



High School Major Classification towards University Students Variable of Score Using Naïve Bayes Algorithm

Usman Sudibyo¹, Yani Parti Astuti², Achmad Wahid Kurniawan³

^{1,2,3}Informatic Technology Departement, FIK, Universitas Dian Nuswantoro, Indonesia
Email: ¹usman.sudibyo@dsn.dinus.ac.id, ²yanipartiasuti@dsn.dinus.ac.id, ³wahid@dsn.dinus.ac.id

Abstract

Completeness of data in each institution, such as major in a university, is necessary. Data of former school has important role in the need of students data. However, there is no relationship between data of former school and variable of students' score. The suitable classification used in this research is data mining technique which is naïve bayes algorithm. This algorithm is able to manage massive data with a relative fast timing. By using this algorithm, the data results 64.77% performances in classifying former major in school towards variable of score. Hence, the researchers optimize selection feature by using Backward Elimination and result 71.71% performances data. It concludes that performance increases with selection feature. The increasing shows that not all variable of score affects the former school major.

Keyword: Naïve Bayes, Classification, Backward Elimination

1. INTRODUCTION

The utilization of data must be maximalized well in the area of education. One of many ways is to know and to learn the marks that are influencing students performance [1]. The mark, which is the background or the high school, is the important tool because it has the impact in the course grade [2]. To learn and to manage the marks is necessary towards the course grade. The challenges of this case are analysing the students performance, noticing the unique mark from each student, and having the strategy and behaviour in the near future [3].

Students data about their former school, high school, can be used to know how efficiency the influence of students former school with the variable of students grade. Those data are processed and one pattern is finally found. The pattern is a saving mode using introduction pattern technique called data mining [4]. Data mining is a process to gain information from data saved in respiratory using technique to gain introduction pattern, statistic, and mathematics [5]. Meanwhile, classification is a process used to gain the function or model that can differ the class of data [6]. The developed model will be used as prediction towards unknown class label. Pandey and Pal [7] observed about Bayes Classification algorithm to estimate students performance whether the students will accomplished the study well. Besides, Bharawaj and Pal [8] observed the students performance using 300 participants from 5 university, and different Bachelor of Coumputer Application (BCA).

The method used in classification is Bayesian using 17 marks. The results show some factors such as grade on final examination in high school, student's address, the learning process, mother qualification, student's activity, family annual income, and student's family status. Those factors are believed to be influencing the students academic achievement.

According to the previous researchers above, the researchers use the Naïve Bayes algorithm as classification techniques. This algorithm is used as the consideration that Naïve Bayes is one of prediction techniques probability which is depended on the regulation and theory of Bayes. Naïve Bayes algorithm is also known as strong independent assumption on the future. It means that the future on data do not include towards existence of other future with the same data [9]. Naïve Bayes algorithm is one of classification algorithm that has high speed computation. It can also solve huge dimension dataset problem [10]. However, it has weakness on the correlation between the marks, so that it decreases the performance of Naïve Bayes classification [9-10].

Optimization using selection feature is necessary to increase that performance of Naïve Bayes such as Backward Elimination and Forward Selection. In this research, the researchers use Backward Elimination to omit the irrelevant marks [5].

2. METHODS

The method applied in this research used the experiment through this following steps as seen on Figure 1 [9]:

