



Hands-On

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

Pertemuan 6

Pertemuan 6 (enam) pada Microcredential Associate Data Scientist 2021 menyampaikan materi mengenai Mengumpulkan Data, Menelaah Data dengan metode Visualisasi

Latihan (1)

Sebelum menelaah data dengan metode visualisasi, kita perlu memanggil modul visualisasi (seaborn & matplotlib) terlebih dahulu.

```
In [8]: # memanggil modul Pandas and Seaborn
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
plt.style.use('fivethirtyeight')

import warnings
warnings.filterwarnings('ignore')
```

```
In [9]: bunga = 'Iris.csv'
df = pd.read_csv(bunga)
```

```
In [10]: # menampilkan 5 baris data
df.head()
```

```
Out[10]:
```

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|----|---------------|--------------|---------------|--------------|-------------|
| 0 | 1 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 2 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 3 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

Latihan (2)

Karena kita tidak membutuhkan kolom "Id" dalam melakukan visualisasi kita dapat menghapus kolom "Id" menggunakan fungsi `.drop()`

```
In [11]: # menghapus kolom "id"
df = df.drop('Id',1)
```

Latihan (3)

Lakukan pengecekan nilai yang hilang (*missing value*) pada dataset. Dengan function `info()`

```
In [12]: # memeriksa missing values pada dataset
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   SepalLengthCm    150 non-null    float64
1   SepalWidthCm     150 non-null    float64
2   PetalLengthCm    150 non-null    float64
3   PetalWidthCm     150 non-null    float64
4   Species          150 non-null    object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
```

```
In [13]: df['Species'].value_counts()
```

```
Out[13]: Iris-virginica      50
Iris-versicolor      50
Iris-setosa      50
Name: Species, dtype: int64
```

```
In [14]: df.head()
```

```
Out[14]:
```

| | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---------------|--------------|---------------|--------------|-------------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

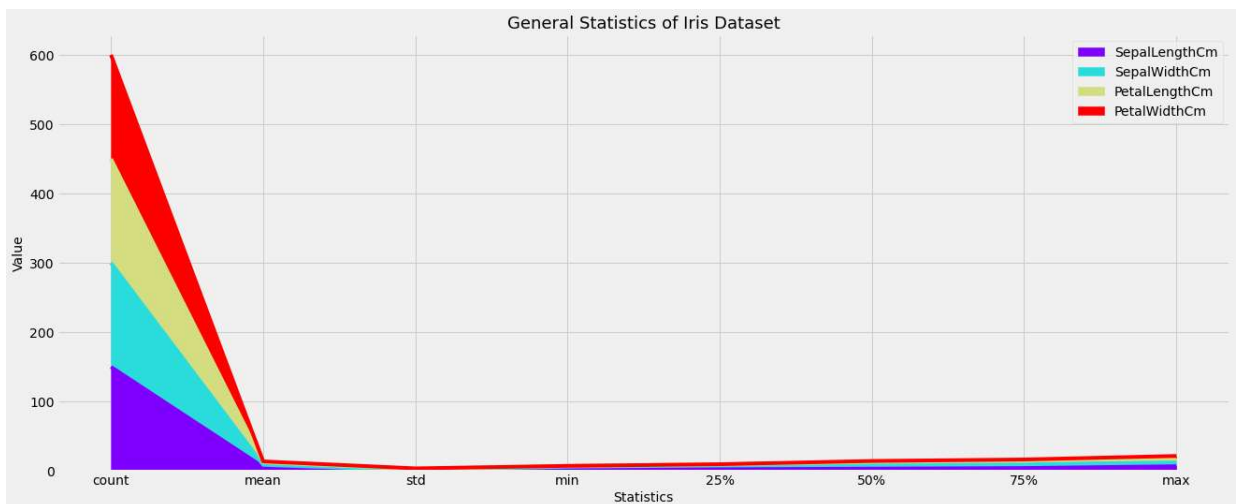
Latihan (4)

Tampilkan visualisasi dari data yang telah menggunakan fungsi describe() untuk mendapatkan informasi umum statistik tentang dataset

```
In [15]: # melakukan visualisasi dari data describe
```

```
In [16]: df.describe().plot(kind = "area", figsize = (20,8), colormap="rainbow")
plt.xlabel('Statistics',fontsize=14)
plt.ylabel('Value',fontsize=14)
plt.title("General Statistics of Iris Dataset",fontsize=18)
```

```
Out[16]: Text(0.5, 1.0, 'General Statistics of Iris Dataset')
```



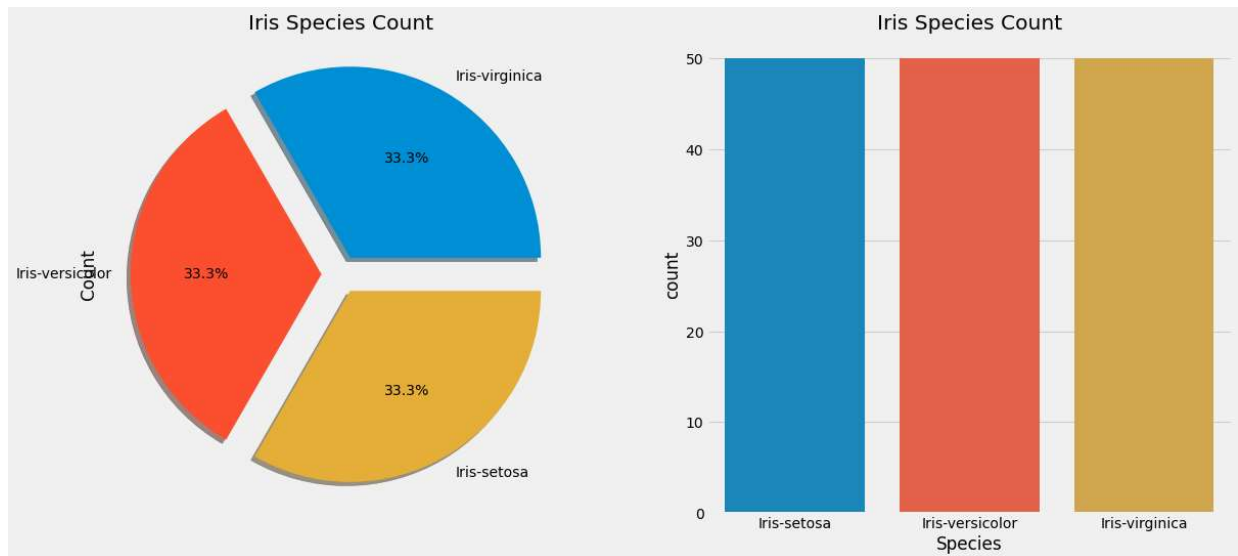
Latihan (5)

Tampilkan visualisasi bar plot dan pie chart untuk menghitung frekuensi setiap species dalam dataset iris

```
In [17]: # visualisasi bar plot dan pie chart
```

```
In [18]: f,ax=plt.subplots(1,2,figsize=(18,8))
df['Species'].value_counts().plot.pie(explode=[0.1,0.1,0.1],autopct='%1.1f%%',ax=
ax[0].set_title('Iris Species Count')
ax[0].set_ylabel('Count')
sns.countplot('Species',data=df,ax=ax[1])
ax[1].set_title('Iris Species Count')
```

