
SATHYA

Easy Visa Project

Objective

In FY 2016, the OFLC processed 775,979 employer applications for 1,699,957 positions for temporary and permanent labor certifications. This was a nine percent increase in the overall number of processed applications from the previous year. The process of reviewing every case is becoming a tedious task as the number of applicants is increasing every year.

The increasing number of applicants every year calls for a Machine Learning based solution that can help in shortlisting the candidates having higher chances of VISA approval. OFLC has hired your firm EasyVisa for data-driven solutions. You as a data scientist have to analyze the data provided and, with the help of a classification model:

- * Facilitate the process of visa approvals.
 - * Recommend a suitable profile for the applicants for whom the visa should be certified or denied based on the drivers that significantly influence the case status.
-

Data Set

The data Set has 25480 samples with 12 columns

No Null values

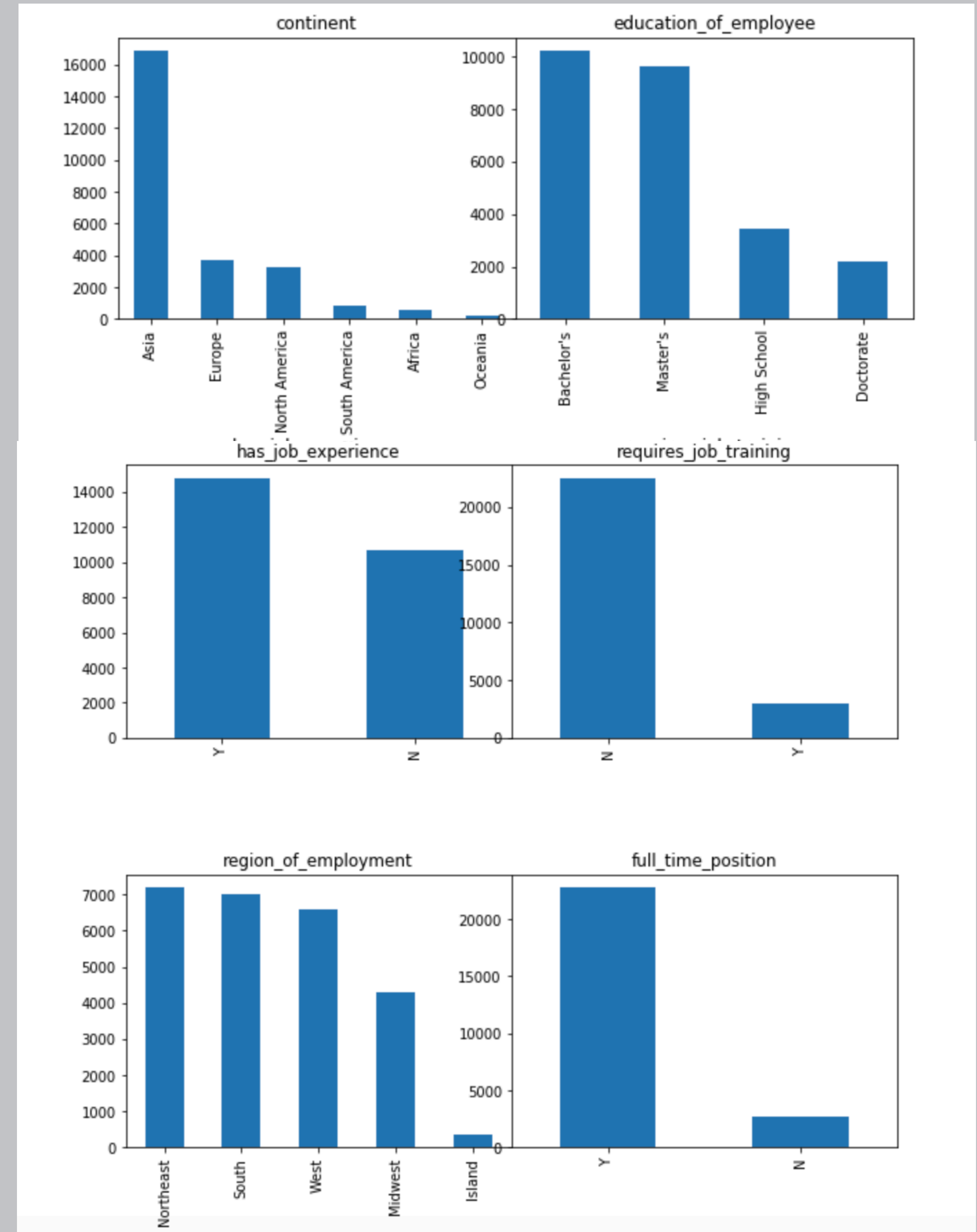
The following are the columns available in the Data set :

- * case_id: ID of each visa application
 - * continent: Information of continent the employee
 - * education_of_employee: Information of education of the employee
 - * has_job_experience: Does the employee has any job experience? Y= Yes; N = No
 - * requires_job_training: Does the employee require any job training? Y = Yes; N = No
 - * no_of_employees: Number of employees in the employer's company
 - * yr_of_estab: Year in which the employer's company was established
 - * region_of_employment: Information of foreign worker's intended region of employment in the US.
 - * prevailing_wage: Average wage paid to similarly employed workers in a specific occupation in the area of intended employment. The purpose of the prevailing wage is to ensure that the foreign worker is not underpaid compared to other workers offering the same or similar service in the same area of employment.
 - * unit_of_wage: Unit of prevailing wage. Values include Hourly, Weekly, Monthly, and Yearly.
 - * full_time_position: Is the position of work full-time? Y = Full Time Position; N = Part Time Position
 - * case_status: Flag indicating if the Visa was certified or denied
-

