

# Using Educational Data Mining (EDM) to Prediction and Classify Students

Samira Talebi<sup>1</sup>, Ali Asghar Sayficar<sup>2</sup>

<sup>1</sup>Islamic Azad University Garmsar Branch, Department of Information Technology,  
University Square, Student Street, Iran  
[samiratlb86@gmail.com](mailto:samiratlb86@gmail.com)

<sup>2</sup>Islamic Azad University Garmsar Branch, Department of Information Technology,  
University Square, Student Street, Iran  
[a\\_sayficar@yahoo.com](mailto:a_sayficar@yahoo.com)

**Abstract:** *The aim of this paper is to predict the students' academic performance. It is useful for identifying weak students at an earlier stage. In this study, we used WEKA open source data mining tool to analyze attributes for predicting students' academic performance. The data set comprised of 180 student records and 21 attributes of students registered between year 2010 and 2013. We chose them from FERDOWSI University of Mashhad. We applied the data set to four classifiers (Naive Bayes, LBR, NBTree, Best-First Decision Tree) and obtained the accuracy of predicting the students' performance into either successful or unsuccessful class. The student's academic performance can be predicted by using past experience knowledge discovered from the existing database. A cross-validation with 10 folds was used to evaluate the prediction accuracy. The result showed that Naive Bayes classifier scored the higher percentage of prediction F-Measure of 83.9%.*

**Keywords:** Data Mining, Prediction, Average, Attributes for predicting students, Educational Data Mining (EDM)

## 1. Introduction

Classification and prediction are of high importance in data mining techniques and used in many fields. Recently, researchers have utilized machine learning in order to make wise career decisions. It is useful for both the students and the instructors getting better in their performances. We got our dataset from the Information system of the biggest virtual university of Iran. We decided to extract the attributes that have significant contribution to the prediction of academic performance. The prediction can be done by using data mining tools such as Weka software.

## 2. Methodology

Many studies were undertaken in order to explain the academic performance or to predict the success or the failure (Kotsiantis *et al.*, 2003; Chamillard, 2006; Minaei-

Bidgoli *et al.*, 2003; Merceron and Yacef, 2005; Romero *et al.*, 2008; Superby *et al.*, 2006; Vandamme *et al.*, 2007; Ardila, 2001; Gallagher, 1996; King, 2000; Minnaert and Janssen, 1999; Parmentier, 1994.) they highlighted a series of explanatory factors associated to the student.

We first considered a set of attributes to be taken into account based on a model used by Parmentier (1994). Secondly, we created a questionnaire allowing us to collect a large amount of interesting information on a certain number of students. We distributed this questionnaire by paper to students in the FERDOWSI University of Mashhad.

We used WEKA open source data mining. It supports many machine learning algorithms and data processing tools. In the data preprocessing step, we collected 205 records of students admitted from year 2010 to 2013 at the FERDOWSI University of Mashhad. According to the total average, the students were classified into four classes:

Grade A (Total Average  $\geq 17$ ),  
Grade B ( $15 \leq$  Total Average  $< 17$ ),  
Grade C ( $13 \leq$  Total Average  $< 15$ ),

Grade D (Total Average<13)

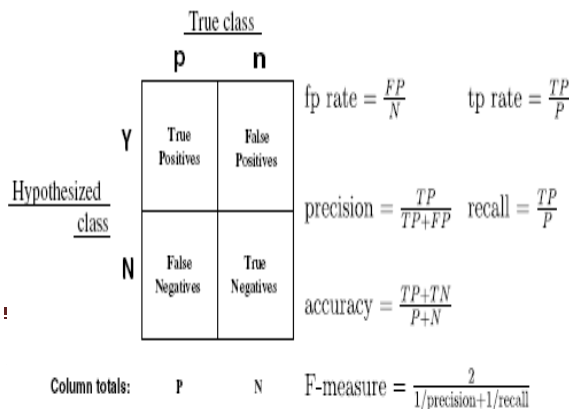
We split the data for training (119 records ~ 66%) and testing (61 records ~ 34%). We used the Naïve Bayes, LBR, NBTree and Best-First Decision Tree classifiers for prediction. Table 1 shows the attributes and their valid values we considered for predicting student's academic performance.

**Table 1:** The attributes used for classification

| Attribute                                        | Value                                         |
|--------------------------------------------------|-----------------------------------------------|
| Sex                                              | Female / male                                 |
| Marital status                                   | Single / married                              |
| Job status                                       | Employed/ unemployed                          |
| City                                             | Mashhad / others                              |
| Right handed or left handed                      | Right hand/ left hand                         |
| The method study                                 | Solo/with the group                           |
| How to study                                     | During the semester/the night before the exam |
| Do my projects                                   | Alone, use the preparation projects           |
| The source of the study                          | Booklet, reference                            |
| Diploma average                                  | A/B/C/D                                       |
| The First university semester average            | A/B/C/D                                       |
| The amount of interest in the field of           | Very high/high/medium /low/very low           |
| Internet accessibility                           | Very high/high/medium /low/very low           |
| Break between high school and university         | Very high/high/medium /low/very low           |
| Mother's level of education                      | Very high/high/medium /low/very low           |
| Type of high school in the pre-university course | Very high/high/medium /low/very low           |
| The number of terms has fallen                   | Very high/high/medium /low/very low           |
| The number of children of the family             | Very high/high/medium /low/very low           |
| The rate of attendance in class                  | Very high/high/medium /low/very low           |
| English language level                           | Very high/high/medium /low/very low           |
| Total Average                                    | A/B/C/D                                       |

### 2.1. Confusion Matrix

A confusion matrix (Kohavi and Provost, 1998) contains information about actual and predicted classifications done by a classification system. Performance of such

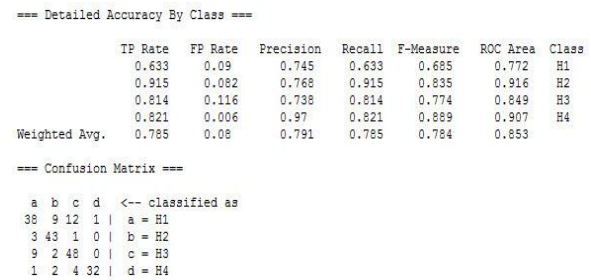


**Fig.1.** Confusion matrix and common performance metrics calculated from it.

### 3. Results

In this paper we used the Naïve Bayes, LBR, NBTree and Best-First Decision Tree to predict student's academic performance. A crossvalidation with 10 folds was used to evaluate the prediction accuracy.

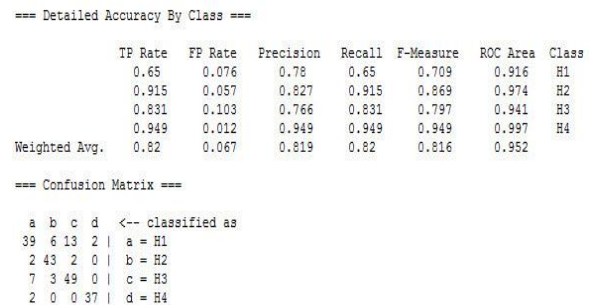
#### 3.1 Best-First Decision Tree



**Fig.2.** Summary of the results of Best-First Decision Tree

As shown in fig 2, the proportion of correct predictions for class H2 are good: 91.5% of the students of class H2 were correctly classified by means of the Naïve Bayes classifier; but the proportion of correct predictions for class H1 are bad, only 63.3% of the students of class H1 were actually classified into class H1. The weighted average of F-Measure is 78.4% and this is not such a good result.

#### 3.2 NBTree



**Fig.3.** Summary of the results of NBTree

As shown in fig 3, the proportion of correct predictions are better than Best-First Decision Tree, 65% of the students of

class H1 were correctly classified by means of the NBTree classifier; and 94.9% of the students of class H4 were actually classified into class H4. The weighted average of F-Measure is 81.6% and this is a good result.

### 3.3 LBR

Fig 4 shows a summary of the results of LBR classifier.

```

=== Detailed Accuracy By Class ===
      TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Class
      0.883    0.069    0.841    0.883    0.862    0.972    H1
      0.851    0.044    0.851    0.851    0.851    0.939    H2
      0.814    0.103    0.762    0.814    0.787    0.939    H3
      0.769    0.012    0.938    0.769    0.845    0.964    H4
Weighted Avg.  0.834    0.062    0.839    0.834    0.835    0.953

=== Confusion Matrix ===
  a  b  c  d  <-- classified as
53  3  4  0  | a = H1
 2  40  5  0  | b = H2
 6  3  48  2  | c = H3
 2  1  6  30  | d = H4

```

**Fig.4.** Summary of the results of LBR classifier

As shown in fig 4, the proportion of correct predictions for class H1 are better than Best-First and LBR classifier: 88.3% of the students of class H1 were correctly classified by means of MLP classifier; and the proportion of correct predictions for class H4 are better than Best-First but is equal to NBTree classifier: 76.9% of the students of class H4 were actually classified into class H4. The weighted average of F-Measure is 83.5% and this is a good result.

### 3.4 Naive Bayes

```

=== Detailed Accuracy By Class ===
      TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Class
      0.9      0.103    0.783    0.9      0.837    0.918    H1
      0.83    0.032    0.886    0.83    0.857    0.961    H2
      0.78    0.068    0.821    0.78    0.8      0.908    H3
      0.846    0.018    0.917    0.846    0.88    0.968    H4
Weighted Avg.  0.839    0.061    0.843    0.839    0.839    0.935

=== Confusion Matrix ===
  a  b  c  d  <-- classified as
54  2  3  1  | a = H1
 4  39  3  1  | b = H2
10  2  46  1  | c = H3
 1  1  4  33  | d = H4

```

**Fig.5.** Summary of the results of Naïve Bayes classifier

As you see in fig 5, the proportion of correct predictions are the best of all: 90% of the students of class H1 were correctly classified by means of Naïve Bayes classifier; and 78% of the students of class H3 were actually classified into class H3. the

weighted average of F-Measure is 83.9% and this is a very good result.

## 4. Conclusion

Identifying the classifiers that contribute the most significant to predict student's academic performance can help to improve the intervention strategies and support services for students who perform poorly in their studies, at an earlier stage. The objective of this study was to introduce and compare some techniques used to predict the student performance at a Azad university of Mashhad. This is important as it provides groundwork for further evaluation of the program. The findings of this study showed that Naïve Bayes classifiers scored the higher percentage of prediction F-Measure of 83.9%. Moreover, the ROC area of LBR classifier is better than other classifiers.

## References

- [1] B.K. Bharadwaj and S. Pal. "Mining Educational Data to Analyze Students' Performance", *International Journal of Advance Computer Science and Applications (IJACSA)*, Vol. 2, No. 6, pp. 63-69, 2011.
- [2] B.K. Bharadwaj and S. Pal. "Data Mining: A prediction for performance improvement using classification", *International Journal of Computer Science and Information Security (IJCSIS)*, Vol. 9, No. 4, pp. 136-140, 2011.
- [3] U.K. Pandey, and S. Pal, "Data Mining: A prediction of performer or underperformer using classification", *(IJCSIT) International Journal of Computer Science and Information Technology*, Vol. 2(2), pp.686-690, ISSN: 0975-9646, 2011.
- [4] U. K. Pandey, and S. Pal, "A Data mining view on class room teaching language", *(IJCSI) International Journal of Computer Science Issue*, Vol. 8, Issue 2, pp. 277-282, ISSN: 1694-0814, 2011.
- [5] Shaeela Ayesha, Tasleem Mustafa, Ahsan Raza Sattar, M. Inayat Khan, "Data mining model for higher education system", *European Journal of Scientific Research*, Vol.43, No.1, pp.24-29, 2010
- [6] Data Mining: A Prediction for Performance Improvement of Engineering Students using Classification *World of Computer Science and Information Technology Journal (WCSIT)* ISSN: 2221-0741 Vol. 2, No. 2, 51-56, 2012

- [7] Kumar, Varun, and Anupama Chadha. Mining Association Rules in Student's Assessment Data. *International Journal of Computer Science Issues* 9. 5:211-216, 2012.