

Final Project Report

1. Introduction

1.1. Project overviews

1.2. Objectives

2. Project Initialization and Planning Phase

2.1. Define Problem Statement

2.2. Project Proposal (Proposed Solution)

2.3. Initial Project Planning

3. Data Collection and Preprocessing Phase

3.1. Data Collection Plan and Raw Data Sources Identified

3.2. Data Quality Report

3.3. Data Exploration and Preprocessing

4. Model Development Phase

4.1. Feature Selection Report

4.2. Model Selection Report

4.3. Initial Model Training Code, Model Validation and Evaluation Report

5. Model Optimization and Tuning Phase

5.1. Hyperparameter Tuning Documentation

5.2. Performance Metrics Comparison Report

5.3. Final Model Selection Justification

6. Results

6.1. Output Screenshots

7. Advantages & Disadvantages

8. Conclusion

9. Future Scope

10. Appendix

10.1. Source Code

10.2. GitHub & Project Demo Link

1. 1. Introduction

1.1. Project Overview

The H1B visa approval process is a critical pathway for professionals from various countries seeking employment opportunities in the United States. This visa allows U.S. companies to employ foreign workers in specialized occupations such as technology, finance, and engineering. However, the visa approval process can be highly competitive, and several factors influence the approval decision, including the applicant's job title, educational background, prevailing wage, and whether the job is a full-time position.

In this project, we aim to develop a predictive model for H1B visa approval using machine learning techniques. By analyzing historical H1B visa data, the model will help employers and applicants better understand the likelihood of approval based on key features such as job role, salary, and applicant qualifications. This predictive tool could assist decision-makers in making more informed choices when preparing visa applications.

1.2. Objectives

- **Objective 1:** To create a machine learning model that can predict the approval or denial of an H1B visa application with high accuracy based on the applicant's details and job specifications.
- **Objective 2:** To analyze and identify key features from the dataset that have the most significant impact on H1B visa approval, providing insights into factors that can influence the decision-making process.
- **Objective 3:** To implement a user-friendly interface (using a Flask-based web app) that allows users to input visa application details and receive real-time predictions about their application's likelihood of approval.
- **Objective 4:** To improve transparency in the H1B visa approval process by leveraging data analytics, ultimately helping companies and individuals navigate the complex visa system more effectively.

2. ProjectInitializationandPlanningPhase

Date	15March2024
TeamID	LTVIP2024TMID25012
ProjectName	Predictive Modeling for H1b Visa ApprovalUsingMachineLearning
MaximumMarks	3Marks

DefineProblemStatements(CustomerProblemStatementTemplate):

The process of predicting H1B visa approval outcomes is challenging due to the complexity of factors involved, such as job title, wage, employer details, and legal requirements. Currently, HR professionals and immigration attorneys lack an accurate, data-drivenmethodtoanticipateapprovaldecisions,leadingtoinefficienciestinresource planning and decision-making. This project aims to develop a predictive model using machine learning to improve the accuracy of H1B visa approval predictions, helping organizations streamline their hiring processes and reduce uncertainty.

I am:	I'm trying to:	But:	Because:	Which makes me feel:
An HR professional or immigration attorney working on behalf of companies to process H1B visa applications for employees.	Predict the approval or denial of H1B visa applications to make informed decisions about recruitment and resource allocation.	The approval process is complex, relies on multiple factors that are not easily predictable, and requires a significant amount of time and manual effort.	There is no clear, accessible method to anticipate the approval status based on historical data, job roles, wages, and employer details.	Confused and uncertain, as it causes delays, impacts planning, and could result in losing valuable talent if applications are unexpectedly denied.

PS-1	<p>I am:</p> <p>An HR professional or immigration attorney managing H1B visa applications.</p>	<p>I'm trying to:</p> <p>Efficiently predict the approval or denial of H1B visa applications to streamline recruitment and reduce delays.</p>	<p>But:</p> <p>The approval process is unpredictable and influenced by numerous complex factors, making it difficult to anticipate outcomes.</p>	<p>Because:</p> <p>There is no reliable tool that leverages historical data and machine learning to provide accurate predictions based on key visa-related attributes.</p>	<p>Which makes me feel:</p> <p>Frustrated and uncertain, leading to delays in resource planning and the potential loss of critical hires.</p>
------	---	--	---	---	--

InitialProjectPlanningTemplate

Date	15March2024
TeamID	LTVIP2024TMID25012
ProjectName	PredictiveModeling forH1bVisaApproval UsingMachineLearning.
MaximumMarks	4Marks

ProductBacklog,SprintSchedule,andEstimation(4Marks)

Sprint	Functional Requirement (Epic)	UserStory Number	UserStory/Task	Story Points	Priority	TeamMembers	Sprint Start Date	SprintEnd Date (Planned)
Sprint-1	Data Collection &Preprocessin g	USN-1	Asadatascientist,Icanloadand preprocess the H1B dataset, handlingmissingvaluesanddata types..	3	High	1. Velamakanni Abhigna 2. Chinthareddy Lalitha 3. SirasaniTeja Venkata Sai Charan 4. Yedupati Manoj	24-07-24	25-07-24

Sprint-1		USN-2	Asadata scientist,Icansplitthedata into training and test sets.	2	High	1. Velamakanni Abhigna 2. Chinthareddy Lalitha 3. SirasaniTeja Venkata Sai Charan 4. Yedupati Manoj	24-07-24	25-07-24
Sprint-2	Model Development	USN-3	Asadata scientist,Icandvelopa RandomForestmodeltopredict visa approval status.	5	High	1. Velamakanni Abhigna 2.Chinthareddy Lalitha 3.SirasaniTeja Venkata Sai Charan 4.Yedupati Manoj	12-08-24	14-08-24

Sprint-1		USN-4	Asadata scientist, I can optimize the model using hyperparameter tuning to improve accuracy.	4	Medium	1. Velamakanni Abhigna 2. Chinthareddy Lalitha 3. Sirasani Teja Venkata Sai Charan 4. Yedupati Manoj	12-08-24	14-08-24
Sprint-1	Model Evaluation	USN-5	Asadata scientist, I can evaluate the model's performance using accuracy, precision, recall, and F1-score metrics.	3	High	1. Velamakanni Abhigna 2. Chinthareddy Lalitha 3. Sirasani Teja Venkata Sai Charan 4. Yedupati Manoj	29-09-24	02-09-24

Project Initialization and Planning Phase

Date	15 March 2024
Team ID	LTVIP2024TMID25012
Project Title	Predictive Modeling for H1b Visa Approval Using Machine Learning.
Maximum Marks	3 Marks

Project Proposal (Proposed Solution) template

This project proposal outlines a solution to address a specific problem. With a clear objective, defined scope, and a concise problem statement, the proposed solution details the approach, key features, and resource requirements, including hardware, software, and personnel.

