

Chapter 21

Predictive Analytics Model of an Engineering and Technology Campus Placement



Sachin Bhoite, Anuradha Kanade, Punam Nikam, and Deepali Sonawane

1 Introduction

According to the requirement of industry, colleges must update their curriculum and provide necessary technical and practical knowledge to the students. It will help in fulfilling the requirement of skilled and qualified students of the industries. DM and machine learning (ML) scholars have studied classification problems most recurrently [1]. In which the value of a dependent variable can be predicted based on the values of other independent variables [2]. This paper aims to determine the features impacting on prediction of placement and also students will get to know the placement status and get help in improving their weaker areas in advance.

Basically, this model will help to make training and placement officers (TPO) work easy and increment the total number of placements. Hence, it will directly lead to an increment in the rank of engineering and technology institutions. As our objective is to predict the placement of a student, in such a way that either he will get placement or not. It is a binary classification problem. To get good accuracy with minimum error, we have experimented with various classification ML algorithms with K-fold cross-validation techniques and trained and tested the data splitting techniques. The Value of K is tested for better results though most of the time it has

S. Bhoite (✉) · A. Kanade · P. Nikam · D. Sonawane
School of Computer Science, MIT-WPU, Pune, Maharashtra, India
e-mail: sachin.bhoite@mitwpu.edu.in

A. Kanade
e-mail: anuradha.kanade@mitwpu.edu.in

P. Nikam
e-mail: punam.nikam@mitwpu.edu.in

D. Sonawane
e-mail: deepali.sonawane@mitwpu.edu.in

considered as 10. Also, we used EL techniques, which are comparatively faster and give better accuracy for classification projects.

2 Related Work

The researchers have studied several connected national and international research papers, thesis to understand datasets, data pre-processing methods, features selection methods, type of algorithms used in the existing studies.

Authors in [3] performed a step-wise analysis based on specific statistical frameworks for the placement. The analysis concluded with student datasets including academic and selection subtleties is important for forecasting future selection possibilities. Authors in [4] proposed the campus placement prediction work using the classification algorithms Decision Tree and Random Forest. The accuracy obtained after analysis for Random Forest is greater than the Decision tree. Authors in [5] used different ML algorithms to analyze students' admission preferences. They found Random Forest classifier is a good classifier as its accuracy is very high. Authors in [6] used different ML models to analyze students' placement, they found AdaBoost classifier along with the Bagging and Decision Tree as Base Classifier gives high accuracy. The student placement analyzer recommendation system, built using classification rules-Naïve Bayes, Fuzzy C Means techniques, to predict the placement status of the student to one of the five categories, viz., Dream Company, Core Company, Mass Recruiters, Not Eligible, and Not Interested in Placement. This model helps weaker students and provides extra care toward improving their performance henceforth [7]. Authors in [8] presented student career prediction using advanced ML techniques. In this paper, Advanced ML algorithms like SVM, Random Forest decision tree, One Hot Encoding, XG boost are used. Out of all, SVM gave more accuracy with 90.3%, and then the XG Boost with 88.33% accuracy.

Authors in [9] presented student placement and skill ranking predictors for programming classes using class attitude, psychological scales, and code metrics. They used Support Vector Machine with RBF Kernel (SVM), Support Vector Machine with Linear Kernel (SVML), Logistic regression (LR), Decision tree (DT), Random Forest (RF) techniques. ML is used to predict placement results and the programming skill level. The researcher created a classification model with precision, recall, and F-measure.

Authors in [10] presented the study on educational data mining for student placement prediction using ML algorithms. ML algorithms are applied in the weka tool and R studio which are J48, Naïve Bayes, Random Forest, Random Tree, Multiple Linear Regression, binomial logistic regression, Recursive Partitioning, Regression Tree, conditional inference tree, Neural Network. In the weka tool, Random Forest and Random Tree algorithms are giving 100% accuracy on the student placement dataset. Authors in [11] presented a survey on placement prediction systems using

ML. The author has suggested ensemble methods, which is a Machine Learning technique that combines several base models in order to produce one optimal predictive model.

3 Research Methodology

The proposed work was carried out by performing experiments on the pass-out student's dataset with various ML algorithms.

3.1 Algorithms Used

The objective of research needs to use classification methods. Hence, researchers have used the following ML classification algorithms.

1. Logistic Regression
2. K-Nearest Neighbors
3. Decision Tree
4. Random Forest
5. Support Vector Machine
6. Naive Bays

Also, used following advanced EL algorithms.

