

An allotment of H1B work visa in USA using machine learning

Pooja Thakur ^{1*}, Mandeep Singh ², Harpreet Singh ³, Prashant Singh Rana ⁴

¹ Chandigarh University, Mohali, India

² Chandigarh University, Mohali, India

*Corresponding author E-mail: Thakur208903@gmail.com

Abstract

H1B work visas are utilized to contract profoundly talented outside specialists at low wages in America which help firms and impact U.S economy unfavorably. In excess of 100,000 individuals for every year apply tight clamp for higher examinations and also to work and number builds each year. Selections of foreigners are done by lottery system which doesn't follow any full proofed method and so results cause a loophole between US-based and foreign workers. We endeavor to examine petitions filled from 2015 to 2017 with the goal that a superior prediction model need to develop using machine learning which helps to foresee the aftereffect of the request of ahead of time which shows whether an appeal to is commendable or not. In this work, we use seven classification models Decision tree, C5.0, Random Forest, Naïve Bayes, Neural Network and SVM which predict the status of a petition as certified, denied, withdrawal or certified withdrawals. The predictions of these models are checked on accuracy parameter. It is found that C5.0 outperform with the best accuracy of 94.62 as a single model but proposed model gives better results of 95.4 accuracies which is built by machine ensemble method and this is validated by 10 fold cross-validation.

Keywords: *H1b; Decision Tree; Random Forest; Naïve Bayes; Neural Network; C5.0; SVM.*

1. Introduction

Visa is the guide of authorization on a travel permit that gives a permit to the holder to move in, leave or stay in the country for a predetermined timeframe. There are distinctive kinds of foreigner visas, the required structures, and the means in the worker visa process contingent upon the nation one needs to move. Moving to America is a vital and complex decision. The U.S of America has numerous classes for settler visas like H1B, L1, and J1 and so on [1]. To be qualified to apply for a worker visa, an outside native must be supported by a USA subject relative, U.S. legitimate perpetual inhabitant, or a planned business, with a couple of special cases. The help begins the movement methodology by recording an interest to for the remote inhabitant's purpose with U.S. Residency and Colonization Facilities (USCIS). Among the better piece of this H-1B are greatly outstanding starting late due to manufactures no of petitions and wrong system for getting consent. H1B is a visa characterization in America under movement and nationality act (INA) [2]. Empowers U.S supervisors to yield outside workers with high degrees and capable of "distinguishing strength occupations". H-1B is a business based non-transient visa gathering for brief remote specialists in the US. For an outside national to apply for H1B visa, a US business must offer an occupation and request to for H-1B visa with the US movement office. This is the most widely recognized visa status connected to and held by universal understudies once they finish school/advanced education (Masters, Ph.D.) and work in a full-time position. The Office of Foreign Labour Certification (OFLC) [3], [4] creates pro- gram information that is helpful data about the movement programs including the H1-B visa. It is intended to carry outside experts with professional educations and specific aptitudes to fill occupations when qualified Americans can't be found. Be that as it may, as of late, worldwide outsourcing organizations have ruled the program,

winning a huge number of visas and pressing out numerous American organizations, including littler new companies. The development in the portrayal of the outside conceived among the US workforce was brought down drastically. In any case, there was a really staggering augmentation in the remote supply of men and women with school preparing in science and building fields [5]. To take one vital case, in India, the quantity of first degrees presented in science and designing rose from 176 thousand of every 1990 to 455 thousand of every 2000. Second, the Act of 1990 set up the H-1B visa package for impermanent labourers in "claim to fame occupations"[4]. The rules defines "claim to fame occupation" as requiring hypothetical and common-sense use of a collection of exceptionally particular learning in a field of human undertaking including, yet not constrained to, design, building, arithmetic, physical sciences, sociologies, solution and wellbeing, instruction, law, bookkeeping, business fortes, religious philosophy, and expressions of the human experience. Furthermore, candidates are required to have achieved a four year certification or it's identical as a base. Firms that desire to contract non-natives on H-1B visas have needs to file a Labour Condition Application (LCA) [4]. In LCA's for H-1B specialists, the business must bear witness to that the firm will pay the non-foreigner the more prominent of the genuine remuneration paid to different representatives in a similar activity or the common pay for that occupation, and the firm will give working conditions to the non- migrant that don't make the working states of alternate workers be unfavourably affected. By then, planned H-1B non-outsiders must exhibit to the US Citizenship and Immigration Services Bureau (USCIS) in America [5]. In spite of the fact that H1B visa contributed a considerable measure to the economy of USA by bringing the skilled non-natives, it additionally influences American work [1]. They lose their employment, as firms incline toward modest work when contrasted with American's. The objective of the H1B program is to connect a work hole in the U.S without influencing U.S specialists. At to

start with, the structure of H1B is to fill work hole however current structure encourages businesses to augment the work hole as there aren't any qualified U.S specialists and they are procuring modest remote labourers as H1B program. There are 3 primary targets of the H1B program: Section 1) To connect a work hole without dislodging U.S specialists forever. Section 2) Review the present structure of H1B program, concentrating on the way toward acquiring H1B visa and what it enables its holders to do. It gives two classifications of the run the show: Qualification of remote specialists.

- a) The framework guarantees that outsiders don't dislodge U.S specialists.

Section 3) Effect of H1B structure on compensation framework. Paying non-natives not exactly or equivalent to U.S representatives to make a disincentive U.S specialists. All these are primary destinations of the H1B program yet it additionally offers to ascend to a few issues like every single talented specialist doesn't get endorsement due to expansive any applications filled by outsourcing firms and the second one is about lost U.S labourers' work because of the procuring of shoddy remote labourers'. These issues are a major provision in the strategy of the H1B visa framework. A profound knowledge is required with the goal that the businesses comprehend the procedure of the visa appeal, to stop the outsourcing firms backtrack applications, amending of the framework, utilizing better techniques like compensation based, justify based or encounter based for conceding of visas. This paper, endeavours to foresee the negative and positive consequence of the applications and finding for which sort of occupation no. of petitions are high or low with the goal that contracting of economical work is extremely dense by utilizing machine learning techniques. The main objectives of current work are as follows:

- 1) Detail investigation of effectively existing machine learning systems and upgraded Machine Learning approaches for a better forecast.
- 2) Approve different models in the wake of observing at on premise insights estimations for using sensible endorsement framework.
- 3) Finally, we prepare a proposed model with the marked data to foresee future petitions as a right one or mishandle and then validate it by using suitable validation technique.

2. Discussion

2.1. Related works

Dhanasekar Sundararaman [1]: In this paper petition filed from 2011-16 analysed using random forest. Prediction of any visa petition in any state is classified as negative or positive. Companies that acquire H1B visa disproportionately with very less wages are identified, which completely conflict the goal of this program. It aimed to secure highly skilled talented employers in case of shortage of same talented people in U.S. decision tree show results with accuracy of 99% by classifying visa petitions

Trim Bach [4]: Demonstrates an investigation of H1B datasets which demonstrates that businesses outsource work and pay fewer wages to the remote specialists when contrasted with U.S based bosses. The structure of low wages impacts inside to the U.S worker's, as firms favour H1B representatives which stagnate U.S specialists to strive for that activity and increment the provision in the H1B framework .

Doran [5]: It includes winning and losing firms in the monetary year 2006 and 2007 lotteries for H1B visas. Winning relates to a direction in the association's general work, prompt a lower normal worker gaining and higher firm benefits and insignificantly affects company's general business, prompts a lower normal representative acquiring and higher firm benefits and insignificantly affects association's licensing and utilization of the examination and experimentation charge credit. They discovered additional H-1B

increment middle firm benefits; diminish in middle income per representative .

