

# Lab1\_Ex1

September 29, 2019

## 1 Laboratorium 1

### 1.1 Ćwiczenie 1 - preprocessing

#### 1.1.1 Skrypt:

```
[137]: import numpy as np
import matplotlib as plt
import pandas as pd
from IPython.display import HTML, display
```

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

```
[139]: display(dataset)
```

	Country	Age	Salary	Purchased
0	France	44.0	72000.0	No
1	Spain	27.0	48000.0	Yes
2	Germany	30.0	54000.0	No
3	Spain	38.0	61000.0	No
4	Germany	40.0	NaN	Yes
5	France	35.0	58000.0	Yes
6	Spain	NaN	52000.0	No
7	France	48.0	79000.0	Yes
8	Germany	50.0	83000.0	No
9	France	37.0	67000.0	Yes

```
[140]: X = dataset.iloc[:, :-1].values
y = dataset.iloc[:, 3].values
```

```
[141]: 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]]
```

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

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

*Ponieważ użycie klasy Imputer wyświetla informację o deprecjacji, użyta jest sugerowana klasa SimpleImputer*

```
[143]: from sklearn.impute import SimpleImputer
```

```
[144]: imputer = SimpleImputer(missing_values=np.nan, strategy="mean")
imputer = imputer.fit(X[:, 1:3])
X[:, 1:3] = imputer.transform(X[:, 1:3])
```

```
[145]: 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 63777.77777777778]
 ['France' 35.0 58000.0]
 ['Spain' 38.77777777777778 52000.0]
 ['France' 48.0 79000.0]
 ['Germany' 50.0 83000.0]
 ['France' 37.0 67000.0]]
```

```
[146]: from sklearn.preprocessing import LabelEncoder
```

```
[147]: labelencoder_X = LabelEncoder()
X[:, 0] = labelencoder_X.fit_transform(X[:, 0])
```

```
[148]: print(X)
```

```
[[0 44.0 72000.0]
 [2 27.0 48000.0]
 [1 30.0 54000.0]
 [2 38.0 61000.0]
 [1 40.0 63777.77777777778]
 [0 35.0 58000.0]
 [2 38.77777777777778 52000.0]
 [0 48.0 79000.0]
 [1 50.0 83000.0]
 [0 37.0 67000.0]]
```

```
[149]: from sklearn.preprocessing import OneHotEncoder
```

```
[ ]: onehotencoder = OneHotEncoder(categorical_features=[0])  
X = onehotencoder.fit_transform(X).toarray()
```

*Tu również wyświetla się informacja o deprecjacji, jednak zmiana jak na razie mi nie wyszła*

```
[151]: print(X)
```

```
[[1.00000000e+00 0.00000000e+00 0.00000000e+00 4.40000000e+01  
 7.20000000e+04]  
[0.00000000e+00 0.00000000e+00 1.00000000e+00 2.70000000e+01  
 4.80000000e+04]  
[0.00000000e+00 1.00000000e+00 0.00000000e+00 3.00000000e+01  
 5.40000000e+04]  
[0.00000000e+00 0.00000000e+00 1.00000000e+00 3.80000000e+01  
 6.10000000e+04]  
[0.00000000e+00 1.00000000e+00 0.00000000e+00 4.00000000e+01  
 6.37777778e+04]  
[1.00000000e+00 0.00000000e+00 0.00000000e+00 3.50000000e+01  
 5.80000000e+04]  
[0.00000000e+00 0.00000000e+00 1.00000000e+00 3.87777778e+01  
 5.20000000e+04]  
[1.00000000e+00 0.00000000e+00 0.00000000e+00 4.80000000e+01  
 7.90000000e+04]  
[0.00000000e+00 1.00000000e+00 0.00000000e+00 5.00000000e+01  
 8.30000000e+04]  
[1.00000000e+00 0.00000000e+00 0.00000000e+00 3.70000000e+01  
 6.70000000e+04]]
```

```
[152]: labelencoder_y = LabelEncoder()  
y = labelencoder_y.fit_transform(y)
```

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

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

```
[154]: from sklearn.model_selection import train_test_split
```

```
[155]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,  
↳random_state=0)
```

```
[156]: print("X_train: ", X_train.shape)  
print("X_test: ", X_test.shape)  
print("y_train: ", y_train.shape)  
print("y_test: ", y_test.shape)
```

```
X_train: (8, 5)
```

```
X_test: (2, 5)
y_train: (8,)
y_test: (2,)
```

```
[157]: from sklearn.preprocessing import StandardScaler
```

```
[158]: sc_X = StandardScaler()
X_train = sc_X.fit_transform(X_train)
X_test = sc_X.transform(X_test)
```

```
[159]: print(X_train)
```

```
[[-1.          2.64575131 -0.77459667  0.26306757  0.12381479]
 [ 1.          -0.37796447 -0.77459667 -0.25350148  0.46175632]
 [-1.          -0.37796447  1.29099445 -1.97539832 -1.53093341]
 [-1.          -0.37796447  1.29099445  0.05261351 -1.11141978]
 [ 1.          -0.37796447 -0.77459667  1.64058505  1.7202972 ]
 [-1.          -0.37796447  1.29099445 -0.0813118  -0.16751412]
 [ 1.          -0.37796447 -0.77459667  0.95182631  0.98614835]
 [ 1.          -0.37796447 -0.77459667 -0.59788085 -0.48214934]]
```

```
[160]: print(X_test)
```

```
[[-1.          2.64575131 -0.77459667 -1.45882927 -0.90166297]
 [-1.          2.64575131 -0.77459667  1.98496442  2.13981082]]
```