

Tugas Mandiri Pertemuan 5_Andrean Yonathan_Institut Teknologi Sepuluh Nopember

Latihan (1)

Silahkan Download sebuah dataset menggunakan API Kaggle

In [1]:	<pre># Install modul kaggle secara inline (di dalam notebook) !pip install kaggle Requirement already satisfied: kaggle in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (1.5.12) Requirement already satisfied: python-slugify in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (from kaggle) (5.0.2) Requirement already satisfied: urllib3 in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (from kaggle) (2.8.1) Requirement already satisfied: certifi in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (from kaggle) (1.26.6) Requirement already satisfied: tqdm in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (from kaggle) (2021.5.30) Requirement already satisfied: requests in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (from kaggle) (1.16.0) Requirement already satisfied: idna<4,>=2.5 in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (from requests->kaggle) (3.2) Requirement already satisfied: colorama in c:\users\user\appdata\local\programs\python\python39\lib\site-packages (from tqdm->kaggle) (0.4.4) WARNING: You are using pip version 21.2.4; however, version 21.3.1 is available. You should consider upgrading via the 'C:\Users\USER\AppData\Local\Programs\Python\Python39\python.exe -m pip install --upgrade pip' command.</pre>
In [2]:	<pre># Mencari dataset yang tersedia di kaggle --> pilih data provider dari UCIML !kaggle datasets list --m Iris ref title size astUpdated downloadCount voteCount usabilityRating ----- uciml/Iris Iris Species 4KB 016-09-27 07:38:05 225795 2679 0.7941176 arshid/Iris-Flower-Dataset Iris Flower Dataset 1010B 018-03-12 15:18:06 40549 370 0.8235294 vikrishnan/Iris-dataset Iris Dataset 999B 017-08-03 16:00:44 2921 26 0.7647059 theroch/Ireland-historical-news Irish Times - Waxy-Many News 52MB 021-09-25 10:52:48 2982 157 1.0 chuckyin/Iris-datasets Iris datasets 1KB 017-03-10 09:35:43 1764 13 0.7352941 Irisan/Iris-dataset-json-version Iris Dataset (JSON Version) 1KB 018-04-06 20:21:31 5629 43 0.75 paruipandey/palmer-archipelago-antarctica-penguin-data Palmer Archipelago (Antarctica) penguin data 17KB 020-06-09 10:14:54 10010 114 0.9705882 conorrot/Irish-weather-hourly-data Irish Weather (hourly data) 61KB 020-06-29 20:15:18 1861 40 0.8235294 saaurabh0007/Iris.csv Iris.csv 1KB 021-11-09 07:34:35 17111 57 0.4117647 jillanisofttech/Iris-dataset-uci Iris dataset uci 1KB 021-11-06 15:11:47 36 12 1.0 fleanend/birds-songa-numeric-dataset Birds' Songa Numeric Dataset 2KB 019-04-01 09:09:46 706 25 0.9411765 kamrankausar/Iris-data Iris data 15MB 017-11-30 10:26:01 1117 9 0.64705884 jeffheaton/Iris-computer-vision Iris Computer Vision 5MB 020-11-24 21:32:29 304 0.875 stuyven/Iris-dataset Iris dataset 1KB 017-11-04 14:10:12 791 8 0.29411766 Irisan14193/Iris-species Iris Species 2KB 020-07-02 06:09:09 57 13 0.5625 oigabelitskaya/flower-color-images Flower Color Images 50MB 020-08-01 22:48:07 8362 161 0.75 naureemohammed/mwu-iris-dataset MWU iris dataset 30MB 020-07-25 18:38:33 644 19 0.5625 rutujavaidya/Iris-dataset Iris dataset 1KB 021-07-25 17:37:14 35 6 0.4117647 shantannuss/Iris-flower-dataset IRIS flower dataset 1KB 020-01-18 19:43:18 195 3 0.9411765 ashishs0ni/Iris-dataset Iris_dataset 1KB 019-08-05 14:26:19 598 7 0.64705884</pre>
In [3]:	<pre># Download dan ekstrak dataset, secara default akan berada dalam satu direktori dengan notebook ini !kaggle datasets download uciml/Iris --unzip Downloading Iris.zip to C:\Users\USER\OneDrive - Universitas Diponegoro\Andrean\Microcredential - Associate Data Science\Pertemuan 5 0% 0.00/3.60k [00:00<?, ?B/s] 100% ##### 3.60k/3.60k [00:00<00:00, 921kB/s]</pre>

