



Neuronowy model klienta bankowego (Bank Marketing Data)

KWD

Grzegorz Podwika

Marek Ziaja

01.02.2021

Streszczenie

PROJEKT MA NA CELU STWORZENIE KLASYFIKATORA REALIZUJĄCEGO PRZEWIDYWANIE CHĘCI UŁOKOWANIA SWOICH ŚRODKÓW W LOKACIE TERMINOWEJ W BANKU. PROGRAM KOMPUTEROWY NA PODSTAWIE WCZEŚNIEJ DOSTARCZONYCH DANYCH, ZBIERANYCH PRZEZ PRACOWNIKÓW BANKU TELEFONUJĄCYCH DO KLIENTÓW, TRENUJE KLASYFIKATOR KTÓRY W KOLEJNYM KROKU ZOSTAJE PODDANY TESTOM. NASZYM ZADANIEM JEST ANALIZA WYNIKÓW PROGRAMU I WYBRANIE ORAZ OPRACOWANIE DAJĄCEGO NAJLEPSZE WYNIKI MODELU KLASYFIKUJĄCEGO.

Spis treści

1	Wprowadzenie	1
1.1	Opis problemu	1
2	Realizacja	2
2.1	Wprowadzenie teoretyczne	2
2.2	Opis danych wejściowych	3
2.3	Przekształcenie zbioru danych	10
2.4	Normalizacja i standaryzacja danych	11
2.5	Ewaluacja uzyskanego modelu	12
2.6	Badania symulacyjne	14
3	Podsumowanie	21
A	Kod programu	22

Rozdział 1

Wprowadzenie

BANKI NA CAŁYM ŚWIECIE GŁOWIĄ SIĘ I TROJĄ, ABY REALIZOWANE PRZEZ NICH STRATEGIE PRZYNOSIŁY JAK NAJWIĘKSZE ZYSKI. JEDNĄ Z OBRANYCH STRATEGII PRZEZ BANK BYŁO KILKU/KILKUNASTOKROTNE PRZEPROWADZANIE ROZMÓW TELEFONICZNYCH PRZEZ KONSULTANTÓW Z KLIENTAMI BANKU Z ZAMIAREM ZAŁOŻENIA KORZYSTNEJ LOKATY TERMINOWEJ. ZBIERAJĄC DANE Z TEGO TYPU PRACY MOŻEMY OPRACOWAĆ MODEL KLIENTA BANKOWEGO, ANALIZOWAĆ GO I ULEPSZAĆ TĄ KONKRETNĄ STRATEGIĘ Z WYMIERNYMI KORZYŚCIAMI DLA BANKU. DZIĘKI TEMU NASTĘPNE KAMPANIE PRZEPROWADZANE PRZEZ BANK BĘDĄ BARDZIEJ EFEKTYWNE I PRZYNOSIĄ KORZYŚCI ZARÓWNO DLA KLIENTÓW JAK I SAMEGO BANKU.

1.1 Opis problemu

PODSTAWOWYM PROBLEMEM JEST STWORZENIE KLASYFIKATORA, KTÓRY PRZY KLASYFIKOWANIU DANYCH, CZYLI W TYM PRZYPADKU CECH KLIENTA BANKU, KTÓRY BĘDZIE POPEŁNIAŁ MOŻLIWIE NAJMNIEJSZĄ ILOŚĆ BŁĘDÓW PRZY PREDYKCJI NASTĘPNYCH DANYCH. POD UWAGĘ NALEŻY WZIĄĆ ODPOWIEDNI DOBÓR DANYCH NA KTÓRYCH NASZ KLASYFIKATOR BĘDZIE TRENOWANY, PONIEWAŻ NA ICH PODSTAWIE BĘDZIE PODEJMOWAĆ DECYZJĘ. MODEL MUSI TAKŻE ZOSTAĆ DOKŁADNIE PRZETESTOWANY NA RÓŻNE SPOSOBY, UWZGLĘDNIAJĄC ZMIANĘ DANYCH TESTOWYCH I TRENUJĄCYCH, TAK BY ZMINIMALIZOWAĆ BŁĄD, A ZMAKSYMALIZOWAĆ WIARYGODNOŚĆ NASZEGO MODELU. UWAGĘ NALEŻY POŚWIĘCIĆ TAKŻE ROZKŁADZIE DANYCH NA TRENINGOWE I TESTUJĄCE. TWORZĄC KAŻDY MODEL NALEŻY ZADBAĆ O TO, BY NIE BYŁ NIEDOUUCZONY, ALE TAKŻE ABY NIE BYŁ PRZEUCZONY, PONIEWAŻ W TAKIM PRZYPADKU MODEL ZA BARDZO PODAŻA ZA BŁĘDAMI/OUTLINERAMI NIŻ ZA OGÓLNYM TRENDEM, KTÓREGO CHCEMY UCHWYCIĆ.

Rozdział 2

Realizacja

2.1 Wprowadzenie teoretyczne

ODMIANĄ REGRESJI UŻYWANĄ DO ROZWIĄZANIA NASZEGO PROBLEMU JEST REGRESJA LOGISTYCZNA. NAJWAŻNIEJSZĄ CECHĄ REGRESJI LOGISTYCZNEJ JEST TO, ŻE ETYKIETA PRZYJMUJE TYLKO DWA WARTOŚCI (NAJCZĘŚCIEJ 0 LUB 1). UŻYWAMY JEJ DO OKREŚLANIA WYSTĄPIENIA JAKIEGOŚ ZDARZENIA LUB ZJAWISKA. REGRESJA LOGISTYCZNA RÓŻNI SIĘ ZASADNICZO OD REGRESJI LINIOWEJ. W PRZYPADKU REGRESJI LOGISTYCZNEJ W PRZECIWIENSTWIE DO REGRESJI LINIOWEJ CELEM NIE JEST PRZEWIDZENIE WARTOŚCI ZMIENNEJ ZALEŻNEJ NA PODSTAWIE WYKORZYSTANIA PREDYKATÓW, ALE PRZEWIDZENIE PRAWDOPODOBIENSTWA WYSTĄPIENIA JAKIEGOŚ ZDARZENIA.

WZÓR NA RÓWNANIE REGRESJI LOGISTYCZNEJ:

$$P(Y = 1 | x_1, x_2, \dots, x_k) = P(X) = \frac{e^{a_0 + \sum_{i=1}^k a_i x_i}}{1 + e^{a_0 + \sum_{i=1}^k a_i x_i}}$$

Gdzie:

$a_i, i = 0, \dots, k$ są współczynnikami regresji,

x_1, x_2, \dots, x_k - zmienne niezależne (ilościowe lub jakościowe).

ANALIZA REGRESJI LOGISTYCZNEJ W PORÓWNANIU DO REGRESJI LINIOWEJ OBARCZONA JEST MNIJSZĄ ILOŚCIĄ ZAŁOŻEŃ. NAJWAŻNIEJSZĄ CECHĄ REGRESJI LOGISTYCZNEJ JEST TO, ŻE PREDYKATEM JEST ZMIENNA DYCHOTOMICZNA 0 - 1.

METODĄ UŻYWANĄ PODCZAS TESTOWANIA NASZEGO MODELU JEST WALIDACJA KRZYŻOWA. JEST TO METODA STATYSTYCZNA POLEGAJĄCA NA

PODZIAŁE PRÓBY STATYSTYCZNEJ NA PODZBIORY, A NASTĘPNIE PRZEWODZENIU WSZELKICH ANALIZ NA ZBIORZE UCZĄCYM, PODCZAS GDY ZBIÓR TESTOWY SŁUŻY DO POTWIERDZENIA WIARYGODNOŚCI WYNIKÓW. METODA TA JEST STOSOWANA W PRZYPADKU NIEDOSTATECZNIE LICZNEGO ZBIORU UCZĄCEGO.

2.2 Opis danych wejściowych

ZBIÓR DANYCH NA KTÓRYCH DZIAŁAĆ BĘDZIE KLASYFIKATOR SKŁADA SIĘ Z 41188 REKORDÓW, Z KTÓRYCH KAŻDY OKREŚLA KONKRETNEGO KLIEN-TA BANKU. W KAŻDEJ KROTCE ZNAJDUJE SIĘ 20 CECH OKREŚLAJĄCYCH KLIENTA. NA ICH PODSTAWIE WYZNACZANA JEST ETYKIETA, KTÓRA OKREŚLA CZY UDAŁO SIĘ NAMÓWIC KLIENTA NA LOKATĘ - NIEPOWODZENIE ("NO"), BĄDŹ POWODZENIE ("YES").

