

# Developing an Efficacious System of Prediction of Loan Eligibility by Combining Artificial Intelligence (AI) and Data Analysis Tools and Techniques

Diksha Choudhary

Indraprastha College for Women  
Civil Lines, New Delhi

## ABSTRACT

*Banks serve the necessities of everybody close to emergency clinics and schools. Individuals contact banks for different purposes. Yet, one of the most widely recognized administrations advertised by banks is credit. In any case, numerous standard individuals should know the financial strategies and advance qualification measures.*

*This study intends to adoptive an AI (ML) model that can foresee whether an individual is qualified for a health loan by examining some essential qualities the client enters. This interaction gathers a dataset of all vital boundaries for a credit application from Kaggle. The collected dataset is then pre-processed using the invalid worth disposal strategy and encoding. At the same time, three ML models were created utilizing three distinct algorithms. They are the Random Forest, Naive Bayes NB, and LR. The pre-processed information will next be used to prepare the models. A couple of boundaries will be contrasted to evaluate the models' adequacy. The examination results demonstrate that the RF algorithm is the best concerning precision and error. The accuracy of the RF algorithm is 91%. What's more, it similarly predicts advanced qualification with lesser error values. The LR model has less accuracy and important error values, making it the most un-proficient algorithm for loan prediction.*

## INTRODUCTION

Banks are one of the most inescapable spots for individuals of all monetary situations. Banks offer different administrations like reserve funds, loans, safe wallets, etc. Out of the relative multitude of advantages, loan assumes a considerable part in banking as it is helpful to individuals of different areas. Various advances include individual, home, health, gem, etc.

Nonetheless, the issues with the financial methods are drawn out and sometimes need clarification for regular users. This might only partially influence different sorts of loans, yet the candidate might confront severe medical problems for a health loan. Assuming the candidate is qualified for that credit, it is ideal that they get it. Because, if not, it prompts pointless exercise in futility and the candidate's health. To avoid such a situation, the candidate should know about all the necessities of the health credit of a specific bank, which is nearly unimaginable, as a rule. To overcome this issue, this study intends to see the best ML algorithm that can be used to predict whether a specific individual is qualified for a health loan or not. This is done by preparing and dissecting different algorithms regarding accuracy and error prediction. The cycle is made sense of obviously in the forthcoming sections.

## MATERIALS AND STRATEGIES

The information comprising fundamental bank subtleties in a table is gathered from Kaggle. The table then, at that point, goes through different cycles. The cycles are pictorially addressed in Fig. 1.

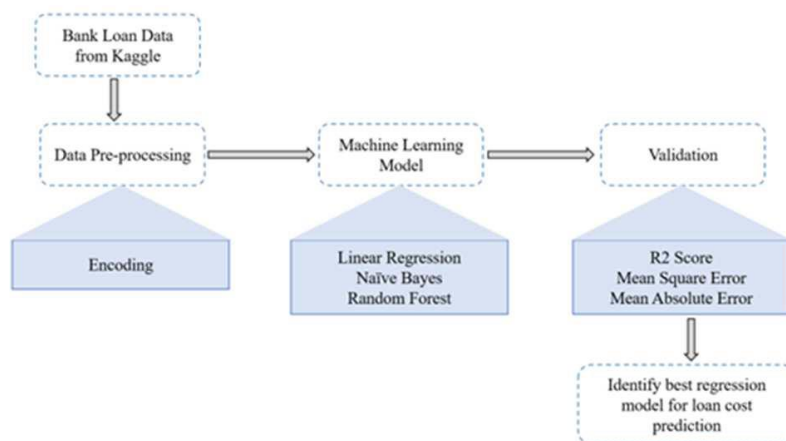


Fig. 1. Workflow of the study

The information gathered from Kaggle is pre-processed to ensure better execution of the ML models. Three ML models were created utilizing three unique algorithms. They are LR, NB, and the RF algorithm. The models will then, at that point, be prepared to use the pre-processed information. The proficiency of the models will then be investigated by contrasting a couple of boundaries.

The best algorithm of the three is picked in light of the research results.

### INFORMATION MINING AND PREPROCESSING

This part makes sense of the information gathered and the techniques utilized to deal with the information.

