

# preprocessing

October 5, 2024

## TUGAS PERTEMUAN 3 DATA MINING PREPROCESSING

---

RAHMADINI CAHYA DEMORA

A11.2022.14464

A11.4509

DATA MINING

Junta Zeniarja, M.Kom

Import Library yang Digunakan

```
[26]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

Import Dataset

```
[27]: dataset = pd.read_csv("Data.csv")
```

```
[28]: x = dataset.iloc[:, :-1].values
y = dataset.iloc[:, -1].values
```

```
[29]: print(x)
```

```
[['France' 44.0 72000.0]
 ['Spain' 27.0 48000.0]
 ['Germany' 30.0 54000.0]
 ['Spain' 38.0 61000.0]
 ['Germany' 40.0 nan]
 ['France' 35.0 58000.0]
 ['Spain' nan 52000.0]
 ['France' 48.0 79000.0]
 ['Germany' 50.0 83000.0]
 ['France' 37.0 67000.0]]
```

```
[30]: print(y)
```

```
['No' 'Yes' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes' 'No' 'Yes']
```

Menghilangkan Missing Value (nan)

```
[32]: from sklearn.impute import SimpleImputer
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
imputer.fit(x[:, 1:3])
x[:, 1:3] = imputer.transform(x[:, 1:3])
```

```
[33]: print(x)
```

```
[['France' 44.0 72000.0]
 ['Spain' 27.0 48000.0]
 ['Germany' 30.0 54000.0]
 ['Spain' 38.0 61000.0]
 ['Germany' 40.0 63777.777777777778]
 ['France' 35.0 58000.0]
 ['Spain' 38.777777777777778 52000.0]
 ['France' 48.0 79000.0]
 ['Germany' 50.0 83000.0]
 ['France' 37.0 67000.0]]
```

Encoding data kategori (Atribut)

```
[34]: from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder
ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [0])],
    ↳remainder='passthrough')
x = np.array(ct.fit_transform(x))
```

```
[35]: print(x)
```

```
[[1.0 0.0 0.0 44.0 72000.0]
 [0.0 0.0 1.0 27.0 48000.0]
 [0.0 1.0 0.0 30.0 54000.0]
 [0.0 0.0 1.0 38.0 61000.0]
 [0.0 1.0 0.0 40.0 63777.777777777778]
 [1.0 0.0 0.0 35.0 58000.0]
 [0.0 0.0 1.0 38.777777777777778 52000.0]
 [1.0 0.0 0.0 48.0 79000.0]
 [0.0 1.0 0.0 50.0 83000.0]
 [1.0 0.0 0.0 37.0 67000.0]]
```

Encoding data kategori (Class/Label)

```
[36]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
y = le.fit_transform(y)
```

```
[37]: print(y)
```

```
[0 1 0 0 1 1 0 1 0 1]
```

Membagi dataset ke dalam training set dan test set

```
[38]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2,
↳ random_state = 1)
```

```
[39]: print(x_train)
```

```
[[0.0 0.0 1.0 38.77777777777778 52000.0]
 [0.0 1.0 0.0 40.0 63777.77777777778]
 [1.0 0.0 0.0 44.0 72000.0]
 [0.0 0.0 1.0 38.0 61000.0]
 [0.0 0.0 1.0 27.0 48000.0]
 [1.0 0.0 0.0 48.0 79000.0]
 [0.0 1.0 0.0 50.0 83000.0]
 [1.0 0.0 0.0 35.0 58000.0]]
```

```
[40]: print(x_test)
```

```
[[0.0 1.0 0.0 30.0 54000.0]
 [1.0 0.0 0.0 37.0 67000.0]]
```

```
[41]: print(y_train)
```

```
[0 1 0 0 1 1 0 1]
```

```
[42]: print(y_test)
```

```
[0 1]
```

Feature Scaling

```
[43]: from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
x_train[:, 3:] = sc.fit_transform(x_train[:, 3:])
x_test[:, 3:] = sc.transform(x_test[:, 3:])
```

```
[44]: print(x_train)
```

```
[[0.0 0.0 1.0 -0.19159184384578545 -1.0781259408412425]
 [0.0 1.0 0.0 -0.014117293757057777 -0.07013167641635372]
 [1.0 0.0 0.0 0.566708506533324 0.633562432710455]
 [0.0 0.0 1.0 -0.30453019390224867 -0.30786617274297867]
 [0.0 0.0 1.0 -1.9018011447007988 -1.420463615551582]
 [1.0 0.0 0.0 1.1475343068237058 1.232653363453549]]
```

```
[0.0 1.0 0.0 1.4379472069688968 1.5749910381638885]  
[1.0 0.0 0.0 -0.7401495441200351 -0.5646194287757332]]
```

```
[45]: print(x_test)
```

```
[[0.0 1.0 0.0 -1.4661817944830124 -0.9069571034860727]  
 [1.0 0.0 0.0 -0.44973664397484414 0.2056403393225306]]
```