



Contents lists available at ScienceDirect

Expert Systems with Applications

journal homepage: www.elsevier.com/locate/eswa

Student data mining solution–knowledge management system related to higher education institutions

Q1 Srečko Natek^{a,*}, Moti Zwilling^b

^a School for Social and Business Studies, Slovenia

^b Netanya Academic College, Israel

ARTICLE INFO

Keywords:

Data mining
Knowledge management system
Student's success rate
Data mining for small data set
Higher education institution
Educational data mining

ABSTRACT

Higher education institutions (HEIs) are often curious whether students will be successful or not during their study. Before or during their courses the academic institutions try to estimate the percentage of successful students. But is it possible to predict the success rate of students enrolled in their courses? Are there any specific student characteristics, which can be associated with the student success rate? Is there any relevant student data available to HEIs on the basis of which they could predict the student success rate? The answers to the above research questions can generally be obtained using data mining tools. Unfortunately, data mining algorithms work best with large data sets, while student data, available to HEIs, related to courses are limited and falls into the category of small data sets. Thus, the study focuses on data mining for small student data sets and aims to answer the above research questions by comparing two different data mining tools. The conclusions of this study are very promising and will encourage HEIs to incorporate data mining tools as an important part of their higher education knowledge management systems.

© 2014 Published by Elsevier Ltd.

1. Introduction

In knowledge management process, data mining technique can be used to extract and discover the valuable and meaningful knowledge from a large amount of data. Nowadays, data mining has given a great deal of concern and attention in the information industry and in society as a whole. This technique is an approach that is currently receiving a great attention in data analysis and it has been recognized as a newly emerging analysis tool (Osei-Bryson, 2010; Park, 2001; Sinha & Zhao, 2008; Tso & Yau, 2007; Wan & Lei, 2009; Zanakis & Fernandez, 2005; Zhuang, Churilov, Q3 Burstein, & Sikaris, 2009).

In the literature, one may notice that they are many areas which adapted this approach to solve their problems such as in finance, medical, marketing, the stock market, telecommunication, manufacturing, health care and customer relationship. However, the data mining application has not attracted much attention from people in relation to educational systems, despite the majority of HEIs admit, that course's students success rates are very important

for their prestige and senior management of the academic institution. Thus the possibilities to deeply understand the reasons and accurately predict the student success have become very valuable (Dermol & Čater, 2013; Marjetič & Lesjak, 2012; Natek & Lesjak, 2013; Rojko, Lesjak, & Vehovar, 2011; Trunk Širca, Babnik, & Breznik, 2013).

OLAP solutions and statistical methods are well established tool to analyze data but data mining enable fresh approach to understand hidden patterns and data prediction. In particular, data mining has turned to be a very popular among researches because many "approachable" standalone or desktop data mining tools are available in the market. From the observation of different tools, one can notice the following as an example: Microsoft Excel, SPSS, Weka, Protégé as Knowledge Acquisition System and Rapid Miner. Some of them (e.g. MS Excel Mining tool) are normally available to HEI's professor and they can benefit from existing knowledge of the Excel. That is the reason to include it in the research. Weka was chosen among other desktop DM tools because of supporting the data mining analysis in very different way, compared with MS Excel. Moreover, Weka is considered as a very legitimate tool analysis where education management data is involved (Cristóbal, Sebastián, & Enrique, 2008; Romero & Ventura, 2006).

HEIs recently promote knowledge management as encouraging data mining environment for their professors and researchers.

Q2 * Corresponding author. Tel.: +386 41630053.

E-mail addresses: srecko.natek@mfdps.si (S. Natek), Moti.Zwilling@gmail.com (M. Zwilling).

was determined as a training data taken from the courses of Informatics during the years 2010–11 and 2011–12 (74 rows). The second group of data sets was determined as predictable data taken from the courses of Informatics during the years 2012–13 (32 rows, without predictable “Final grade” column). In these courses, there 106 students who were included in the research (rows) participated. Regarding the Informatics 2012–13 course the actual “Final grade” was available for testing the results (32 row of actual “Final grade” data). The described data set from small number of examples but their content was estimated to be relevant for machine learning (Han, Kamber, & Pei, 2012).

In the next technological step the Data Mining Technology were chosen. Where, the Microsoft company generally offered three possibilities of data mining level for analysis. The basic level included MS Excel Table tools that include: data mining features analysis. The intermediate level included MS Excel Data Mining Adds-in features. The Expert level included MS SQL Server Data Mining capabilities. All levels of data mining are using algorithms from MS SQL Server Data mining, with different user interface, different technique and number of parameters, used to manage data mining process. For research purpose the basic level were chosen.