7. Adaptive Boosting,
8. Extreme Gradient Boosting (XGBoost) and
9. Grid Search CV

4 Steps in Building Predictive Models Using ML

We followed the Cross-Industry Standard Process (CRISP) methodology.

Understanding of problem and objectives of the research: Understanding dataset of already placed students and selection of the appropriate features for placement prediction.

Data Understanding: Data of already placed students were collected. All the attributes of the dataset were analyzed based on their importance and relevance based on the placement prediction. Point 5, About the dataset of this topic explains details about the dataset.

Feature Engineering: In this phase, the data from multiple data sources were integrated into one dataset. The next step is that the data were cleaned by removing unwanted columns, handling missing values, creating unique classes, performing transformation for numerical data, and all the cleaning activities on the data. Point 6, Feature engineering of this topic explains details about the same.

Table 1 Univariate feature selection for placement prediction

| Feature name | Feature score | Feature name | Feature score |
|-------------------------|---------------|-----------------------------|---------------|
| Sem_IV_Aggregate Marks | 311.517737 | Sem_VI_Pending Back Papers | 22.920021 |
| Aggregate Present Marks | 223.533006 | Sem_V_Pending Back Papers | 20.066178 |
| Sem_III_Aggregate Marks | 198.255768 | Sem_III_Back Papers | 16.854853 |
| Sem_VI_Aggregate Marks | 151.002450 | Back Papers | 12.203902 |
| Sem_V_Aggregate Marks | 147.823502 | Pending Back Papers | 7.822886 |
| Sem_II_Aggregate Marks | 142.692722 | Sem_I_Back Papers | 5.458014 |
| Sem_I_Aggregate Marks | 138.860021 | Sem_II_Back Papers | 2.265504 |
| College Name | 55.624319 | Sem_II_Pending Back Papers | 0.701740 |
| Sem_VI_Back Papers | 53.893101 | Sem_IV_Pending Back Papers | 0.558951 |
| SSC Aggregate Marks | 42.917186 | Sem_III_Pending Back Papers | 0.200665 |
| Defense Type | 27.490582 | Sem_I_Pending Back Papers | 0.124713 |
| Category | 27.385665 | Gender | 12.518480 |
| 12th/Diploma marks | 26.351629 | Branch | 0.103342 |

Experimenting: A number of ML algorithms were tested and experimented with parameter tuning mentioned in Table 2 and 3 to predict the college, and its results are discussed in point 9, result and discussion.

Evaluation: Models developed were evaluated based on their performance for accuracy metric [13]. More information is presented in point 8.

Result and Discussion: Result and discussion are discussed in point 9.

Implementation: Once the model is evaluated, it is used to evaluate unseen data, which is discussed in point 10.

5 About the Dataset

Researchers have collected 16 engineering colleges' 9766 data records. A number of columns in the dataset per college was varied from 20 to 46. We merged the dataset by considering common and important columns from our objective point of view in the excel file format, and then converted it into a CSV file. Which is essential to read by the python code to implement ML algorithms.

6 Feature Engineering

In general, every ML algorithm takes some input data to generate desired outputs. These input data are called features, which are usually presented in structured columns. As per goal or objective algorithms require input features with some specific

Table 2 List of experiments with model combinations

| Sr No | Name of Algorithm | Data splitting method used | Datasplitting folds/ratio | | | Parameter tuned | No. of parameter Tested |
|-------|------------------------------|----------------------------|---------------------------|-------|-------|-----------------------|-------------------------|
| | | | | | | | |
| 1 | Logistic Regression | K-FCV | 3 | 5 | 10 | label encoding | 6-10 |
| | | | | | | onehot encoding | 6-10 |
| | | T-TS | 70:30 | 80:20 | 90:10 | label encoding | 6-10 |
| | | | | | | onehot encoding | 6 to 10 |
| 2 | Support Vector Machine (SVC) | K-FCV | 3 | 5 | 10 | estimator | 6-10 |
| | | | | | | param_grid | 6-10 |
| | | T-TS | 70:30 | 80:20 | 90:10 | estimator | 6-10 |
| | | | | | | param_grid | 6-10 |
| 3 | Decision Tree | K-FCV | 3 | 5 | 10 | max_depth | 6-10 |
| | | | | | | min_impurity_decrease | 6-10 |
| | | | | | | max_leaf_nodes | 6-10 |
| | | | | | | min_leaf_nodes | 6-10 |
| | | | | | | max_features | 6-10 |
| | | T-TS | 70:30 | 80:20 | 90:10 | max_depth | 6-10 |
| | | | | | | min_impurity_decrease | 6-10 |
| | | | | | | max_leaf_nodes | 6-10 |
| | | | | | | min_leaf_nodes | 6-10 |
| | | | | | | max_features | 6-10 |
| 4 | Random Forest | K-FCV | 3 | 5 | 10 | max_depth | 6-10 |
| | | | | | | min_impurity_decrease | 6-10 |
| | | | | | | max_leaf_nodes | 6-10 |
| | | | | | | min_leaf_nodes | 6-10 |
| | | | | | | max_features | 6-10 |
| | | T-TS | 70:30 | 80:20 | 90:10 | max_depth | 6-10 |
| | | | | | | min_impurity_decrease | 6-10 |
| | | | | | | | |

