

Combining the Unsupervised

by Puput Wanarti

Submission date: 02-Jun-2021 02:51PM (UTC+0700)

Submission ID: 1598904888

File name: yamasari2020.pdf (223.97K)

Word count: 3396

Character count: 17883

Combining the Unsupervised Discretization Method and the Statistical Machine Learning on the Students' Performance

1
Yuni Yamasari
Department of Informatics
Universitas Negeri Surabaya
Surabaya, Indonesia
yuniyamasari@unesa.ac.id

Anita Qoiriah
Department of Informatics
Universitas Negeri Surabaya
Surabaya, Indonesia
anitaqoiriah@unesa.ac.id

1
Naim Rochmawati
Department of Informatics
Universitas Negeri Surabaya
Surabaya, Indonesia
naimrochmawati@unesa.ac.id

Wiyli Yustanti
Department of Informatics
Universitas Negeri Surabaya
Surabaya, Indonesia
wiyliyustanti@unesa.ac.id

Hapsari P. A. Tjahyaningtjas
Department of Electrical Engineering
Universitas Negeri Surabaya
Surabaya, Indonesia
hapsaripeni@unesa.ac.id

Puput W. Rusimanto
Department of Electrical Engineering
Universitas Negeri Surabaya
Surabaya, Indonesia
puputwanarti@unesa.ac.id

Abstract— the suitability of the data with the method in the process of data mining is very important to increase the process performance. However, In Educational Data Mining (EDM), not much research has focused on this field. Therefore, this study proposes to combine an unsupervised discretization method called "equal width interval" and logistic regression as statistical machine learning to improve the performance of the model relating to students' performance. The discretization method is performed on student data with several intervals, namely: 3-interval, 4-interval, and 5-interval. Then, these intervals are combined with logistic regression in two regularizations, namely: lasso and ridge. Evaluation is carried out on all combinations. The experimental results indicate that the highest performance, in terms of the accuracy level, is achieved by the model combining a 3-interval and logistic regression on all regularization. This combination can increase the model performance based on the average accuracy level of about 4.08-8.49 on the ridge regularization and about 4.28-8.6 on the lasso regularization.

Keywords—students' performance, data mining, machine learning, logistic regression, discretization

8 I. INTRODUCTION

Nowadays, almost all education institutions explore information and communication technology to enhance their process, for example learning process [1][2][3], evaluation process [4][5] despite its security [6]. This situation generates massive data pushing research about Educational Data Mining (EDM) [7]. The one of popular tasks in EDM is the classification of students' performance [8]. Here, student data is mined to get information about the students' performance. Before the mining process is done, methods sometimes are applied in the pre-processing stage [9]. They are the normalization, feature extraction, feature selection, the discretization method, etc. One of the objectives is to improve the performance of the system built [10].

Relating to the discretization method, it is used to transform the numerical method to a categorical method. Also, this method changes the non-standard probability distribution to the standard probability distribution. This method is applied with many reasons as follows: algorithms of machine learning requiring categorical or ordinal variables [11], the non-standard probability distributions causing the performance degradation of machine learning, the result of mapping smoothing out the relationships between observations because of providing a high-order ranking of values [12]. There are

two mainstreams of the discretization methods, namely: the supervised methods and the unsupervised methods. The methods included in the unsupervised discretization methods are equal width interval, equal frequency interval, etc. For the supervised method, they are adaptive quantizes, chi merge, predictive value max, etc. [13]

15
The logistic regression is one of the methods in machine learning requiring the categorical variables. This method is grouped in the statistical machine learning [14]. As we have known, statistics have a very important role in the development of other sciences to draw conclusions, test hypotheses or theories, understand phenomena, analyze experiments, determine decisions, and so forth. Meanwhile, Machine learning, which is one branch of artificial intelligence (artificial intelligence) is currently continuing to experience growth and increasingly popular. The development of statistics and machine learning is of course because it can not be separated from the main factor, namely data. Machine learning has at least two main objectives, namely: solving problems in predicting the future (unobserved event) and gaining knowledge (knowledge discovery). Statistical machine learning refers to techniques for predicting the future and getting knowledge from data rationally. To be able to get these goals, statistical machine learning can be the right tool or method. Statistics acts as a learning base that utilizes statistical theory to inference and interprets the models, while machine learning focuses on the use of models to predict new data. Statistics and machine learning form a concept called Statistical machine learning using logistic regression models. Logistic regression models are among the models that are often used by machine learning practitioners [15].