Out[18]: Text(0.5, 1.0, 'Iris Species Count')



Visualisasi jointplot digunakan untuk menganalisis dua variable dan menggambarkan distribusi pada plot

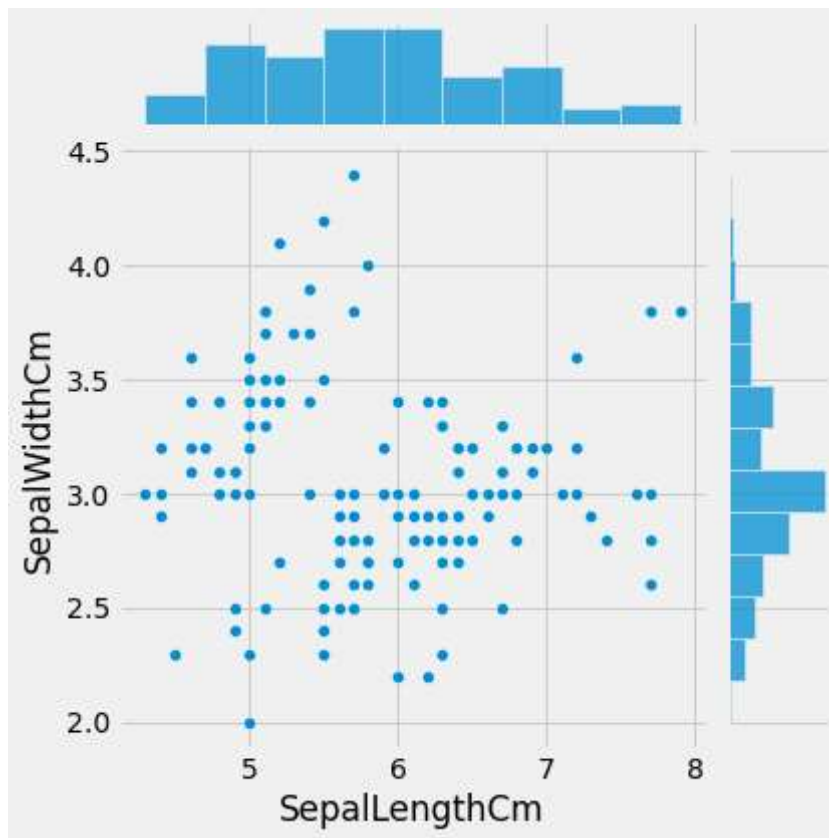
Tampilkan visualisasi jointplot menggunakan fitur 'SepalLengthCm' dan 'SepalWidthCm' dalam visualisasi jointplot.

Latihan (6)

Gunakan jenis plot residual

```
In [19]: sns.jointplot(x='SepalLengthCm',y='SepalWidthCm', data=df)
```

```
Out[19]: <seaborn.axisgrid.JointGrid at 0x1f28c72a910>
```

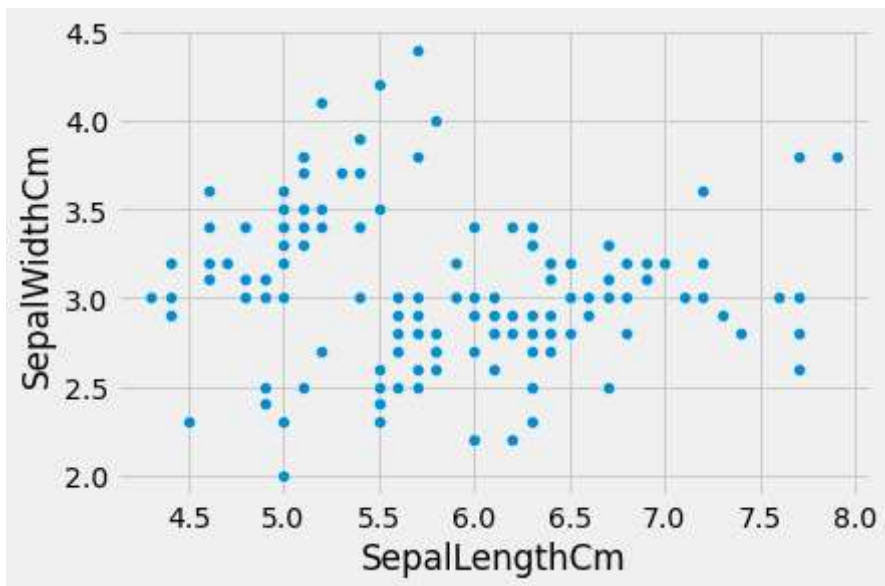


Latihan (7)

Gunakan jenis scatter plot

```
In [28]: sns.scatterplot(x='SepalLengthCm',y='SepalWidthCm',data=df)
```

```
Out[28]: <AxesSubplot:xlabel='SepalLengthCm', ylabel='SepalWidthCm'>
```

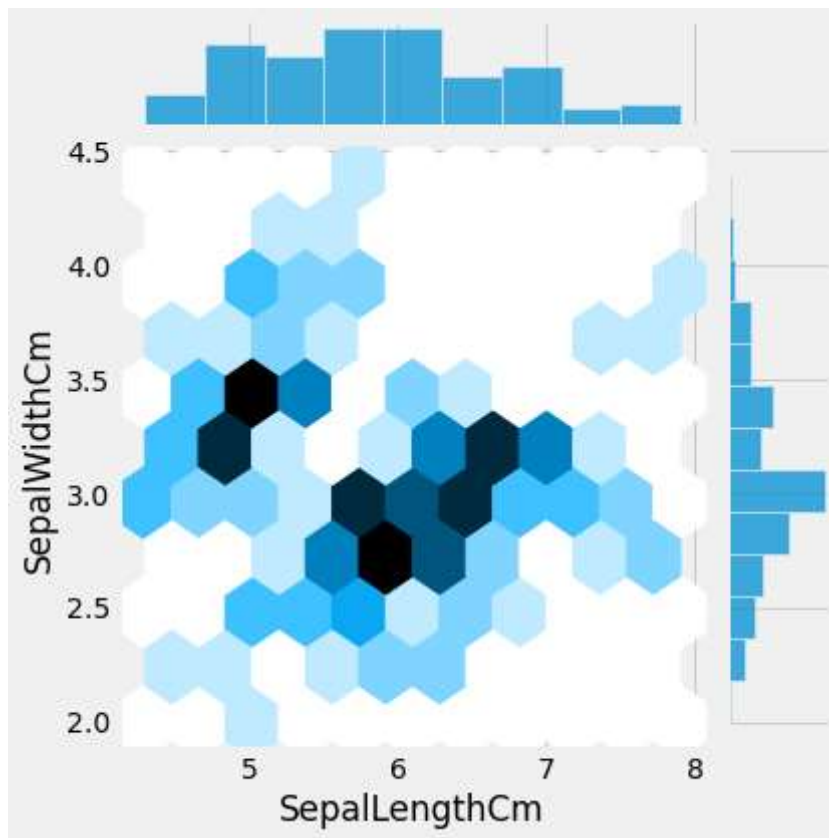


Latihan (8)

