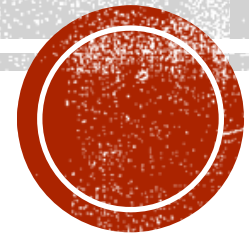


# **PENGOLAHAN DATA DAN VISUALISASI DENGAN PYTHON**

Week 2

Dosen: Yusuf Fadlila Rachman, S.Kom., M.Kom



# INSTALASI LIBRARY PYTHON

- -> Open Anaconda Prompt / CMD
- Pip install pandas
- Pip install numpy
- Pip install matplotlib
- -> Open Jupyter Notebook



# ARRAY, LIST, DATAFRAME

Fitur	Array	List	DataFrame
Definisi	Struktur data yang berisi elemen dengan tipe data yang sama.	Koleksi elemen yang dapat berisi tipe data berbeda.	Struktur data dua dimensi dengan baris dan kolom, biasanya digunakan untuk data tabular.
Tipe Data	Homogen (semua elemen harus bertipe sama).	Heterogen (elemen bisa memiliki tipe berbeda).	Homogen untuk setiap kolom, tetapi bisa berbeda antar kolom.
Ukuran	Ukuran tetap setelah deklarasi.	Ukuran bisa dinamis, bisa bertambah atau berkurang.	Ukuran bisa dinamis, tetapi lebih terstruktur dan besar.



# ARRAY, LIST, DATAFRAME

Fitur	Array	List	DataFrame
Operasi Matematika	Mendukung operasi vektor dan matriks secara langsung (terutama pada <b>NumPy</b> ).	Tidak mendukung operasi matematika langsung.	Mendukung operasi matematika dan statistik, terutama menggunakan <b>pandas</b> .
Fungsi dan Metode	Terbatas pada fungsi dasar yang disediakan oleh pustaka seperti <b>NumPy</b> .	Memiliki beberapa fungsi dasar untuk manipulasi data.	Memiliki banyak fungsi untuk manipulasi dan analisis data, seperti <b>groupby</b> , <b>merge</b> , <b>pivot</b> , dll.
Akses Elemen	Indeks berbasis integer.	Indeks berbasis integer atau elemen bisa diakses dengan indeks atau nilai.	Akses elemen bisa menggunakan nama kolom atau indeks baris.



# ARRAY, LIST, DATAFRAME

Fitur	Array	List	DataFrame
Kecepatan	Lebih cepat dalam pemrosesan matematis, terutama untuk data numerik.	Relatif lebih lambat jika digunakan untuk operasi numerik kompleks.	Sedikit lebih lambat dibanding array atau list dalam operasi numerik, tetapi sangat efisien untuk data tabular besar.
Penggunaan Umum	Digunakan untuk operasi numerik, komputasi ilmiah, dan algoritma berbasis angka.	Digunakan untuk koleksi data heterogen atau data yang tidak terstruktur.	Digunakan untuk analisis data dan statistik, terutama dengan data tabular besar.



# NUMPY

## Code Examples :

```
import numpy as np
```

```
data = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
```

```
std_dev = np.std(data)  
variance = np.var(data)
```

```
print(f"Standar Deviasi: {std_dev}")  
print(f"Varians: {variance}")
```

Fungsi / Metode	Deskripsi	Contoh Penggunaan
<code>np.mean()</code>	Menghitung rata-rata dari array.	<code>np.mean([1, 2, 3, 4, 5])</code> akan menghasilkan <code>3.0</code> .
<code>np.median()</code>	Menghitung median dari array.	<code>np.median([1, 2, 3, 4, 5])</code> akan menghasilkan <code>3.0</code> .
<code>np.std()</code>	Menghitung deviasi standar dari array (penyebaran data).	<code>np.std([1, 2, 3, 4, 5])</code> akan menghasilkan <code>1.41</code> .
<code>np.var()</code>	Menghitung varians dari array.	<code>np.var([1, 2, 3, 4, 5])</code> akan menghasilkan <code>2.0</code> .
<code>np.min()</code> dan <code>np.max()</code>	Menghitung nilai minimum dan maksimum dari array.	<code>np.min([1, 2, 3, 4, 5])</code> akan menghasilkan <code>1</code> dan <code>np.max()</code> menghasilkan <code>5</code> .
<code>np.percentile()</code>	Menghitung persentil ke-n dari array.	<code>np.percentile([1, 2, 3, 4, 5], 50)</code> akan menghasilkan <code>3.0</code> .
<code>np.corrcoef()</code>	Menghitung koefisien korelasi antara	<code>np.corrcoef([1, 2, 3], [4, 5, 6])</code> akan



