

Tugas 2: Praktikum & Praktikum Mandiri 2

Pandu Linggar Kumara - 0110221277,
Link GitHub - https://github.com/PanduLgg/M_Learning.git

¹ Teknik Informatika, STT Terpadu Nurul Fikri, Depok

*E-mail: pandulinggar1@gmail.com

Abstract. Kegiatan praktikum mencakup pembacaan dataset, perhitungan nilai-nilai statistik dasar (seperti mean, median, modus, variansi, standar deviasi, dan kuartil), serta analisis korelasi antarvariabel. Selain itu, dilakukan pula visualisasi data menggunakan boxplot, histogram, dan scatter plot untuk memahami distribusi serta hubungan antar variabel.

1. Connecting Google Colab & Drive

1.1 Menghubungkan lingkungan Google Colab dengan akun Google Drive

Sel ini berfungsi untuk menghubungkan lingkungan Google Colab dengan akun Google Drive

```
# menghubungkan colab dengan google drive
from google.colab import drive
drive.mount('/content/gdrive')

Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive.mount("/content/gdrive", force_remount=True).
```

Gambar 1.1. Proses ini hanya perlu dilakukan satu kali per sesi.

1.2 Memanggil Data set dari Gdrive dan Membaca file .CSV menggunakan Pandas

Sel ini menggunakan library Pandas untuk membaca file data, yang diinginkan yaitu 'data/500_Person_Gender_Height_Weight_Index.csv'.

```
# memanggil data set lewat gdrive
path = "/content/gdrive/MyDrive/praktikum_ml/praktikum02/"

# membaca file csv menggunakan pandas
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

df = pd.read_csv(path + 'data/500_Person_Gender_Height_Weight_Index.csv')
df
```

	Gender	Height	Weight	Index
0	Male	174	96	4
1	Male	189	87	2
2	Female	185	110	4
3	Female	195	104	3
4	Male	149	61	3
...
495	Female	150	153	5
496	Female	184	121	4
497	Female	141	136	5
498	Male	150	95	5
499	Male	173	131	5

500 rows x 4 columns

Gambar 1.2. Proses ini hanya perlu dilakukan satu kali per sesi.

1.3 Mencari informasi data yang ada pada file

Sel ini menampilkan informasi yang ada di dalam file dari mulai tipe data nama kolom, dsb.

```
# Mencari info data pada file (tipe datanya, non nul count data, nama kolom)
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Gender      500 non-null    object
1   Height      500 non-null    int64
2   Weight      500 non-null    int64
3   Index       500 non-null    int64
dtypes: int64(3), object(1)
memory usage: 15.8+ KB
```

Gambar 1.3. Mencari info data pada file

1.4 Menghitung Ukuran Tendensi Sentral dan Ukuran Penyebaran Data

Sel ini menghitung mean, median, modus, variansi, dan standard deviasi pada file data

```
# Menghitung mean semua kolom numerik
df['Height'].mean()

np.float64(169.944)

# Menghitung Median semua kolom numerik
df['Height'].median()

170.5

# Mencari Modus bisa lebih dari satu
df['Height'].mode()

Height
0      188
dtype: int64

# Menghitung Variansi & Standard Deviasi
df.var(numeric_only=True)

Height      268.149162
Weight     1048.633267
Index        1.836168
dtype: float64

# Menghitung Standar Deviasi
df.std(numeric_only=True)

Height      16.375261
Weight      32.382607
Index        1.355053
dtype: float64
```

Gambar 1.4. Menghitung Ukuran Tendensi Sentral dan Ukuran Penyebaran Data

1.5 Menghitung Quartil dan Membuat Deskripsi pada Data

Sel ini menghitung q1, q2, dan IQR, lalu membuat deskripsi pada type data Integer, lalu menghitung matriks korelasi dalam semua kolom numerik

```
[ ] # Menghitung Kuartil Pertama (Q1)
q1 = df['Height'].quantile(0.25)
print("Q1 : ", q1)

#Hitung Kuartil Ketiga (Q3)
q3 = df['Height'].quantile(0.75)
print("Q3 : ", q3)

#Hitung IQR (interquartile Range)
iqr = q3 - q1
print("IQR : ", iqr)
```

Q1 : 156.0
Q3 : 184.0
IQR : 28.0

```
[ ] # Untuk membuat statistika deskripsi pada type data int
df.describe()
```

