

Data Collection and Preprocessing Phase

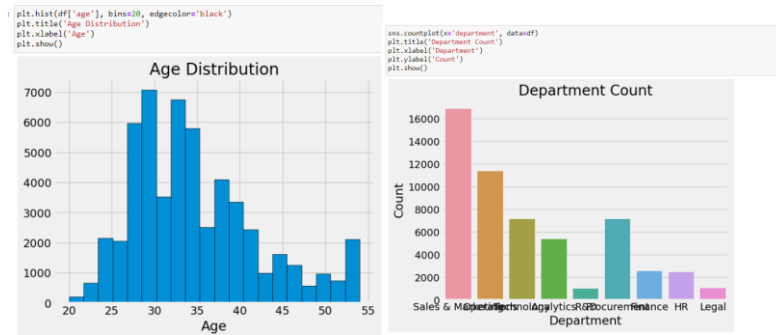
Date	09 July 2024
Team ID	SWTID1720455879
Project Title	Human Resource Management: Predicting Employee Promotions Using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

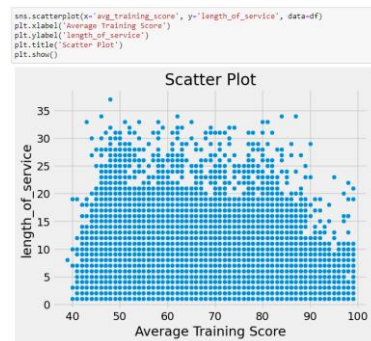
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description																																																															
Data Overview	<div><pre>df.describe()</pre></div> <table><thead><tr><th></th><th>employee_id</th><th>no_of_trainings</th><th>age</th><th>previous_year_rating</th><th>length_of_service</th><th>KPI</th></tr></thead><tbody><tr><td>count</td><td>54808.000000</td><td>54808.000000</td><td>54808.000000</td><td>50684.000000</td><td>54808.000000</td><td>54808.000000</td></tr><tr><td>mean</td><td>39195.830627</td><td>1.253011</td><td>34.803915</td><td>3.329256</td><td>5.865512</td><td>0.300000</td></tr><tr><td>std</td><td>22586.581449</td><td>0.609264</td><td>7.660169</td><td>1.259993</td><td>4.265094</td><td>0.400000</td></tr><tr><td>min</td><td>1.000000</td><td>1.000000</td><td>20.000000</td><td>1.000000</td><td>1.000000</td><td>0.000000</td></tr><tr><td>25%</td><td>19669.750000</td><td>1.000000</td><td>29.000000</td><td>3.000000</td><td>3.000000</td><td>0.000000</td></tr><tr><td>50%</td><td>39225.500000</td><td>1.000000</td><td>33.000000</td><td>3.000000</td><td>5.000000</td><td>0.000000</td></tr><tr><td>75%</td><td>58730.500000</td><td>1.000000</td><td>39.000000</td><td>4.000000</td><td>7.000000</td><td>1.000000</td></tr><tr><td>max</td><td>78298.000000</td><td>10.000000</td><td>60.000000</td><td>5.000000</td><td>37.000000</td><td>1.000000</td></tr></tbody></table> <div><div></div></div> <div><pre>: df.shape</pre><pre>: (54808, 14)</pre></div>		employee_id	no_of_trainings	age	previous_year_rating	length_of_service	KPI	count	54808.000000	54808.000000	54808.000000	50684.000000	54808.000000	54808.000000	mean	39195.830627	1.253011	34.803915	3.329256	5.865512	0.300000	std	22586.581449	0.609264	7.660169	1.259993	4.265094	0.400000	min	1.000000	1.000000	20.000000	1.000000	1.000000	0.000000	25%	19669.750000	1.000000	29.000000	3.000000	3.000000	0.000000	50%	39225.500000	1.000000	33.000000	3.000000	5.000000	0.000000	75%	58730.500000	1.000000	39.000000	4.000000	7.000000	1.000000	max	78298.000000	10.000000	60.000000	5.000000	37.000000	1.000000
		employee_id	no_of_trainings	age	previous_year_rating	length_of_service	KPI																																																									
	count	54808.000000	54808.000000	54808.000000	50684.000000	54808.000000	54808.000000																																																									
	mean	39195.830627	1.253011	34.803915	3.329256	5.865512	0.300000																																																									
	std	22586.581449	0.609264	7.660169	1.259993	4.265094	0.400000																																																									
	min	1.000000	1.000000	20.000000	1.000000	1.000000	0.000000																																																									
	25%	19669.750000	1.000000	29.000000	3.000000	3.000000	0.000000																																																									
	50%	39225.500000	1.000000	33.000000	3.000000	5.000000	0.000000																																																									
	75%	58730.500000	1.000000	39.000000	4.000000	7.000000	1.000000																																																									
	max	78298.000000	10.000000	60.000000	5.000000	37.000000	1.000000																																																									

