



Klasyfikator wspomagający diagnozę raka (Breast Cancer Wisconsin)

KWD

Michał Mamla

Dominik Pepaś

01.02.2021

Streszczenie

WYŻEJ WYMIENIONY PROJEKT MA NA CELU STWORZENIE KLASYFIKATORA REALIZUJĄCEGO PRZEWIDYWANIE DIAGNOZY DLA PACJENTA KLINIKI ONKOLOGICZNEJ. PROGRAM KOMPUTEROWY NA PODSTAWIE WCZEŚNIEJ OPRACOWANYCH DANYCH, UTWORZONYCH NA PRZYKŁAD W WYWIADZIE PROWADZONYM PRZEZ LEKARZA, TRENUJE KLASYFIKATOR KTÓRY W KOLEJNYM KROKU ZOSTAJE PODDANY TESTOM. ZADANIEM NASZYM, CZYLI PROGRAMISTÓW JEST ANALIZA WYNIKÓW PROGRAMU I WYBRANIE ORAZ OPRACOWANIE NAJBARDZIEJ WIARYGODNEGO 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	Normalizacja i standaryzacja danych	8
2.4	Analiza jakości uzyskanego modelu	8
2.5	Badania symulacyjne	10
3	Podsumowanie	17
A	Kod programu	18

Rozdział 1

Wprowadzenie

NOWOTWÓR PIERSI, CZYLI RAK JEST OBECNIE DRUGĄ GŁÓWNĄ PRZYCZYNĄ ZGONÓW NA ŚWIECIE. W NASZYM KRAJU CO ROKU Z POWODU ZACHOROWAŃ NA ZŁOŚLIWE NOWOTWORY PIERSI UMIERA OKOŁO 6 TYSIĘCY OSÓB, NA ŚWIECIE PONAD 600 TYSIĘCY. EKSPERCI ALARMUJĄ, ŻE STAŁE ZWIĘKSZA SIĘ LICZBA PACJENTÓW ONKOLOGICZNYCH. WCZESNE WYKRYCIE ZMIAN NOWOTWOROWYCH W CIELE PACJENTA MOŻE ZNACZĄCO ZWIĘKSZYĆ JEGO PROGNOZY, CZYLI SZANSE NA WYLECZENIE. NATOMIAST ZIDENTYFIKOWANIE ZMIANY JAKO ŁAGODNEJ MOŻE UCHRONIĆ PACJENTA PRZED PODJĘCIEM NIEPOTRZEBNEGO, WYCZERPUJĄCEGO LECZENIA. ZNALEZIENIE METODY WYKRYWANIA ZMIAN NOWOTWOROWYCH I KLASYFIKOWANIA ICH NA RÓŻNE KATEGORIE MOŻE ZNACZĄCO PRZYSPIESZYĆ CAŁY PROCES LECZENIA, DLATEGO JEST TO AKTUALNIE BARDZO POPULARNY ORAZ POŻĄDANY TEMAT. NARZĘDZIEM, KTÓRE SPRAWDZA SIĘ W TEGO TYPU ZADANIACH JEST UCZENIE MASZYNOWE. JEGO METODY KLASYFIKOWANIA DANYCH SĄ CZĘSTO WYKORZYSTYWANE W MEDYCYNIE, DO DIAGNOZOWANIA PRZYPADKÓW I PODEJMOWANIA DECYZJI.

1.1 Opis problemu

PODSTAWOWYM PROBLEMEM JEST STWORZENIE KLASYFIKATORA, KTÓRY PRZY KLASYFIKOWANIU DANYCH, CZYLI W TYM PRZYPADKU CECH PACJENTA KLINIKI ONKOLOGICZNEJ, BĘDZIE NIEOMYLNÝ LUB TEŻ BĘDZIE POPEŁNIAŁ MOŻLIWIE NAJMNIEJSZĄ ILOŚĆ BŁĘDÓW. POD UWAGĘ NALEŻY WZIĄĆ ODPOWIEDNI DOBÓR DANYCH NA KTÓRYCH KLASYFIKATOR NAS BĘDZIE TRENOWANY. NA ICH PODSTAWIE NASTĘPNIE BĘDZIE PODEJMOWAĆ DECYZJĘ. MODEL MUSI TAKŻE ZOSTAĆ BARDZO DOKŁADNIE PRZETESTOWANY, UWZGLĘDNIAJĄC ZMIANĘ DANYCH TESTOWYCH I TRENUJĄCYCH, TAK BY WYKLUCZYĆ JAK NAJWIĘKSZĄ ILOŚĆ POMYŁEK. UWAGĘ NALEŻY POŚWIĘCIĆ TAKŻE ROZKŁADZIE DANYCH NA TRENINGOWE I TESTUJĄCE. TWORZĄC MODEL DBAMY O TO, BY NIE BYŁ PRZEUCZONY.

Rozdział 2

Realizacja

2.1 Wprowadzenie teoretyczne

ODMIANĄ ANALIZY 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 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 NA WYSTĄPIENIE 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 MNIEJSZĄ ILOŚCIĄ ZAŁOŻEŃ. NAJWAŻNIEJSZĄ CHARAKTERYSTYKĄ 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 PODZIALE PRÓBY STATYSTYCZNEJ NA PODZBIORY, A NASTĘPNIE PRZEPROWADZENIU WSZELKICH ANALIZ NA ZBIORZE UCZĄCYM, PODCZAS GDY ZBIÓR TESTOWY SŁUŻY DO POTWIERDZENIA WIARYGODNOŚCI WYNIKÓW. METODA TA JEST STOSOWANA W PRZYPADKU NIEDOSTATECZNIE LICZNEGO ZBIORU UCZĄCEGO. POWTARZAJĄC PROCES KROSWALIDACJI K RAZY I OBLICZAJĄC STATYSTYKI NA ZBIORACH WALIDACYJNYCH TAK ZBUDOWANYCH MODELI MOŻLIWA JEST ESTYMACJA BŁĘDU I GENERALIZACJI MODELU.