Latihan (2)

Lakukan import Library Pandas dan Library Numpy

In [4]:	<pre>#Latihan(2) #Import Library Pandas import pandas as pd #Import Library Numpy import numpy as np</pre>
---------	---

Latihan (3)

Panggil file (load dataset) dengan format .csv untuk dataset mengenai bunga Iris yang sudah peserta unduh dari Kaggle, dan akan disimpan di dalam dataframe df. Lalu tampilkan 5 baris awal dataset dengan function head()

```
#Hitung ukuran (jumlah baris dan kolom) dari dataset
df_size = df.shape
df_size
```

Out[8]: (150, 6)

Latihan (7)

Berapakah jumlah baris, dan jumlah kolom pada dataset? (silakan diisi pada cell di bawah ini)

```
In [9]: #Latihan (7)

#Jumlah Baris pada dataset adalah = <isikan jawaban di sini>

#Jumlah kolom pada dataset adalah = <isikan jawaban di sini>

print('Jumlah baris pada dataset adalah = %d \nJumlah kolom pada dataset adalah = %d' %df_size)
```

```
Jumlah baris pada dataset adalah = 150
Jumlah kolom pada dataset adalah = 6
```

Latihan (8)

Tampilkan data yang hanya berisi kolom "Id" dan kolom "Species" dalam bentuk dataframe.

```
In [10]: #Latihan (8)

#Tampilkan data untuk kolom "Id" dan kolom "Species" dalam bentuk dataframe

df_id_species = df[['Id', 'Species']]
df_id_species.head()
```

```
Out[10]:
```

	Id	Species
0	1	Iris-setosa
1	2	Iris-setosa
2	3	Iris-setosa

Latihan (4)

Tampilkan tipe data dari kolom yang ada pada dataset

In [6]:	<pre>#Latihan(4) #Tampilkan tipe data dari kolom yang ada pada dataset dfType = df.dtypes dfType</pre>														
Out[6]:	<table><tbody><tr><td>Id</td><td>int64</td></tr><tr><td>SepalLengthCm</td><td>float64</td></tr><tr><td>SepalWidthCm</td><td>float64</td></tr><tr><td>PetalLengthCm</td><td>float64</td></tr><tr><td>PetalWidthCm</td><td>float64</td></tr><tr><td>Species</td><td>object</td></tr><tr><td>dtype:</td><td>object</td></tr></tbody></table>	Id	int64	SepalLengthCm	float64	SepalWidthCm	float64	PetalLengthCm	float64	PetalWidthCm	float64	Species	object	dtype:	object
Id	int64														
SepalLengthCm	float64														
SepalWidthCm	float64														
PetalLengthCm	float64														
PetalWidthCm	float64														
Species	object														
dtype:	object														

Latihan (5)

Apakah tipe Data dari kolom berikut ini: (silakan diisi pada cell di bawah ini)

In [7]:	<pre>#Latihan (5) #Tipe Data dari kolom yang ada di dataset #Kolom "Id" memiliki tipe data = <isikan jawaban di sini> #Kolom "SepalLengthCm" memiliki tipe data = <isikan jawaban di sini> #Kolom "Species" memiliki tipe data = <isikan jawaban di sini> dfType_list = dfType.index for i in range(6): df_list = dfType_list[i] df_type = dfType[i] print(f'Kolom {df_list} memiliki tipe data = {df_type}')</pre>
	<p>Kolom Id memiliki tipe data = int64 Kolom SepalLengthCm memiliki tipe data = float64 Kolom SepalWidthCm memiliki tipe data = float64 Kolom PetalLengthCm memiliki tipe data = float64 Kolom PetalWidthCm memiliki tipe data = float64 Kolom Species memiliki tipe data = object</p>