Our paper focuses on the exploration of one of the discretization method called "equal width interval" and the statistical machine learning, namely: logistic regression on the students' performance domain.

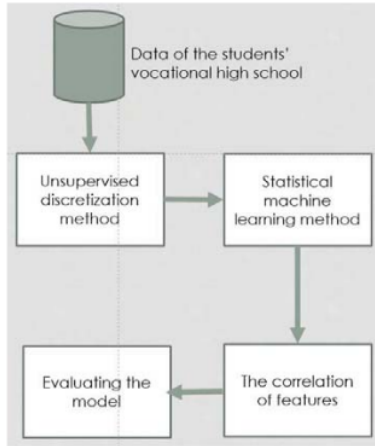


Fig. 1. The proposed architecture

II. METHOD

In this section, this proposed method is depicted in Fig. 1. Here, this architecture consists of many stages. They are as follows:

Stage 1: data of students' vocational high school

This research mine the student data of the previous work [16]. The student data is collected when the students join the evaluation process in the e-learning system. The data consists of 5 features that are extracted from 101 features of the raw dataset. They are as follows: Done, PercentTrue, Time, Hint, and Score which having the numeric data type.

Stage 2: the exploration of the unsupervised discretization method

In this stage, we apply the unsupervised discretization method called "equal width interval" with formula as follows:

$$w = \frac{(value_{max} - value_{min})}{k} \quad (1)$$

Where, w = width of an interval, k = the number of intervals which can be determined manually. For the range threshold on $value_{min} + iw$, where, $i = 1, \dots, k - 1$, each range can be defined as follows:

$$value_{min} + w, value_{min} + 2w, \dots, value_{min} + (k - 1)w \quad (2)$$

We divide student data into some intervals, namely: 3, 4, and 5. Then we definite as 3-interval, 4-interval, and 5-interval. This stage is done to know how many of the best intervals are implemented on our student data. This information is very important to the labeling process relating to the students' performance level.

Stage 3: logistic regression

We combine the result of the previous stage with logistic regression in this stage. Also, there are two regularizations in the logistic regression, namely: Lasso [17] and Ridge [18][19]. We explore all regularizations to be experimented in our research to reach the optimum result.

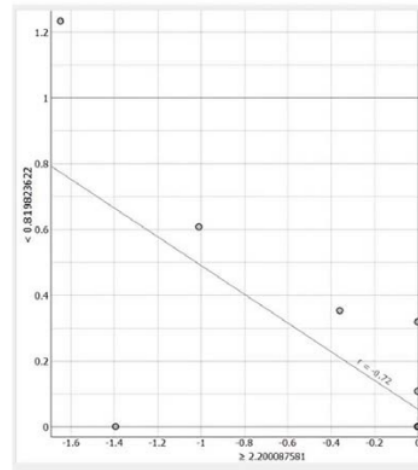


Fig. 2. Regression line with the highest correlation on discretization-3-interval-lasso

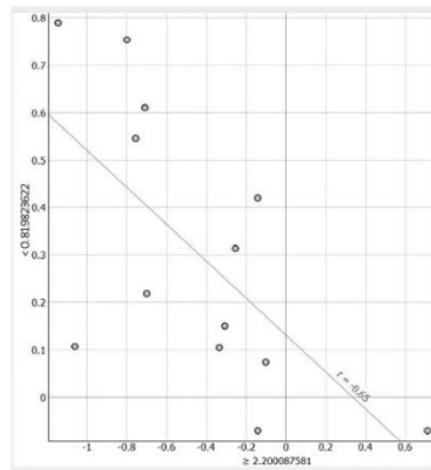


Fig. 3. Regression line with the highest correlation on discretization-3-interval-ridge

Stage 4: analyzing the correlation of features

In this stage, we analyze the correlation on all features. We do one by one of the features to observe the correlation between the one features with the others. This stage also can be used to evaluate the relevant features to determine the target. Further, we also do the visualization for this step.