2.2 Opis danych wejściowych

ZBIÓR DANYCH NA KTÓRYCH DZIAŁAĆ BĘDZIE KLASYFIKATOR SKŁADA SIĘ Z 569 REKORDÓW, Z KTÓRYCH KAŻDY OKREŚLA ODDZIELNEGO PACJENTA KLINIKI. W KAŻDEJ LINII ZNAJDUJE SIĘ 30 CECH OKREŚLAJĄCYCH KONKRETNEGO PACJENTA. NA ICH PODSTAWIE WYZNACZANA JEST ETYKIETA, KTÓRA OKREŚLA RODZAJ ZMIANY NOWOTWOROWEJ - ZŁOŚLIWEJ (0), BĄDŹ ŁAGODNEJ (1).

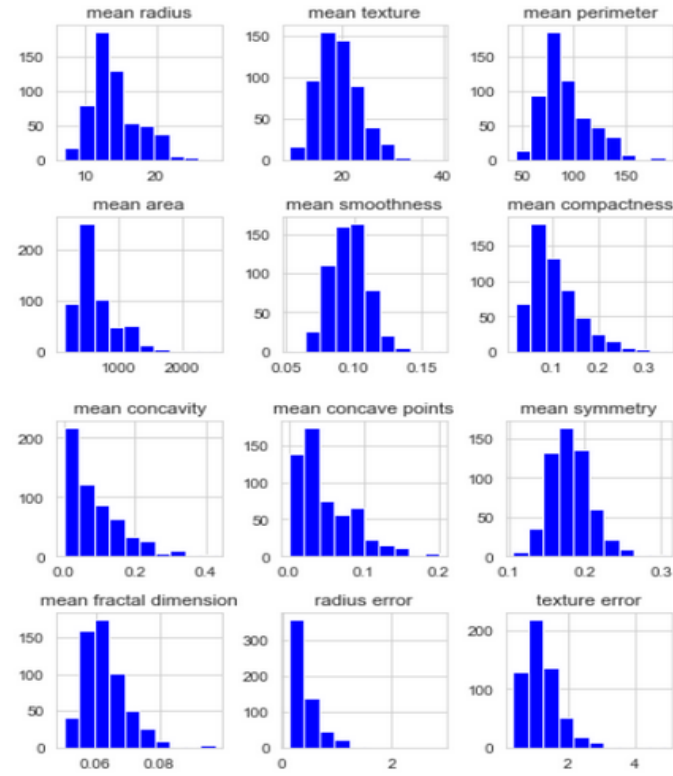
SPECYFIKACJA OKREŚLAJĄCA ZBIÓR

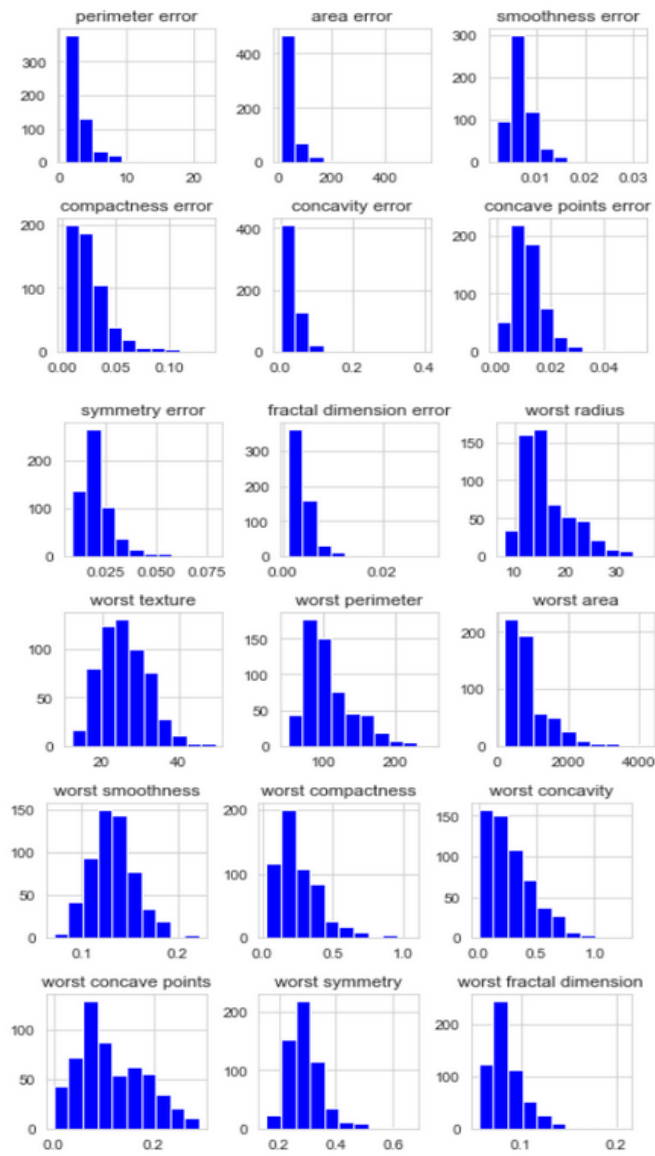
Classes	2
Samples per class	212(M),357(B)
Samples total	569
Dimensionality	30
Features	real, positive

