

### Лабораторна робота 3 - Логістична регресія.

У цій роботі Ви побудуєте логістичну регресію для класифікації зображень рукописних символів за датасетом MNIST.

Реалізуйте методи з позначкою #TODO класу LogisticRegression:

Метод `preprocess` повинен додавати колонку з одиниць у матрицю  $X$ .  
Опціонально – додайте поліноміальні або будь-які інші нелінійні ознаки.

Метод `onehot` повинен виконувати onehot-перетворення:

$$\text{onehot} : \mathbb{R} \rightarrow \mathbb{R}^c$$

$$\overline{\text{onehot}(y_i)}_j = \begin{cases} 1, & j = y_i \\ 0, & j \neq y_i \end{cases}$$

де  $c$  – кількість класів. Метод `h` - гіпотеза:

$$h(X) = \sigma(X\theta)$$

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

Метод `J` повинен обчислювати оціночну функцію логістичної регресії:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^c (-y_{i,j} \cdot \log(h(x_i)_j) - (1 - y_{i,j}) \cdot \log(1 - h(x_i)_j)) + \alpha_1 \sum_{i=1}^N \sum_{j=1}^c |\hat{\theta}_{i,j}| + \alpha_2 \sum_{i=1}^N \sum_{j=1}^c \hat{\theta}_{i,j}^2$$

Метод `grad` має обчислювати градієнт  $\frac{\partial J}{\partial \theta}$ :

$$\frac{\partial J}{\partial \theta} = -\frac{1}{m} X^T (Y - h(X)) + \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \times (\alpha_1 \text{sign}(\theta) + 2\alpha_2 \theta)$$

Метод `moments` має повертати вектор-рядки  $\mu, \sigma$  для середнього і стандартного відхилення кожної колонки. Пам'ятайте, що колонку з одиницями не потрібно нормалізувати, тож відповідні середнє і стандартне відхилення для неї вкажіть рівними 0 і 1 відповідно. Можна використовувати функції `np.mean` і `np.std`.

Метод `normalize` має виконувати нормалізацію  $X$  на основі статистик  $\mu, \sigma$ , що повернув метод `moments`. Для того щоб уникнути ділення на 0, можете до  $\sigma$  додати маленьку величину, наприклад  $10^{-8}$ .

Метод `get_batch` має повертати матриці  $X_b, Y_b$  з довільно обраних  $b$  елементів вибірки ( $b$  у коді - `self.batch_size`).

Метод `fit` виконує оптимізацію  $J(\theta)$ . Для кращої збіжності реалізуйте алгоритм оптимізації **Momentum**:

$$v_t = \gamma v_{t-1} + \alpha \nabla J(\theta_{t-1})$$

$$\theta_t = \theta_{t-1} - v_t$$

де  $\gamma$  встановіть рівним 0.9 (можете поекспериментувати з іншими величинами),  $v_1 = [0]_{N,c}$ .