Gunakan jenis hexagons plot

```
In [21]: sns.jointplot('SepalLengthCm', 'SepalWidthCm', kind="hex", data=df)
```

```
Out[21]: <seaborn.axisgrid.JointGrid at 0x1f28c82d220>
```

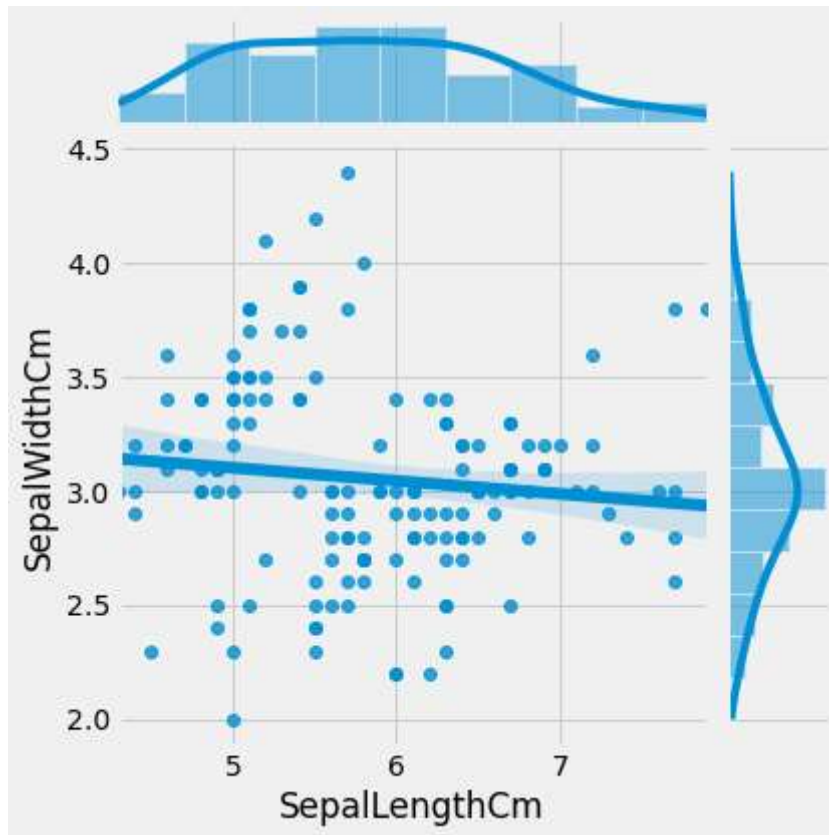


Latihan (9)

Gunakan jenis Linear regression line plot

```
In [26]: sns.jointplot('SepalLengthCm', 'SepalWidthCm', data=df, kind="reg")
```

```
Out[26]: <seaborn.axisgrid.JointGrid at 0x1f28e395070>
```

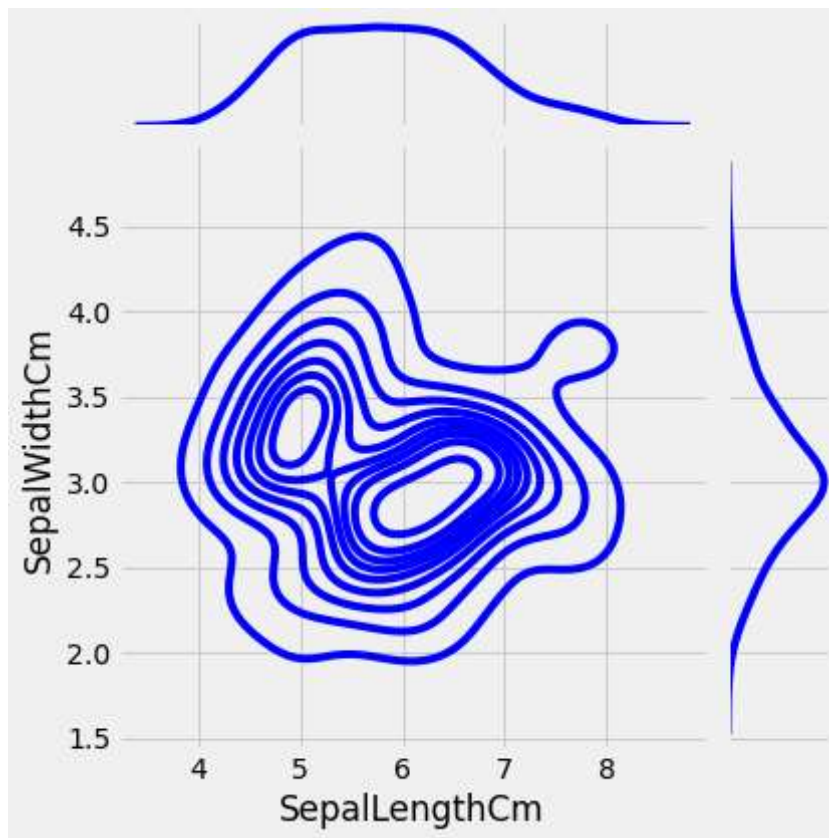


Latihan (10)

Gunakan jenis kernel density estimate plot


```
In [23]: sns.jointplot("SepalLengthCm", "SepalWidthCm", data=df, kind="kde", color='b')
```

```
Out[23]: <seaborn.axisgrid.JointGrid at 0x1f28c964b80>
```

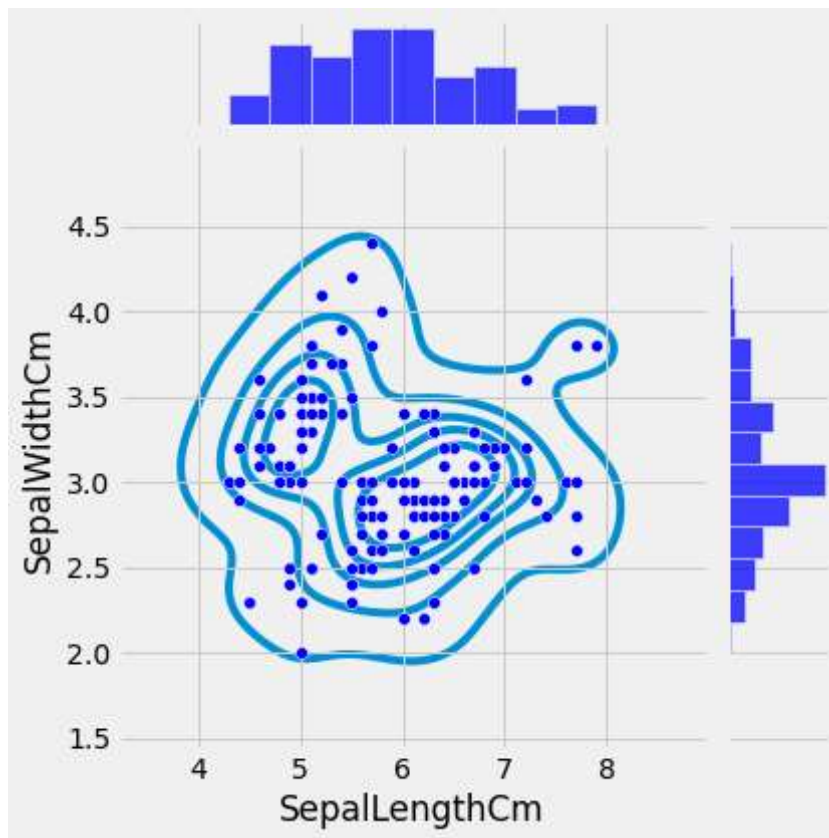


Latihan (11)

