

# Przemek Wolczacki, czerwiec 2021

## Projekt ze Sztucznej inteligencji: Grupowanie i klasyfikacja danych

Podstawowe dane dotyczące komputera, na którym wykonano zadania:

Laptop ASUS TUF GAMING A15,  
CPU: Intel(R) Core(TM) i7-1065G7,  
CPU: AMD Ryzen z serii 4000,  
RAM: 16 GB DDR4,  
 dysk: SSD 512 GB,  
System: Windows 10, 64-bit.

Projekt wykonany w Python Notebook Jupyter.

In [1]:

```
import numpy as np
import pandas as pd
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")

import matplotlib.pyplot as plt
%matplotlib inline

from IPython.display import Markdown, display
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error as mse
from sklearn.metrics import accuracy_score
from sklearn.inspection import permutation_importance
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
from sklearn.cluster import DBSCAN
```

## Klasyfikacja

Zbiór danych użyty do klasyfikacji: *titanic.csv*  
<https://www.kaggle.com/c/titanic/data?select=test.csv> (1309 osób i 14-ście informacji o każdej z nich)

Krótki opis zbioru:

Zbiór zawiera informacje na temat pasażerów legendarnego Titanic.

Opis najważniejszych zmiennych:

**pclass** - typ klasy pasażerskiej (1-najlepsza),  
**survived** - czy przeżył (1-tak, 0-nie),  
**sex** - płeć,  
**age** - wiek,  
**sibbsp** - rodzzeń/swoiomałżonkowie na pokładzie (tak/nie),  
**parch** - liczba rodziców/dzieci na pokładzie.

Klasyfikacja będzie dotyczyła tego czy osoba płynąca Titanic przeżyła katastrofę.

Atrybut decyzyjny - kolumna **survived**.

Many dane modułowości: 0 - osoba nie przeżyła, 1 - osoba przeżyła

Algorytm klasyfikacji wybrane przeze mnie: **XGBoost** i **RandomForest**

Przebieg ze zbioru usuniętych wartości z kolumny Age, w których brakowało wartości. Następnie, pozostałe informacje o wieku pasażera zamieniam na zero jeżeli przetrwał (w zależności od tego czy osoba miała portret), czy powyżej 18 lat). Typem **male** i **female** z kolumny mówiącej o płci - również przypisalem 0 lub 1, żeby opierać na samych zmiennych numerycznych. Następnie reszcie kolumn mających braki w danych i kolumn według mnie mało istotnych, takich jak np. cel podróży (ang. home.dest) - wyrzucilem.

In [2]:

```
df = pd.read_csv('titanic2.csv')
display(df)

df = df.dropna(subset=['age'])
df['age'] = df['age'].apply(lambda x: int(x == 18))
df['sex'] = df['sex'].apply(lambda x: int(x == 'male'))
print(df.isnull().any())
df = df.dropna(axis=1)
df.drop(['name', 'ticket'], axis=1, inplace=True)
X = df.drop(['survived'], axis=1)
y = df['survived']
X = df.drop(['survived'], axis=1)
display(df)
```

```
pclass survived name sex age sibsp parch ticket fare cabin embarked boat body home.dest
0 1 1 Allen, Miss Elizabeth Walton female 29.0000 0 0 34100 211.3375 85 S 2 NaN Montreal, PQ / Chesapeake, ON
1 1 1 Allison, Master Hudson Trevor male 0.9167 1 2 113781 151.5500 C22 C26 S 11 NaN Montreal, PQ / Chesapeake, ON
2 1 0 Allison, Mrs. Helen Louise female 2.0000 1 2 113781 151.5500 C22 C26 S NaN NaN Montreal, PQ / Chesapeake, ON
3 1 0 Allison, Mr. Hudson Joshua Creighton male 30.0000 1 2 113781 151.5500 C22 C26 S NaN 125.0 Montreal, PQ / Chesapeake, ON
4 1 0 Allison, Mrs. Hudson J C (Bessie Wade Daniels) female 25.0000 1 2 113781 151.5500 C22 C26 S NaN NaN Montreal, PQ / Chesapeake, ON
...
```

```
1304 3 0 Zabour, Miss Helen female 14.5000 1 0 2665 14.4542 NaN C NaN 328.0 NaN
1305 3 0 Zabour, Miss Theresia female NaN 1 0 2665 14.4542 NaN C NaN NaN
1306 3 0 Zakarian, Mr. Napredoter male 26.5000 0 0 2656 7.2250 NaN C NaN 304.0 NaN
1307 3 0 Zakarian, Mr. Otin male 27.0000 0 0 2670 7.2250 NaN C NaN NaN
1308 3 0 Zimmerman, Mr. Leo male 29.0000 0 0 315082 7.8750 NaN S NaN NaN
```