#### A. Information Mining

A dataset containing all the fundamental information about banking necessities for applying for home credit is gathered from Kaggle. This dataset will be as a table. It has age, sex, BMI, kids, smoker, locale, and charges. An example of the gathered dataset is displayed in Fig. 2. The total dataset comprised of 1,000 300 38 lines of information with every one of the segments filled. Nonetheless, the dataset should be pre-processed before preparing the ML models with the dataset. The preprocessing strategies incorporate invalid worth end and encoding.

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southeast	1725.55230
2	28	male	33.000	3	no	southeast	4449.46200
3	33	male	22.705	0	no	northwest	21984.47061
4	32	male	28.880	0	no	northwest	3866.85520

Fig. 2. Sample dataset

#### B. Elimination of Null Values

Sometimes, a dataset may contain invalid entries in some sections. These invalid qualities will bring about fake forecasts [11]. Consequently, the dataset must be dissected for weak attributes, and if a unit comprises multiple invalid grades, the segment can be removed. Be that as it may, this invalid worth examination won't consider '0' is a null

worth as a candidate might not have children. At the point when a value is missing, it is viewed as weak. The consequences of the invalid worth end are displayed in the figure 3.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  -
0   age          1338 non-null   int64
1   sex          1338 non-null   object
2   bmi          1338 non-null   float64
3   children     1338 non-null   int64
4   smoker       1338 non-null   object
5   region       1338 non-null   object
6   charges     1338 non-null   float64
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
```

Fig. 3. Column details after null value elimination

From Fig. 3, no cells contain invalid values, i.e., every dataset segment incorporates a particular esteem. In this way, all parts are held together.

**Encoding**

Encoding is a cycle where values in sections with extremely low cardinality, like three and four, are changed into numbers for more straightforward handling [12]. For instance, by and large, an individual can be a smoker or not. It implies that the cardinality of the smoker segment is only two. Here, if they are a smoker, the worth is changed to 0; if the individual is a nonsmoker, the value is changed to 1. This change is made for three of the eight segments, and they are sex, smoker, and district. The qualities changed are displayed in Table 1. From Table 1, the rates are changed from strings to mathematical attributes. Here, the community has the most elevated cardinality of four.

TABLE 1. ENCODED DATA

Sex	Female	0
	Male	1
Smoker	No	0
	Yes	1
Region	North East	0
	North West	1
	South East	2
	South West	3

**DEVELOPMENT OF ML MODELS**

The ML model utilized for bank credit prediction is defined below:

**A. Linear Regression**

The LR algorithm is an ML procedure taken on from the statistics field. This algorithm, by and large, predicts the result esteem by investigating input values. Two distinct factors, named the ward and free factors, assume a considerable part in the expectation by this algorithm [13]. The models created by this algorithm are of two kinds. Basic LR model and Various LR models.