Lengkapi potongan code penggabungan plot

```
In [24]: sns.jointplot("SepalLengthCm", "SepalWidthCm", data=df, color="b").plot_joint(sns.
```

```
Out[24]: <seaborn.axisgrid.JointGrid at 0x1f28e2c3700>
```



Latihan (12)

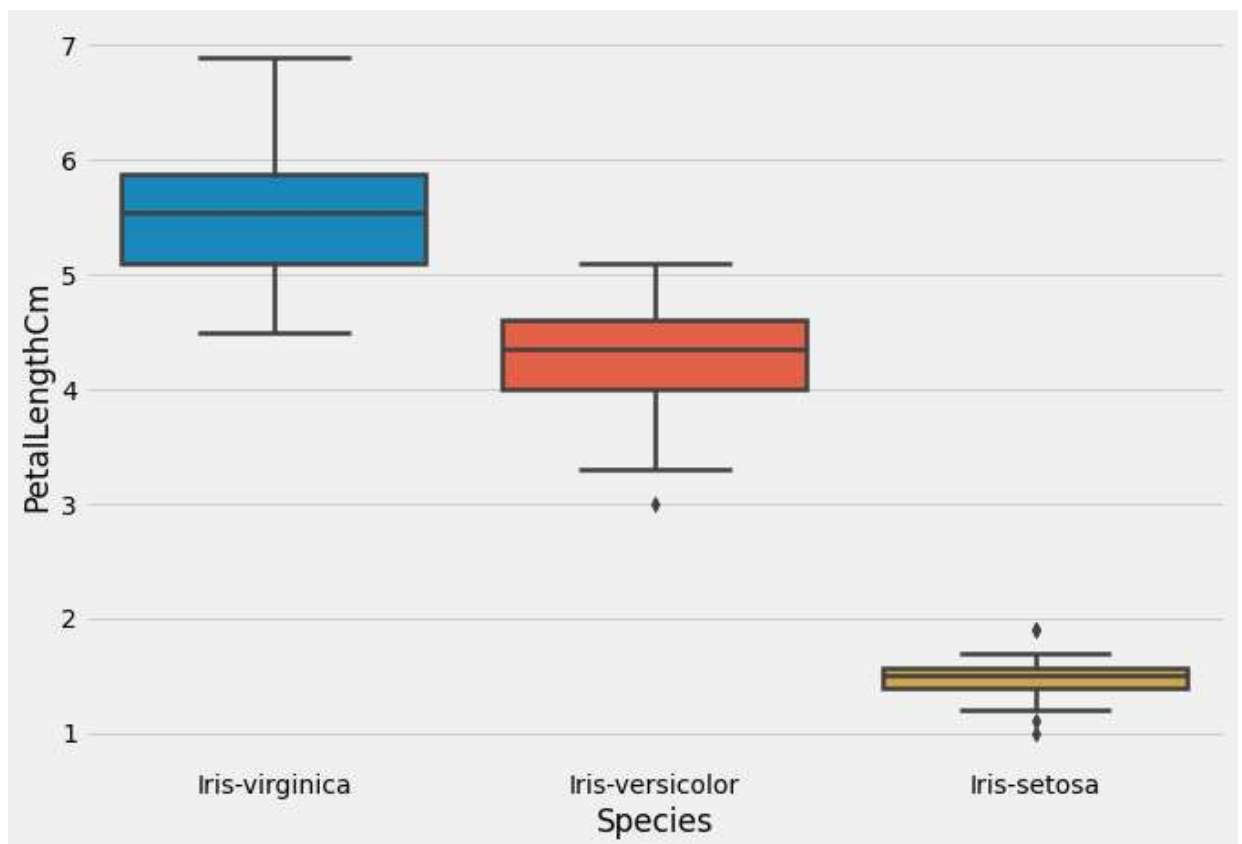
Visualisasi Boxplot untuk memberikan ringkasan statis dari fitur yang diplot.

- Garis atas mewakili nilai maksimal
- Tepi atas kotak adalah Kuartil ketiga
- Tepi tengah adalah median,
- Tepi bawah adalah nilai kuartil pertama.
- Garis paling bawah adalah nilai minimum.
- Ketinggian kotak disebut sebagai rentang Interkuartil.
- Titik-titik hitam pada plot adalah nilai outlier dalam data.

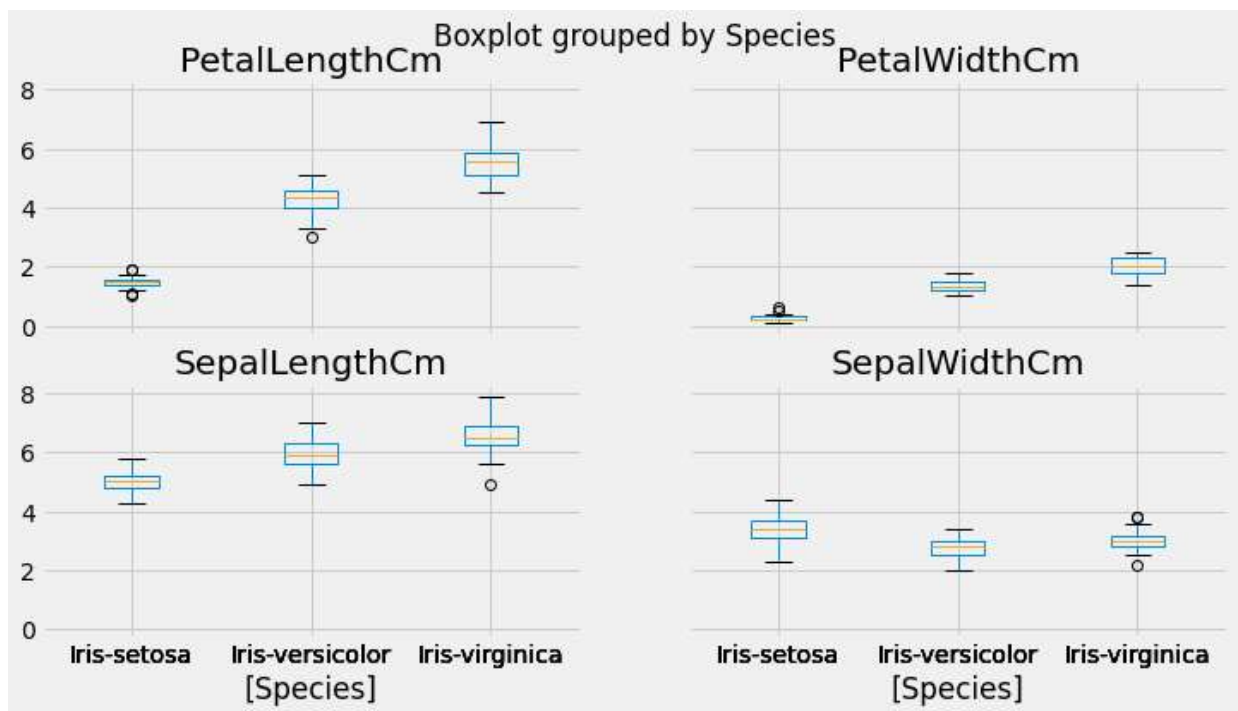
Tampilkan visualisasi boxplot menggunakan kolom "Species" dan "PetalLengthCm" dalam dataset iris

```
In [ ]: # visualisasi Boxplot
```

```
In [29]: fig=plt.gcf()
fig.set_size_inches(10,7)
fig=sns.boxplot(x='Species',y='PetalLengthCm',data=df,order=['Iris-virginica','Ir
```



```
In [32]: # visualisasi Boxplot yang di kelompokkan berdasarkan "Species"
df.boxplot(by="Species", figsize=(12, 6))
pass
```