1309 rows x 14 columns

```
pclass survived sex age sibsp parch
0 1 1 0 0 1 0 0
1 1 1 1 1 0 1 2
2 1 0 0 0 1 1 2
3 1 0 1 1 1 1 2
4 1 0 0 1 1 1 2
...
```

```
1304 3 0 0 1 1 0 0
1305 3 0 0 0 1 0 0
1306 3 0 1 1 0 0 0
1307 3 0 1 1 0 0 0
1308 3 0 1 1 0 0 0
1309 rows x 6 columns
```

In [3]:

```
seed=27
train_X, test_X, train_Y, test_Y = train_test_split(X, Y,
                                                    test_size = 0.2, random_state = seed)

def model_train(model, train_X, test_X, train_Y, test_Y):
    model.fit(train_X, train_Y)
    pred = model.predict(test_X)
    print ("Accuracy: ", accuracy_score(test_Y, pred))

print("XGBoost:")
model_train(xgb.XGBClassifier(use_label_encoder=False, eval_metric='nlogloss'), train_X, test_X, train_Y, test_Y)
print("RandomForest:")
model_train(RandomForestClassifier(), train_X, test_X, train_Y, test_Y)
```

XGBoost:  
Accuracy: 0.8265714285714286  
RandomForest:  
Accuracy: 0.8190476190476191

Przeanalizowałem również działanie klasyfikatorów na zbiorze bez zmiany wieku na 0/1

In [4]:

```
df = pd.read_csv('titanic2.csv')
df = df.dropna(subset=['age'])
df['sex'] = df['sex'].apply(lambda x: int(x == 'male'))
df['age'] = df['age'].apply(lambda x: int(x == 18))
df = df.dropna(axis=1)
df.drop(['name', 'ticket'], axis=1, inplace=True)
y = df['survived']
X = df.drop(['survived'], axis=1)
train_X, test_X, train_Y, test_Y = train_test_split(X, Y,
                                                    test_size = 0.2)
```

```
print("XGBoost:")
model_train(xgb.XGBClassifier(use_label_encoder=False, eval_metric='nlogloss'), train_X, test_X, train_Y, test_Y)
print("RandomForest:")
model_train(RandomForestClassifier(), train_X, test_X, train_Y, test_Y)

XGBoost:
Accuracy: 0.8095238962389623
RandomForest:
Accuracy: 0.7619047619047619
```

**Generowanie z wiekiem zero-jedynkowymi wyniki algorytmów są lepsze. Poza tym oba algorytmy mają przybliżoną dokładność.**

Stosunek zbioru treningowego do testowego ustawiony na 7:3

In [5]:

```
train_X, test_X, train_Y, test_Y = train_test_split(X, Y,
                                                    test_size = 0.4, random_state = seed)

def model_train(model, train_X, test_X, train_Y, test_Y):
    model.fit(train_X, train_Y)
    pred = model.predict(test_X)
    print ("Accuracy: ", accuracy_score(test_Y, pred))

print("XGBoost:")
model_train(xgb.XGBClassifier(use_label_encoder=False, eval_metric='nlogloss'), train_X, test_X, train_Y, test_Y)
print("RandomForest:")
model_train(RandomForestClassifier(), train_X, test_X, train_Y, test_Y)
```

XGBoost:  
Accuracy: 0.792993885732485  
RandomForest:  
Accuracy: 0.789808911074523