	Height	Weight	Index
count	500.000000	500.000000	500.000000
mean	169.944000	106.000000	3.748000
std	16.375261	32.382607	1.355053
min	140.000000	50.000000	0.000000
25%	156.000000	80.000000	3.000000
50%	170.500000	106.000000	4.000000
75%	184.000000	136.000000	5.000000
max	199.000000	160.000000	5.000000

```
[ ] # Menghitung matriks korelasi untuk semua kolom numerik
correlation_matrix = df.corr(numeric_only=True)

# Menampilkan Matriks Korelasi
print("Matriks Korelasi:")
print(correlation_matrix)
```

Matriks Korelasi:

	Height	Weight	Index
Height	1.000000	0.000446	-0.422223
Weight	0.000446	1.000000	0.804569
Index	-0.422223	0.804569	1.000000

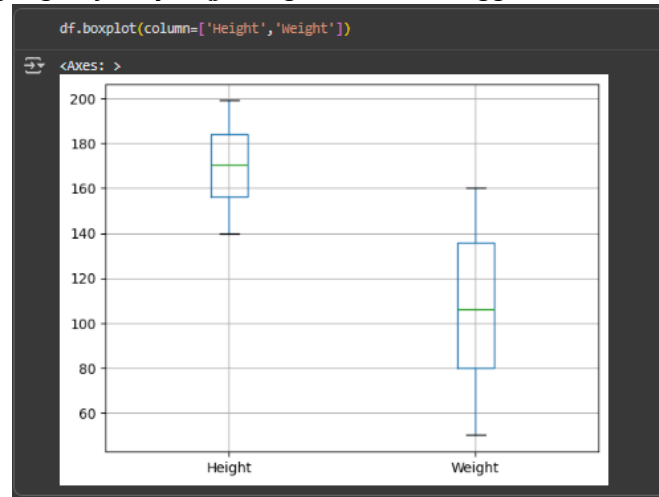
Gambar 1.5. Menghitung Quartil dan Membuat Deskripsi pada Data

2. Visualisasi Data

Memvisualisasikan Data menggunakan matplotlib, dengan metode Box Plot, Histogram, Scatter Plot

2.1 Visualisasi Data menggunakan Box Plot

Sel ini menggunakan fungsi `df.boxplot()` mengambil data Tinggi dan Berat Badan



Gambar 2.1. Menampilkan Box Plot

2.2 Visualisasi Data menggunakan Histogram

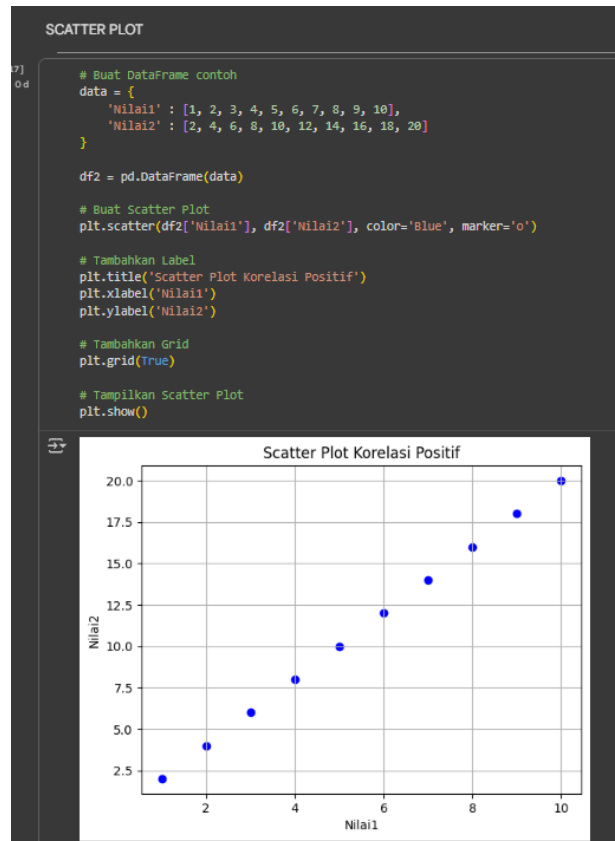
Sel ini membuat histogram untuk kolom 'Height'. Histogram adalah grafik yang menunjukkan distribusi frekuensi data