After this step where the Data Mining modeling Technique were chosen, MS Table tools were utilized to offer several Data Mining Techniques: Analyze key influencers, Detect Categories, Fill from Examples, Forecast, Highlight Exceptions, Scenario Analysis, Prediction Calculation and Shopping Basket Analysis. Some analysis are not suitable for a given student data sets, e.g. Forecast, Scenario Analysis and Shopping Basket due to data or content limitation. Among other techniques, the most usable technique is the Key Influencer and Fill from Examples, which were dynamically used for the following research.

Regarding the Model training, a Fill from example technique was used. Analysis was repeated several times, by choosing the “Final grade” as a prediction column (e.g. selecting column containing the examples). After choosing different combination of columns the authors found out the most relevant attributes by omitting student number (not relevant at all) and Exam points columns (directly defined the Final grade). The research dynamically explore the Microsoft Excel Data mining tool to build data mining model and thus relevant research datasets.

Finally the best model was chosen and the results were evaluated. By comparison the results of prediction for student’s during the year 2012/13 according to the prediction class attribute: “Final grade”, with the actual Final grade of the same student the authors could choose the best parameters for the Fill from example data mining technique. The final Data mining Fill from example predictive model was ready for prediction as a ‘future student Final grade’, which was based on already known instances. The model could be used by confidence, until several assumptions were obtained as True (Berry & Linoff, 2000). The first assumption assumed that the past results were considered as good predictors of the future (if the student generation is changed, therefore; their attitude and behavior may vary from the past and create different student pattern). The next assumption assumed that the data was available at hand. The researchers assumed that the data will be always available for HEIs. The last assumption assumed that the data contained what the researchers aimed to predict – the research data mining model as a result contains the relevant data for analysis.

2.2. Building student data mining model with Weka

The data was also built using the WEKA Data Mining tool, were the following steps were conducted: (Fig. 1). In general, the first step in working with data mining analysis tool is a creation of data sets (Step 1). In this work, the data was saved as described above, except from the last column (“the predicted class”) – which was

attributed as the – ‘Final grade’. This attribute was transformed into an ordinal one using the following index: To values between 8 and 10 a “High” value was assigned, to values between 6 and 7 a “Medium” value was assigned and to lower values than 6 a “Low” value was assigned. (In the case of decision tree M5P the class attributes remained the same as numerical).

In the next step the Data Mining Technology were chosen, and the “Weka” tool was found as the one that could be used for the analysis phase (Step 2).

Next the specific technique which was utilized using the “Weka” tool was conducted (Step 3). In this phase three various decision trees techniques were evaluated: J48, M5P and RepTree. These techniques were chosen as best suited for data analysis that was implemented in the last phase (Step 4). (In the case of M5P no conversion of values related to the class attribute was necessary and the Final grade values were remained as numbers).

Finally, in the last phase (Step 5), the data was analyzed by each of the models, and the best one was chosen. In other words – the model that best predicted the class attribute value, with a high accuracy rate was selected. (Each model was evaluated on a training set and its accuracy was evaluated again on the test set respectively). As much as the ‘gap’ between the trained predicted value and the actual data was lower the model was found as a better predictor for analysis regarding the actual data.

3. Results

3.1. Key influencer for student Final grade

At the beginning of the study, the researches aimed to find out which data (variables) are most powerful influencers for the Final grade prediction. Key Influencer Analysis is the right and useful tool for the job. Table 2 exhibits the most powerful relative impact for different variables. For example, the last study year (2012–13), where “Final grade” is empty (the authors aimed to predict it) shows “100” relative impact to favor empty column. Similar, high final points strongly favor high “Final grade”, and low “Activity points” or “Exam points” favor low “Final grade”. The results suggested to eliminate some of these columns from final prediction or at least to try several combination of different columns to produce relevant results.