Stosunek zbioru treningowego do testowego ustawiony na 6:4

In [6]:

```
train_X, test_X, train_Y, test_Y = train_test_split(X, Y,
                                                    test_size = 0.4, random_state = seed)

def model_train(model, train_X, test_X, train_Y, test_Y):
    model.fit(train_X, train_Y)
    pred = model.predict(test_X)
    print ("Accuracy: ", accuracy_score(test_Y, pred))

print("XGBoost:")
model_train(xgb.XGBClassifier(use_label_encoder=False, eval_metric='nlogloss'), train_X, test_X, train_Y, test_Y)
print("RandomForest:")
model_train(RandomForestClassifier(), train_X, test_X, train_Y, test_Y)
```

XGBoost:  
Accuracy: 0.787594898060826  
RandomForest:  
Accuracy: 0.7947494833422887

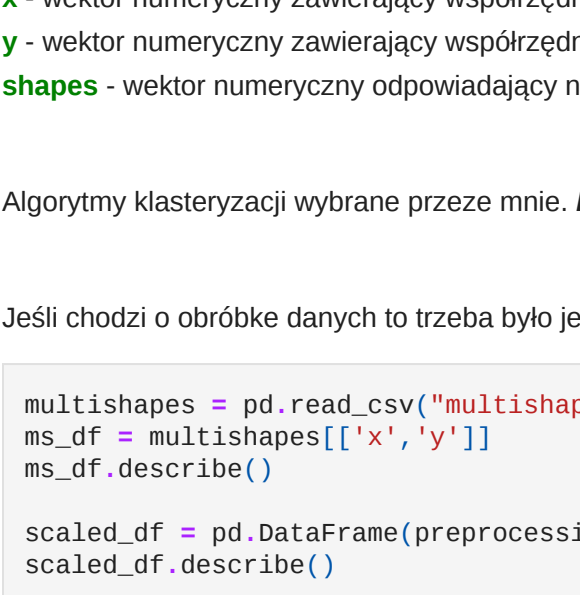
**W zależności od stosunku zbioru treningowego do testowego, ale też i od obserwacji wykazanych do poszczególnych grup - wyniki klasyfikacji oscylują w okolicy 80% dokładności (arg. accuracy).**

Poza klasyfikacją przeanalizowałem cały zbiór za pomocą wykresów.

Ile osób przeżyło, a ile zginęło:

In [7]:

```
sns.set_palette("set1")
sns.catplot(data=df,
            x='pclass',
            kind='count',
            hue='survived')
plt.show()
```

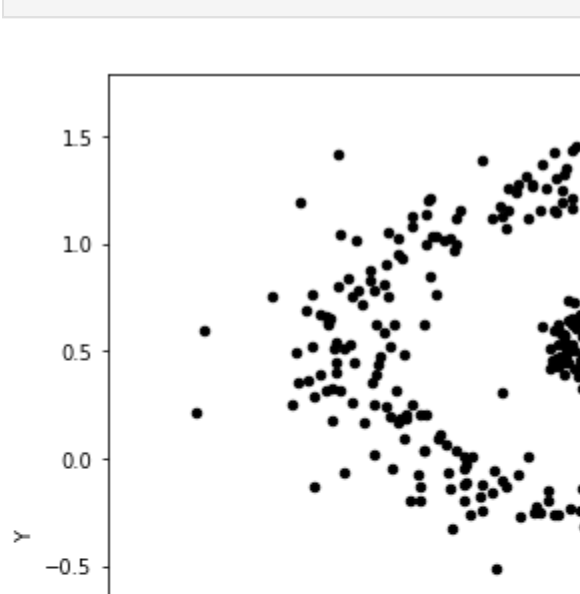


Czy klasa którą podróżowały osoby, miała wpływ na przeżycie?

1 - klasa najwyższa (premium). Być może, osoby podróżujące tą klasą, miały priorytet przy ewakuacji, analogicznie do osób wsiadających do samolotu na lotnisku, a być może ich bagaży, były zaskładowane w specjalnych części statku. Uważając lub utrudniając ewakuację.