DANE OKREŚLAJĄCE PACJENTA:

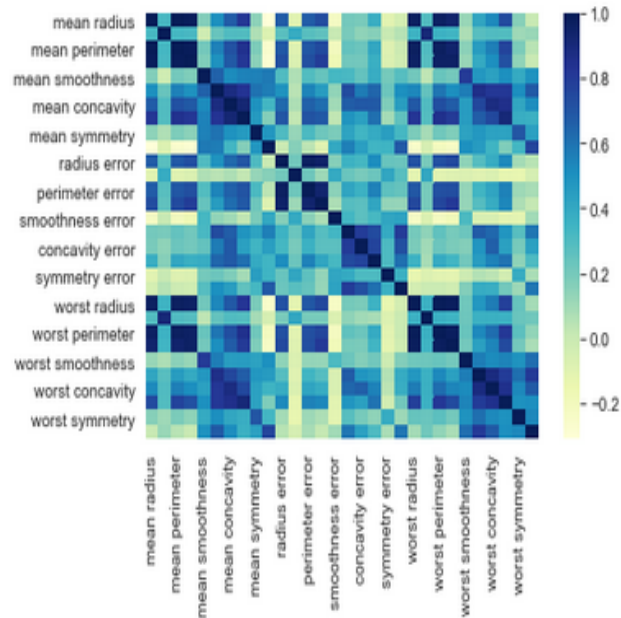
	Min	Max
radius (mean):	6.981	28.11
texture (mean):	9.71	39.28
perimeter (mean):	43.79	188.5
area (mean):	143.5	2501.0
smoothness (mean):	0.053	0.163
compactness (mean):	0.019	0.345
concavity (mean):	0.0	0.427
concave points (mean):	0.0	0.201
symmetry (mean):	0.106	0.304
fractal dimension (mean):	0.05	0.097
radius (standard error):	0.112	2.873
texture (standard error):	0.36	4.885
perimeter (standard error):	0.757	21.98
area (standard error):	6.802	542.2
smoothness (standard error):	0.002	0.031
compactness (standard error):	0.002	0.135
concavity (standard error):	0.0	0.396
concave points (standard error):	0.0	0.053
symmetry (standard error):	0.008	0.079
fractal dimension (standard error):	0.001	0.03
radius (worst):	7.93	36.04
texture (worst):	12.02	49.54
perimeter (worst):	50.41	251.2
area (worst):	185.2	4254.0
smoothness (worst):	0.071	0.223
compactness (worst):	0.027	1.058
concavity (worst):	0.0	1.252
concave points (worst):	0.0	0.291
symmetry (worst):	0.156	0.664
fractal dimension (worst):	0.055	0.208

ROZKŁAD DANYCH Z PODZIAŁEM NA CECHY PACJENTÓW

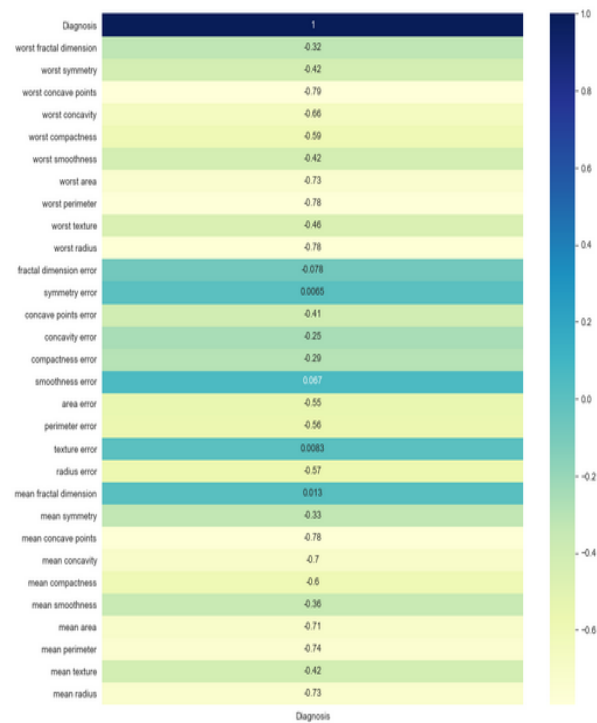




WZAJEMNA KORELACJA CECH



WZAJEMNA KORELACJA CECH Z DIAGNOZĄ



2.3 Normalizacja i standaryzacja danych

BIORĄC POD UWAGĘ ZRÓŻNICOWANIE ZAKRESU DANYCH W ZBIORZE, W OBLICZENIACH KONIECZNE JEST UWZGLĘDNIENIE NORMALIZACJI ORAZ STANDARYZACJI. WIELE KLASYFIKATORÓW OBLICZA ODLEGŁOŚĆ MIĘDZY DWO-
MA PUNKTAMI NA PODSTAWIE ODLEGŁOŚCI EUKLIDESOWEJ. JEŚLI JEDNA
Z CECH MA SZEROKI ZAKRES WARTOŚCI, ODLEGŁOŚĆ BĘDZIE ZALEŻAŁA
OD TEJ KONKRETNEJ CECHY. DLATEGO ZAKRES WSZYSTKICH CECH POWI-
NIEN ZOSTAĆ ZNORMALIZOWANY, TAK ABY KAŻDA CECHA MIAŁA UDZIAŁ
W PRZYBLIŻENIU PROPORCJONALNIE DO KOŃCOWEJ ODLEGŁOŚCI. INNYM
POWODEM STOSOWANIA SKALOWANIA CECH JEST TO, ŻE ZEJŚCIE GRA-
DIENTU ZBIEGA SIĘ ZNACZNIE SZYBCIEJ ZE SKALOWANIEM CECH NIŻ BEZ
NIEGO. STANDARYZACJA ODNOSI SIĘ DO PRZESUNIĘCIA ROZKŁADU KAŻDE-
GO ATRYBUTU W CELU UZYSKANIA ŚREDNIEJ RÓWNEJ ZERO I ODCHYLENIA
STANDARDOWEGO WYNOSZĄCEGO JEDEN.