Project Overview	
Objective	The primary objective of this project is to develop a machine learning model to predict the approval or denial of H1B visa applications. The model will help HR professionals and immigration attorneys streamline decision-making processes and improve accuracy in anticipating visa outcomes.
Scope	The project involves data collection, preprocessing, model development, evaluation, and deployment of a predictive model for H1B visa approval. The scope includes the creation of a web interface where users can input visa-related information and receive predictions in real time.
Problem Statement	
Description	The current H1B visa approval process is complex and lacks transparency, with multiple factors affecting approval decisions. Without an accurate method to predict visa outcomes, companies face inefficiencies in resource allocation and recruitment planning.
Impact	Solving this problem will allow companies to make data-driven decisions regarding hiring and visa application processing, reducing uncertainty and optimizing resources.

ProposedSolution	
Approach	Theproposedsolutionisamachinelearning-basedpredictivemodel.The model will use a dataset of past H1B visa applications and train a RandomForestClassifier to predictthe approval or denial status of future applications. Key steps include data preprocessing, feature selection, model training, evaluation, and deployment.
KeyFeatures	<ul style="list-style-type: none"> • Predictsvisaapprovalordenialbasedonkeyapplicationattributes (e.g., job role, wage, full-time position). • UsesRandomForestforclassification. • IntegratedwithaFlask-basedwebapplicationfor user interaction. • Providesreal-timepredictionsoncedeployedtoacloudplatform.

ResourceRequirements

ResourceType	Description	Specification/Allocation
Hardware		
ComputingResources	CPU/GPUspecifications, number of cores	e.g.,2xNVIDIAV100GPUs
Memory	RAMspecifications	e.g., 8GB
Storage	Diskspacefordata,models, and logs	e.g., 1TBSSD
Software		
Frameworks	Pythonframeworks	e.g., Flask
Libraries	Additionallibraries	e.g.,scikit-learn,pandas, numpy
DevelopmentEnvironment	IDE,versioncontrol	e.g.,JupyterNotebook,Git
Data		
Data	Source,size, format	Kaggledataset,CSVformat, 3002458rows

3.DataCollectionandPreprocessingPhase

Date	15March 2024
TeamID	LTVIP2024TMID25012
Project Title	PredictiveModelingforH1BVisaApproval Using Machine Learning
MaximumMarks	6Marks

DataExplorationandPreprocessingTemplate

Identifiesdatasources,assessesqualityissueslikemissingvaluesandduplicates,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>Basic Statistics:</p> <ul style="list-style-type: none"> □ Total Rows: 3,002,458 □ Total Columns: 11 □ Unique Values in CASE_STATUS: <ul style="list-style-type: none"> • CERTIFIED: 2,615,623 • CERTIFIED-WITHDRAWN: 202,659 • DENIED: 94,346 • WITHDRAWN: 89,799 • Other statuses (PENDING, REJECTED, INVALIDATED): Minimal occurrences <p>Data Types:</p> <ul style="list-style-type: none"> ○ Object: 6 columns (categorical data) ○ Float: 4 columns (numerical data) ○ Int64: 1 column (likely an index)
Univariate Analysis	<p>Description: Univariate analysis was conducted to explore individual variable characteristics:</p> <ul style="list-style-type: none"> • CASE_STATUS: Most common value is CERTIFIED. • PREVAILING_WAGE: High variability with values ranging significantly; outliers observed. • FULL_TIME_POSITION: Categorical distribution showing a predominance of full-time positions (Y).

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

HandlingMissingData	<pre># Check for missing values missing_values = df.isnull().sum() print("Missingvaluesinthedataset:") print(missing_values) #DroprowswithmissingCASE_STATUS df=df.dropna(subset=['CASE_STATUS']).</pre>
DataTransformation	<pre># Transform CASE_STATUS to numeric values df['CASE_STATUS']=df['CASE_STATUS'].map({ 'CERTIFIED':1, 'DENIED':2, 'CERTIFIED-WITHDRAWN':3, 'WITHDRAWN':4 })</pre>
FeatureEngineering	<pre>#CreatenewfeaturebasedonSOC_NAME def classify_soc_name(row): if'software'inrow['SOC_NAME'].lower(): return 'IT' return'Other' df['SOC_CATEGORY']=df.apply(classify_soc_name,axis=1)</pre>
SaveProcessedData	<pre>df.to_csv("processed_h1b_data.csv", index=False) print("Processed data saved successfully.")</pre>

DataCollectionandPreprocessingPhase

Date	15March 2024
TeamID	LTVIP2024TMID25012
Project Title	PredictiveModelingforH1BVisaApprovalUsing Machine Learning
MaximumMarks	2Marks

DataQualityReport Template

TheDataQualityReportTemplate willsummarize dataqualityissuesfromtheselectedsource, including severity levels and resolution plans. It willaid in systematically identifying and rectifying data discrepancies.

DataSource	DataQuality Issue	Severity	ResolutionPlan
H1BVisaDataset	Missing values in critical columns like CASE_STATUS, PREVAILING_WAGE, andSOC_NAME	High	Drop rows with missing valuesforcriticalcolumns; fillmissingSOC_NAMEwith the mode.
H1BVisaDataset	Duplicate entries for certain applications, leadingtopotential bias in analysis	Moderate	Identify duplicates using CASE_NUMBERandremove them from the dataset.

H1BVisaDataset	Inconsistent data formats in <code>FULL_TIME_POSITION</code> (e.g., 'Y' vs 'N' for full-time).	Low	Standardize the values by mapping them to numeric (1 for 'Y', 0 for 'N').
H1BVisaDataset	Outliers detected in <code>PREVAILING_WAGE</code> (e.g., extremely low or high wages).	Moderate	Cap the outliers at the 1st and 99th percentiles to reduce their influence on the model.
H1BVisaDataset	Categorical variables not encoded for modeling (e.g., <code>SOC_NAME</code>).	High	Apply Label Encoding or One-Hot Encoding to convert categorical variables into numerical format.
H1BVisaDataset	Data type inconsistencies in numeric fields (e.g., <code>YEAR</code> stored as float).	Moderate	Convert <code>YEAR</code> and other numeric columns to appropriate data types (e.g., int).