# NUMPY

## Code Example :

```
import numpy as np
import pandas as pd
from scipy import stats
```

```
data = [5, 7, 8, 9, 9]
```

```
mean = np.mean(data)
median = np.median(data)
mode = stats.mode(data)[0][0]
```

```
print(f"Mean: {mean}")
print(f"Median: {median}")
print(f"Mode: {mode}")
```

<code>np.corrcoef()</code>	Menghitung koefisien korelasi antara dua variabel.	<code>np.corrcoef([1, 2, 3], [4, 5, 6])</code> akan menghasilkan matriks korelasi.
<code>np.random.normal()</code>	Membuat data acak yang terdistribusi normal.	<code>np.random.normal(0, 1, 1000)</code> menghasilkan data acak dengan distribusi normal.
<code>np.random.binomial()</code>	Membuat data acak berdasarkan distribusi binomial.	<code>np.random.binomial(10, 0.5, 1000)</code> menghasilkan distribusi binomial dengan n=10 dan p=0.5.
<code>np.random.uniform()</code>	Membuat data acak dengan distribusi seragam.	<code>np.random.uniform(0, 1, 1000)</code> menghasilkan angka acak antara 0 dan 1.
<code>np.random.choice()</code>	Memilih elemen secara acak dari array atau list.	<code>np.random.choice([1, 2, 3, 4, 5], size=3)</code> menghasilkan 3 elemen acak dari array tersebut.
<code>np.random.seed()</code>	Menetapkan nilai awal (seed) untuk generator bilangan acak, memungkinkan replikasi hasil.	<code>np.random.seed(42)</code> memastikan hasil acak dapat direplikasi.



# PANDAS

## Code Example :

```
import pandas as pd

data = {'Nama': ['Andi', 'Budi', 'Cici', 'Diana', 'Eka'],
        'Usia': [22, 23, 21, 25, 24]}

df = pd.DataFrame(data)

rata_rata = df['Usia'].mean()
varians = df['Usia'].var()

print("Rata-rata Usia:", rata_rata)
print("Varians Usia:", varians)
```

Fungsi / Metode	Deskripsi	Contoh Penggunaan
<code>df.mean()</code>	Menghitung rata-rata dari setiap kolom dalam DataFrame.	<code>df['age'].mean()</code> akan menghasilkan rata-rata usia dari kolom <code>age</code> .
<code>df.median()</code>	Menghitung median dari setiap kolom dalam DataFrame.	<code>df['age'].median()</code> akan menghasilkan median usia dari kolom <code>age</code> .
<code>df.std()</code>	Menghitung deviasi standar dari setiap kolom dalam DataFrame.	<code>df['age'].std()</code> akan menghasilkan deviasi standar usia dari kolom <code>age</code> .
<code>df.var()</code>	Menghitung varians dari setiap kolom dalam DataFrame.	<code>df['age'].var()</code> akan menghasilkan varians usia dari kolom <code>age</code> .
<code>df.min()</code> dan <code>df.max()</code>	Menghitung nilai minimum dan maksimum dari setiap kolom dalam DataFrame.	<code>df['age'].min()</code> menghasilkan nilai usia minimum, <code>df['age'].max()</code> menghasilkan nilai usia maksimum.
<code>df.describe()</code>	Memberikan ringkasan statistik dari setiap kolom numerik dalam DataFrame (mean, std, min, max, dll).	<code>df.describe()</code> akan memberikan ringkasan statistik untuk kolom numerik di DataFrame.