EDA

Observation:

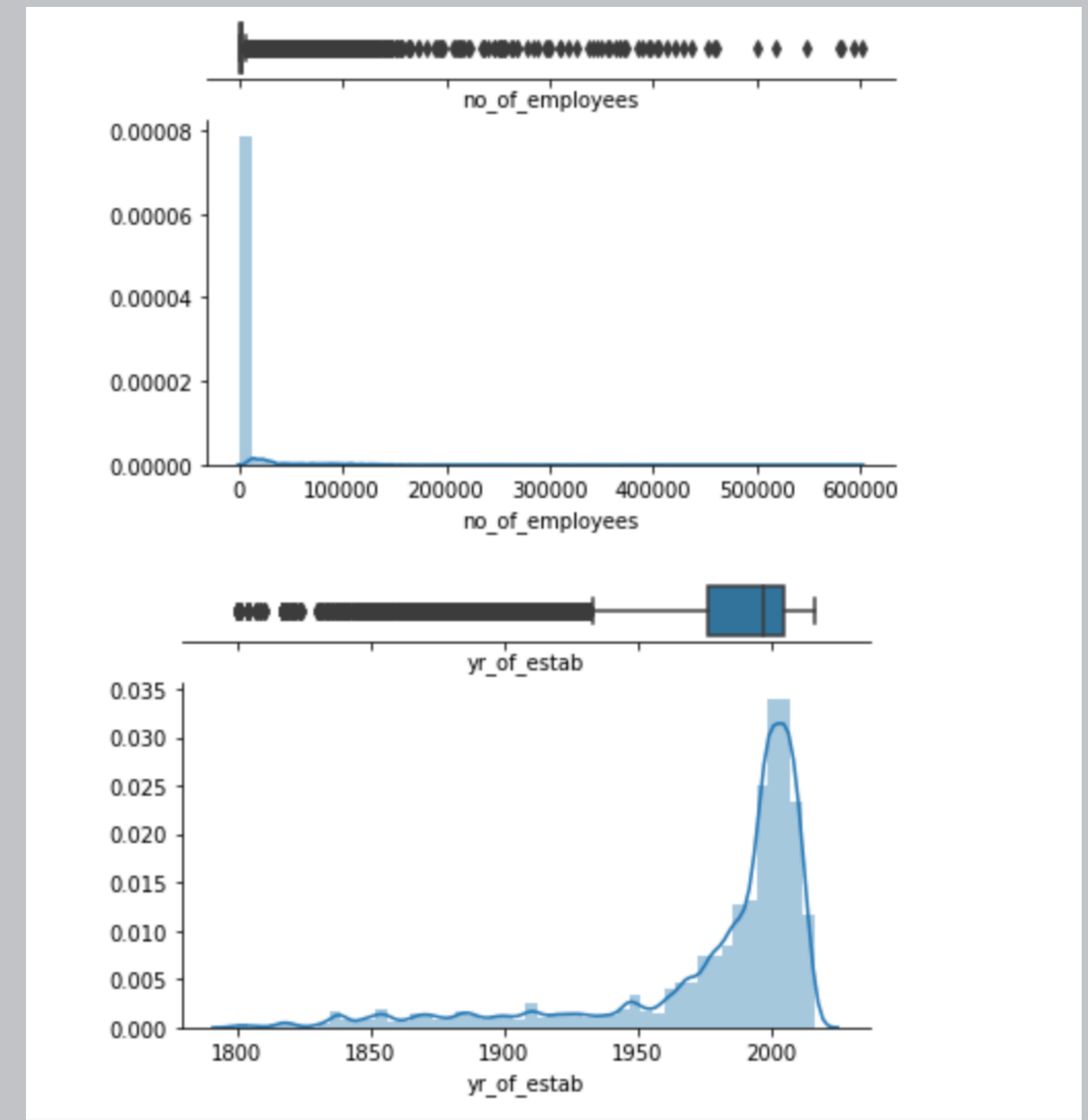
- 1) High number of applicants from Asia Continent
- 2) Mostly Educations of applicants are Bachelor & Masters
- 3) Applicants dont need job training , applicants are with the prior experience for the job
- 4) The regions of the employer are from Noreast & South
- 5) Applicants are mostly applied for the full time position



EDA

Observation: (Employers Profile)

- 1) Median value for Years of establishments 1997
- 2) No of employees with mean 5667.043210 and median 2109.00