The quantity of reliant and free factors decides the sort. It works by finding the best line which can be fitted given the reliant factors and seeing where the free factor will be available. One of the critical benefits of the LR algorithm is its effortlessness. It is more straightforward to carry out than some other ML algorithm and can foresee the result with only a few information factors.

### B. Naive Bayes

There are two kinds of ML algorithms - directed and unaided. The NB algorithm goes under the classification of a regulated ML algorithm. This was planned given the hypothesis named Bayes hypothesis. The NB algorithm is primarily utilized in applications that require separating or grouping algorithms. The Bayes hypothesis clarifies how to decide the value of an obscure variable by investigating the known factors. As this algorithm depends on the Bayes hypothesis, it works in basically the same manner. It utilizes various probabilities: the Back, Probability, Earlier, and Peripheral [14]. During forecast characterization, this algorithm works in three stages. The dataset will initially be changed over into tables known as recurrence tables. By tracking down the probabilities of the recurrence table, another table is named the probability table.

Utilizing the upsides of the probability tables, the back likelihood of the dataset is anticipated using the Bayes hypothesis.

### C. Random Forest

The RF technique is a clustering algorithm for the hardware. A few decision trees are used to decide the best outcome. Various groups of evaluating and regression trees prepared on informational collections that are equivalent to the training sets are the groundwork of the RF approach [15].

Due to its superior or predominant exhibition across various issues, including regression, arrangement, and prescient demonstrating, this strategy is undoubtedly the most well-known and utilized ML. The RF algorithm's sacking process is one of its fundamental advantages. Since each decision tree is fitted to a somewhat particular informational index and displays some variety in execution, stowing is a legitimate gathering strategy. Compared to regular model decision trees, trees are undiscovered and overfit to the preparation dataset. The outfit utilizes the tree. Congestion may likewise influence the dynamic trees. By making arbitrary subsets of traits and more modest trees using these gatherings, the RF regularly keeps away from this.

## RESULT AND CONVERSATION

A dataset with data about a credit candidate's fundamental boundaries is gathered from Kaggle. This dataset then uses two or three strategies to guarantee an accurate forecast.

The systems are invalid and data encoding.

TABLE 2. PERFORMANCE OF THE ALGORITHMMS

Regression Model	MSE	MAE	R2
LR	5800.763	4032.558884007944	0.824
NB	5123.142	3762.164	0.887
RF	4921.26	3537.794	0.911

Three ML models are created utilizing the LR, NB, and the RF algorithm. The models are then prepared using the pre-processed dataset. In the wake of training, a couple of boundaries are broken down to track down the best algorithm.