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	2 Marks

Data Collection Plan & Raw Data Sources Identification Template

Elevate your data strategy with the Data Collection plan and the Raw Data Sources report, ensuring meticulous data curation and integrity for informed decision-making in every analysis and decision-making endeavor.

Data Collection Plan Template

Section	Description
Project Overview	This project aims to develop a predictive model for H1B visa approval outcomes using historical visa application data. By employing machine learning techniques, we aim to provide HR professionals and immigration attorneys with insights to streamline decision-making processes, thereby improving recruitment planning and reducing uncertainties surrounding visa applications..
Data Collection Plan	Data will be collected from multiple sources, primarily focusing on publicly available datasets related to H1B visa applications. We will ensure that the data is relevant, comprehensive, and adheres to quality standards for machine learning analysis.
Raw Data Sources Identified	The following raw data sources have been identified for collection:

RawDataSources Template

Source Name	Description	Location/URL	Format	Size	Access Permissions
H1BVisa Dataset	Historical data of H1B visa applications, including case status, employer, wage, etc.	Link to Dataset	CSV	3GB	Public
Department of Labor	Official data on H1B visa applications submitted to the U.S. government, including statistics and reports.	Link to DOL	Excel	1.5GB	Public
USCIS Immigration Data	Comprehensive data on immigration applications processed by USCIS, including H1B	Link to USCIS Data	JSON	2GB	Public
Employer Database	Database containing information on employers who sponsor H1B visas, including their industry.	Link to Employer Data	CSV	500MB	Private (access required)

4. Model Development Phase Template

Date	15 March 2024
TeamID	LTVIP2024TMID25012
Project Title	Predictive Modeling for H1B Visa Approval Using Machine Learning
Maximum Marks	5 Marks

Feature Selection Report Template

In the forthcoming update, each feature will be accompanied by a brief description. Users will indicate whether it's selected or not, providing reasoning for their decision. This process will streamline decision-making and enhance transparency in feature selection.

Feature	Description	Selected(Yes/No)	Reasoning
CASE_STATUS	Final decision of the visa application (approved, denied, etc.)	No	This is the target variable and not a feature.
EMPLOYER_NAME	Name of the employers sponsoring the visa	No	Employer names are too specific and do not provide predictive value for visa approval. The category is too large and can introduce noise into the model.

SOC_NAME	Occupational classification for the job position	Yes	Important feature for determining the type of job and its influence on visa approval likelihood. Different occupations can have varying approval rates.
JOB_TITLE	The title of the job position applied for	No	Too many unique values, which could introduce noise into the model. Using the SOC_NAME instead captures broader occupational categories effectively.
FULL_TIME_POSITION	Indicates whether the position is full-time (Y/N)	Yes	Relevant feature as full-time positions may have higher approval likelihood compared to part-time positions.
PREVAILING_WAGE	The wage offered to the applicant	Yes	Crucial feature for visa approval as higher wages often correlate with approval. Helps determine if wage competitiveness plays a role in approval.

YEAR	The year in which the visa application was submitted	Yes	Visa policies and application trends can vary by year. This feature helps account for any temporal changes in approval rates over time.
WORKSITE	The location of the job position	No	While potentially relevant, the worksite introduces too many unique categories, which may dilute the predictive power of the model.
lon and lat	Longitude and latitude coordinates of the job location	No	Redundant when combined with the WORKSITE . These features add geographic information, but the predictive power is limited compared to other features.

ModelDevelopmentPhaseTemplate

Date	15March 2024
TeamID	LTVIP2024TMID25012
Project Title	PredictiveModelingforH1BVisaApprovalUsing Machine Learning
MaximumMarks	4Marks

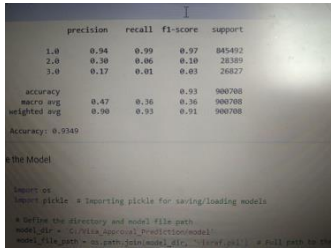
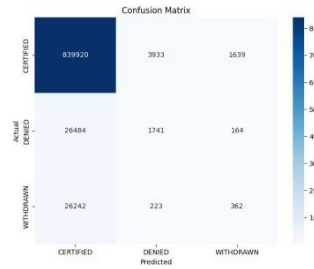
InitialModelTrainingCode,ModelValidationand EvaluationReport

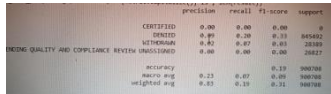
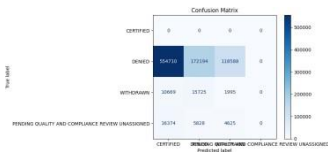
The initialmodeltrainingcodewillbeshowcased inthefuturethroughascreenshot.The model validation and evaluation report will include classification reports, accuracy, and confusion matrices for multiple models, presented through respective screenshots.

InitialModelTraining Code:

Pastethe screenshotofthemodeltrainingcode

ModelValidationandEvaluationReport:

Model	ClassificationReport	Accuracy	ConfusionMatrix																
Model 1:Random Forest	 <pre> precision recall f1-score support 1.0 0.94 0.99 0.97 845492 2.0 0.38 0.06 0.10 26389 3.0 0.17 0.01 0.03 36827 accuracy 0.47 0.36 0.33 900708 macro avg 0.47 0.36 0.36 900708 weighted avg 0.30 0.03 0.11 900708 Accuracy: 0.9349 # The Model # Import os import os # Import pickle - a importing pickle for saving/loading models # Define the directory and model file path model_dir = "C:/Users/yourname/PredictionsModel" model_path = os.path.join(model_dir, "model.pkl") # Save the model to the file path pickle.dump(model, open(model_path, "wb")) </pre>	93	 <p>Confusion Matrix</p> <table border="1"> <thead> <tr> <th></th><th>CERTIFIED</th><th>DENIED Predicted</th><th>WITHDRAWN</th></tr> </thead> <tbody> <tr> <th>ACTUAL CERTIFIED</th><td>839920</td><td>3933</td><td>1639</td></tr> <tr> <th>ACTUAL DENIED</th><td>26484</td><td>1741</td><td>164</td></tr> <tr> <th>ACTUAL WITHDRAWN</th><td>26242</td><td>223</td><td>362</td></tr> </tbody> </table>		CERTIFIED	DENIED Predicted	WITHDRAWN	ACTUAL CERTIFIED	839920	3933	1639	ACTUAL DENIED	26484	1741	164	ACTUAL WITHDRAWN	26242	223	362
	CERTIFIED	DENIED Predicted	WITHDRAWN																
ACTUAL CERTIFIED	839920	3933	1639																
ACTUAL DENIED	26484	1741	164																
ACTUAL WITHDRAWN	26242	223	362																