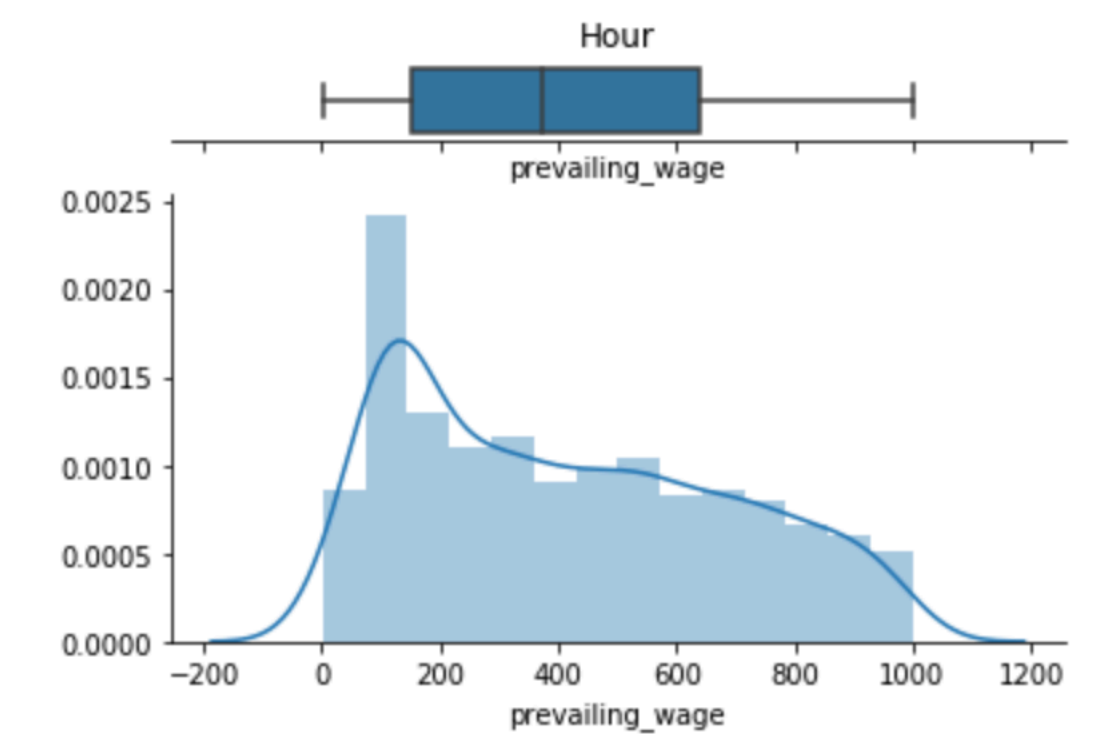
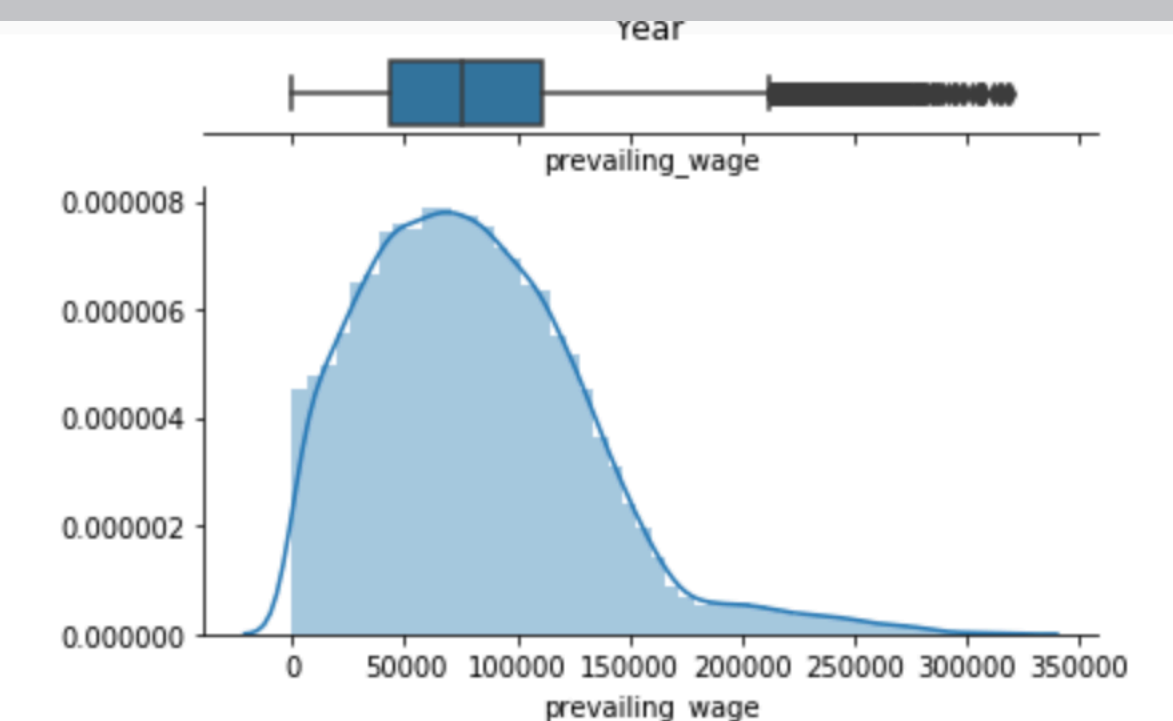