Univariate Analysis



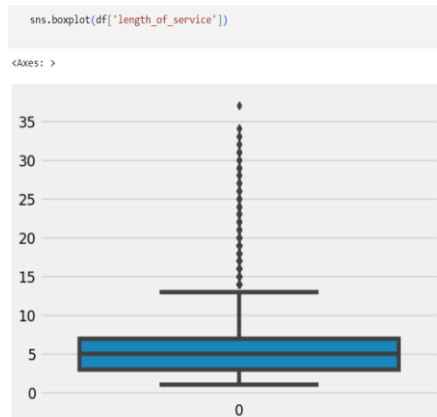
Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies



	<pre># Handle outliers with capping numerical_cols = ['no_of_trainings', 'age', 'previous_year_rating', 'length_of_service', 'avg_training_score'] for col in numerical_cols: Q1 = df[col].quantile(0.25) Q3 = df[col].quantile(0.75) IQR = Q3 - Q1 lower_bound = Q1 - 1.5 * IQR upper_bound = Q3 + 1.5 * IQR df[col] = np.where(df[col] < lower_bound, lower_bound, df[col]) df[col] = np.where(df[col] > upper_bound, upper_bound, df[col]) q1=np.quantile(df['length_of_service'],0.25) q3=np.quantile(df['length_of_service'],0.75) IQR= q3-q1 upper_bound=(1.5*IQR)+q3 lower_bound=(1.5*IQR)-q1 print("Skewed data:",len(df[df['length_of_service']>upper_bound])) pd.crosstab([df['length_of_service']>upper_bound],df['is_promoted']) df['length_of_service']=[upper_bound if x>upper_bound else x for x in df['length_of_service']] pd.crosstab([df['length_of_service']<lower_bound],df['is_promoted']) df['length_of_service']=[upper_bound if x<lower_bound else x for x in df['length_of_service']]</pre>
Data Preprocessing Code Screenshots	
Loading Data	<pre># Load your dataset df = pd.read_csv("emp_promotion.csv") df</pre>
Handling Missing Data	<pre># Drop unwanted features df = df.drop(['employee_id', 'region', 'gender', 'recruitment_channel'], axis=1) df print(df.isnull().sum()) print(df['education'].value_counts()) df['education'] = df['education'].fillna(df['education'].mode()[0]) print(df['previous_year_rating'].value_counts()) df['previous_year_rating'] = df['previous_year_rating'].fillna(df['previous_year_rating'].mode()[0]) print(df.isnull().sum())</pre>
Data Transformation	<pre>le = LabelEncoder() df['department'] = le.fit_transform(df['department']) df['education'] = le.fit_transform(df['education']) df.head()</pre>
Handling Imbalanced Data	<pre>from imblearn.over_sampling import SMOTE sm=SMOTE() X_new,y_new=sm.fit_resample(X,y) count_0 = np.count_nonzero(y_new == 0) count_1 = np.count_nonzero(y_new== 1) print(f"Number of 0s after sampling: {count_0}") print(f"Number of 1s after sampling: {count_1}")</pre>
Splitting data	<pre>X_train,X_test,y_train,y_test=train_test_split(X_new,y_new,test_size=0.3,random_state=42)</pre>