



Hands-On

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

Pertemuan 9

Pertemuan 9 (sembilan) pada Microcredential Associate Data Scientist 2021 menyampaikan materi mengenai Mengkonstruksi Data

Pada Tugas Mandiri Pertemuan 9

silakan Anda kerjakan Latihan 1 s/d 10. Output yang anda lihat merupakan panduan yang dapat Anda ikuti dalam penulisan code :)

Latihan (1)

Melakukan import library yang dibutuhkan

```
In [1]: 1 # import library pandas
2
3 import pandas as pd
4
5 # Import library scipy
6
7 import scipy
8
9 # Import library winsorize dari scipy
10
11 from scipy.stats.mstats import winsorize
12
13 # Import library trima dari scipy
14
15 from scipy.stats.mstats import trima
16
17 # Import library RandomSampleImputer dari feature engine imputation
18
19 from feature_engine.imputation import RandomSampleImputer
20
21 # import library StandardScaler dari sklearn
22
23 from sklearn.preprocessing import StandardScaler
```

Latihan (2)

Menghitung nilai null pada dataset :

1. Load dataset Iris_Unclean
2. Tampilkan dataset
3. Hitung jumlah nilai null pada dataset

```
In [2]: 1 # Load dataset Iris_Unclean
2
3 dfUnclean = pd.read_csv('Iris_unclean.csv')
```

```
In [3]: 1 # tampilkan dataset
2
3 dfUnclean
```

```
Out[3]:
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	NaN	3.5	1.4	0.2	Iris-setosa
1	4.9	2000.0	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
...
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

150 rows x 5 columns

```
In [4]: 1 # hitung jumlah nilai null pada dataset
2
3 dfUnclean.isnull().sum()
```

```
Out[4]: SepalLengthCm    2
SepalWidthCm          0
PetalLengthCm         0
PetalWidthCm          0
Species              0
dtype: int64
```

Latihan (3)

Melakukan handle missing value dengan Imputasi Mean:

1. Load dataset Iris_Unclean
2. Ambil 10 data teratas "SepalLengthCm", kemudian tampilkan
3. Mengganti missing value Imputasi dengan mean, kemudian masukkan ke variable
4. Tampilkan 10 data teratas "SepalLengthCm" setelah handle missing value dengan Imputasi mean()

```
In [5]: 1 # Load dataset Iris_Unclean
2
3 dfUnclean = pd.read_csv('Iris_unclean.csv')
```

```
In [6]: 1 # ambil 10 data teratas SepallengthCm, kemudian tampilkan
        2
        3 df = dfUnclean['SepallengthCm'][:10]
        4 df
```

```
Out[6]: 0      NaN
        1      4.9
        2      4.7
        3      4.6
        4      5.0
        5      5.4
        6      NaN
        7      5.0
        8      4.4
        9      4.9
        Name: SepallengthCm, dtype: float64
```

```
In [7]: 1 # mengganti missing value dengan mean(), kemudian masukkan ke variabel
        2
        3 df = df.fillna(df.mean())
```

```
In [8]: 1 # tampilkan 10 data teratas SepallengthCm setelah handle missing value dengan imputasi mean
        2
        3 df
```

```
Out[8]: 0      4.8625
        1      4.9000
        2      4.7000
        3      4.6000
        4      5.0000
        5      5.4000
        6      4.8625
        7      5.0000
        8      4.4000
        9      4.9000
        Name: SepallengthCm, dtype: float64
```

Latihan (4)

Melakukan handle missing value dengan nilai suka-suka (Arbitrary):

1. Load dataset Iris_Unclean
2. Ambil 10 data teratas "SepallengthCm", kemudian tampilkan
3. Mengganti missing value dengan imputasi nilai suka-suka (Arbitrary), kemudian masukkan ke variabel
4. Tampilkan 10 data teratas "SepallengthCm" setelah handle missing value dengan nilai suka-suka

```
In [9]: 1 # Load dataset Iris_Unclean
        2
        3 dfUnclean = pd.read_csv('Iris_unclean.csv')
```

```
In [10]: 1 # ambil 10 data teratas SepallengthCm, kemudian tampilkan
        2
        3 df = dfUnclean['SepallengthCm'][:10]
        4 df
```

```
Out[10]: 0      NaN
         1      4.9
         2      4.7
         3      4.6
         4      5.0
         5      5.4
         6      NaN
         7      5.0
         8      4.4
         9      4.9
         Name: SepallengthCm, dtype: float64
```

```
In [11]: 1 # melakukan imputasi nilai suka-suka (Arbitrary), masukkan ke dalam variabel
         2
         3 df = df.fillna(99)
```

```
In [12]: 1 # tampilkan 10 data teratas SepallengthCm setelah handle missing value dengan nilai suka-suka (arbitrary)
         2
         3 df
```

```
Out[12]: 0      99.0
         1      4.9
         2      4.7
         3      4.6
         4      5.0
         5      5.4
         6      99.0
         7      5.0
         8      4.4
         9      4.9
         Name: SepallengthCm, dtype: float64
```

Latihan (5)

Melakukan handle missing value dengan frequent category / modus:

1. Load dataset Iris_Unclean
2. Ambil 10 data teratas "SepallengthCm", kemudian tampilkan
3. Mengganti missing value dengan frequent category / modus
4. Tampilkan hasil imputasi "SepallengthCm" setelah handle dengan frequent category / modus

```
In [13]: 1 # Load dataset Iris_Unclean
         2
         3 dfUnclean = pd.read_csv('Iris_unclean.csv')
```

```
In [14]: 1 # tampilkan 10 data teratas kolom SepallengthCm
         2
         3 dfUnclean['SepallengthCm'][:10]
```

```
Out[14]: 0      NaN
         1      4.9
         2      4.7
         3      4.6
         4      5.0
         5      5.4
         6      NaN
         7      5.0
         8      4.4
         9      4.9
         Name: SepallengthCm, dtype: float64
```

```
In [15]: 1 # Import SimpleImputer dari sklearn.impute
         2
         3 from sklearn.impute import SimpleImputer
         4
         5 # Mengatasi missing value dengan frequent category / modus
         6
         7 imp = SimpleImputer(strategy='most_frequent')
```

```
In [16]: 1 # Tampilkan hasil imputasi "SepallengthCm"
         2
         3 imp.fit_transform(dfUnclean[['SepallengthCm']])
```

```
Out[16]: array([[5. ],
                [4.9],
                [4.7],
                [4.6],
                [5. ],
                [5.4],
                [5. ],
                [5. ],
                [4.4],
                [4.9]])
```

```
[5.7],  
[5.4],  
[4.8],  
[4.8],  
[4.3],  
[5.8],  
[5.7],  
[5.4],  
[5.1],  
[5.7],
```

Latihan (6)

Melakukan handle missing value dengan Imputasi Random Sample:

1. Load dataset Iris_Unclean
2. Tampilkan 10 data teratas
3. Membuat imputer random sample dengan random state = 5
4. Cocokkan imputer ke data
5. Ubah data dengan imputer masukkan ke dalam variable
6. Tampilkan hasil imputasi data "SepalLengthCm"

```
In [17]: 1 # Load dataset Iris_Unclean  
2  
3 dfUnclean = pd.read_csv('Iris_unclean.csv')
```

```
In [18]: 1 # tampilkan 10 data teratas SepalLengthCm  
2  
3 dfUnclean['SepalLengthCm'][:10]
```

```
Out[18]: 0      NaN  
1      4.9  
2      4.7  
3      4.6  
4      5.0  
5      5.4  
6      NaN  
7      5.0  
8      4.4  
9      4.9  
Name: SepalLengthCm, dtype: float64
```

```
In [19]: 1 # Membuat imputer random sample dengan random state = 5  
2  
3 imputer = RandomSampleImputer(random_state=5)  
4  
5 # Cocokkan imputer ke data  
6  
7 imputer.fit(dfUnclean)  
8  
9 # Ubah data dengan imputer masukkan ke dalam variable  
10  
11 test_t = imputer.transform(dfUnclean)
```

```
In [20]: 1 # Tampilkan data hasil imputasi data "SepalLengthCm"  
2  
3 test_t['SepalLengthCm'][:10]
```

```
Out[20]: 0      5.8  
1      4.9  
2      4.7  
3      4.6  
4      5.0  
5      5.4  
6      6.9  
7      5.0  
8      4.4  
9      4.9  
Name: SepalLengthCm, dtype: float64
```

Latihan (7)

Melakukan Winsorizing

1. Import library winsorize dari scipy
2. Load data Iris_AfterClean
3. Ambil 10 data teratas "SepalLengthCm", kemudian masukkan ke dalam variabel datan tampilkan
4. Winsorize data dengan batas nilai terendah 10% dan batas nilai tinggi 20%
5. Tampilkan hasil winsorize

```
In [21]: 1 # Import library scipy  
2  
3 import scipy
```

```
In [22]: 1 # Load data Iris_AfterClean  
2 data = pd.read_csv('Iris_AfterClean.csv')  
3  
4 # Ambil 10 data teratas "SepalLengthCm", kemudian masukkan ke dalam variabel datan tampilkan  
5 a = data['SepalLengthCm'][:10]  
6 a
```

```
Out[22]: 0      4.6  
1      5.0  
2      5.4  
3      4.9  
4      5.4  
5      4.8  
6      4.8  
7      4.3  
8      5.8  
9      5.4  
Name: SepalLengthCm, dtype: float64
```

```
In [23]: 1 # Winsorize data dengan batas nilai terendah 10% dan batas nilai tinggi 20%  
2  
3 wins = winsorize(a, limits=[0.1, 0.2])  
4  
5 # Tampilkan hasil winsorize  
6 print(wins)
```

```
[4.6 5.  5.4 4.9 5.4 4.8 4.8 4.6 5.4 5.4]
```

Latihan (8)

Melakukan Trimming

1. Import library trima dari scopy
2. Load data Iris_AfterClean
3. Ambil 10 data teratas "SepalLengthCm", kemudian masukkan ke dalam variabel datan tampilkan
4. Trimming data dengan batas nilai terendah 2 dan batas nilai tinggi 5
5. Tampilkan hasil trimming

```
In [24]: 1 # Import library trima dari scopy  
2  
3 from scipy.stats.mstats import trim
```

```
In [25]: 1 # Load data Iris_AfterClean  
2 data = pd.read_csv('Iris_AfterClean.csv')  
3  
4 # Ambil 10 data teratas "SepalLengthCm", kemudian masukkan ke dalam variabel datan tampilkan  
5 a = data['SepalLengthCm'][:10]  
6 a
```

```
Out[25]: 0    4.6
1    5.0
2    5.4
3    4.9
4    5.4
5    4.8
6    4.8
7    4.3
8    5.8
9    5.4
Name: SepalLengthCm, dtype: float64
```

```
In [26]: 1 # Trimming data dengan batas nilai terendah 2 dan batas nilai tinggi 5
2
3 trims = trima(a, limits=(2,5))
4
5 # Tampilkan hasil trimming
6
7 print(trims)

[4.6 5.0 -- 4.9 -- 4.8 4.8 4.3 -- --]
```

Latihan (9)

Melakukan Scaling: Normalisasi

1. Load data Iris_AfterClean
2. Ambil 10 data teratas SepalLengthCm dan SepalWidthCm
3. Menghitung mean data
4. Menghitung max - min pada data
5. Menerapkan transformasi ke data
6. Tampilkan hasil scaling

```
In [27]: 1 # Load data Iris_AfterClean
2
3 data = pd.read_csv('Iris_AfterClean.csv')
4
5 # Ambil 10 data teratas SepalLengthCm dan SepalWidthCm
6
7 data = data[['SepalLengthCm', 'SepalWidthCm']][:10]
8
9 data
```

```
Out[27]:
```

	SepalLengthCm	SepalWidthCm
0	4.6	3.1
1	5.0	3.6
2	5.4	3.9
3	4.9	3.1
4	5.4	3.7
5	4.8	3.4
6	4.8	3.0
7	4.3	3.0
8	5.8	4.0
9	5.4	3.9

```
In [28]: 1 # Menghitung mean
2 means = data.mean()
3
4 # menghitung max - min
5 max_min = data.max() - data.min()
6
7 # menerapkan transformasi ke data
8 train_scaled = (data - means) / max_min
```

```
In [29]: 1 # Tampilkan hasil scaling
2
3 train_scaled
```

```
Out[29]:
```

	SepalLengthCm	SepalWidthCm
0	-0.293333	-0.37
1	-0.026667	0.13
2	0.240000	0.43
3	-0.093333	-0.37
4	0.240000	0.23
5	-0.160000	-0.07
6	-0.160000	-0.47
7	-0.493333	-0.47
8	0.506667	0.53
9	0.240000	0.43

Latihan (10)

Melakukan Scaling: Standardisasi

1. Load data Iris_AfterClean
2. Ambil 10 data teratas SepalLengthCm dan SepalWidthCm
2. Import library StandardScaler dari sklearn
3. Membuat objek scaler
4. Sesuaikan scaler dengan data
5. Mengubah data
6. Tampilkan hasil scaling dengan standarisasi

```
In [30]: 1 # Load data Iris_AfterClean
2
3 data = pd.read_csv('Iris_AfterClean.csv')
4
5 # Ambil 10 data teratas SepalLengthCm dan SepalWidthCm
6
7 data = data[['SepalLengthCm', 'SepalWidthCm']][:10]
8
9 data
```

```
Out[30]:
```

	SepalLengthCm	SepalWidthCm
0	4.6	3.1
1	5.0	3.6
2	5.4	3.9
3	4.9	3.1
4	5.4	3.7
5	4.8	3.4
6	4.8	3.0
7	4.3	3.0
8	5.8	4.0
9	5.4	3.9

```
In [31]: 1 # import library StandardScaler dari sklearn
2 from sklearn.preprocessing import StandardScaler
3
4 # Buat objek scaler
5 scaler = StandardScaler()
6
7 # Sesuaikan scaler dengan data
```

```
7 # Sesuaikan scaler dengan data
8 scaler.fit(data)
9
10 # Mengubah data
11 train_scaled = scaler.transform(data)
```

```
In [32]: 1 # Tampilkan hasil
         2
         3 train_scaled
```

```
Out[32]: array([[ -1.02464215, -0.97469723],
                 [ -0.09314929,  0.34246119],
                 [  0.83834358,  1.13275625],
                 [ -0.3260225 , -0.97469723],
                 [  0.83834358,  0.60589288],
                 [ -0.55889572, -0.18440218],
                 [ -0.55889572, -1.23812892],
                 [ -1.7232618 , -1.23812892],
                 [  1.76983644,  1.39618793],
                 [  0.83834358,  1.13275625]])
```

```
In [ ]: 1
```