      36669
Master's & above 14925
Below Secondary  885
Name: count, dtype: int64
```

```
#Replacing nan with mode
```

```
print(df['previous_year_rating'].value_counts())
df['previous_year_rating']=df['previous_year_rating'].fillna(df['previous_year_rating'].mode()[0])
```

```
previous_year_rating
3.0    18618
5.0    11741
4.0     9877
1.0     6223
2.0     4225
Name: count, dtype: int64
```

Data Transformation

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler
```

```
# Select numerical columns
```

```
numerical_columns = ['no_of_trainings', 'age', 'previous_year_rating',
                    'length_of_service', 'KPIs_met >80%', 'awards_won?',
                    'avg_training_score']
```

```
# Standard scaling
```

```
scaler = StandardScaler()
```

```
data_standard_scaled = pd.read_csv('C:\Dataset\emp_promotion.csv')
```

```
data_standard_scaled[numerical_columns] = scaler.fit_transform(data_standard_scaled[numerical_columns])
```

```
# Min-Max normalization
```

```
min_max_scaler = MinMaxScaler()
```

```
data_min_max_scaled = pd.read_csv('C:\Dataset\emp_promotion.csv')
```

```
data_min_max_scaled[numerical_columns] = min_max_scaler.fit_transform(data_min_max_scaled[numerical_columns])
```

```
# Display the first few rows of the scaled and normalized data
```

```
print("Standard Scaled Data:\n", data_standard_scaled.head())
```

```
print("\nMin-Max Normalized Data:\n", data_min_max_scaled.head())
```

```
Standard Scaled Data:
```

```
employee_id  department  region  education gender \
0    65438  Sales & Marketing  region_7  Master's & above  f
1    65141  Operations      region_22  Bachelor's      m
2    7513   Sales & Marketing  region_19  Bachelor's      m
3    2542   Sales & Marketing  region_23  Bachelor's      m
4    48945  Technology      region_26  Bachelor's      m
```




```
recruitment_channel  no_of_trainings  age  previous_year_rating \
0      sourcing      -0.415276  0.025598      1.326008
1      other      -0.415276  -0.627135      1.326008
2      sourcing      -0.415276  -0.104948      -0.261318
3      other      1.226063  0.547785      -1.848645
4      other      -0.415276  1.331064      -0.261318
```

```
length_of_service  KPIs_met >80%  awards_won?  avg_training_score \
0    0.500460      1.356878      -0.154018      -1.075931
1    -0.437395      -0.736986      -0.154018      -0.253282
2    0.265996      -0.736986      -0.154018      -1.001145
3    0.969387      -0.736986      -0.154018      -1.001145
4    -0.906322      -0.736986      -0.154018      0.718939
```

```
is_promoted
```

```
0    0
1    0
2    0
3    0
4    0
```