Latihan (6)

Hitunglah ukuran (jumlah baris dan kolom) dari dataset. Dengan menggunakan method function

In [8]:	<pre>#Latihan (6) #Hitung Ukuran (jumlah baris dan kolom) dari dataset df_size = df.shape df_size</pre>
Out[8]:	(150, 6)

Latihan (7)

Berapakah jumlah baris, dan jumlah kolom pada dataset? (silakan diisi pada cell di bawah ini)

In [9]:	<pre>#Latihan (7) #Jumlah Baris pada dataset adalah = <isikan jawaban di sini> #Jumlah kolom pada dataset adalah = <isikan jawaban di sini> print(f'Jumlah baris pada dataset adalah = {d}\nJumlah kolom pada dataset adalah = {d} \ndf_size=')</pre>
	<p>Jumlah baris pada dataset adalah = 150 Jumlah kolom pada dataset adalah = 6</p>

Latihan (8)

Tampilkan data yang hanya berisi kolom "Id" dan kolom "Species" dalam bentuk dataframe.

df_corr()

Out[15]:

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
Id	1.000000	0.716676	-0.397729	0.882747	0.899759
SepalLengthCm	0.716676	1.000000	-0.109369	0.871754	0.817954
SepalWidthCm	-0.397729	-0.109369	1.000000	-0.420516	-0.356544
PetalLengthCm	0.882747	0.871754	-0.420516	1.000000	0.962757
PetalWidthCm	0.899759	0.817954	-0.356544	0.962757	1.000000

Latihan (14)

Berdasarkan pada perhitungan korelasi di Latihan (13), apakah yang dapat Bapak/Ibu simpulkan sementara? Silakan tuliskan simpulan sementara Bapak/Ibu pada cell di bawah ini.

Latihan (9)

Tampilkan data dengan dataframe, dan data yang ditampilkan adalah data pada baris dengan indeks 0 (nol) sampai dengan indeks 9 (sembilan)

Hitung korelasi untuk kolom berikut ini: `PetalLengthCm`, `PetalWidthCm`

```

In [16]: #Latihan (15)
          #Hitung korelasi untuk kolom PetalLengthCm, PetalWidthCm

          df[['PetalLengthCm', 'PetalWidthCm']].corr()

```

```

Out[16]:
          PetalLengthCm  PetalWidthCm
PetalLengthCm  1.000000    0.962757
PetalWidthCm   0.962757    1.000000

```

Latihan (16)

Method "describe" secara otomatis melakukan komputasi statistik untuk semua continous variable. Secara default "describe" melakukan ignore terhadap variabel bertipe objek.

Latihan (10)

Tampilkan data hanya kolom "Id" dan kolom "Species" dengan dataframe, dan yang ditampilkan adalah data pada baris dengan indeks 11 (sebelas) sampai dengan indeks 15 (limabelas)

out[17]:

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	150.000000	150.000000	150.000000	150.000000	150.000000
mean	75.500000	5.843333	3.054000	3.758667	1.198667
std	43.453568	0.828066	0.433594	1.764420	0.763161
min	1.000000	4.300000	2.000000	1.000000	0.100000
25%	38.250000	5.100000	2.800000	1.600000	0.300000
50%	75.500000	5.800000	3.000000	4.350000	1.300000
75%	112.750000	6.400000	3.300000	5.100000	1.800000
max	150.000000	7.900000	4.400000	6.900000	2.500000

Latihan (17)

Latihan (11)

Pada DataFrame dapat menampilkan beberapa baris pertama/terakhir dari dataset yang di load. Gunakan Method head() dan tail().