Latihan (13)

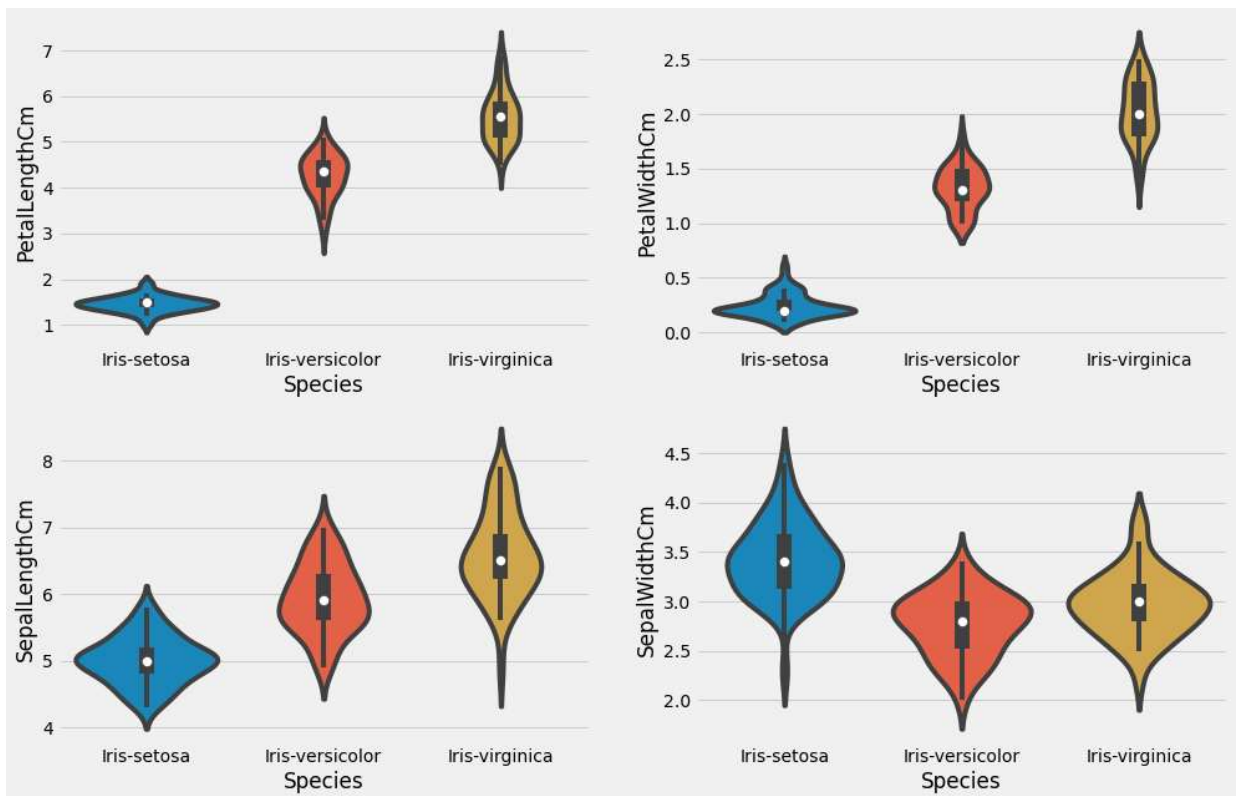
Visualisasi Violin Plot untuk memvisualisasikan sebaran data dan distribusi probabilitas.

- Bilah hitam tebal di tengah mewakili rentang interkuartil
- Garis hitam tipis yang memanjang darinya mewakili interval kepercayaan 95%
- Titik putih adalah median.

Tampilkan visualisasi Violin Plot dengan menggunakan setiap kolom yang ada untuk melihat sebaran data terhadap kolom "Species" dalam dataset iris

```
In [33]: plt.figure(figsize=(15,10))
plt.subplot(2,2,1)
sns.violinplot(x='Species',y='PetalLengthCm',data=df)
plt.subplot(2,2,2)
sns.violinplot(x='Species',y='PetalWidthCm',data=df)
plt.subplot(2,2,3)
sns.violinplot(x='Species',y='SepalLengthCm',data=df)
plt.subplot(2,2,4)
sns.violinplot(x='Species',y='SepalWidthCm',data=df)
```

```
Out[33]: <AxesSubplot:xlabel='Species', ylabel='SepalWidthCm'>
```



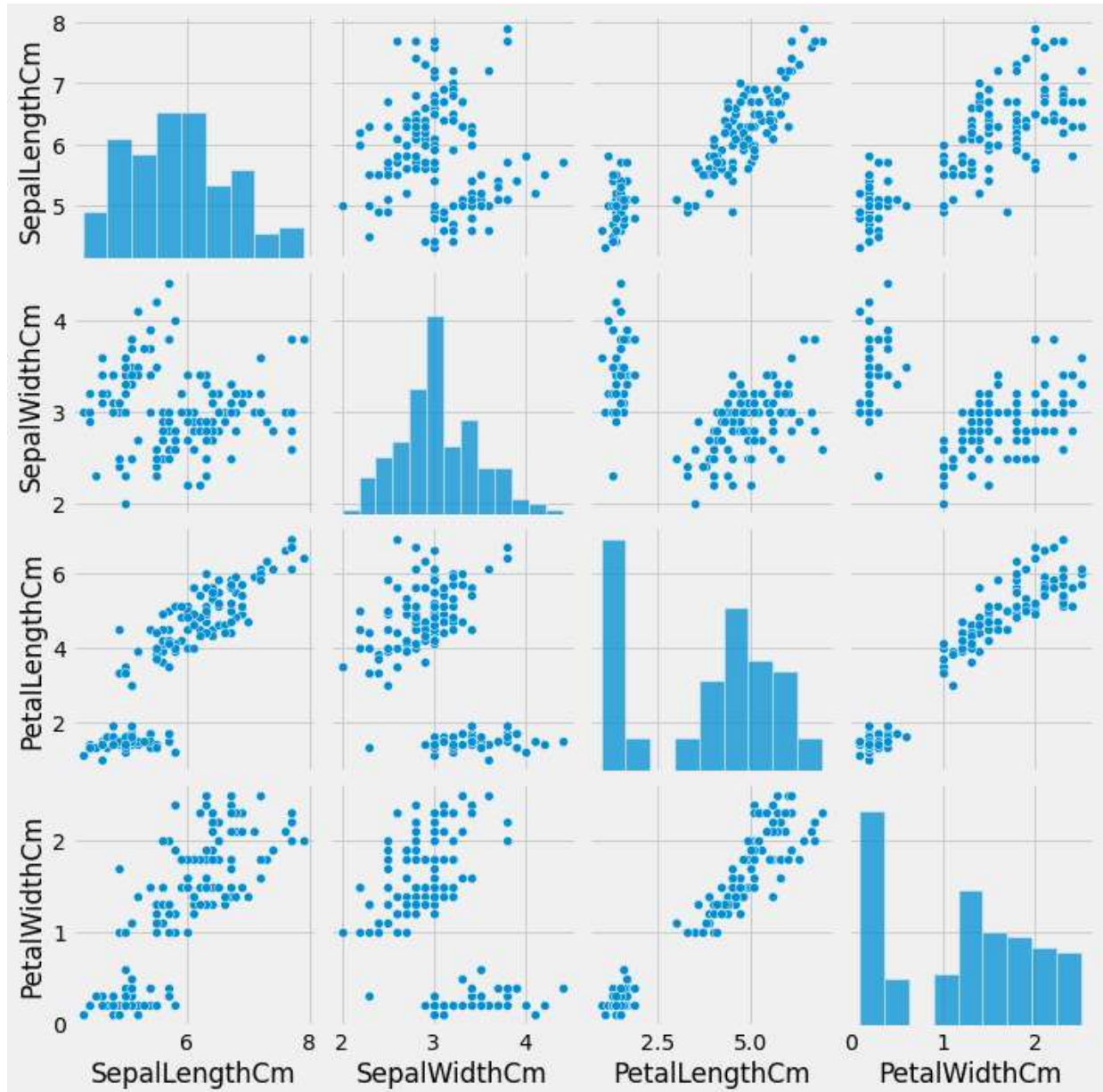
Latihan (14)