Bound [6]: One of the current work which connotes the benefits of H-1B on US economy and the lessening of the costs of PC related innovation and expanded the yield. It likewise tells how firms earned plenty of benefits utilizing these visas, set up of local individuals .

Mithas et al. [7]: finds that the pay for outsiders and those on labour permits vary in light of supply shocks made by tops on new H-1B visas. Lower and completely used tops outcomes in a higher pay for foreigners and those on work visas. These days, these tops are filled by lottery which in a circuitous way proposes that an excessive number of undeserving individuals apply .

Despite the fact that there are different hypothetical and overview chips away at focal points of H-1B and how these visas are abused in the current circumstances, there isn't a strong reasonable execution to discover the utilization of these visas.

We endeavour to demonstrate it by taking a down to earth H-1B dataset and discover the organizations, which pay lower than the greater part of the others, and organizations, which record a large portion of these applications and term them negative.

2.2. Dataset and attributes

The data set for foreseeing examination on H-1B visa are gathered from OFLC database. It is an association which has the obligations of handling of work confirmation and work validation applications. The Office of Foreign Labour Certification (OFLC) produces program data that is important both for centre valuation of program proficiency and furnishing the Department's outside financial specialists with helpful measurements about the migration programs oversight by OFLC [3]. This association incorporates database data organized into 3 primary classifications:

- 1) OFLC's yearly reports for developing sign of program data and data;
- 2) Measurements by Program for giving quarterly data of significant migration programs in charge of depiction perspectives of the OFLC programs [8, 9]; and
- 3) Quarter and yearly issues of program introduction data to add to outside Investigation and database assessment [3, 9]. H1B dataset from OFLC contained about 40 attributes out of which, some are expelled in data cleaning technique. Original datasets are shown in table 1.

Unimportant properties are evacuated by profoundly comprehension of dataset and by finding a connection between these qualities. Traits like Case_Number, Employer_Phone, Employer_Address, and Employer phone_Ext are having one of a kind esteem, henceforth are immaterial in foreseeing target i.e. case status, so these characteristics are evaluated. The greater part of the properties is in content configuration which isn't the suitable arrangement for machine learning procedure so pre-processing of data is required. Sampled dataset taken as input are shown in table 2 and table 3 shows output variables. Pairs .panels (figure 9: Pearson correlation between attributes) demonstrate a scramble plot of networks (SPLOM), with bivariate, diffuse plots beneath the corner to corner, histograms on the inclining and the Pearson correlation over the planting which is proven valuable for unmistakable measurements of data sets. On the bases of correlation and significant investigation of dataset just 20 qualities approach as pertinent, and thus just these are operated as a part of further work. At last all relevant qualities are changed over into some numeric incentive by producing recipe simply on some scientific idea at that point last dataset has 528136 esteems with 20 qualities, out of which CASE_STATUS is the objective characteristic. Distributions of some attributes are shown in figure 1.

Table 1: Attributes of Dataset

Attribute	Description
Submitted_Date	Date and Time the application was submitted
Case_No	Case number
Program_Designation	Type of Visa classes
	H1B
	E-3 Australian
EMPLOYER_Name	H1B1 Chile
	H1B1 Singapore
	Employer's name
Employer_Address 1	Employer's address
Employer_Address 2	Employer address 2
.	.
.	.
.	.
Employer_postal code	Employer's postal code
SOC_NAME	Name of wage sources
TOTAL_WORKERS	Total no of workers from one firm
PREVAILING_WAGE	Prevailing wage rate
PW_UNIT_OF_PAY	Unit of pay for wages.
PW_SOURCE	Collective bargaining
PW_SOURCE_YEAR	Year of prevailing wages
PW_SOURCE_OTHER	Others sources for prevailing wages
WAGE_RATE_OF_PAY_FROM	Starting of pay
WAGE_RATE_OF_PAY_TO	Maximum of wages
WAGE_UNIT_OF_PAY	Year of wages
H1B_DEPENDENT	Candidates depend on Visa classes
WILLFUL_VIOATOR	Violated candidates
CASE_STATUS	Status of candidates

The experts in charge of visas alongside the administration are taking activities to keep this and guarantee these visas go to the merited individuals. A portion of the proposed philosophies incorporate legitimacy based and pay based allowing for the visa. In this paper, we endeavour to dissect the instance of visa manhandle state astute by finding the situations for which the quantity of petitions is high, and the compensation is low in each state.

To compress, :

- 1) First, we assemble all the data about a visa request of like compensation, position, manager state shrewd for every one of the 50 states.
- 2) Second, we perform essential factual tasks like mean, highest and least pay for which the request of was petitioned for each state.

Table 2: Sample Dataset Taken As Inputs

Visa clas- s	Em- ployer- name	Em- ploy- er- State	Em- ployer- coun- try	H1B- de- pen- dent	Works ite- state	Works ite- postal- code	
0	43436	26	2	.	1	24	10469
0	46027	48	2	.	2	46	5912
.
0	17079	40	2	.	2	38	2671
0	27211	54	2	.	1	24	8875

Table 3: Sample Dataset Taken as Output

s. no	Case_Status
1	1
2	1
.....
528134	2
528135	3

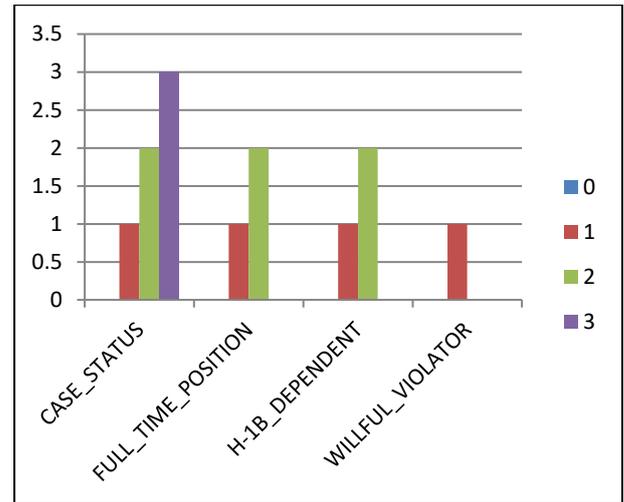


Fig. 1: Distribution of Different Attributes.

- 1) Third, we utilize grouping to locate the total pay incentives for each state, underneath which an excessive number of petitions are recorded to be delegated manhandle, i.e., the end-point esteem.

After an accumulation of data and pre-processing subsequent stage go for targets that to acquire substantial data set by standardization utilizing different systems. By finding such outright pay esteem, which candled the cut esteem, it is helpful to plan and effective visa framework in which visas are allowed for individuals who in any event draw The cut-off esteems pay. Finally, visas are allowed to the merited individuals, and use of visas to employ inferior work is definitely diminished. The highlights are utilized by seven machine learning models in particular decision tree, linear model, C5.0, Support Vector Machine, Naïve Bayes, Neural Network and Random forest for providing better prediction model for H1B Visa petitions without any involvement from firms or petitioners. Different statistics measures are used for measuring the best prescient model. Rest of the paper is sorted out as takes after. A brief diagram of the thought about highlights, data set, and correlation model; furthermore, machine learning models are exhibited in section 2. The proposed discrimination of all models are anticipated and portrayed in Section 3. Demonstrate assessment is displayed in Section 4. Segment 5 portrays tests, results, and talk. At last, conclusion is exhibited in Section 6.

2.3. Exploratory data analysis

After collection and pre-processing of data, dataset assembled state wise to get insight of data, due to the serious fluctuation of pay crosswise over states e.g., compensation for a product designer is observed to be much lower in Florida than that of California. Presently subsequent to gathering the data state astute, few data examination assignments and en- vision are done for investigation. The accompanying is the data investigation errands required for a state. Distributions of different visa classes are shown in figure 2.