3.2. Fill from example – student “Final grade” prediction for 2012/13

The next step of study aimed to predict the Student “Final grade” for the missing 2012/13 data by using Fill from example

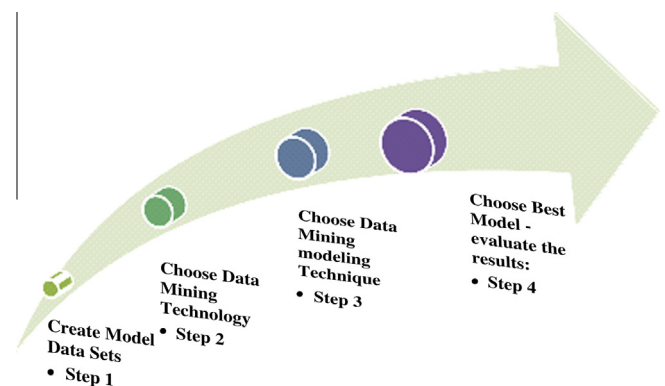


Fig. 1. Data mining process using the weka tool – the process is divided into 4 steps: model creation and data sets preparation (Step 1), choosing the technology for data mining analysis [in this research either excel/weka] (Step 2), choose of modeling technique – such as decision trees or regression (Step 3), choose the best model to evaluate the results (Step 4).

data mining technique, based on machine learning using previous instances from 2010/11 and 2011/12. For more relevant result “Student”, “Exam points” and “Final points” were eliminated from analysis. The Pattern Report in Table 3 shows a diversity of influences. The results justified the author’s decision to eliminate powerful key influencer from analysis.

Prediction for student “Final grade” for 2012/13 was the final part of data mining analysis, using a Fill from Example Data Mining technique. The results for 32 students (rows) are shown in Table 4. The prediction based on data - knowledge, hidden in the student data. Different algorithms were used to exploit data from year 2010/11 and 2011/12 for data mining prediction model (training data) and were aimed to predict the missing student “Final grade” data for the year 2012/13. As noticed, the entire student’s “Final grade” are positive (above 6) because the majority of previous results were found as positive too.

3.3. Data interpretation

Interpretation is considered as the final but vital phase of data mining (Negnevitsky, 2011). To evaluate the relevance of data mining student “Final grade” prediction the study compared the actual student’s “Final grade” with the predicted one (Table 5). The following conclusions were calculated and discussed:

Table 2

Key influencer report for final grade (the table contains the following values: year of study, favors, relative impact and final/activities points).

Filter by 'column' or 'favors' to see how various columns influence 'Final grade 10'			
Column	Value	Favors	Relative impact
Study year	2012–13	(Empty)	100
Final points _100_	92	10	100
Activities points _50_	39–42	7	100
Final points _100_	77	8	100
Final points _100_	85	9	100
Exam points _50_	9	2	100
Final points _100_	53	2	100
Exam points _50_	13	2	100
Final points _100_	51	2	100
Exam points _50_	8	4	100
Final points _100_	41	4	100
Exam points _50_	21	4	100
Final points _100_	66	6	100

Table 3

Fill from example analyze for student “Final grade” (the table contains several values/attributes taken for analysis such as: registration [first/repeat], activities points, etc.).

Filter by 'column' or 'favors' to see how various columns influence 'Final grade 10'			
Column	Value	Favors	Relative impact
Type of study	Part time	10	13
Activities points _50_	43,539–47,000	10	13
Registration	Repeat	2	100
Study year	2010–11	2	13
Type of study	Part time	4	17
Gender	Male	6	33
Study year	2011–12	6	21
Activities points _50_	33,000–38,902	7	43
Status_sport_	Yes	7	36
Employment	Yes	7	31
Activities points _50_	38,902–41,221	7	18
Year of birth	1,967,148–1,982,240	8	25
Activities points _50_	43,539–47,000	8	21
Activities points _50_	43,539–47,000	9	36
Type of study	Part time	9	24
Employment	No	9	10
Activities points _50_	41,221–43,539	9	10

In the list of 32 students from Table 5, only 20 students already approached the exam (actual Final grade in third column). Using the data mining approaches the models predicted all 32 students “Final grade”, included 12 (37.5%) without actual “Final grade”. With zero “Final grade” tolerance, the prediction was successful for 7 “Final grade” (21.8%). With +– 1 “Final grade” tolerance the prediction was successful for 11 “Final grade” (43.3%). The positive “Final grade” prediction was successful for 13 of 32 students (40.6%). And finally if 12 students (without actual “Final grade”) were omitted from the data, the positive “Final grade” prediction was considered as successful for 13 of 20 students (65%), which could be evaluated as generally positive results.

The study’s results are not very clear. From observation of the final interpretation, it seems that the positive “Final grade” prediction was successful for 13 of 20 students (65%) indicating that the potential of data mining technique even for small student data sets is still of importance. However, if a larger data set or even more relevant student data was available the result relevance could be improved, particularly when parameters and columns are carefully combined in data mining analysis pre-process.