# PANDAS

## Code Example Menghitung Korelasi Antar Kolom :

```
data = {'Usia': [22, 23, 21, 25, 24],  
        'Jam_Belajar': [15, 18, 10, 20, 17]}
```

```
df = pd.DataFrame(data)
```

```
# Menghitung korelasi antara  
Usia dan Jam Belajar  
korelasi = df.corr()
```

```
print("Matriks Korelasi:")  
print(korelasi)
```

<code>df.corr()</code>	Menghitung koefisien korelasi antara semua kolom numerik dalam DataFrame.	<code>df.corr()</code> menghasilkan matriks korelasi antar kolom numerik.
<code>df.cov()</code>	Menghitung kovarians antar kolom dalam DataFrame.	<code>df.cov()</code> menghasilkan matriks kovarians antar kolom numerik.
<code>df.value_counts()</code>	Menghitung jumlah nilai unik dalam suatu kolom.	<code>df['gender'].value_counts()</code> menghitung frekuensi nilai unik dalam kolom <code>gender</code> .
<code>df.pivot_table()</code>	Membuat pivot table untuk menganalisis data berdasarkan agregasi.	<code>df.pivot_table(values='age', index='gender', aggfunc='mean')</code> akan menghitung rata-rata usia berdasarkan gender.
<code>df.groupby()</code>	Mengelompokkan data berdasarkan satu atau lebih kolom dan menerapkan agregasi statistik.	<code>df.groupby('gender')['age'].mean()</code> menghitung rata-rata usia berdasarkan gender.
<code>df.quantile()</code>	Menghitung kuantil ke-n (misalnya kuartil) dari suatu kolom.	<code>df['age'].quantile(0.25)</code> menghasilkan kuartil pertama dari kolom <code>age</code> .



# PANDAS

<code>df.sample()</code>	Mengambil sampel acak dari DataFrame.	<code>df.sample(5)</code> menghasilkan 5 sampel acak dari DataFrame.
<code>df.isna()</code>	Mengidentifikasi nilai yang hilang (missing values) dalam DataFrame.	<code>df.isna()</code> menghasilkan DataFrame boolean yang menunjukkan nilai yang hilang.
<code>df.fillna()</code>	Mengisi nilai yang hilang dengan nilai yang ditentukan.	<code>df['age'].fillna(df['age'].mean())</code> mengisi nilai hilang pada kolom <code>age</code> dengan rata-ratanya.
<code>df.random_state</code>	Menetapkan nilai acak pada operasi seperti sampling atau pemisahan data untuk reproduksi.	<code>df.sample(5, random_state=42)</code> memastikan pengambilan sampel yang dapat direplikasi.

Jika kita ingin mengetahui kuartil pertama (Q1), kuartil ketiga (Q3), dan persentil ke-90 dari data usia mahasiswa, kita bisa menggunakan fungsi `quantile()`.

## Code Example Hitung Quantile

```
Q1 = df['Usia'].quantile(0.25)
Q3 = df['Usia'].quantile(0.75)
persentil_90 = df['Usia'].quantile(0.90)
print("Kuartil 1:", Q1)
print("Kuartil 3:", Q3)
print("Persentil ke-90:", persentil_90)
```



# MATPLOTLIB

## Code Examples :

```
import matplotlib.pyplot as plt

# Data contoh
data = [2, 3, 1, 4, 3, 2, 4, 1, 3, 2]

# Membuat histogram
plt.hist(data)

# Menambahkan label dan judul
plt.xlabel('Nilai')
plt.ylabel('Frekuensi')
plt.title('Histogram Distribusi Nilai')

# Menampilkan grafik
plt.show()
```

Fungsi / Metode	Deskripsi	Contoh Penggunaan
<code>plt.plot()</code>	Membuat plot garis untuk data.	<code>plt.plot([1, 2, 3, 4], [10, 20, 25, 30])</code> menghasilkan plot garis antara dua array.
<code>plt.hist()</code>	Membuat histogram untuk visualisasi distribusi data.	<code>plt.hist([1, 2, 2, 3, 3, 3, 4], bins=4)</code> menghasilkan histogram dari data dengan 4 bin.
<code>plt.boxplot()</code>	Membuat diagram kotak (boxplot) untuk menggambarkan distribusi data dan outliers.	<code>plt.boxplot([1, 2, 3, 4, 5])</code> menghasilkan boxplot dari data tersebut.
<code>plt.scatter()</code>	Membuat plot sebar (scatter plot) untuk menggambarkan hubungan antar dua variabel.	<code>plt.scatter([1, 2, 3], [10, 20, 25])</code> menghasilkan scatter plot antara dua array.
<code>plt.bar()</code>	Membuat diagram batang untuk data kategorikal.	<code>plt.bar(['A', 'B', 'C'], [5, 7, 3])</code> menghasilkan diagram batang berdasarkan kategori.
<code>plt.errorbar()</code>	Membuat plot dengan garis error untuk menampilkan ketidakpastian atau variasi dalam data.	<code>plt.errorbar([1, 2, 3], [10, 20, 30], yerr=[1, 2, 1])</code> menghasilkan plot dengan error bar pada sumbu y.



