# Implementation Regression and Naïve Bayes to Predict and Classify Data Asset at Educational Institutions

Handy Noviyarto

Faculty of Computer Science, Mercu Buana University, Indonesia

*Abstract*— *Regional assets represent regional assets which in essence belong to the respective provincial government. These government assets can play a role as collateral for regional development. The preparation of asset documents aims to safeguard assets from the aspect of regional administration. In this study, to process prediction analysis use the regression method, and to process classification use the Naive Bayes method. The purpose of this study was to predict and classify data asset that will be used and categorized according to regional planning using Regression and Naïve Bayes Method. This research was conducted using the Python programming language and the Visual Studio code.*

*Keywords*— *Regression, naïve bayes, clustering, data mining.*

## I. INTRODUCTION

Local assets are essentially regional wealth is owned by the provincial government each - each. One is a regional asset is an asset not move. As for which is included in the fixed assets to which such land or land, buildings, and so forth. In this aspect, it can play a role of government assets as collateral development in the region. Preparation of the document aims to secure the assets of the assets of the administrative aspects of the area.

According to the Government Accounting Standards (2016) assets are economic resources controlled or owned by the government as a result of past events and from which economic and social benefits in the future is expected can be obtained either by the government or the public, as well as can be measured, including non-financial resources needed to provide services to the general public and resources maintained for historical and cultural reasons.

Asset security aims to keep local assets do not change hands illegally and facilitate local authorities in managing further. Safeguarding assets is absolutely necessary by completing the assets referred to as legal documents. In addition, a regional asset wealth can act as a guarantee of regional development.

A common problem of the government's assets, which is not yet completed the document, even none at all. Not infrequently, the region's assets lost due to various reasons. Such as the lack of accuracy of the value of the assets being managed, the unclear status of the assets being managed, and others.

Based on the background of the issue, so in this study was taken the title "Implementation Regression and Naïve Bayes To Predict And Classify Data Asset".

## II. PLATFORM THEORY

### 2.1 Definition of Data Mining

Data mining is the process to obtain useful information from large data base warehouse. Techniques in Data Mining: how to search for the data that is to build a model. The model was used to identify the pattern of other data that are not in the data base stored.

### 2.2. Regression Analysis

Regression analysis in statistics is one method for determining the causal relationship between one variable and another variable (s). "Cause" variables are referred to by various terms: explanatory variables, explanatory variables, independent variables, or independently, variable X (because it is often depicted on the graph as abscissa, or the X-axis). Variables affected as a result are known as influenced variables, dependent variables, dependent variables, or Y variables. Both of these variables can be random variables (random), but the variables affected must always be random variables.

Regression analysis is one of the most popular and widely used analyzes. Regression analysis is widely used to make predictions and forecasts, with uses that complement each other in the field of machine learning. This analysis is also used to understand which independent variables are related to the dependent variable, and to find out the forms of the relationship.

### 2.3 Naïve Bayes

According Thomas Bayes, The Naive Bayes algorithm is a classification method using probability and statistical methods. The Naive Bayes algorithm predicts future opportunities based on past experience so it is known as the Bayes Theorem. The main characteristic of Naïve Bayes Classifier is a very strong assumption of independence from each condition / event.

The advantage of using this method is only requires a small amount of training data to determine the estimated parameters needed in the classification process. Because it is assumed to be an independent variable, only the variance of a variable in a class is needed to determine the classification, not the whole of the covariance matrix.

## III. RESULTS AND DISCUSSION

### 3.1 Regression Analysis Method

a.  Data Preprocessing

Before using the code for data processing, first input library that will be used:



We have made, then uflood some parts, such as the name of the dataset is loaded:



Figure 1. Data Pre Processing Regression Analysis

b.  Fitting Simple Linear Regression in Training Set



Figure 2. Fitting Simple Linear Regression in Training Set

Creating value coefficient and the intercept on predictive data tabulation.