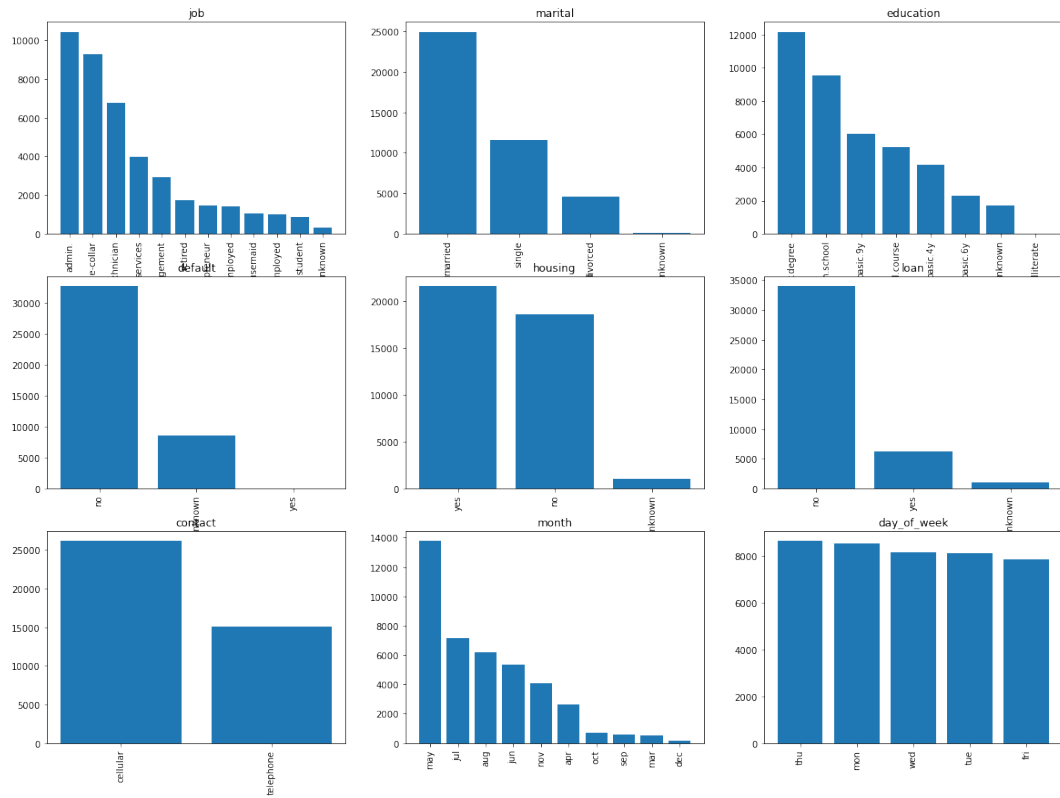
INFO O ZBIORZE DANYCH:

```
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   age                    41188 non-null  int64
1   job                    41188 non-null  object
2   marital                41188 non-null  object
3   education              41188 non-null  object
4   default                41188 non-null  object
5   housing                41188 non-null  object
6   loan                   41188 non-null  object
7   contact                41188 non-null  object
8   month                  41188 non-null  object
9   day_of_week            41188 non-null  object
10  duration                41188 non-null  int64
11  campaign                41188 non-null  int64
12  pdays                  41188 non-null  int64
13  previous                41188 non-null  int64
14  poutcome                41188 non-null  object
15  emp.var.rate            41188 non-null  float64
16  cons.price.idx           41188 non-null  float64
17  cons.conf.idx            41188 non-null  float64
18  euribor3m                41188 non-null  float64
19  nr.employed              41188 non-null  float64
20  y                        41188 non-null  object
dtypes: float64(5), int64(5), object(11)
```

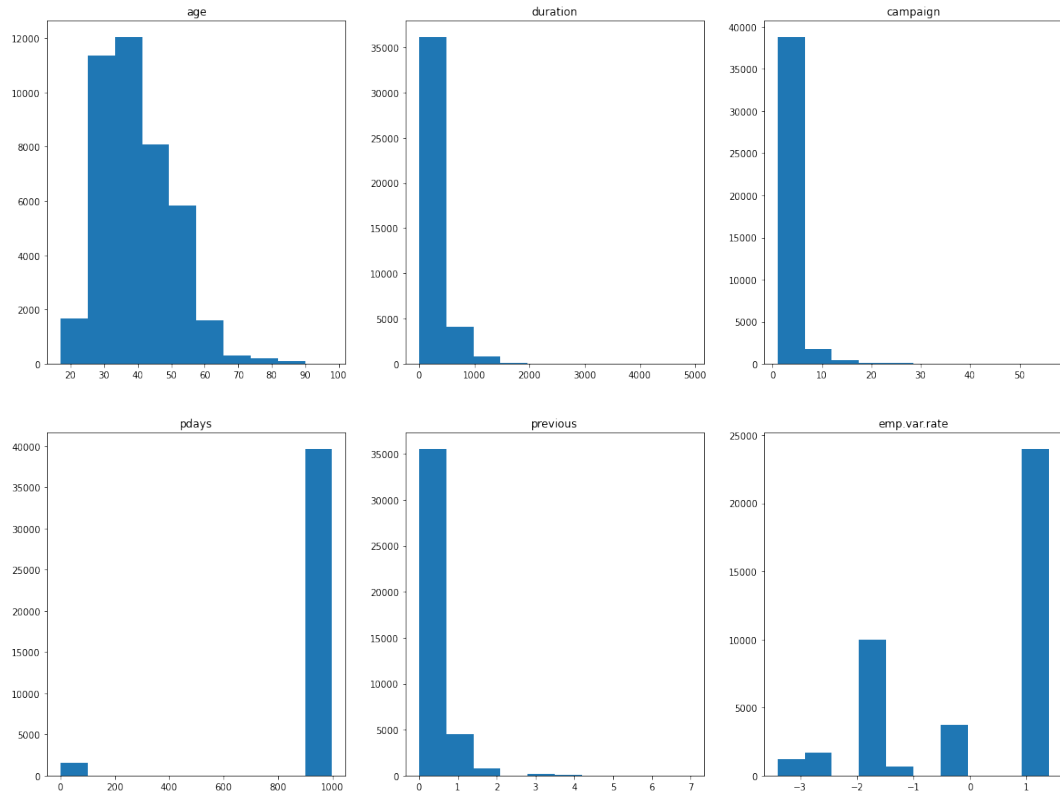
DANE OKREŚLAJĄCE KLIENTA:

	0	1	2	3	4
age	56	57	37	40	56
job	housemaid	services	services	admin.	services
marital	married	married	married	married	married
education	basic.4y	high.school	high.school	basic.6y	high.school
default	no	unknown	no	no	no
housing	no	no	yes	no	no
loan	no	no	no	no	yes
contact	telephone	telephone	telephone	telephone	telephone
month	may	may	may	may	may
day_of_week	mon	mon	mon	mon	mon
duration	261	149	226	151	307
campaign	1	1	1	1	1
pdays	999	999	999	999	999
previous	0	0	0	0	0
poutcome	nonexistent	nonexistent	nonexistent	nonexistent	nonexistent
emp.var.rate	1.1	1.1	1.1	1.1	1.1
cons.price.idx	93.994	93.994	93.994	93.994	93.994
cons.conf.idx	-36.4	-36.4	-36.4	-36.4	-36.4
euribor3m	4.857	4.857	4.857	4.857	4.857
nr.employed	5191	5191	5191	5191	5191