Gambar 2.2. Menampilkan Histogram

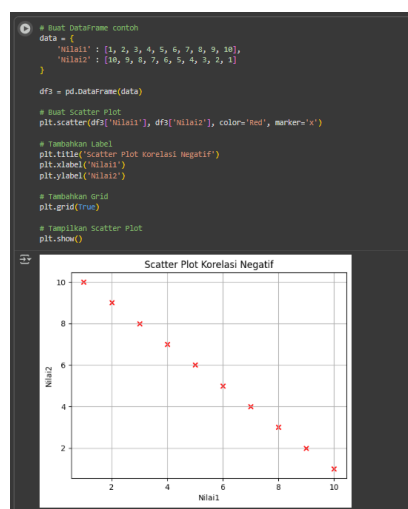
2.3 Visualisasi Data menggunakan Scatter Plot

Sel ini membuat scatter plot atau diagram pencar untuk memvisualisasikan hubungan antara dua variabel numerik. Berisi dua blok kode untuk membuat dua scatter plot terpisah. Kode pertama menggunakan data Nilai1 dan Nilai2 yang memiliki korelasi positif



Gambar 2.3. Menampilkan Scatter Plot Positif

sedangkan kode kedua memiliki korelasi negatif.



Gambar 2.4. Menampilkan Scatter Plot Negatif

Hasil:

1. Plot Korelasi Positif: Titik-titik data membentuk pola yang naik ke atas. Ini menunjukkan bahwa ketika Nilai1 meningkat, Nilai2 juga cenderung meningkat.
2. Plot Korelasi Negatif: Titik-titik data membentuk pola yang menurun. Ini menunjukkan bahwa ketika Nilai1 meningkat, Nilai2 justru cenderung menurun.

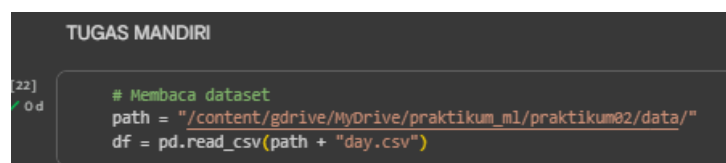
3. Praktikum Mandiri

Buat program untuk membagi dataset day.csv menjadi tiga bagian, yaitu:

- (a) Data Training: 80% dari total dataset
- (b) Data Validation: 10% dari data training
- (c) Data Testing: 20% dari total dataset

3.1 Membaca Dataset File

Sel ini membaca dataset yang diinginkan yaitu day.csv



```
[22]:  
# Membaca dataset  
path = "/content/gdrive/MyDrive/praktikum_ml/praktikum02/data/"  
df = pd.read_csv(path + "day.csv")
```

Gambar 3.1. Proses ini hanya perlu dilakukan satu kali per sesi.

3.2 Membaca Dataset File

Sel ini membagi data Training 80%, Data Testing 20%, dan Data validation 10%, lalu menampilkan jumlah datanya



```
[23]:  
# Split 80% training dan 20% testing  
train_data, test_data = train_test_split(df, test_size=0.2, random_state=42)  
  
# Split 10% validation dari data training  
train_data, val_data = train_test_split(train_data, test_size=0.1, random_state=42)  
  
# Tampilkan jumlah data  
print("Jumlah data total: {len(df)}")  
print("Training: {len(train_data)} data")  
print("Validation: {len(val_data)} data")  
print("Testing: {len(test_data)} data")  
  
Jumlah data total: 731  
Training: 525 data  
Validation: 59 data  
Testing: 147 data
```

Gambar 3.2. Menampilkan Jumlah Data hasil pembagian persentase

Referensi:

- Munir, S., Seminar, K. B., Sudradjat, Sukoco, H., & Buono, A. (2022). The Use of Random Forest Regression for Estimating Leaf Nitrogen Content of Oil Palm Based on Sentinel 1-A Imagery. *Information*, 14(1), 10. <https://doi.org/10.3390/info14010010>
- Seminar, K. B., Imantho, H., Sudradjat, Yahya, S., Munir, S., Kaliana, I., Mei Haryadi, F., Noor Baroroh, A., Supriyanto, Handoyo, G. C., Kurnia Wijayanto, A., Ijang Wahyudin, C., Liyantono, Budiman, R., Bakir Pasaman, A., Rusiawan, D., & Sulastri. (2024). PreciPalm: An Intelligent System for Calculating Macronutrient Status and Fertilizer Recommendations for Oil Palm on Mineral Soils Based on a Precision Agriculture Approach. *Scientific World Journal*, 2024(1). <https://doi.org/10.1155/2024/1788726>