Visualisasi pairplot yang juga dikenal dengan scatterplot adalah visualisasi sebaran data yang menunjukkan keterkaitan antar kolom.

Tampilkan visualisasi pairplot dari setiap kolom yang ada untuk melihat sebaran data dalam dataset iris

```
In [34]: # visualisasi pairplot
sns.pairplot(data=df, kind='scatter')
```

```
Out[34]: <seaborn.axisgrid.PairGrid at 0x1f28fada670>
```



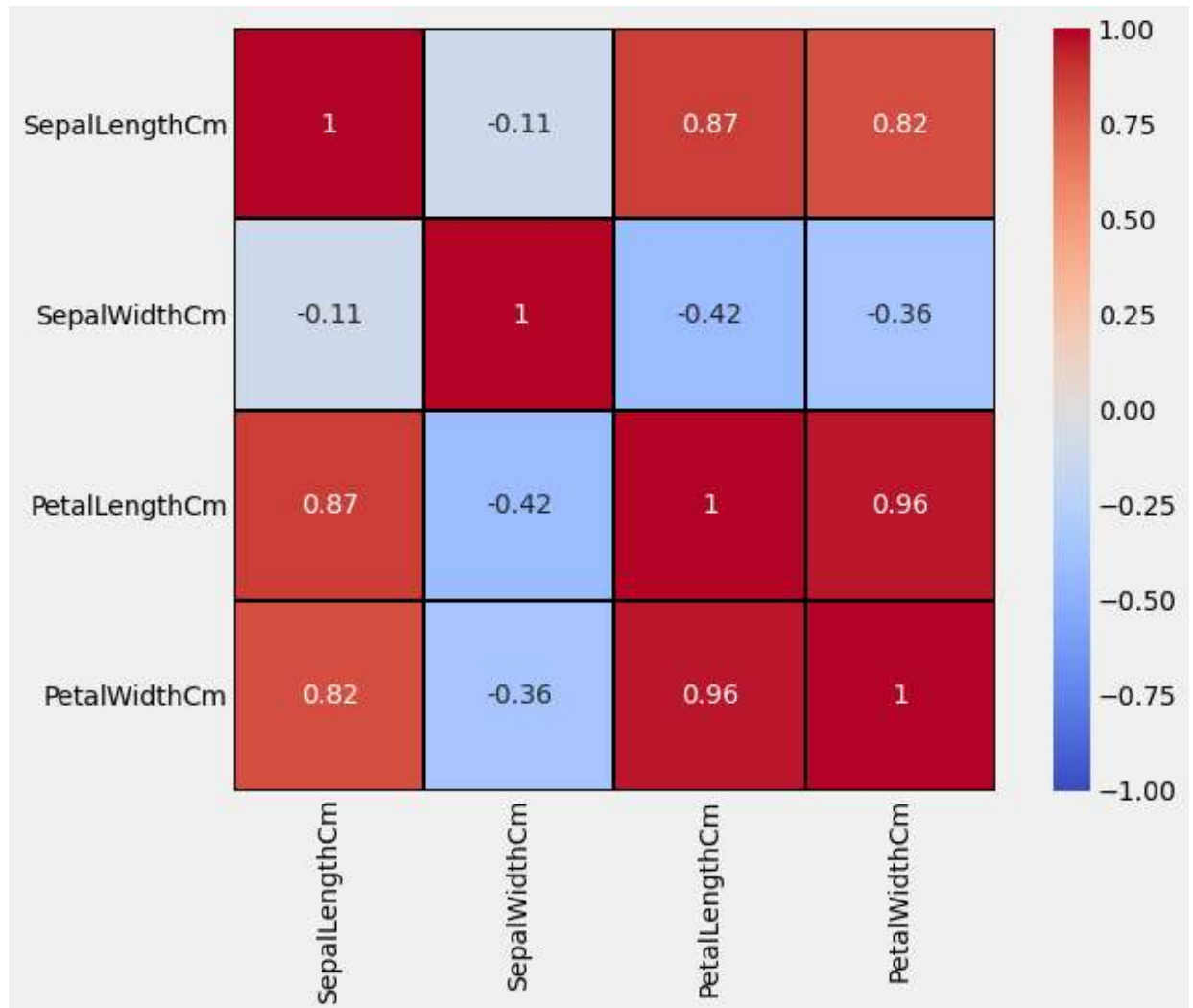
Latihan (15)

Visualisasi heatmap digunakan untuk mengetahui korelasi setiap kolom dalam dataset. Nilai positif atau negatif yang tinggi menunjukkan bahwa fitur tersebut memiliki korelasi yang tinggi. Hal ini membantu kita memilih parameter untuk machine learning.