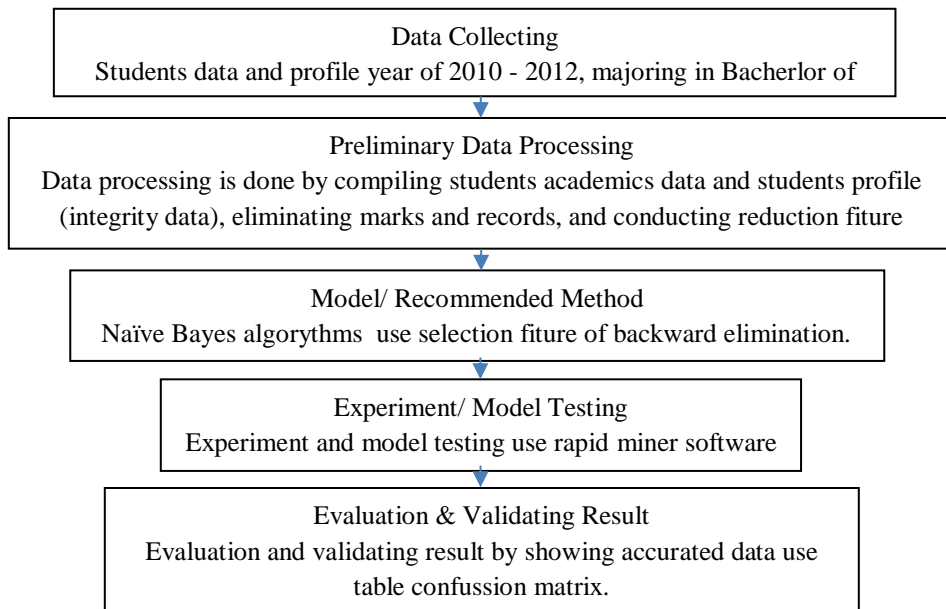


Figure 1. Research Methods Process

2.1. Data Collecting

Data was taken from the students academic section such as courses score and admission section about students former school. Data was taken from 1030 students year of 2010-2012.

2.2. Preliminary Data Processing

In this step, the researchers simplified data to define data and to be applied in algorithm or recommended method. These are the following steps:

1. Data integration was saving the media by compiling them. Identifier data and academic data were compiled in one saved medium.
2. Data reduction was the taken data with less records and marks which were less necessary. Then, they were decreased the unnecessary mark and record such as unidentified former school data inputed manually.

In data reduction, unnecessary marks were eliminated such as Grade Point Average (GPA), general competence subjects, and supporting competence subjects. Record like unidentified former data which are not inputed manually by students were also eliminated because variable of former school were used as class label in the classification. From preliminary data processing, it showed valid data with 489 records contained 25 marks, consisted of course grade and supporting competence subject grade, and former school as variable of label.

2.3. Recommended Method

2.3.1. Naïve Bayes Algorithm

Naïve Bayes is an applied theory of Bayes. Naïve Bayes algorithm is based on assumption to simplify marks which are not related among the marks [11]. Bayes is also statistic classification used to predict probability in sub-class [5]. Bayes has high accurated level and high speed computation when it is applied in dimension data.

Bayes means that each future do not have relation. Bayes predictions are based on Bayes theory with standart formula in Formula(1) [9]:

$$P(X | H) = \frac{P(X|H) * P(H)}{P(X)} \quad (1)$$

Remarks:

X: data of unidentified class

H: Hypothesa on spesific class of data x

$P(X | H)$: Probability that X gained by the value of H

$P(H)$: Pre-probability (chance) on hypothesa H without other marks

$P(X)$: Pre-probability (chance) on X gained without other hypothesis

Formula of Naïve Bayes used as classification is showed as follow:

$$P(Y|X) = \frac{P(Y)\prod_{i=1}^q P(X_i|Y)}{P(X)} \quad (2)$$

$P(Y|X)$ is probability data of variable X on label class Y. $P(Y)$ is pre-probability on class Y. $\prod_{i=1}^q P(X_i|Y)$ is bound probability which is independent on label class Y from variable future X. $P(X)$ is defined as prediction and it finally is formulated on $P(Y)\prod_{i=1}^q P(X_i|Y)$ by choosing the biggest part as prediction result class. Meanwhile, the independent probability $\prod_{i=1}^q P(X_i|Y)$ is the result of each future data in each class and formulated as follow:

$$P(X/Y=y) = \prod_{i=1}^q P(X_i|Y = y) \quad (3)$$

The group future $X = \{X_1, X_2, X_3, \dots, X_q\}$ consist of q mark or q dimensionon.

