## Data Collection and Preprocessing Phase
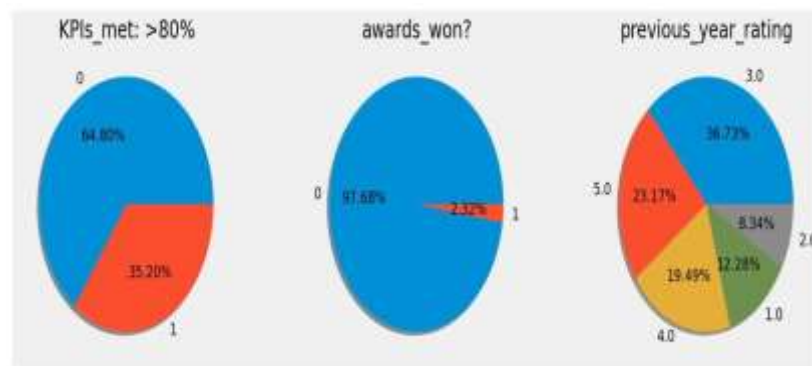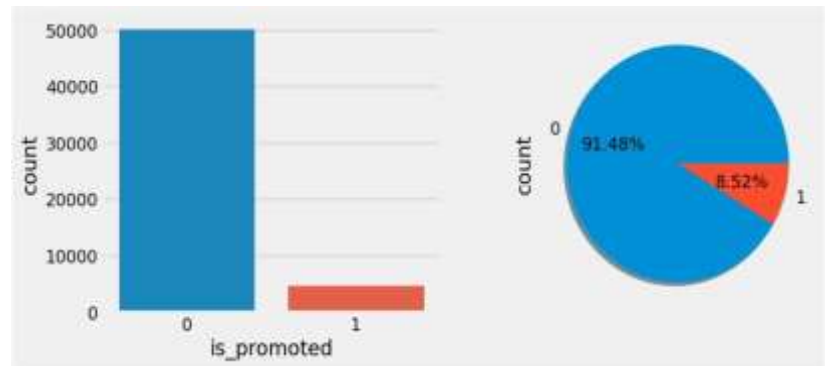
| | |
|---|---|
| Date | 10 July 2024 |
| Team ID | xxxxxx |
| Project Title | Human Resource Management: Predicting Employee Promotions Using Machine Learning |
| Maximum Marks | 6 Marks |

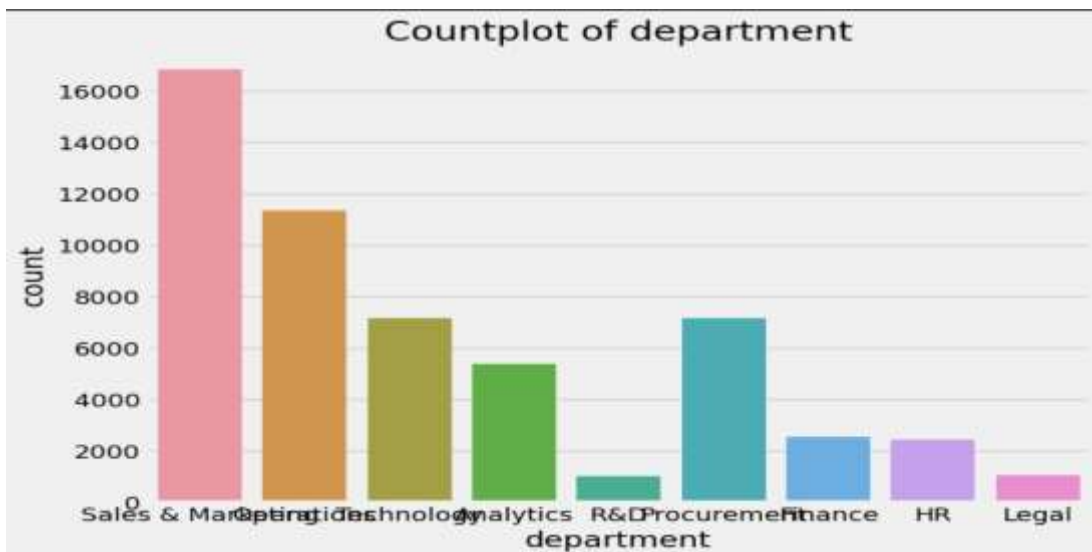Data Exploration and Preprocessing Report

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 | Dimensions:<br>54808 rows × 14 columns    Descriptive statistics:<br> |

| | |
|---|---|
| Univariate Analysis |  |
| Bivariate Analysis |  |

| | |
|---|---|
| Multivariate Analysis |  |

| | |
|---|---|
| Outliers and Anomalies |  |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data |  |

| Handling Missing Data | ```
# Replacing nan with mode

print(df['education'].value_counts())
df['education']=df['education'].fillna(df['education'].mode()[0])
# 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])
``` |
|---|---|
| Data Transformation | ```
# feature mapping is done on education column

df['education']=df['education'].replace(("Below Secondary","Bachelor's","Master's & above"),(1,2,3))

lb = LabelEncoder()
df['department']=lb.fit_transform(df['department'])
``` |
| Feature Engineering | Attached the codes in final submission |

| Save Processed Data | - |
|---|---|