# Data Collection and Preprocessing Phase

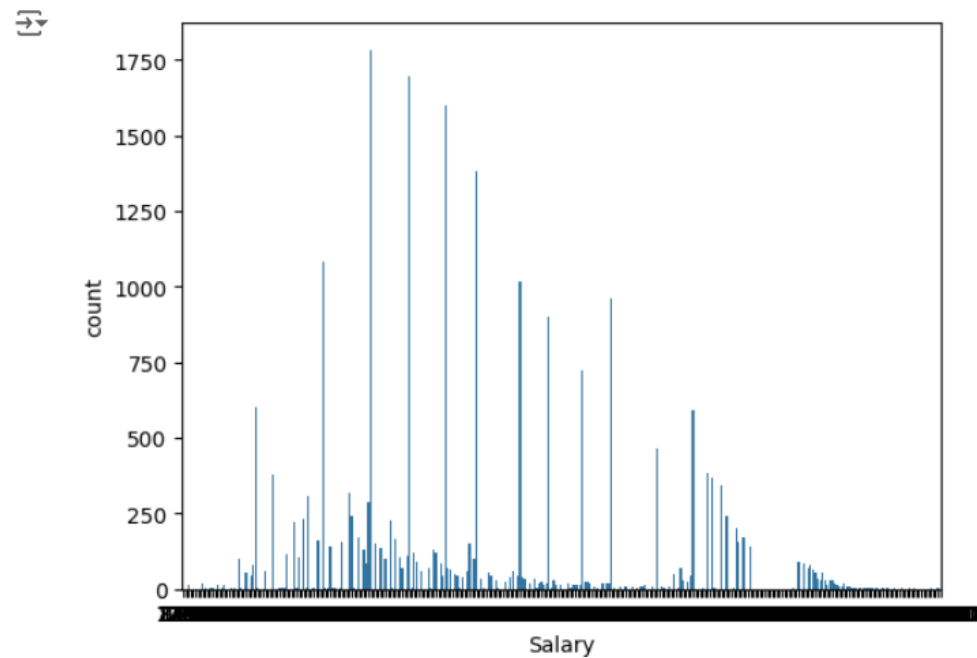| | |
|---|---|
| Date | 31 June 2024 |
| Team ID | 739854 |
| Project Title | Software Salary Prediction |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing**

Dataset gets statistically analysed to identify patterns and outliers. Data preprocessing addresses missing values, improving data quality for further analysis and modelling, and forming a strong foundation for insights and predictions.

| Section | Description |
|---|---|
| Data Overview | Dimensions: 64rows*8 Columns<br>Descriptive Statistics:<br><br>`[ ] df.describe()` |

<table>
<thead>
<tr><th></th><th>Rating</th><th>Company Name</th><th>Job Title</th><th>Salary</th><th>Salaries Reported</th><th>Location</th><th>Employment Status</th><th>Job Roles</th></tr>
</thead>
<tbody>
<tr><td>count</td><td>22770.000000</td><td>22770.000000</td><td>22770.000000</td><td>2.277000e+04</td><td>22770.000000</td><td>22770.000000</td><td>22770.000000</td><td>22770.000000</td></tr>
<tr><td>mean</td><td>3.918213</td><td>5478.825209</td><td>597.435968</td><td>6.953872e+05</td><td>1.855775</td><td>3.150812</td><td>1.071322</td><td>5.465086</td></tr>
<tr><td>std</td><td>0.519675</td><td>3224.603280</td><td>348.305504</td><td>8.843990e+05</td><td>6.823668</td><td>3.529116</td><td>0.342450</td><td>3.221968</td></tr>
<tr><td>min</td><td>1.000000</td><td>0.000000</td><td>0.000000</td><td>2.112000e+03</td><td>1.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td></tr>
<tr><td>25%</td><td>3.700000</td><td>2756.000000</td><td>237.000000</td><td>3.000000e+05</td><td>1.000000</td><td>0.000000</td><td>1.000000</td><td>3.000000</td></tr>
<tr><td>50%</td><td>3.900000</td><td>5317.500000</td><td>753.000000</td><td>5.000000e+05</td><td>1.000000</td><td>2.000000</td><td>1.000000</td><td>7.000000</td></tr>
<tr><td>75%</td><td>4.200000</td><td>8336.000000</td><td>850.000000</td><td>9.000000e+05</td><td>1.000000</td><td>8.000000</td><td>1.000000</td><td>8.000000</td></tr>
<tr><td>max</td><td>5.000000</td><td>11260.000000</td><td>1079.000000</td><td>9.000000e+07</td><td>361.000000</td><td>9.000000</td><td>3.000000</td><td>10.000000</td></tr>
</tbody>
</table>