2.3.2. Backward Elimination

Backward Elimination is defined as eliminator for irrelevant marks [5]. Backward Elimination Algorithms is based on linear regression. The steps of Backward Elimination are as follow:

- Regressing variable of class Y with variable of course grade as predictor, such as X_1, X_2, \dots, X_k .
- Choosing one predictor with the smallest F_{partial} . The significant value of partial (formulated as P_{partial}) is bigger than P_{out} . After that, it is excluded from the model. Value of $F_{\text{partial}} =$ quadratic statistic value on testing T. Value of P_{partial} is the extension on the right side of F_{partial} towards the distribution of $F(1, v)$, using v as error or residual degree of freedom. Value of P_{out} is written as α_{out} resulted 0.1 from the distribution table $F(1, v)$; if $F_{\text{partial}} < F(1, v, \alpha_{\text{out}})$, predictor is omitted from model.
- Regressing Y with the rest $k-1$ predictor.
- Repeating the step 2.
- The output in predictor and regression value Y towards the rest predictor is done until predictor has significant value bigger than P_{out} .
- If there is no predictor contained value $F_{\text{partial}} < F(1, v, \alpha_{\text{out}})$, the model is chosen as the best model.

2.4. Experiment and Model Testing

Both experiment and model testing done in this research are as follow:

- Preparing data that will be used in experiment.
- Pre-processing by omitting the unidentified data in coloumn former school.
- Implementing data mining using Rapid Miner Software to develop classification model of Naïve bayes algorithms by Backward Elimination. Rapid Miner is free software for data mining and machine learning.
- Conducting model testing of Naïve bayes algorithms by getting accuracy value in the classification through confusion matrix such as Table 1.

Table 1. Confusion matrix: case of two class models

CLASSIFICATION		PREDICTED CLASS	
OBSERVED CLASS	Class = YES	Class = YES a (true positive-TP)	Class = NO b (false negative-FN)
	Class = NO	c (false positive-FP)	d (true negative-TN)

Column a and d are the correct classification in which classification will exactly predict in real result. However, column b is wrong classification because model is predicted as no (negative), but the result is yes (positif). Column c is also wrong classification because the prediction is yes (positif), but the result is no {negative} [5]. According to the prediction class, the formula from confusion matrix shows the accuracy as follow:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \tag{4}$$

- e. Analysing the result using Naïve bayes algorithm with Backward Elimination selection figure.

3. RESULTS AND DISCUSSION

Data taken from pre-processing is tested its validity using corss validation k=10. Before the data is validated using cross validation, the marks will be reduced using selection figure, backward elimination. Data before tested in backward elimination is shown in Figure 2.

ExampleSet (488 examples, 1 special attribute, 26 regular attributes)