3.4. Key influencer for student “Final grade attribute” using the Weka tool

Designing and offering new courses to students is a very important task to higher education decision makers, especially in dynamic and changed environment, where data sets regarding students are limited and does not contain substantial information. In this case “weka” was found by the authors to be used as a data mining analysis tool with no obstacle. In addition, this tool was found (during analysis) as a good prediction on the success of students in their studies along the years, where the Final grade attribute could be used as a key influencer for prediction alongside other attributes such as: Type of studies (Part time/Full time), Age, Employment and the Type of study taken by the student. This finding, that using only several attributes for analysis, where the Final grade is used as a key player in the over whole prediction process was found as very interesting for implication in the higher education system from the aspect of design, offer and budget planning for higher educational programs.

3.5. REPTree model

Training phase indicated 97.0588% where 66 instances were correctly classified and 2.9412% (2 instances) incorrectly classified (Fig. 2).

Testing phase indicated 100% (20/20 instances) correctly classified Instances.

3.6. J48 model

Training phase indicated 98.5294% (67 instances) which were correctly classified where 1.4706% (1 instance) was incorrectly classified, as presented in Fig. 3.

Testing phase indicated 90% (18/20 instances) correctly classified instances (see Fig. 4).

3.7. M5P model

M5P model indicated correlation coefficient of 0.9358 and relative absolute error of 32.2884%.

Evaluation of results related to Table 6, revealed that 18 out of 20 instances were correctly classified by the method J48 and fully classified (20 of 20) where the REPTree method was taken. Therefore when using a small data set of data along with structured data

Table 4

“Final grade” prediction for 2012/13 (the table contains the following attributes: study year, student #, Final grade out of 10, and Final grade prediction received by data mining analysis using the Excel tool).

Study year	Student	Gender	Year of birth	Employment	Status (sport...)	Registration	Type of study	Exam Condition	Activities points (50)	Exam points (50)	Final Points (100)	Final grade (10)	Final grade _10_Extended
2012–13	75	Female	1990	No	No	First	Full time	Yes	38	22	60	7	
2012–13	76	Male	1990	No	No	First	Full time	Yes	38			7	
2012–13	77	Male	1992	No	No	First	Full time	Yes	44			8	
2012–13	78	Male	1992	No	No	First	Full time	Yes	38	30	68	7	
2012–13	79	Male	1990	No	No	First	Full time	Yes	43	30	73	8	
2012–13	80	Male	1983	No	No	First	Full time	Yes	43	19	62	8	
2012–13	81	Female	1991	No	No	First	Full time	Yes	43	49	92	7	
2012–13	82	Female	1973	Yes	No	First	Full time	Yes	44	42	86	8	
2012–13	83	Female	1989	No	No	First	Full time	Yes	42			7	
2012–13	84	Female	1989	No	No	First	Full time	Yes	38			7	
2012–13	85	Female	1990	No	No	First	Full time	Yes	38			7	
2012–13	86	Female	1976	Yes	No	First	Full time	Yes	44	36	80	8	
2012–13	87	Female	1987	No	No	First	Full time	Yes	44	32	76	8	
2012–13	88	Male	1992	No	No	First	Full time	Yes	38	39	77	7	
2012–13	89	Female	1992	No	No	First	Full time	Yes	43	24	67	7	
2012–13	90	Female	1992	No	No	First	Full time	Yes	44			9	
2012–13	91	Male	1992	No	No	First	Full time	Yes	38	30	68	7	
2012–13	92	Male	1992	No	No	First	Full time	Yes	38	6	44	7	
2012–13	93	Male	1992	No	No	First	Full time	Yes	38	24	62	7	
2012–13	94	Female	1992	No	No	First	Full time	Yes	38			7	
2012–13	95	Male	1991	No	No	First	Full time	Yes	43	36	79	8	
2012–13	96	Male	1987	No	No	First	Full time	Yes	38	32	70	7	
2012–13	97	Female	1990	No	No	Repeat	Full time	Yes	41			7	
2012–13	98	Female	1991	No	No	First	Full time	Yes	38			7	
2012–13	99	Female	1968	Yes	No	First	Full time	Yes	44	38	82	8	
2012–13	100	Female	1986	No	No	First	Full time	Yes	44	46	90	8	
2012–13	101	Male	1989	No	No	First	Full time	Yes	38	17	55	7	
2012–13	102	Male	1990	No	No	First	Part time	No				7	
2012–13	103	Female	1988	No	No	First	Part time	Yes	38	25	63	7	
2012–13	104	Female	1969	No	No	First	Part time	Yes	44			8	
2012–13	105	Female	1988	No	No	First	Part time	Yes	44			9	
2012–13	106	Female	1990	No	No	First	Part time	Yes	44	30	74	9	

Table 5

Actual and predicted “Final grade” comparison (study year 2012–13) using 2 decision tree methods. Analysis was conducted by “Weka”. The Final grade prediction (out of 10) is compared with the actual Final grade. Other attributers taken are: Student # & year of study.