Stage 5: evaluating the model

The last stage evaluates the model built by the combination of the discretization method and logistic regression using the accuracy level metric. The stage observes the highest performance of the model.

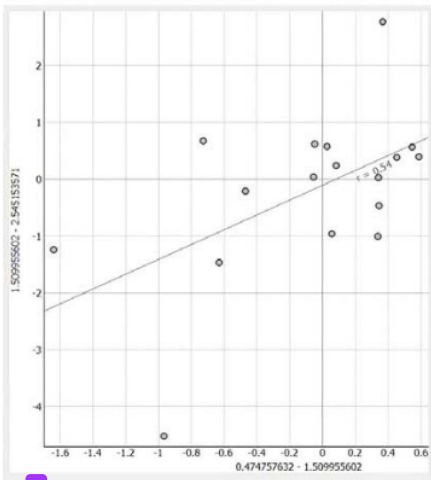


Fig. 4. Regression line with the highest correlation on discretization-4-interval-lasso

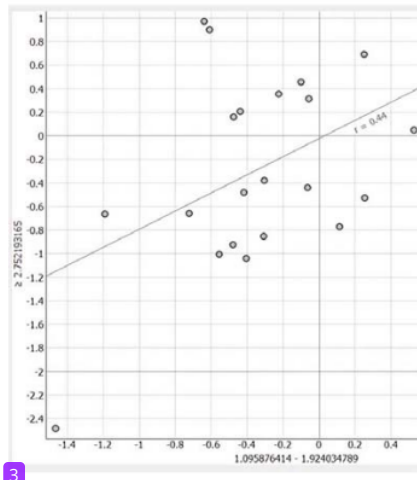


Fig. 6. Regression line with the highest correlation on discretization-5-interval-lasso

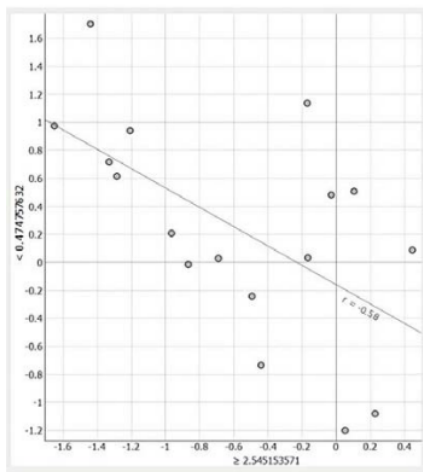


Fig. 5. Regression line with the highest correlation on discretization-4-interval-ridge

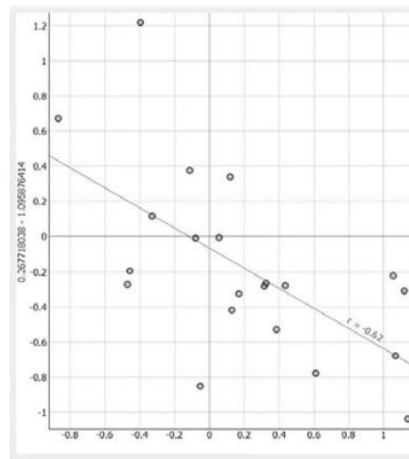


Fig. 7. Regression line with the highest correlation on discretization-5-interval-ridge

III. RESULT AND DISCUSSION

In this section, the proposed architecture is executed and then the result is analyzed. There are 2 sub-section explained, namely: the features correlation and the performance measurement. The performance of the model is evaluated using the accuracy level.

A. The correlation of features

The first session describes the correlation of features. After the discretization method is applied, we do combine the results of the discretization method using the logistic regression on all regulations.

For that, there are 6 combinations, namely: 3-interval-lasso, 3-interval-ridge, 4-interval-lasso, 4-interval-ridge, 5-interval-lasso and 5-interval-ridge.

Then, we compute the correlation of features on each combination which is stated with the r symbol. In here, we only visualize the highest results on the all combination depicted in the Fig.2-Fig.7. The general, the combinations of the discretization method, and the logistic regression on the ridge regulation correlate higher than others in all cases. In detail, the highest correlation is reached by 3-interval-lasso, namely: around $r = -0.72$. Contrarily, the lowest correlation of about $r = 0.44$ is achieved by 5-interval-lasso.