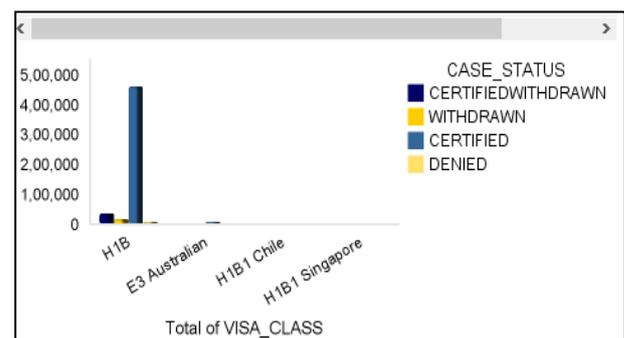


Fig. 2: Distribution of Different Type of Visa Classes.

- The mean pay offered in a state for a specific year .this pattern is assembled and we find that Computer Occupations are the most paid posts with as shown in figure 4.
- And Analyst is second in the rundown. Most denied and Pulled back cases are from Scientist as presented in figure 5.
- Computer Occupations are at the top for both years 2016 and 2017 (as shown in figure 6 and 7).
- The visa petitions are assembled in light of the status of the visa i.e. denied, certified, certified withdrawn and withdrawn (as shown in figure 3).

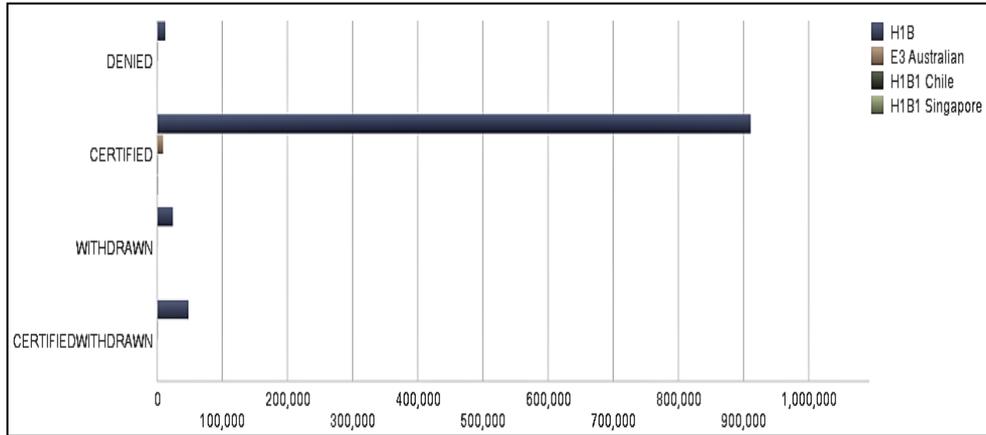


Fig. 3: No of Petitions as Per Visa Status.

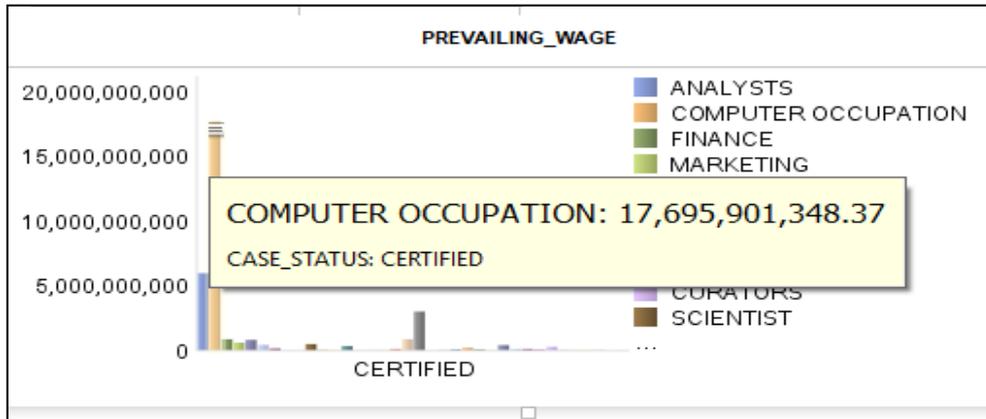


Fig. 4: Certified Cases as Per Jobs.

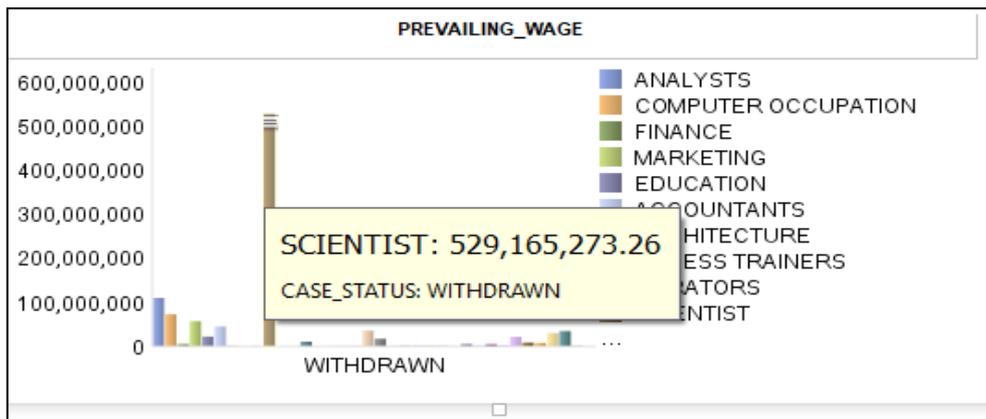


Fig. 5: Withdrawal Cases as Per Jobs.

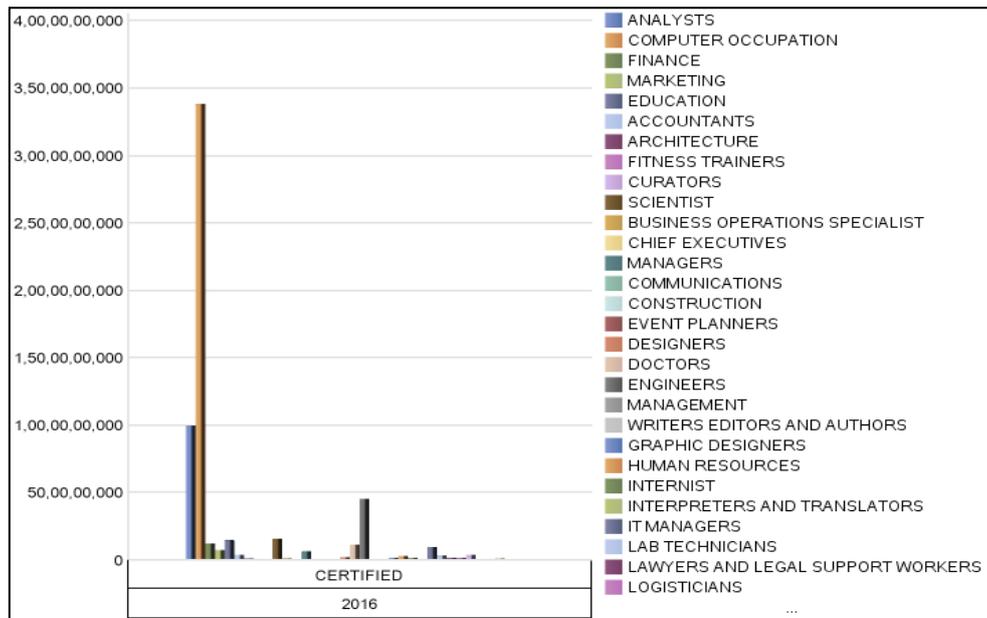


Fig. 6: Certified Cases in Year 2016.