R..	jur	m_nim	JA..	IN..	KR..	KA..	FI..	DA..	PE..	KA..	FI..	P...A..	PR...M..	MA..	S...S..	O...PE..	O...B..	R...R..	SI..	SI..	K...K..	DA..					
1	NON-EKSA	A11.2010	3	3	1	3	2	4	3	4	3	3	2	3	4	2	4	2	4	4	4	3	4	3	4	2	3
2	EKSAK	A11.2010	2	3	2	3	4	2	3	3	3	4	4	2	3	3	2	3	3	2	3	3	3	3	4	2	1
3	EKSAK	A11.2010	4	4	3	3	3	3	4	3	4	3	4	2	3	3	4	4	4	4	3	3	3	3	4	3	1
4	EKSAK	A11.2010	3	3	3	3	4	3	4	4	3	2	4	4	4	3	3	4	4	3	4	4	3	3	4	4	1
5	EKSAK	A11.2010	4	3	3	3	3	3	3	3	3	2	4	4	4	3	3	4	4	4	4	3	3	3	3	4	1
6	EKSAK	A11.2010	3	3	3	4	4	3	4	3	3	4	3	4	4	3	3	4	4	3	2	3	2	3	3	4	3
7	EKSAK	A11.2010	3	3	4	4	4	4	3	3	3	3	3	4	4	3	3	4	4	3	4	3	3	4	3	4	4
8	EKSAK	A11.2010	3	3	4	3	3	3	4	3	3	3	4	3	3	3	3	4	3	3	4	3	3	3	3	3	3
9	EKSAK	A11.2010	4	4	3	4	4	3	4	3	4	4	3	3	4	4	4	4	4	4	4	4	4	3	4	4	4
10	EKSAK	A11.2010	4	4	3	3	3	4	4	4	3	3	4	3	4	3	4	4	3	4	4	4	3	3	3	3	1
11	EKSAK	A11.2010	3	3	1	4	4	3	3	4	4	2	3	3	4	4	4	3	4	3	2	3	3	3	3	4	1
12	EKSAK	A11.2010	4	3	4	4	4	4	3	4	4	3	4	3	4	4	2	3	3	3	2	3	4	4	4	4	1
13	EKSAK	A11.2010	4	3	2	3	4	3	4	3	3	3	3	2	4	2	4	3	3	2	3	3	4	3	3	1	
14	EKSAK	A11.2010	2	4	3	4	3	2	3	3	3	4	3	3	3	3	4	3	3	3	4	3	2	4	3	3	3
15	EKSAK	A11.2010	3	2	1	3	4	3	3	4	3	3	3	4	3	3	2	3	3	3	4	3	2	4	3	2	2
16	NON-EKSA	A11.2010	3	4	3	4	4	3	3	4	3	3	3	3	4	3	3	3	2	3	3	3	3	3	4	3	1
17	EKSAK	A11.2010	3	4	1	4	4	3	3	4	2	4	4	3	4	3	3	4	4	4	4	3	3	3	4	3	1
18	NON-EKSA	A11.2010	2	4	1	4	4	3	3	4	3	3	3	2	4	3	2	3	3	3	4	2	3	4	3	3	1
19	EKSAK	A11.2010	4	4	3	4	3	3	3	4	3	2	2	3	3	3	2	3	3	4	4	3	4	4	3	3	3

Figure 2. Students Data before tested by selection figure

The following step is reducing data using backward elimination, resulted reduction data as shown in Figure 3.

ExampleSet (488 examples, 1 special attribute, 26 regular attributes)																												
R...	jur	m_nim	JA..	IN..	KR..	KA..	FI..	DA..	PE..	KA..	FI..	P..	A..	PR..	M..	MA..	S..	S..	O..	PE..	O..	B..	R..	SI..	K..	DA..		
1	NON-EKSA	A11.2010	3	3	1	3	2	4	3	4	3	3	2	3	4	2	4	2	4	4	4	3	4	3	4	2	3	
2	EKSAK	A11.2010	2	3	2	3	4	2	3	3	4	4	2	3	3	3	2	3	3	2	3	3	3	3	3	4	2	1
3	EKSAK	A11.2010	4	4	3	3	3	3	4	3	4	3	4	2	3	3	4	4	4	4	3	3	3	3	4	3	1	
4	EKSAK	A11.2010	3	3	3	4	3	4	4	3	2	4	4	4	3	3	4	4	3	4	4	3	4	3	3	4	4	1
5	EKSAK	A11.2010	4	3	3	3	3	3	3	3	3	2	4	4	4	3	3	4	4	4	4	3	3	3	3	4	1	
6	EKSAK	A11.2010	3	3	3	4	4	3	4	3	3	4	3	4	4	3	3	4	4	3	2	3	2	3	3	4	3	
7	EKSAK	A11.2010	3	3	4	4	4	4	3	3	3	3	4	4	3	3	4	4	3	4	3	4	3	3	3	4	4	
8	EKSAK	A11.2010	3	3	4	3	3	4	3	3	4	3	3	4	3	3	3	4	3	3	4	3	3	3	3	3	3	
9	EKSAK	A11.2010	4	4	3	4	4	3	4	3	4	4	3	3	4	4	4	4	4	4	4	4	4	4	4	4	4	
10	EKSAK	A11.2010	4	4	3	3	3	4	4	4	3	3	4	3	4	3	4	4	3	4	4	4	4	3	3	3	1	
11	EKSAK	A11.2010	3	3	1	4	4	3	3	4	4	2	3	3	4	4	4	3	4	3	2	3	3	3	3	4	1	
12	EKSAK	A11.2010	4	3	4	4	4	3	4	4	3	4	4	3	4	4	2	3	3	3	2	3	4	4	4	4	1	
13	EKSAK	A11.2010	4	4	3	2	3	4	3	4	3	3	3	3	2	4	2	4	3	3	2	3	3	4	3	3	1	
14	EKSAK	A11.2010	2	4	3	4	3	2	3	3	3	4	3	3	3	3	4	3	3	3	4	3	3	4	2	4	3	3
15	EKSAK	A11.2010	3	2	1	3	4	3	3	4	3	3	3	4	3	3	2	3	3	4	3	2	4	3	2	4	2	
16	NON-EKSA	A11.2010	3	4	3	4	4	3	3	4	3	3	3	3	4	3	3	2	3	3	3	3	3	3	3	4	3	1
17	EKSAK	A11.2010	3	4	1	4	4	3	3	4	2	4	4	3	4	3	3	4	4	4	4	3	3	3	4	3	1	
18	NON-EKSA	A11.2010	2	4	1	4	4	3	3	4	3	3	2	4	3	2	3	3	3	4	2	3	4	3	3	3	1	
19	EKSAK	A11.2010	4	4	3	4	3	3	3	4	3	2	2	3	3	3	2	3	3	4	4	3	4	4	3	3	3	

