

# IMPLEMENTATION REGRESSION AND NAÏVE BAYES TO PREDICT AND CLASSIFY DATA ASSET AT EDUCATIONAL INSTITUTIONS

# Handy Noviyarto

# 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

## **1. 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 nonfinancial 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".

## 2. 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.

Faculty of Computer Science, Mercu Buana University, Indonesia

#### HOW TO CITE THIS ARTICLE:

Handy Noviyarto (2020). Implementation Regression And Naïve Bayes To Predict And Classify Data Asset At Educational Institutions, International Educational Journal of Science and Engineering (IEJSE), Vol: 3, Issue: 3, 01-05

Copyright© 2020, IEJSE. This open-access article is published under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License which permits Share (copy and redistribute the material in any medium or format) and Adapt (remix, transform, and build upon the material) under the Attribution-NonCommercial terms.

## 2.3 Naive 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.

## **3 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:

OUT[19]:	Ка	de_Barang	KIB				Jenis_Barang	Luas	Satuan	Alamat	Tahun
	Û	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	5.190	M2	JI. Cikini Raya No 87	1950
	1	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	1.800	M2	JALAN PERUK NO. 32	1951
	2	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	4.180	M2	JI.Perwira No.10 Rt.002 / Rw.008	1951
	3	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	9.700	M2	JL PALMERAH BARAT NO. 59	1959
	4	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	775.000	M2	JI. Masjid Nur No.33 Rt.002 / 010	1962
	5	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	2.400	M2	Jin. Medan Merdeka Timur 14	1962
	6	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	4.808	M2	jl.Pertanian Klender	1963
	7	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	2.363	M2	JI Bulak Timur I/7 Klender	1963
	8	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	2.900	M2	JI. F. No.1 Kebon Baru	1969
	9	1011104002	KIB D	Tanah Ba	ngunan Pendidika	an Dan Li	atihan (sekolah)	1.730	M2	JI. Utama Raya No. 41	1969
n [20]: M	datas	et.descri	be()								
Out[20]:		Kode_Bara	ng	Luas	Tahun						
	count	1.000000e+	63 10	00.00000.00	1000.000000						
	mean	1.011084e+	-09	21.916909	1987.657000						
	std	1.407996e+	-05 1	13.472694	8.608685						
	min	1.010101e+	-09	1.057000	1950.000000						
	25%	1.011104e+	-09	2.129500	1984.000000						
	50 %	1.011104e+	-09	3.035000	1985.000000						

max 1.011301e+09 976.000000 2017.000000

Figure 1: Data Pre Processing Regression Analysis

b. Fitting Simple Linear Regression in Training Set

3.	Library Train - Test Dataset
In [38]: M	msk = np.random.rand(lan(df)) < 0.5 train = zaset[msk] test = zaset[msk]
4.	Fitting Simple Linear Regression pada Training-Set
In [72]: H	<pre>from sklearn import linear_model regr = linear_model.LinearRegression() train_x = np.asanyarray(train[['Tahun']]) train_y = np.asanyarray(train[['Kde_Bararg']]) regr.fit (train_x, train_y)</pre>
Out[72]:	LinearRepression(copy_X=True, fit_intercept=True, n_jobs=Wone, normalize=False)
In [73]: H	<pre>plt.scatter(train.Tahun, train.Kode_Barang, color='green') plt.plot(train_x, regr.coef_[0][0]*train_x + regr.intercept_[0], '-r') plt.slabe1("Kode_Barang") plt.stile("Kode_Barang")</pre>
Out[73]:	Text(0.5, 1.0, 'Analisa Data Train')
	109 Analisa Data Train
	10112 -
	10110 -
	5 10108 ·
	8 10106 -
	1.0104 -
	1.0102 -
	1950 1960 1970 1980 1990 2000 2010 2020

Figure 2: Fitting Simple Linear Regression in Training Set

Creating value coefficient and the intercept on predictive data tabulation.