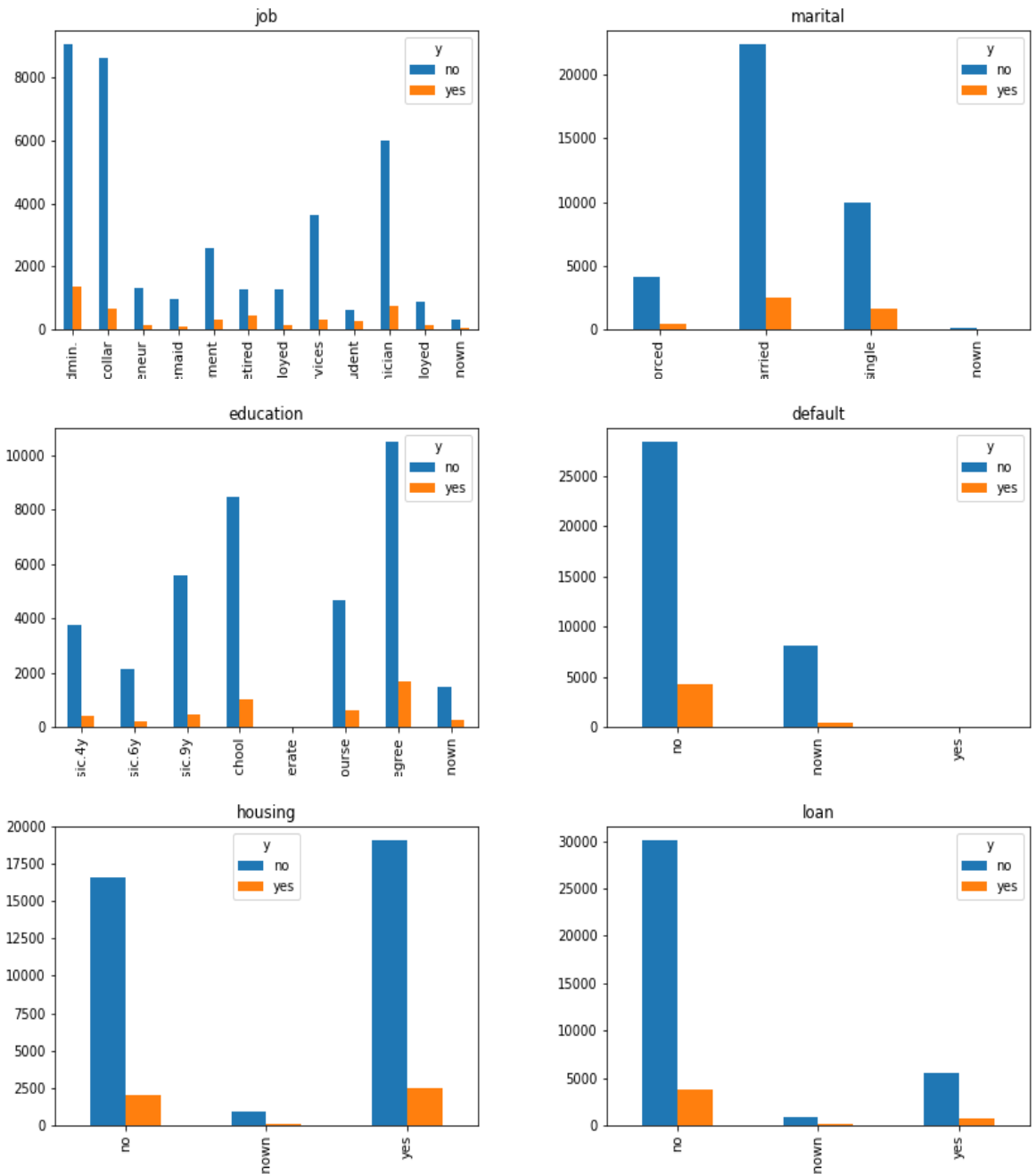
EKSPLORACJA KOLUMN DANYCH KATEGORYCZNYCH:

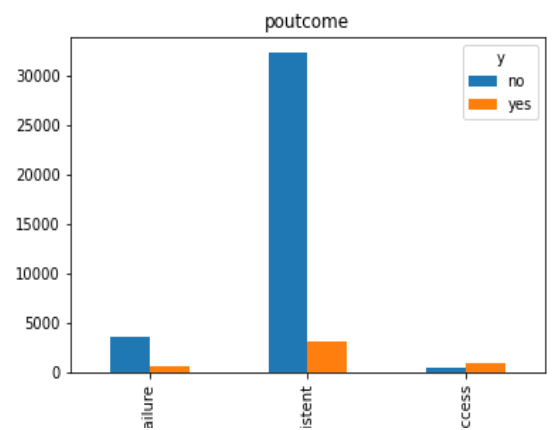
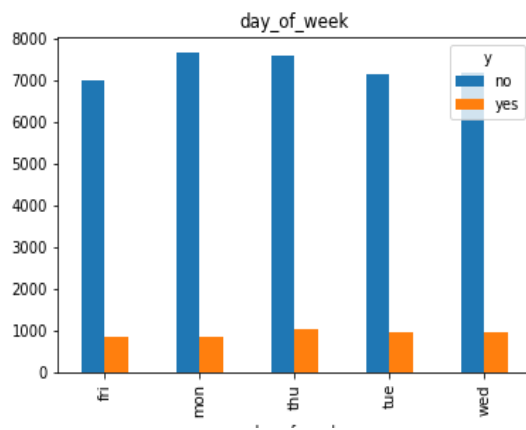
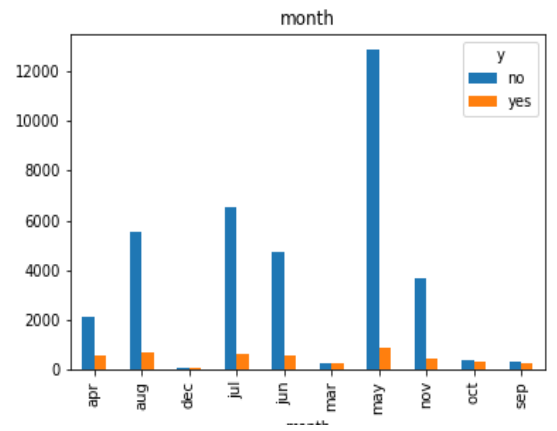
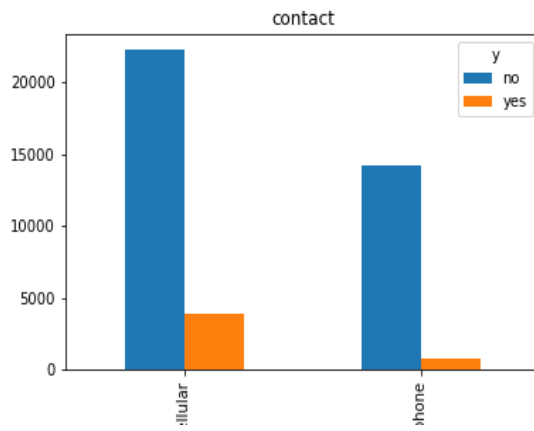


EKSPLORACJA KOLUMN DANYCH NUMERYCZNYCH:



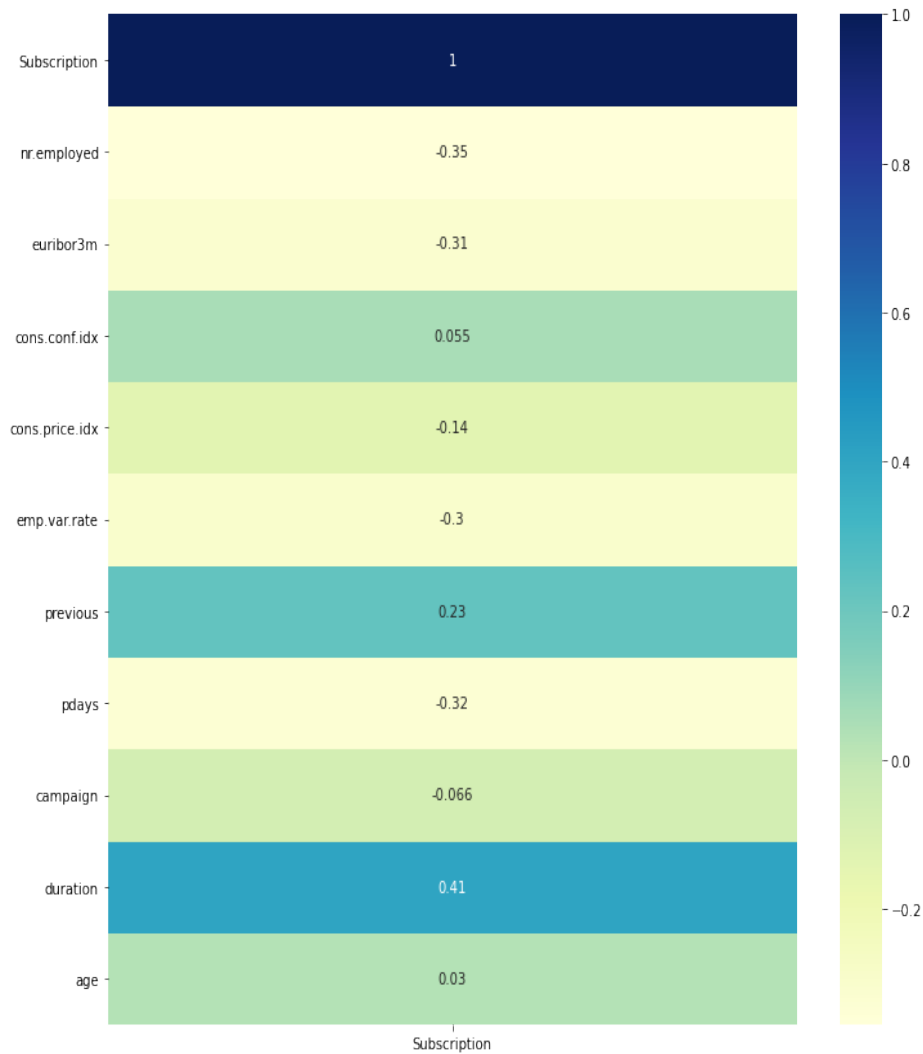
ROZKŁAD DANYCH KATEGORYCZNYCH WRAZ Z ETYKIETAMI





Heatmap showing the correlation matrix for the variables: age, duration, campaign, pdays, previous, emp.var.rate, cons.price.idx, cons.conf.idx, euribor3m, and nr.employed. The color scale ranges from -0.4 (yellow) to 1.0 (dark blue).

WZAJEMNA KORELACJA CECH Z SUBSKRYPCJĄ



2.3 Przekształcenie zbioru danych

ZAPOZNAWSZY SIĘ Z POWYŻSZYMI ZDJĘCIAMI OBRAZUJĄCYMI ZBIÓR DANYCH KLIENTÓW BANKU MOŻEMY OD RAZU ZAUWAŻYĆ, ŻE CO PO NIEKTÓRE KOLUMNY SĄ ZAPISANE W POSTACI KATEGORYCZNYCH CO UNIEMOŻLIWIA UCZENIE NASZEGO MODELU, WIĘC MUSIMY ZAKODOWAĆ CECHY TYCH KOLUMN W POSTACI NUMERYCZNEJ. TO SAMO MUSIMY UCZYNIĆ Z ETYKIETAMI, PRZEKSZTAŁCAJĄC JE Z "YES" ORAZ "NO" NA ODPOWIEDNIO 1 I 0.

ENKODOWANIE DANYCH:

```
# Transform 'yes' and 'no' into (0,1) dependency
bank_data['y']=[1 if i=='yes' else 0 for i in bank_data['y']]
```

```
# Create dummy variable for the categorical data and drop first column
# to avoid dummy variable trap
bank_data = pd.get_dummies(bank_data, columns = ['job', 'marital', 'education',
                                                'default', 'housing', 'loan',
                                                'month', 'contact',
                                                'day_of_week', 'outcome'],
                           drop_first = True)
```

2.4 Normalizacja i standaryzacja danych

BIORĄC POD UWAGĘ ZRÓŻNICOWANIE ZAKRESU DANYCH W ZBIORZE, W OBLICZENIACH KONIECZNE JEST UWZGLĘDNIENIE NORMALIZACJI ORAZ STANDARYZACJI. STANDARYZACJA DANYCH DO MODELI MACHINE LEARNING POLEGA NA PRZEKSZTAŁCENIU DANYCH PIERWOTNYCH, ABY ICH ROZKŁAD MIAŁ ŚREDNIĄ WARTOŚĆ RÓWNĄ 0 I ODCHYLENIE STANDARDOWE RÓWNE 1. OD KAŻDEJ WARTOŚĆ W KOLUMNIE DANYCH BĘDZIE ODEJMOWANA ŚREDNIA WARTOŚĆ KOLUMNY, A NASTĘPNIE TO CO WYJDZIE BĘDZIE PODZIELONA PRZEZ ODCHYLENIE STANDARDOWE KOLUMNY DANYCH. OPISANY PROCES DOTYCZY KAŻDEJ KOLUMNY ODDZIELNIE. WIELE KLASYFIKATORÓW OBLICZA ODLEGŁOŚĆ MIĘDZY DWOMA PUNKTAMI NA PODSTAWIE ODLEGŁOŚCI EUKLIDESOWEJ. JEŚLI JEDNA Z CECH MA SZEROKI ZAKRES WARTOŚCI, ODLEGŁOŚĆ BĘDZIE ZALEŻAŁA OD TEJ KONKRETNEJ CECHY. DLATEGO ZAKRES WSZYSTKICH CECH POWINIEN ZOSTAĆ ZNORMALIZOWANY, TAK ABY KAŻDA CECHA MIAŁA UDZIAŁ W PRZYBLIŻENIU PROPORCJONALNIE DO KOŃCOWEJ ODLEGŁOŚCI.

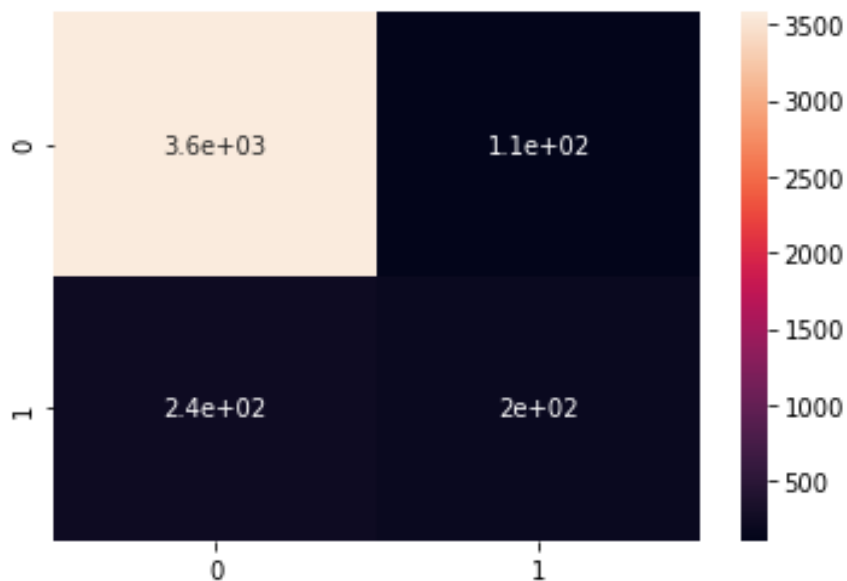
STANDARYZACJA DANYCH DLA MODELU:

```
# scale data
sc = StandardScaler()
bank_train_data = sc.fit_transform(bank_train_data)
bank_test_data = sc.transform(bank_test_data)
```

2.5 Ewaluacja uzyskanego modelu

WYTRENOWANY PRZEZ NAS MODEL ZA POMOCĄ REGRESJI LOGISTYCZNEJ DLA PRZESKALOWANYCH DANYCH DAJE NASTĘPUJĄCE WYNIKI:

```
Basic Logistic Regression with scaled data
Basic model accuracy of scaled data= 0.9169701383831027
[[3579  107]
 [ 235  198]]
```



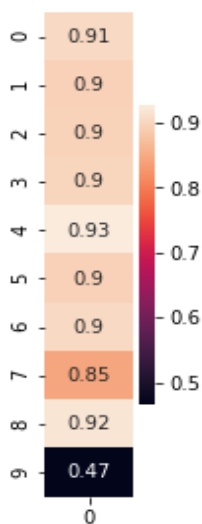
WYNIK MODELU JEST PRZEDSTAWIONY ZA POMOCĄ WSPÓŁCZYNNIKA DOKŁADNOŚCI - ACCURACY SCORE. OBLICZAMY GO DZIELĄC LICZBĘ POPRAWNYCH DOPASOWAŃ PRZEZ WSZYSTKIE PRÓBY DOPASOWANIA. W NASZYM PRZYPADKU DOKŁADNOŚĆ WYNIOSŁA OKOŁO 91 %, CO JEST PRZYZWOITYM WYNIKIEM. MACIERZ KONFUZJI PRZEDSTAWIA ILOŚĆ POPRAWNYCH