Figure 3. Students data after reduced by backward elimination

According to figure 2 and 3, it is shown that there are some reduction marks from 25 to 21 marks. It can be conclude that the 4 reduced marks are irrelevant. Those reduced variable are grade of Calculus 1, Probability and Statistics, Theory of Language and Automata, and Artificial Science. The following step is testing the validity of data using cross validity k=10 as shown in Figure 4.

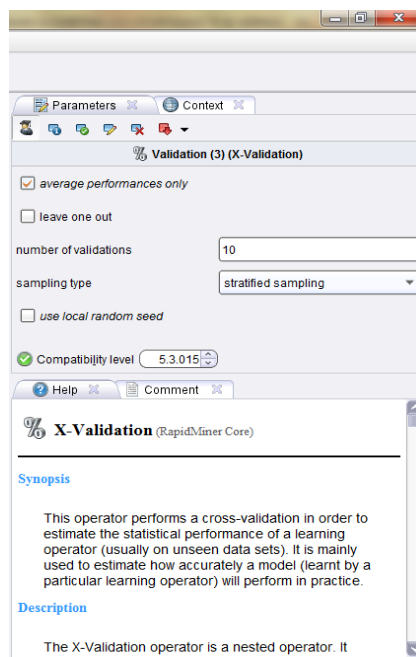


Figure 4. The usage of cross validity k=10

Data validated by validity $k=10$ is tested using naïve bayes with table of confusion matrix shown in Figure 5.

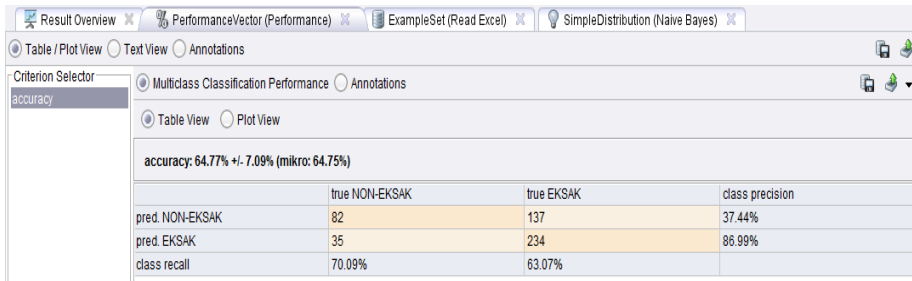


Figure 5. Confusion matrix algorithms of naïve bayes

According to figure 5, the accuracy shown from table is 64.77%. It can be concluded that naïve bayes algorithm is not optimal. So that it is recommended to use the selection figure. The figure used for data is backward elimination resulted as Figure 6.