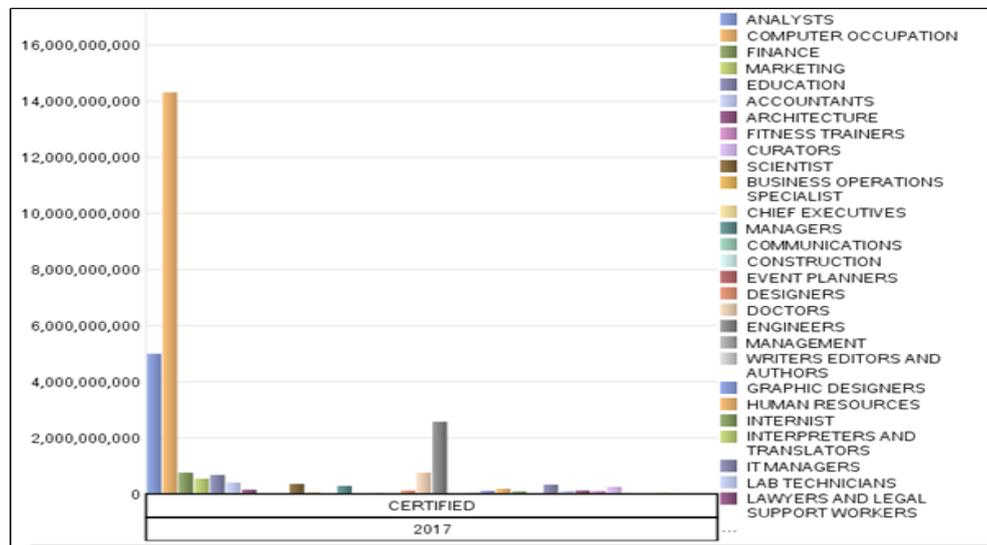


Fig. 7: Certified Cases in Year 2017.

3. Methodology

This segment contains Methodology which is shown in figure 8. In the initial step, the H1B visa dataset is taken from The Office of Foreign Labor Certification (OFLC) [3] official site with 40 attrib-

utes. It is likewise accessible on Kaggle site [8]. The second step is data cleaning and pre-handling under which initially copy esteems and missing esteems are evacuated. We used Microsoft Excel to clean adjust and standardize our dataset. The initial phase in cleaning the dataset was to evacuate accentuations to sustain a tactical

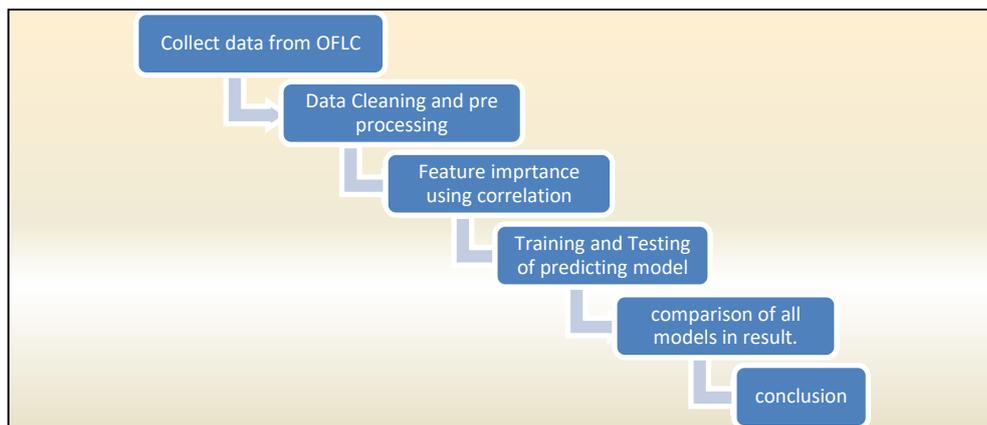


Fig. 8: Methodology Used.

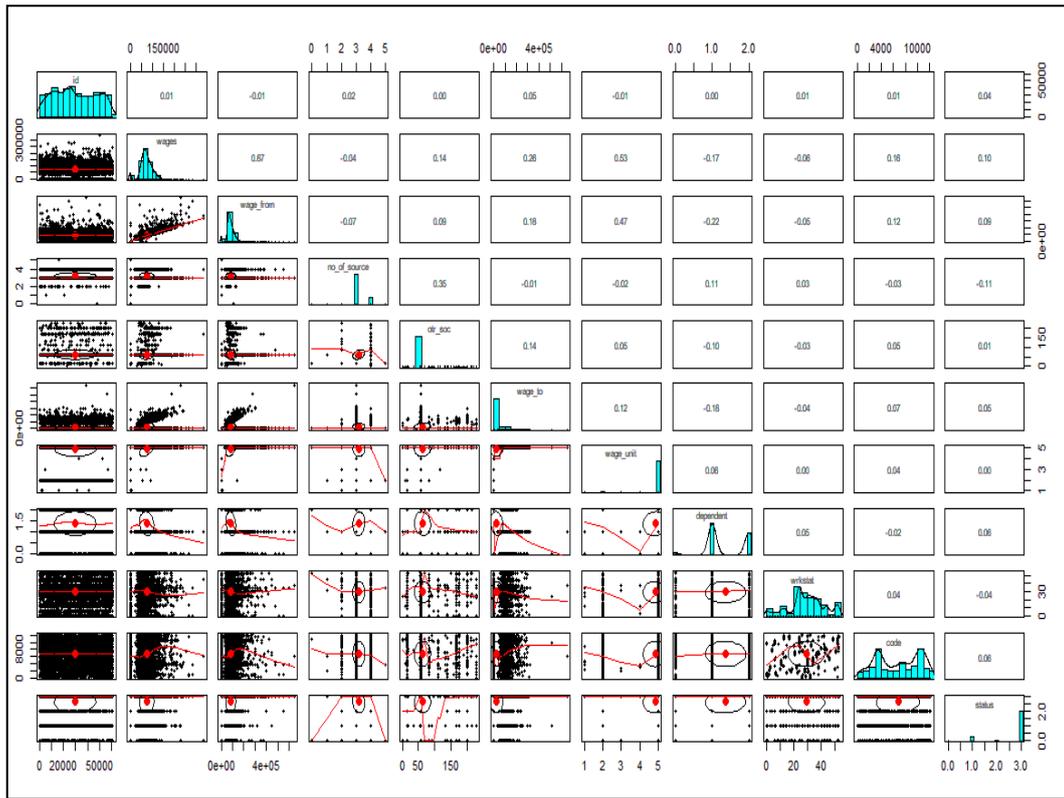


Fig. 9: Pearson Correlations between Attributes.

Table 4: Correlation between Features

	ID	Wages	Wages to	Wages from	Wage unit of pay	H1B dependent	Other sources	Will violator	Wage unit
ID	1.00	0.01	0.50	0.34	0.01	0.03	0.00	-0.34	-0.01
Wages	0.01	1.00	0.26	0.87	0.55	-0.17	0.14	-0.05	0.53
Wages to	0.50	0.26	1.00	0.18	0.20	-0.18	0.34	-0.10	0.20
Wages from	0.34	0.87	0.18	1.00	0.47	-0.22	0.09	-0.03	0.47
Wages unit of pay	0.01	0.55	0.20	0.47	1.00	0.06	0.12	0.47	0.47
H1B dependent	0.03	-0.07	-0.18	-0.22	0.06	1.00	-0.10	0.47	0.06
Other sources	0.00	0.14	0.14	0.09	0.12	0.12	1.00	0.00	-0.47
Will violator	-0.34	-0.05	-0.10	-0.03	0.47	0.47	0.00	1.00	0.76
Wage unit	-0.01	0.55	0.20	0.47	0.47	0.06	-0.47	0.76	1.00

