

# Hands-On

Hands-On ini digunakan pada kegiatan Microcredential Associate Data Scientist 2021

## Tugas Mandiri Pertemuan 10

Pertemuan 10 (sepuluh) pada Microcredential Associate Data Scientist 2021 menyampaikan materi mengenai Membangun Model (Dasar Regresi dan Regresi Linier). silakan Anda kerjakan Latihan 1 s/d 20. Output yang anda lihat merupakan panduan yang dapat Anda ikuti dalam penulisan code :)

## Latihan (1)

### Melakukan import library yang dibutuhkan

```
In [34]: # Import library pandas
import pandas as pd

# Import library numpy
import numpy as np

# Import library matplotlib and seaborn untuk visualisasi
import matplotlib.pyplot as plt
import seaborn as sns

# Import Module LinearRegression digunakan untuk memanggil algoritma Linear Regression.
from sklearn.linear_model import LinearRegression

# Import Module train_test_split digunakan untuk membagi data kita menjadi training dan testing set.
from sklearn.model_selection import train_test_split

# Import modul mean_absolute_error dari library sklearn
from sklearn.metrics import mean_absolute_error

# Import math agar program dapat menggunakan semua fungsi yang ada pada modul math. (math.sqrt)
import math

# Menon aktifkan peringatan pada python
import warnings
warnings.filterwarnings('ignore')
```

### Load Dataset

```
In [35]: # Panggil file (load file bernama CarPrice_Assignment.csv) dan simpan dalam DataFrame Lalu tampilkan 10 baris as
data = pd.read_csv('CarPrice_Assignment.csv')
dataset = pd.DataFrame(data)
dataset.head(10)
```

	car_ID	symboling	CarName	fueltpe	aspiration	doornumber	carbody	drivewheel	engineLocation	wheelbase	...	engineize	fuely
0	1	3	alfa-romero guilia	gas	std	two	convertible	rwd	front	88.6	...	130	
1	2	3	alfa-romero Quadrifoglio	gas	std	two	convertible	rwd	front	88.6	...	130	
2	3	4	alfa-romero Quadrifoglio	gas	std	two	hatchback	rwd	front	99.5	...	102	
3	4	2	audi 100 l	gas	std	four	sedan	fwd	front	94.8	...	159	
4	5	2	audi 100ls	gas	std	four	sedan	4wd	front	99.4	...	136	
5	6	2	audi fox	gas	std	two	sedan	fwd	front	99.8	...	136	
6	7	1	audi 100ls	gas	std	four	sedan	fwd	front	105.8	...	136	
7	8	1	audi 5000	gas	std	four	wagon	fwd	front	105.8	...	136	
8	9	1	audi 4000	gas	turbo	four	sedan	fwd	front	105.8	...	131	
9	10	0	audi 5000s (diesel)	gas	turbo	two	hatchback	4wd	front	99.5	...	131	

10 rows × 26 columns

## Latihan (2)

### Review Dataset

```
In [36]: # melihat jumlah baris dan jumlah kolom (bentuk data) pada data df dengan Fungsi .shape
dataset.shape
```

Out [36]: (205, 26)

Data kita mempunyai 26 kolom dengan 205 baris.

```
In [37]: # Melihat Informasi lebih detail mengenai struktur DataFrame dapat dilihat menggunakan fungsi info()
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 26 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   car_ID              205 non-null    int64
 1   symboling           205 non-null    int64
 2   CarName             205 non-null    object
 3   fueltpe             205 non-null    object
 4   aspiration          205 non-null    object
 5   doornumber         205 non-null    object
 6   carbody            205 non-null    object
 7   drivewheel         205 non-null    object
 8   engineLocation     205 non-null    object
 9   wheelbase          205 non-null    float64
10   carlength          205 non-null    float64
11   carwidth           205 non-null    float64
12   carheight          205 non-null    float64
13   curbweight         205 non-null    object
14   enginetype         205 non-null    object
15   cylindernumber     205 non-null    object
16   engineize          205 non-null    int64
17   fuelsystem         205 non-null    object
18   boreratio          205 non-null    float64
19   stroke            205 non-null    float64
20   compressionratio   205 non-null    float64
21   horsepower         205 non-null    int64
22   peakrpm           205 non-null    int64
23   citympg           205 non-null    int64
24   highwaympg        205 non-null    int64
25   price             205 non-null    float64
dtypes: float64(8), int64(8), object(10)
memory usage: 41.4+ KB
```

```
In [72]: # melihat statistik data untuk data numeric seperti count, mean, standard deviation, maximum, minimum, dan qu
dataset.describe()
```

	car_ID	symboling	wheelbase	carlength	carwidth	carheight	curbweight	engineize	boreratio	stroke	compression
count	205.000000	205.000000	205.000000	205.000000	205.000000	205.000000	205.000000	205.000000	205.000000	205.000000	205.000000
mean	103.000000	0.834146	98.756585	174.049268	65.907805	53.724878	2555.565854	126.907317	3.329756	3.254515	0.101428
std	59.322565	1.245307	6.021776	12.337289	2.145204	2.443522	520.680204	41.642693	0.270844	0.313597	0.159215
min	1.000000	-2.000000	86.600000	141.100000	60.300000	47.800000	1488.000000	61.000000	2.540000	2.070000	0.000000
25%	52.000000	0.000000	94.500000	166.300000	64.100000	52.000000	2145.000000	97.000000	3.150000	3.110000	8.600000
50%	103.000000	1.000000	97.000000	173.200000	65.500000	54.100000	2414.000000	120.000000	3.310000	3.290000	9.000000
75%	154.000000	2.000000	102.400000	183.100000	66.900000	55.500000	2995.000000	141.000000	3.580000	3.410000	9.400000
max	205.000000	3.000000	120.900000	208.100000	72.300000	59.800000	4066.000000	326.000000	3.940000	4.170000	23.000000

```
In [39]: # cek nilai yang hilang / missing values di dalam data
dataset.isna().sum()
```

```
Out [39]: car_ID              0
          symboling         0
          CarName           0
          fueltpe           0
          aspiration        0
          doornumber        0
          carbody           0
          drivewheel        0
          engineLocation    0
          wheelbase         0
          carlength         0
          carwidth          0
          carheight         0
          curbweight        0
          enginetype        0
          cylindernumber    0
          engineize         0
          fuelsystem        0
          boreratio         0
          stroke            0
          horsepower        0
          peakrpm           0
          citympg           0
          highwaympg        0
          price             0
          dtype: int64
```

Ternyata data kita tidak ada missing values.

Simple linear regression atau regresi linier sederhana merupakan jenis regresi yang paling sederhana karena hanya melibatkan satu variabel bebas atau variabel independen X.

## Visualisasi data untuk pemilihan fitur / variabel independen X

1. Variabel y atau variabel dependet adalah 'price'
2. Lakukan Visualisasi dalam penerapannya agar dapat terlihat jelas / mempermudah dalam membaca data tsb
3. Untuk dapat menentukan variabel X yaitu dapat melihat korelasi antar variabel dengan variabel y / kolom 'price'

## Latihan (3)

untuk dapat menentukan lebih detail / akurat dalam pemilihan fitur dapat dilihat dari hubungan korelasi nya dengan function corr()

```
In [40]: dataset.corr()
```

	car_ID	symboling	wheelbase	carlength	carwidth	carheight	curbweight	engineize	boreratio	stroke	compression
car_ID	1.000000	-0.151621	0.129729	0.170636	0.052387	0.255960	0.071962	-0.033930	0.260064	-0.160824	0.151500
symboling	-0.151621	1.000000	-0.531954	-0.357612	-0.232919	-0.541038	-0.227691	-0.105790	-0.130051	-0.008735	-0.171500
wheelbase	0.129729	-0.531954	1.000000	0.874587	0.795144	0.589435	0.776386	0.569329	0.488750	0.160959	0.240100
carlength	0.170636	-0.357612	0.874587	1.000000	0.841118	0.491029	0.877728	0.683360	0.606454	0.129533	0.151500
carwidth	0.052387	-0.232919	0.795144	0.841118	1.000000	0.279210	0.867032	0.735433	0.559150	0.162942	0.181500
carheight	0.255960	-0.541038	0.589435	0.491029	0.279210	1.000000	0.295572	0.067439	0.171071	-0.055307	0.261500
curbweight	0.071962	-0.227691	0.776386	0.877728	0.867032	0.295572	1.000000	0.850594	0.648480	0.168790	0.151500
engineize	-0.033930	-0.105790	0.569329	0.683360	0.735433	0.067439	0.850594	1.000000	0.583774	0.203129	0.002100
boreratio	0.260064	-0.130051	0.488750	0.606454	0.559150	0.171071	0.648480	0.583774	1.000000	-0.055909	0.000100
stroke	0.160824	-0.008735	0.160959	0.129533	0.182942	-0.055307	0.168790	0.203129	-0.055909	1.000000	0.181500
compression	0.150726	-0.178515	0.249786	0.158414	0.181129	0.261214	0.151362	0.028971	0.005197	0.186110	1.000000
horsepower	-0.015006	0.0070873	0.353294	0.552623	0.640732	-0.108802	0.750739	0.809769	0.573677	0.080940	-0.200100
peakrpm	-0.020389	-0.057969	0.063469	0.087728	-0.287242	-0.220012	-0.320411	-0.266243	-0.244660	-0.067964	-0.431500
citympg	0.015940	-0.035823	-0.470414	-0.070909	-0.642704	-0.048640	-0.757414	-0.653658	-0.584532	-0.042435	-0.332100
highwaympg	0.011255	0.034606	-0.540882	-0.704662	-0.677218	-0.107358	-0.797465	-0.677470	-0.587012	-0.439391	0.261500
price	-0.109093	-0.079978	0.577816	0.682920	0.759325	0.119336	0.835305	0.874145	0.553173	0.079443	0.061500

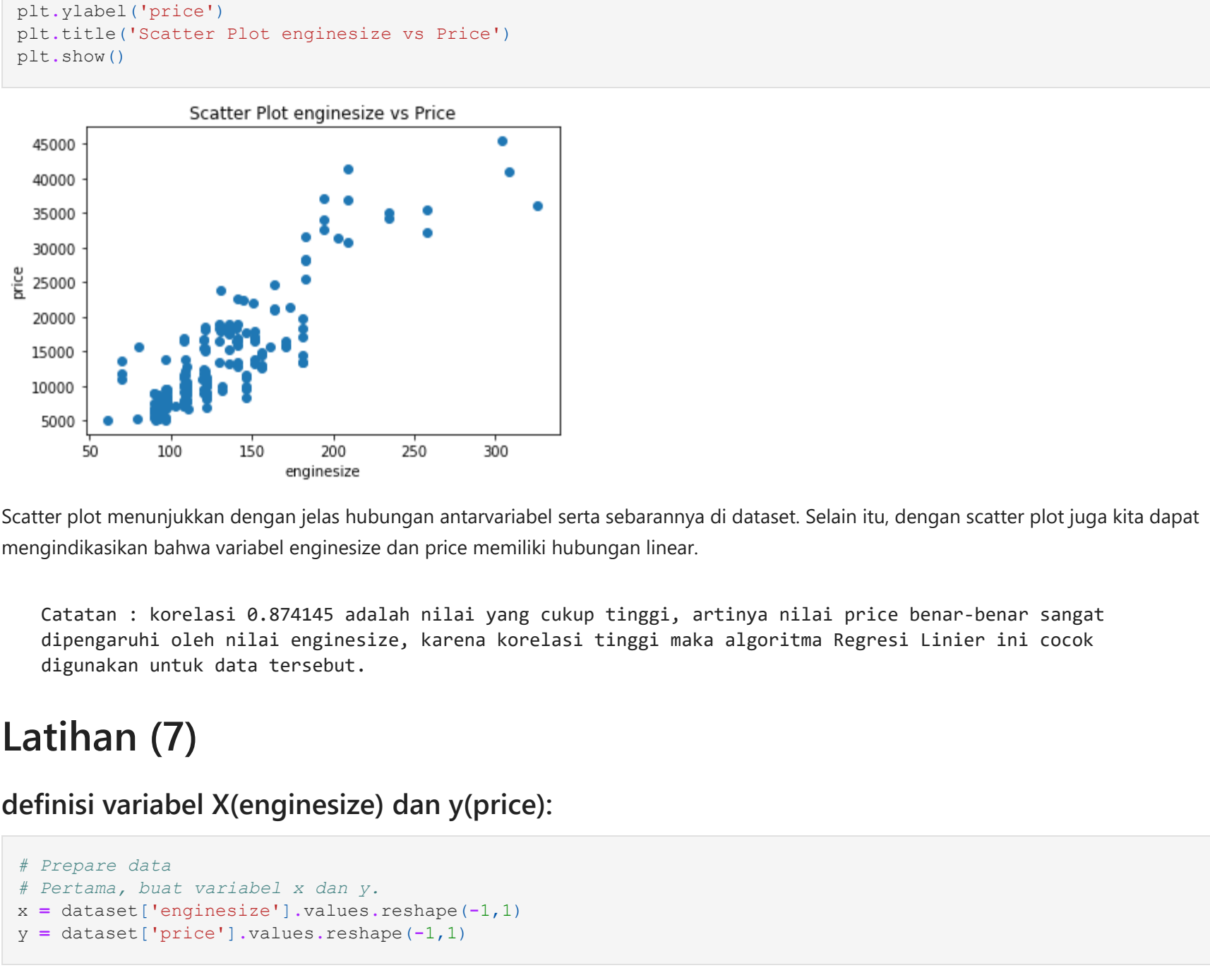
ternyata engineize, boreratio, horsepower memiliki korelasi yang signifikan dengan harga/price.

## Latihan (4)

Buat Visualisasi scatter plot dari kolom:

'engineize', 'boreratio', 'stroke', 'compressionratio', 'horsepower', 'peakrpm', 'wheelbase', 'citympg', 'highwaympg'

```
In [71]: df = pd.DataFrame(dataset)
sns.pairplot(dataset, x_vars=[x, y, z], y_vars='price', size=4, aspect=1, diag_kind=None)
plt.show()
```

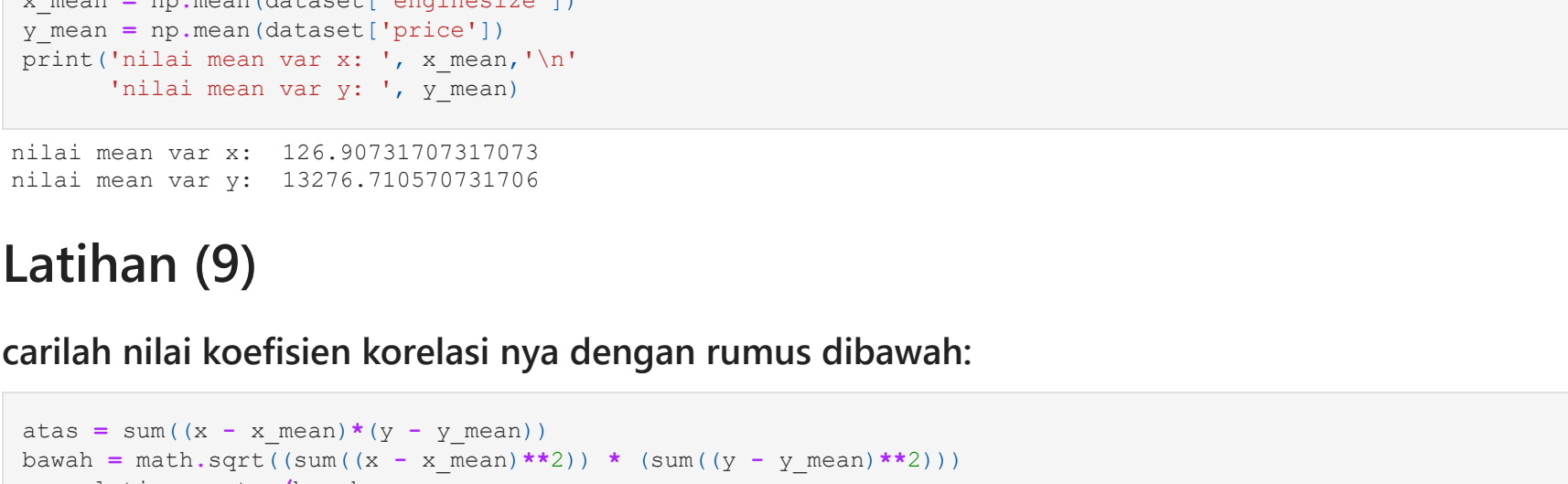


## Latihan (5)

Buat Visualisasi Heatmap dari kolom:

'engineize', 'boreratio', 'stroke', 'compressionratio', 'horsepower', 'peakrpm', 'wheelbase', 'citympg', 'highwaympg'

```
In [43]: plt.figure(figsize=(8,8))
data_fitur = dataset[['engineize', 'boreratio', 'stroke', 'compressionratio', 'horsepower', 'peakrpm', 'wheelbase', 'citympg', 'highwaympg', 'price']]
sns.heatmap(data_fitur.corr(), annot=False, fmt="")
plt.show()
```



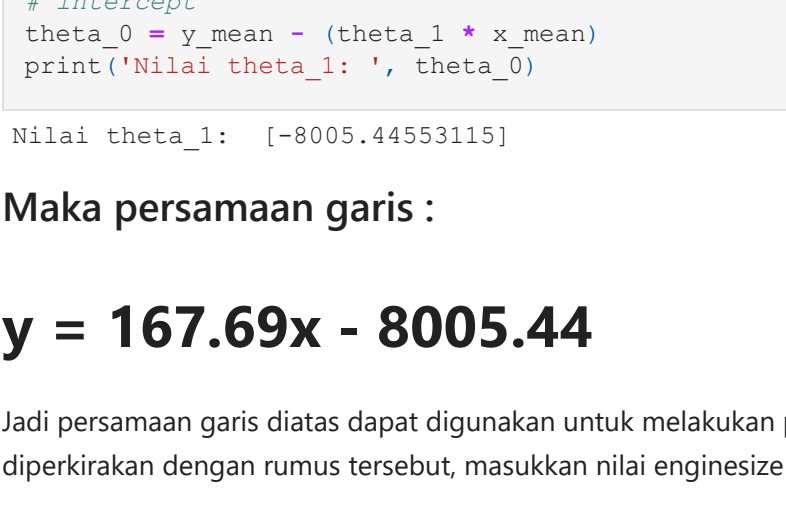
Dari hasil visualisasi diatas bahwa fitur/kolom engineize memiliki korelasi yang tinggi terhadap kolom price / variabel dependet sehingga kita mengambil fitur/kolom engineize untuk di training

- Independent variabel adalah engineize.
- Dependent variabel adalah price.

## Latihan (6)

Buat Visualisasi Scatter Plot antara calon variabel X(engineize) dan y(price):

```
In [44]: plt.scatter(dataset['engineize'], dataset['price'])
plt.xlabel('engineize')
plt.ylabel('price')
plt.title('Scatter Plot engineize vs Price')
plt.show()
```



Scatter plot menunjukkan dengan jelas hubungan antarvariabel serta sebarannya di dataset. Selain itu, dengan scatter plot juga kita dapat mengindikasikan bahwa variabel engineize dan price memiliki hubungan linier.

Catatan : korelasi 0.874145 adalah nilai yang cukup tinggi, artinya nilai price benar-benar sangat dipengaruhi oleh nilai engineize, karena korelasi tinggi maka algoritma Regresi Linier ini cocok digunakan untuk data tersebut.

## Latihan (7)

definisi variabel X(engineize) dan y(price):

```
In [45]: # Prepare data
# Persiapan buat variabel x dan y.
x = dataset['engineize'].values.reshape(-1,1)
y = dataset['price'].values.reshape(-1,1)
```

Formula Regresi Linear

"Jika kita melihat formula regresi linear di atas, kita pasti ingat rumus persamaan garis yang pernah dipelajari di bangku sekolah, yaitu  $y = mx + c$ , dimana  $m$  merupakan gradien atau kemiringan garis dan  $c$  merupakan konstanta."

- $y$  from  $scat$
- $y = ax + b$  atau  $y = wx + w0$  atau  $y = mx + c$
- $x = input$
- $y = output$
- $b$  atau  $w0 = intercept / bias$
- $a$  atau  $w1 = slope / gradient / coefficient$

## Latihan (8)

definisi variabel nilai mean/rata-rata X(engineize) dan nilai mean/rata-rata y(price):

```
In [46]: x_mean = np.mean(dataset['engineize'])
y_mean = np.mean(dataset['price'])
print('nilai mean var x: ', x_mean, '\n'
      'nilai mean var y: ', y_mean)
```

nilai mean var x: 126.90731707317073  
nilai mean var y: 13276.710570731706

## Latihan (9)

carilah nilai koefisien korelasi nya dengan rumus dibawah:

```
In [47]: atas = sum((x - x_mean)*(y - y_mean))
bawah = math.sqrt((sum((x - x_mean)**2)) * (sum((y - y_mean)**2)))
correlation = atas/bawah
print('Nilai Correlation Coefficient: ', correlation)
```

Nilai Correlation Coefficient: [0.874144]

carilah nilai parameter theta 1 dan theta 0 dengan rumus dibawah:

```
theta_1 = ((111-104.11) * (13495-13276.71)) + ... + ((114-104.11) * (22625-13276.71)) / ((111-104.11)**2
+ ... + (114-104.11)**2)
```

## Latihan (10)

carilah nilai theta\_1 atau nilai slope

```
In [48]: # slope
# Slope adalah tingkat kemiringan garis, intercept
# adalah jarak titik y pada garis dari titik 0
variance = sum((x - x_mean)**2)
covariance = sum((x - x_mean)*(y - y_mean))
theta_1 = covariance/variance
print('Nilai theta_1: ', theta_1)
```

Nilai theta\_1: [167.69841639]

## Latihan (11)

carilah nilai theta\_0 atau nilai intercept

```
In [49]: # intercept
theta_0 = y_mean - (theta_1 * x_mean)
print('Nilai theta_0: ', theta_0)
```

Nilai theta\_0: [-8005.4453115]

Maka persamaan garis:

**y = 167.69x - 8005.44**

Jadi persamaan garis diatas dapat digunakan untuk melakukan prediksi apabila kita memiliki data engineize yang baru, price dapat diperkirakan dengan rumus tersebut, masukkan nilai engineize baru ke x, maka perkiraan nilai y (price) akan didapat.

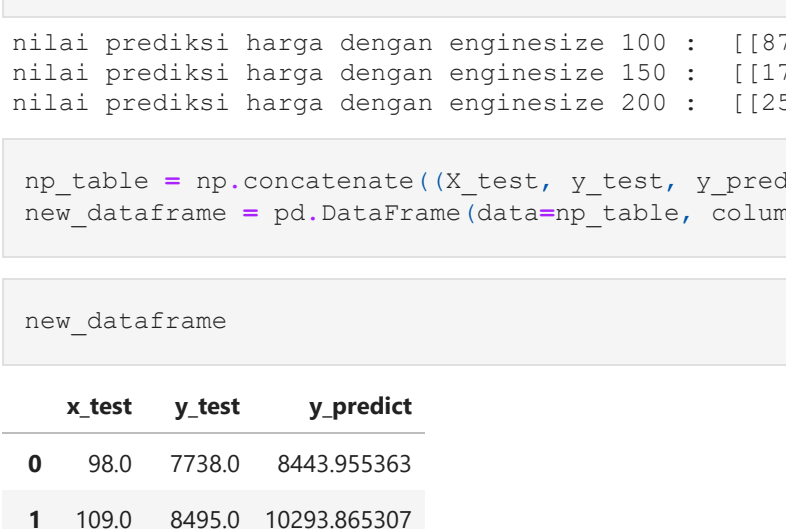
## Latihan (12)

carilah nilai prediksi secara manual dan buatlah visualisasi scatter plot nya

```
In [50]: # prediction manual
y_pred = theta_0 + (theta_1 * x)
print(y_pred)
```

[13795.34859997]

```
In [73]: # visualisasi prediksi dengan scatter plot
y_pred = theta_0 + (theta_1 * x)
plt.scatter(dataset['engineize'], dataset['price'])
plt.plot(x, y_pred, 'r')
plt.xlabel('engineize')
plt.ylabel('Price')
plt.title('Plot engineize vs Price')
plt.show()
```



Linier Regression digunakan untuk Prediksi dengan mencari pola garis terbaik antara variabel independen dan dependen

Pros:

- Mudah diimplementasikan
- Digunakan untuk memprediksi nilai numerik/continous /data jenis interval dan ratio

Cons:

- Cenderung mudah Overfitting
- Tidak dapat digunakan bila relasi antara variabel independen dan dependen tidak linier atau korelasi variabel rendah

## Linier Regression dengan menggunakan library sklearn

1. Pertama yang kita lakukan adalah split data. Train/test split adalah salah satu metode yang dapat digunakan untuk mengevaluasi performa model machine learning. Metode evaluasi model ini membagi dataset menjadi dua bagian yakni bagian yang digunakan untuk training data dan untuk testing data dengan proporsi tertentu. Train data digunakan untuk fit model machine learning, sedangkan test data digunakan untuk mengevaluasi hasil fit model tersebut.

Python memiliki library yang dapat mengimplementasikan train/test split dengan mudah yaitu Skit-Learn. Untuk menggunakannya, kita perlu mengimport Skit-Learn terlebih dahulu, kemudian setelah itu kita dapat menggunakan fungsi train\_test\_split().

## Latihan (13)

split data train dan test dengan function train\_test\_split() dengan train\_size=0.8, test\_size=0.2 dan random\_state=100

```
In [52]: X_train, X_test, y_train, y_test = train_test_split(x, y, train_size=0.8, test_size=0.2, random_state=100)
```

- X\_train: Untuk menampung data source yang akan dilatih.
- X\_test: Untuk menampung data target yang akan dilatih.
- y\_train: Untuk menampung data source yang akan digunakan untuk testing.
- y\_test: Untuk menampung data target yang akan digunakan untuk testing.

X dan y adalah nama variabel yang digunakan saat mendefinisikan data source dan data target. Parameter test\_size digunakan untuk mendefinisikan ukuran data testing. Dalam contoh di atas, test\_size=0.2 berarti data yang digunakan sebagai data testing adalah sebesar 20% dari keseluruhan dataset.

Perlu diketahui bahwa metode ini akan membagi train set dan test set secara random atau acak. Jadi, jika kita mengulang proses running, maka tentunya hasil yang didapat akan berubah-ubah. Untuk mengatasinya, kita dapat menggunakan parameter random\_state

## Latihan (14)

buat object variabel linier regression

```
In [53]: regressor = LinearRegression()
```

## Latihan (15)

training the model menggunakan training data yang sudah displit sebelumnya.

```
In [54]: regressor.fit(X_train, y_train)
```

Out [54]: LinearRegression()

## Latihan (16)

carilah nilai slope/koeffisien (m) dan intercept (b), dengan menggunakan function dari library sklearn -> LinierRegression

```
In [55]: print(regressor.coef_)
print(regressor.intercept_)
```

[[168.17363122]]  
[-8037.06049611]

Dari nilai m dan b diatas, kalau dimasukkan ke dalam rumus persamaan menjadi:

**y = 168.17x - 8037.06**

## &lt;