Study year	Student	Final grade (10)	Prediction
2012–13	75	4	7
2012–13	76		7
2012–13	77		8
2012–13	78	7	7
2012–13	79	7	8
2012–13	80	3	8
2012–13	81	10	7
2012–13	82	9	8
2012–13	83		7
2012–13	84		7
2012–13	85		7
2012–13	86	8	8
2012–13	87	8	8
2012–13	88	8	7
2012–13	89	4	7
2012–13	90		9
2012–13	91	7	7
2012–13	92	1	7
2012–13	93	4	7
2012–13	94		7
2012–13	95	8	8
2012–13	96	7	7
2012–13	97		7
2012–13	98		7
2012–13	99	8	8
2012–13	100	9	8
2012–13	101	3	7
2012–13	102		7
2012–13	103	4	7
2012–13	104		8
2012–13	105		9
2012–13	106	7	9

mining algorithms (decision trees) one can notice that the over whole prediction rate is found as high (above 80%).

4. Discussion

The study explores the possibility to predict the success rate of students enrolled to an academic course using a contemporary data mining tools normally available to HEIs. The research clearly exhibits that available desktop data mining tools have matured in terms of their usability and ease of use, and provide usable results without extensive investment.

The results of the study yield that student data, available to higher education decision makers (such as professors) via export/import features, carries enough student-specific characteristics (or in other words – information) in the sense of hidden knowledge which can be successfully associated with student success rates.

Despite the well-known fact that data mining algorithms work best on large data sets, the study focused on student data available to HEIs which is limited and clearly falls into the category of a small data set. Results show that small student data sets in the specific data mining analysis did not limit the use of data mining tools.

Moreover the Weka tool, that was used as a comparative analysis to MS Excel, showed that by using decision tree models a high prediction accuracy (especially with the REPTree model) is obtained (the accuracy was verified during the test phase). Moreover, M5P regression tree did not perform as well as the other decision trees since it required more assumptions than the formers. In fact, The REPTree was found as a fast decision tree learner which builds a decision/regression tree using information gain as the splitting criterion, and prunes it using reduced error pruning. The model only sorted values for numeric attributes once. Missing values were dealt with using C4.5's method of using fractional