Space from mistakes. We expelled the accompanying accentuations from our dataset: (, @ # \$) " ' : / ? ' ~) Subsequent to evacuating accentuations, we also divided the dataset and continue to our next ventures to cleaning it. Irrelevant qualities can delude the classifier into building an off-base model for foreseeing exactness. Qualities with exceptional esteem, for example, ID can assemble a model that would be based on that trait with special esteem and anticipate with most extreme exactness. It would not help us in examining the informational index. Such a characteristic is additionally called a False Predictor. The properties with one of kind esteems are Case_Number, Employer_Address, Employer_Phone, and Employer_Phone_Ext. Case_Number is a novel identifier that is doled out to each case. Employer address and telephone number are pointless traits in anticipating the Case_Status. Along these lines, we evacuated superfluous properties Case_Number, Employer_Address, Employer_Phone, and Employer_Phone_Ext on the grounds that they have special esteems. For Labor Condition Application, business time can't be over three years. In our dataset,

there are two factors Employment_Start_Date and Employment_End_Date. We computed the distinction in time between those two properties and over 90% of the distinction was 3 years for our dataset. That would make them irrelevant traits in anticipating our class variable, so they were expelled. We additionally had monotonous qualities, for example, Employer_City, Worksite_City, Work-site_County, Employer_Province, and Employer_Postal_Code that were giving a similar data. Those traits were expelled too. We kept Employer_State and Worksite_State since there are 50 states and 5 US regions, we can examine class variable in light of that. Job_Title, SOC_Code, and SOC_Name were all giving the same word related data. In this way, we expelled Job_Title and SOC_Code and kept SOC_Name. Dataset had 20 irrelevant properties that were expelled. This progression gives an aggregate of 528316 estimations of each property. The significance of highlights is checked by Pearson correlation model (as shown in figure 9) [30].

In the next step, forecasting models are prepared and tried on the date set with insights measure that are utilized for models correlations and help in anticipating best model out of the considerable number of models. In this paper, we utilize seven machine learning models. For models, R instrument is utilized which is an open source apparatus and after that unimportant esteems are evacuated and finally dataset stay with just 20 traits. The correlations between different attributes are given in table 4.

Built-in models are available in tool and are utilized straight forwardly on the dataset just by introducing required model from the R library.

3.1. Proposed multilevel ensemble model

A gathering is utilized to manage the most doubtful scenario of model forecast. In the present work, the emphasis is on the false expectation and in addition, a genuine forecast of the model and multilevel outfit model is utilized to manage false and genuine predictions. Seven models i.e. Decision tree, RF, SVM, Naive Bayes, Neural network, C5.0 and Linear Model are consolidated to show signs of improvement precision as shown in figure 10. Every one of the models is prepared on 70% of the data set and

30% is utilized for testing. The future model is separated into three stages and all stages are clarified underneath:

Stage I: The decision tree, SVM, Random Forest and neural network model is pre- pared with 70% of the dataset and create forecasts from 30% of the dataset.

Stage II: The bogus expectations of two models (decision tree and SVM) from stage I are utilized to prepare the Naive Bayes model. The bogus forecasts of two models (neural network and RF) from stage I are utilized to prepare the linear model

Stage III: The false expectations from stage II and genuine forecasts from Phase 1 are combined. This joined new dataset is utilized to prepare C5.0 model that gives last predictions.

In this approach, genuine expectations, and also false expectations, are refined to get a precise-proposed model. The reason for utilizing genuine expectation as the contribution of different models is to manage false positive outcomes (Non-antigenic is considered as antigenic). The information is gone through seven models as a result of this model superbly takes in the information to give solid and exact outcomes.

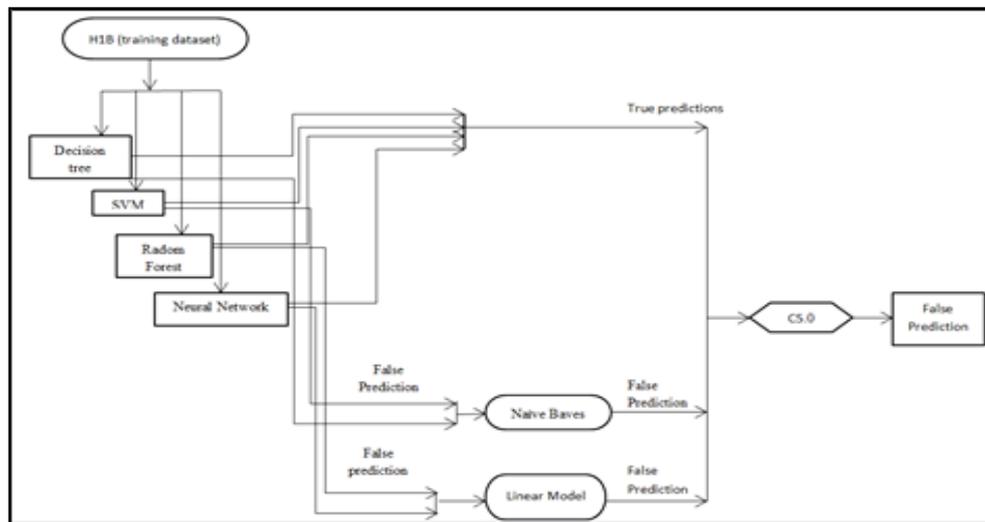


Fig. 10: Ensemble Model.

3.2. Machine learning methods

A brief of the models is exhibited underneath and table 5 shows the parameters used by these models.

- 1) Decision Trees: This model is an augmentation of C5.0 grouping calculations portrayed by Quinlan [10], [11], and [12]. Decision tree calculation is a standout between the most vital classification measures in information mining. Decision tree classifier as one sort of classifier is a stream diagram like a tree structure, where each inside hub indicates a test on a characteristic, each branch speaks to a result of the test, and each leaf hub speaks to a class. The technique that a decision tree demonstrate is utilized to group a record is to discover a way that from root to a leaf by estimating the characteristics test, and the trait on the leaf is classification result. The decision tree is the fundamental innovation utilized for classification and expectation. Decision tree learning is a regular inductive calculation in light of the case, which centers around classification rules showing as decision trees induced from a gathering of turmoil and sporadic case. In a best down the recursive way, it thinks about characteristics between inside hubs of decision tree judge the descending branches as indicated by an alternate property of the hub and reach an inference from leaf hubs in the decision tree. So from a root to a leaf hub relates to a conjunctive

manage, and the whole tree compares to a gathering of disjunctive articulation rules [31].

- 2) C5.0: C5.0 is a new decision tree algorithm developed from C4.5. The idea of construction of a decision tree in C5.0 is similar to C4.5. Keeping all the functions of C4.5, C5.0 introduces more new technologies. One of the important technologies is boosting, and another one is the construction of a cost-sensitive tree [12], [13], and [14]. Take the decision tree as a Boolean capacity. The contribution of the capacity is the protector all property of circumstance, and the yield is the "yes" or "no" decision esteem. In the decision tree, each tree hub relates to a property test, each leaf hub compares to a Boolean esteem, and each branch speaks to one of the conceivable estimations of the testing trait. The most common decision tree learning framework is ID3, which began with the idea learning framework CLS, lastly advanced into C4.5 (C5.0), which can manage consistent qualities. It is a learning guide, in light of a decision tree makes out of preparing subsets. In the event that the tree neglects to group accurately all the given preparing subset, pick other preparing subset adding to the first subset, rehash it until the point when the right decision set. To prepare various preparing occurrence classification, a decision tree which can characterize an obscure case classification in view of the particular event of property estimation sets. Utilizing the decision tree to characterize cases, you can test bit by bit the estimation of the items' properties beginning at roots, and after that go-

ing down the branch until achieving a leaf hub, in which class is the class of the protest. The decision tree is a broadly utilized technique for classification. There are different decision tree strategies, for example, ID3, C4.5, PUBLIC, CART, CN2, SLIQ, SPRINT and so forth. The most created decision tree is a variation of the center algorithm [32].

