



Pelatihan Microcredential CERTIFICATION untuk Associate Data Scientist

1 November - 10 Desember 2021



Hands-On

Hands-On 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 [49]: # import library pandas
import pandas as pd

# Import library scipy
import scipy as sp

# Import library winsorize dari scipy
from scipy.stats.mstats import winsorize

# Import library trim dari scipy
from scipy.stats.mstats import trim

# Import library RandomSampleImputer dari feature engine imputation
from feature_engine.imputation import RandomSampleImputer

# import library StandardScaler dari sklearn
from sklearn.preprocessing import StandardScaler
```

Latihan (2)

Menghitung nilai null pada dataset :

```
In [ ]: # load dataset Iris_Unclean
df = pd.read_csv('Iris_unclean.csv')

In [ ]: # tampilkan dataset
df

In [ ]: # hitung jumlah nilai null pada dataset
df.isna().sum()
```

Latihan (3)

Melakukan handle missing value dengan Imputasi Mean:

```
In [ ]: # Load dataset Iris_Unclean
        1. Ambil 10 data teratas "SepallengthCm", kemudian tampilkan
        2. Mengganti missing value "SepallengthCm" dengan mean, kemudian masukkan ke variabel
        3. Tampilkan 10 data teratas "SepallengthCm" setelah handle missing value dengan imputasi mean()

In [ ]: # load dataset Iris_Unclean
df = pd.read_csv('Iris_unclean.csv')

In [ ]: # ambil 10 data teratas SepallengthCm, kemudian tampilkan
df = df['SepallengthCm'][0:10]
df

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

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

Latihan (4)

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

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

In [ ]: # load dataset Iris_Unclean
df = pd.read_csv('Iris_unclean.csv')

In [ ]: # ambil 10 data teratas SepallengthCm, kemudian tampilkan
df = df['SepallengthCm'][0:10]
df

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

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

Latihan (5)

Melakukan handle missing value dengan frequent category / modus:

```
In [2]: # Load dataset Iris_Unclean
        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 [3]: # tampilkan 10 data teratas kolom SepallengthCm
df = df['SepallengthCm'][0:10]
df

Out[3]: 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 [4]: df = pd.DataFrame(df)
```

```
Out[4]: SepallengthCm
```

```
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 [5]: # Import SimpleImputer dari sklearn.impute
from sklearn.impute import SimpleImputer

# Mengganti missing value dengan frequent category / modus
imp = SimpleImputer(strategy='most_frequent')
```

```
In [6]: # Tampilkan hasil imputasi "SepallengthCm"
imp.fit_transform(df)
```

```
Out[6]: array([[4.9],
       [4.9],
       [4.7],
       [4.6],
       [5.0],
       [5.4],
       [NaN],
       [5.0],
       [4.4],
       [4.9]])
```

Latihan (6)

Melakukan handle missing value dengan Imputasi Random Sample:

```
In [7]: # Load dataset Iris_Unclean
df = pd.read_csv('Iris_unclean.csv')

In [8]: # tampilkan 10 data teratas kolom SepallengthCm
df = df['SepallengthCm'][0:10]
df

Out[8]: 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]: df = pd.DataFrame(df)
```

```
Out[11]: SepallengthCm
```

```
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 [12]: # Membuat imputer random sample dengan random state = 5
imputer = RandomSampleImputer(random_state = 5)

# Cokok imputer ke data
imputer.fit(df)

# Ubah data dengan imputer masukkan ke dalam variable
test_t = imputer.transform(df)
```

```
In [13]: # Tampilkan data hasil imputasi data "SepallengthCm"
test_t
```

```
Out[13]: SepallengthCm
```

```
0      5.0
```

```
1      4.9
```

```
2      4.7
```

```
3      4.6
```

```
4      5.0
```

```
5      5.4
```

```
6      4.4
```

```
7      5.0
```

```
8      4.4
```

```
9      4.9
```

Latihan (7)

Melakukan Winsorizing

```
In [14]: # Import library winsorize dari scipy
import scipy as sp

In [25]: # Load data Iris_AfterClean
data = pd.read_csv('Iris_AfterClean.csv')

# Ambil 10 data teratas "SepallengthCm", kemudian masukkan ke dalam variabel data tampilkan
a = data['SepallengthCm'][0:10]
```

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

```
In [26]: # Winsorize data dengan batas nilai terendah 10% dan batas nilai tinggi 20%
wins = winsorize(a, limits=[0.1, 0.2])

# Tampilkan hasil winsorize
print(wins)
```

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

Latihan (8)

Melakukan Scaling: Normalisasi

```
In [27]: # Load dataset Iris_AfterClean
df = pd.read_csv('Iris_AfterClean.csv')

# Ambil 10 data teratas "SepallengthCm", kemudian masukkan ke dalam variabel data tampilkan
a = df['SepallengthCm'][0:10]
```

```
Out[27]: 0    4.6
1    5.0
2    4.9
3    4.8
4    5.4
5    4.8
6    5.0
7    4.3
8    5.8
9    5.4
Name: SepallengthCm, dtype: float64
```

```
In [28]: # Trimming data dengan batas nilai terendah 2 dan batas nilai tinggi 5
trims = trim(a, limits=(2, 5))

# Tampilkan hasil trimming
print(trim)
```

```
[4.6 5.0 -- 4.8 4.9 4.8 4.6 5.4 5.4]
```

Latihan (9)

Melakukan Scaling: Standardisasi

```
In [29]: # Load data Iris_AfterClean
data = pd.read_csv('Iris_AfterClean.csv')

# Ambil 10 data teratas "SepallengthCm", kemudian masukkan ke dalam variabel data tampilkan
a = data[['SepallengthCm']][0:10]
```

```
Out[29]: SepallengthCm
```

```
0      4.6
1      5.0
2      4.9
3      4.8
4      5.4
5      4.8
6      5.0
7      4.3
8      5.8
9      5.4
Name: SepallengthCm, dtype: float64
```

```
In [30]: # Menghitung mean
means = a.mean(axis = 0)

# menghitung max - min
max_min = a.max(axis = 0) - a.min(axis = 0)

# menerapkan transformasi ke data
train_scaled = (a - means) / max_min
```

```
In [31]: # Tampilkan hasil scaling
train_scaled
```

```
Out[31]: SepallengthCm
```

```
0      0.293333
1      0.506667
2      0.47
3      0.46
4      0.54
5      0.48
6      0.58
7      0.43
8      0.58
9      0.54
Name: SepallengthCm, dtype: float64
```

```
In [32]: # Import library winsorize dari scipy
import scipy as sp
```

```
In [33]: # Load data Iris_AfterClean
data = pd.read_csv('Iris_AfterClean.csv')

# Ambil 10 data teratas "SepallengthCm", kemudian masukkan ke dalam variabel data tampilkan
a = data['SepallengthCm'][0:10]
```

```
Out[33]: SepallengthCm
```

```
0      4.6
1      5.0
2      4.9
3      4.8
4      5.4
5      4.8
6      5.0
7      4.3
8      5.8
9      5.4
Name: SepallengthCm, dtype: float64
```

```
In [34]: # Winsorize data dengan batas nilai terendah 10% dan batas nilai tinggi 20%
wins = winsorize(a, limits=[0.1, 0.2])

# Tampilkan hasil winsorize
print(wins)
```

```
[4.6 5.0 -- 4.8 4.9 4.8 4.6 5.4 5.4]
```

Latihan (10)

Melakukan Scaling: Standardisasi

```
In [35]: # Load data Iris_AfterClean
data = pd.read_csv('Iris_AfterClean.csv')

# Ambil 10 data teratas "SepallengthCm", kemudian masukkan ke dalam variabel data tampilkan
a = data[['SepallengthCm']][0:10]
```

```
Out[35]: SepallengthCm
```

```
0      4.6
1      5.0
2      4.9
3      4.8
4      5.4
5      4.8
6      5.0
7      4.3
8      5.8
9      5.4
Name: SepallengthCm, dtype: float64
```

```
In [36]: # Import library StandardScaler dari sklearn
from sklearn.preprocessing import StandardScaler
```

```
In [37]: # Buat objek scaler
scaler = StandardScaler()

#sesuaikan scaler dengan data
scaler.fit(a)

#tambahkan data kereta
train_scaled = scaler.transform(a)

# Tampilkan hasil
train_scaled
```

```
Out[37]: SepallengthCm
```

```
0      0.293333
1      0.506667
2      0.47
3      0.46
4      0.54
5      0.48
6      0.58
7      0.43
8      0.58
9      0.54
Name: SepallengthCm, dtype: float64
```

```
In [38]: # Import library winsorize dari scipy
import scipy as sp
```

```
In [39]: # Load data Iris_AfterClean
data = pd.read_csv('Iris_AfterClean.csv')

# Ambil 10 data teratas "SepallengthCm", kemudian masukkan ke dalam variabel data tampilkan
a = data['SepallengthCm'][0:10]
```

```
Out[39]: SepallengthCm
```

```
0      4.6
1      5.0
2      4.9
3      4.8
4      5.4
5      4.8
6      5.0
7      4.3
8      5.8
9      5.4
Name: SepallengthCm, dtype: float64
```

```
In [40]: # Menghitung mean
means = a.mean(axis = 0)
```

```
In [41]: # menghitung max - min
max_min = a.max(axis = 0) - a.min(axis = 0)

# menerapkan transformasi ke data
train_scaled = (a - means) / max_min
```

```
In [42]: # Tampilkan hasil scaling
train_scaled
```

```
Out[42]: SepallengthCm
```

```
0      0.293333
1      0.506667
2      0.47
3      0.46
4      0.54
5      0.48
6      0.58
7      0.43
8      0.58
9      0.54
Name: SepallengthCm, dtype: float64
```

```
In [43]: # Import library winsorize dari scipy
import scipy as sp
```

```
In [44]: # Load data Iris_AfterClean
data = pd.read_csv('Iris_AfterClean.csv')

# Ambil 10 data teratas "SepallengthCm", kemudian masukkan ke dalam variabel data tampilkan
a = data[['SepallengthCm']][0:10]
```

```
Out[44]: SepallengthCm
```

```
0      4.6
1      5.0
2      4.9
3      4.8
4      5.4
5      4.8
6      5.0
7      4.3
8      5.8
9      5.4
Name: SepallengthCm, dtype: float64
```

```
In [45]: # Import library winsorize dari scipy
```