TABLE I. THE ACCURACY LEVEL OF THE COMBINATION OF LOGISTIC REGRESSION AND DISCRETIZATION METHOD WITH 3-INTERVAL

| Training set size | 3-interval-lasso | | | | | 3-interval-ridge | | | | |
|-------------------|------------------|-------------|-------------|-------------|-------------|------------------|-------------|-------------|------|------|
| | Repeat train | | | | | | | | | |
| | 2 | 3 | 5 | 10 | 20 | 2 | 3 | 5 | 10 | 20 |
| 10% | 85.3 | 85.3 | 85.3 | 85.3 | 85.3 | 84.8 | 85 | 84.9 | 84.7 | 84.5 |
| 20% | 84.6 | 84.6 | 84.6 | 84.6 | 84.6 | 84.6 | 84.6 | 84.6 | 84.4 | 84.4 |
| 30% | 85 | 85 | 85 | 85 | 85 | 85 | 85 | 84.8 | 84.9 | 84.9 |
| 40% | 85.3 | 85.3 | 85.3 | 85.3 | 85.3 | 85.3 | 85.3 | 85.3 | 85.1 | 85.2 |
| 50% | 84.2 | 84.2 | 84.2 | 84.2 | 84.2 | 84.2 | 83 | 85.3 | 83.9 | 83.9 |
| 60% | 84.8 | 84.8 | 84.8 | 84.8 | 84.8 | 84.8 | 84.8 | 84.8 | 84.8 | 84.8 |

TABLE II. THE ACCURACY LEVEL OF THE COMBINATION OF LOGISTIC REGRESSION AND DISCRETIZATION METHOD WITH 4-INTERVALS

| Training set size | 4-interval-lasso | | | | | 4-interval-ridge | | | | |
|-------------------|------------------|-------------|-------------|-------------|-------------|------------------|-------------|------|-----------|-----------|
| | Repeat train | | | | | | | | | |
| | 2 | 3 | 5 | 10 | 20 | 2 | 3 | 5 | 10 | 20 |
| 10% | 79.4 | 80.1 | 80.6 | 81 | 81.2 | 79.9 | 80.4 | 80.8 | 80.5 | 80.3 |
| 20% | 80.2 | 80.6 | 80.9 | 81.1 | 81.1 | 80.8 | 81 | 80.7 | 81 | 80.9 |
| 30% | 78.7 | 79.6 | 80.2 | 80.7 | 80.9 | 80.6 | 80.4 | 80.5 | 80.9 | 80.7 |
| 40% | 79.4 | 79.9 | 79.7 | 80 | 80.4 | 79.4 | 79.4 | 80 | 80.1 | 80.5 |
| 50% | 78.9 | 79.5 | 80 | 80.4 | 80.3 | 80.7 | 80.7 | 80.7 | 80.5 | 80.4 |
| 60% | 82.6 | 82.6 | 82.6 | 82.6 | 82.6 | 80.4 | 81.2 | 81.7 | 82 | 82 |

TABLE III. THE ACCURACY LEVEL OF THE COMBINATION OF LOGISTIC REGRESSION AND DISCRETIZATION METHOD WITH 5-INTERVALS

| Training set size | 5-interval-lasso | | | | | 5-interval-ridge | | | | |
|-------------------|------------------|-------------|------|-------------|------|------------------|-------------|------|------|------|
| | Repeat train | | | | | | | | | |
| | 2 | 3 | 5 | 10 | 20 | 2 | 3 | 5 | 10 | 20 |
| 10% | 77.5 | 77.5 | 77.3 | 77.5 | 77.3 | 76 | 76.8 | 76.5 | 76.7 | 76.5 |
| 20% | 75.3 | 73.6 | 74.9 | 75.9 | 76.3 | 75.8 | 75.1 | 75.8 | 76.3 | 76.5 |
| 30% | 76.9 | 76.7 | 76 | 75.7 | 75.7 | 76.9 | 76.7 | 76.7 | 76.5 | 76.4 |
| 40% | 76.5 | 76.5 | 75.6 | 76 | 76.2 | 76.5 | 76.5 | 75.6 | 75.9 | 75.9 |
| 50% | 77.2 | 76 | 76.5 | 76.3 | 76.8 | 76.3 | 76 | 76.1 | 76.1 | 76.7 |
| 60% | 76.1 | 76.1 | 76.1 | 76.1 | 76.1 | 76.1 | 76.1 | 76.1 | 75.9 | 75.9 |