Figure 3. Coefficient and Intercept

Coefficients and interception in a simple linear regression fit the parameters of the line. Given that this is a simple linear regression, with only two parameters, and knowing that the parameter is the intercept and the slope of the line, can sklearn direct estimate of the data.

c.  Predicting Results of Test-Set

And to know the quality of the data is carried out also for testing against test data



Figure 4. Predicting Result

And included the form of the plot



Figure 5. Regression Analysis

3.2  *Naïve Bayes Method*

a.  Selection Data

23

Before using the code for data processing, first input library that will be used:



Figure 6. Selection Data

b. Preprocessing/Cleaning Data

Before the data is used, cleansing data will be taken:



Figure 7. Preprocessing/Cleaning Data

c. Data Transformation



Figure 8. Data Transformation

Classification Data Naïve Bayes Method
1. Determining the predictive data



Figure 9. Determining Predictive Data

2. Determine the probability of data



Figure 10. Determine the probability data

3. Determining Matrix Model:



Figure 11. Determining Matrix Model

d.   Specifying Process Data Mining
Apply an algorithm to classify the data.

```
[29]  1 # Import train_test_split function
      2 from sklearn.model_selection import train_test_split
      3 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 123)
```

```
[30]  1 # Import Gaussian Naive Bayes model
      2 from sklearn.naive_bayes import GaussianNB
      3
      4 # Mengaktifkan/memanggil/membuat fungsi klasifikasi Naive bayes
      5 modelnb = GaussianNB()
      6
      7 # Memasukkan data training pada fungsi klasifikasi naive bayes
      8 nbtrain = modelnb.fit(x_train, y_train)
      9 nbtrain.class_count_
```

```
array([  1.,   2.,   1.,   1.,   4.,   1.,   2.,   2.,   2.,   6.,   3.,
         5.,   6.,   5.,   6.,   9.,  11.,  10.,  83., 234.,  93.,  64.,
        44.,  20.,   5.,   3.,   1.,   6.,   7.,   4.,   1.,  11.,  26.,
        26.,  16.,   3.,   4.,   4.,   9.,   6.,   3.,   2.,   7.,  24.,
         1.,   5.,   3.,   3.,   1.,   2.,   2.])
```

|       | precision | recall | f1-score | support |
|-------|-----------|--------|----------|---------|
| 1962  | 0.00      | 0.00   | 0.00     | 2       |
| 1963  | 0.00      | 0.00   | 0.00     | 1       |
| 1975  | 0.00      | 0.00   | 0.00     | 1       |
| 1976  | 0.00      | 0.00   | 0.00     | 3       |
| 1977  | 0.00      | 0.00   | 0.00     | 2       |
| 1978  | 0.00      | 0.00   | 0.00     | 3       |
| 1979  | 0.00      | 0.00   | 0.00     | 3       |
| 1981  | 0.00      | 0.00   | 0.00     | 2       |
| 1982  | 0.00      | 0.00   | 0.00     | 1       |
| 1983  | 0.00      | 0.00   | 0.00     | 27      |
| 1984  | 0.00      | 0.00   | 0.00     | 56      |
| 1985  | 0.00      | 0.00   | 0.00     | 11      |
| 1986  | 0.08      | 1.00   | 0.16     | 16      |
| 1987  | 0.00      | 0.00   | 0.00     | 10      |
| 1988  | 0.00      | 0.00   | 0.00     | 5       |
| 1989  | 0.00      | 0.00   | 0.00     | 3       |
| 1990  | 0.00      | 0.00   | 0.00     | 0       |
| 1991  | 0.00      | 0.00   | 0.00     | 2       |
| 1993  | 0.00      | 0.00   | 0.00     | 1       |
| 1994  | 0.00      | 0.00   | 0.00     | 4       |
| 1996  | 0.00      | 0.00   | 0.00     | 5       |
| 1997  | 0.00      | 0.00   | 0.00     | 6       |
| 1998  | 0.00      | 0.00   | 0.00     | 9       |
| 1999  | 0.00      | 0.00   | 0.00     | 2       |
| 2000  | 0.00      | 0.00   | 0.00     | 2       |
| 2001  | 0.00      | 0.00   | 0.00     | 1       |
| 2003  | 0.00      | 0.00   | 0.00     | 3       |
| 2004  | 0.00      | 0.00   | 0.00     | 8       |
| 2005  | 0.00      | 0.00   | 0.00     | 2       |
| 2008  | 0.00      | 0.00   | 0.00     | 5       |
| 2010  | 0.00      | 0.00   | 0.00     | 3       |
| 2014  | 0.00      | 0.00   | 0.00     | 1       |
| accuracy     |      |        | 0.08     | 200     |
| macro avg    | 0.00 | 0.03   | 0.00     | 200     |
| weighted avg | 0.01 | 0.08   | 0.01     | 200     |