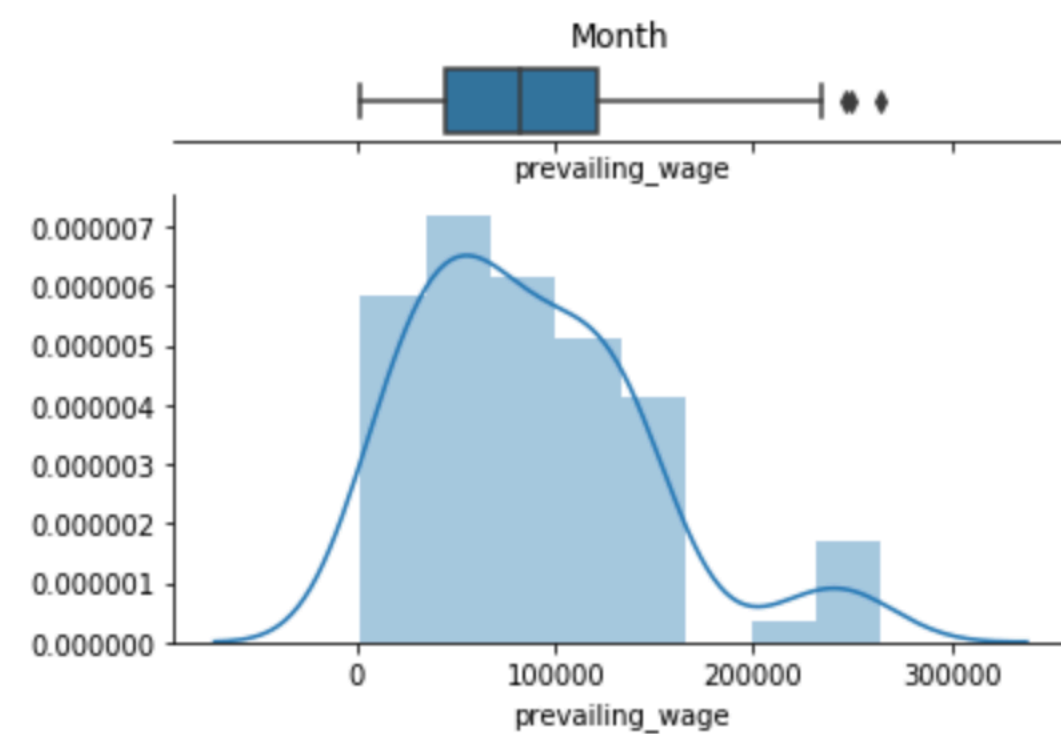
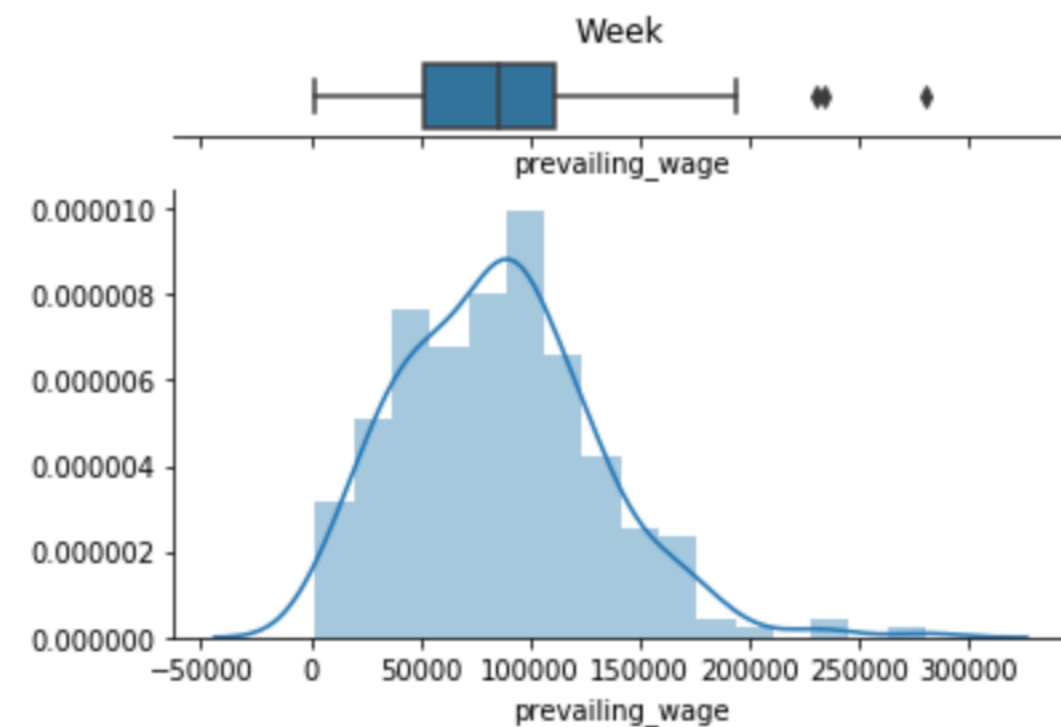