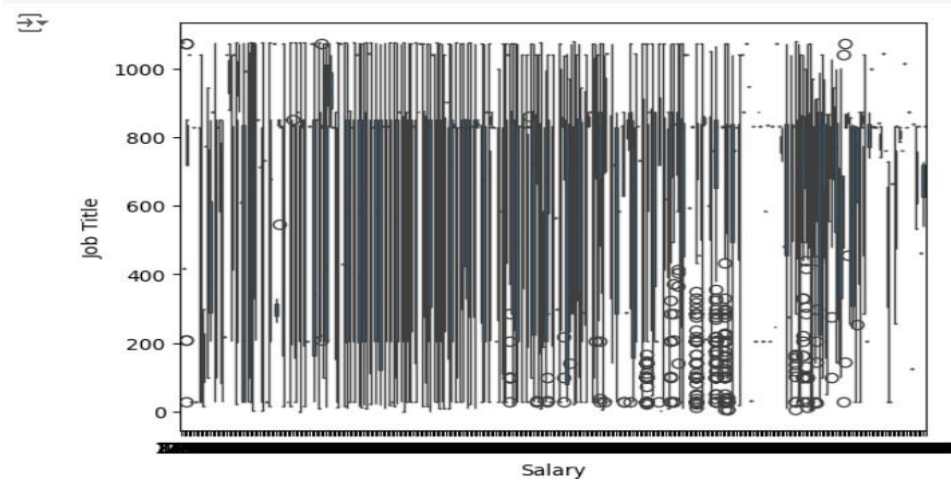| | |
|---|---|
| Univariate Analysis | ```
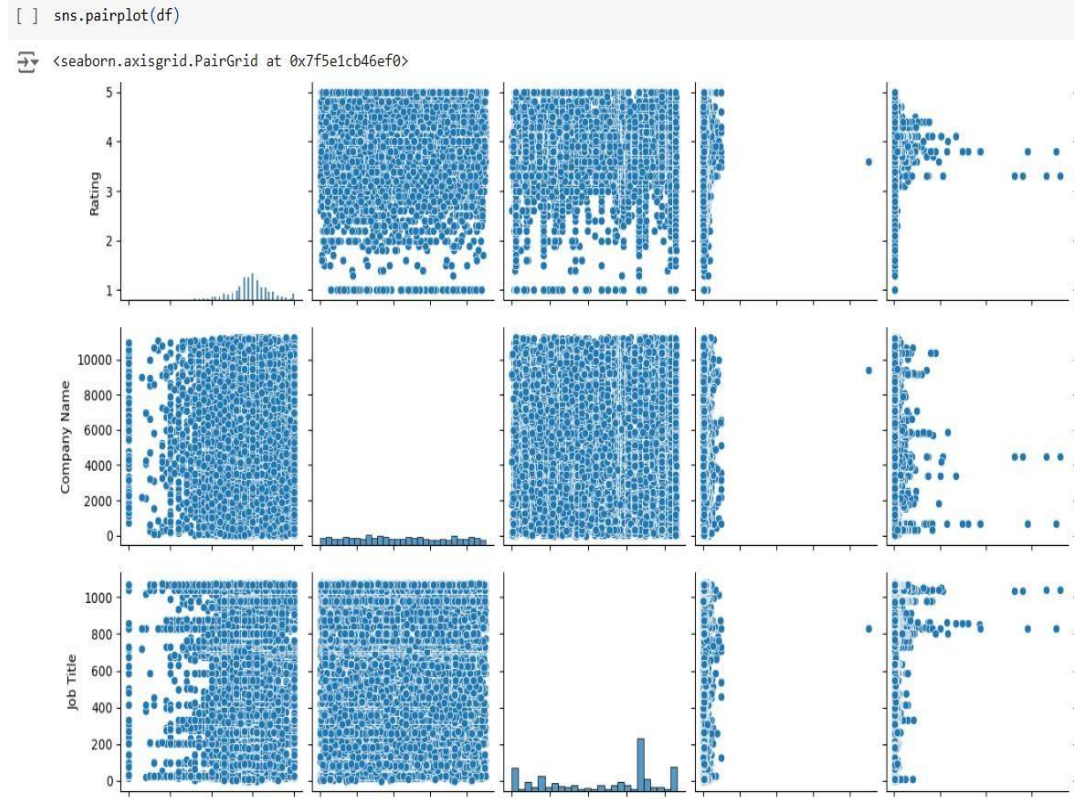[ ] sns.countplot(x='Salary',data=df)
    plt.show()
```<br> |

| | |
|---|---|
| Bivariate Analysis | ```
[ ] sns.boxplot(x='Salary',y='Job Title', data=df)
    plt.show()
```<br> |

| Multivariate Analysis | ```[ ] sns.pairplot(df)```<br>```<seaborn.axisgrid.PairGrid at 0x7f5e1cb46ef0>```<br> |
| --- | --- |
| Outliers and Anomalies | -------------- |

**Data Preprocessing**
**Code Screenshots**

| Loading Data | ```df = pd.read_csv('Salary_Dataset_with_Extra_Features.csv')```<br>```df.head()```<br> |
| --- | --- |

For the Loading Data table image:

|   | Rating | Company Name | Job Title | Salary | Salaries Reported | Location | Employment Status | Job Roles |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 0 | 3.8 | Sasken | Android Developer | 400000 | 3 | Bangalore | Full Time | Android |
| 1 | 4.5 | Advanced Millennium Technologies | Android Developer | 400000 | 3 | Bangalore | Full Time | Android |
| 2 | 4.0 | Unacademy | Android Developer | 1000000 | 3 | Bangalore | Full Time | Android |
| 3 | 3.8 | SnapBizz Cloudtech | Android Developer | 300000 | 3 | Bangalore | Full Time | Android |
| 4 | 4.4 | Appoids Tech Solutions | Android Developer | 600000 | 3 | Bangalore | Full Time | Android |

| Handling Missing Values | ```
[ ] df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 22770 entries, 0 to 22769
Data columns (total 8 columns):
 #   Column             Non-Null Count   Dtype
---  ------             --------------   -----
 0   Rating             22770 non-null   float64
 1   Company Name       22769 non-null   object
 2   Job Title          22770 non-null   object
 3   Salary             22770 non-null   int64
 4   Salaries Reported  22770 non-null   int64
 5   Location           22770 non-null   object
 6   Employment Status  22770 non-null   object
 7   Job Roles          22770 non-null   object
dtypes: float64(1), int64(2), object(5)
memory usage: 1.4+ MB
``` |
|---|---|

| Data Transformation | -- |
|---|---|
| Feature Engineering | -- |
| Save Processed Data | -- |