Further, the 3-interval is illustrated in Fig.2-3. The correlation exploration of all features on the 3-interval is found that the highest correlation on logistic regression with lasso regulation. It is around $r = -0.72$. This condition occurs on the value of axis-x ≥ 2.200087581 and the value of axis-y < 0.819823672 . For logistic regression with ridge regulation, the highest correlation is about $r = -0.65$ with the value of axis-x ≥ 2.200087581 and on the value of axis-y < 0.8198234672 .

The 4-interval is illustrated in Fig.4-5. The correlation exploration of all features on the 4-interval is found that the highest correlation on logistic regression with lasso regulation. It is around $r = 0.54$. This condition occurs on the value of axis-x $0.474757632-1.509955602$ and the value of axis-y $= 1.509955602-2.545153571$. For logistic regression with ridge regulation, the highest correlation is about $r = -0.58$ with the value of axis-x ≥ 2.545153571 and axis-y < 0.474757632 .

The 5-interval is illustrated in Fig.6-7. The correlation exploration of all features on the 4-interval is found that the highest correlation on logistic regression with lasso regulation. It is around $r = 0.44$. This condition occurs on the value of axis-x $0.474757632-1.509955602$ and the value of axis-y $= 1.509955602-2.545153571$. For logistic regression with ridge regulation, the highest correlation is about $r = -0.62$ with the value of axis-x ≥ 2.545153571 and axis-y < 0.474757632 .

B. The performance measurement

The next discussion is about the model performance, in terms of level accuracy. Discretization methods at all intervals and all regulations are presented in TABLE I-III. This step is carried out to analyze which intervals and regulations have the best performance in logistic regression. Trials are conducted on a set of training sizes multiples of 10 from 10% -60% and repetition of train 2, 3, 5, 10, and 20.

The trial results of a combination of 3-interval discretization and logistic regression regulation of the lasso and ridge are presented in TABLE I. The table shows the highest accuracy level of 85.3% occurred in the trial scenario training size of 10% and 40% with all the repeat train settings for lasso regulation. Whereas in ridge regulation, the test scenario is the training set size of 40% with repeat train 2, 3, 5, and the training set size 50% with repeat train 5. Conversely, in discretization 3 intervals, logistic regression with lasso regulation experiences the lowest accuracy level of 84.2 % occurs in the trial scenario the training set size is 50% for all repeat trains. Whereas logistic regression with ridge regulation experienced the lowest accuracy level of 83% in the training set size of 60% and repeat train 3.

The experimental result of a combination of 4-interval and logistic regression with the regulation on the lasso and ridge is presented in TABLE II. The table shows the highest level of accuracy in both regulations that occurred in the training set= 60%. For lasso, the highest accuracy level is 82.6% on all repeat trains. The ridge reaches the highest accuracy level of

10 The authors wish to thank you for the support of the Department of Informatics, Universitas Negeri Surabaya, Indonesia on this research.

REFERENCES

- [1] R. Yunis and K. Telaumbanua, "Pengembangan E-Learning Berbasis LMS untuk Sekolah, Studi Kasus SMA/SMK di Sumatera Utara," *J. Nas. Tek. Elektro dan Teknol. Inf.*, vol. 6, no. 1, Feb. 2017.
- [2] P. Wanarti, E. Ismayanti, H. Peni, and Y. Yamasari, "The Enhancement of Teaching-Learning Process Effectiveness through The Development of Instructional Media Based on E-learning of Surabaya's Vocational Student," in *Proceedings of the 6th International Conference on Educational, Management, Administration and Leadership*, 2016, pp. 342–346.
- [3] U. Aich and S. Banerjee, "Application of teaching learning based optimization procedure for the development of SVM learned EDM process and its pseudo Pareto optimization," *Appl. Soft Comput.*, vol. 39, pp. 64–83, Feb. 2016.
- [4] Y. Yamasari, S. M. S. Nugroho, K. Yoshimoto, H. Takahashi, and M. H. Pumomo, "Identifying Dominant Characteristics of Students' Cognitive Domain on Clustering-based Classification," *Int. J. Intell. Eng. Syst.*, vol. 13, no. 1, 2020.
- [5] M. W. Rodrigues, L. E. Zárate, and S. Isotani, "Educational Data Mining: A review of evaluation process in the e-learning," *Telemat. Informatics*, May 2018.
- [6] P. Manirih and T. Ahmad, "Information hiding scheme for digital images using difference expansion and modulus function," *J. King Saud Univ. - Comput. Inf. Sci.*, vol. 31, no. 3, pp. 335–347, Jul. 2019.
- [7] L. C. Liñán, Á. Alejandro, and J. Pérez, "Educational Data Mining and Learning Analytics: differences, similarities, and time evolution Learning Analytics: Intelligent Decision Support Systems for Learning Environments," *RUSC. Univ. Knowl. Soc. J.*, vol. 12, no. 3, pp. 98–112, 2015.
- [8] A. Peña-Ayala, *Educational data mining: Applications and trends*. 2014.
- [9] N. T. Pang, M. Steinbach, and V. Kumar, *Introduction to Data mining*. 2006.
- [10] Y. Yamasari, S. M. S. Nugroho, R. Harimurti, and M. H. Pumomo, "Improving the cluster validity on student's psychomotor domain using feature selection," in *2018 International Conference on Information and Communications Technology (ICOLACT)*, 2018, pp. 460–465.
- [11] I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, 2016.
- [12] M. Kuhn and K. Johnson, *Feature Engineering and Selection: a Practical Approach for Predictive Models*. CRC Press LLC, 2019.
- [13] J. Dougherty, R. Kohavi, and M. Sahami, "Supervised and Unsupervised Discretization of Continuous Features," Morgan Kaufmann Publishers, 1995.
- [14] M. Sugiyama, *Introduction to Statistical Machine Learning*. Elsevier Inc., 2015.
- [15] A. Genkin, D. D. Lewis, and D. Madigan, "Large-Scale Bayesian Logistic Regression for Text Categorization," *Technometrics*, vol. 49, no. 3, pp. 291–304, Aug. 2007.
- [16] Y. Yamasari, S. Mardi, S. Nugroho, K. Yoshimoto, H. Takahashi, and M. H. Purnomo, "Expanding Tree-Based Classifiers Using Meta-Algorithm Approach: An Application for Identifying Students' Cognitive Level," *Int. J. Innov. Comput.*, vol. 15, no. 6, pp. 2085–2107, 2019.
- [17] R. Tibshirani, "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58. WileyRoyal Statistical Society, pp. 267–288, 1996.
- [18] B. F. Swindel, "Geometry of Ridge Regression Illustrated," *Am. Stat.*

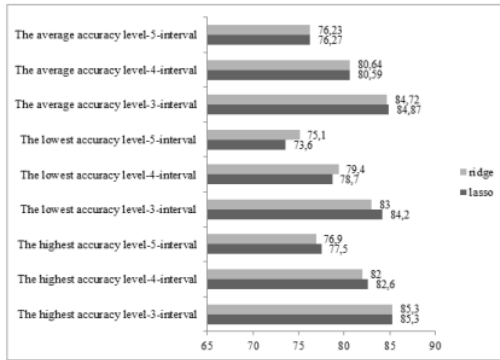


Fig. 8. The comparison of accuracy level on all intervals

82% on repeat trains 10 and 20. Meanwhile, the lowest accuracy level of 78.7% is experienced by Lasso regulations in the training set size of 30% with repeat train 2. For ridge, the lowest level of accuracy 79.4% is experienced on the training set size of 40% with repeat train 2 and 3.

The experiment result of a combination of 5-interval and logistic regression with the regulation on lasso and ridge is represented in TABLE III. The table shows the highest accuracy level of 77.5% occurring on the training set size of 10% in almost all repeat train settings for lasso regulation. Meanwhile, in ridge regulation, the highest accuracy level of 76.9% is achieved with a training set size of 30% and repeat train 2. Conversely, logistic regression with the Lasso regulation experience the lowest accuracy level of 73.6% occurred in the training set size of 20% and repeat train 3. While in logistic regression with ridge regulation, the lowest accuracy level is 75.1% in the training set size of 20% and repeat train 2.

The overall results of the experiment presented in TABLE I-III show special conditions, namely: achieving quite high accuracy with a small training set size (10%). This is possible because the data of students selected in the construction of the model are sufficiently representative. Moreover, the data selected for this training set is most likely to have the most influence on the target. Thus, this data is sufficient to represent the data as a whole although only slightly in size. To further clarify the results obtained, the average accuracy level, lowest accuracy level, and highest accuracy level are presented in Fig.8. The discretization method for 3-interval dominates the highest results among the others, and then followed by the discretization method for 4-interval and finally 5-interval in all regulations of logistic regression

IV. CONCLUSION

The Equal width interval can be combined with logistic regression with lasso and ridge regulation in the students' performance data. Among the intervals that have been carried out, the combination of discretization of 3 intervals in this realm and logistic regression of all regulations has been proven to achieve the best results, in terms of the highest level of accuracy.

vol. 35, no. 1, p. 12, Feb. 1981.

no. 4, p. 451, Nov. 1979.

- [19] N. R. Draper and R. C. van Nostrand, "Ridge Regression and James-Stein Estimation: Review and Comments," *Technometrics*, vol. 21,

Combining the Unsupervised

ORIGINALITY REPORT

12%

SIMILARITY INDEX

7%

INTERNET SOURCES

11%

PUBLICATIONS

5%

STUDENT PAPERS

PRIMARY SOURCES

| | | |
|---|--|----|
| 1 | repository.unair.ac.id Internet Source | 3% |
| 2 | Submitted to South Bank University Student Paper | 2% |
| 3 | Eero Helle. "Aerial census of ringed seals Pusa hispida basking on the ice of the Bothnian Bay, Baltic", Ecography, 1980 Publication | 1% |
| 4 | Yuni Yamasari, Anita Qoiriah, Hapsari P. A. Tjahyaningtjas, Ricky E. Putra, Agus Prihanto, Asmunin. "Improving the Quality of the Clustering Process on Students' Performance using Feature Selection", 2020 International Seminar on Application for Technology of Information and Communication (iSemantic), 2020 Publication | 1% |
| 5 | Yuni Yamasari, Naim Rochmawati, Anita Qoiriah, Dwi F. Suyatno, Tohari Ahmad. "Chapter 18 Reducing the Error Mapping of the Students' Performance Using Feature | 1% |

Selection", Springer Science and Business
Media LLC, 2021

Publication

| | | |
|----|--|------|
| 6 | repository.unesa.ac.id Internet Source | 1 % |
| 7 | Yuni Yamasari, Muhammad H. Garry, Supeno Mardi Susiki Nugroho, Mauridhi Hery Purnomo. "Enhancing the Classification Performance of Students Behavior on Serious Game using Discretization-based k-NN", 2019 IEEE International Conference on Engineering, Technology and Education (TALE), 2019 Publication | 1 % |
| 8 | Yuni Yamasari, Supeno M. S. Nugroho, I N. Sukajaya, Mauridhi H. Purnomo. "Features extraction to improve performance of clustering process on student achievement", 2016 International Computer Science and Engineering Conference (ICSEC), 2016 Publication | 1 % |
| 9 | Submitted to Universitas Negeri Surabaya The State University of Surabaya Student Paper | <1 % |
| 10 | Y. Yamasari, S. M. S. Nugroho, R. Harimurti, M. H. Purnomo. "Improving the cluster validity on student's psychomotor domain using feature selection", 2018 International | <1 % |

Conference on Information and Communications Technology (ICOIACT), 2018

Publication

11

Alfred Ultsch. "Optimizing time series discretization for knowledge discovery", Proceeding of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining - KDD 05 KDD 05, 2005

Publication

<1 %

12

ns2.thinkmind.org

Internet Source

<1 %

13

orca.cf.ac.uk

Internet Source

<1 %

14

dc.uwm.edu

Internet Source

<1 %

15

dokumen.pub

Internet Source

<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography On