

Nama : Miladyna Fauzia

NPM : 41155050210023

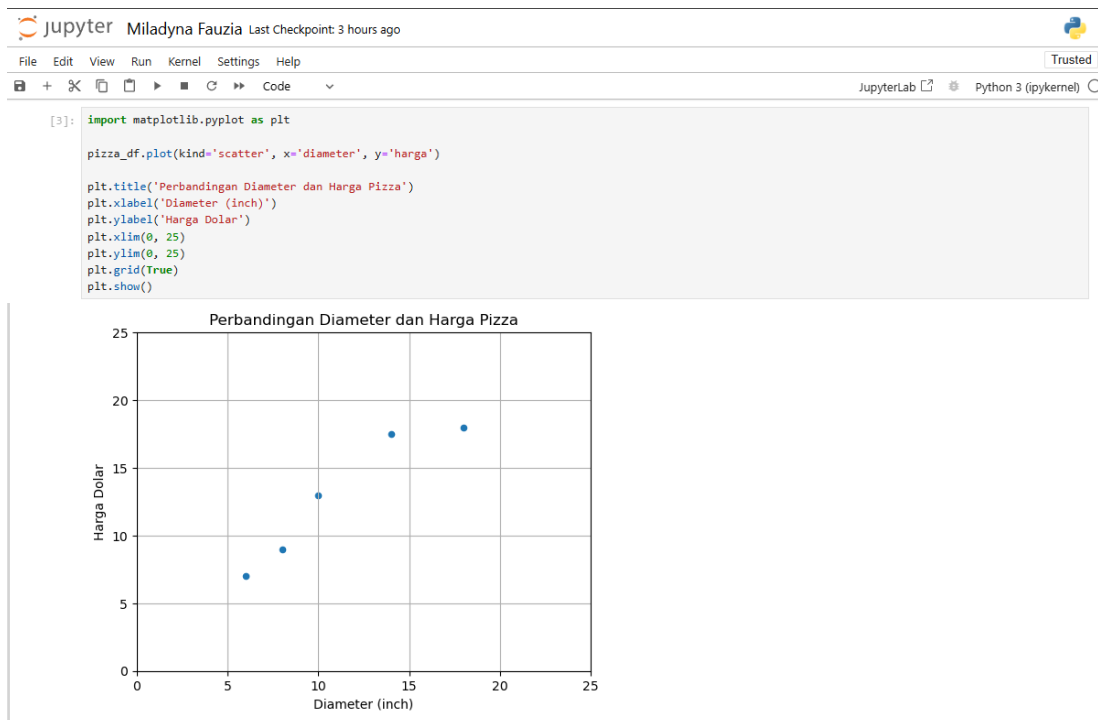
Kelas : INF A1

Pertemuan 2

1. Lakukan praktek dari <https://youtu.be/lcj7-2zMSA?si=f4jWJR6lY8y0BZKl> dan buat screen shot hasil run dengan nama anda pada hasil run tersebut. Praktek tersebut yaitu
 - 1.1. Sample dataset



1.2. Visualisasi dataset



1.3. Transformasi dataset

```
Jupyter Miladyna Fauzia Last Checkpoint: 3 hours ago
File Edit View Run Kernel Settings Help Trusted
JupyterLab Python 3 (ipykernel)

[5]: import numpy as np

X = np.array(pizza_df['diameter'])
y = np.array(pizza_df['harga'])

print(f'X: {X}')
print(f'y: {y}')

X: [ 6  8 10 14 18]
y: [ 7.  9. 13. 17.5 18. ]

[9]: X = X.reshape(-1, 1)
X.shape

[9]: (5, 1)

[11]: X

[11]: array([[ 6],
           [ 8],
           [10],
           [14],
           [18]], dtype=int64)
```

1.4. Training Simple Linear Regression Model

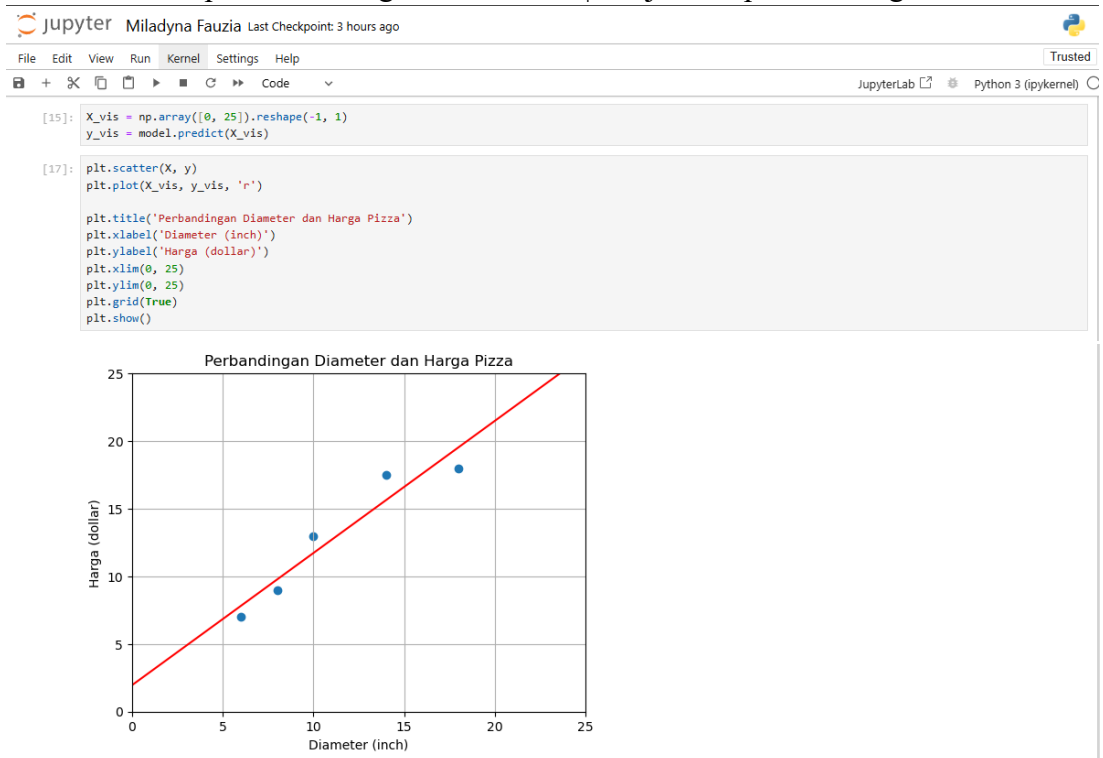
```
Jupyter Miladyna Fauzia Last Checkpoint: 3 hours ago
File Edit View Run Kernel Settings Help Trusted
JupyterLab Python 3 (ipykernel)

[13]: from sklearn.linear_model import LinearRegression

model = LinearRegression()
model.fit(X, y)

[13]: LinearRegression
LinearRegression()
```

1.5. Visualisasi Simple Linear Regression Model | Penjelasan persamaan garis linear



Formula Linear Regression $y = \alpha + \beta x$

- y : response variable
- x : explanatory variable
- α : intercept
- β : slope

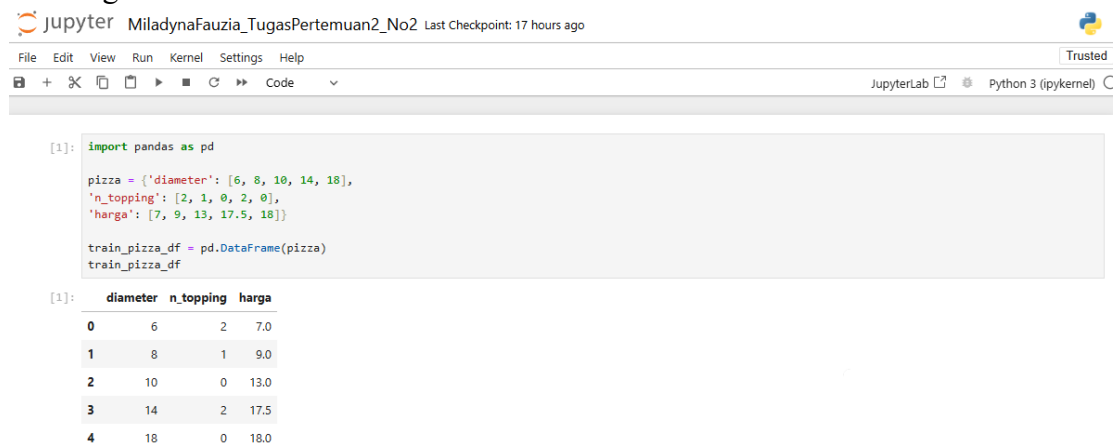
```
intercept: 1.965517241379315
slope:[0.9762931]
```



2. Lakukan praktek dari <https://youtu.be/nWJUJenAyB8?si=BQDzWwrMnr8jtzpV> dan buat screen shot hasil run dengan nama anda pada hasil run tersebut. Praktek tersebut yaitu:

2.1. Sample dataset

Training Dataset



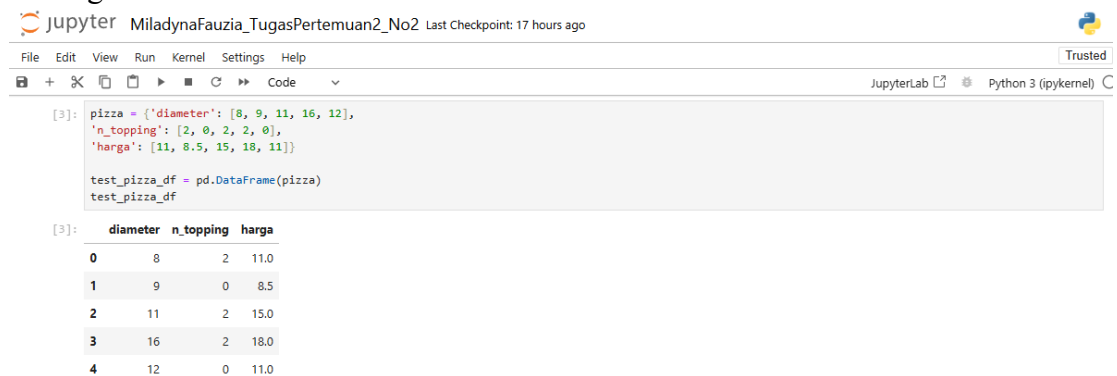
```
[1]: import pandas as pd

pizza = {'diameter': [6, 8, 10, 14, 18],
        'n_topping': [2, 1, 0, 2, 0],
        'harga': [7, 9, 13, 17.5, 18]}

train_pizza_df = pd.DataFrame(pizza)
train_pizza_df
```

	diameter	n_topping	harga
0	6	2	7.0
1	8	1	9.0
2	10	0	13.0
3	14	2	17.5
4	18	0	18.0

Testing Dataset



```
[3]: pizza = {'diameter': [8, 9, 11, 16, 12],
        'n_topping': [2, 0, 2, 2, 0],
        'harga': [11, 8.5, 15, 18, 11]}

test_pizza_df = pd.DataFrame(pizza)
test_pizza_df
```

	diameter	n_topping	harga
0	8	2	11.0
1	9	0	8.5
2	11	2	15.0
3	16	2	18.0
4	12	0	11.0

2.2. Preprocessing dataset



```
[5]: import numpy as np

X_train = np.array(train_pizza_df[['diameter', 'n_topping']])
y_train = np.array(train_pizza_df['harga'])

print(f'X_train:\n{X_train}\n')
print(f'y_train: {y_train}')

X_train:
[[ 6  2]
 [ 8  1]
 [10  0]
 [14  2]
 [18  0]]

y_train: [ 7.  9. 13. 17.5 18. ]

[13]: X_test = np.array(test_pizza_df[['diameter', 'n_topping']])
y_test = np.array(test_pizza_df['harga'])

print(f'X_test:\n{X_test}\n')
print(f'y_test: {y_test}')

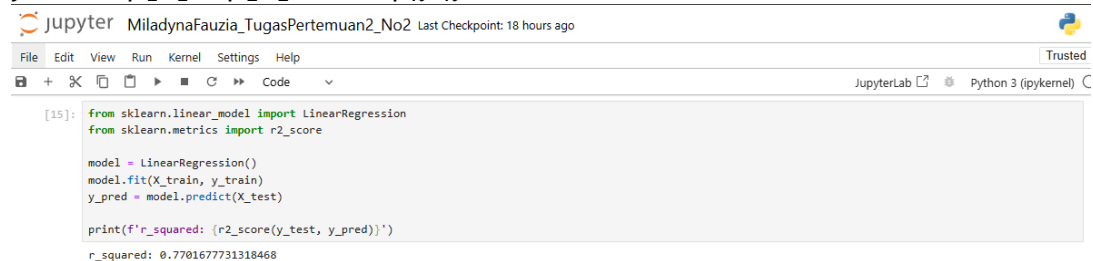
X_test:
[[ 8  2]
 [ 9  0]
 [11  2]
 [16  2]
 [12  0]]

y_test: [11.  8.5 15. 18. 11. ]
```

2.3. Pengenalan Multiple Linear Regression | Apa itu Multiple Linear Regression?

Multiple Linear Regression merupakan generalisasi dari Simple Linear Regression yang memungkinkan untuk menggunakan beberapa explanatory variables.

$$y = a + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$



```
[15]: from sklearn.linear_model import LinearRegression
      from sklearn.metrics import r2_score

      model = LinearRegression()
      model.fit(X_train, y_train)
      y_pred = model.predict(X_test)

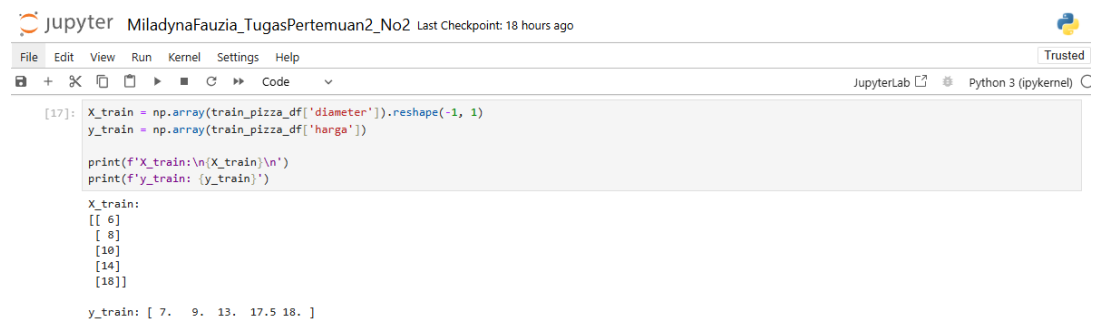
      print(f'r_squared: {r2_score(y_test, y_pred)}')
```

r_squared: 0.7701677731318468

2.4. Pengenalan Polynomial Regression | Apa itu Polynomial Regression?

Polynomial Regression memodelkan hubungan antara independent variable x dan dependent variable y sebagai polynomial dalam x

Proses Dataset



```
[17]: X_train = np.array(train_pizza_df['diameter']).reshape(-1, 1)
      y_train = np.array(train_pizza_df['harga'])

      print(f'X_train:\n{X_train}\n')
      print(f'y_train: {y_train}')
```

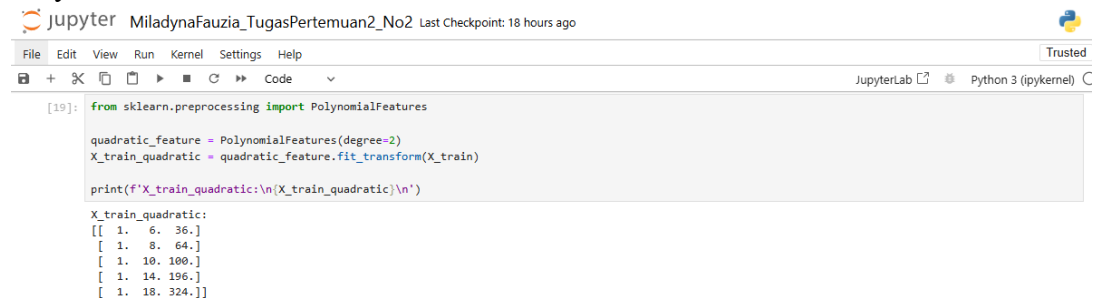
X_train:
[[6]
 [8]
 [10]
 [14]
 [18]]

y_train: [7. 9. 13. 17.5 18.]

2.5. Quadratic Polynomial Regression

$$y = a + \beta_1 x + \beta_2 x^2$$

Polynomial Features



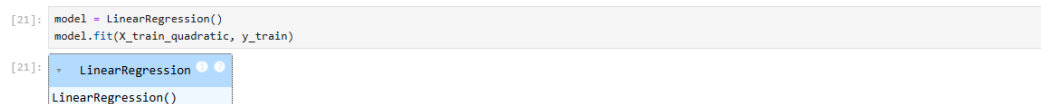
```
[19]: from sklearn.preprocessing import PolynomialFeatures

      quadratic_feature = PolynomialFeatures(degree=2)
      X_train_quadratic = quadratic_feature.fit_transform(X_train)

      print(f'X_train_quadratic:\n{X_train_quadratic}\n')
```

X_train_quadratic:
[[1. 6. 36.]
 [1. 8. 64.]
 [1. 10. 100.]
 [1. 14. 196.]
 [1. 18. 324.]]

Training Model

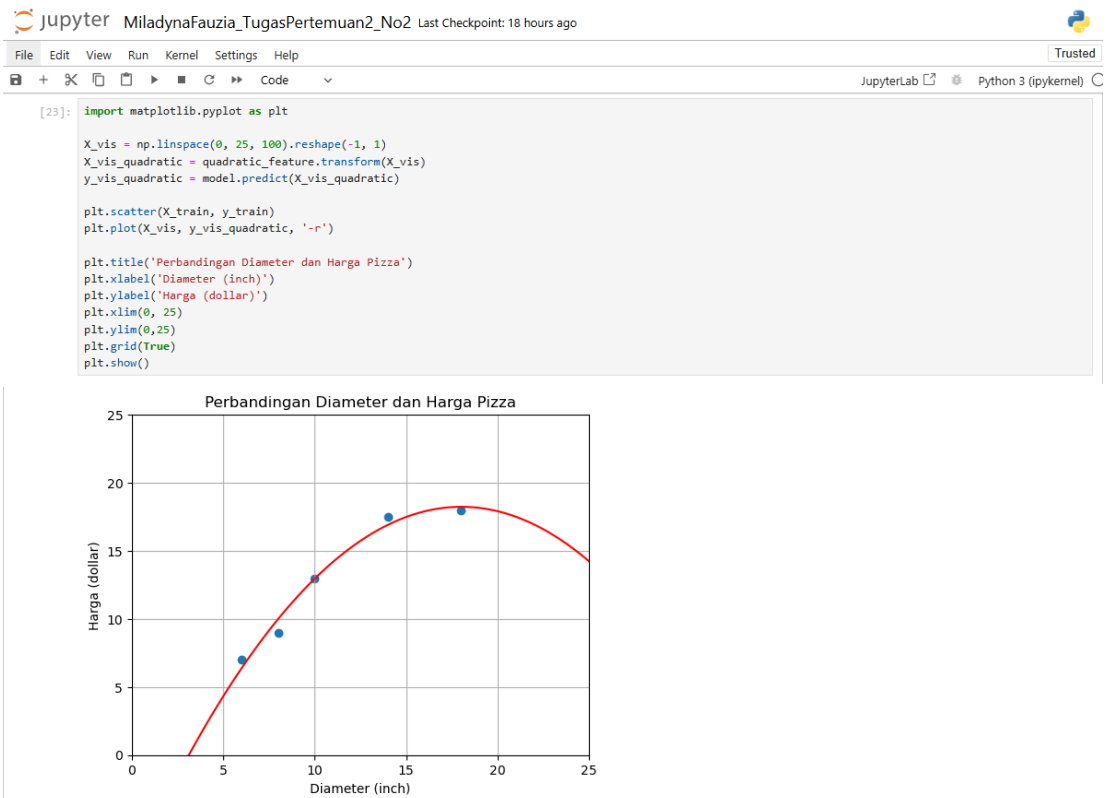


```
[21]: model = LinearRegression()
      model.fit(X_train_quadratic, y_train)
```

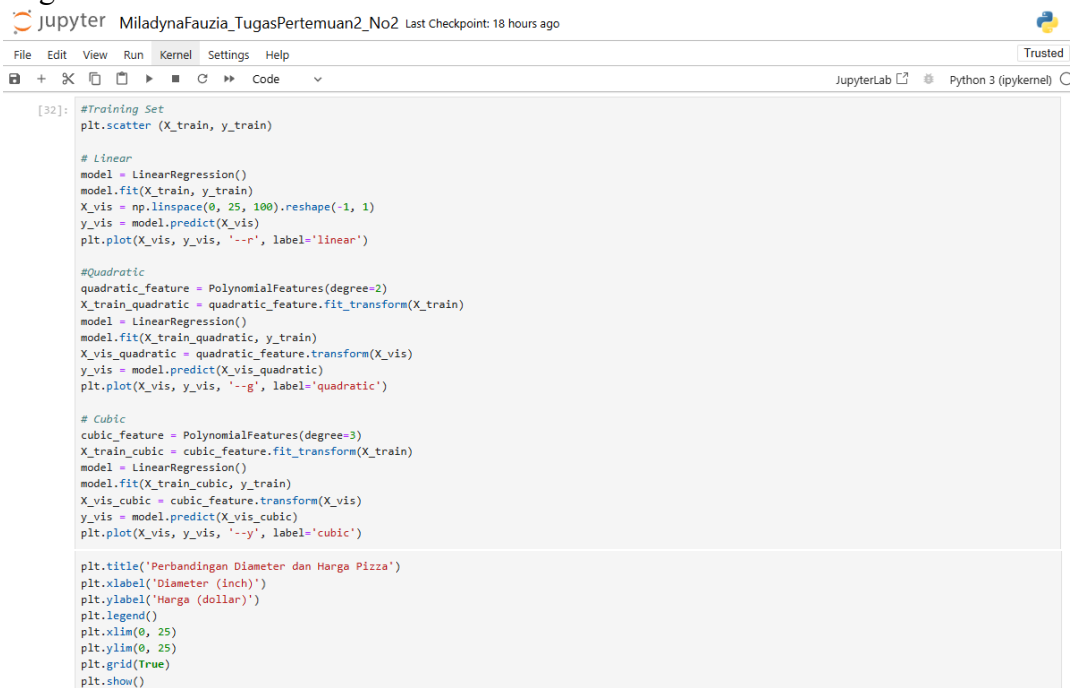
[21]: LinearRegression

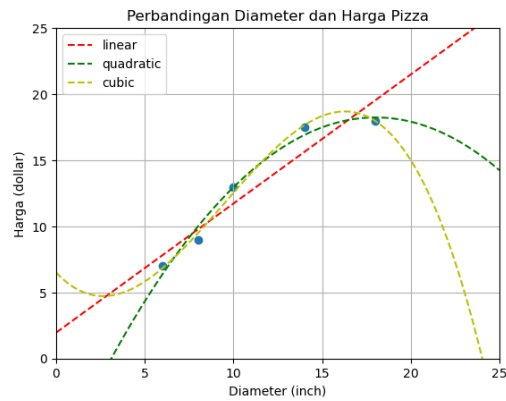
LinearRegression()