I NIEPOPRAWNYCH DOPASOWAŃ. W PIERWSZYM WIERSZU ZNAJDUJĄ SIĘ INFORMACJE O 3569 POPRAWNYCH DOPASOWANIACH BRAKU SUBSKRYPCJI (0), ORAZ O 98 POMYŁKACH, GDZIE MODEL PRZYPISAŁ BŁĘDNE ZDECYDOWANIE SIĘ NA OTWARCIE LOKATY, A W RZECZYWISTOŚCI KLIENT NIE ZDECYDOWAŁ SIĘ NA TEN KROK. W DRUGIM WIERSZU ZNAJDUJĄ SIĘ INFORMACJĘ O 206 POPRAWNYCH DOPASOWANIACH SUBSKRYPCJI LOKATY (1), ORAZ O 246 POMYŁKACH KIEDY MODEL PRZYPISAŁ BRAK SUBSKRYPCJI, A W RZECZYWISTOŚCI KLIENT ZDECYDOWAŁ SIĘ NA ZAŁOŻENIE LOKATY. PODCZAS TWORZENIA MODELU 90% DANYCH BYŁO ZBIOREM TRENUJĄCYM, NATOMIAST 10% DANYCH BYŁO ZBIOREM TESTUJĄCYM.

```
# split data and target sets into 4 sets for train and test purposes
bank_train_data, bank_test_data, \
bank_train_target, bank_test_target = \
train_test_split(data_np, target_np, test_size=0.1)
```

DLA NASZYCH DANYCH PRZEPROWADZILIŚMY RÓWNIEŻ KROSVALIDACJĘ, KTÓREJ WYNIKI SĄ ZBLIŻONE DO MODELU WYTRENOWANEGO PRZEZ NAS.

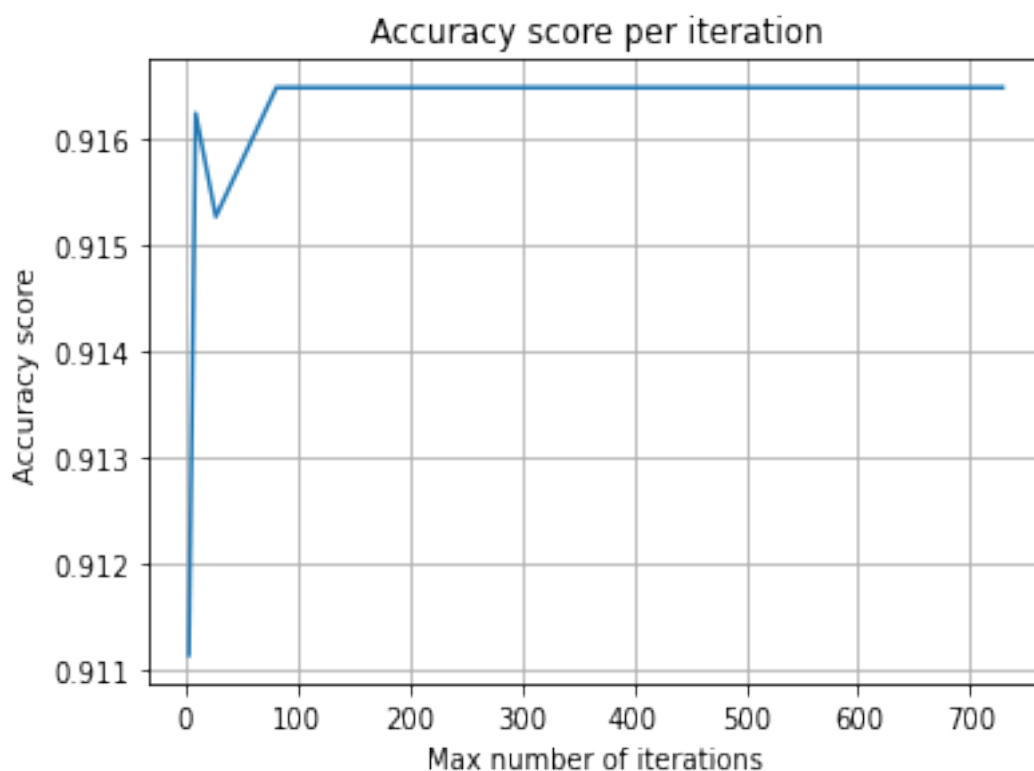
WYNIKI Z KROSVALIDACJI WYGLĄDAJĄ NASTĘPUJĄCO:



2.6 Badania symulacyjne

BADANIA WYNIKU ZMIENNOŚCI NASZEGO MODELU ROZPOCZĘLIŚMY OD OKREŚLANIA MAKSYMALNEJ LICZBY ITERACJI DLA ALGORYTMU BUDUJĄCEGO KLASYFIKATOR.

WYKRES PRZEDSTAWIAJĄCY DOKŁADNOŚĆ MODELU W ZALEŻNOŚCI OD MAKSYMALNEJ LICZBY ITERACJI:



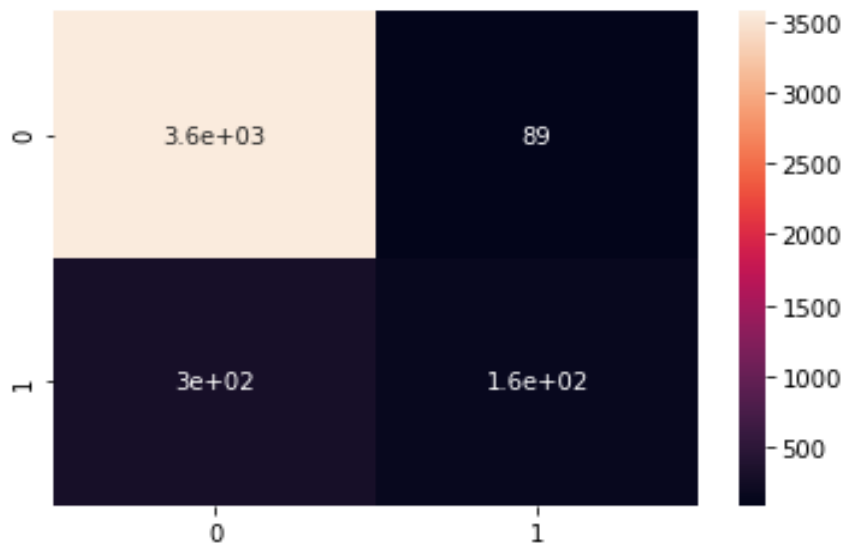
MAKSYMALNE ITERACJE SĄ POTĘGAMI LICZBY 3 DO 729 WŁĄCZNIE.