In [8]:

```
sns.catplot(data=df,
            x='pclass',
            kind='count',
            hue='survived')
plt.show()
```

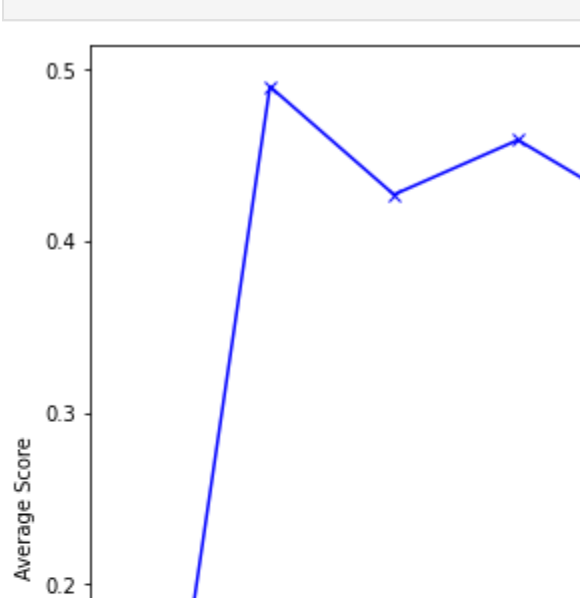


Czy wiek osoby, miał wpływ na przeżycie?

Widzimy, że większość dzieci przeżyła. Nie ma widocznej różnicy, dla osób po 20 roku życia. Zapewne dzieci były traktowane priorytetowo.

In [9]:

```
sns.distplot((df['age'].notnull() & (df['survived']==1))['age'],
            kde_kws={"label": "Survived"},
            bins=10)
sns.distplot((df['age'].notnull() & (df['survived']==0))['age'],
            kde_kws={"label": "Not Survived"},
            bins=18)
plt.show()
```

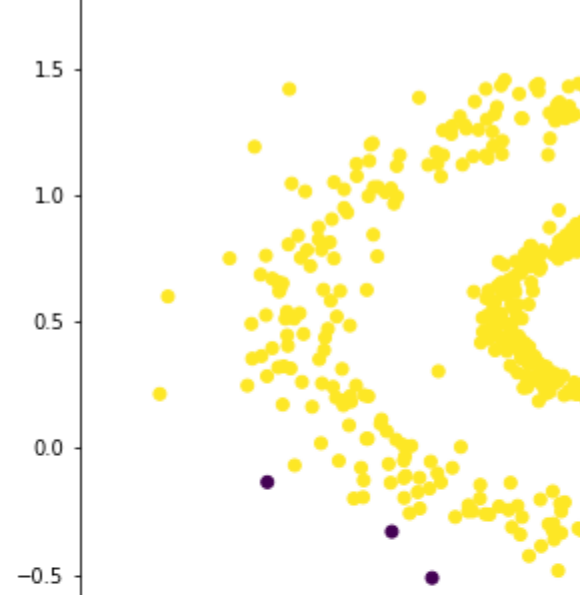


Czy płeć osoby, miała wpływ na przeżycie?

Kobiety miały większe szanse na przeżycie. Prawdopodobnie miały pierwszeństwo

In [10]:

```
sns.catplot(data=df,
            x='sex',
            kind='count',
            hue='survived',
            kind='count')
plt.show()
```



## Klasteryzacja

Zbiór danych użyty do klasteryzacji: *multishapes.csv*  
<https://www.datacollection.org/packages/factototat/versions/0.7.0/topics/multishapes> (1100 obserwacji i 3 zmienne)

Krótki opis zbioru:

Dane zawierające klasty dowolnych kształtów.

Opis zmiennych:

**x** - wektor numeryczny zawierający współrzędne obserwacji x,  
**y** - wektor numeryczny zawierający współrzędne obserwacji y,  
**shapes** - wektor numeryczny odpowiadający numerowi klasty każdej obserwacji.

Algorytm klasteryzacji wybrane przeze mnie: **KMeans** i **DBSCAN**

Jeśli chodzi o obróbkę danych to trzeba było je przekształcić, aby poszczególne pomary nie dominowały w późniejszych obliczeniach.