Latihan: Tampilkan data pada 8 (delapan) baris pertama dari dataset, dengan dataframe.

| In [13]: | ``` #Latihan (11) #Tampilkan data pada 8 (delapan) baris pertama dari dataset, dengan dataframe df.head(8) ``` |
| Out[13]: | | | 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 | | 5 | 6 | 5.4 | 3.9 | 1.7 | 0.4 | Iris-setosa | | 6 | 7 | 4.6 | 3.4 | 1.4 | 0.3 | Iris-setosa | | 7 | 8 | 5.0 | 3.4 | 1.5 | 0.2 | Iris-setosa | |

Latihan (12)

Tampilkan data pada 3 (tiga) baris terakhir dari dataset, dengan dataframe.

| In [14]: | ``` #Latihan (12) #Tampilkan data pada 3 (tiga) baris terakhir dari dataset, dengan dataframe df.tail(3) ``` |
| Out[14]: | | | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species | |-----|-----|---------------|--------------|---------------|--------------|----------------| | 147 | 148 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica | | 148 | 149 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica | | 149 | 150 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica | |

Latihan (13)

Hitung korelasi dari dataset. Dengan menggunakan method function

Perhatikan bahwa describe pada dataset yang sudah di load akan berisi type data
 df.describe(include = 'all')

out[19]:

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
count	150.000000	150.000000	150.000000	150.000000	150.000000	150
unique	NaN	NaN	NaN	NaN	NaN	3
top	NaN	NaN	NaN	NaN	NaN	Iris-virginica
freq	NaN	NaN	NaN	NaN	NaN	50
mean	75.500000	5.843333	3.054000	3.758667	1.198667	NaN
std	43.445368	0.828066	0.433594	1.764420	0.763161	NaN
min	1.000000	4.300000	2.000000	1.000000	0.100000	NaN
25%	36.250000	5.100000	2.800000	1.600000	0.300000	NaN
50%	75.500000	5.800000	3.000000	4.350000	1.300000	NaN
75%	113.750000	6.400000	3.200000	6.100000	1.900000	NaN

Latihan (14)

Berdasarkan pada perhitungan korelasi di Latihan (13), apakah yang dapat Bapak/Ibu simpulkan sementara? Silakan tuliskan simpulan sementara Bapak/Ibu pada cell di bawah ini.

Simpulan Sementara Hasil Korelasi di latihan (13)

variabel SepalLengthCm, PetalLengthCm, dan PetalWidthCm ketiganya memiliki korelasi positif kuat antar variabel

Latihan (15)

Hitung korelasi untuk kolom berikut ini: PetalLengthCm, PetalWidthCm

Latihan (20)

Hitung nilai mean dari dataset untuk kolom `Petal.LengthCm`.

```
In [21]: #Latihan (20)
#Hitung nilai Mean untuk kolom Petal.LengthCm

petalLengthmean = df['Petal.LengthCm'].mean()
print(f'Nilai mean untuk kolom Petal.LengthCm adalah {petalLengthmean}')
```

Nilai mean untuk kolom `Petal.LengthCm` adalah 3.7586666666666693

Latihan (21)

Carilah nilai minimal dari dataset untuk kolom `Sepal.WidthCm`.

```
In [22]: #Latihan (21)
```

Latihan (16)

Method "describe" secara otomatis melakukan komputasi statistik untuk semua continous variable. Secara default "describe" melakukan ignore terhadap variabel bertipe objek.

Komputasi statistik yang dilakukan terdiri dari: count, mean, std, min, max, 25%, 75%, max.