STANDARYZACJA DANYCH W PROGRAMIE:

```
from sklearn.preprocessing import StandardScaler  
scaler = StandardScaler()  
scaled_data_cancer = scaler.fit_transform(data_cancer['data'])
```

2.4 Analiza jakości uzyskanego modelu

UZYSKANY PRZEZ NAS MODEL ZA POMOCĄ REGRESJI LOGISTYCZNEJ DAJE
NASTĘPUJĄCE WYNIKI:

```
Model accuracy is 0.96  
[[15  2]  
 [ 0 40]]
```



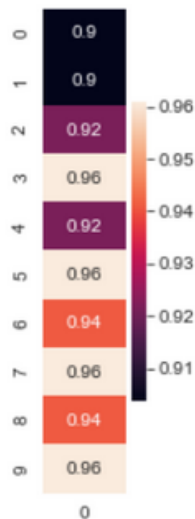
WYNIK MODELU JEST PRZEDSTAWIONY ZA POMOCĄ WSPÓŁCZYNNIKA DO-
KŁADNOŚCI. OBLICZAMY GO DZIELĄC LICZBĘ POPRAWNYCH DOPASOWAŃ
PRZEZ WSZYSTKIE PRÓBY DOPASOWANIA. W NASZYM PRZYPADKU DO-
KŁADNOŚĆ WYNIOSŁA 96 %, CO JEST DOSYĆ ZADOWALAJĄCYM WYNIKIEM.
MACIERZ PRZEDSTAWIA ILOŚĆ POPRAWNYCH I NIEPOPRAWNYCH DOPA-
SOWAŃ. W PIERWSZYM WIERSZU ZNAJDUJĄ SIĘ INFORMACJE O 15 PO-
PRAWNYCH DOPASOWANIACH DIAGNOZY ZMIANY ZŁOŚLIWEJ (0), ORAZ O
2 POMYŁKACH, GDZIE MODEL PRZYPISAŁ ZMIANIE ZŁOŚLIWEJ ETYKIETĘ
ZMIANY ŁAGODNEJ. W DRUGIM WIERSZU ZNAJDUJĄ SIĘ INFORMACJĘ O 40
POPRAWNYCH DOPASOWANIACH DIAGNOZY ZMIANY ŁAGODNEJ (1), ORAZ
O BRAKU POMYŁEK. PODCZAS TWORZENIA MODELU 10% DANYCH BYŁO
ZBIOREM TESTUJĄCYM, NATOMIAST 90% ZBIOREM TRENUJĄCYM.

```
#Rozdzielenie danych na zbiór trenujący i testujący (90% do 10%)
cancer_train_data, cancer_test_data, \
cancer_train_target, cancer_test_target = \
train_test_split(scaled_data_cancer, data_cancer.target, test_size=0.1)
```

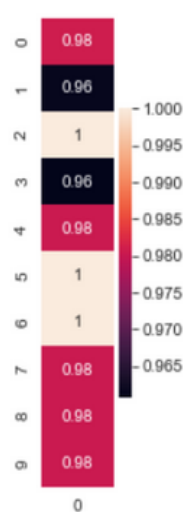
DLA NASZYCH DANYCH ZOSTAŁA RÓWNIEŻ PRZEPROWADZONA SPRAWDZE-
NIE KRZYŻOWE, CZYLI KROSWALIDACJA, KTÓREJ WYNIKI SĄ ZBLIŻONE DO
MODELU WYTRENOWANEGO PRZEZ NAS.

WYNIKI PRZED SKALOWANIEM DANYCH I PO KROSWALIDACJI WYGLĄDAJĄ
NASTĘPUJĄCO:

PRZED:



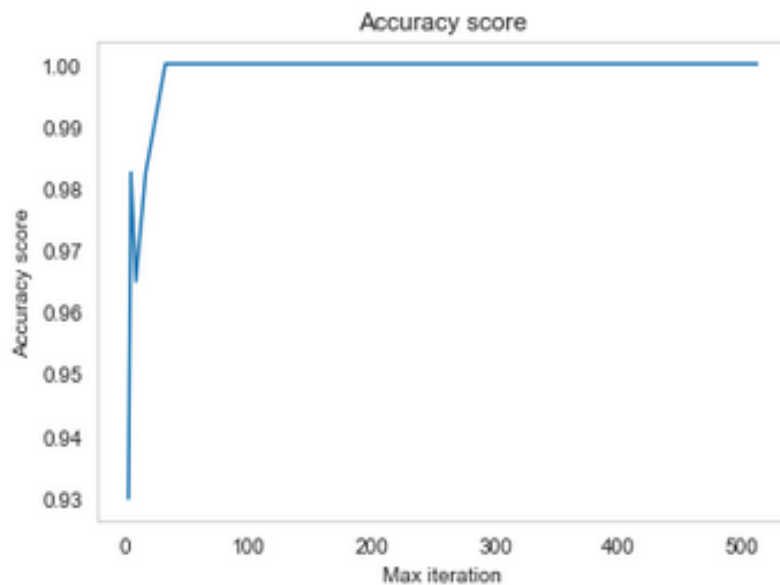
PO:



2.5 Badania symulacyjne

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

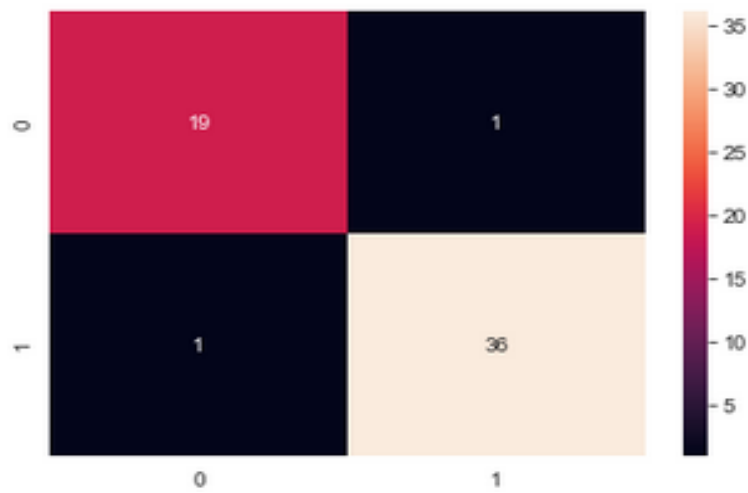
WYKRES PRZEDSTAWIAJĄCY DOKŁADNOŚĆ MODELU W ZALEŻNOŚCI OD MAKSYMALNEJ LICZBY ITERACJI:



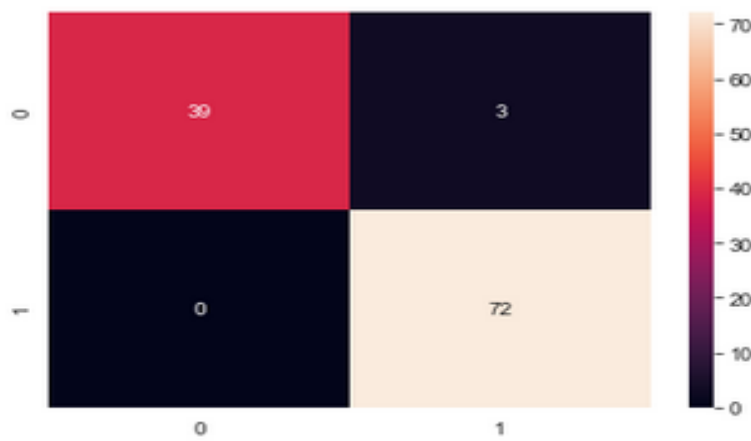
MAKSYMALNE ITERACJE SĄ POTĘGAMI LICZBY 2 DO 512 WŁĄCZNIE.

PO PRZEANALIZOWANIU WYKRESU MOŻNA WYWNIOSKOWAĆ, ŻE DOKŁADNOŚĆ WZRASTA WRAZ ZE WZROSTEM LICZBY ITERACJI, AŻ OSIĄGA STU PROCENTOWĄ SKUTECZNOŚĆ OKOŁO 32 MAKSYMALNYCH ITERACJI. PRZEPROWADZONE ZOSTAŁY RÓWNIEŻ BADANIA WPŁYWU ROZKŁADU DANYCH NA ZBIORY TESTUJĄCE I TRENUJĄCE:

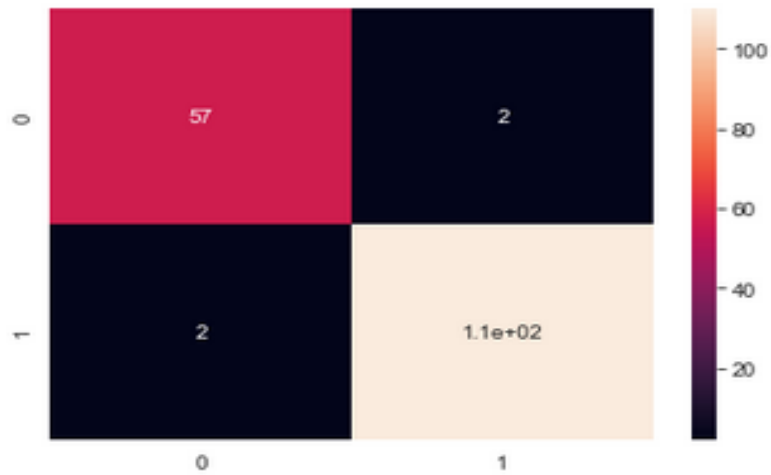
```
Model LogisticRegression basic - Test data = 10 %  
Model accuracy is 0.96  
[[19  1]  
 [ 1 36]]
```



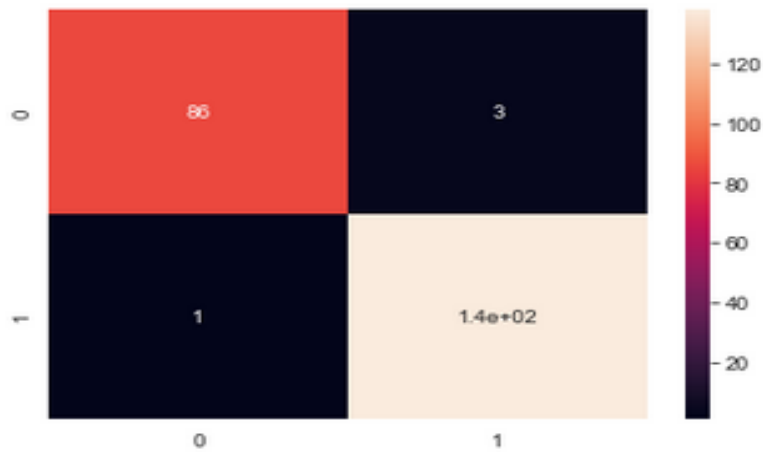
```
Model LogisticRegression basic - Test data = 20 %  
Model accuracy is 0.97  
[[39  3]  
 [ 0 72]]
```



```
Model LogisticRegression basic - Test data = 30 %  
Model accuracy is 0.98  
[[ 57  2]  
 [  2 110]]
```



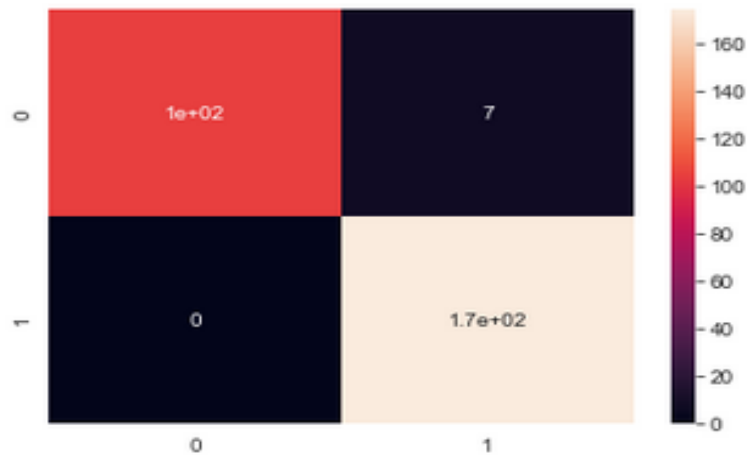
```
Model LogisticRegression basic - Test data = 40 %  
Model accuracy is 0.98  
[[ 86  3]  
 [  1 138]]
```



```

Model LogisticRegression basic - Test data = 50 %
Model accuracy is 0.98
[[104  7]
 [  0 174]]

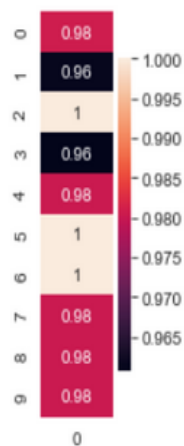
```



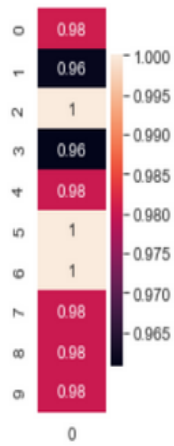
ZWIĘKSZAJĄC ZBIÓR DANYCH TESTUJĄCYCH WYNIK DLA NASZEGO MODELU POZOSTAJE BEZ ZNACZĄCYCH ZMIAN.

SPRAWDZONY ZOSTAŁ TAKŻE WPŁYW RODZAJU ALGORYTMU OPTYMALIZACYJNEGO NA REZULTAT DZIAŁANIA MODELU:

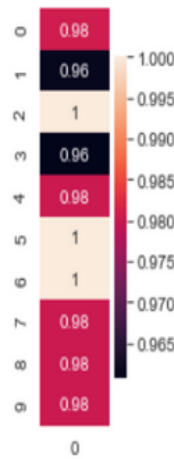
solver:lbfgs



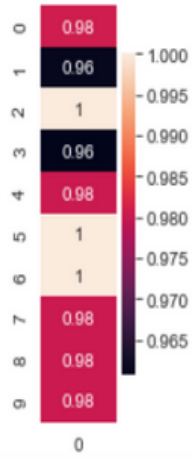
Solver:liblinear



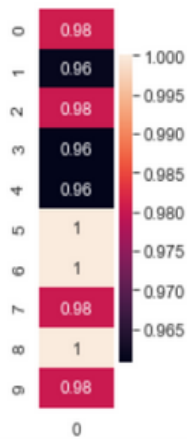
Solver:newton-cg



Solver:sag



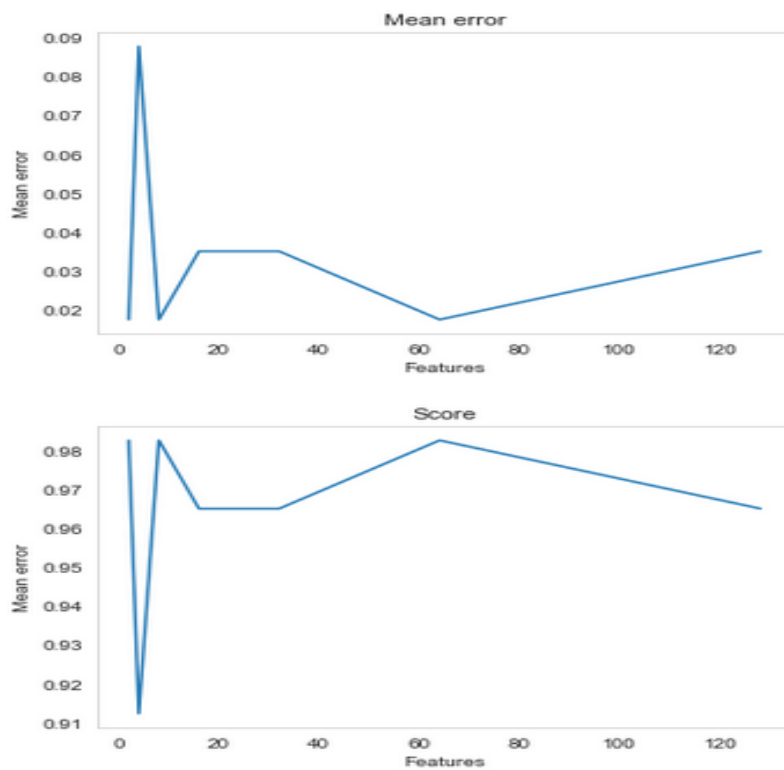
Solver:saga



ZMIANA ALGORYTMU OPTYMALIZACYJNEGO NIE WPŁYNĘŁA ZNACZĄCO NA WYNIK MODELU.

ZESTAW DANYCH ZOSTAŁ RÓWNIEŻ PODDANY TESTOM REKURENCYJNEJ ELIMINACJI CECH. JEST TO METODA OPAKOWANIA, POLEGAJĄCA NA SZKOLENIU ESTYMATORA NA POCZĄTKOWYM ZESTAWIE ATRYBUTÓW, A NASTĘPNIE USUWANIU NAJMNIEJ WAŻNYCH ATRYBUTÓW. PROCEDURA TA JEST POWTARZANA REKURENCYJNIE, AŻ DO MOMENTU OSIĄGNIĘCIA POTRZEBNEJ NAM ILOŚCI ATRYBUTÓW DO WYBORU. NA WYKRESACH PRZEDSTAWIONA JEST ZALEŻNOŚĆ ŚREDNIEGO BŁĘDU I WYNIKU OD LICZBY WYBRANYCH ATRYBUTÓW. W OBU PRZYPADKACH MOŻNA ZAUWAŻYC SPADEK JAKOŚCI W PRZYPADKU WYBRANIA WIĘCEJ NIŻ 60 ATRYBUTÓW.

Model LogisticRegression with PF and RFE:



Rozdział 3

Podsumowanie

WYTRENOWANY KLASYFIKATOR BARDZO POPRAWNIE RADZI SOBIE Z KLASYFIKOWANIEM ANALIZOWANYCH PRZEZ NAS DANYCH. SKALOWANIE ZBIORU DANYCH ZNACZĄCO POLEPSZYŁO WYNIK MODELU. OPTYMALNĄ ILOŚCIĄ MAKSYMALNYCH ITERACJI ALGORYTMU OPTYMALIZUJĄCEGO PODCZAS BUDOWANIA NASZEGO MODELU JEST 100. RODZAJ ALGORYTMU OPTYMALIZUJĄCEGO NIE MA WIĘKSZEGO ZNACZENIA NA SKUTECZNOŚĆ KLASYFIKATORA. MIMO ZWIĘKSZANIA ZBIORU TESTUJĄCEGO I ZMNIEJSZANIA TRENUJĄCEGO LINIOWA REGRESJA RADZI SOBIE BARDZO DOBRZE Z DANymi KLINIKI. REKURENCYJNA ELIMINACJA CECH NIE JEST KONIECZNA W TRENOWANIU NASZEGO MODELU. NALEŻY PAMIĘTAĆ O TYM, ŻE NAWET JEŚLI SKUTECZNOŚĆ MODELU SIĘGA STU PROCENT, W DAŁSZYM CIĄGU DIAGNOZA STAWIANA PRZEZ ALGORYTM UCZENIA MASZYNOWEGO POWINNA STANOWIĆ TYLKO JEDEN Z KROKÓW W OCENIE PACJENTA.

Dodatek A

Kod programu

```
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import cross_val_score
from sklearn.preprocessing import StandardScaler
import numpy as np
import matplotlib.pyplot as plt

#Cross walidacja dla naszych danych
def cross_valid_data(iteration):
    print("\nCross walidacja:")
    scores = cross_val_score(LogisticRegression(), cancer_train_data,
                              cancer_train_target, cv=iteration)
    plt.figure(figsize=(1, 5))
    sns.heatmap(scores[:, np.newaxis],annot=True)
    plt.show()

#Przetestowanie naszego modelu dla określonych przez nas elementów
#- porównanie wyniku modelu z rzeczywistym wynikiem
def compare_model_resault_to_real_value(first_item, last_item):
    for i in range(first_item, last_item):

        prediction = logistic_regression.predict(cancer_test_data[i,:].reshape(1,-1))
        print("\nModel predicted for client {0} value {1}".format(i, prediction))
        print("Real value for client \"{0}\" is {1}".format(i, cancer_test_target[i]))
        #Wyświetlenie
        print(logistic_regression.predict_proba(cancer_test_data[i,:].reshape(1,-1)))
```

```
#Załadowanie danych
data_cancer = load_breast_cancer()

#Skalowanie daty
scaler = StandardScaler()
scaled_data_cancer = scaler.fit_transform(data_cancer['data'])

#Rozdzielenie danych na zbiór trenujący i testujący (90\% do 10\%)
cancer_train_data, cancer_test_data, \
cancer_train_target, cancer_test_target = \
train_test_split(scaled_data_cancer, data_cancer.target, test_size=0.1)

#Informacje o rozmiarze podzielonych danych
print("Training dataset:")
print("Train_data:", cancer_train_data.shape)
print("Train_target:", cancer_train_target.shape)

print("Testing dataset:")
print("Test_data:", cancer_test_data.shape)
print("Test_target:", cancer_test_target.shape)

#Trenowanie modelu z pomocą regresji logistycznej
logistic_regression = LogisticRegression()
logistic_regression.fit(cancer_train_data, cancer_train_target)

#Skuteczność naszego modelu
print("Model LogisticRegression basic:")
acc = accuracy_score(cancer_test_target, logistic_regression \
.predict(cancer_test_data))
print("Model accuracy is {0:0.2f}".format(acc))
conf_matrix = confusion_matrix(cancer_test_target, logistic_regression \
.predict(cancer_test_data))
print(conf_matrix)

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

compare_model_resault_to_real_value(0,10)
cross_valid_data(10)

-----

from sklearn.datasets import load_breast_cancer
import pandas as pd
```

```

import seaborn as sns
import matplotlib.pyplot as plt

data_cancer = load_breast_cancer()
pd_data = pd.DataFrame(data_cancer['data'], columns= data_cancer.feature_names)

fig = plt.figure()

for i,b in enumerate(list(bcdf.columns[0:30])):
    i +=1
    ax = fig.add_subplot(5,6,i)
    ax.hist(pd_data[b], label = 'Value', stacked = True, color= 'b')
    ax.set_title(b)

plt.tight_layout()
plt.show()

-----

\%matplotlib inline

import matplotlib.pyplot as plt
import seaborn as sb

sb.pairplot(pd_data, diag_kind="kde")

-----

# POLY FEATURES AND RFE
from sklearn.preprocessing import PolynomialFeatures
from sklearn.feature_selection import RFE
from sklearn.metrics import mean_squared_error

x = []
y_mean = []
y_score = []

def PolynomialFeatures_for_LogisticRegression(degree):
    print("\nModel LogisticRegression with PF and RFE:")

    pt = PolynomialFeatures(degree, )

    cancer_train_poly = pt.fit_transform(cancer_train_data)
    cancer_test_poly = pt.fit_transform(cancer_test_data)

```

```

#logistic_regression_poly = LogisticRegression(solver='lbfgs', max_iter=300)
#logistic_regression_poly.fit(credit_clients_train_poly, cancer_train_target)

i=2

while i <= 128:
    sel_ = RFE(estimator=LogisticRegression(solver='lbfgs', max_iter=300), \
        n_features_to_select=i)
    sel_.fit(cancer_train_poly, cancer_train_target)
    x.append(i);
    y_mean.append(mean_squared_error(cancer_test_target, \
        sel_.predict(cancer_test_poly)))
    y_score.append(sel_.score(cancer_test_poly, cancer_test_target))
    i*=2

#acc = accuracy_score(cancer_test_target, sel_.predict(credit_clients_test_poly))
#print("Model accuracy is {0:0.2f}".format(acc))

#conf_matrix = confusion_matrix(cancer_test_target
#, sel_.predict(credit_clients_test_poly))
#print(conf_matrix)

PolynomialFeatures_for_LogisticRegression(2)

plt.plot(x,y_mean)
plt.title("Mean error")
plt.grid()
plt.xlabel("Features")
plt.ylabel("Mean error")
plt.show()

plt.plot(x,y_score)
plt.title("Score")
plt.grid()
plt.xlabel("Features")
plt.ylabel("Mean error")
plt.show()

-----

\%matplotlib inline

```

```

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

for sol in tab:
    scores = cross_val_score(LogisticRegression(solver=sol), \
        cancer_train_data, cancer_train_target, cv=10)
    plt.figure(figsize=(1,4))
    sns.heatmap(scores[:, np.newaxis], annot=True)
    print("Solver:" + sol)
    plt.show()

-----

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

pd_data = pd.DataFrame(data_cancer['data'], columns= data_cancer.feature_names)

sns.heatmap(pd_data.corr(), cmap="YlGnBu")
plt.show()

-----

pd_data['Diagnosis'] = data_cancer.target
fig, ax = plt.subplots(figsize=(12,12))
sns.heatmap(pd_data.corr()[['Diagnosis']], cmap='YlGnBu', annot=True);
ax.invert_yaxis()

-----

for size in range (1,6):
    cancer_train_data, cancer_test_data, \
    cancer_train_target, cancer_test_target = \
    train_test_split(scaled_data_cancer, data_cancer.target, test_size=(size/10))
    #Trenowanie modelu z pomocą regresji logistycznej
    logistic_regression = LogisticRegression()
    logistic_regression.fit(cancer_train_data, cancer_train_target)

    #Skuteczność naszego modelu
    print("Model LogisticRegression basic - Test data = {} \%".format(size*10))
    acc = accuracy_score(cancer_test_target, \
        logistic_regression.predict(cancer_test_data))
    print("Model accuracy is {0:0.2f}".format(acc))

```



```
conf_matrix = confusion_matrix(cancer_test_target, \
    logistic_regression.predict(cancer_test_data))
print(conf_matrix)

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

-----

i = 2
x = []
y_acc = []

while i <= 512:
    x.append(i)
    logistic_regression = LogisticRegression(max_iter=i)
    logistic_regression.fit(cancer_train_data, cancer_train_target)
    acc = accuracy_score(cancer_test_target, \
        logistic_regression.predict(cancer_test_data))
    y_acc.append(acc)
    i *= 2

plt.plot(x, y_acc)
plt.title("Accuracy score")
plt.grid()
plt.xlabel("Max iteration")
plt.ylabel("Accuracy score")
```