EDA

Observations :

The wages are at different unit. Hours , week , month, Year

The high number of unit of wages are Year and it follows normal distribution of data. Except the Hourly wage type , others follows normal distribution.

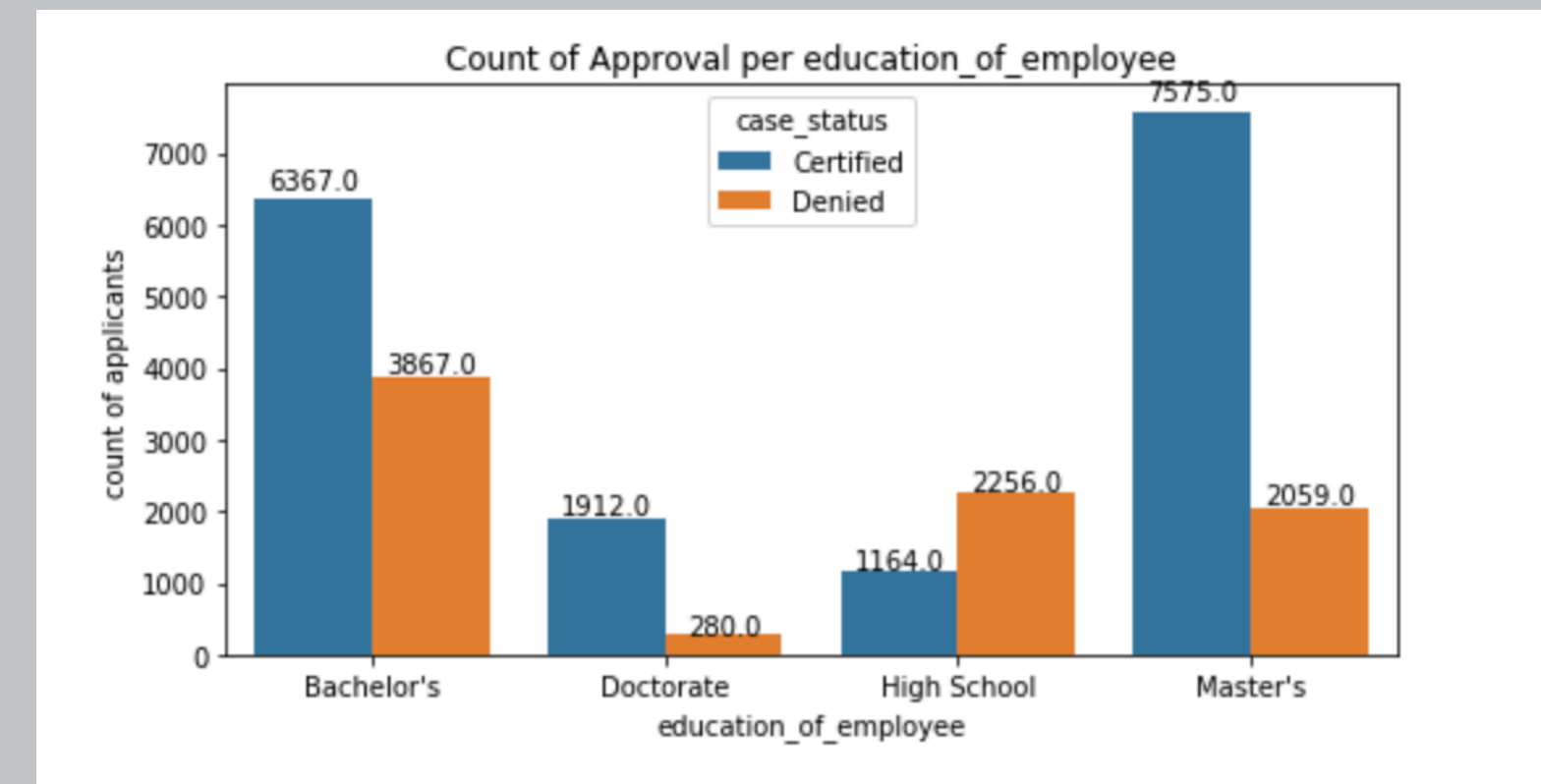


EDA

Observation :

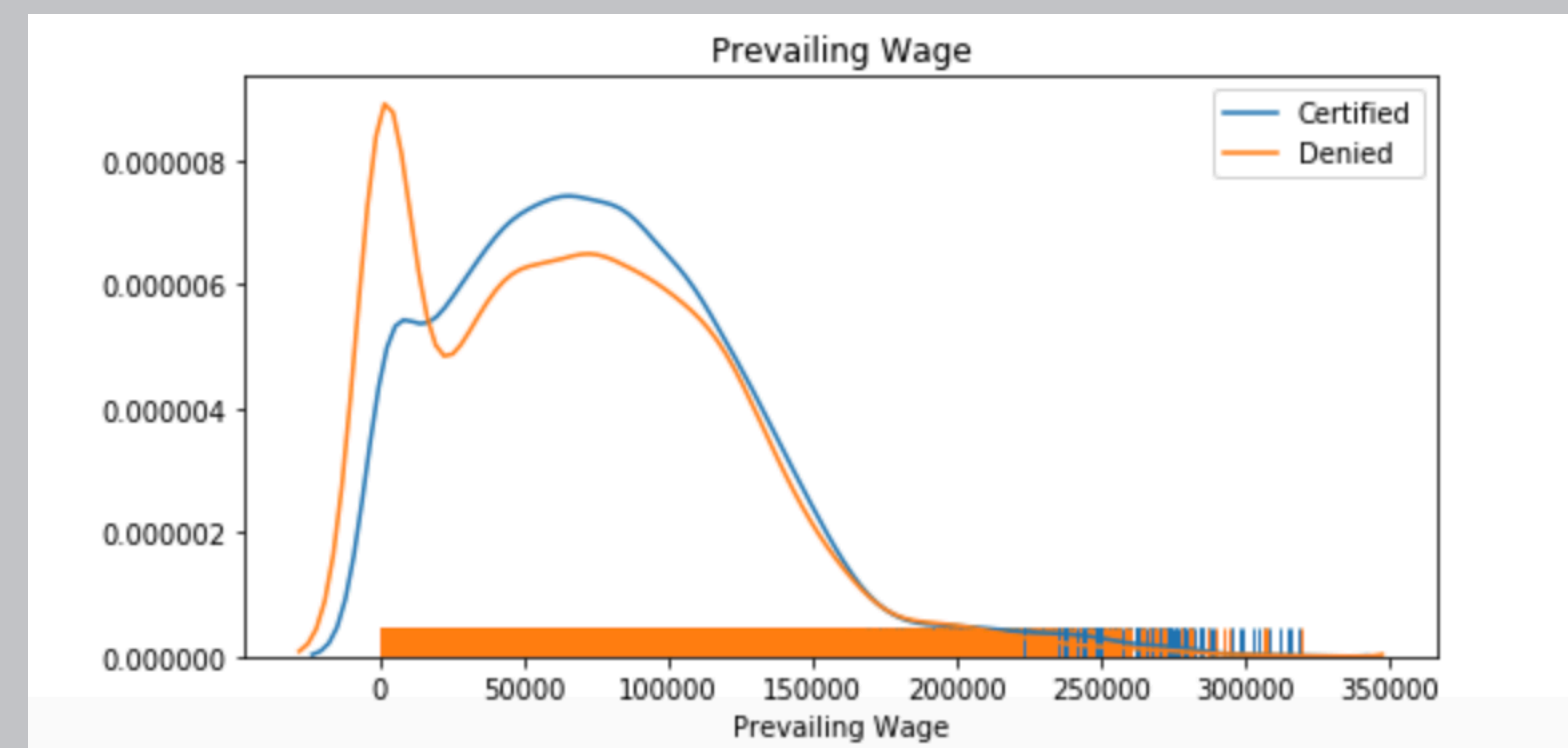
The visa status is rejected mostly with the applicants of High School.