Latihan: Gunakan method describe pada dataset yang sudah di load untuk semua continous variabel. (Dataset Iris.csv)

```
In [23]: #Latihan (22)
#Hitung nilai mean dari dataset untuk SepalLengthCm per Species dengan metode groupby

dfSpecies = df.groupby('Species')['SepalLengthCm'].mean()
dfSpecies
```

```
Out[23]: Species
Iris-setosa      5.006
Iris-versicolor  5.936
Iris-virginica   6.588
Name: SepalLengthCm, dtype: float64
```

Latihan (23)

Hitunglah frekuensi pada kolom 'Species' dengan menggunakan metode value_counts().

```
In [24]: #Latihan (23)
#Hitung frekuensi pada kolom 'Species' dengan menggunakan metode value_counts()
```

Latihan (17)

Gunakan method describe pada dataset yang sudah di load untuk data bertipe objek. (Dataset Iris.csv)

In [18]:	<pre>#Latihan (17) #Gunakan method describe pada dataset yang sudah di load untuk data bertipe objek df.describe(include = 'object')</pre>										
Out[18]:	<table><thead><tr><th></th><th>Species</th></tr></thead><tbody><tr><td>count</td><td>150</td></tr><tr><td>unique</td><td>3</td></tr><tr><td>top</td><td>Iris-virginica</td></tr><tr><td>freq</td><td>50</td></tr></tbody></table>		Species	count	150	unique	3	top	Iris-virginica	freq	50
	Species										
count	150										
unique	3										
top	Iris-virginica										
freq	50										

Latihan 18

Gunakan method describe pada dataset yang sudah di load untuk semua type data (continous variabel dan type objek).

1.4	12
5.1	8
4.5	8
1.6	7
1.3	7
5.6	6
4.0	5
4.9	5
4.7	5
5.0	4
1.7	4
4.8	4
4.4	4
4.2	4
4.1	3
5.7	3
5.5	3
6.1	3
3.9	3
4.6	3
5.8	3
5.2	2

Latihan (19)

Hitunglah nilai mean dari dataset.

In [20]:	<pre>#Latihan (19) #Hitung nilai Mean dari dataset print('Mean dari dataset adalah:\n') df.mean()</pre>												
	<p>Mean dari dataset adalah:</p>												
Out[20]:	<table><tbody><tr><td>Id</td><td>75.500000</td></tr><tr><td>SepalLengthCm</td><td>5.843333</td></tr><tr><td>SepalWidthCm</td><td>3.054000</td></tr><tr><td>PetalLengthCm</td><td>3.758667</td></tr><tr><td>PetalWidthCm</td><td>1.198667</td></tr><tr><td>dtype:</td><td>float64</td></tr></tbody></table>	Id	75.500000	SepalLengthCm	5.843333	SepalWidthCm	3.054000	PetalLengthCm	3.758667	PetalWidthCm	1.198667	dtype:	float64
Id	75.500000												
SepalLengthCm	5.843333												
SepalWidthCm	3.054000												
PetalLengthCm	3.758667												
PetalWidthCm	1.198667												
dtype:	float64												

Latihan (20)

Hitung nilai mean dari dataset untuk kolom PetalLengthCm.

In [21]:	<pre>#Latihan (20) #Hitung nilai Mean untuk kolom PetalLengthCm petalLengthmean = df['PetalLengthCm'].mean() print(f'Nilai mean untuk kolom PetalLengthCm adalah {petalLengthmean}')</pre>
	<p>Nilai mean untuk kolom PetalLengthCm adalah 3.758666666666693</p>

Latihan (21)

Carilah nilai minimal dari dataset untuk kolom SepalWidthCm.

In [22]:	<pre>#Latihan (21) #Cari nilai minimal untuk kolom SepalWidthCm minSepalWidth = df['SepalWidthCm'].min() print(f'Nilai minimal untuk kolom SepalWidthCm (minSepalWidth) ')</pre>
	<p>Nilai minimal untuk kolom SepalWidthCm 2.0</p>

Latihan (22)

Hitunglah nilai mean dari dataset untuk kolom SepalLengthCm per Species dengan menggunakan metode groupby.