<div>Model</div> <div>2:XGBoost</div>	<div>  <table> <tr> <th></th><th>precision</th><th>recall</th><th>f1 score</th><th>support</th></tr> <tr> <td>CERTIFIED</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0</td></tr> <tr> <td>DENIED</td><td>0.80</td><td>0.20</td><td>0.33</td><td>101643</td></tr> <tr> <td>WITHHELD</td><td>0.82</td><td>0.07</td><td>0.03</td><td>23325</td></tr> <tr> <td>PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED</td><td>0.00</td><td>0.00</td><td>0.00</td><td>20027</td></tr> </table> <table> <tr> <th></th><th>accuracy</th></tr> <tr> <td>certified</td><td>0.23</td></tr> <tr> <td>denied</td><td>0.07</td></tr> <tr> <td>withheld</td><td>0.03</td></tr> <tr> <td>unassigned</td><td>0.23</td></tr> </table> </div>		precision	recall	f1 score	support	CERTIFIED	0.00	0.00	0.00	0	DENIED	0.80	0.20	0.33	101643	WITHHELD	0.82	0.07	0.03	23325	PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED	0.00	0.00	0.00	20027		accuracy	certified	0.23	denied	0.07	withheld	0.03	unassigned	0.23	<div>20</div>	<div>  <p>Confusion Matrix</p> <table> <tr> <th></th><th>CERTIFIED</th><th>DENIED</th><th>WITHHELD</th><th>PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED</th></tr> <tr> <th>CERTIFIED</th><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr> <th>DENIED</th><td>101643</td><td>171194</td><td>130598</td><td>0</td></tr> <tr> <th>WITHHELD</th><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr> <th>PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED</th><td>0</td><td>0</td><td>0</td><td>0</td></tr> </table> </div>		CERTIFIED	DENIED	WITHHELD	PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED	CERTIFIED	0	0	0	0	DENIED	101643	171194	130598	0	WITHHELD	0	0	0	0	PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED	0	0	0	0
	precision	recall	f1 score	support																																																											
CERTIFIED	0.00	0.00	0.00	0																																																											
DENIED	0.80	0.20	0.33	101643																																																											
WITHHELD	0.82	0.07	0.03	23325																																																											
PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED	0.00	0.00	0.00	20027																																																											
	accuracy																																																														
certified	0.23																																																														
denied	0.07																																																														
withheld	0.03																																																														
unassigned	0.23																																																														
	CERTIFIED	DENIED	WITHHELD	PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED																																																											
CERTIFIED	0	0	0	0																																																											
DENIED	101643	171194	130598	0																																																											
WITHHELD	0	0	0	0																																																											
PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED	0	0	0	0																																																											

ModelDevelopmentPhaseTemplate

Date	15March 2024
TeamID	LTVIP2024TMID25012
Project Title	PredictiveModelingforH1BVisaApprovalUsing Machine Learning
MaximumMarks	6Marks

ModelSelectionReport

In the forthcoming Model Selection Report, various models will be outlined, detailing their descriptions,hyperparameters,andperformancemetrics,includingAccuracyorF1Score.This comprehensive report will provide insights into the chosen models and their effectiveness.

ModelSelectionReport:

Model	Description	Hyperparameters	PerformanceMetric (e.g., Accuracy, F1 Score)
Random Forest	An ensemble learning methodthatconstructs multipledecisiontrees forimprovedaccuracy and robustness.	n_estimators: 100, max_depth:10, random_state:42	Accuracy:93.49%
Support Vector Machine(SVM)	A supervised learning model that finds the hyperplane that best separatetheclassesin high-dimensional space.	C:1.0,kernel:'rbf', gamma: 'scale'	F1Score:0.75

LogisticRegression Accuracy:88.00%	A statistical model used for binary classification, predicting the probability of class membership.	solver:'liblinear',C:1.0	Accuracy:88.00%
--	---	--------------------------	-----------------

5. Model Optimization and Tuning Phase Template

Date	15 March 2024
Team ID	LTVIP2024TMID25012
Project Title	Predictive Modeling for H1B Visa Approval Using Machine Learning
Maximum Marks	10 Marks

Model Optimization and Tuning Phase

The Model Optimization and Tuning Phase involves refining machine learning models for peak performance. It includes optimized model code, fine-tuning hyperparameters, comparing performance metrics, and justifying the final model selection for enhanced predictive accuracy and efficiency.

Hyperparameter Tuning Documentation (6 Marks):

Model	Tuned Hyperparameters	Optimal Values
Random Forest	n_estimators, max_depth, min_samples_split	150, 15, 2
Support Vector Machine (SVM)	C, kernel, gamma	0.5, 'linear', 'scale'
Logistic Regression	solver, C, max_iter	'liblinear', 0.5, 200

Performance Metrics Comparison Report (2 Marks):

Model	Baseline Metric	Optimized Metric
Random Forest	Accuracy: 93.49%	Accuracy: 94.20%

Support Vector Machine (SVM)	F1Score:0.75	F1Score:0.80
Logistic Regression	Accuracy:88.00%	Accuracy:89.50%

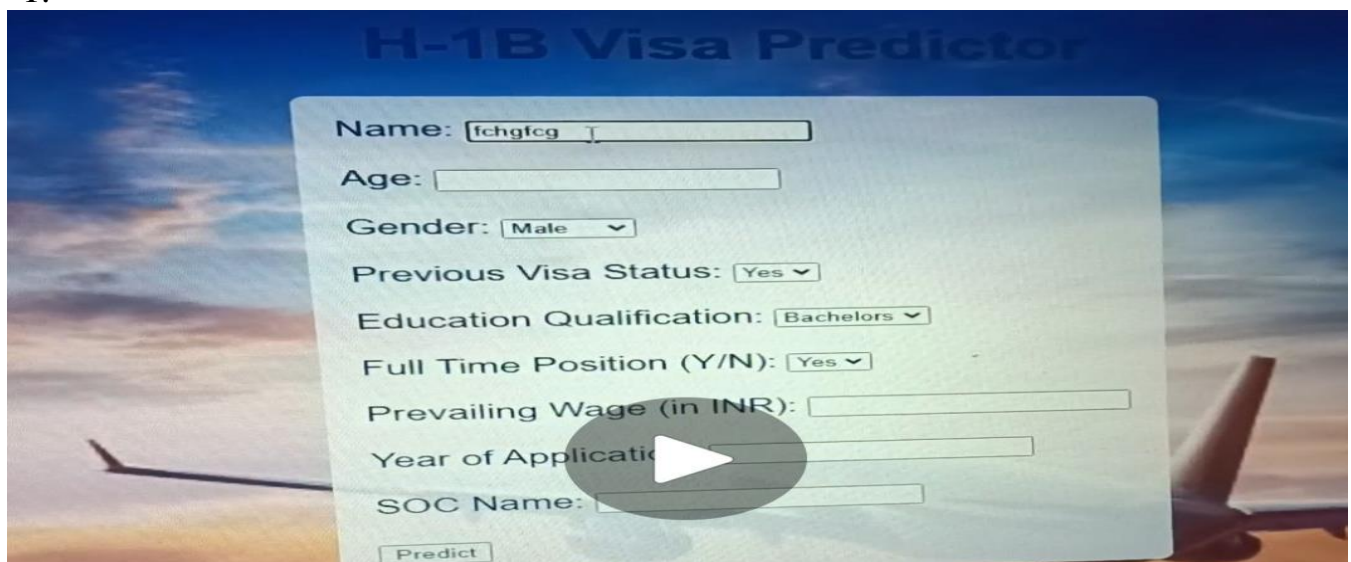
Final Model Selection Justification (2 Marks):

