

### Data Collection and Preprocessing Phase

Date	15 March 2024
Team ID	LTVIP2024TMID25012
Project Title	Predictive Modeling for H1B Visa Approval 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	<p>The dataset consists of 3,002,458 entries and 11 columns, including key attributes such as CASE_STATUS, EMPLOYER_NAME, SOC_NAME, JOB_TITLE, FULL_TIME_POSITION, and PREVAILING_WAGE.</p> <p><b>Basic Statistics:</b></p> <ul style="list-style-type: none"> <li>□ Total Rows: 3,002,458</li> <li>□ Total Columns: 11</li> <li>□ Unique Values in CASE_STATUS: <ul style="list-style-type: none"> <li>• CERTIFIED: 2,615,623</li> <li>• CERTIFIED-WITHDRAWN: 202,659</li> <li>• DENIED: 94,346</li> <li>• WITHDRAWN: 89,799</li> </ul> </li> <li>• Other statuses (PENDING, REJECTED, INVALIDATED): Minimal occurrences</li> </ul> <p><b>Data Types:</b></p> <ul style="list-style-type: none"> <li>○ Object: 6 columns (categorical data)</li> <li>○ Float: 4 columns (numerical data)</li> <li>○ Int64: 1 column (likely an index)</li> </ul>
Univariate Analysis	<p><b>Description:</b></p> <p>Univariate analysis was conducted to explore individual variable characteristics:</p> <ul style="list-style-type: none"> <li>• CASE_STATUS: Most common value is CERTIFIED.</li> <li>• PREVAILING_WAGE: High variability with values ranging significantly; outliers observed.</li> <li>• FULL_TIME_POSITION: Categorical distribution showing a predominance of full-time positions (Y).</li> </ul>

Bivariate Analysis	<p><b>Description:</b></p> <ul style="list-style-type: none"> <li>○ <b>Relationships:</b> The correlation between PREVAILING_WAGE and CASE_STATUS shows that higher wages correlate positively with approval.</li> <li>○ <b>Scatter Plots:</b> Generated to visualize the relationship between wage and approval status, revealing trends in the data.</li> </ul>
Multivariate Analysis	<p><b>Description:</b> Patterns involving multiple variables were analyzed:</p> <ul style="list-style-type: none"> <li>• A combination of FULL_TIME_POSITION, PREVAILING_WAGE, and SOC_NAME appears to provide a more robust prediction model for visa approval outcomes.</li> <li>• <b>Visualization:</b> PCA (Principal Component Analysis) was applied to reduce dimensionality and identify key features.</li> </ul>
Outliers and Anomalies	<p><b>Description:</b> Outliers were identified in the PREVAILING_WAGE column. For instance:</p> <ul style="list-style-type: none"> <li>○ Extreme low values (e.g., below \$20,000) and high values (e.g., above \$250,000) were capped at the 1st and 99th percentiles to maintain model integrity.</li> </ul>
<b>Data Preprocessing Code Screenshots</b>	
Loading Data	<pre>import pandas as pd  df = pd.read_csv("path/to/h1b_dataset.csv") print(df.shape) print(df.head())</pre>

Handling Missing Data	<pre># Check for missing values missing_values = df.isnull().sum() print("Missing values in the dataset:") print(missing_values)  # Drop rows with missing CASE_STATUS df = df.dropna(subset=['CASE_STATUS']).</pre>
Data Transformation	<pre># Transform CASE_STATUS to numeric values df['CASE_STATUS'] = df['CASE_STATUS'].map({     'CERTIFIED': 1,     'DENIED': 2,     'CERTIFIED-WITHDRAWN': 3,     'WITHDRAWN': 4 })</pre>
Feature Engineering	<pre># Create new feature based on SOC_NAME def classify_soc_name(row):     if 'software' in row['SOC_NAME'].lower():         return 'IT'     return 'Other'  df['SOC_CATEGORY'] = df.apply(classify_soc_name, axis=1)</pre>
Save Processed Data	<pre>df.to_csv("processed_h1b_data.csv", index=False) print("Processed data saved successfully.")</pre>