# MATPLOTLIB

## Code Examples :

```
import matplotlib.pyplot as plt
```

```
# Data contoh
x = [1, 2, 3, 4, 5]
y = [2, 4, 6, 8, 10]
```

```
# Membuat plot garis
plt.plot(x, y)
```

```
# Menambahkan label dan judul
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Plot Garis Hubungan X-Y')
```

```
# Menampilkan grafik
plt.show()
```

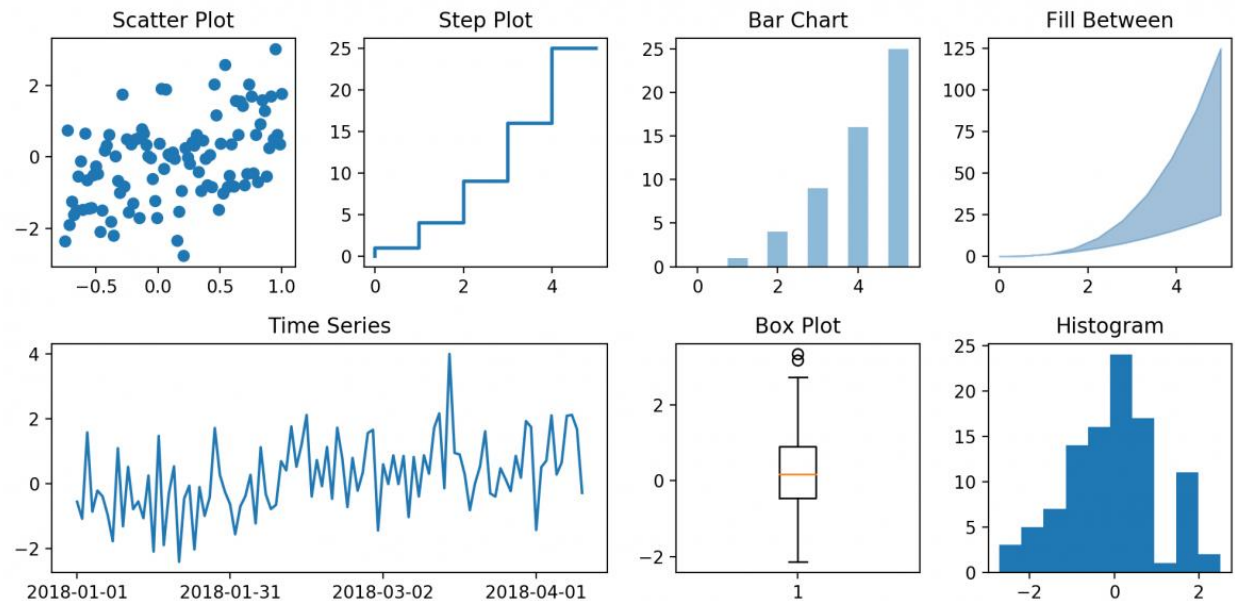
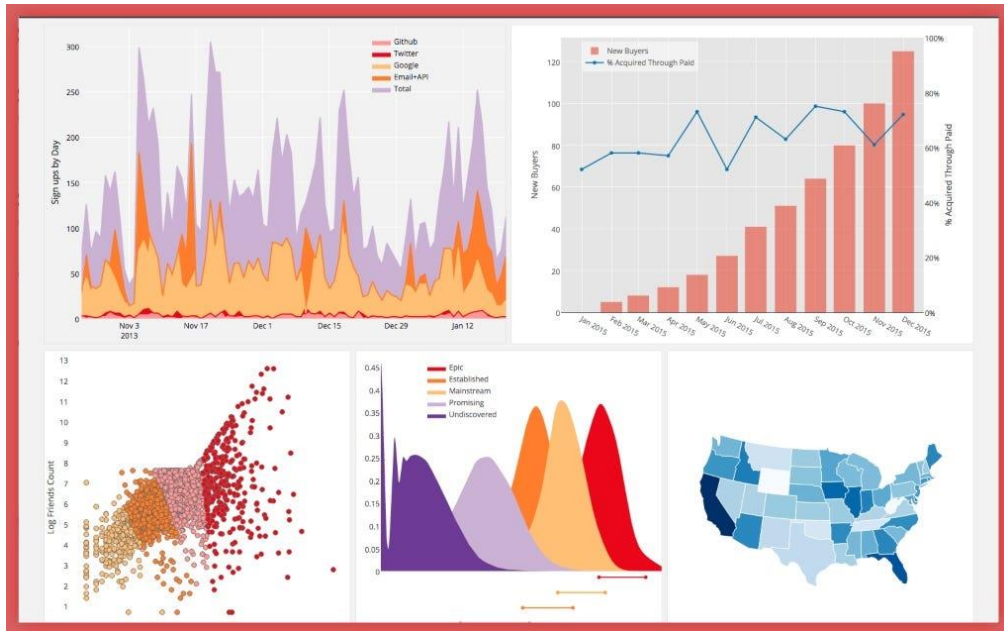
<code>plt.pie()</code>	Membuat diagram lingkaran (pie chart) untuk menggambarkan proporsi data.	<code>plt.pie([30, 30, 40], labels=['A', 'B', 'C'], autopct='%1.1f%%')</code> menghasilkan pie chart dengan tiga kategori.
<code>plt.subplot()</code>	Membuat beberapa subplot dalam satu figure untuk menampilkan beberapa grafik sekaligus.	<code>plt.subplot(1, 2, 1)</code> membuat subplot pertama dari dua, yang dapat digunakan untuk plot terpisah.
<code>plt.xlim()</code> dan <code>plt.ylim()</code>	Mengatur batas sumbu x dan y dalam plot.	<code>plt.xlim(0, 10); plt.ylim(0, 100)</code> mengatur batas sumbu x antara 0 dan 10 serta sumbu y antara 0 dan 100.
<code>plt.title()</code>	Menambahkan judul pada plot.	<code>plt.title("Distribusi Data Usia")</code> menambahkan judul pada plot.
<code>plt.grid()</code>	Menambahkan grid pada plot untuk mempermudah pembacaan data.	<code>plt.grid(True)</code> menambahkan grid pada plot yang sedang dibuat.
<code>plt.fill_between()</code>	Mengisi area antara dua kurva pada plot, biasanya digunakan untuk visualisasi distribusi.	<code>plt.fill_between([1, 2, 3], [10, 20, 15], [5, 10, 5], color='blue', alpha=0.3)</code> mengisi area antara dua garis dengan warna biru.