In [11]:

```
multishapes = pd.read_csv('multishapes.csv')
ms_df = multishapes[['x', 'y']]
ms_df.describe()

scaled_df = pd.DataFrame(preprocessing.scale(ms_df), index=multishapes['shape'], columns = ms_df.columns)
scaled_df.describe()
```

Out[11]:

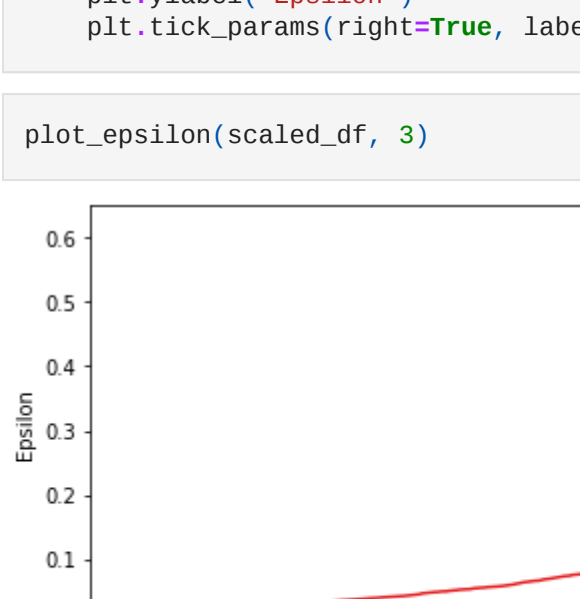
```
shape
1 1.32749 -0.197016
1 1.448907 0.944692
1 1.564203 0.298101
1 -1.043483 0.960859
1 1.222429 1.221561
...
```

1100 rows x 2 columns

Przedstawienie zbioru w formie graficznej

In [12]:

```
msplot = scaled_df.plot.scatter(x='x', y='y', c='Black', title='Multishapes data', figsize=(11, 8.5))
msplot.set_xlabel('x')
msplot.set_ylabel('y')
plt.show()
```



## KMeans

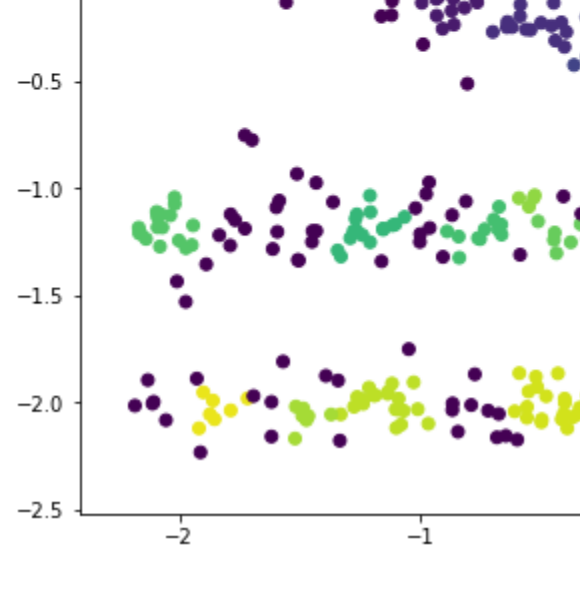
Do znalezienia optymalnej liczby klastów - ręcznie zapiełem algorytm KMeans, aby później przedstawić jego wyniki na wykresie i wybrać najlepsze n.

In [13]:

```
from sklearn.metrics import silhouette_score

scores = []
for i in range(2, 11):
    fit = KMeans(n_clusters=i, init='random', n_init=5, random_state=109).fit(scaled_df)
    distances, indices = fitted_neighbors.kneighbors(fit)
    scores.append(score)

plt.figure(figsize=(11, 8.5))
plt.plot(range(2, 11), np.array(scores), 'b-')
plt.xlabel('Number of clusters Ks')
plt.ylabel('Average Score')
plt.show()
```

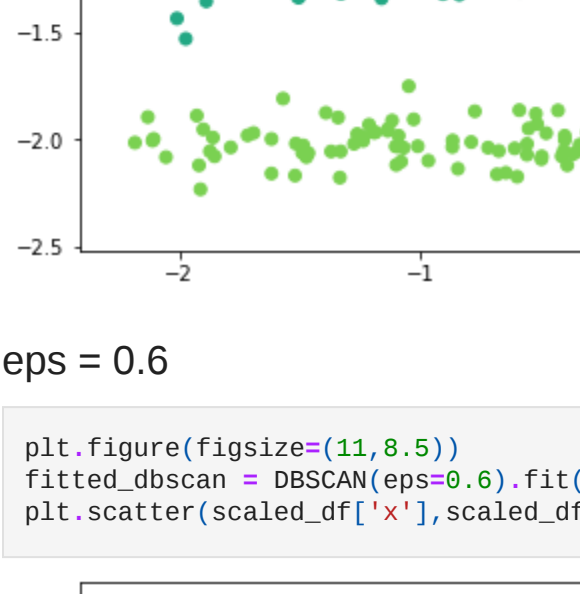


Wykres sugeruje że n przyjąć 2.

