# Data Collection and Preprocessing Phase

| Date | 15 March 2024 |
|------|---------------|
| Team ID | 740071 |
| Project Title | Work Force Retention System |
| Maximum Marks | 6 Marks |

## Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

| Section | Description |
|---------|-------------|

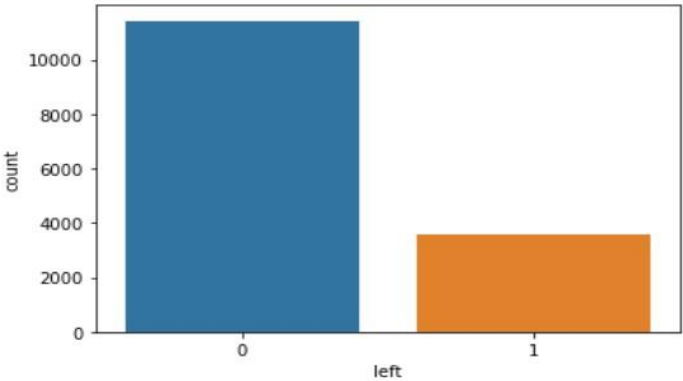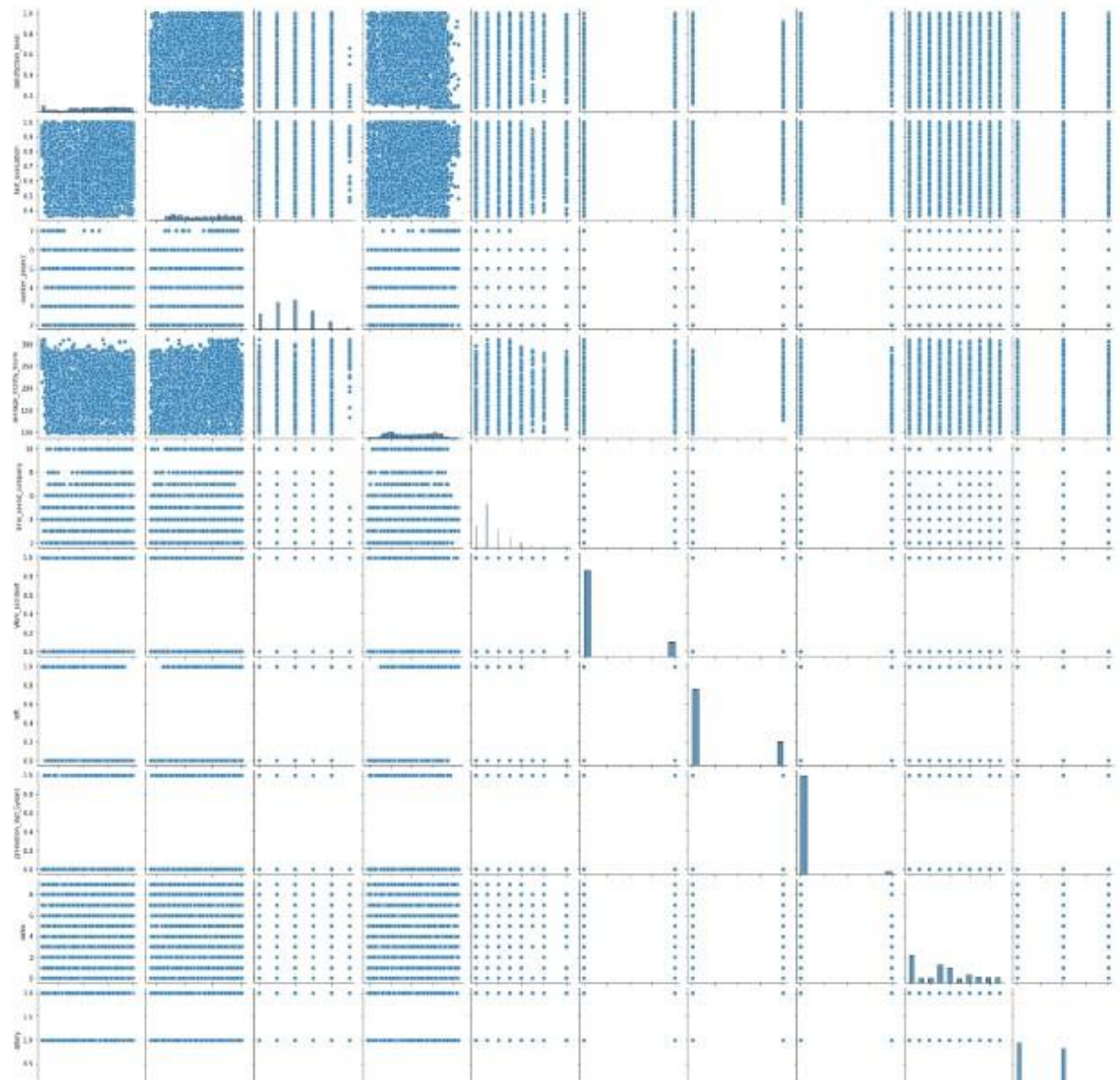| | |
|---|---|
| Data Overview | Dimension:<br> 14999 rows × 10 columns<br>Descriptive statistics:<br><br>`df.describe()`<br><br>| | satisfaction_level | last_evaluation | number_project | average_montly_hours | time_spend_company | Work_accident | left | promotion_last_5years |<br>|---|---|---|---|---|---|---|---|---|<br>| count | 14999.000000 | 14999.000000 | 14999.000000 | 14999.000000 | 14999.000000 | 14999.000000 | 14999.000000 | 14999.000000 |<br>| mean | 0.612834 | 0.716102 | 3.803054 | 201.050337 | 3.498233 | 0.144610 | 0.238083 | 0.021268 |<br>| std | 0.248631 | 0.171169 | 1.232592 | 49.943099 | 1.460136 | 0.351719 | 0.425924 | 0.144281 |<br>| min | 0.090000 | 0.360000 | 2.000000 | 96.000000 | 2.000000 | 0.000000 | 0.000000 | 0.000000 |<br>| 25% | 0.440000 | 0.560000 | 3.000000 | 156.000000 | 3.000000 | 0.000000 | 0.000000 | 0.000000 |<br>| 50% | 0.640000 | 0.720000 | 4.000000 | 200.000000 | 3.000000 | 0.000000 | 0.000000 | 0.000000 |<br>| 75% | 0.820000 | 0.870000 | 5.000000 | 245.000000 | 4.000000 | 0.000000 | 0.000000 | 0.000000 |<br>| max | 1.000000 | 1.000000 | 7.000000 | 310.000000 | 10.000000 | 1.000000 | 1.000000 | 1.000000 | |
| Univariate Analysis | |
| | `sns.countplot(df['left'])`<br><br>`<AxesSubplot:xlabel='left', ylabel='count'>`<br><br> |

| Bivariate Analysis | To find the relation between two features we use bivariate analysis. You can use seaborn package to plot visualisation uisng two variables of the datase |
| --- | --- |

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

```
sns.pairplot(df)
```

```
<seaborn.axisgrid.PairGrid at 0x2ab22b4e5b0>
```

| Loading Data | **Loading Data**<br><br>`: df = pd.read_csv('HR_comma_sep.csv')`<br><br>`: df`<br><br>`:`<br><br>| | satisfaction_level | last_evaluation | number_project | average_montly_hours | time_spend_company | Work_accident | left | promotion_las |<br>|---|---|---|---|---|---|---|---|---|<br>| 0 | 0.38 | 0.53 | 2 | 157 | 3 | 0 | 1 | |<br>| 1 | 0.80 | 0.86 | 5 | 262 | 6 | 0 | 1 | |<br>| 2 | 0.11 | 0.88 | 7 | 272 | 4 | 0 | 1 | |<br>| 3 | 0.72 | 0.87 | 5 | 223 | 5 | 0 | 1 | |<br>| 4 | 0.37 | 0.52 | 2 | 159 | 3 | 0 | 1 | |<br>| ... | ... | ... | ... | ... | ... | ... | ... | |<br>| 14994 | 0.40 | 0.57 | 2 | 151 | 3 | 0 | 1 | |<br>| 14995 | 0.37 | 0.48 | 2 | 160 | 3 | 0 | 1 | |<br>| 14996 | 0.37 | 0.53 | 2 | 143 | 3 | 0 | 1 | |<br>| 14997 | 0.11 | 0.96 | 6 | 280 | 4 | 0 | 1 | |<br>| 14998 | 0.37 | 0.52 | 2 | 158 | 3 | 0 | 1 | |<br><br>14999 rows × 10 columns |

| Handling missing Values | ```
df.shape
``` |
| --- | --- |

```
(14999, 10)
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14999 entries, 0 to 14998
Data columns (total 10 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
 0   satisfaction_level     14999 non-null  float64
 1   last_evaluation        14999 non-null  float64
 2   number_project         14999 non-null  int64
 3   average_montly_hours   14999 non-null  int64
 4   time_spend_company     14999 non-null  int64
 5   Work_accident          14999 non-null  int64
 6   left                   14999 non-null  int64
 7   promotion_last_5years  14999 non-null  int64
 8   sales                  14999 non-null  object
 9   salary                 14999 non-null  object
dtypes: float64(2), int64(6), object(2)
memory usage: 1.1+ MB
```

| | |
|---|---|
| **Feature Engineering**<br><br>**Save processed data** | Attached the codes in final submission<br><br><br>- |
| | |
| | |