# Data Collection and Preprocessing Phase

| | |
|---|---|
| Date | 10 July 2024 |
| Team ID | 739779 |
| Project Title | Predictive Modelling for H1b Visa Approval Using Machine Learning |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Report**

Dataset variables will be statistically analysed 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 modelling, and forming a strong foundation for insights and predictions.

| Section | Description |
|---------|-------------|
| | **Dimension:** <br> 923 rows × 49 columns Descriptive statistics: |
| Data Overview | |

| | Unnamed: 0 | CASE_STATUS | EMPLOYER_NAME | SOC_NAME | JOB_TITLE | FULL_TIME_POSITION | PREVAILING_WAGE | YEAR | WORKSITE |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | CERTIFIED-WITHDRAWN | UNIVERSITY OF MICHIGAN | BIOCHEMISTS AND BIOPHYSICISTS | POSTDOCTORAL RESEARCH FELLOW | N | 36067.0 | 2016.0 | ANN ARBOR, MICHIGAN |
| 1 | 2 | CERTIFIED-WITHDRAWN | GOODMAN NETWORKS, INC. | CHIEF EXECUTIVES | CHIEF OPERATING OFFICER | Y | 242674.0 | 2016.0 | PLANO, TEXAS |
| 2 | 3 | CERTIFIED-WITHDRAWN | PORTS AMERICA GROUP, INC. | CHIEF EXECUTIVES | CHIEF PROCESS OFFICER | Y | 193066.0 | 2016.0 | JERSEY CITY, NEW JERSEY |
| 3 | 4 | CERTIFIED-WITHDRAWN | GATES CORPORATION, A WHOLLY-OWNED SUBSIDIARY O... | CHIEF EXECUTIVES | REGIONAL PRESIDEN, AMERICAS | Y | 220314.0 | 2016.0 | DENVER, COLORADO |
| 4 | 5 | WITHDRAWN | PEABODY INVESTMENTS CORP. | CHIEF EXECUTIVES | PRESIDENT MONGOLIA AND INDIA | Y | 157518.4 | 2016.0 | ST. LOUIS, MISSOURI |

| | |
|---|---|
| Outliers and Anomalies | ```python
df = df [df['PREVAILING_WAGE'] <= 500000]
by_emp_year = df [['EMPLOYER_NAME', 'YEAR', 'PREVAILING_WAGE']] [df['EMPLOYER_NAME'].isin(top_emp)]
# Group by the columns and reset the index to bring the grouping columns back as regular columns.
by_emp_year = by_emp_year.groupby(['EMPLOYER_NAME', 'YEAR']).mean().reset_index()
print(by_emp_year['EMPLOYER_NAME'])
``` |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data | ```python
df = pd.read_csv("h1b_kaggle.csv")
df.shape
df.head()
```

| | Unnamed: 0 | CASE_STATUS | EMPLOYER_NAME | SOC_NAME | JOB_TITLE | FULL_TIME_POSITION | PREVAILING_WAGE | YEAR | WOR |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | CERTIFIED-WITHDRAWN | UNIVERSITY OF MICHIGAN | BIOCHEMISTS AND BIOPHYSICISTS | POSTDOCTORAL RESEARCH FELLOW | N | 36067.0 | 2016.0 | ANN A MICH |
| 1 | 2 | CERTIFIED-WITHDRAWN | GOODMAN NETWORKS, INC. | CHIEF EXECUTIVES | CHIEF OPERATING OFFICER | Y | 242674.0 | 2016.0 | PLANO, |
| 2 | 3 | CERTIFIED-WITHDRAWN | PORTS AMERICA GROUP, INC. | CHIEF EXECUTIVES | CHIEF PROCESS OFFICER | Y | 193066.0 | 2016.0 | JERSEY NEW J |
| 3 | 4 | CERTIFIED-WITHDRAWN | GATES CORPORATION, A WHOLLY-OWNED SUBSIDIARY O... | CHIEF EXECUTIVES | REGIONAL PRESIDEN, AMERICAS | Y | 220314.0 | 2016.0 | DE COLO |
| 4 | 5 | WITHDRAWN | PEABODY INVESTMENTS CORP. | CHIEF EXECUTIVES | PRESIDENT MONGOLIA AND INDIA | Y | 157518.4 | 2016.0 | ST. MIS. |
|

| | |
|---|---|
| Handling Missing Data | ```python
df.isnull().sum()
```

```
Unnamed: 0            0
CASE_STATUS          0
EMPLOYER_NAME       42
SOC_NAME         17698
JOB_TITLE           26
FULL_TIME_POSITION   0
PREVAILING_WAGE      0
YEAR                 0
WORKSITE             0
lon             107089
lat             107089
dtype: int64
```

```python
df['SOC_NAME'] = df['SOC_NAME'].fillna(df['SOC_NAME'].mode()[0])
```

```python
df['CASE_STATUS'] = df['CASE_STATUS'].map({'CERTIFIED':0, 'CERTIFIED-WITHDRAWN': 1, 'DENIED': 2, 'WITHDRAWN': 3, 'PENDING QUALITY AND COMPLIANCE REVIEW
'REJECTED': 5, 'INVALIDATED': 6})
``` |

| Data Transformation |  |
| --- | --- |

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