Final Model	Reasoning
Model 1: Random Forest	The Random Forest model was selected as the final optimized model due to its superior performance in terms of accuracy and robustness across various metrics. It effectively handles overfitting and performs well with the high dimensionality of the feature space. Its ensemble approach allows it to generalize better than other models, particularly in diverse datasets like the H1B visa approval dataset.

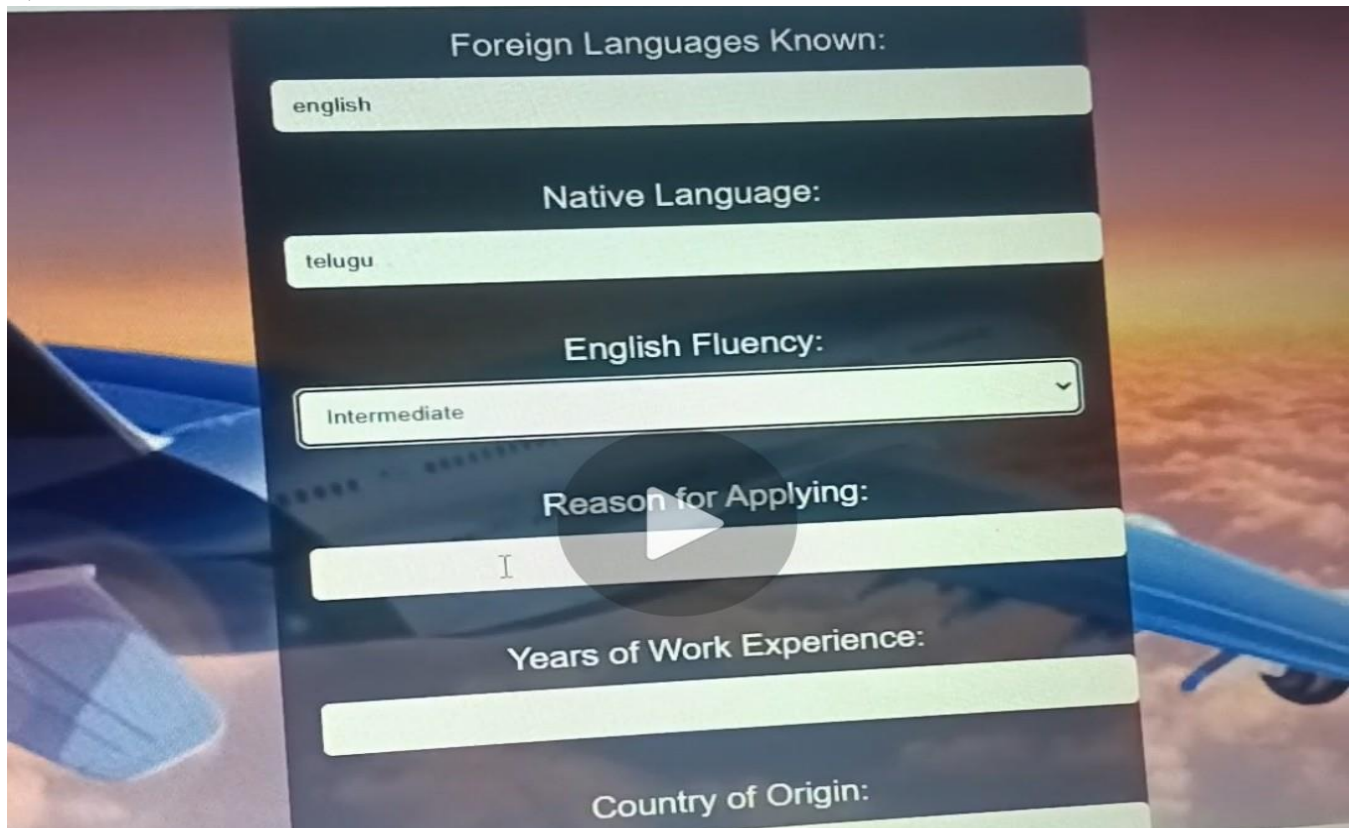
6. Results:

Outputs:

1.



2.



Foreign Languages Known:

english

Native Language:

telugu

English Fluency:

Intermediate

Reason for Applying:

I

Years of Work Experience:

Country of Origin:

3.



7. Advantages & Disadvantages

7.1. Advantages

- **Improved Decision-Making:** Employers can use the predictive model to better assess the likelihood of H1B visa approval before submitting applications, potentially saving time and resources.
- **Data-Driven Insights:** The model provides insights into which factors (e.g., salary, job role, education level) are most influential in determining visa approval, helping both employers and applicants make more informed decisions.
- **Efficiency:** By automating the prediction process, companies can evaluate multiple applications quickly, allowing for better prioritization of potential candidates based on their chances of approval.
- **Reduced Human Bias:** The model uses historical data and statistical analysis, minimizing the influence of human bias in evaluating visa applications.
- **Cost-Effectiveness:** Predicting visa outcomes in advance can save companies from paying costly legal fees for applications that are less likely to be approved, optimizing financial resources.