(continued)

Table 2 (continued)

| Sr No | Name of Algorithm | Data splitting method used | Datasplitting folds/ratio | | | Parameter tuned | No. of parameter Tested |
|-------|------------------------|----------------------------|---------------------------|-------|-------|-----------------|-------------------------|
| | | | | | | max_leaf_nodes | 6-10 |
| | | | | | | min_leaf_nodes | 6-10 |
| | | | | | | max_features | 6-10 |
| 5 | Gaussian NB | K-FCV | 3 | 5 | 10 | | 6-10 |
| | | T-TS | 70:30 | 80:20 | 90:10 | | 6-10 |
| 6 | K Neighbors Classifier | K-FCV | 3 | 5 | 10 | leaf_size | 6-10 |
| | | | | | | n_neighbors | 6-10 |
| | | T-TS | 70:30 | 80:20 | 90:10 | leaf_size | 6-10 |
| | | | | | | n_neighbors | 6-10 |

Table 3 List of experiments with advanced algorithms

| Sr. No | Name of the Algorithm | Data splitting method used |
|--------|--|----------------------------|
| 1 | Ada Boost Classifier (DT) | T-TS |
| 2 | Extreme Gradient Boosting (XGBoost) Classifier | T-TS |
| 3 | Grid Search CV | T-TS |

characteristic to get the desired output. Hence, there is a need of feature engineering. Feature engineering efforts mainly have two goals:

1. Generating the proper input dataset, as per the requirement of the ML algorithm.
2. Improving the performance of ML models.

As per the experience of the researcher, we need to spend more than 70% of the time on data preparation. The following steps are carried out to achieve the same.

1. Missing Values
2. Handling categorical data (Label Encoder)
3. Change the data type
4. Drop columns

7 Feature Selection

Every time domain experts may not be available to decide independent features to predict the category of the target feature. Hence, before fitting model, we must make sure that all the features that we have selected are contributing to the model properly and weights assigned to it are good enough so that our model gives satisfactory accuracy. For that, we have used 3 feature selection techniques: Univariate Selection,

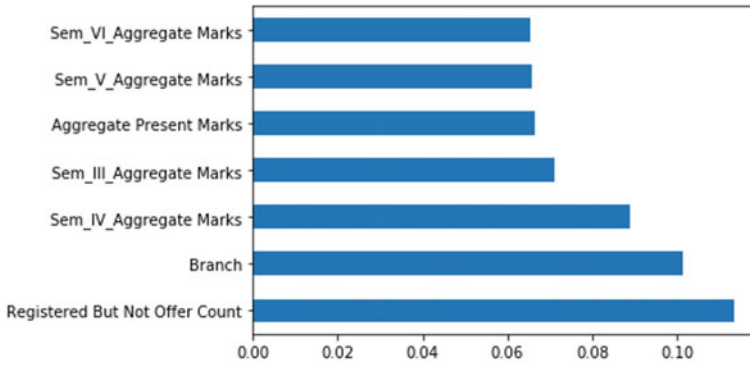


Fig. 1 Feature selection using feature importance for placement prediction

Recursive Features Importance, and Feature importance. We used the python scikit-learn library to implement it.

The Univariate Selection method shows the highest score for the following features (Table 1).

While using the Recursive Feature Importance method, the following features are selected, and the remaining are rejected.

Selected Features: ['Pending Back Papers', 'Sem_III_Pending Back Papers', 'Sem_IV_Aggregate Marks', 'Sem_IV_Pending Back Papers', 'Sem_V_Pending Back Papers', 'Sem_VI_Back Papers'].

Inbuilt class Feature importance comes with Tree based Classifiers; we used Extra Tree Classifier from python scikit-learn library for extracting the top 7 features of the dataset (Fig. 1).