Tampilkan visualisasi heatmap dari korelasi setiap fitur dalam dataset iris

```
In [ ]: # visualisasi heatmap
```

```
In [35]: fig=plt.gcf()  
fig.set_size_inches(10,7)  
fig=sns.heatmap(df.corr(),annot=True,cmap='coolwarm',linewidths=1,linecolor='k',s
```



Latihan (16)

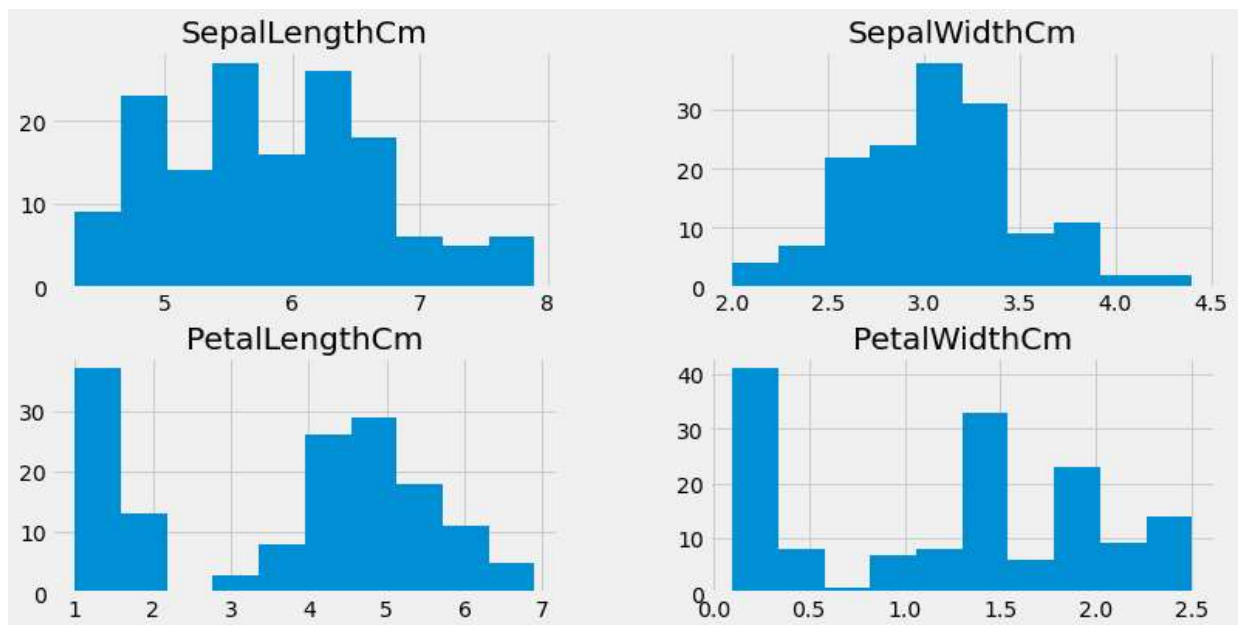
Visualisasi Distribution plot untuk membandingkan jangkauan dan distribusi untuk dataset numerik. Visualisasi Distribution plot tidak relevan untuk analisis data yang terperinci karena berkaitan dengan ringkasan distribusi data.

Tampilkan visualisasi Distribution plot setiap fitur dalam dataset iris

```
In [ ]: # visualisasi distribution plot
```



```
In [36]: df.hist(linewidth=1.2)
fig=plt.gcf()
fig.set_size_inches(12,6)
```



Latihan (17)

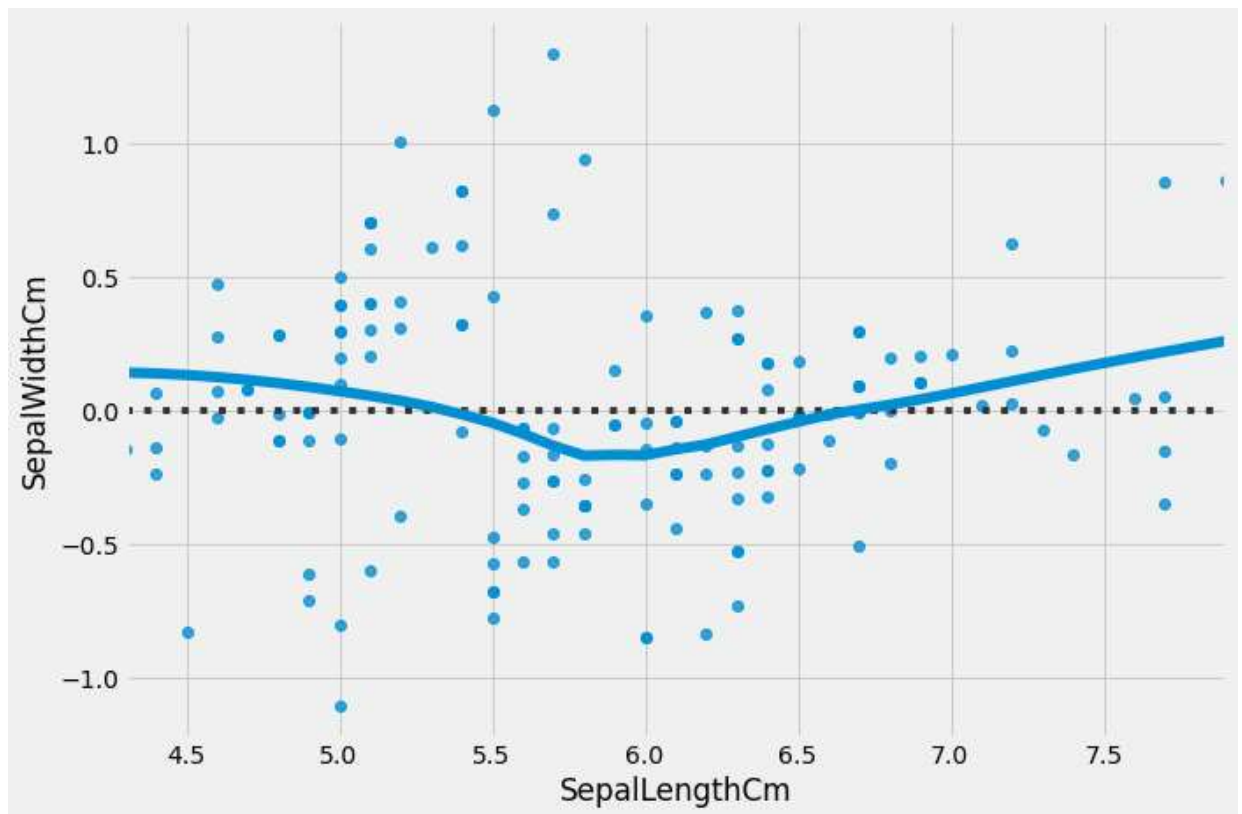
Visualisasi Residual Plot adalah visualisasi untuk memplot residu dengan nilai prediksi pada sumbu x, dan residu pada sumbu y. Jarak dari garis 0 adalah seberapa buruk prediksi untuk nilai yang di plot.