- 3) Support Vector Machine: SVM is a capable strategy for general (nonlinear) order and exceptions recognition with a natural model portrayal [15]. SVMs depend on the guideline of basic hazard minimization, which plans to limit the genuine mistake rate. SVMs work by finding a linear hyper plane that isolates the positive and negative cases with a most extreme Interclass separation or edge by taking care of the accompanying improvement issue [26], [28].

$$\text{Min } 1/2 \|w\|^2 + c \sum_{i=1}^l \xi_i \quad (1)$$

Subject to $y_i (w(x_i) + b) \geq 1 - \xi_i$, $\xi_i \geq 0$, Where ξ_i is a slack variable to represent preparing mistake, and C is a regularization parameter characterizing the exchange off between the preparation blunder and the edge. For nonlinearly distinct information, SVMs outline information to be linearly distinguishable in a high-dimensional component space utilizing bit capacities and a Lagrange multiplier α . Thus, the double of the enhancement issue progresses toward becoming.

$$\text{Max } w(\alpha) = \sum_{i=1}^l \alpha_i - 1/2 \sum_{j=1}^l (\alpha_j y_j - K(x_i, y_j))$$

Designs with nonzero α s are called SVs. The isolating hyper plane is totally characterized by the SVs, and they are the main examples that add to the grouping decision [26].

- 4) Naive Bayes: In machine learning, naive Bayes classifiers are a cluster of direct "probabilistic classifiers" in sight of relating Bayes' theory with dense (innocent) autonomy presumptions between the highlights. In the learning procedure of this classifier with the known structure, class probabilities and restrictive probabilities are computed utilizing preparing information, and after that estimation of these probabilities are utilized to arrange new perceptions [16], [17]. The discourse so far has inferred the free component model, that is, the guileless Bayes likelihood model. The guileless Bayes classifier consolidates this model with a choice run the show. One basic control is to pick the speculation that is most plausible; this is known as the greatest a posteriori or map decision run the show. The relating classifier, a Bayes classifier, is the capacity that doles out a class ($y=C_k$) mark for some k as takes after:

$$Y^* = \text{argumentmax}_p C^k; i=1 \prod_{j=1}^n (x_j^i / C^k) \quad (2)$$

- 5) Neural Network: An ANN depends on a gathering of associated units or hubs called manufactured neurons (an improved rendition of natural neurons in a creature cerebrum). Every association (an improved rendition of a neurotransmitter) between manufactured neurons can transmit a flag starting with one then onto the next. The counterfeit neuron that gets the flag can practice it and after that flag fake neurons associated with it [19], [20].
- 6) Random forest: Random forests are a group learning technique for arrangement, degeneration, and different undertakings, that work by building a large number of decision trees at preparing time and yielding the class that is the method of the classes or mean expectation (relapse) of the individual trees [22] [23]. Random decision forests modify predictions for decision trees' tendency for over fitting to their training dataset. [21]. The random forest contains a few trees and chooses random subsets of a number of various Indicators tried at each node [27].

To develop the trees, a deterministic procedure is utilized to choose each tree from a random arrangement of qualities. The hub is used to break down the hubs for diminishing the aggregate number through depiction accessible for investigation. The inspected random vector and standard random forest are blended as an estimator for each tree by utilizing comparative dissemination for all trees. It contains input vector $X = x^1; x^2; \dots; x^p$, where a p -dimensional info vector works in building a forest. Inside the forest an arrangement of k trees $T^1(x); T^2(x); \dots; T^k(x)$, the yield of each tree evaluates the genuine esteem.

$$\hat{a}Y^1 = T^1(x); \hat{a}Y^m = T^m(x), \text{ where } m = 1: k.$$

The final consequence of it is the mean of the considerable number of qualities anticipated by different trees.

$$\text{Estimate}^{RF}(X) = 1/k \sum_{k=1}^k Y_k(x). \quad (3)$$

The preparation dataset is freely taken from the information and yield the equation $D = D^1; D^2; \dots; D^n = (x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)$, where $x_i; i = 1; \dots; n$, is the preparing dataset for input vector and $y_i; i = 1; \dots; n$, is preparing the dataset for yield vector.

- 7) Linear Model: It utilizes straight models to do relapse, single stratum examination of change and examination of covariance [24]. On the off chance that the information includes vector to the classifier is a genuine vector x , at that point the yield score is:

$$y = f(w \cdot x) = f(\sum_{j=1}^n w_j x_j), \quad (4)$$

Where ω is a genuine vector of weights and f is a capacity that changes over the speck result of the two vectors into the coveted yield. The weight vector ω is gained from an arrangement of named preparing tests. Regularly f is a straightforward capacity that maps all qualities over a specific edges to the top of the line and every single other incentive to the inferior. A more mind-boggling f may give the likelihood that a thing has a place with a specific class [29].

3.3. Model valuation

There are hundreds of classifiers and depending on nature of problem different classifiers are used as performance measurement depending on nature of considered application. In the paper we do classification and for that 6 classification models are used. Dataset is distributed into preparation and investigation dataset, then target and inputs are defined. There are 1 target object i.e. CASE_STATUS and 19 input objects. The formula is same for all ML models i.e.

$$\text{CASE_STATUS} \sim \text{VISA_CLASS} + \text{EMPLOYER_NAME} + \text{EMPLOYER_STATE} + \text{EMPLOYER_COUNTRY} + \text{SOC_NAME} + \text{TOTAL_WORKERS} + \text{FULL_TIME_POSITION} + \text{PREVAILING_WAGE} + \text{PW_UNIT_OF_PAY} + \text{PW_SOURCE} + \text{PW_SOURCE_YEAR} + \text{PW_SOURCE_OTHER} + \text{WAGE_RATE_OF_PAY_FROM} + \text{WAGE_RATE_OF_PAY_TO} + \text{WAGE_UNIT_OF_PAY} + \text{H.IB_DEPENDENT} + \text{WILLFUL_VIOLATOR} + \text{WORK_SITE_STATE} + \text{WORKSITE_POSTAL_CODE}.$$

After defining formula Models are formed for seven classifiers. In this work we take accuracy and time factors comparing work of different classifiers. To decide the precision and time taken, an error or confusion network is framed demonstrating the data about genuine and predicted arrangements are done by a classifier. The diagonal components of disarray or mistake framework speak to

the quantity of items for which the anticipated mark is equivalent to the genuine name, while off-corner to corner components are those that are mislabelled by the classifier.

Table 5: Classification Models

Models	Parameters used	Library installed
Decision trees	Min spilt = 20, max depth = 30, Min Bucket =7.	rPart
C5.0	No of samples = 30000, predictors =19	C5.0
SVM	Kernel used rbfdot and Anovadot	e1071, ksvm.
Naïve Bayes	Naïve Bayes	naive Bayes
Neural Network	size=10, Linout =TRUE, max newts =10000.	neural net
Linear Model	Multinomial	Car, nnet.
Random Forest	No of trees =200	Random forest

The higher the corner to corner estimations of the matrix, better the exactness. The accuracy is calculated by using formula i.e.

$$Acc \leftarrow \text{round}(\text{mean}(\text{Actual} == \text{Predicted}) * 100, 2) \tag{5}$$

Here Actual is the actual values of the Target object and Predicted contain those Values of target object that is predicted by ML model. The models are assessed on different parameters as said in Table 5. From the outcomes, it is reasoned that the accuracy of a proposed show is expanded when contrasted with the single model

accuracy. Table 6 portrays the accuracy and other metrics of the proposed display. The accuracy has been recorded by applying 10-overlap cross validation 3 times. For cross-validation, 70% of a dataset is utilized for preparing and 30% utilized for testing. The Figure 12 portray the accuracy of proposed demonstrate 3 times in 10 runs and demonstrates the consistency in the accuracy of the proposed display.