# MATPLOTLIB

berguna dalam visualisasi data statistik dan probabilitas, membantu dalam menggambarkan distribusi data, hubungan antar variabel, dan memberikan representasi visual yang intuitif untuk analisis lebih lanjut.

<code>plt.legend()</code>	Menambahkan legenda pada plot untuk menunjukkan arti dari elemen yang ditampilkan.	<code>plt.plot([1, 2, 3], [10, 20, 25], label="Data 1"); plt.legend()</code> menambahkan legenda dengan label "Data 1".
<code>plt.barh()</code>	Membuat diagram batang horizontal untuk visualisasi data kategorikal.	<code>plt.barh(['A', 'B', 'C'], [5, 7, 3])</code> menghasilkan diagram batang horizontal berdasarkan kategori.



## PLOT EXAMPLES



# TUGAS

1. Diberikan data distribusi normal dengan mean 50 dan standar deviasi 5. Buatlah program Python untuk menghasilkan **100 nilai acak** dari distribusi normal ini dan hitung **mean**, **median**, dan **standar deviasi** dari data yang dihasilkan.
2. Diberikan DataFrame df dengan dua kolom x dan y yang memiliki data sebagai berikut: x: [5, 10, 15, 20, 25] y: [10, 20, 20, 30, 30]. Hitunglah mean, median, dan modus untuk setiap kolom x dan y menggunakan Pandas !
3. Visualisasikan soal nomor 1 dan 2 menggunakan Matplotlib !