PO PRZEANALIZOWANIU WYKRESU MOŻNA WYWNIOSKOWAĆ, ŻE DOKŁADNOŚĆ WZRASTA WRAZ ZE WZROSTEM LICZBY ITERACJI, AŻ ZATRZYMUJE SIĘ NA POZIOMIE 91,64% I UTRZYMUJE SIĘ POMIMO ZWIĘKSZANIA LICZBY ITERACJI. PRZEPROWADZONE ZOSTAŁY RÓWNIEŻ BADANIA WPŁYWU ROZKŁADU DANYCH NA ZBIORY TESTUJĄCE I TRENUJĄCE:

Model LogisticRegression basic - Test data = 10 %

Model accuracy is 0.91

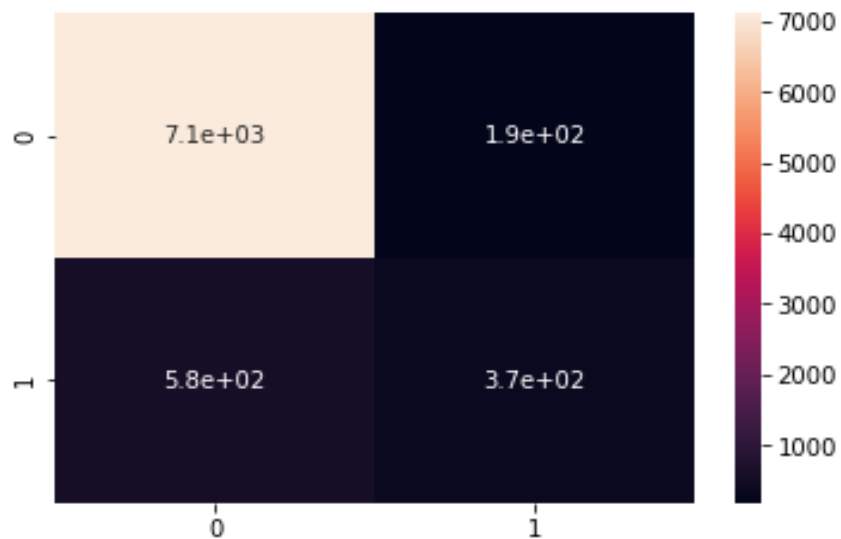
```
[[3570  89]  
 [ 295 165]]
```



Model LogisticRegression basic - Test data = 20 %

Model accuracy is 0.91

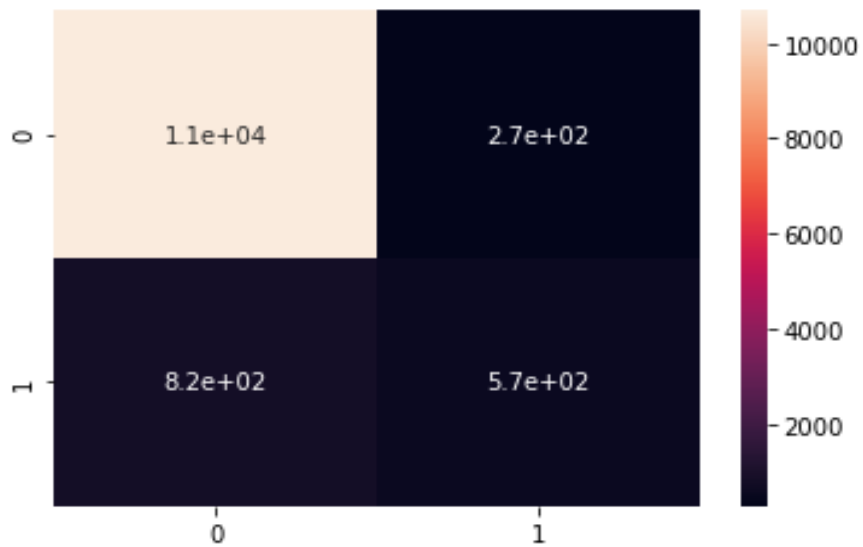
```
[[7107 189]  
 [ 575 367]]
```



Model LogisticRegression basic - Test data = 30 %

Model accuracy is 0.91

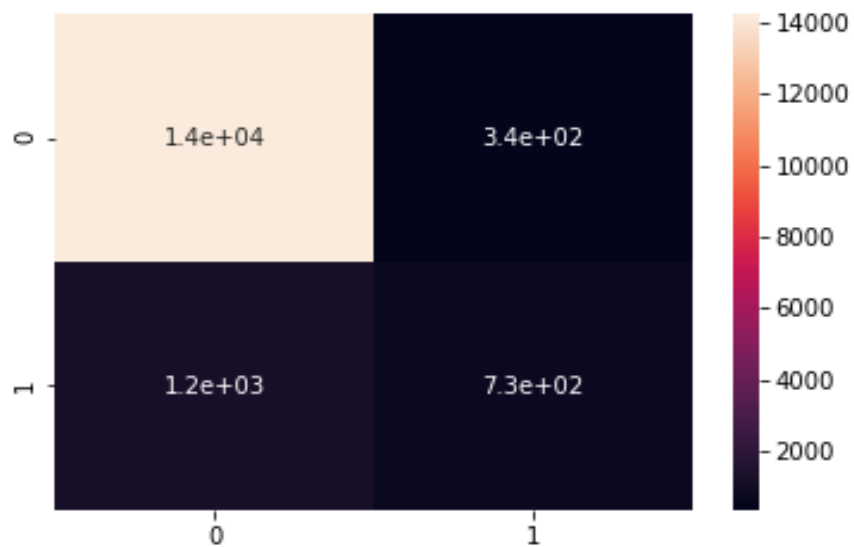
```
[[10689  274]
 [  820  574]]
```



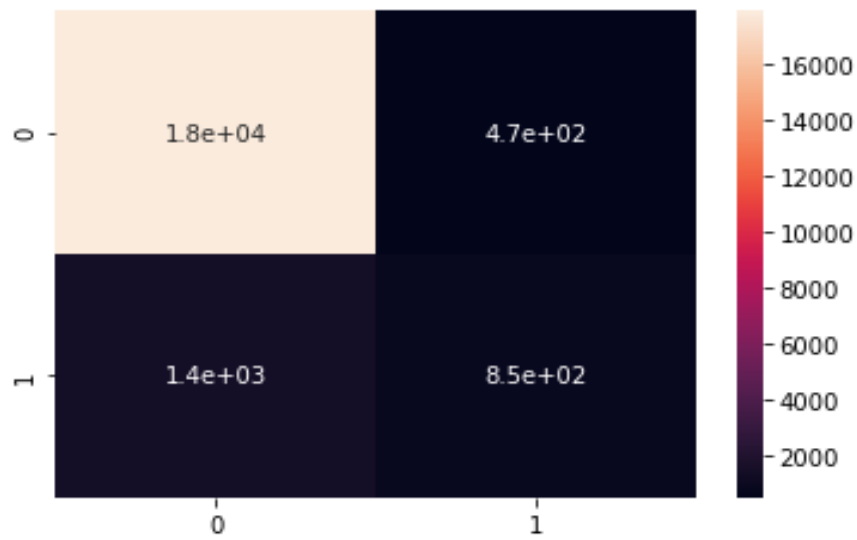
Model LogisticRegression basic - Test data = 40 %

Model accuracy is 0.91

```
[[14241  341]
 [ 1167  727]]
```



```
Model LogisticRegression basic - Test data = 50 %  
Model accuracy is 0.91  
[[17869  473]  
 [ 1399   853]]
```



ZWIĘKSZAJĄC ZBIÓR DANYCH TESTUJĄCYCH WYNIK POZOSTAJE BEZ ZMIAN.
SPRAWDZONY ZOSTAŁ TAKŻE WPŁYW RODZAJU ALGORYTMU OPTYMALIZA-

CYJNEGO NA REZULTAT DZIAŁANIA MODELU:



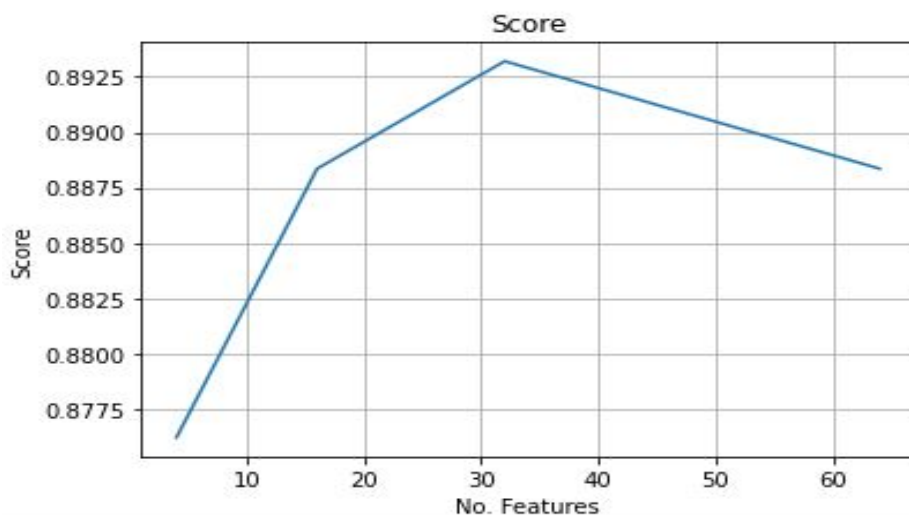
ZMIANA ALGORYTMU OPTYMALIZACYJNEGO NIE WPŁYNEŁA NA WYNIK MODELU.

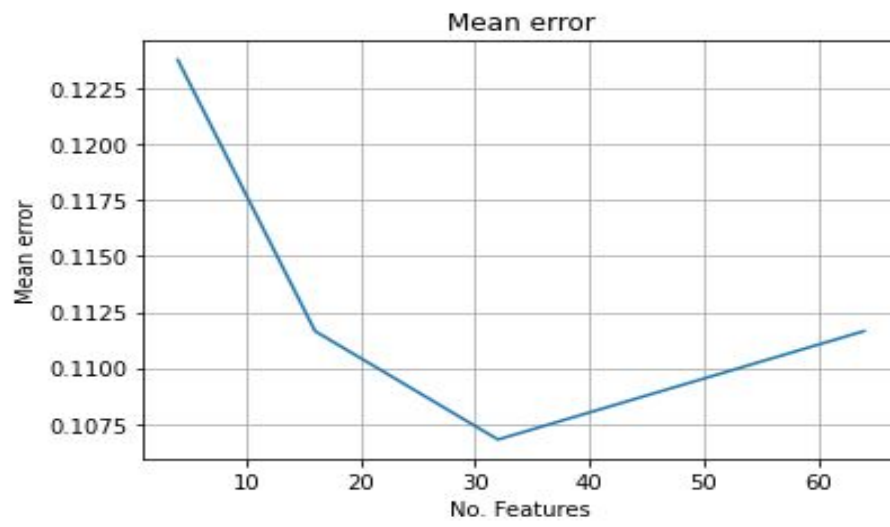
DO TESTÓW Z WYKORZYSTANIEM POLYNOMIAL FEATURES I RFE (REKURENCYJNEJ ELIMINACJI CECH) UŻYLIŚMY MNIEJSZEGO ZESTAWU DANYCH, PONIEWAŻ CZAS WYKONYWANIA TESTU ROŚNIE WYKŁADNICZO WRAZ Z ROZMIAREM ZBIORU.

ZA POMOCĄ POLYNOMIAL FEATURES ZESTAW CECH JEST PODNOSZONY DO POTĘGI (W TYM PRZYPADKU JEST TO 2). REKURENCYJNA ELIMINACJA ATRYBUTÓW POLEGA NA WYBIERANIU PODZBIORU NAJWAŻNIEJSZYCH CECH DLA ZBIORU DANYCH. MNIEJSZA LICZBA CECH MOŻE UMOŻLIWIĆ BARDZIEJ WYDAJNE DZIAŁANIE ALGORYTMÓW UCZENIA MASZYNOWEGO (MNIEJSZA ZŁOŻONOŚĆ PRZESTRZENNA LUB CZASOWA) I WIĘKSZĄ SKUTECZNOŚĆ. NIEKTÓRE ALGORYTMY UCZENIA MASZYNOWEGO MOGĄ BYĆ WPROWADZANE W BŁĄD PRZEZ NIEISTOTNE POŁĄCZENIA ATRYBUTÓW.

NA PONIŻSZYCH WYKRESACH MOŻNA DOSTRZEC ZALEŻNOŚĆ WYNIKU (DOKŁADNOŚCI) ORAZ ŚREDNIEGO BŁĘDU. TEST BYŁ PRZEPROWADZONY PODAJĄC RÓŻNE ILOŚCI ATRYBUTÓW (4, 16, 32, 64). MOŻNA ZAUWAŻYĆ MINIMUM BŁĘDU PRZY 32 ATRYBUTACH, DALEJ DLA 64 ATRYBUTÓW NIEWIELKI WZROST. NALEŻY MIEĆ NA UWADZE JEDNAK WYSKALOWANIE WYKRESU - RÓŻNICE BŁĘDÓW ZALEŻNE OD ILOŚCI WYBRANYCH ATRYBUTÓW (CO ZA TYM IDZIE TEŻ WYNIKÓW) SĄ NIEWIELKIE, ZAUWAŻALNE ZMIANY NA DRUGIM MIEJSCU PO PRZECINKU, CZYLI PROCENTOWY RZĄD WIELKOŚCI.

ANALIZUJĄC WYNIKI TESTÓW NALEŻY STWIERDZIĆ, ŻE REKURENCYJNA ELIMINACJA CECH NIE POPRAWIA WYNIKU TRENOWANIA MODELU.





Rozdział 3

Podsumowanie

KLASYFIKATOR WYTRENOWANY PRZEZ NAS DOBRZE RADZI SOBIE Z ANALIZĄ ZADANYCH DANYCH. NAJLEPSZYM UZYSKANYM WYNIKIEM DOKŁADNOŚCI BYŁO OKOŁO 92 PROCENT, CO JEST ZADOWALAJĄCE.

SKALOWANIE ZBIORU NIEZNACZNIE POLEPSZYŁO WYNIK MODELU.

DLA PODSTAWOWEJ REGRESJI LOGISTYCZNEJ OPTYMALNĄ ILOŚCIĄ ITERACJI BYŁO OKOŁO 80, ZWIĘKSZANIE ITERACJI NIE DAWAŁO LEPSZEGO EFEKTU.

MANIPULOWANIE ROZMIARAMI ZBIORÓW TESTOWYCH I TRENUJĄCYCH NIE NIESIE ZA SOBĄ ZMIAN W DOKŁADNOŚCI.

ALGORYTM OPTYMALIZUJĄCY RÓWNIEŻ NIE MA WPŁYWU NA SKUTECZNOŚĆ KLASYFIKATORA.

RFE POSKUTKOWAŁO NIEZNACZNYM SPADKIEM DOKŁADNOŚCI KLASYFIKATORA.

MIMO DOSKONALENIA PROCESU UCZENIA MASZYNOWEGO, NALEŻY ZAWSZE BRAĆ POD UWAGĘ TZW. "CZYNNIK LUDZKI", CZYLI MOŻLIWOŚĆ POJEDYNCZYCH ODSTĘPSTW OD OGÓLNIE WYPRACOWANEGO MODELU.

Dodatek A

Kod programu

```
#imports
import pandas as pd
import numpy as np
import seaborn as sns
from matplotlib import pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from sklearn.preprocessing import PolynomialFeatures
from sklearn.feature_selection import RFE

-----

# load data set
bank_data = pd.read_csv("bank-additional-full.csv", sep=';')
print("Shape of DataFrame is ", bank_data.shape)

-----

# Closer look at data
bank_data.head().T

bank_data.info()

# Statistical summary of our dataset
bank_data.describe().T
```



```

# iterated through bank_data and stored data with datatype as 'object'
#to new variable cat_col

cat_col = [n for n in bank_data.columns if bank_data[n].dtypes == 'object']
cat_col.remove('y')

-----

fig, axs = plt.subplots(3, 3, sharex=False, sharey=False, figsize=(20, 15))

counter = 0
for col in cat_col:
    value_counts = bank_data[col].value_counts()

    trace_x = counter // 3
    trace_y = counter % 3
    x_pos = np.arange(0, len(value_counts))

    axs[trace_x, trace_y].bar(x_pos, value_counts.values, tick_label
                             = value_counts.index)

    axs[trace_x, trace_y].set_title(col)

    for tick in axs[trace_x, trace_y].get_xticklabels():
        tick.set_rotation(90)

    counter += 1

plt.show()
#plt.savefig('dataset_cat_columns_exploration.png')

-----

num_columns = ['age', 'duration', 'campaign', 'pdays', 'previous', 'emp.var.rate',
               'cons.price.idx', 'cons.conf.idx', 'euribor3m', 'nr_employed']

fig, axs = plt.subplots(2, 3, sharex=False, sharey=False, figsize=(20, 15))

counter = 0
for num_column in num_columns:

    trace_x = counter // 3
    trace_y = counter % 3

```

```

    axs[trace_x, trace_y].hist(bank_data[num_column])

    axs[trace_x, trace_y].set_title(num_column)

    counter += 1

plt.show()
plt.savefig("dataset_num_columns_exploration.png")

-----

# column y has unbalanced data
bank_data['y'].value_counts().plot.bar()

-----

for col in cat_col:
    pd.crosstab(bank_data[col], bank_data.y).plot(kind = 'bar')
    plt.title(col)
    #plt.savefig(col + ".png")

-----

plt.figure(figsize = (10, 6))
plt.grid(b = True)
sns.distplot(a = bank_data['age'], kde = False)

-----

data_heatmap = bank_data.loc[:, bank_data.columns != 'y']
sns.heatmap(data_heatmap.corr(), cmap="YlGnBu")
plt.show()

-----

target_heatmap = bank_data.loc[:, bank_data.columns == 'y']
data_heatmap['Subscription'] = target_heatmap
fig, ax = plt.subplots(figsize=(12,12))
sns.heatmap(data_heatmap.corr()[['Subscription']], cmap='YlGnBu', annot=True);
ax.invert_yaxis()

-----

# Transform 'yes' and 'no' into (0,1) dependency
bank_data['y']=[1 if i=='yes' else 0 for i in bank_data['y']]

```

```
-----

# Create dummy variable for the categorical data and drop first column
# to avoid dummy varibale trap
bank_data = pd.get_dummies(bank_data, columns = ['job', 'marital', 'education',
                                                'default', 'housing', 'loan',
                                                'month', 'contact',
                                                'day_of_week', 'poutcome'],
                           drop_first = True)

bank_data.shape
bank_data.head().T

-----

# split dataset into data and target sets
data = bank_data.loc[:, bank_data.columns != 'y']
target = bank_data.loc[:, bank_data.columns == 'y']

data_np = np.array(data)
target_np = np.array(target)
#target_np = target_np.ravel()

print(data_np.shape)
print(target_np.shape)

-----

# split data and target sets into 4 sets for train and test purposes
bank_train_data, bank_test_data, \
bank_train_target, bank_test_target = \
train_test_split(data_np, target_np, test_size=0.1)

bank_train_data_basic = np.copy(bank_train_data)
bank_test_data_basic = np.copy(bank_test_data)

-----

# scale data
sc = StandardScaler()
bank_train_data = sc.fit_transform(bank_train_data)
bank_test_data = sc.transform(bank_test_data)

-----
```

```

# Training our model
logistic_regression = LogisticRegression()
logistic_regression.fit(bank_train_data, bank_train_target)

log_reg_basic = LogisticRegression()
log_reg_basic.fit(bank_train_data_basic, bank_train_target.ravel())

-----

# Evaluate accuracy of our model
print("Basic Logistic Regression")
accu_score = accuracy_score(bank_test_target,
log_reg_basic.predict(bank_test_data_basic))
print("Basic model accuracy = {}".format(accu_score))

-----

# Evaluate accuracy of our model
print("Basic Logistic Regression with scaled data")
accu_score = accuracy_score(bank_test_target,
logistic_regression.predict(bank_test_data))
print("Basic model accuracy of scaled data= {}".format(accu_score))

# confusion matrix
conf_matrix = confusion_matrix(bank_test_target,
logistic_regression.predict(bank_test_data))
print(conf_matrix)

sns.heatmap(conf_matrix, annot=True)
plt.show()

-----

def cross_valid_data(numOfIteration):
    print("\nCROSS walidacja:")
    scores = cross_val_score(LogisticRegression(),
data_np, target_np.ravel(), cv=numOfIteration)
    plt.figure(figsize=(1, 5))
    sns.heatmap(scores[:, np.newaxis], annot=True)
    plt.show()

def compare_given_result_with_real(left_bound, right_bound):
    for i in range(left_bound, right_bound):
        pred = logistic_regression.predict(bank_test_data[i, :].reshape(1, -1))

```

```

        print("\nModel predicted for client {0} value {1}".format(i, pred))
        print("Real value for client \"{0}\" is {1}".format(i, bank_test_target[i]))
        #Wyświetlenie
        print(logistic_regression.predict_proba(bank_test_data[i,:].reshape(1,-1)))

compare_given_result_with_real(0, 10)
cross_valid_data(10)

-----

%matplotlib inline
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

bank_data2 = pd.read_csv("bank-additional.csv", sep=';')
data2 = bank_data2.loc[:, bank_data2.columns != 'y']

sns.pairplot(data2, diag_kind = "kde")

-----

# Mean squared error and variance score of a basic logistic model
from sklearn.metrics import mean_squared_error

m_squared_error = mean_squared_error(bank_test_target,
logistic_regression.predict(bank_test_data))
print("Mean squared error = %.2f" % m_squared_error)

var_score = logistic_regression.score(bank_test_data, bank_test_target)
print("Logistic Regression variance score: %.2f" % var_score)

-----

i = 3
x = []
y_acc = []

while i <= 729:
    x.append(i)
    logistic_regression = LogisticRegression(max_iter=i)
    logistic_regression.fit(bank_train_data, bank_train_target.ravel())

```

```

        acc = accuracy_score(bank_test_target,
                              logistic_regression.predict(bank_test_data))
        y_acc.append(acc)
        i *= 3

plt.plot(x, y_acc)
plt.title("Accuracy score per iteration")
plt.grid()
plt.xlabel("Max number of iterations")
plt.ylabel("Accuracy score")

-----

for size in range (1,6):
    target_np = target_np.ravel()

    bank_train_data, bank_test_data, \
    bank_train_target, bank_test_target = \
    train_test_split(data_np, target_np, test_size=(size/10))

    log_reg = LogisticRegression()
    log_reg.fit(bank_train_data, bank_train_target)

    #Skuteczność naszego modelu
    print("Model LogisticRegression basic - Test data = {} {}".format(size*10))
    acc = accuracy_score(bank_test_target, log_reg.predict(bank_test_data))
    print("Model accuracy is {}".format(acc))
    conf_matrix = confusion_matrix(bank_test_target,
                                    log_reg.predict(bank_test_data))
    print(conf_matrix)
    sns.heatmap(conf_matrix,annot=True)
    plt.show()
    plt.savefig("dataset" + str(size*10) + ".png")

-----

solvers = ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']

for sol in solvers:
    scores = cross_val_score(LogisticRegression(solver=sol),
                              bank_train_data, bank_train_target, cv=10)
    plt.figure(figsize=(1,4))
    sns.heatmap(scores[:, np.newaxis],annot=True)

```

```

    print("Solver:" + sol)
    plt.show()

-----

from sklearn.preprocessing import PolynomialFeatures
from sklearn.feature_selection import RFE
from sklearn.metrics import mean_squared_error

bank_data2 = pd.read_csv("bank-additional.csv", sep=';')

bank_data2['y']=[1 if i=='yes' else 0 for i in bank_data2['y']]

bank_data2 = pd.get_dummies(bank_data2,columns =
['job','marital','education','default',
 'housing','loan','month', 'contact',
 'day_of_week','poutcome'], drop_first = True)

data2 = bank_data2.loc[:, bank_data2.columns != 'y']
target2 = bank_data2.loc[:, bank_data2.columns == 'y']

data2 = np.array(data2)
target2 = np.array(target2)

bank_train_data, bank_test_data, \
bank_train_target, bank_test_target = \
train_test_split(data2, target2, test_size=0.1)

pf = PolynomialFeatures(2, )

bank_train_poly = pf.fit_transform(bank_train_data)
bank_test_poly = pf.fit_transform(bank_test_data)

print("train poly size {}".format(bank_train_poly.size))
print("test poly size {}".format(bank_test_poly.size))

-----

x = [4, 16, 32, 64]
y_mean_error = []
y_r2_score = []

for feat in x:
    sel_ = RFE(estimator=LogisticRegression(solver='lbfgs',

```

```
max_iter=300), n_features_to_select=feat)
sel_.fit(bank_train_poly, bank_train_target)
y_mean_error.append(mean_squared_error(bank_test_target,
sel_.predict(bank_test_poly)))
y_r2_score.append(sel_.score(bank_test_poly, bank_test_target))

plt.plot(x, y_mean_error)
plt.title("Mean error")
plt.grid()
plt.xlabel("No. Features")
plt.ylabel("Mean error")
plt.show()

plt.plot(x,y_r2_score)
plt.title("Score")
plt.grid()
plt.xlabel("No. Features")
plt.ylabel("Score")
plt.show()
```