4. Results

In present work, we prognosis the consequences of all the seven machine learning characterization models on the testing dataset. All the seven approaches are running on the parameters as appeared in Table 5. The accuracy which is computed by utilizing Eq. (5) and Other statistical measures are also present in table 6. These measures are computed by using R function mmetric which computes the classification error metrics. In this work we use metrics like Accuracy, Precision, Total Positive Rate (TPR), Total Negative Rate (TNR), Classification error (CE) and F1 score (F1). Time and accuracy are used for comparison between these seven models. Table 7 shows different results of accuracy and time on 50-50%, 60-40%, 70-30% and 80-20% partition of train and test datasets. In cross-validation, models are accomplished n number of times and accuracy is recorded if accuracy is very fluctuating then that model is over fitted/under fitted/one-sided. In exhibit work, rehashed K overlay cross-validation is utilized that portrays the consistency in the accuracy which implies proposed to demonstrate isn't influenced by these issues. For authentication of the projected show, benchmark dataset is utilized and contrast and the current model by utilizing different parameters, for example, accuracy, TNP, Precision, and F1 esteems as portrayed in Table 6. The outcome finishes up two things about the proposed show. In the main place, the proposed show is free.

Table 6: Classification Error Metrics of Models

Methods	Time (in sec)	Accuracy	Precision	CE	TPR	TNR	F1
1) Decision tree							
a) class	50.62	92.81	99.10	7.192	93.94	86.87	96.43
b) Anova	35.28	83.27	-	16.73	-	-	-
c) poisson	11.98	83.15	-	16.85	-	-	-
d) exp	13.39	83.18	-	16.82	-	-	-
2) C5.0	66.46	94.62	98.78	5.37	96.04	87.52	9.39
2) SVM (kernels							
a) rbfdot	7.92	81.84	96.48	18.16	97.94	56.89	91.04
b) anovadot	126.03	10.76	78.92	89.4	64.78	89.99	89.44
c) other kernels	Very large time	Less accuracy	---	----	-----	-----	-----
4) Naïve Bayes	119.72	75.21	92.13	24.79	80.29	46.04	85.78
5) Neural Network	113.8	88.33	66.7	11.67	97.81	88.44	80.78
6) Linear Model	169.9	90.22	95.78	9.77	90.40	91.31	94.91
7) Random Forest	162.92	93.90	91.25	6.1	89.09	69.97	79.89
8) Proposed model	158.4	95.4	94.45	4.6	88.08	79.98	95.89

Table 7: Performance Analysis of All Seven Models on Different Training and Testing Dataset Partitions

Model	Partitions of training and testing datasets 50-50% (T,A) (T,A)	60-40% (T,A) (T,A)	70-30% (T,A) (T,A)	80-30%
Decision Tree	(93.8, 92.75)	(112.1, 92.75)	(104.4, 92.75)	(82.76, 92.73)
C5.0	(114.2, 94.50)	(232.9, 94.59)	(129.9, 94.77)	(136.3, 94.71)
SVM	(143.14, 90.45)	(436.2, 90.12)	(342.6, 90.17)	(219.1, 90.05)
Naïve Bayes	(163.6, 82.88)	(145.3, 74.88)	(107.0, 84.41)	(72.39, 71.05)
Neural Network	(181.3, 88.51)	(255.6, 88.44)	(684.4, 88.01)	(1222.1, 88.8)
Random forest	(15.43, 67.12)	(15.31, 66.12)	(015.00, 67.14)	(16.67, 88.51)
Linear Model	(162.3, 93.07)	(113.05, 94.19)	(105.44, 93.21)	(104.4, 93.27)

from overfitted/under fitted/one-sided issues. Second, the result of the proposed display is enhanced when contrasted with the current procedure. Table 5 portrays the machine learning models that are prepared on the dataset with ideal tuning parameters. The dataset is separated into two sections 70% and 30%. The readied models are obscure to the 30% of the dataset. The proposed demonstrate is a blend of seven models that make it a multilevel outfit show as talked about in Section 3.2. From the outcomes, it is reasoned that the accuracy of the proposed show is expanded when contrasted with the single model accuracy. Table 8 portrays the accuracy of the proposed demonstrate. The accuracy has been recorded by applying 10-overlay cross validation 3 times. For cross-validation, 70% of the dataset is utilized for preparing and 30% utilized for testing. The Figure 12 depict the accuracy of proposed demonstrate 3 times in 10 runs and demonstrates the consistency in the accuracy of the proposed display. All results are shown in table 7 and it is clear that C5.0 has highest Accuracy Rate with less time taken (114.2, 94.50%), (232.9, 94.59%), (129.9, 94.77%) and (136.3, 94.71%) among all single models. Table 8 shows the result of 10 fold cross validation of proposed model.

5. Conclusion

H1B visa category is one of the most applied categories among other visas categories. It is designed to overcome the shortage of skilled workers in America but it affects the hiring of American workers and no. of foreign workers increased day by day. So in current work we investigate the machine learning arrangement

models with 20 properties to foresee the real H1B visa solicitors with no contribution from any external sources. Based on classification error metrics proposed model give better accuracy of 95.45 as compared single models. The projected model increases the accuracy rate as validated by 10 fold cross validations. There are chances to train data on some other classification models which may give better results. The work can be stretched out to more properties with a better relationship and other computational strategies to upgrade the execution of machine learning techniques. Some pros and cons are that we select important features by using Pearson correlation only, in future other researchers can also go for other feature selection methods which may predict better results by training new models under new conditions. The dataset and source code utilized as a part of the examination are accessible at <https://www.kaggle.com/nsharan/h-1b-visa>.

Acknowledgement

The Satisfaction and Euphoria that go with the effective culmination of an assignment would be fragmented without the specifying of the general population whose steady direction and consolation made it conceivable. I consequently entire heartedly recognize the kind undertakings' of the general people who have been associated with me particularly Mr. Harpreet Singh and Dr. Prashant Singh Rana.

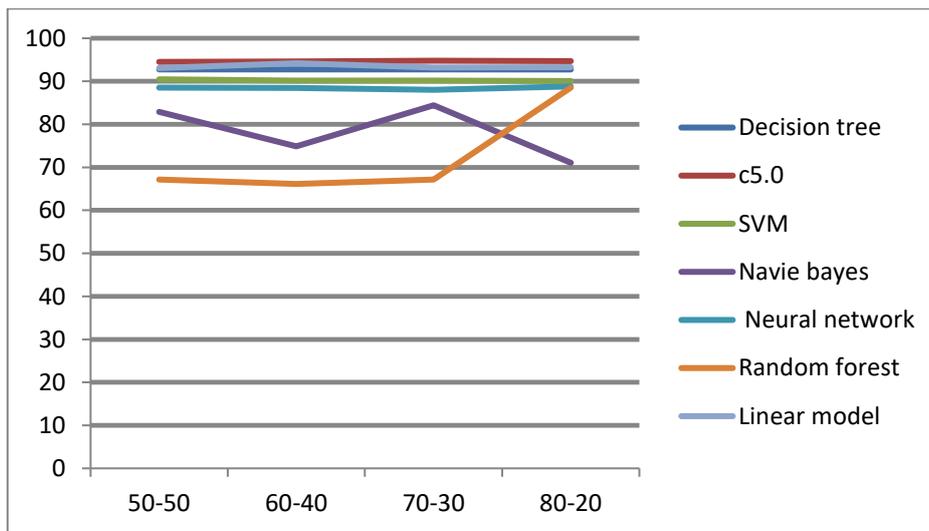


Fig. 11: Evaluation Graphs of Models on the Basis of Accuracy.