The limitations are the Mean Square Error (MSE), Mean Absolute Error (MAE) and Precision (R2) scores. MSE is a contraction of Mean Squared Error. This worth is utilized to figure out considerable mistakes in forecasts. The MAE estimation is the mean contrast between the genuine and anticipated values, while the R2 score is the reliant variable. Table 2 comprises the MAE, MSE, and R2 scores of all the e-algorithms [16].

The LR algorithm individually has the most elevated MSE and MAE values of 5800 and 4032. The second-biggest upside of MAE and MSE is the NB algorithm. The MSE of this algorithm is 5123, and the MAE esteem is 3762. The RF algorithm is the best, with the most minor MSE worth of 4921. This MAE is significantly less than the most noteworthy MAE value [17]. The MSE of the RF algorithm is additionally lesser than the other two. The MSE of this algorithm is 3537. Concerning the R2 score, the RF algorithm has the most noteworthy worth of 0.911, and the LR algorithm has the lowest value of 0.824.

The error investigation of every one of the three algorithms is graphically made sense of in Fig. 4.

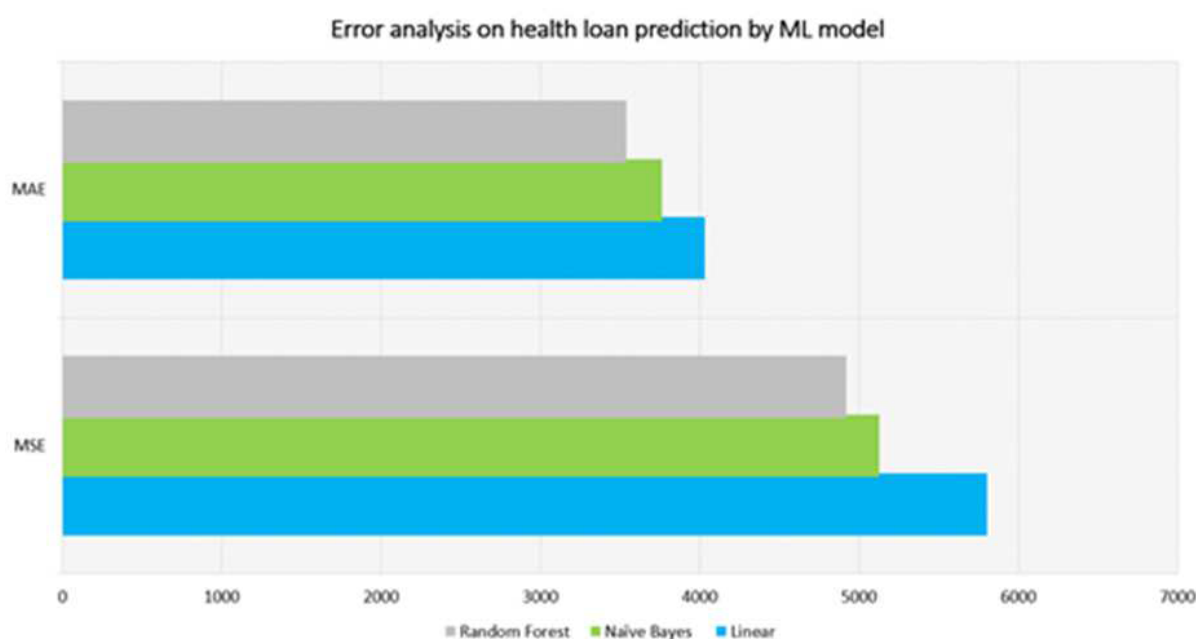


Figure 4. Error analysis

Crops yield must be shielded from pesticides and irresistible. Harm to the harvests must look at the use of quality appraisal. Temperature variety should be assessed utilizing soil clamminess to recognize the grasshopper attack.

Grasshopper plagues happen in vacant regions where wet soil is present. Bug synthetic substances cause crop harm, which is hurtful to the yields. These preventive measures must decrease food shortage. The AI approach connects with the dataset, which has more secret layers of neurons in light of classification [18]. A protection measure for the insect plagues must be started by SVM, Rf and LR techniques. These strategies are utilized to compute the derivate of examination. AI has a subset, a deep learning approach that uses more than adequate information space with intricacy in the algorithm. A few algorithms for the dataset near investigation are Backing vector machines, Irregular Woods Algorithm and Strategic Relapse. The CNN model is expressed to distinguish which conspicuous insects utilize their field state pictures.

Intently checking the problem areas previously impacted by the beetle swarm.

The graphical investigation again demonstrates that the RF algorithm has the least mistake values. However, the proficiency of an ML model can't be resolved exclusively by investigating the mistake values. The ML model's

expectation accuracy is just as imperative as the mistake values. The accuracy worth of each of the three algorithms is displayed as a 3D diagram in Figure 5.

From Fig. 5, the precision of the RF algorithm is the most noteworthy among the three algorithms. The LR algorithm had the most minimal precision value [19].

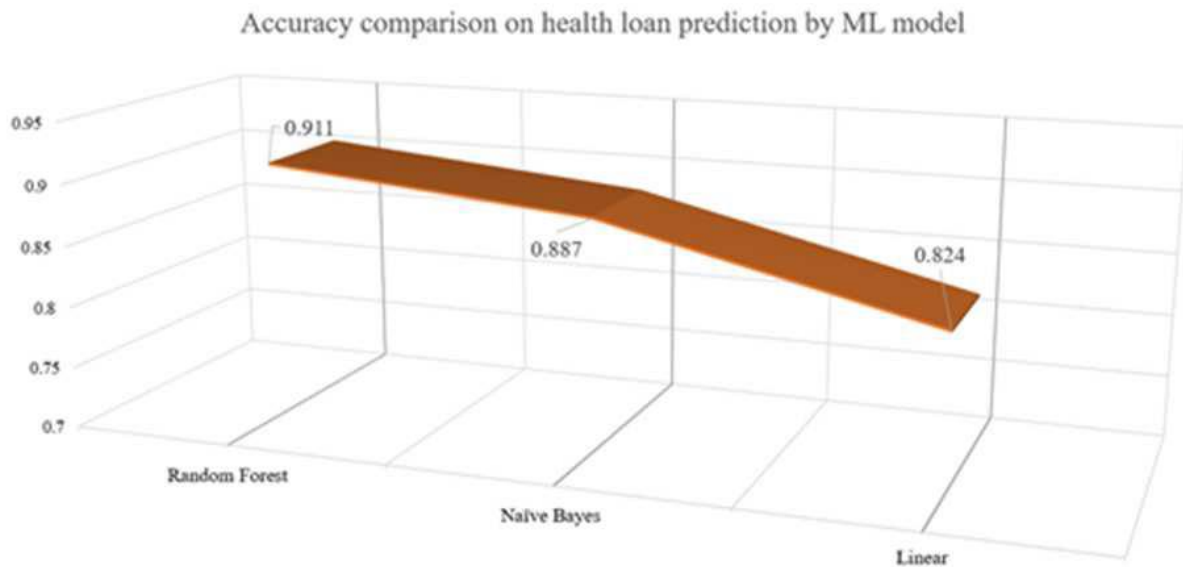


Fig. 5. Accuracy analysis results

