

## Zad 1

Zaimplementowano kod SGD:

```
1 # Linear Regression With Stochastic Gradient Descent for  
2 Wine Quality  
3 from random import seed  
4 from random import randrange  
5 from csv import reader  
6 from math import sqrt  
7  
8 # Load a CSV file  
9 def load_csv(filename):  
10     dataset = list()  
11     with open(filename, 'r') as file:  
12         csv_reader = reader(file)  
13         for row in csv_reader:  
14             if not row:  
15                 continue  
16             dataset.append(row)  
17     return dataset  
18  
19 # Convert string column to float  
20 def str_column_to_float(dataset, column):  
21     for row in dataset:  
22         row[column] = float(row[column].strip())  
23  
24 # Find the min and max values for each column  
25 def dataset_minmax(dataset):  
26     minmax = list()  
27     for i in range(len(dataset[0])):  
28         col_values = [row[i] for row in dataset]  
29         value_min = min(col_values)  
30         value_max = max(col_values)  
31         minmax.append([value_min, value_max])  
32     return minmax  
33  
34 # Rescale dataset columns to the range 0-1  
35 def normalize_dataset(dataset, minmax):  
36     for row in dataset:  
37         for i in range(len(row)):  
38             row[i] = (row[i] - minmax[i][0]) / (minmax[i][1] - minmax[i][0])
```

```

39 # Split a dataset into k folds
40 def cross_validation_split(dataset, n_folds):
41     dataset_split = list()
42     dataset_copy = list(dataset)
43     fold_size = int(len(dataset) / n_folds)
44     for i in range(n_folds):
45         fold = list()
46         while len(fold) < fold_size:
47             index = randrange(len(dataset_copy))
48             fold.append(dataset_copy.pop(index))
49         dataset_split.append(fold)
50     return dataset_split
51
52 # Calculate root mean squared error
53 def rmse_metric(actual, predicted):
54     sum_error = 0.0
55     for i in range(len(actual)):
56         prediction_error = predicted[i] - actual[i]
57         sum_error += (prediction_error ** 2)
58     mean_error = sum_error / float(len(actual))
59     return sqrt(mean_error)
60
61 # Evaluate an algorithm using a cross validation split
62 def evaluate_algorithm(dataset, algorithm, n_folds, *args
63 ):
64     folds = cross_validation_split(dataset, n_folds)
65     scores = list()
66     for fold in folds:
67         train_set = list(folds)
68         train_set.remove(fold)
69         train_set = sum(train_set, [])
70         test_set = list()
71         for row in fold:
72             row_copy = list(row)
73             test_set.append(row_copy)
74             row_copy[-1] = None
75         predicted = algorithm(train_set, test_set, *args)
76         actual = [row[-1] for row in fold]
77         rmse = rmse_metric(actual, predicted)
78         scores.append(rmse)
79     return scores
80
81 # Make a prediction with coefficients
82 def predict(row, coefficients):
83     yhat = coefficients[0]
84     for i in range(len(row)-1):
85         yhat += coefficients[i + 1] * row[i]
86     return yhat

```

```

86
87 # Estimate linear regression coefficients using
   stochastic gradient descent
88 def coefficients_sgd(train, l_rate, n_epoch):
89     coef = [0.0 for i in range(len(train[0]))]
90     for epoch in range(n_epoch):
91         for row in train:
92             yhat = predict(row, coef)
93             error = yhat - row[-1]
94             coef[0] = coef[0] - l_rate * error
95             for i in range(len(row)-1):
96                 coef[i + 1] = coef[i + 1] - l_rate *
error * row[i]
97                 # print(l_rate, n_epoch, error)
98     return coef
99
100 # Linear Regression Algorithm With Stochastic Gradient
   Descent
101 def linear_regression_sgd(train, test, l_rate, n_epoch):
102     predictions = list()
103     coef = coefficients_sgd(train, l_rate, n_epoch)
104     for row in test:
105         yhat = predict(row, coef)
106         predictions.append(yhat)
107     return(predictions)
108
109 # Linear Regression on wine quality dataset
110 seed(1)
111 # load and prepare data
112 filename = 'winequality-white.csv'
113 dataset = load_csv(filename)
114 for i in range(len(dataset[0])):
115     str_column_to_float(dataset, i)
116
117
118 # normalize
119 minmax = dataset_minmax(dataset)
120 normalize_dataset(dataset, minmax)
121 # evaluate algorithm
122 n_folds = 5
123
124 # stała uczenia
125 l_rate = 0.01
126

```

```

127 # epoki
128 n_epoch = 500
129 scores = evaluate_algorithm(dataset,
    linear_regression_sgd, n_folds, l_rate, n_epoch)
130 print('Scores: %s' % scores)
131 print('Mean RMSE: %.3f' % (sum(scores)/float(len(scores)
    ))))

```

Przeprowadzono testy na różnych wartościach epok i szybkości uczenia się. Wyniki przedstawiono w Tabeli 1.

*Tabela 1: Wyniki testów*

Lp	Wartość epoki	Szybkość uczenia się	RMSE
1	50	0,01	0,126
2	50	0,1	0,131
3	100	0,1	0,131
4	100	0,01	0,126
5	500	0,01	0,126

Analizując Tabelę możemy zauważyć, że nie udało się zejść z wynikiem RMSE poniżej 0.126. Oznacza, że algorytmy są do siebie zbliżone i zwiększenie wartości epoki czy szybkości uczenia się, nie wpłynie na poprawę uzyskiwanych wyników. W celu uzyskania lepszych wyników RMSE prawdopodobnie należałoby zmienić algorytm na inny.