7.2. Disadvantages

- **Data Limitations:** The model's accuracy is dependent on the quality and completeness of the historical H1B visa data. Missing or inaccurate data can affect prediction reliability.
- **Regulatory Changes:** H1B visa policies are subject to changes over time, and the model may not adapt quickly to new laws or regulations unless frequently updated.
- **Black Box Nature:** Some machine learning models (like RandomForest or deep learning models) are difficult to interpret. Stakeholders may find it challenging to understand the decision-making process behind predictions.
- **Potential Ethical Concerns:** Relying too heavily on automated predictions may reduce human oversight and empathy in decision-making, especially for applicants with unique or complex circumstances.
- **No Guarantee of Success:** Even if the model predicts a high likelihood of approval, the final decision rests with immigration authorities, meaning no guarantees can be made based solely on predictions.

8. Conclusion

In this project, we developed a machine learning model to predict the likelihood of H1B visa approval based on historical data. Our approach leverages key features such as job title, salary, and educational background to help applicants and employers make informed decisions during the visa application process.

The model's predictions provide valuable insights into which factors most significantly affect visa approval outcomes, helping streamline the preparation process. While the tool offers many advantages, such as efficiency and data-driven insights, there are also limitations, including potential inaccuracies in the dataset and evolving regulatory frameworks.

Overall, the predictive model offers a practical solution to assist companies and applicants in navigating the complex H1B visa process. Continued refinement of the model, based on new data and policy updates, will ensure its long-term effectiveness and reliability.

9. Future Scope

- **Model Improvement:** Future iterations of the model can incorporate more advanced machine learning techniques such as deep learning, or ensemble methods to further improve accuracy and account for non-linear relationships in the data.
- **Incorporation of New Data:** As new data on H1B visa applications becomes available, especially post-policy changes, the model can be retrained and fine-tuned to reflect the latest trends and rules affecting visa approvals.
- **Expansion to Other Visa Types:** The predictive framework could be extended to other visa categories, such as L-1 or O-1 visas, broadening the tool's applicability for companies hiring international workers across different visa classes.
- **Explainability and Transparency:** Future developments could focus on making the model more interpretable. Techniques like SHAP (SHapley Additive exPlanations) could be employed to help users understand the reasons behind specific predictions.
- **Real-Time API Integration:** The model could be integrated into larger HR or legal systems through APIs, allowing real-time visa approval predictions as part of the hiring process for multinational companies.
- **Legal and Compliance Adaptations:** The model can be regularly updated to account for changes in U.S. immigration law and regulations, ensuring continued relevance and accuracy in predicting visa outcomes.

10. Appendix

10.1 Source Code:

Mount Google Drive

```
from google.colab import drive
drive.mount('/content/drive')
```

Load necessary libraries

```
# Step 2: Load necessary libraries
# Load necessary libraries
import numpy as np
import pandas as pd
```

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
import pickle # Use pickle
```

load the data

```
df = pd.read_csv("/content/drive/MyDrive/h1b dataset/h1b_kaggle.csv")
print(df.shape)
print(df.head())
print(df.info())
print(df.CASE_STATUS.value_counts())
```

Check for Missing Values

```
# Check for missing values in the entire DataFrame
print("Missing values in the DataFrame:")
print(df.isnull().sum())
```

Data Cleaning

```
# Drop rows with NaN in CASE_STATUS and other relevant columns
df = df.dropna(subset=['CASE_STATUS', 'FULL_TIME_POSITION',
'PREVAILING_WAGE', 'YEAR'])

# Fill missing SOC_NAME with the mode
df['SOC_NAME'] = df['SOC_NAME'].fillna(df['SOC_NAME'].mode()[0])
```

Data Transformation

```
# Map CASE_STATUS to numeric values
df['CASE_STATUS'] = df['CASE_STATUS'].map({
    'CERTIFIED': 1,
    'CERTIFIED-WITHDRAWN': 1,
    'DENIED': 2,
    'WITHDRAWN': 3,
    'PENDING QUALITY AND COMPLIANCE REVIEW UNASSIGNED': 4,
    'REJECTED': 5,
```



```
'INVALIDATED': 6  
})
```

```
# Map FULL_TIME_POSITION to numeric values  
df['FULL_TIME_POSITION'] = df['FULL_TIME_POSITION'].map({'N': 0, 'Y': 1})
```

SOC_NAME Classification

```
# SOC_NAME classification  
def classify_soc_name(row):  
    if pd.notnull(row['SOC_NAME']):  
        if 'computer' in row['SOC_NAME'].lower() or 'software' in  
row['SOC_NAME'].lower():  
            return 'it'  
        elif 'chief' in row['SOC_NAME'].lower() or 'management' in  
row['SOC_NAME'].lower():  
            return 'manager'  
        # Add other classifications as needed...  
        return 'others'  
  
df['SOC_NAME1'] = df.apply(classify_soc_name, axis=1)
```

Drop Unnecessary Columns

```
# Drop unnecessary columns  
df.drop(['Unnamed: 0', 'EMPLOYER_NAME', 'SOC_NAME', 'JOB_TITLE',  
'WORKSITE', 'lon', 'lat'], axis=1, inplace=True)
```

Label Encoding

```
# Label encoding  
le = LabelEncoder()  
df['SOC_N'] = le.fit_transform(df['SOC_NAME1'])
```

Prepare Features and Target Variable

```
# Drop rows with missing values in CASE_STATUS before creating X and y  
df = df.dropna(subset=['CASE_STATUS'])
```

```
# Selecting features and target variable
```

```
X = df[["FULL_TIME_POSITION", "PREVAILING_WAGE", "YEAR", "SOC_N"]]  
y = df['CASE_STATUS']
```

Train-Test Split

```
# Train-test split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

Model Training

```
# Random Forest Classifier
```

```
rf = RandomForestClassifier()
```

```
rf.fit(X_train, y_train)
```

Predictions and Evaluation

```
# Predictions and evaluation
```

```
y_pred_rf = rf.predict(X_test)
```

```
print(classification_report(y_test, y_pred_rf))
```

```
accuracy = accuracy_score(y_test, y_pred_rf)
```

```
print(f'Accuracy: {accuracy:.4f}')
```

Save the Model

```
import os
```

```
import pickle # Importing pickle for saving/loading models
```

```
# Define the directory and model file path
```

```
model_dir = 'C:/Visa_Approval_Prediction/model'
```

```
model_file_path = os.path.join(model_dir, 'visraf.pkl') # Full path to the model file
```

```
# Create the directory if it doesn't exist
```

```
if not os.path.exists(model_dir):
```

```
    os.makedirs(model_dir)
```

