

LAPORAN PERTEMUAN 4

Nama : Muhammad Salman Imamwan Abdillah

Nim : 231011401032

Kelas : 05TPLE017

1. Kode ini digunakan untuk membaca data mahasiswa dari file CSV, menampilkan informasi dan isi datanya, serta memastikan datanya siap digunakan untuk analisis lebih lanjut.

```
import pandas as pd
df = pd.read_csv("kelulusan_mahasiswa.csv")
print(df.info())
print(df.head())
```

[2] ✓ 29.5s

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
#   Column             Non-Null Count  Dtype  
---  --
0    IPK                 10 non-null    float64
1    Jumlah_Absensi      10 non-null    int64  
2    Waktu_Belajar_Jam   10 non-null    int64  
3    Lulus               10 non-null    int64  
dtypes: float64(1), int64(3)
memory usage: 448.0 bytes
None
```

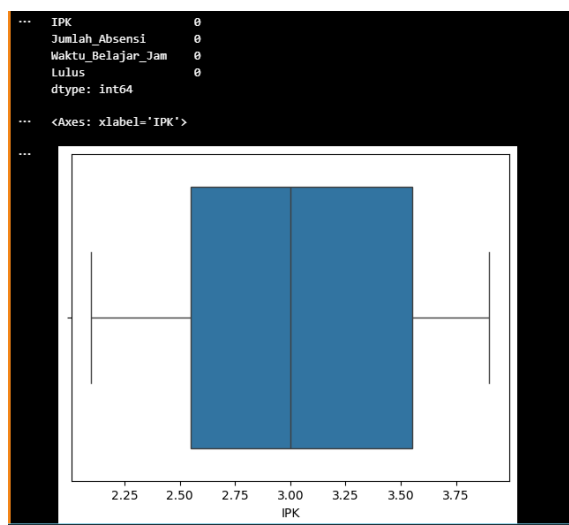
	IPK	Jumlah_Absensi	Waktu_Belajar_Jam	Lulus
0	3.8	3	10	1
1	2.5	8	5	0
2	3.4	4	7	1
3	2.1	12	2	0
4	3.9	2	12	1

2. Kode ini digunakan untuk memastikan data tidak memiliki nilai kosong maupun data ganda, serta untuk melihat distribusi nilai IPK menggunakan grafik boxplot. Hasilnya menunjukkan bahwa data bersih, lengkap, dan nilai IPK mahasiswa terdistribusi normal tanpa adanya outlier.

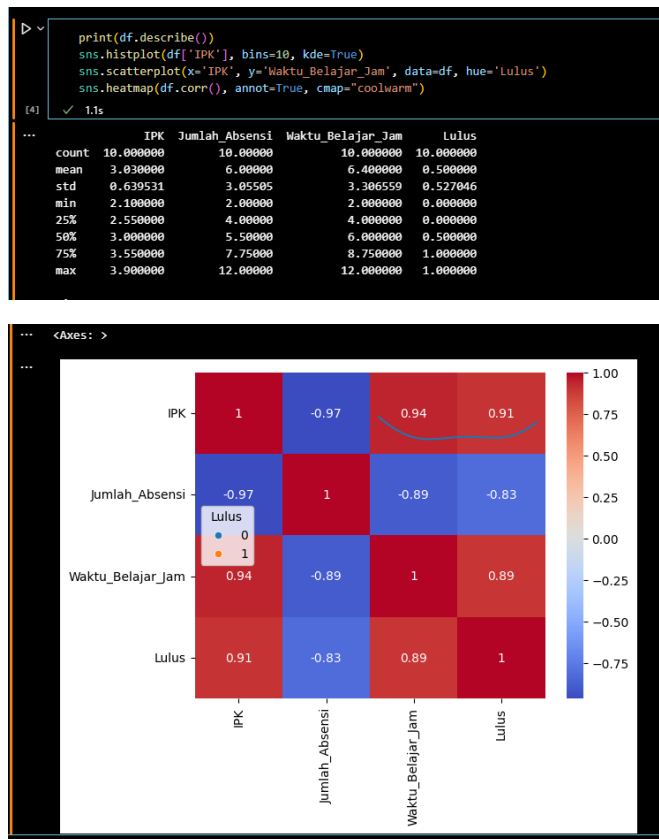
```
print(df.isnull().sum())
df = df.drop_duplicates()

import seaborn as sns
sns.boxplot(x=df['IPK'])
```

[3] ✓ 47.7s



- Analisis dilakukan untuk melihat hubungan antara IPK, jumlah absensi, waktu belajar, dan kelulusan mahasiswa. Berdasarkan hasil statistik deskriptif dan visualisasi data, diketahui bahwa IPK dan waktu belajar memiliki korelasi positif yang kuat terhadap kelulusan, sedangkan jumlah absensi berkorelasi negatif. Hal ini menunjukkan bahwa mahasiswa dengan IPK tinggi dan waktu belajar lebih banyak cenderung lulus, sementara absensi yang tinggi menurunkan peluang kelulusan.



- Pada tahap ini dilakukan proses pembuatan fitur baru dan pembagian dataset untuk keperluan pelatihan model. Dua fitur tambahan, yaitu *Rasio_Absensi* (hasil pembagian jumlah absensi dengan total 14 pertemuan) dan *IPK_x_Study* (perkalian antara IPK dan waktu belajar), ditambahkan untuk memperkaya informasi dalam data. Selanjutnya, data disimpan ke file *processed_kelulusan.csv* dan dibagi menjadi tiga bagian menggunakan metode *train_test_split*, yaitu data latih (60%), data validasi (20%), dan data uji (20%). Pembagian dilakukan secara stratified agar proporsi kelas tetap seimbang.