[60]: M	<pre># The coefficients print ('Coefficients: ', regr.coef_) print ('Intercept: ',regr.intercept_)</pre>	
	Coefficients: [[-304.93274936]] Intercept: [1.0116871e+09]	

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

In

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

In [80]:	<pre>M g_pred = regr.predict(test_x) print(g_pred)</pre>
	[[1.01107519++09]
	[1.91163917-469]
	[1.9190917499]
	[1.0110000 eres]
	[1.01190205-00]
	[1.01109205-+09]
	1.01109264e+091
	[1.01109264e+09]
	[1.01109322=+09]
	(1.01109322=+09)
	[1.01109379=+09]
	[1.01109379=+09]
	[1.01109495e+09]
	[1.01109553=+09]
	[1.0109553=+09]
	[1.011095536409]
	[1.0109508409]
	[1.01105016405]
	[1.0109000009]
In [81]:	<pre>M x_pred = regr.predict(test_y) print(x_pred)</pre>
	[3.020000000+11]
	[5.85206296e+11]
	[5.85206296+11]
	[5.85296296e+11]
	[5.85296296+11]
	5. 55265296411
	5.85266296411
	[5.05202/04/11] [5.07202014/04/11]
	[5.02/02/02/02/11] [5.02/02/02/02/11]
	[5.05260260011]
	[5:0530557345573451]
	[5,8526626541]
	[5,85296296e+11]
	5.85286296+11
	LE DEDUCATE ALL

**Figure 4 : Predicting Result** 

## And included the form of the plot



**Figure 5: Regression Analysis** 

## 3.2 Naïve Bayes Method

#### a. Selection Data

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

[12]	1 import pandas as pd 2 import numpy as np
[13]	<pre>1 from google.colab import drive 2 drive.mount('/content/drive')</pre>

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

[14] 1 data = pd.read\_csv('/content/drive/Ny Drive/Project/Data/Data\_aset\_II\_3.csv', sep=';')

#### **Figure 6: Selection Data**

#### b. Preprocessing/Cleaning Data

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

[15] 1 data.head()

Ŀ		Kode_Barang	KIB	Jenis_Barang	Ukuran	Satuan	Alamat	Tahun
	0	1011104002	KIB D	Tanah Bangunan Pendidikan Dan Latihan (sekolah)	5.19	M2	JI. Cikini Raya No 87	1950
	1	1011104002	KIB D	Tanah Bangunan Pendidikan Dan Latihan (sekolah)	1.80	M2	JALAN PERUK NO. 32	1951
	2	1011104002	KIB D	Tanah Bangunan Pendidikan Dan Latihan (sekolah)	4.18	M2	JI.Perwira No.10 Rt.002 / Rw.008	1951
	3	1011104002	KIB D	Tanah Bangunan Pendidikan Dan Latihan (sekolah)	9.70	M2	JL PALMERAH BARAT NO. 59	1959
	4	1011104002	KIB D	Tanah Bangunan Pendidikan Dan Latihan (sekolah)	775.00	M2	JI. Masjid Nur No.33 Rt.002 / 010	1962
[57]	1	# Mengecek ap data.emotv	bakah a	da deret yang kosong				
Ŀ	Fals	se						
[58]	1	# Melihat uku data.size	uran da	ri data				

C+ 7000



## c. Data Transformation

[59]	1 2 3	# Menetapkan x = data.drop x.head()	Variabel D(['Tahun	independen ', 'KIB', 'Satu	uan','Jenis_	Barang', '	Alamat'], a	exis = 1)		
C•		Kode_Barang	Ukuran							
	0	1011104002	5.19							
	1	1011104002	1.80							
	2	1011104002	4.18							
	3	1011104002	9.70							
	4	1011104002	775.00							
[60]	1 2 3 4	# Menetapkan y = data['Tał y.head()	Variabel nun']	independen						
C•	0 1 2 3 4 Nam	1950 1951 1951 1959 1962 e: Tahun, dty	/pe: int6	4						

## **Figure 8: Data Transformation**

## **Classification Data Naïve Bayes Method**

- 1. Determining the predictive data
- [31] 1 # Menentukan hasil prediksi dari x\_test
  - 2 y\_pred = nbtrain.predict(x\_test)
  - 3 y\_pred
- array([1986, 2010, 1986, 1986, 1986, 1986, 1982, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1982, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1990, 1986, 1986, 1986, 1986, 1990, 1986, 1986, 1986, 1982, 1986, 1990, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1990, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1982, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1986, 1990, 1986, 1986, 1990, 1986])

**Figure 9: Determining Predictive Data** 

## 2. Determine the probability of data

# [32] 1 # Menentukan probabilitas hasil prediksi 2 nbtrain.predict\_proba(x\_test)

C•	array([[0.00501279, 0.01233699, 0.00127184,, 0.	, 0.01218608,	
	0.01033674],		
	[0.00566488, 0.0119682 , 0.00196036,, 0.	, 0.01160836,	
	0.01085605],		
	[0.00493076, 0.01237628, 0.00120303,, 0.	, 0.01225602,	
	0.01027391],		
	,		
	[0.00636245, 0.01145526, 0.00309225,, 0. 0.01146746],	, 0.01093496,	
	[0.00538868, 0.01213673, 0.00163482,, 0.	, 0.01185809,	
	0.01063127],		
	[0.00725134, 0.0105622 , 0.00567666,, 0.	, 0.00992861,	
	0.01239659]])		

## Figure 10: Determine the probability data

#### 3. Determining Matrix Model:

[33]	<pre>1 # import confusion_matrix model 2 from sklearn.metrics import confusion_matrix 3 confusion_matrix(y_test, y_pred)</pre>	
Ŀ	array([[0, 0, 0,, 0, 0, 0], [0, 0, 0,, 0, 0, 0], [0, 0, 0,, 0, 0, 0], , [0, 0, 0,, 0, 0, 0], [0, 0, 0,, 0, 0, 0], [0, 0, 0,, 0, 0, 0]])	
[37]	<pre>1 # Merapikan hasil confusion matrix 2 y_actual1 = pd.Series([1, 0,1,0,1,0,1,0,1,0,0,1,1,0,0], name ='actual 3 y_pred1 = pd.Series([1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1], 4 df_confusion = pd.crosstab(y_actual1, y_pred1) 5 df_confusion</pre>	') name ='prediction')

## C\* prediction 0 1

actual

0 7 2

1 1 8

#### **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/menanggil/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_
[p 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.,
25., 15., 3., 4., 4., 9., 6., 3., 2., 7., 24.,
1., 5., 3., 3., 1., 2., 2.])
precision recall f1-score support
```

1	962	0.00	0.00	0.00	2
1	963	0.00	0.00	0.00	1
1	975	0.00	0.00	0.00	1
1	976	0.00	0.00	0.00	з
1	977	0.00	0.00	0.00	2
1	978	0.00	0.00	0.00	3
1	979	0.00	0.00	0.00	3
1	981	0.00	0.00	0.00	2
1	982	0.00	0.00	0.00	1
1	983	0.00	0.00	0.00	27
1	984	0.00	0.00	0.00	56
1	985	0.00	0.00	0.00	11
1	986	0.08	1.00	0.16	16
1	987	0.00	0.00	0.00	10
1	988	0.00	0.00	0.00	5
1	989	0.00	0.00	0.00	3
1	990	0.00	0.00	0.00	0
1	991	0.00	0.00	0.00	2
1	993	0.00	0.00	0.00	1
1	994	0.00	0.00	0.00	4
1	996	0.00	0.00	0.00	5
1	997	0.00	0.00	0.00	6
1	998	0.00	0.00	0.00	9
1	999	0.00	0.00	0.00	2
2	000	0.00	0.00	0.00	2
2	001	0.00	0.00	0.00	1
2	003	0.00	0.00	0.00	3
2	004	0.00	0.00	0.00	8
2	005	0.00	0.00	0.00	2
2	008	0.00	0.00	0.00	5
2	010	0.00	0.00	0.00	3
2	014	0.00	0.00	0.00	1
accur	acy			0.08	200
macro	avg	0.00	0.03	0.00	200
ighted	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="0")
3 plt.xticks(data['Tahun'],rotation=45)
4 plt.xlabel('Tahun',fontsize=18)
5 plt.ylabel('Tahun',fontsize=18)
6 plt.show()



e. Application of Interpretation / Evaluation



3 print(classification\_report(y\_test,y\_pred))

#### **Figure 14: Interpretation/Evaluation**

#### 4. 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%.

#### **5. REFERENCE**

- 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
- M.W.Berry and J.Kogan. Text Mining Analysis and Theory. Wiley:United Kingdom.2010
- O. Maimon and L. Rokach, Data Mining and Knowledge Discovery Handbo- ok. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2005.
- S. Russell and P. Norvig, Artificial Intelligence A Modern Approach. Upper Saddle River, New Jersey 07458: Pearson Education, Inc., 3 ed., 2010
- 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
- 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
- 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.