The higher education level has high probability of visa approval



Observation :

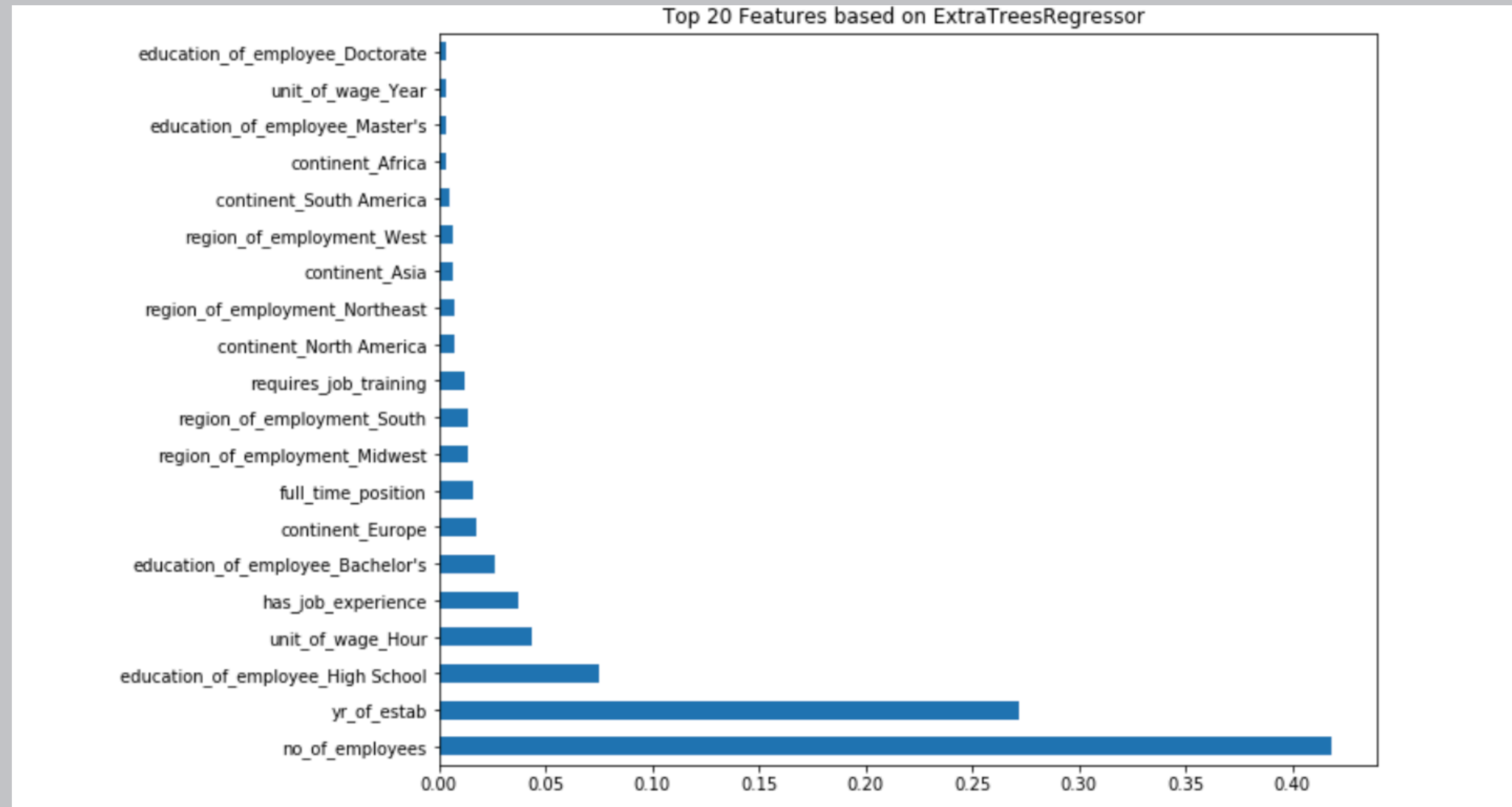
the prevailing wage is more for certified applicants than the Denied applicants