Hence, as per all the above methods and also as per domain our knowledge, we have chosen 25 important features which are as follows to predict target feature 'Job Offer'.

['Branch', 'Aggregate Present Marks', 'Back Papers', 'Pending Back Papers', 'Sem_I_Aggregate Marks', 'Sem_I_Back Papers', 'Sem_I_Pending Back Papers', 'Sem_II_Aggregate Marks', 'Sem_II_Back Papers', 'Sem_II_Pending Back Papers', 'Sem_III_Aggregate Marks', 'Sem_III_Back Papers', 'Sem_III_Pending Back Papers', 'Sem_IV_Aggregate Marks', 'Sem_IV_Back Papers', 'Sem_IV_Pending Back Papers', 'Sem_V_Aggregate Marks', 'Sem_V_Back Papers', 'Sem_V_Pending Back Papers', 'Sem_VI_Aggregate Marks', 'Sem_VI_Back Papers', 'Sem_VI_Pending Back Papers', '12th/Diploma_Aggre_marks', 'SSC Aggregate Marks'].

8 Experimentation

There are adequate models that are studied and tested for the objective with optimal values for K-fold cross-validation (K-FCV), Train-Test split (T-TS), parameters tuning, and testing. In this process, the python sklearn library has played a very important role. So detail is mentioned in the table below.

Apart from the above methods while doing parameter tuning, we have used the following ensemble algorithms.

After the discussion of the accuracy results researcher has suggested a web module named 'Free guide to notify the campus placement status (FGNCPS)' through which students will get to know their placement status in advance and also come to know to work more on weaker areas.

9 Result and Discussion

After implementing data cleaning process, removing all the noise, selecting relevant features and encoded it into ML form, the next step is building a predictive model by applying various ML techniques to find out the best model which gives us more accuracy for train data and test data.

Model selection for placement prediction: After implementing all the above methods mentioned in Table 4 and 5, we found XGBoost classifier is the best classifier to predict campus placement.

Table 4 Results of placement prediction using ML techniques with K-fold cross validation

| Sr. No | Name of Algorithm | Train Accuracy | Test Accuracy |
|--------|--------------------------|--------------------|---------------------|
| 1 | Logistic Regression | 0.7251336898395722 | 0.7371794871794872 |
| 2 | Support Vector Machine | 0.7235294117647059 | 0.7393162393162394 |
| 3 | Decision Tree Classifier | 0.8165775401069518 | 0.782051282051282 |
| 4 | Random Forest Classifier | 0.823663101604278 | 0.7162393162393162 |
| 5 | Gaussian NB | 0.5294117647058824 | 0.49145299145299143 |
| 6 | K Neighbors Classifier | 0.8235294117647058 | 0.7606837606837606 |

Table 5 Results of placement prediction using Ensemble Learning

| Sr. No | Algorithm | Train Accuracy | Test Accuracy |
|--------|------------------------|----------------|------------------|
| 1 | AdaBoostClassifier(DT) | 0.85 | 0.82 |
| 2 | XGBoost | 0.88 | 0.84 |
| 3 | GridSearchCV | 0.851336898395 | 0.82478632478632 |

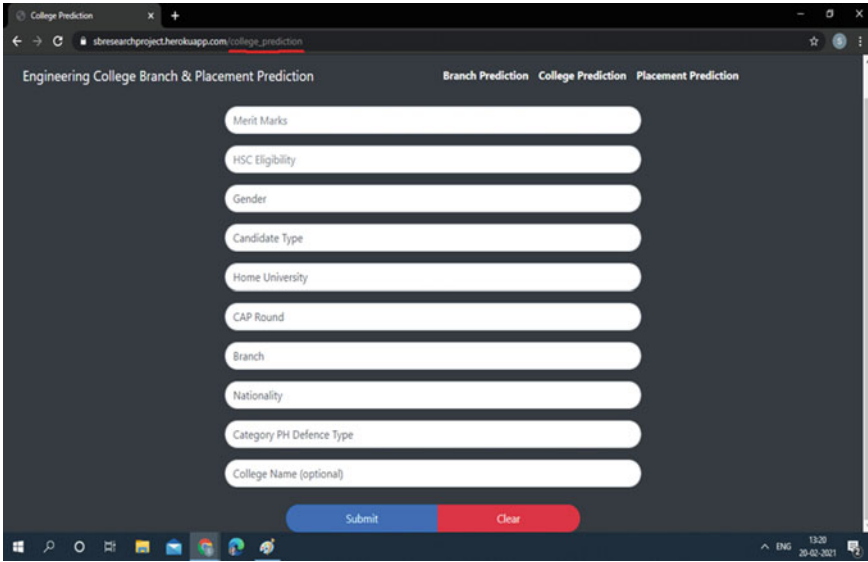


Fig. 2 Placement prediction web module

We can see the result of placement prediction using ensemble classifier XGBoost with 0.88 training accuracy and with 0.84 testing accuracy, which is comparatively very high. Hence, we have chosen the XGBoost classifier to implement the model.

10 Implementation

10.1 A Free Guide to Notify the Campus Placement Status (FGNCPS)

While predicting the campus placement of Engineering and Technology students, we have proposed the following FGNCPS web module. The aspirant student has to submit some basic information which is nothing but selected input features to predict their placement status in the early stage of academics (Fig. 2).

11 Conclusion

In this research, to predict the campus placement of Engineering and Technology students, all the ML model building steps are rigorously implemented on the dataset. Python, various libraries played a vital role during whole this process. In this study,

25 input features are selected out of the existing 46 features of the dataset. These features are very important, according to Univariate Selection, Recursive Features Importance, Lasso feature selection methods, and researchers' domain knowledge. To predict the campus placement, suit of ML and EL methods are experimented and compared. This suit contains Logistic Regression, K-Nearest Neighbors', Decision Tree Classifier, Random Forest Classifier, Naive Bayes, and Support Vector Machine classifiers. Under EL, we have experimented with Adaptive Boosting, Gradient Boosting, and GridSearchCV methods. After a comparison of all algorithms' accuracy, we found that the XGBoost classifier has greater accuracy for this project. Also, it has been observed that feature engineering is a very important step in a model building because, after it, results have been more improved. At the end researchers have suggested, 'A free guide to notify the campus placement status (FGNCPS)' web module for placement aspirant students.

References

1. Kabakchieva D, Stefanova K, Kisimov V (2011) Analyzing university data for determining student profiles and predicting performance. In: 4th International conference on educational data mining (EDM 2011). The Netherlands, pp 347–348
2. Nie M, Yang L, Sun J, Su H, Xia H, Lian D, Yan K (2017) Advanced forecasting of career choices for college students based on campus big data. Higher Education Press and Springer-Verlag Berlin Heidelberg
3. Kumar N, Singh AS, Thirunavukkarasu K, Rajesh E (2020) Campus placement predictive analysis using machine learning. In: 2nd International conference on advances in computing, communication control and networking (ICACCCN), ISBN: 978-1-7281-8337-4/20/\$31.00 ©2020. IEEE
4. Pothuganti M, Swaroopa N (2019) Campus placement prediction using supervised machine learning techniques. *Int J Appl Eng Res* 14(9):2188–2191, ISSN 0973-4562
5. Kalathiya D, Padalkar R, Shah R, Bhoite S (2019) Engineering college admission preferences based on student performance. *Int J Comput Appl Technol Res* 8(09):379–384, ISSN:-2319–8656
6. Khandale S, Bhoite S (2019) Campus placement analyzer: using supervised machine learning algorithms. *Int J Comput Appl Technol Res* 8(09):379–384, ISSN:- 2319–8656, 358–362
7. Apoorva Rao R, Deeksha KC, Vishal Prajwal R, Vrushak K, Nandini (2018) Student placement analyzer: a recommendation system using machine learning. *IJARIIIE* 4(3), ISSN(O)-2395-4396
8. Roy KS, Roopkanth K, Teja VU, Bhavana V, Priyanka J (2018) Student career prediction using advanced machine learning techniques. *Int J Eng Technol* 7:26–29
9. Ishizue R, Sakamoto K, Washizaki H, Fukazawa Y (2018) Student placement and skill ranking predictors for programming classes using class attitude, psychological scales, and code metrics. *Res Pract Technol Enhanced Learn* 13. <https://doi.org/10.1186/s41039-018-0075-y>
10. Sreenivasa Rao K, Swapna N, Praveen Kumar P (2017) Educational data mining for student placement prediction using machine learning algorithms. *Int J Eng Technol*, [S.l.] 7(1.2):43–46, ISSN 2227-524X
11. Bangale M, Bavane S, Gunjal A, Dandhare R, Salunkhe SD (2019) A survey on placement prediction system using machine learning. *IJSART* 5(2), ISSN [ONLINE]: 2395-1052