Visualisasi Model



2.6. Linear Regression vs Quadratic Polynomial Regression vs Cubic Polynomial Regression





3. Lakukan praktek dari <https://youtu.be/oe7DW4rSH1o?si=H-PZJ9rs9-Kab-Ln> dan buat screen shot hasil run dengan nama anda pada hasil run tersebut. Praktek tersebut yaitu: Logistic Regression pada Binary Classification Task

3.1. Formula dasar pembentuk Logistic Regression | Fungsi Sigmoid

Simple Linear Regression

- $y = a + \beta x$
- $g(x) = a + \beta x$

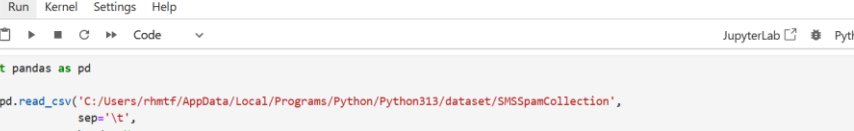
Multiple Linear Regression

- $y = a + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$
- $g(X) = a + \beta X$

Logistic Regression

- $g(X) = \text{sigmoid}(a + \beta X)$
- $\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$

3.2. Persiapan dataset | SMS Spam Collection Dataset



The screenshot shows a JupyterLab window with a single code cell. The menu bar at the top includes File, Edit, View, Run, Kernel, Settings, and Help. The toolbar below the menu contains icons for saving, opening, and running files, as well as a code editor icon. The code cell contains the following Python code:

```
[4]: import pandas as pd

df = pd.read_csv('C:/Users/rhmtf/AppData/Local/Programs/Python/Python313/dataset/SMS SpamCollection',
                 sep='\t',
                 header=None,
                 names=['label', 'sms'])

df.head()
```

The output of the code cell is displayed below the code, showing the first five rows of the DataFrame:

```
[4]:
```

	label	sms
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...

Below the output table, the next code cell is visible, showing the command to count the values in the 'label' column:

```
[5]: df['label'].value_counts()
```

The output for this command is shown at the bottom of the screenshot:

```
[5]: label
ham    4825
spam    747
Name: count, dtype: int64
```


3.3. Pembagian training dan testing dataset

```
jupyter Miladyna Fauzia Last Checkpoint: 1 hour ago
File Edit View Run Kernel Settings Help
JupyterLab Python 3 (ipykernel) Trusted

[7]: from sklearn.preprocessing import LabelBinarizer

X = df['sms'].values
y = df['label'].values

lb = LabelBinarizer()
y = lb.fit_transform(y).ravel()
lb.classes_

[7]: array(['ham', 'spam'], dtype='<U4')

[9]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X,
                                                    y,
                                                    test_size=0.25,
                                                    random_state=0)

print(X_train, '\n')
print(y_train)

['Its going good...no problem..but still need little experience to understand american customer voice...'
'U have a secret admirer. REVEAL who thinks U R So special. Call 09065174042. To opt out Reply REVEAL STOP. 1.50 per msg recd. Cust care 07821238901'
'Ok...' ...
'For ur chance to win a £250 cash every wk TXT: ACTION to 80608. T's&C's www.movietrivia.tv custcare 08712405022, 1x150p/wk'
'R U 8SAM P IN EACHOTHER. IF WE MEET WE CAN GO 2 MY HOUSE'
'Mm feeling sleepy. today itself i shall get that dear']

[0 1 0 ... 1 0 0]
```

3.4. Feature extraction dengan TF-IDF

```
jupyter Miladyna Fauzia Last Checkpoint: 1 hour ago
File Edit View Run Kernel Settings Help
JupyterLab Python 3 (ipykernel) Trusted

[13]: from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer = TfidfVectorizer(stop_words='english')

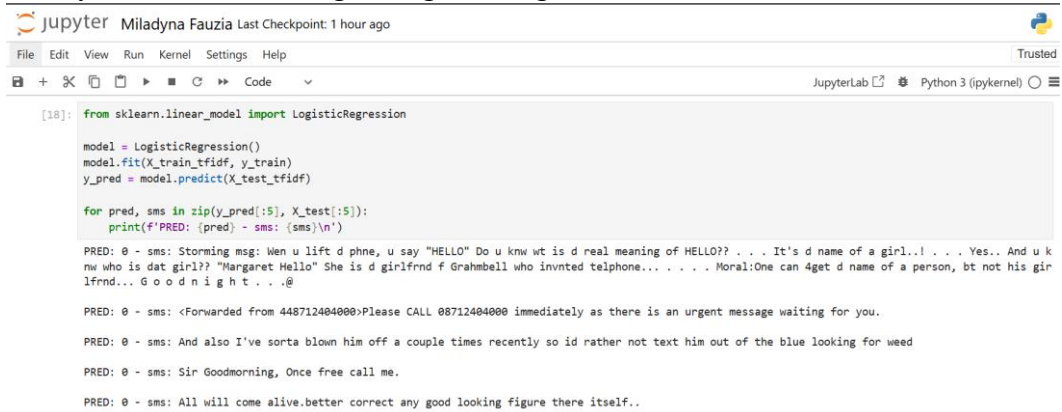
X_train_tfidf = vectorizer.fit_transform(X_train)
X_test_tfidf = vectorizer.transform(X_test)

print(X_train_tfidf)

<Compressed Sparse Row sparse matrix of dtype 'float64'
with 32656 stored elements and shape (4179, 7287)>

Coords      Values
(0, 2997)    0.23173982975834367
(0, 3007)    0.21421364386658514
(0, 5123)    0.308974289326673
(0, 4453)    0.2297719954323795
(0, 3926)    0.3126721340000456
(0, 2554)    0.3825278811525034
(0, 6739)    0.3546359942838148
(0, 900)     0.4114867709157148
(0, 2006)    0.2890802580285881
(0, 6903)    0.3591386422223876
(1, 5642)    0.24344998442381355
(1, 799)     0.25048918791828574
(1, 5441)    0.5009783758205715
(1, 6472)    0.24039776602646504
(1, 6013)    0.20089911182610476
(1, 216)     0.28902673040368515
(1, 4677)    0.24039776602646504
(1, 5394)    0.16464655071448758
(1, 6131)    0.16142609035094446
(1, 532)     0.20186022353306565
(1, 4358)    0.17341410292348694
...
(1, 5301)    0.2711077935907125
(1, 2003)    0.2711077935907125
(1, 1548)    0.18167737976542422
(1, 36)      0.28902673040368515
:
:
(4176, 6792) 0.1407604617250961
(4176, 6693) 0.16491299289150899
(4176, 6684) 0.22114159453800114
(4176, 7083) 0.19523751585154273
(4176, 1569) 0.18895085073406012
(4176, 7195) 0.17892283441772988
(4176, 779)  0.2811868572055718
(4176, 1612) 0.21138425595332702
(4176, 365)  0.2380005587702937
(4176, 7114) 0.4512018097459442
(4176, 637)  0.29968668460649284
(4176, 4350) 0.29968668460649284
(4176, 2004) 0.25589560236817055
(4176, 107)  0.29968668460649284
(4176, 343)  0.2811868572055718
(4177, 3319) 0.43046342221720785
(4177, 4177) 0.3636187667918345
(4177, 5565) 0.5506066649743346
(4177, 2362) 0.6158854885899457
(4178, 2068) 0.3055766821331892
(4178, 2641) 0.3993042639531407
(4178, 6555) 0.2897850627168302
(4178, 5720) 0.3963527249882828
(4178, 4279) 0.4530624713751054
(4178, 5883) 0.548491137555895
```

3.5. Binary Classification dengan Logistic Regression



```
[18]: from sklearn.linear_model import LogisticRegression

model = LogisticRegression()
model.fit(X_train_tfidf, y_train)
y_pred = model.predict(X_test_tfidf)

for pred, sms in zip(y_pred[:5], X_test[:5]):
    print(f'PRED: {pred} - sms: {sms}\n')
```

PRED: 0 - sms: Storming msg: Wen u lift d phne, u say "HELLO" Do u knw wt is d real meaning of HELLO?? . . . It's d name of a girl!..! . . . Yes.. And u k nw who is dat girl?? "Margaret Hello" She is d girlfrnd f Grahmbell who invnted telphone... . . . Moral:One can 4get d name of a person, bt not his gir lfrnd... G o o d n i g h t . . .@

PRED: 0 - sms: <Forwarded from 448712404000>Please CALL 08712404000 immediately as there is an urgent message waiting for you.

PRED: 0 - sms: And also I've sorta blown him off a couple times recently so id rather not text him out of the blue looking for weed

PRED: 0 - sms: Sir Goodmorning, Once free call me.

PRED: 0 - sms: All will come alive.better correct any good looking figure there itself..

3.6. Evaluation Metrics pada Binary Classification Task

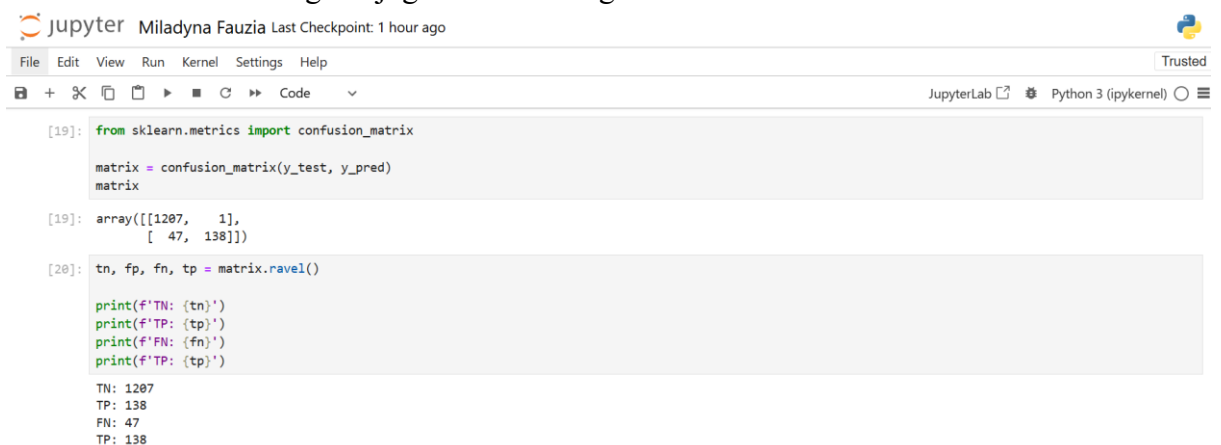
- Confusion Matrix
- Accuracy
- Precision dan Recall
- F1 Score
- ROC

Terminologi Dasar

- True Positive (TP)
- True Negative (TN)
- False Positive (FP)
- False Negative (FN)

3.7. Pengenalan Confusion Matrix

Confusion Matrix seringkali juga dikenal sebagai error matrix



```
[19]: from sklearn.metrics import confusion_matrix

matrix = confusion_matrix(y_test, y_pred)
matrix

[19]: array([[1207,  1],
        [ 47, 138]])

[20]: tn, fp, fn, tp = matrix.ravel()

print(f'TN: {tn}')
print(f'TP: {tp}')
print(f'FN: {fn}')
print(f'FP: {fp}')
```

TN: 1207
TP: 138
FN: 47
TP: 138



3.8. Pengenalan Accuracy Score

Accuracy mengukur porsi dari hasil prediksi yang tepat

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} = \frac{correct}{total}$$

Jupyter Miladyna Fauzia Last Checkpoint: 1 hour ago

```
[22]: from sklearn.metrics import accuracy_score

accuracy_score(y_test, y_pred)
```

[22]: 0.9655419956927495

3.9. Pengenalan Precision dan Recall

Selain menggunakan accuracy, performa dari suatu classifier umumnya juga diukur berdasarkan nilai Precision dan Recall.

Precision or Positive Predictive Value (PPV)

$$Precision = \frac{TP}{TP + FP}$$

```
[18]: from sklearn.metrics import precision_score

precision_score(y_test, y_pred)
```

[18]: np.float64(0.9928857553956835)

Recall or True Positive Rate (TPR) or Sensitivity

$$Recall = \frac{TP}{TP + FN}$$

```
[23]: from sklearn.metrics import recall_score

recall_score(y_test, y_pred)
```

[23]: np.float64(0.745945945945946)

3.10. Pengenalan F1 Score | F1 Measure

F1-score atau F1-measure adalah harmonic mean dari precision dan recall.

$$F1\ score = \frac{precision \times recall}{precision + recall}$$

```
[24]: from sklearn.metrics import f1_score

f1_score(y_test, y_pred)
```

[24]: np.float64(0.8518518518518519)

3.11. Pengenalan ROC | Receiver Operating Characteristic

Jupyter Miladyna Fauzia Last Checkpoint: 1 hour ago

```
[25]: from sklearn.metrics import roc_curve, auc

prob_estimates = model.predict_proba(X_test_tfidf)

fpr, tpr, threshold = roc_curve(y_test, prob_estimates[:, 1])
nilai_auc = auc(fpr, tpr)

plt.plot(fpr, tpr, 'b', label=f'AUC={nilai_auc}')
plt.plot([0,1], [0,1], 'r--', label='Random Classifier')

plt.title('ROC: Receiver Operating Characteristic')
plt.xlabel('Fallout or False Positive Rate')
plt.ylabel('Recall or True Positive Rate')
plt.legend()
plt.show()
```