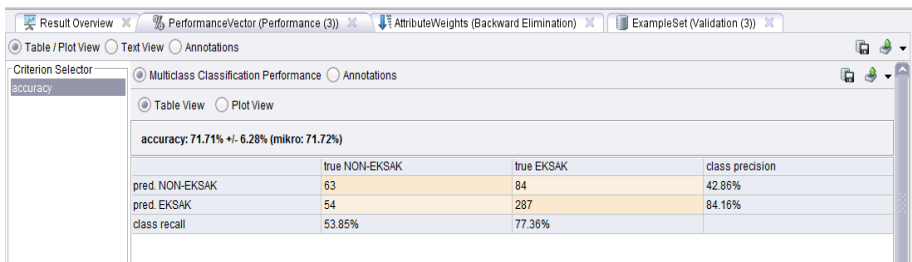


Figure 6. Confusion matrix of naïve bayes algorithm – backward elimination

According to Figure 6, it can be seen that backward elimination on naïve bayes algorithm shows 71.71% which means that it shows increasing performance by using optimization 6.94%. The data results Tabel 2.

Table 2. Comparison of testing result

Algoritma	Waktu	Akurasi
Naïve Bayes	1 detik	64,77%
Naïve Bayes – Backward Elimination	6 detik	71,71%

4. CONCLUSION

Some course grade is influenced by the major in former school or high school or students education background. It is showed by the increasing performance from naïve bayes algorithm testing using backward elimination. The increasing performance shows 6.94%.

5. REFERENCES

- [1] Kurniawan, D., Wicaksono, W. and Astuti, Y.P. 2016. Pemanfaatan Educational Data Mining (EDM) Untuk Memprediksi Masa Studi Mahasiswa Menggunakan Algoritma C4. 5 (Studi Kasus: TI-S1 Udinus). *Momentum*, Vol. 3(2), pp. 48-52.
- [2] Andriani, P. 2010. Pengaruh Asal Sekolah dan Jurusan Terhadap Hasil Belajar Pengantar Dasar Matematika Mahasiswa Fakultas Tarbiyah IAIN Mataram. *Beta Jurnal Tadris Matematika*, Vol. 3(2), pp.118-133.
- [3] Fauzan, A. 2010. *Ide-ide penelitian pendidikan matematika. Proceedings seminar Nasional Pendidikan Matematika*. UIN Syarif Hidayatullah Jakarta, November 27, 2010.
- [4] Larose D. T. 2006. *Discovering Knowledge In Data*. United States of America: John Wiley & Sons, Inc.
- [5] Larose, D. T. 2006. *Data Mining Methods and Models*, Hoboken, New Jersey, United State of America: John Wiley & Sons, Inc.
- [6] Han, J., & Kamber, M. 2007. *Data Mining Concepts and Techniques (2nd ed.)*. San Francisco, United State America: Morgan Kaufmann Publishers.
- [7] Pandey, U., & Pal, S. 2011. 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.
- [8] Yadav, S. K., & Pal, S. 2012. Data Mining A Prediction for Performance Improvement of Engineering Students using Classification. *World of Computer Science and Information Technology Journal (WCSIT)*, Vol. 2(2), pp.51-56.
- [9] Prasetyo E. 2012. *Data Mining: Konsep dan Aplikasi menggunakan MATLAB*, 1st ed. Yogyakarta,Indonesia: Andi.
- [10] Bose, I. and Chen, X. 2009. Quantitative models for direct marketing: A review from systems perspective. *European Journal of Operational Research*, Vol. 195(1), pp.1-16.
- [11] Santoso B. 2007. *Data Mining Teknik Pemanfaatan Data untuk Keperluan Bisnis*, 1st ed. Yogyakarta, Indonesia.