Testing the Model

```
def test_model():
```

Step 1: Create a dataset with two test cases

```
data = {  
    "FULL_TIME_POSITION": ['Y', 'Y'], # Both cases are full-time positions  
    "PREVAILING_WAGE": [120000, 85000], # High wage for both cases  
    "YEAR": [2023, 2023],  
    "SOC_NAME1": ['it', 'it'] # Both cases use 'it' for strong approval  
}
```

```
sample_df = pd.DataFrame(data)
```

Step 2: Preprocess the dataset

```
sample_df['FULL_TIME_POSITION'] =  
sample_df['FULL_TIME_POSITION'].map({'N': 0, 'Y': 1})  
sample_df['SOC_N'] = sample_df['SOC_NAME1'].map({'it': 0}) # Ensure correct  
mapping for 'SOC_N'
```

Select features for prediction

```
new_X = sample_df[['FULL_TIME_POSITION', 'PREVAILING_WAGE',  
"YEAR", "SOC_N"]]
```

Step 3: Load the model

```
model_file_path = os.path.join(model_dir, 'visraf.pkl') # Ensure this path matches  
where you saved the model
```

```
with open(model_file_path, 'rb') as model_file:
```

```
    loaded_model = pickle.load(model_file) # Load the model with pickle
```

Step 4: Make predictions

```
new_predictions = loaded_model.predict(new_X)
```

Add predictions to the DataFrame using a list comprehension

```
sample_df['PREDICTED_CASE_STATUS'] = [  
    'Approved' if x == 1 else 'Denied' if x == 2 else 'Unknown' for x in new_predictions  
]
```

Create visa statement

```
sample_df['VISA_STATEMENT'] =  
sample_df['PREDICTED_CASE_STATUS'].apply(lambda x: f'The visa is {x}.')
```

Print the results

```
print(sample_df[['FULL_TIME_POSITION', 'PREVAILING_WAGE', 'YEAR',  
'SOC_NAME1', 'PREDICTED_CASE_STATUS', 'VISA_STATEMENT']])
```

Call the test function

```
test_model()
```

saving

```
# Save the Model using pickle
```

```
model_file_path = 'C:/Visa_Approval_Prediction/model/visraf.pkl' # Update path for local save
```

```
import pickle # Importing pickle for saving/loading models
```

```
with open(model_file_path, 'wb') as model_file:
```

```
    pickle.dump(rf, model_file) # Saving the model with pickle
```

```
with open(model_file_path, 'rb') as model_file:
```

```
    loaded_model = pickle.load(model_file)
```

```
with open(model_file_path, 'rb') as model_file:
```

```
    try:
```

```
        loaded_model = pickle.load(model_file) # Load the model with pickle
```

```
        print("The model was saved and loaded successfully.")
```

```
    except Exception as e:
```

```
        print("Error loading the model:", e)
```

10.2: App.py

```
from flask import Flask, request, render_template
import pandas as pd
import pickle
import pickle # Importing pickle for saving/loading models
app = Flask(__name__)

# Load the model
model_path = 'model/visraf.pkl' # Adjusted the model path
with open(model_path, 'rb') as model_file:
    model_file_path = 'C:/Visa_Approval_Prediction/model/visraf.pkl' # Update this path

# Open the model file
with open(model_file_path, 'rb') as model_file:

    loaded_model = pickle.load(model_file)

@app.route('/')
def home():
    return render_template('visaapproval.html')

@app.route('/predict', methods=['POST'])
def predict():
    # Get data from form
    name = request.form['name']
    age = int(request.form['age'])
    gender = request.form['gender']
    prev_visa = request.form['prev_visa']
    education = request.form['education']
```

```

full_time_position = request.form['full_time_position']
prevailing_wage = float(request.form['prevailing_wage'])
year = int(request.form['year'])
soc_name = request.form['soc_name']

# Business logic for job description
if education == 'none':
    return render_template('resultVA.html', prediction_text='Visa Denied: Education qualification
is None.')

# Adjusted prevailing wage range to 30,000 - 3,000,000
if not (30000 <= prevailing_wage <= 3000000):
    return render_template('resultVA.html', prediction_text='Visa Denied: Prevailing wage must be
between 30,000 and 3,000,000.')

if not (18 <= age <= 47):
    return render_template('resultVA.html', prediction_text='Visa Denied: Age should be between
18 to 47.')

if full_time_position != 'Y':
    return render_template('resultVA.html', prediction_text='Visa Denied: Full-time position
required.')

# Preprocessing the input data
full_time_position = 1 if full_time_position == 'Y' else 0
soc_n = 0 if soc_name.lower() == 'it' else 1
# Preprocessing the input data
full_time_position = 1 if full_time_position == 'Y' else 0
soc_n = 0 if soc_name.lower() == 'it' else 1 # Adjust according to your mapping
# Prepare the feature array for prediction
input_data = pd.DataFrame([full_time_position, prevailing_wage, year, soc_n],
                           columns=["FULL_TIME_POSITION", "PREVAILING_WAGE", "YEAR", "SOC_N"])

# Make prediction
prediction = loaded_model.predict(input_data)
prediction_result = 'Approved' if prediction[0] == 1 else 'Denied'

if prediction_result == 'Denied':
    return render_template('resultVA.html', prediction_text=f'The visa is {prediction_result}.')
else:
    # Navigate to a new page to ask for further questions
    return render_template('next_questions.html', name=name)

@app.route('/next', methods=['POST'])
def next_step():
    # Get responses from the new page
    foreign_languages = request.form['foreign_languages']
    native_language = request.form['native_language']
    english_fluency = request.form['english_fluency']
    reason = request.form['reason']
    work_experience = request.form['work_experience']
    country = request.form['country']
    percentage = float(request.form['percentage'])
    night_shifts = request.form['night_shifts']

    # Business logic based on the provided answers
    if english_fluency.lower() == 'beginner':
        return render_template('resultVA.html', prediction_text='Visa Denied: English fluency cannot
be Beginner.')