	<pre>Min-Max Normalized Data:</pre> <table><thead><tr><th></th><th>employee_id</th><th>department</th><th>region</th><th>education</th><th>gender</th><th>\</th></tr></thead><tbody><tr><td>0</td><td>65438</td><td>Sales & Marketing</td><td>region_7</td><td>Master's & above</td><td>f</td><td></td></tr><tr><td>1</td><td>65141</td><td>Operations</td><td>region_22</td><td>Bachelor's</td><td>m</td><td></td></tr><tr><td>2</td><td>7513</td><td>Sales & Marketing</td><td>region_19</td><td>Bachelor's</td><td>m</td><td></td></tr><tr><td>3</td><td>2542</td><td>Sales & Marketing</td><td>region_23</td><td>Bachelor's</td><td>m</td><td></td></tr><tr><td>4</td><td>48945</td><td>Technology</td><td>region_26</td><td>Bachelor's</td><td>m</td><td></td></tr></tbody></table> <table><thead><tr><th></th><th>recruitment_channel</th><th>no_of_trainings</th><th>age</th><th>previous_year_rating</th><th>\</th></tr></thead><tbody><tr><td>0</td><td>sourcing</td><td>0.000000</td><td>0.375</td><td>1.0</td><td></td></tr><tr><td>1</td><td>other</td><td>0.000000</td><td>0.250</td><td>1.0</td><td></td></tr><tr><td>2</td><td>sourcing</td><td>0.000000</td><td>0.350</td><td>0.5</td><td></td></tr><tr><td>3</td><td>other</td><td>0.111111</td><td>0.475</td><td>0.0</td><td></td></tr><tr><td>4</td><td>other</td><td>0.000000</td><td>0.625</td><td>0.5</td><td></td></tr></tbody></table> <table><thead><tr><th></th><th>length_of_service</th><th>KPIs_met >80%</th><th>awards_won?</th><th>avg_training_score</th><th>\</th></tr></thead><tbody><tr><td>0</td><td>0.194444</td><td>1.0</td><td>0.0</td><td>0.166667</td><td></td></tr><tr><td>1</td><td>0.083333</td><td>0.0</td><td>0.0</td><td>0.350000</td><td></td></tr><tr><td>2</td><td>0.166667</td><td>0.0</td><td>0.0</td><td>0.183333</td><td></td></tr><tr><td>3</td><td>0.250000</td><td>0.0</td><td>0.0</td><td>0.183333</td><td></td></tr><tr><td>4</td><td>0.027778</td><td>0.0</td><td>0.0</td><td>0.566667</td><td></td></tr></tbody></table> <table><thead><tr><th></th><th>is_promoted</th></tr></thead><tbody><tr><td>0</td><td>0</td></tr><tr><td>1</td><td>0</td></tr><tr><td>2</td><td>0</td></tr><tr><td>3</td><td>0</td></tr><tr><td>4</td><td>0</td></tr></tbody></table>		employee_id	department	region	education	gender	\	0	65438	Sales & Marketing	region_7	Master's & above	f		1	65141	Operations	region_22	Bachelor's	m		2	7513	Sales & Marketing	region_19	Bachelor's	m		3	2542	Sales & Marketing	region_23	Bachelor's	m		4	48945	Technology	region_26	Bachelor's	m			recruitment_channel	no_of_trainings	age	previous_year_rating	\	0	sourcing	0.000000	0.375	1.0		1	other	0.000000	0.250	1.0		2	sourcing	0.000000	0.350	0.5		3	other	0.111111	0.475	0.0		4	other	0.000000	0.625	0.5			length_of_service	KPIs_met >80%	awards_won?	avg_training_score	\	0	0.194444	1.0	0.0	0.166667		1	0.083333	0.0	0.0	0.350000		2	0.166667	0.0	0.0	0.183333		3	0.250000	0.0	0.0	0.183333		4	0.027778	0.0	0.0	0.566667			is_promoted	0	0	1	0	2	0	3	0	4	0												
	employee_id	department	region	education	gender	\																																																																																																																																					
0	65438	Sales & Marketing	region_7	Master's & above	f																																																																																																																																						
1	65141	Operations	region_22	Bachelor's	m																																																																																																																																						
2	7513	Sales & Marketing	region_19	Bachelor's	m																																																																																																																																						
3	2542	Sales & Marketing	region_23	Bachelor's	m																																																																																																																																						
4	48945	Technology	region_26	Bachelor's	m																																																																																																																																						
	recruitment_channel	no_of_trainings	age	previous_year_rating	\																																																																																																																																						
0	sourcing	0.000000	0.375	1.0																																																																																																																																							
1	other	0.000000	0.250	1.0																																																																																																																																							
2	sourcing	0.000000	0.350	0.5																																																																																																																																							
3	other	0.111111	0.475	0.0																																																																																																																																							
4	other	0.000000	0.625	0.5																																																																																																																																							
	length_of_service	KPIs_met >80%	awards_won?	avg_training_score	\																																																																																																																																						
0	0.194444	1.0	0.0	0.166667																																																																																																																																							
1	0.083333	0.0	0.0	0.350000																																																																																																																																							
2	0.166667	0.0	0.0	0.183333																																																																																																																																							
3	0.250000	0.0	0.0	0.183333																																																																																																																																							
4	0.027778	0.0	0.0	0.566667																																																																																																																																							
	is_promoted																																																																																																																																										
0	0																																																																																																																																										
1	0																																																																																																																																										
2	0																																																																																																																																										
3	0																																																																																																																																										
4	0																																																																																																																																										