Figure 12. Specifying Process Data Mining

```
[56]  1 plt.figure(figsize = (17,6))
      2 plt.scatter(y_test,y_pred, c='red', s=300, alpha=0.1 , marker="o")
      3 plt.xticks(data['Tahun'],rotation=45)
      4 plt.xlabel('Tahun',fontsize=18)
      5 plt.ylabel('Tahun',fontsize=18)
      6 plt.show()
```
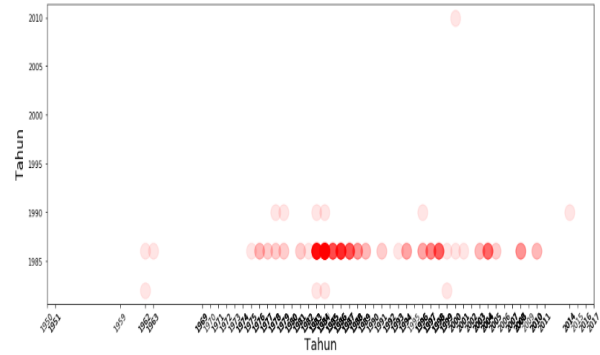


Figure 13. K-Means

e.   Application of Interpretation / Evaluation

```
[38]  1 # Menghitung nilai akurasi dari klasifikasi naive bayes
      2 from sklearn.metrics import classification_report
      3 print(classification_report(y_test,y_pred))
```

Figure 14. Interpretation/Evaluation

## IV.   CONCLUSION

Based on the discussion above, it can be concluded that the asset data influences the utilization of the performance program. Naïve Bayes method successfully classifies 900 data from 1000 data tested. So the Naïve Bayes method succeeded in making the accuracy of the data percentage of accuracy 71.42%.

## REFERENCES

[1]   B.Liu. Sentiment Analysis and Opinion Mining. San Rafael : Morgan and Claypool Publishers. 2012
[2]   Han, J., Kamber, M., & Pei, J. Data Mining Concepts and Techniques (Third Edition)Elsevier Inc. 2012
[3]   M.W.Berry and J.Kogan. Text Mining Analysis and Theory. Wiley:United Kingdom.2010
[4]   O. Maimon and L. Rokach, Data Mining and Knowledge Discovery Handbo- ok. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2005.
[5]   S. Russell and P. Norvig, Artificial Intelligence A Modern Approach. Upper Saddle River, New Jersey 07458: Pearson Education, Inc., 3 ed., 2010
[6]   Tan, M. Steinbach, and V. Kumar, Introduction to Data Mining, (First Edition). Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 2005.
[7]   X. Wu and V. Kumar, eds., The Top Ten Algorithms in Data Mining.Chapman and Hall, 2009
[8]   Nia Rahma Kurnianda & Yunita Sartika Sari. *Analysis and Design of Information System for Self-Journal on Food Based Dietary Assessment Record for Diabetes Patients*. International Research Journal of Computer Science (IRJCS). Volume 06 Issue 5. 2018
[9]   Suhendra & Ranggadara, Indra. Naïve Bayes Algorithm with Chi Square and NGram Feature for Reviewing Laptop Product on Amazon Site. International Research Journal of Computer Science, Vol 4, issue 12.2017
[10]  Triana, Yaya Sudarya, and Astari Retnowardhani. "Enhance interval width of crime forecasting with ARIMA model-fuzzy alpha cut." Telkomnika 17.3 (2019): 1193-1201.