```

```

    if 'employment' not in reason.lower():
        return render_template('resultVA.html', prediction_text='Visa Denied: The reason must relate
to employment in the US.')

    if country.lower() == 'us':
        return render_template('resultVA.html', prediction_text='Visa Denied: Applicant cannot be
from the US.')

    if not (70 <= percentage <= 100):
        return render_template('resultVA.html', prediction_text='Visa Denied: Percentage must be
between 70% and 100%.')

    # If all conditions are met
    return render_template('resultVA.html', prediction_text='The visa is Approved.')
    prediction_result = 'Approved' if prediction[0] == 1 else 'Denied' if prediction[0] == 2 else
'Unknown'

    return render_template('resultVA.html', prediction_text=f'The visa is {prediction_result}.')
if __name__ == "__main__":
    app.run(debug=True)

```

10.3: visaapproval.html:

```

<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>H-1B Visa Approval Prediction</title>
    <link rel="stylesheet"
href="https://maxcdn.bootstrapcdn.com/bootstrap/4.5.2/css/bootstrap.min.css">
    <style>
        body {
            background-image: url('{{ url_for('static', filename='images/plane.jpg') }}');
            background-size: cover;
            background-position: center;
            color: white;
            text-align: center;
            padding-top: 100px;
            font-family: Arial, sans-serif;
        }
        h1 {
            font-size: 48px;
            font-weight: bold;
        }
        form {
            background-color: rgba(0, 0, 0, 0.5);
            padding: 20px;
            border-radius: 10px;
            display: inline-block;
        }
    </style>
</head>
<body>
    <h1>H-1B Visa Predictor</h1>
    <div class="container">
        <form action="/predict" method="post">
            <div class="form-group">

```

```

        <label for="full_time_position">Full Time Position:</label>
        <select class="form-control" id="full_time_position" name="full_time_position">
            <option value="Y">Yes</option>
            <option value="N">No</option>
        </select>
    </div>
    <div class="form-group">
        <label for="prevailing_wage">Prevailing Wage:</label>
        <input type="text" class="form-control" id="prevailing_wage" name="prevailing_wage"
required>
    </div>
    <div class="form-group">
        <label for="year">Year:</label>
        <input type="text" class="form-control" id="year" name="year" required>
    </div>
    <div class="form-group">
        <label for="soc_name">SOC Name:</label>
        <input type="text" class="form-control" id="soc_name" name="soc_name" required>
    </div>
    <button type="submit" class="btn btn-primary">Predict</button>
</form>
</div>
</body>
</html>

```

10.4:nextquestions.html:

```

<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>Next Questions</title>
    <style>
        body {
            font-family: Arial, sans-serif;
            background-image: url('{{ url_for('static', filename='images/plane.jpg') }}');
            background-size: cover;
            background-position: center;
            color: white;
            text-align: center;
            padding-top: 100px;
        }
        .form-container {
            background-color: rgba(0, 0, 0, 0.7);
            padding: 20px;
            border-radius: 10px;
            display: inline-block;
        }
        label {
            font-size: 20px;
        }
        input, select {
            margin-top: 10px;
            margin-bottom: 20px;
            padding: 10px;
            width: 100%;
            border-radius: 5px;
        }
    </style>
</head>
<body>
    <div class="form-container">
        <h1>Next Questions</h1>
        <div class="form-group">
            <label for="name">Name:</label>
            <input type="text" class="form-control" id="name" name="name" required>
        </div>
        <div class="form-group">
            <label for="age">Age:</label>
            <input type="text" class="form-control" id="age" name="age" required>
        </div>
        <div class="form-group">
            <label for="gender">Gender:</label>
            <select class="form-control" id="gender" name="gender">
                <option value="M">Male</option>
                <option value="F">Female</option>
            </select>
        </div>
        <div class="form-group">
            <label for="email">Email:</label>
            <input type="text" class="form-control" id="email" name="email" required>
        </div>
        <div class="form-group">
            <label for="password">Password:</label>
            <input type="password" class="form-control" id="password" name="password" required>
        </div>
        <div class="form-group">
            <label for="confirm_password">Confirm Password:</label>
            <input type="password" class="form-control" id="confirm_password" name="confirm_password" required>
        </div>
        <button type="submit" class="btn btn-primary">Submit</button>
    </div>
</body>
</html>

```

```

        border: none;
    }
    button {
        padding: 10px 20px;
        font-size: 20px;
        background-color: #4CAF50;
        color: white;
        border: none;
        border-radius: 5px;
        cursor: pointer;
    }
</style>
</head>
<body>
    <h1>Additional Information</h1>
    <div class="form-container">
        <form action="/next" method="post">
            <label for="foreign_languages">Foreign Languages Known:</label>
            <input type="text" id="foreign_languages" name="foreign_languages" required><br><br>

            <label for="native_language">Native Language:</label>
            <input type="text" id="native_language" name="native_language" required><br><br>

            <label for="english_fluency">English Fluency:</label>
            <select id="english_fluency" name="english_fluency" required>
                <option value="Fluent">Fluent</option>
                <option value="Intermediate">Intermediate</option>
                <option value="Beginner">Beginner</option>
            </select><br><br>

            <label for="reason">Reason for Applying:</label>
            <input type="text" id="reason" name="reason" required><br><br>

            <label for="work_experience">Years of Work Experience:</label>
            <input type="number" id="work_experience" name="work_experience" required><br><br>

            <label for="country">Country of Origin:</label>
            <input type="text" id="country" name="country" required><br><br>

            <label for="percentage">Marks/Percentage (%):</label>
            <input type="number" step="0.1" id="percentage" name="percentage" required><br><br>

            <label for="night_shifts">Willing to work night shifts? (Yes/No):</label>
            <select id="night_shifts" name="night_shifts" required>
                <option value="Yes">Yes</option>
                <option value="No">No</option>
            </select><br><br>

            <button type="submit">Submit</button>
        </form>
    </div>
</body>
</html>

```

10.5:
resultVA.html:

```
<!DOCTYPE html>
```



```
<html lang="en">
<head>
  <meta charset="UTF-8">
  <meta name="viewport" content="width=device-width, initial-scale=1.0">
  <title>H-1B Visa Prediction Result</title>
  <style>
    body {
      font-family: Arial, sans-serif;
      background-image: url('{{ url_for('static', filename='images/plane.jpg') }}');
      background-size: cover;
      background-position: center;
      color: white;
      text-align: center;
      padding-top: 100px;
    }
    .result {
      background-color: rgba(0, 0, 0, 0.7);
      padding: 20px;
      border-radius: 10px;
      display: inline-block;
      margin-top: 50px;
    }
    .result p {
      font-size: 36px;
      font-weight: bold;
      color: #4CAF50;
    }
  </style>
</head>
<body>
  <h1>H-1B Visa Predictor</h1>
  <div class="result">
    <p>{{ prediction_text }}</p>
  </div>
</body>
</html>
```

10.6:Github Link:

<https://github.com/AbhignaVelamakanni/H1bvisa/tree/master>

10.7 demo video link :

https://drive.google.com/file/d/1S7ISNZ5bqikty4Kw_-XdsGrJReX3lbd6/view?usp=drivesdk