**n\_clusters = 2**

In [14]:

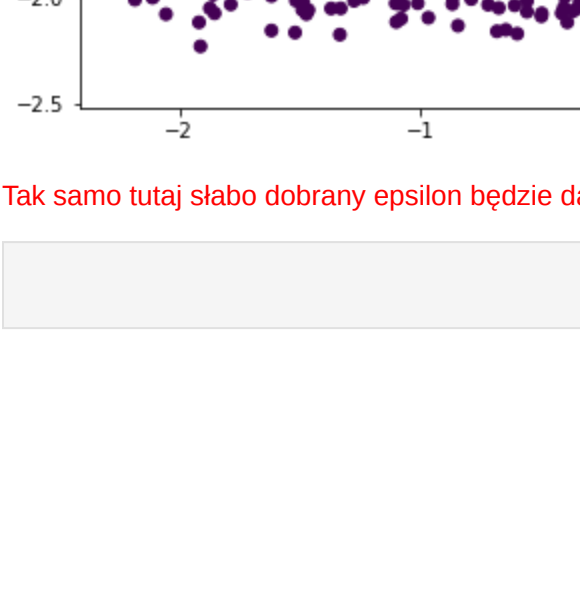
```
ms_kmeans = KMeans(n_clusters=2, init='random', n_init=2, random_state=109).fit(scaled_df)
plt.figure(figsize=(10, 10))
plt.scatter(scaled_df['x'], scaled_df['y'], c=ms_kmeans.labels_)
plt.scatter(ms_kmeans.cluster_centers_[0], ms_kmeans.cluster_centers_[1], c='r', marker='h', s=100)
```



**n\_clusters = 3**

In [15]:

```
ms_kmeans = KMeans(n_clusters=3, init='random', n_init=3, random_state=109).fit(scaled_df)
plt.figure(figsize=(10, 10))
plt.scatter(scaled_df['x'], scaled_df['y'], c=ms_kmeans.labels_)
plt.scatter(ms_kmeans.cluster_centers_[0], ms_kmeans.cluster_centers_[1], c='r', marker='h', s=100)
```



**n\_clusters = 5**

In [16]:

```
ms_kmeans = KMeans(n_clusters=5, init='random', n_init=5, random_state=109).fit(scaled_df)
plt.figure(figsize=(10, 10))
plt.scatter(scaled_df['x'], scaled_df['y'], c=ms_kmeans.labels_)
plt.scatter(ms_kmeans.cluster_centers_[0], ms_kmeans.cluster_centers_[1], c='r', marker='h', s=100)
```



**n\_clusters = 6**

In [17]:

```
plt.figure(figsize=(11, 8.5))
fitted_dbSCAN = DBSCAN(eps=0.5).fit(scaled_df)
plt.scatter(scaled_df['x'], scaled_df['y'], c=fitted_dbSCAN.labels_)
```



**eps = 0.2**

In [18]:

```
plt.figure(figsize=(11, 8.5))
fitted_dbSCAN = DBSCAN(eps=0.2).fit(scaled_df)
plt.scatter(scaled_df['x'], scaled_df['y'], c=fitted_dbSCAN.labels_)
```



**eps = 0.6**

In [19]:

```
plt.figure(figsize=(11, 8.5))
fitted_dbSCAN = DBSCAN(eps=0.6).fit(scaled_df)
plt.scatter(scaled_df['x'], scaled_df['y'], c=fitted_dbSCAN.labels_)
```



**Tak samo zbyt słabo dobrany epsilon będzie dawał słabe wyniki, tak więc warto znać funkcję, dzięki której z łatwością określmy najlepsze parametry.**

In [ ]: