

Tugas LAB WEEK 4 Christopher Darren

Jupyter 00000054804_Christopher Darren_Week4 Last Checkpoint: 5 minutes ago (autosaved) Python 3 (ipykernel)

File Edit View Insert Cell Kernel Widgets Help

Code

Data Mining Methodology 2

```
In [1]: 1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
```

Linear Regression

```
In [2]: 1 #independent_var merupakan variable yang mempengaruhi dependent_var
2 #independent_var sebagai data x, dependent_var sebagai data y
3
4 independent_var = np.array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
5 dependent_var = np.array([1, 3, 2, 5, 7, 8, 8, 9, 10, 12])
```

```
In [3]: 1 #number of observations/points
2 n = np.size(independent_var)
```

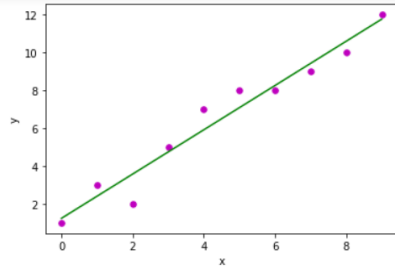
```
In [4]: 1 #menghitung rata - rata (mean) untuk menghitung cross deviation dan squared deviation
2
3 # mean of x and y vector
4 mean_indep_var = np.mean(independent_var)
5 mean_dep_var = np.mean(dependent_var)
6
7 #fungsi cross deviation dan squared deviation adalah untuk membuat garis linear yang menggambarkan posisi
8 #dependent variable yang akan diprediksi
9
10 #calculating cross-deviation and deviation about x
11 cross_deviation = np.sum(dependent_var * independent_var) - n*mean_dep_var*mean_indep_var
12 squared_deviation = np.sum(independent_var*independent_var) - n*mean_indep_var*mean_indep_var
```

Gambar 1.

```
In [5]: 1 #calculating regression coefficients
2 b_1 = cross_deviation / squared_deviation
3 b_0 = mean_dep_var - b_1*mean_indep_var
4
5 #plotting the actual points as scatter plot
6 plt.scatter(independent_var, dependent_var, color = 'm', marker = "o", s = 30)
7
8 # predicted response vector
9 y_pred = b_0 + b_1*independent_var
```

```
In [6]: 1 # plotting the regression Line
2 plt.plot(independent_var, y_pred, color = 'g')
3
4 # putting labels
5 plt.xlabel('x')
6 plt.ylabel('y')
7
8 #function to show plot
9 plt.scatter(independent_var, dependent_var, color = "m", marker = "o", s = 30)
10 plt.show()
```

Gambar 2.



1. Penerapan Linear Regression pada Weather untuk memprediksi temperatur suhu

```
In [7]: 1 from sklearn.linear_model import LinearRegression
        2 from sklearn.model_selection import train_test_split
```

```
In [8]: 1 #reading data
        2 dataset = pd.read_csv(r"D:\SEMESTER 4\IS411 Data Modelling\LAB\Bahan Modul 4\weather.csv", header=0)
        3 dataset.info()
        4 dataset
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 366 entries, 0 to 365
Data columns (total 22 columns):
#   Column                Non-Null Count  Dtype
---  -
0   MinTemp                366 non-null   float64
1   MaxTemp                366 non-null   float64
2   Rainfall               366 non-null   float64
3   Evaporation            366 non-null   float64
4   Sunshine               366 non-null   float64
5   WindGustDir            366 non-null   object
6   WindGustSpeed          366 non-null   float64
7   WindDir9am            366 non-null   object
8   WindDir3pm            366 non-null   object
9   WindSpeed9am          366 non-null   float64
10  Humidity3pm           366 non-null   int64
11  Pressure9am           366 non-null   float64
12  Pressure3pm           366 non-null   float64
13  Cloud9am              366 non-null   int64
14  Cloud3pm              366 non-null   int64
15  Temp9am               366 non-null   float64
16  Temp3pm               366 non-null   float64
17  RainToday             366 non-null   object
18  RISK_MM               366 non-null   float64
19  RainTomorrow          366 non-null   object
dtypes: float64(12), int64(5), object(5)
memory usage: 63.0+ KB
```

Gambar 3.

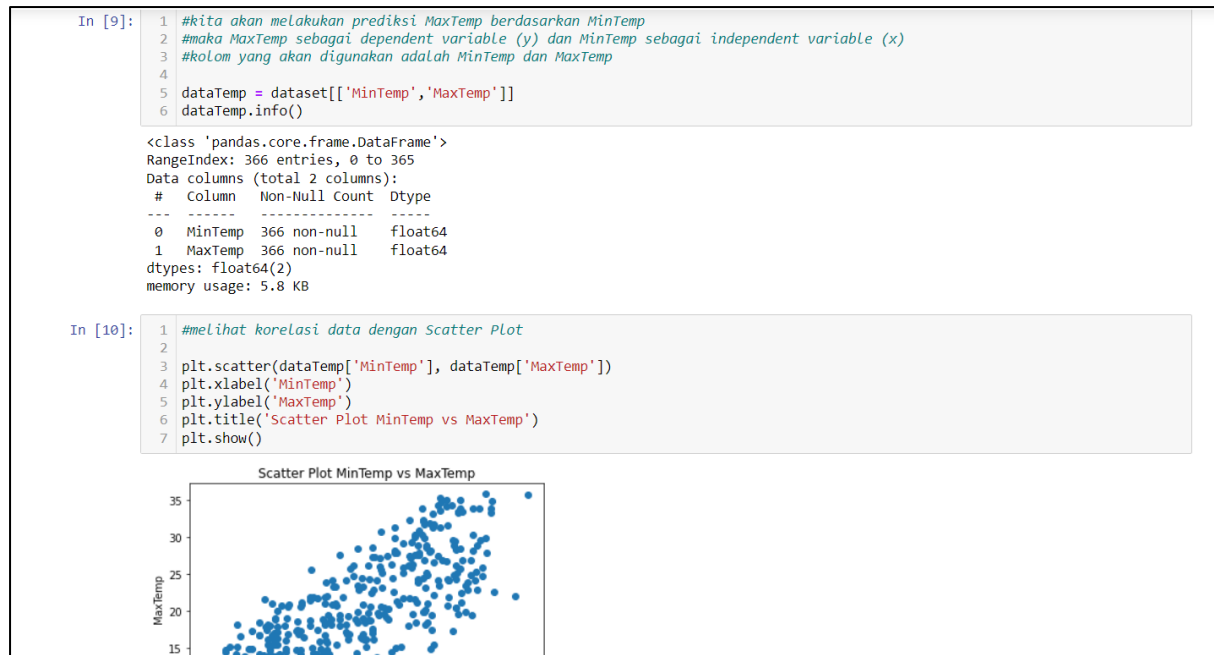
```
10 WindSpeed3pm 366 non-null int64
11 Humidity9am 366 non-null int64
12 Humidity3pm 366 non-null int64
13 Pressure9am 366 non-null float64
14 Pressure3pm 366 non-null float64
15 Cloud9am 366 non-null int64
16 Cloud3pm 366 non-null int64
17 Temp9am 366 non-null float64
18 Temp3pm 366 non-null float64
19 RainToday 366 non-null object
20 RISK_MM 366 non-null float64
21 RainTomorrow 366 non-null object
dtypes: float64(12), int64(5), object(5)
memory usage: 63.0+ KB
```

Out[8]:

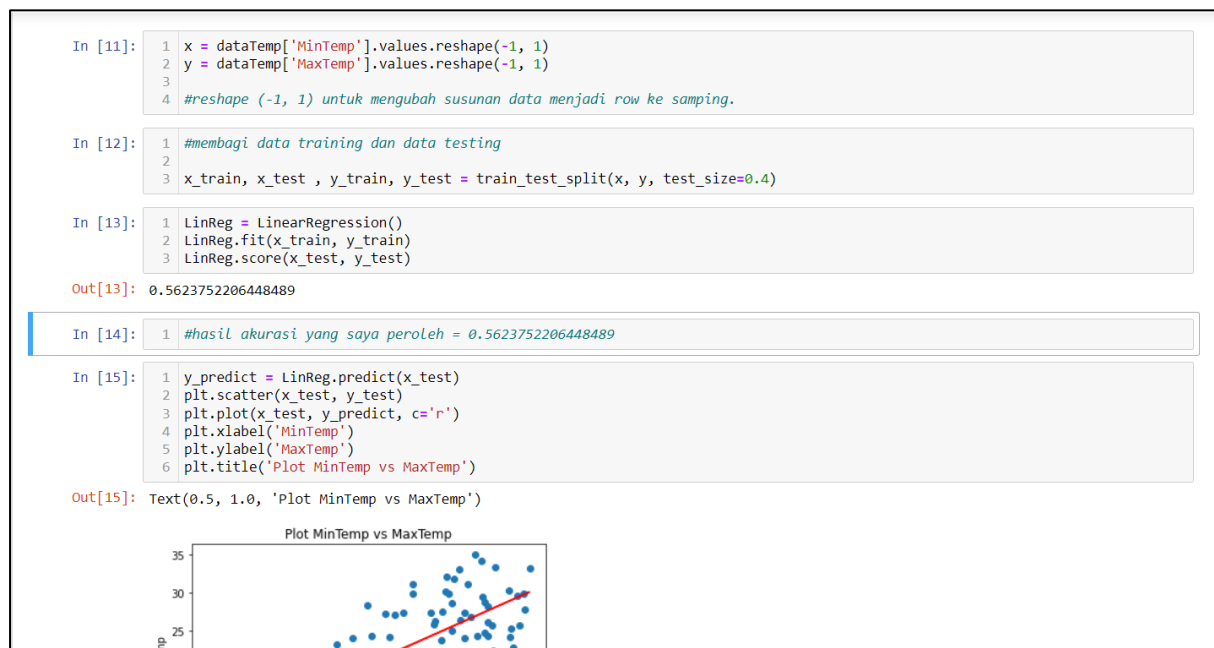
	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed9am	...	Humidity3pm	Pressure9am
0	8.0	24.3	0.0	3.4	6.3	NW	30.0	SW	NW	6.0	...	29	10
1	14.0	26.9	3.6	4.4	9.7	ENE	39.0	E	W	4.0	...	36	10
2	13.7	23.4	3.6	5.8	3.3	NW	85.0	N	NNE	6.0	...	69	10
3	13.3	15.5	39.8	7.2	9.1	NW	54.0	WNW	W	30.0	...	56	10
4	7.6	16.1	2.8	5.6	10.6	SSE	50.0	SSE	ESE	20.0	...	49	10
...
361	9.0	30.7	0.0	7.6	12.1	NNW	76.0	SSE	NW	7.0	...	15	10
362	7.1	28.4	0.0	11.6	12.7	N	48.0	NNW	NNW	2.0	...	22	10
363	12.5	19.9	0.0	8.4	5.3	ESE	43.0	ENE	ENE	11.0	...	47	10
364	12.5	26.9	0.0	5.0	7.1	NW	46.0	SSW	WNW	6.0	...	39	10
365	12.3	30.2	0.0	6.0	12.6	NW	78.0	NW	WNW	31.0	...	13	10

366 rows x 22 columns

Gambar 4.



Gambar 5



Gambar 6.

Logistic Regression

```
In [16]: 1 #copy dataset
2 dataLogistic = dataset.copy()
3 dataLogistic.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 366 entries, 0 to 365
Data columns (total 22 columns):
#   Column              Non-Null Count  Dtype
---  -
0   MinTemp              366 non-null   float64
1   MaxTemp              366 non-null   float64
2   Rainfall             366 non-null   float64
3   Evaporation          366 non-null   float64
4   Sunshine             363 non-null   float64
5   WindGustDir          363 non-null   object
6   WindGustSpeed        364 non-null   float64
7   WindDir9am           335 non-null   object
8   WindDir3pm           365 non-null   object
9   WindSpeed9am         359 non-null   float64
10  WindSpeed3pm         366 non-null   int64
11  Humidity9am          366 non-null   int64
12  Humidity3pm          366 non-null   int64
13  Pressure9am          366 non-null   float64
14  Pressure3pm          366 non-null   float64
15  Cloud9am             366 non-null   int64
16  Cloud3pm             366 non-null   int64
17  Temp9am              366 non-null   float64
18  Temp3pm              366 non-null   float64
19  RainToday            366 non-null   object
20  RISK_MM              366 non-null   float64
21  RainTomorrow         366 non-null   object
dtypes: float64(12), int64(5), object(5)
memory usage: 63.0+ KB
```

Gambar 7.

```
In [17]: 1 #mengatasi null
2 dataLogistic['WindGustSpeed'] = dataLogistic['WindGustSpeed'].fillna(0)
3 dataLogistic['WindSpeed9am'] = dataLogistic['WindSpeed9am'].fillna(0)
4 dataLogistic['Sunshine'] = dataLogistic['Sunshine'].fillna(0)
5 dataLogistic.info()
6 dataLogistic

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 366 entries, 0 to 365
Data columns (total 22 columns):
#   Column              Non-Null Count  Dtype
---  -
0   MinTemp              366 non-null   float64
1   MaxTemp              366 non-null   float64
2   Rainfall             366 non-null   float64
3   Evaporation          366 non-null   float64
4   Sunshine             366 non-null   float64
5   WindGustDir          363 non-null   object
6   WindGustSpeed        366 non-null   float64
7   WindDir9am           335 non-null   object
8   WindDir3pm           365 non-null   object
9   WindSpeed9am         366 non-null   float64
10  WindSpeed3pm         366 non-null   int64
11  Humidity9am          366 non-null   int64
12  Humidity3pm          366 non-null   int64
13  Pressure9am          366 non-null   float64
14  Pressure3pm          366 non-null   float64
15  Cloud9am             366 non-null   int64
16  Cloud3pm             366 non-null   int64
17  Temp9am              366 non-null   float64
18  Temp3pm              366 non-null   float64
19  RainToday            366 non-null   object
20  RISK_MM              366 non-null   float64
21  RainTomorrow         366 non-null   object
dtypes: float64(12), int64(5), object(5)
memory usage: 63.0+ KB
```

Gambar 8

Out[17]:

	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed9am	...	Humidity3pm	Pressure!
0	8.0	24.3	0.0	3.4	6.3	NW	30.0	SW	NW	6.0	...	29	10
1	14.0	26.9	3.6	4.4	9.7	ENE	39.0	E	W	4.0	...	36	10
2	13.7	23.4	3.6	5.8	3.3	NW	85.0	N	NNE	6.0	...	69	10
3	13.3	15.5	39.8	7.2	9.1	NW	54.0	WNW	W	30.0	...	56	10
4	7.6	16.1	2.8	5.6	10.6	SSE	50.0	SSE	ESE	20.0	...	49	10
...
361	9.0	30.7	0.0	7.6	12.1	NNW	76.0	SSE	NW	7.0	...	15	10
362	7.1	28.4	0.0	11.6	12.7	N	48.0	NNW	NNW	2.0	...	22	10
363	12.5	19.9	0.0	8.4	5.3	ESE	43.0	ENE	ENE	11.0	...	47	10
364	12.5	26.9	0.0	5.0	7.1	NW	46.0	SSW	WNW	6.0	...	39	10
365	12.3	30.2	0.0	6.0	12.6	NW	78.0	NW	WNW	31.0	...	13	10

366 rows x 22 columns

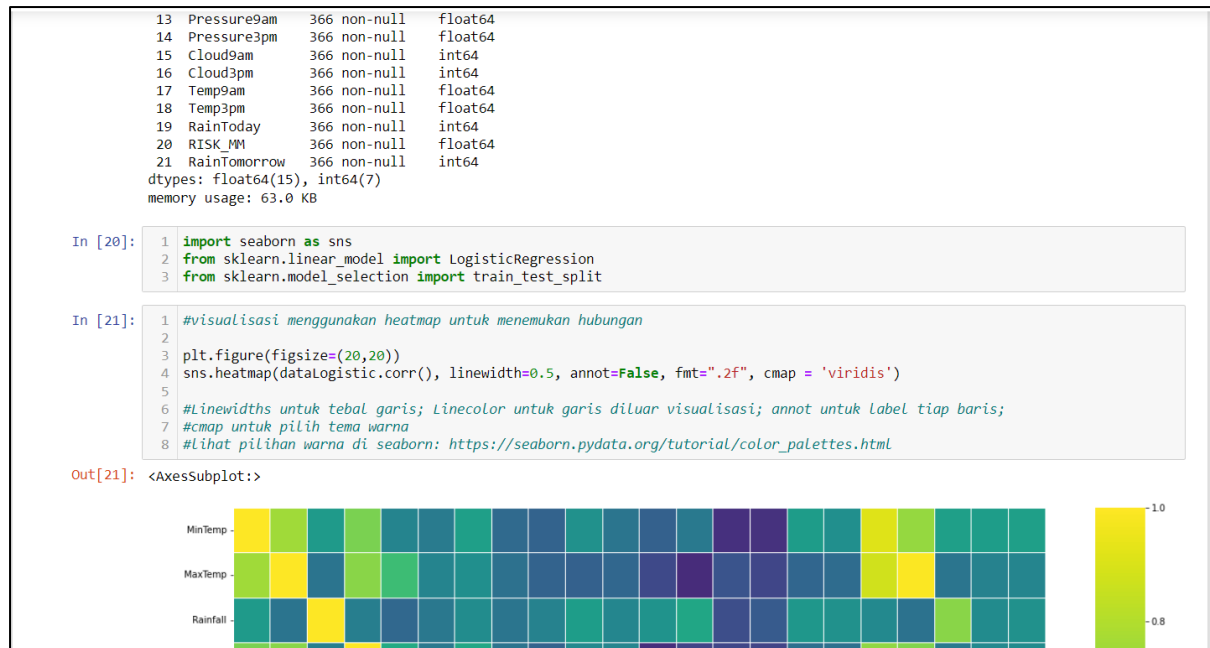
In [18]:

```
1 #encoding
2
3 def encode_data(feature_name):
4     """
5     This function takes feature name as a parameter and returns mapping dictionary to replace(or map) categorical data with
6     """
7     mapping_dict = {}
8
9     unique_values = list(dataLogistic[feature_name].unique())
10
11
12     for idx in range(len(unique_values)):
```

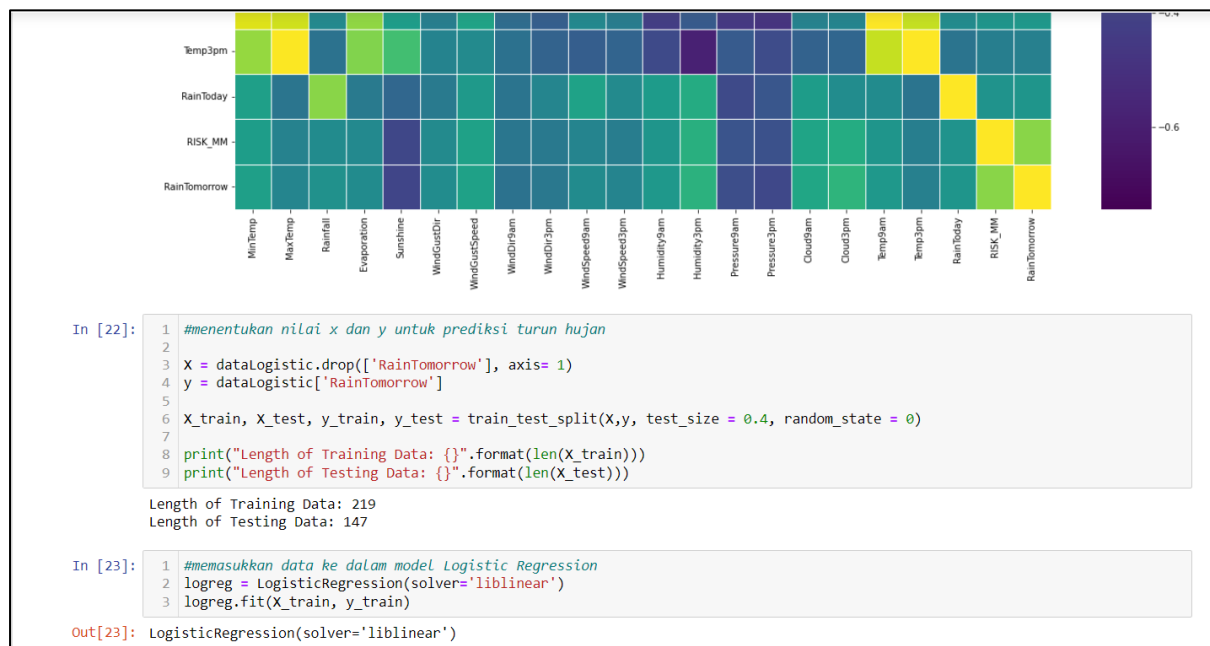
Gambar 9.

	<pre> return mapping_dict dataLogistic['RainToday'].replace({'No':0,'Yes': 1}, inplace = True) dataLogistic['RainTomorrow'].replace({'No':0,'Yes': 1}, inplace = True) dataLogistic['WindGustDir'].replace(encode_data('WindGustDir'), inplace = True) dataLogistic['WindDir9am'].replace(encode_data('WindDir9am'), inplace = True) dataLogistic['WindDir3pm'].replace(encode_data('WindDir3pm'), inplace = True) </pre>
In [19]:	<pre> #checking data dataLogistic dataLogistic.info() </pre> <div> <div><class 'pandas.core.frame.DataFrame'></div> <div>RangeIndex: 366 entries, 0 to 365</div> <div>Data columns (total 22 columns):</div> <div> <div>#</div> <div>Column</div> <div>Non-Null Count</div> <div>Dtype</div> </div> <div> <div>0</div> <div>MinTemp</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>1</div> <div>MaxTemp</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>2</div> <div>Rainfall</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>3</div> <div>Evaporation</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>4</div> <div>Sunshine</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>5</div> <div>WindGustDir</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>6</div> <div>WindGustSpeed</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>7</div> <div>WindDir9am</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>8</div> <div>WindDir3pm</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>9</div> <div>WindSpeed9am</div> <div>366 non-null</div> <div>float64</div> </div> <div> <div>10</div> <div>WindSpeed3pm</div> <div>366 non-null</div> <div>int64</div> </div> <div> <div>11</div> <div>Humidity9am</div> <div>366 non-null</div> <div>int64</div> </div> <div> <div>12</div> <div>Humidity3pm</div> <div>366 non-null</div> <div>int64</div> </div> </div>

Gambar 10.



Gambar 11.



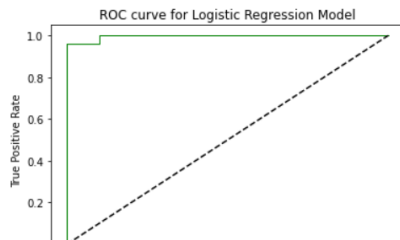
Gambar 12.

```

In [24]: 1 #evaluating data using ROC Curve
2 #hasil yang kemungkinan terjadi adalah
3
4 #True Positive: prediksi hujan, dan turun hujan
5 #True Negative: prediksi tidak ada hujra, dan tidak turun hujan
6
7 print("Accuracy : ", logreg.score(X_test, y_test))
8
9 from sklearn.metrics import roc_curve
10
11 y_pred_logreg_proba = logreg.predict_proba(X_test)[:,-1]
12 fpr, tpr, thresholds = roc_curve(y_test, y_pred_logreg_proba)
13 plt.figure(figsize=(6, 4))
14 plt.plot(fpr, tpr, '-g', linewidth=1)
15 plt.plot([0,1],[0,1], 'k--')
16 plt.title('ROC curve for Logistic Regression Model')
17 plt.xlabel("False Positive Rate")
18 plt.ylabel('True Positive Rate')
19 plt.show()

```

Accuracy : 0.9931972789115646



Gambar 13.

==Exercise Linear Regression dan Logistic Regression Terhadap harga_mobil==

2. Penerapan Linear Regression pada harga_mobil dalam memprediksi harga

```

In [25]: 1 #reading data
2 dataset2 = pd.read_csv(r"D:\SEMESTER 4\IS411 Data Modelling\LAB\Bahan Modul 4\harga_mobil.csv", header=0)
3 dataset2.info()
4 dataset2.head(5)

```

```

23 city-mpg      201 non-null   int64
24 highway-mpg   201 non-null   int64
25 price         201 non-null   int64
dtypes: float64(10), int64(6), object(10)
memory usage: 41.0+ KB

```

Out[25]:

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	...	engine-size	fuel-system	bore	stroke	compression-ratio	horsepower
0	3	NaN	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111
1	3	NaN	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111
2	1	NaN	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	9.0	154
3	2	164.0	audi	gas	std	four	sedan	fwd	front	99.8	...	109	mpfi	3.19	3.40	10.0	102

Gambar 14.

```
In [26]: 1 dataset2.describe()

Out[26]:
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke	compression-ratio	horsepower
count	201.000000	164.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	197.000000	197.000000	201.000000	199.000000
mean	0.840796	122.000000	98.797015	174.200995	65.889055	53.766667	2555.666667	126.875622	3.330711	3.256904	10.164279	103.396985
std	1.254802	35.442168	6.066366	12.322175	2.101471	2.447822	517.296727	41.546834	0.270793	0.319256	4.004965	37.553843
min	-2.000000	65.000000	86.600000	141.100000	60.300000	47.800000	1488.000000	61.000000	2.540000	2.070000	7.000000	48.000000
25%	0.000000	94.000000	94.500000	166.800000	64.100000	52.000000	2169.000000	98.000000	3.150000	3.110000	8.600000	70.000000
50%	1.000000	115.000000	97.000000	173.200000	65.500000	54.100000	2414.000000	120.000000	3.310000	3.290000	9.000000	95.000000
75%	2.000000	150.000000	102.400000	183.500000	66.600000	55.500000	2926.000000	141.000000	3.590000	3.410000	9.400000	116.000000
max	3.000000	256.000000	120.900000	208.100000	72.000000	59.800000	4066.000000	326.000000	3.940000	4.170000	23.000000	262.000000

```
In [27]: 1 dataset2.isnull().sum()

Out[27]:
```

symboling	0
normalized-losses	37
make	0
fuel-type	0
aspiration	0
num-of-doors	2
body-style	0
drive-wheels	0
engine-location	0
wheel-base	0
length	0
width	0
height	0
curb-weight	0
engine-type	0

Gambar 15.

```
num-of-doors      2
body-style        0
drive-wheels      0
engine-location   0
wheel-base       0
length           0
width            0
height           0
curb-weight       0
engine-type       0
num-of-cylinders  0
engine-size       0
fuel-system       0
bore              4
stroke           4
compression-ratio 0
horsepower        2
peak-rpm          2
city-mpg          0
highway-mpg       0
price            0
dtype: int64

In [28]: 1 def clean_dataset(dataset2):
2         assert isinstance(dataset2, pd.DataFrame), "dataset2 needs to be a pd.DataFrame"
3         dataset2.dropna(inplace=True)
4         indices_to_keep = ~dataset2.isin([np.nan, np.inf, -np.inf]).any(axis=1)
5         return dataset2[indices_to_keep].astype(np.float64)

In [29]: 1 #get rid of infinite values.
2         dataset2.replace([np.inf, -np.inf], np.nan, inplace=True)

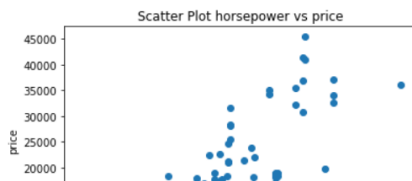
In [30]: 1 #fill missing values
2         dataset2.fillna(0, inplace=True)
```

Gambar 16.

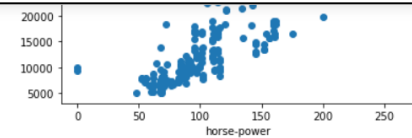

```
In [31]: 1 #kita akan melakukan prediksi price berdasarkan engine-size(cc)
2 #maka price sebagai dependent variable (y) dan horsepower sebagai independent variable (x)
3 #kolom yang akan digunakan adalah horsepower dan price
4
5 dataMobil = dataset2[['horsepower','price']]
6 dataMobil.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 201 entries, 0 to 200
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   horsepower  201 non-null    float64
1   price       201 non-null    int64
dtypes: float64(1), int64(1)
memory usage: 3.3 KB
```

```
In [32]: 1 #melihat korelasi data dengan Scatter Plot
2
3 plt.scatter(dataMobil['horsepower'], dataMobil['price'])
4 plt.xlabel('horse-power')
5 plt.ylabel('price')
6 plt.title('Scatter Plot horsepower vs price')
7 plt.show()
```



Gambar 17.



```
In [33]: 1 x = dataMobil['horsepower'].values.reshape(-1, 1)
2 y = dataMobil['price'].values.reshape(-1, 1)
3
4 #reshape (-1, 1) untuk mengubah susunan data menjadi row ke sampling.
```

```
In [34]: 1 #membagi data training dan data testing untuk harga_mobil
2
3 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
```

```
In [35]: 1 LinReg2 = LinearRegression()
2 LinReg2.fit(x_train, y_train)
3 LinReg2.score(x_test, y_test)
```

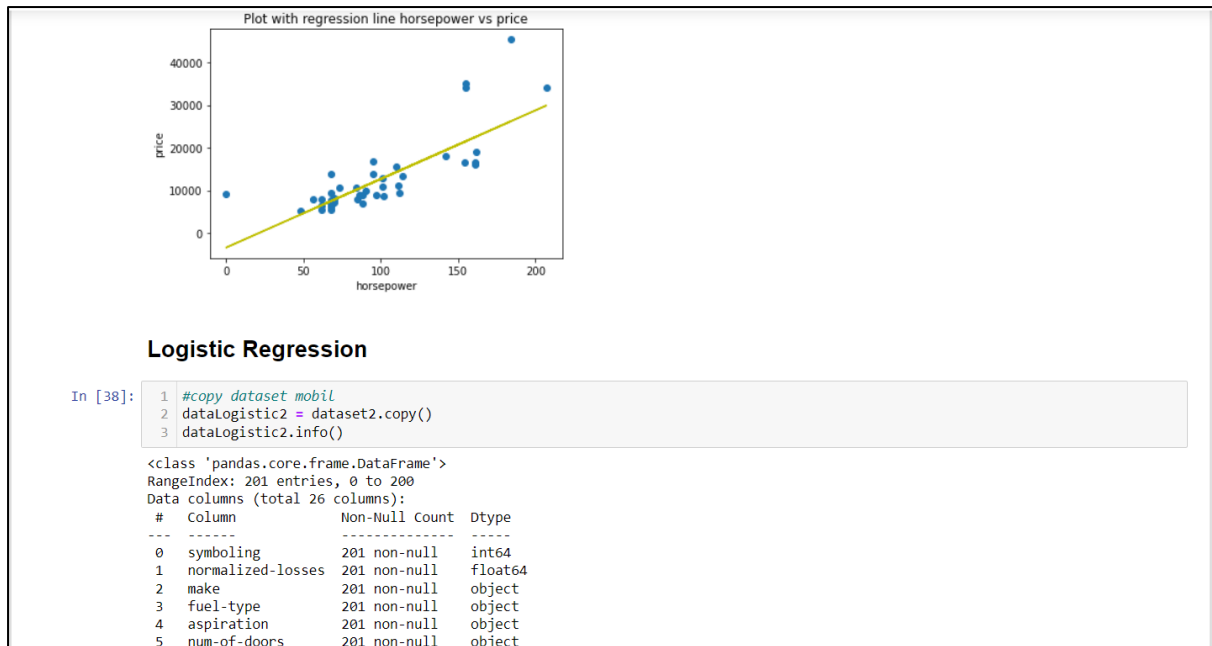
Out[35]: 0.6253679708642874

```
In [36]: 1 #hasil akurasi train data mobil yang saya peroleh = 0.6253679708642874
```

```
In [37]: 1 y_predict = LinReg2.predict(x_test)
2 plt.scatter(x_test, y_test)
3 plt.plot(x_test, y_predict, c='y')
4 plt.xlabel('horsepower')
5 plt.ylabel('price')
6 plt.title('Plot with regression line horsepower vs price')
```

Out[37]: Text(0.5, 1.0, 'Plot with regression line horsepower vs price')

Gambar 18.



Gambar 19.

```
6   body-style           201 non-null    object
7   drive-wheels         201 non-null    object
8   engine-location      201 non-null    object
9   wheel-base           201 non-null    float64
10  length               201 non-null    float64
11  width                201 non-null    float64
12  height               201 non-null    float64
13  curb-weight          201 non-null    int64
14  engine-type           201 non-null    object
15  num-of-cylinders      201 non-null    object
16  engine-size           201 non-null    int64
17  fuel-system           201 non-null    object
18  bore                 201 non-null    float64
19  stroke               201 non-null    float64
20  compression-ratio     201 non-null    float64
21  horsepower            201 non-null    float64
22  peak-rpm              201 non-null    float64
23  city-mpg              201 non-null    int64
24  highway-mpg           201 non-null    int64
25  price                 201 non-null    int64
dtypes: float64(10), int64(6), object(10)
memory usage: 41.0+ KB
```

```
In [39]: 1 #encoding
          2
          3 def encode_data(feature_name):
          4     ...
          5     This function takes feature name as a parameter and returns mapping dictionary to replace(or map) categorical data with
          6     ...
          7     mapping_dict = {}
          8
          9     unique_values = list(dataLogistic2[feature_name].unique())
          10
          11
          12     for idx1 in range(len(unique_values)):
          13
```

Gambar 20.

```

13     mapping_dict[unique_values[idx1]] = idx1
14
15     return mapping_dict
16
17
18 #dataLogistic2['RainToday'].replace({'No':0,'Yes': 1}, inplace = True)
19
20 #dataLogistic2['RainTomorrow'].replace({'No':0,'Yes': 1}, inplace = True)
21
22 dataLogistic2['make'].replace(encode_data('make'), inplace = True)
23
24 dataLogistic2['fuel-type'].replace(encode_data('fuel-type'), inplace = True)
25
26 dataLogistic2['aspiration'].replace(encode_data('aspiration'), inplace = True)
27
28 dataLogistic2['num-of-doors'].replace(encode_data('num-of-doors'), inplace = True)
29
30 dataLogistic2['body-style'].replace(encode_data('body-style'), inplace = True)
31
32 dataLogistic2['drive-wheels'].replace(encode_data('drive-wheels'), inplace = True)
33
34 dataLogistic2['engine-location'].replace(encode_data('engine-location'), inplace = True)
35
36 dataLogistic2['engine-type'].replace(encode_data('engine-type'), inplace = True)
37
38 dataLogistic2['num-of-cylinders'].replace(encode_data('num-of-cylinders'), inplace = True)
39
40 dataLogistic2['fuel-system'].replace(encode_data('fuel-system'), inplace = True)
41

```

In [40]:

```

1 #checking the changes
2 dataLogistic2.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 201 entries, 0 to 200
Data columns (total 26 columns):

```

Gambar 21.

```

# Column Non-Null Count Dtype
---
0 symboling 201 non-null int64
1 normalized-losses 201 non-null float64
2 make 201 non-null int64
3 fuel-type 201 non-null int64
4 aspiration 201 non-null int64
5 num-of-doors 201 non-null int64
6 body-style 201 non-null int64
7 drive-wheels 201 non-null int64
8 engine-location 201 non-null int64
9 wheel-base 201 non-null float64
10 length 201 non-null float64
11 width 201 non-null float64
12 height 201 non-null float64
13 curb-weight 201 non-null int64
14 engine-type 201 non-null int64
15 num-of-cylinders 201 non-null int64
16 engine-size 201 non-null int64
17 fuel-system 201 non-null int64
18 bore 201 non-null float64
19 stroke 201 non-null float64
20 compression-ratio 201 non-null float64
21 horsepower 201 non-null float64
22 peak-rpm 201 non-null float64
23 city-mpg 201 non-null int64
24 highway-mpg 201 non-null int64
25 price 201 non-null int64
dtypes: float64(10), int64(16)
memory usage: 41.0 KB

```

In [41]:

```

1 def clean_dataset(dataLogistic2):
2     assert isinstance(dataLogistic2, pd.DataFrame), "dataset2 needs to be a pd.DataFrame"
3     dataLogistic2.dropna(inplace=True)
4     indices_to_keep = ~dataLogistic2.isin([np.nan, np.inf, -np.inf]).any(axis=1)
5     return dataLogistic2[indices_to_keep].astype(np.float64)

```

Gambar 22.

```

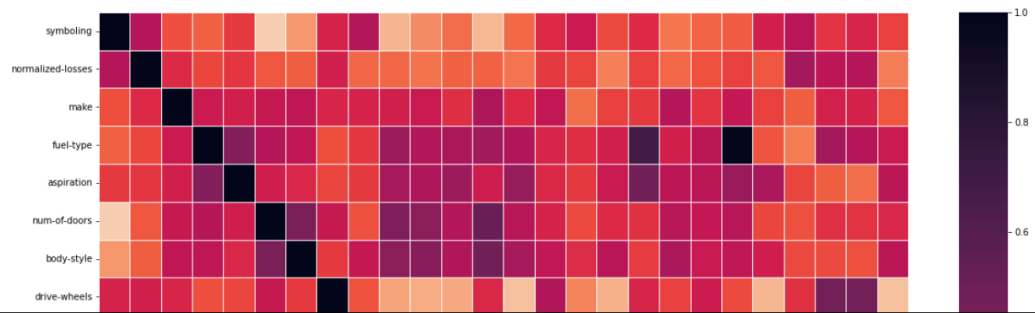
In [42]: 1 #get rid of infinite values.
          2 dataLogistic2.replace([np.inf, -np.inf], np.nan, inplace=True)

In [43]: 1 #fill missing values
          2 dataLogistic2.fillna(0, inplace=True)

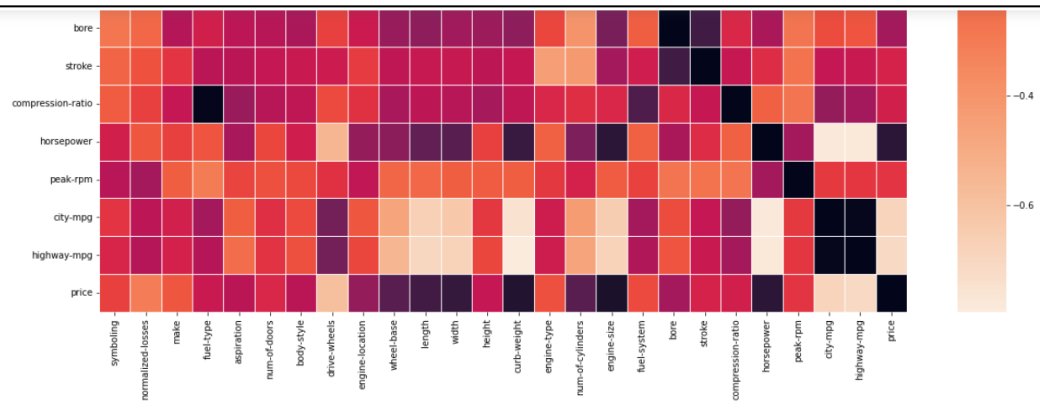
In [44]: 1 #visualisasi menggunakan heatmap untuk menemukan hubungan
          2
          3 plt.figure(figsize=(20,20))
          4 sns.heatmap(dataLogistic2.corr(), linewidth=0.5, annot=False, fmt=".2f", cmap = 'rocket_r')
          5
          6 #Linewidths untuk tebal garis; Linecolor untuk garis diluar visualisasi; annot untuk label tiap baris;
          7 #cmap untuk pilih tema warna
          8 #lihat pilihan warna di seaborn: https://seaborn.pydata.org/tutorial/color_palettes.html

```

Out[44]: <AxesSubplot>



Gambar 23.



```

In [45]: 1 #menentukan nilai x dan y untuk prediksi harga mobil
          2
          3 X = dataLogistic2.drop(['price'], axis= 1)
          4 y = dataLogistic2['price']
          5
          6 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size = 0.2, random_state = 0)
          7
          8 print("Length of Training Data: {}".format(len(X_train)))
          9 print("Length of Testing Data: {}".format(len(X_test)))

```

Length of Training Data: 160
Length of Testing Data: 41

Gambar 24.

```

In [46]: 1 #memasukkan data ke dalam model Logistic Regression
2 logreg2 = LogisticRegression(solver='liblinear')
3 logreg2.fit(X_train, y_train)

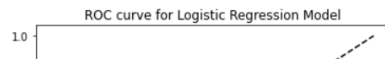
Out[46]: LogisticRegression(solver='liblinear')

In [47]: 1 #evaluating data using ROC Curve
2 #hasil yang kemungkinan terjadi adalah
3
4 #True Positive: prediksi harga naik, dan harga turun
5 #True Negative: prediksi harga tidak naik, dan harga tidak turun
6
7 print("Accuracy : ", logreg2.score(X_test, y_test))
8
9 from sklearn.metrics import roc_curve
10
11 y_pred_logreg2_proba = logreg2.predict_proba(X_test)[:,-1]
12 fpr1, tpr1, thresholds = roc_curve(y_test, y_pred_logreg2_proba, pos_label=1)
13 plt.figure(figsize=(6, 4))
14 plt.plot(fpr1, tpr1, 'g', linewidth=1)
15 plt.plot([0,1],[0,1], 'k--')
16 plt.title('ROC curve for Logistic Regression Model')
17 plt.xlabel("False Positive Rate")
18 plt.ylabel("True Positive Rate")
19 plt.show()

```

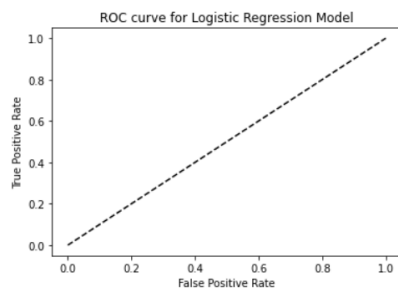
Accuracy : 0.024390243902439025

C:\Users\Darren\anaconda3\lib\site-packages\sklearn\metrics_ranking.py:999: UndefinedMetricWarning: No positive samples in y_true, true positive value should be meaningless
warnings.warn(



Gambar 25.

C:\Users\Darren\anaconda3\lib\site-packages\sklearn\metrics_ranking.py:999: UndefinedMetricWarning: No positive samples in y_true, true positive value should be meaningless
warnings.warn(



In []: 1

Gambar 26.