Correlation Matrix

no_of_employees	1	-0.02	-0.01	-0.02	-0.01	-0.02	-0.01	0	-0.02	-0.01	-0.01	-0.01	-0.04	-0.01	-0.02	0.05	-0.01	0.01	-0.01	-0	-0	0	-0.01	0.01	0.0
yr_of_estab	-0.02	1	-0.01	0.03	-0.02	-0.01	-0	-0.02	0.01	-0	-0.01	-0	-0.04	-0.01	-0.02	0.05	-0.01	0.01	-0.01	-0	-0	0	-0.01	0.01	0.0
continent_Africa	0.01	-0.01	1	-0.21	-0.06	-0.06	-0.01	-0.03	-0.04	0.01	-0.01	0.04	-0.01	0	0.01	0.01	-0	-0	0.03	-0.03	-0	0.01	-0.02	0.02	0.0
continent_Asia	-0.02	0.03	-0.21	1	-0.58	-0.54	-0.12	-0.26	0.07	-0.16	0.01	0.02	-0.03	-0	-0.02	0.04	-0.02	0.12	-0.03	0.03	-0.09	-0	-0.05	-0.04	-0.0
continent_Europe	0.01	-0.02	-0.06	-0.58	1	-0.16	-0.04	-0.08	-0.05	0.21	-0	-0.07	-0.07	-0.01	-0.02	0.08	-0.02	-0.1	-0.02	-0.01	0.13	-0	0.11	0.1	0.1
continent_North America	0.02	-0.01	-0.06	-0.54	-0.16	1	-0.03	-0.07	-0.02	-0.01	-0.01	0.04	0.09	0.01	0.05	-0.1	0.03	-0.05	0.05	-0.01	-0	0.01	-0.05	-0.08	-0.0
continent_Oceania	-0.01	-0	-0.01	-0.12	-0.04	-0.03	1	-0.02	-0.01	0.01	0.01	-0	0	0.01	0.01	-0.01	0	0.01	-0.01	-0.01	0.01	-0.01	0.01	0	-0.0
continent_South America	-0	-0.02	-0.03	-0.26	-0.08	-0.07	-0.02	1	-0	0.01	0.01	-0.01	0.07	-0	0.01	-0.06	0.02	-0.02	0.01	-0.01	0	-0.01	0.03	0.04	-0.0
education_of_employee_Bachelor's	-0.02	0.01	-0.04	0.07	-0.05	-0.02	-0.01	-0	1	-0.25	-0.32	-0.64	0.03	0.01	0.01	-0.04	-0.01	-0.09	-0	0.03	0.05	-0.01	0.02	0.11	-0.0
education_of_employee_Doctorate	0.02	-0	0.01	-0.16	0.21	-0.01	0.01	0.01	-0.25	1	-0.12	-0.24	-0.05	-0.01	-0.02	0.05	-0.01	-0.04	0.01	-0.02	0.05	-0.01	0.07	0.07	0.1
education_of_employee_High School	-0.01	-0.01	-0.01	0.01	-0	-0.01	0.01	0.01	-0.32	-0.12	1	-0.31	0.04	0	-0.01	-0.04	0.01	0.05	-0.02	-0	-0.03	0.01	0.01	0.06	-0.2
education_of_employee_Master's	0.02	-0	0.04	0.02	-0.07	0.04	-0	-0.01	-0.64	-0.24	-0.31	1	-0.04	-0.01	0	0.04	0.01	0.08	0.01	-0.02	-0.06	0.01	-0.07	-0.19	0.2
unit_of_wage_Hour	-0.02	-0.04	-0.01	-0.03	-0.07	0.09	0	0.07	0.03	-0.05	0.04	-0.04	1	-0.02	-0.03	-0.92	-0.02	-0.08	0.08	0.03	-0.04	-0.13	-0.08	0.1	-0.2
unit_of_wage_Month	-0.01	-0.01	0	-0	-0.01	0.01	0.01	-0	0.01	-0.01	0	-0.01	-0.02	1	-0.01	-0.18	-0.01	-0	0.01	0.01	-0.01	0.01	-0.01	0.01	-0.0
unit_of_wage_Week	0	-0.02	0.01	-0.02	-0.02	0.05	0.01	0.01	0.01	-0.02	-0.01	0	-0.03	-0.01	1	-0.31	-0	-0.01	0.01	-0.01	0.01	0.03	-0.03	0.02	-0.0
unit_of_wage_Year	0.02	0.05	0.01	0.04	0.08	-0.1	-0.01	-0.06	-0.04	0.05	-0.04	0.04	-0.92	-0.18	-0.31	1	0.02	0.08	-0.08	-0.03	0.04	0.11	0.09	-0.1	0.2
region_of_employment_Island	0.01	-0.01	-0	-0.02	-0.02	0.03	0	0.02	-0.01	-0.01	0.01	0.01	-0.02	-0.01	-0	0.02	1	-0.06	-0.08	-0.08	-0.07	0	-0.02	-0.01	-0.0
region_of_employment_Midwest	-0.02	0.01	-0	0.12	-0.1	-0.05	0.01	-0.02	-0.09	-0.04	0.05	0.08	-0.08	-0	-0.01	0.08	-0.06	1	-0.28	-0.28	-0.27	-0.03	-0.09	0	0.0
region_of_employment_Northeast	0.01	-0.01	0.03	-0.03	-0.02	0.05	-0.01	0.01	-0	0.01	-0.02	0.01	0.08	0.01	0.01	-0.08	-0.08	-0.28	1	-0.39	-0.37	0.01	-0.06	-0.01	-0.0
region_of_employment_South	0.01	-0	-0.03	0.03	-0.01	-0.01	-0.01	-0.01	0.03	-0.02	-0	-0.02	0.03	0.01	-0.01	-0.03	-0.08	-0.28	-0.39	1	-0.36	0	0.09	-0.05	0.0
region_of_employment_West	0	-0	-0	-0.09	0.13	-0	0.01	0	0.05	0.05	-0.03	-0.06	-0.04	-0.01	0.01	0.04	-0.07	-0.27	-0.37	-0.36	1	0.01	0.05	0.07	-0.0
has_job_experience	0.01	0	0.01	-0	-0	0.01	-0.01	-0.01	-0.01	-0.01	0.01	0.01	-0.13	0.01	0.03	0.11	0	-0.03	0.01	0	0.01	1	-0.11	0.04	0.1
requires_job_training	-0.01	-0.01	-0.02	-0.05	0.11	-0.05	0.01	0.03	0.02	0.07	0.01	-0.07	-0.08	-0.01	-0.03	0.09	-0.02	-0.09	-0.06	0.09	0.05	-0.11	1	0.1	0.0
full_time_position	-0.01	0.01	0.02	-0.04	0.1	-0.08	0	0.04	0.11	0.07	0.06	-0.19	0.1	0.01	0.02	-0.1	-0.01	0	-0.01	-0.05	0.07	0.04	0.1	1	-0.0
case_status	-0.01	-0.01	-0.02	-0.04	0.11	-0.04	-0.01	-0.04	-0.08	-0.13	-0.27	-0.2	-0.21	-0.01	-0.01	-0.2	-0.02	-0.08	-0.05	-0.04	-0.06	-0.15	-0.01	-0.01	1
no_of_employees	no_of_employees	yr_of_estab	continent_Africa	continent_Asia	continent_Europe	continent_North America	continent_Oceania	continent_South America	education_of_employee_Bachelor's	education_of_employee_Doctorate	education_of_employee_High School	education_of_employee_Master's	unit_of_wage_Hour	unit_of_wage_Month	unit_of_wage_Week	unit_of_wage_Year	region_of_employment_Island	region_of_employment_Midwest	region_of_employment_Northeast	region_of_employment_South	region_of_employment_West	has_job_experience	requires_job_training	full_time_position	case_status

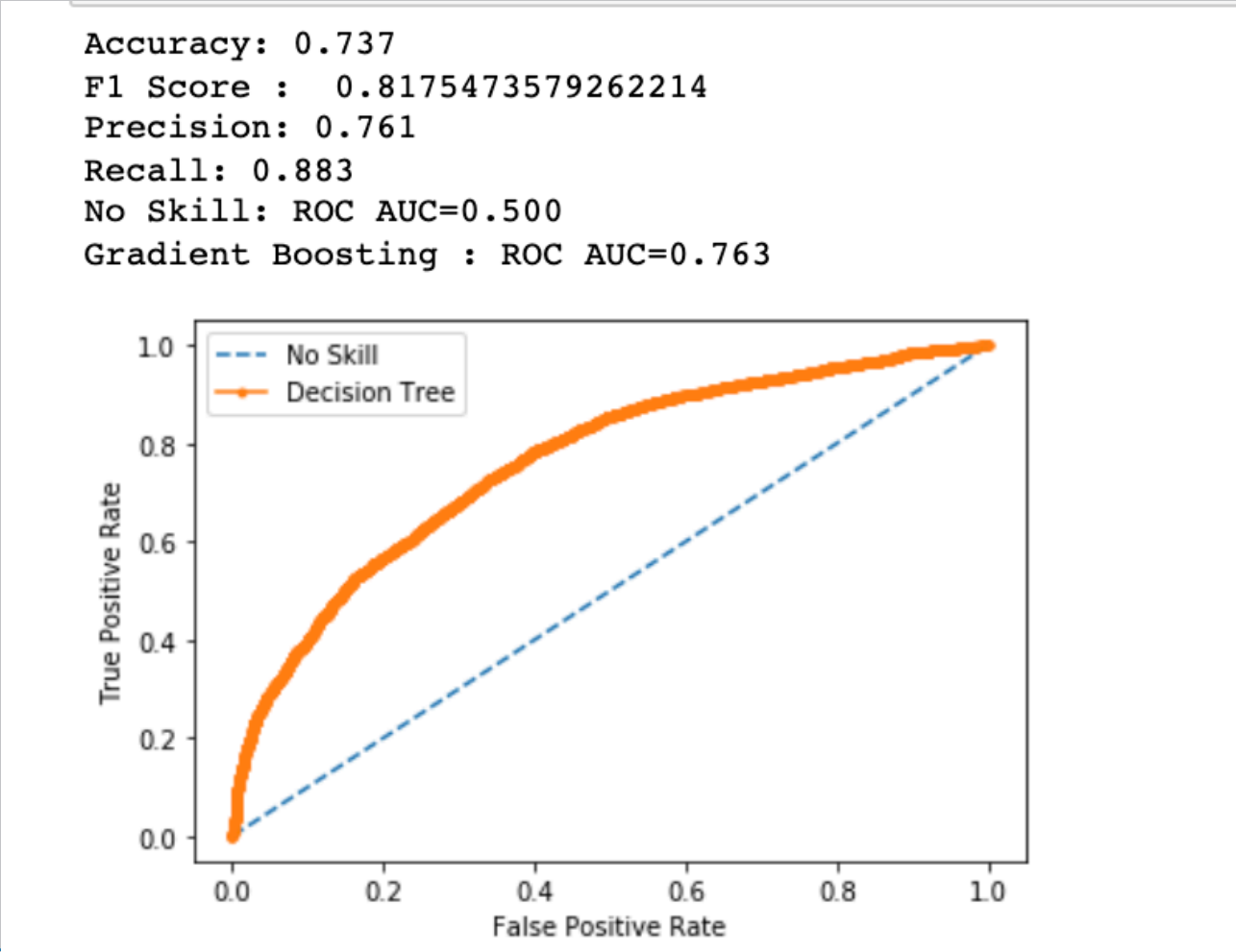
Top 20 Features



Model Performance And Evaluation

Sno	Model	Accuracy	Precision	Recall	F1_Score
1	Random Forest	0.690	0.756	0.790	0.7727229192606072
2	Decision Tree	0.659	0.748	0.738	0.7431952662721892
3	Logistic Regression	0.668	0.668	1.0	0.8008471252647267
4	Gradient Boosting	0.737	0.761	0.883	0.8175473579262214
5	Tuned Random Forest	0.701	0.761	0.804	0.7820952380952382

Out of the 5 Models the Gradient Boosting model performs better with F1 Score as 0.81



The ROC curve of the final selection Gradient Boosting Model

Conclusion

The Following are the actionable Insights:

- > We can add additional features to the data set like Relevant Years of Experience
- > Field of Study can be included as part of education