Feature Engineering	<pre>import pandas as pd # Load the CSV file data = pd.read_csv('C:\Dataset\emp_promotion.csv') # Create a new feature 'is_high_performer' data['is_high_performer'] = data.apply(lambda row: 1 if row['KPIs_met >80%'] == 1 and row['awards_won?'] == 1 else 0, axis=1) # Create an 'age_group' feature bins = [20, 30, 40, 50, 60] labels = ['20-29', '30-39', '40-49', '50-59'] data['age_group'] = pd.cut(data['age'], bins=bins, labels=labels, right=False) # Modify the 'education' feature to fewer categories data['education'] = data['education'].replace({ "Master's & above": "Postgraduate", "Bachelor's": "Undergraduate", "Below Secondary": "Secondary" }) # Display the first few rows to verify changes print(data.head())</pre> <table><thead><tr><th></th><th>employee_id</th><th>department</th><th>region</th><th>education</th><th>gender</th><th>\</th></tr></thead><tbody><tr><td>0</td><td>65438</td><td>Sales & Marketing</td><td>region_7</td><td>Postgraduate</td><td>f</td><td></td></tr><tr><td>1</td><td>65141</td><td>Operations</td><td>region_22</td><td>Undergraduate</td><td>m</td><td></td></tr><tr><td>2</td><td>7513</td><td>Sales & Marketing</td><td>region_19</td><td>Undergraduate</td><td>m</td><td></td></tr><tr><td>3</td><td>2542</td><td>Sales & Marketing</td><td>region_23</td><td>Undergraduate</td><td>m</td><td></td></tr><tr><td>4</td><td>48945</td><td>Technology</td><td>region_26</td><td>Undergraduate</td><td>m</td><td></td></tr></tbody></table> <table><thead><tr><th></th><th>recruitment_channel</th><th>no_of_trainings</th><th>age</th><th>previous_year_rating</th><th>\</th></tr></thead><tbody><tr><td>0</td><td>sourcing</td><td>1</td><td>35</td><td>5.0</td><td></td></tr><tr><td>1</td><td>other</td><td>1</td><td>30</td><td>5.0</td><td></td></tr><tr><td>2</td><td>sourcing</td><td>1</td><td>34</td><td>3.0</td><td></td></tr><tr><td>3</td><td>other</td><td>2</td><td>39</td><td>1.0</td><td></td></tr><tr><td>4</td><td>other</td><td>1</td><td>45</td><td>3.0</td><td></td></tr></tbody></table> <table><thead><tr><th></th><th>length_of_service</th><th>KPIs_met >80%</th><th>awards_won?</th><th>avg_training_score</th><th>\</th></tr></thead><tbody><tr><td>0</td><td>8</td><td>1</td><td>0</td><td>49</td><td></td></tr><tr><td>1</td><td>4</td><td>0</td><td>0</td><td>60</td><td></td></tr><tr><td>2</td><td>7</td><td>0</td><td>0</td><td>50</td><td></td></tr><tr><td>3</td><td>10</td><td>0</td><td>0</td><td>50</td><td></td></tr><tr><td>4</td><td>2</td><td>0</td><td>0</td><td>73</td><td></td></tr></tbody></table> <table><thead><tr><th></th><th>is_promoted</th><th>is_high_performer</th><th>age_group</th></tr></thead><tbody><tr><td>0</td><td>0</td><td>0</td><td>30-39</td></tr><tr><td>1</td><td>0</td><td>0</td><td>30-39</td></tr><tr><td>2</td><td>0</td><td>0</td><td>30-39</td></tr><tr><td>3</td><td>0</td><td>0</td><td>30-39</td></tr><tr><td>4</td><td>0</td><td>0</td><td>40-49</td></tr></tbody></table>		employee_id	department	region	education	gender	\	0	65438	Sales & Marketing	region_7	Postgraduate	f		1	65141	Operations	region_22	Undergraduate	m		2	7513	Sales & Marketing	region_19	Undergraduate	m		3	2542	Sales & Marketing	region_23	Undergraduate	m		4	48945	Technology	region_26	Undergraduate	m			recruitment_channel	no_of_trainings	age	previous_year_rating	\	0	sourcing	1	35	5.0		1	other	1	30	5.0		2	sourcing	1	34	3.0		3	other	2	39	1.0		4	other	1	45	3.0			length_of_service	KPIs_met >80%	awards_won?	avg_training_score	\	0	8	1	0	49		1	4	0	0	60		2	7	0	0	50		3	10	0	0	50		4	2	0	0	73			is_promoted	is_high_performer	age_group	0	0	0	30-39	1	0	0	30-39	2	0	0	30-39	3	0	0	30-39	4	0	0	40-49
	employee_id	department	region	education	gender	\																																																																																																																																					
0	65438	Sales & Marketing	region_7	Postgraduate	f																																																																																																																																						
1	65141	Operations	region_22	Undergraduate	m																																																																																																																																						
2	7513	Sales & Marketing	region_19	Undergraduate	m																																																																																																																																						
3	2542	Sales & Marketing	region_23	Undergraduate	m																																																																																																																																						
4	48945	Technology	region_26	Undergraduate	m																																																																																																																																						
	recruitment_channel	no_of_trainings	age	previous_year_rating	\																																																																																																																																						
0	sourcing	1	35	5.0																																																																																																																																							
1	other	1	30	5.0																																																																																																																																							
2	sourcing	1	34	3.0																																																																																																																																							
3	other	2	39	1.0																																																																																																																																							
4	other	1	45	3.0																																																																																																																																							
	length_of_service	KPIs_met >80%	awards_won?	avg_training_score	\																																																																																																																																						
0	8	1	0	49																																																																																																																																							
1	4	0	0	60																																																																																																																																							
2	7	0	0	50																																																																																																																																							
3	10	0	0	50																																																																																																																																							
4	2	0	0	73																																																																																																																																							
	is_promoted	is_high_performer	age_group																																																																																																																																								
0	0	0	30-39																																																																																																																																								
1	0	0	30-39																																																																																																																																								
2	0	0	30-39																																																																																																																																								
3	0	0	30-39																																																																																																																																								
4	0	0	40-49																																																																																																																																								
Save Processed Data	<pre>import pickle pickle.dump(rf, open('promotion_model.pkl', 'wb'))</pre>																																																																																																																																										

	Today			
	 promotion_model.pkl	14-07-2024 17:33	PKL File	1,11,024 KB
	 final	14-07-2024 17:35	Jupyter Source File	2,287 KB
	 emp_promotion	14-07-2024 16:50	Microsoft Excel Co...	3,672 KB