In [23]:	<pre>#Latihan (22) #Hitung nilai mean dari dataset untuk SepalLengthCm per Species dengan metode groupby dfSpecies = df.groupby('Species')['SepalLengthCm'].mean() dfSpecies</pre>										
Out[23]:	<table><tbody><tr><td>Species</td><td>5.006</td></tr><tr><td>Iris-setosa</td><td>5.006</td></tr><tr><td>Iris-versicolor</td><td>5.936</td></tr><tr><td>Iris-virginica</td><td>6.588</td></tr><tr><td>Name: SepalLengthCm, dtype: float64</td><td></td></tr></tbody></table>	Species	5.006	Iris-setosa	5.006	Iris-versicolor	5.936	Iris-virginica	6.588	Name: SepalLengthCm, dtype: float64	
Species	5.006										
Iris-setosa	5.006										
Iris-versicolor	5.936										
Iris-virginica	6.588										
Name: SepalLengthCm, dtype: float64											

Latihan (23)

Hitunglah frekuensi pada kolom 'Species' dengan menggunakan metode value_counts().

In [24]:	<pre>#Hitung frekuensi pada kolom 'Species' dengan menggunakan metode value_counts() df['Species'].value_counts()</pre>						
Out[24]:	<table><tbody><tr><td>Iris-virginica</td><td>50</td></tr><tr><td>Iris-versicolor</td><td>50</td></tr><tr><td>Iris-setosa</td><td>50</td></tr></tbody></table>	Iris-virginica	50	Iris-versicolor	50	Iris-setosa	50
Iris-virginica	50						
Iris-versicolor	50						
Iris-setosa	50						

Latihan (24)

Tampilkan perhitungan frekuensi pada kolom 'Species' dengan menggunakan metode value_counts() dalam bentuk dataframe.

In [27]:	<pre>#Latihan (24) #Hitung frekuensi pada kolom 'Species' dengan menggunakan metode value_counts() dalam bentuk dataframe df['Species'].value_counts().to_frame()</pre>								
Out[27]:	<table><thead><tr><th></th><th>Species</th></tr></thead><tbody><tr><td>Iris-virginica</td><td>50</td></tr><tr><td>Iris-versicolor</td><td>50</td></tr><tr><td>Iris-setosa</td><td>50</td></tr></tbody></table>		Species	Iris-virginica	50	Iris-versicolor	50	Iris-setosa	50
	Species								
Iris-virginica	50								
Iris-versicolor	50								
Iris-setosa	50								

Latihan (25)

Hitunglah frekuensi pada kolom 'PetalLengthCm' dengan menggunakan metode value_counts() dan dalam bentuk dataframe.

In [31]:	<pre>#Latihan (25) #Hitung frekuensi pada kolom 'PetalLengthCm' dengan menggunakan metode value_counts() df['PetalLengthCm'].value_counts()</pre>																																
Out[31]:	<table><tbody><tr><td>1.5</td><td>14</td></tr><tr><td>1.4</td><td>12</td></tr><tr><td>5.1</td><td>8</td></tr><tr><td>4.5</td><td>8</td></tr><tr><td>1.6</td><td>7</td></tr><tr><td>1.3</td><td>7</td></tr><tr><td>5.6</td><td>6</td></tr><tr><td>4.0</td><td>5</td></tr><tr><td>4.9</td><td>5</td></tr><tr><td>4.7</td><td>5</td></tr><tr><td>5.0</td><td>4</td></tr><tr><td>1.7</td><td>4</td></tr><tr><td>4.8</td><td>4</td></tr><tr><td>4.4</td><td>4</td></tr><tr><td>4.2</td><td>4</td></tr><tr><td>4.1</td><td>3</td></tr></tbody></table>	1.5	14	1.4	12	5.1	8	4.5	8	1.6	7	1.3	7	5.6	6	4.0	5	4.9	5	4.7	5	5.0	4	1.7	4	4.8	4	4.4	4	4.2	4	4.1	3
1.5	14																																
1.4	12																																
5.1	8																																
4.5	8																																
1.6	7																																
1.3	7																																
5.6	6																																
4.0	5																																
4.9	5																																
4.7	5																																
5.0	4																																
1.7	4																																
4.8	4																																
4.4	4																																
4.2	4																																
4.1	3																																