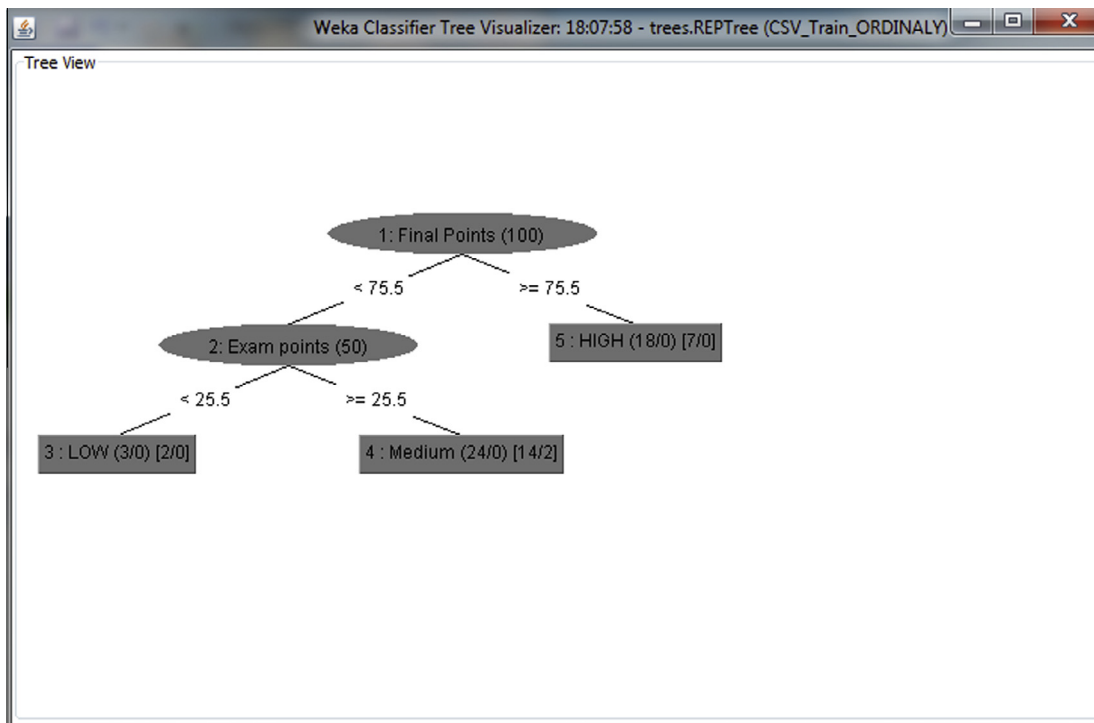


Fig. 2. "REPTree Decision Tree Model (constructed in the train phase) – the root contains the final points attribute where the classification of Exam points and final points (low, medium, high) is represented in the leafs.

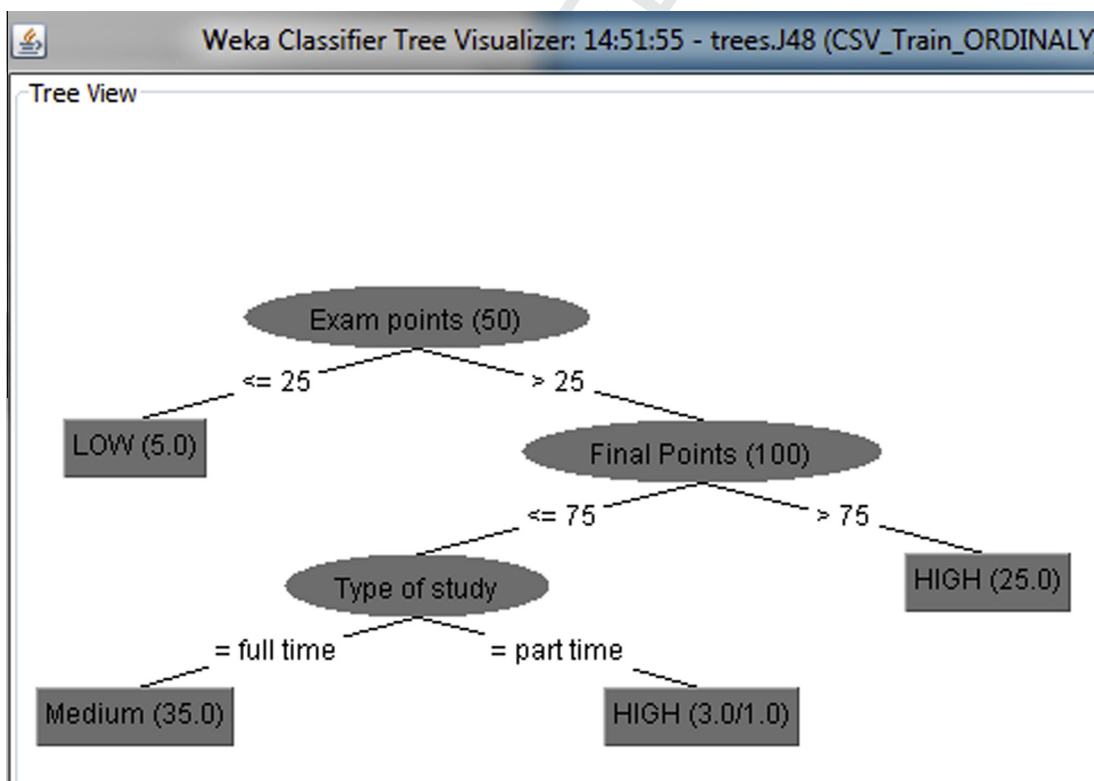


Fig. 3. "J48 Decision Tree Model (constructed in the train phase) – the root contains the exam points attribute where the classification of final points and type of study (low, medium, high) is represented in the leafs.

instances. From the above, it may be deduced that since the predicted attribute was transformed into an ordinary value, REPTree was less sensitive to missing values than J48 thus prediction accuracy on the test set was better performed, however on the training

set it was shown that a slightly less performance than J48 (~97% vs. ~98%) prediction rate is achieved. Hence, J48 is a slightly modified C4.5 algorithm that generates a classification-decision tree for the given data-set by recursive partitioning of data (the decision is

394
395
396
397