Tampilkan visualisasi Residual Plot dari fitur "SepalLengthCm" dan SepalWidthCm dalam dataset iris

```
In [ ]: # visualisasi Residual Plot
```



```
In [37]: fig=plt.gcf()
fig.set_size_inches(10,7)
fig=sns.residplot('SepalLengthCm', 'SepalWidthCm',data=df,lowess=True)
```



Latihan (18)

Visualisasi Stacked Histogram digunakan untuk menunjukkan bagaimana fitur yang lebih besar dibagi menjadi fitur yang lebih kecil dan menunjukkan hubungan masing-masing fitur terhadap jumlah total

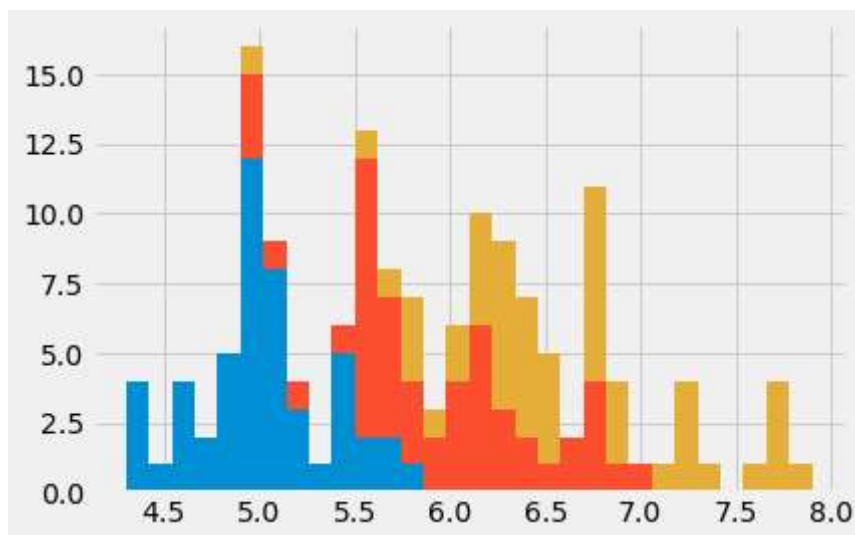
Tampilkan visualisasi Stacked Histogram dari fitur "Species" dengan mengubah tipe datanya menjadi *category* (astype) dalam dataset iris

```
In [191]: # visualisasi stacked histogram
```

```
In [39]: df['Species'] = df['Species'].astype('category')
#df.head()
list1=list()
mylabels=list()
for gen in df.Species.cat.categories:
    list1.append(df[df.Species==gen].SepalLengthCm)
    mylabels.append(gen)

h=plt.hist(list1,bins=30,stacked=True,rwidth=1,label=mylabels)
plt.hist
```

```
Out[39]: <function matplotlib.pyplot.hist(x, bins=None, range=None, density=False, weights=None, cumulative=False, bottom=None, histtype='bar', align='mid', orientation='vertical', rwidth=None, log=False, color=None, label=None, stacked=False, *, data=None, **kwargs)>
```



Latihan (19)

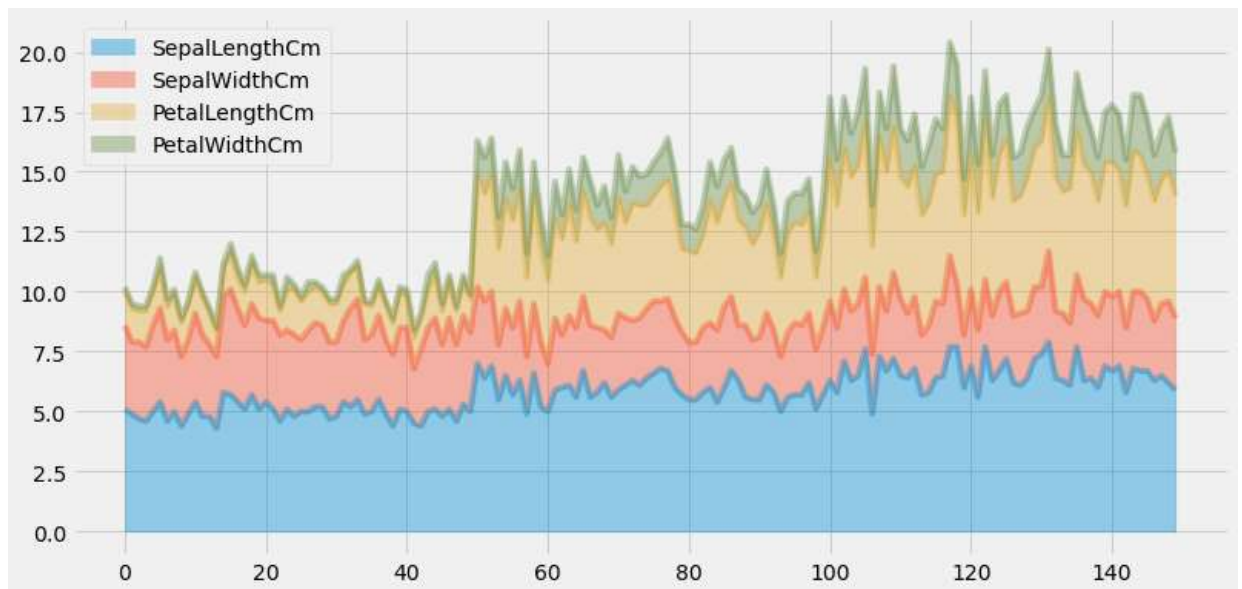
Visualisasi Area Plot memberi kita representasi visual dari Berbagai dimensi dalam dataset.

Tampilkan visualisasi Area Plot dari fitur

'SepalLengthCm','SepalWidthCm','PetalLengthCm','PetalWidthCm' dalam dataset iris

```
In [ ]: # visualisasi Area Plot
```

```
In [43]: df.plot.area(y=['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm'], ax=)
```



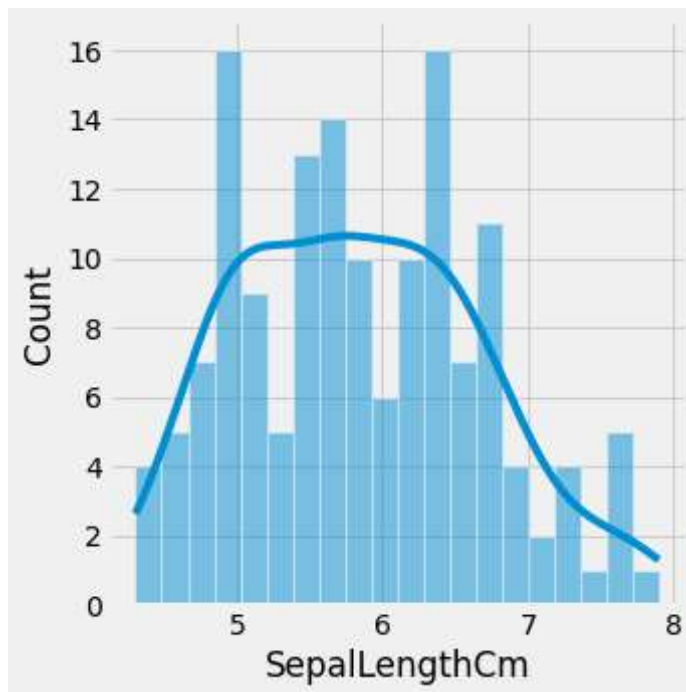
Latihan (20)

Visualisasi distplot membantu untuk melihat distribusi variabel tunggal. Kde menunjukkan kepadatan distribusi

Tampilkan visualisasi distplot dari fitur 'SepalLengthCm' dengan menggunakan *Kde* (kind) untuk menunjukkan kepadatan distribusi dalam dataset iris

```
In [ ]: # visualiasi distplot
```

```
In [44]: sns.displot(df['SepalLengthCm'],kde=True,bins=20);
```



```
In [ ]: #Friska Andalusia_Universitas Telkom
```