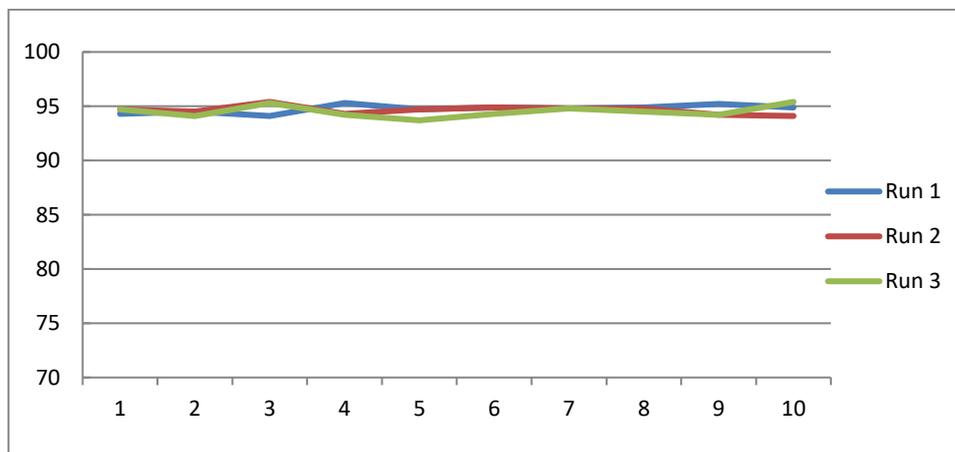


Fig. 12: K Fold Cross Validation of C5.0 Model.

Table 8: 10 Fold Cross Validation of the Proposed Model

Folds	1	2	3	4	5	6	7	8	9	10
Run 1	94.3	94.5	94.1	94.3	94.7	93.9	94.8	94.9	95.2	94.9
Run 2	94.7	94.5	95.4	94.3	94.7	94.9	94.8	94.8	94.2	94.1
Run 3	94.7	94.1	95.3	94.2	93.7	94.3	94.8	94.5	94.2	95.4

References

- [1] Dhanasekar Sundararaman , Nabarun Pal , Aashish Kumar Misraa ,(2017),” An analysis of nonimmigrant work visas in the USA using Machine Learning” , International Journal of Computer Science and Security(IJCSS), Vol. 6,
- [2] <https://www.foreignlaborcert.doleta.gov/performance/cfm>.
- [3] UNITED STATES DEPARTMENT OF LABOR. (2009, January 15). OFLC Performance Data. (www.dol.gov) Retrieved September 09, 2017, from UNITED STATES DEPARTMENT OF LABOR Employment & Training Administration:
- [4] Trim Bach, S., (2016), Giving the Market a Microphone: Solutions to the Ongoing Displacement of US Workers through the H1B Visa Program. *Nw. J. Int'l L. & Bus.*, 37, p.275.
- [5] Doran, K., Gelber, A. and Isen, A., 2014. The effects of high-skilled immigration policy on firms: Evidence from H-1B visa lotteries (No. w20668). National Bureau of Economic Research. <https://doi.org/10.3386/w20668>.
- [6] Bound, J., Khanna, G. and Morales, N., (2017). Understanding the Economic Impact of the H-1B Program on the US. In *High-Skilled Migration to the United States and its Economic Consequences*. University of Chicago Press.
- [7] Mithas, S. and Lucas Jr, H.C., 2010. Are foreign IT workers cheaper? US visa policies and compensation of information technology professionals. *Management Science*, 56(5), pp.745-765. <https://doi.org/10.1287/mnsc.1100.1149>.
- [8] Kaggle H-1B dataset, <https://www.kaggle.com/nsharan/h-1b-visa> 13 outsourcing companies took nearly one-third of all H-1B visas in 2014, <https://www.nytimes.com/interactive/2015/11/06/us/outourcing-companies-dominate-h1b-visas.html>.
- [9] Disney 'forced 250 of its American IT workers to train up the Indian workers who replaced them', <http://www.dailymail.co.uk/news/article-4037392/Disney-fired-250-American-workers-replaced-Indian-staff-visas-suit-says.html>. 2016
- [10] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V. and Vanderplas, J., Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12(Oct), pp.2825-2830, 2011
- [11] Qing-yun Dai, Chun-ping Zhang, Hao Wu, “Research of Decision Tree Classification Algorithm in Data Mining”, Vol.9, No.5 (2016), pp.1-8
- [12] J.Ross Quinlan, “C4.5: Programs for machine learning”, Elsevier, 2014.
- [13] Sohag Sundar Nanda, Soumya Mishra, Sanghamitra Mohan-ty, Oriya Language Text Mining Using C5.0 Algorithm, (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 2 (1) , 2011
- [14] PANG Su-lin, GONG Ji-zhang C5.0 Classification Algorithm and Application on Individual Credit Evaluation of Banks; Volume 29, Issue 12, February 2009.
- [15] Sona Taheri, Musa Mammadov, Learning the Naive Bayes Classifier with Optimization Models, Vol. 23, No. 4, 787–795, 2013.
- [16] Chang, C. and Lin, C. (2001). LIBSVM: A library for support vector machines, <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [17] S.S. Keerthi and E.G. Gilbert. Convergence of a generalized SMO algorithm for SVM classifier design. *Machine Learning*, 46(1):351–360, 2002. <https://doi.org/10.1023/A:1012431217818>.
- [18] CTomM.Mitchel, McGrawHil, Decision Tree Learning, Lecture slides for textbook Machine Learning,, 197
- [19] JürgenSchmidhuber, “Deep learning in neural networks: An overview”, Elsevier, Volume 61, January 2015, Pages 85-117.
- [20] Gao Huang, Guang-Bin Huang, Shiji Song, Keyou You, “Trends in extreme learning machines: A review”, Elsevier, Volume 61, January 2015, Pages 32-48.
- [21] Andy Liaw and Matthew Wiener,” Classification and Regression by random Forest”. *R News*, 2(3):18–22, 2002.
- [22] Arun Pretorius, Surette Bierman, Sarel J.steel (2017.).” A meta-analysis of research in random forests for classification” IEEE Conference, 16 January.
- [23] Eesha Goel, Er. Abhilasha, “Random Forest: A Review”, *ijarcsse*, Volume 7, Issue 1, January 2017.
- [24] L. Breiman, “Random Forest”. October 2001, Volume 45, Issue 1, pp. 5–32 45
- [25] Christian Heumann, Michael Schomaker, Shalabh (2016), “Introduction of statistics and data analysis”, Springer.
- [26] C. Burges, (1998.), “A Tutorial on Support Vector Machines for Pattern Recognition,” *Data Mining and Knowledge Discovery*, Kluwer Academic Publishers.
- [27] A. Liaw, M. Wiener, “Classification and regression by random forest”, *R news* 2 (3), (2002) 18–22.
- [28] Nahla H. Barakat and Andrew P. Bradley,”Rule Extraction from Support Vector Machines: A Sequential Covering Approach”, VOL. 19, NO. 6, JUNE 2007.
- [29] Guo-Xun Yuan; Chia-Hua Ho; Chih-Jen Lin "Recent Advances of Large-Scale Linear Classification". *IEEE*. 100, (2012).
- [30] C. K. Williams, A. Engelhardt, T. Cooper, Z. Mayer, A. Ziem, L. Scrucca, Y. Tang, C. Candan, M. M. Kuhn, Package caret.
- [31] V. W. Aalst, “Exterminating the Dynamic Change Bug: A Concrete Approach to Support Workflow Change”, Eindhoven, UK: Eindhoven University of Technology, (2000).
- [32] S. Rinderle, M. Reichert and P. Dadam, “Correctness Criteria for Dynamic Changes in Workflow Systems”, *Data & Knowledge Engineering*, vol. 50, no. 1, pp. 9-34. <https://doi.org/10.1016/j.datak.2004.01.002>.