## CONCLUSION

From Kaggle, a dataset including insights regarding a credit candidate's central attributes is accumulated. Then, at that point, to ensure precise expectations, this dataset goes through barely any cycles. Information encoding and invalid value disposal are the techniques. Three algorithms — LR, NB, and RF are utilized to make three ML models. The pre-processed dataset is accordingly used to prepare the models. A couple of boundaries are inspected following preparation to decide the best algorithm. The MSE, MAE, and R2 scores are the boundaries. After examination, it was tracked down that the RF algorithm had the most minor mistake values and the most elevated precision esteem, making it the best ML algorithm to foresee qualification for health loans. While the RF algorithm is the best, the NB algorithm is the subsequent best, and the LR algorithm is the most awful among the three algorithms. The best algorithm, i.e., the RF algorithm, can be utilized daily by sending it to a site or a product application. Along these lines, it tends to be used by individuals with medical problems and incapability to predict, regardless of whether they are qualified to apply for a health credit. Thus, the time wasted looking out for the bank premises to realize that the individual is ineligible for a health loan can be forestalled.

## REFERENCES

1. Y. Xie, Y. Li, Z. Xia and R. Yan, "An Improved Forward Regression Variable Selection Algorithm for High-Dimensional Linear Regression Models," in *IEEE Access*, vol. 8, pp. 129032-129042, 2020, doi: 10.1109/ACCESS.2020.3009377.
2. G. Qiu, X. Gui and Y. Zhao, "Privacy-Preserving Linear Regression on Distributed Data by Homomorphic Encryption and Data Masking," in *IEEE Access*, vol. 8, pp. 107601-107613, 2020, doi: 10.1109/ACCESS.2020.3000764.
3. W. Deng, Y. Guo, J. Liu, Y. Li, D. Liu and L. Zhu, "A missing power data filling method based on improved random forest algorithm," in *Chinese Journal of Electrical Engineering*, vol. 5, no. 4, pp. 33-39, Dec. 2019, doi: 10.23919/CJEE.2019.000025.

4. F. Dufrenois and D. Hamad, "Sparse and online null proximal discriminant analysis for one class learning in large- scale datasets" 2019 International Joint Conference on Neural Networks (IJCNN), 2019, pp. 1-8, doi: 10.1109/IJCNN.2019.8852312.
5. V. Chouhan and S. K. Peddoju, "Investigation of Optimal Data Encoding Parameters Based on User Preference for Cloud Storage," in IEEE Access, vol. 8, pp. 75105-75118, 2020, doi: 10.1109/ACCESS.2020.2987999.
6. Ghosh, S., Rana, A., & Kansal, V. A benchmarking framework using nonlinear manifold detection techniques for software defect prediction. International Journal of Computational Science and Engineering, 2020, 21(4), 593-614.
7. Rana, A., & Sharma, S., Mechanism of sphingosine-1-phosphate induced cardioprotection against I/R injury in diabetic rat heart: Possible involvement of glycogen synthase kinase 3 $\beta$  and mitochondrial permeability transition pore. Clinical and Experimental Pharmacology and Physiology, 2016, 43(2), 166-173.
8. Kunwar, V., Agarwal, N., Rana, A., Pandey, J.P. (2018). Load Balancing in Cloud—A Systematic Review. In: Aggarwal, V., Bhatnagar, V., Mishra, D. (eds) Big Data Analytics. Advances in Intelligent Systems and Computing, vol 654. Springer, Singapore. [https://doi.org/10.1007/978-981-10-6620-7\\_56](https://doi.org/10.1007/978-981-10-6620-7_56)
9. M. S. Raghavendra, P. Chawla and A. Rana, "A Survey of Optimization Algorithms for Fog Computing Service Placement," 2020 8th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO), Noida, India, 2020, pp. 259-262, doi: 10.1109/ICRITO48877.2020.9197885.