## Код класу LogisticRegression:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.ticker as ticker
4 from sklearn.metrics import confusion_matrix
5 import pandas as pd
6 import random
7 from sklearn.utils import shuffle
8 import os
9 from datetime import datetime
10 from datetime import timedelta
11 class LogisticRegression:
12     def __init__(self,
13         items=[],
14         # alpha1
15         # alpha2,
16         # learning_rate,
17         # batch_size,
18         # train_steps
19     ):
20         self.l = items
21         # self.alpha1 = alpha1
22         # self.alpha2 = alpha2
23         # self.learning_rate = learning_rate
24         # self.batch_size = batch_size
25         # self.train_steps = train_steps
26
27     def onehot(self, y):
28         OnehotY=np.zeros((len(y),10))
29         for i in range(len(y)):
30             OnehotY[i][y[i]]=1
31         return OnehotY
32     def preprocess(self, x, poly_deg): # TODO
33         PolynomialX = [x]
34         for degree in range(2, poly_deg + 1):
35             PolynomialX.append(x ** degree)
36         newx_poly = np.concatenate((PolynomialX[0], np.ones((x.shape[0], 1))), axis = 1)
37         newx_poly[:, :1:] = 1
38         return newx_poly
39         # previous preprocess (it works slowly)
40         """CountOfRows=len(x)
41         CountOfColumns=len(x[0])
42         NewX=np.zeros((CountOfRows,CountOfColumns*poly_deg+1))
43         for i in range(CountOfRows):
44             NewX[i][0]=1
45             l=1
46             for j in range(CountOfColumns):
47                 NewX[i][j+1]=x[i][j]
48                 k=CountOfColumns+1
49                 for j in range(2,poly_deg+1):
50                     for l in range(CountOfColumns):
51                         NewX[i][k]=x[i][l]**j
52                         k+=1
53         return NewX"""
54
55     def normalize(self, x): # TODO
56         # Z-масштабування даних на основі середнього значення та стандартного відхилення:
57         # ділення різниці між змінною та середнім значенням на стандартне відхилення.
58         CountOfRows=len(x)
59         CountOfColumns=len(x[0])
60         VerySmallNumber=pow(10,-8)
61         for i in range(CountOfRows):
62             for j in range(CountOfColumns):
63                 x[i][j]=((x[i][j]-self.mu[j]))/(self.sigma[j]+VerySmallNumber)
64         return x
65     def moments(self, x): # TODO
66         CountOfRows=len(x)
67         CountOfColumns=len(x[0])
68         MeanDeviations=[0]
69         StandardDeviations=[1]
70         for i in range(1,CountOfColumns):
71             column=[]
72             for j in range(CountOfRows):
73                 column.append(x[j][i])
74             MeanDeviations.append(np.mean(column,axis=0))
75             StandardDeviations.append(np.std(column,axis=0))
76         return [MeanDeviations,StandardDeviations]
```

```

74 def get_batch(self, x, y, batch_size): # TODO
75     RandomIndexes=np.random.randint(len(x),size=batch_size)
76     return x[RandomIndexes],y[RandomIndexes]
77     # previous get_batch 2 (it works faster)
78     """XSize=len(x)
79     RandomIndexes=np.array([i for i in range(XSize)])
80     random.shuffle(RandomIndexes)
81     return np.array([x[RandomIndexes[i]] for i in
82     range(batch_size)],np.array([y[RandomIndexes[i]] for i in range(batch_size)]))"""
83     # previous get_batch 1 (it works slowly)
84     """XSize=len(x)
85     YSize=len(y)
86     XBatch=np.zeros((batch_size,len(x[0])))
87     YBatch=np.zeros((batch_size,len(y[0])))
88     RandomIndexes=[]
89     for i in range(XSize):
90         RandomIndexes.append(i);
91     for i in range(XSize):
92         j=random.randrange(0,XSize)
93         t=RandomIndexes[i]
94         RandomIndexes[i]=RandomIndexes[j]
95         RandomIndexes[j]=t
96     for i in range(batch_size):
97         XBatch[i]=x[RandomIndexes[i]]
98         YBatch[i]=y[RandomIndexes[i]]
99     return [XBatch,YBatch]"""
100 def PrepareX(self,x,poly_deg):
101     x = self.preprocess(x,poly_deg)
102     self.mu, self.sigma = self.moments(x)
103     return self.normalize(x)
104
105 def sigmoid(self,x):
106     return 1/(1+np.exp(-x))
107 def h(self, x, theta): # TODO
108     return self.sigmoid(x@theta)
109 def grad(self, x, y, theta, alpha1, alpha2): # TODO
110     return ((1/len(x)*-1)*x.T@(y-self.h(x,theta)))+(alpha1*np.sign(theta)+2*alpha2*theta)
111 def fit(self, x, y, alpha1, alpha2, learning_rate, batch_size, train_steps):
112     (m, n), (_, c) = x.shape, y.shape
113     theta = np.zeros(shape=(n, c))
114     gamma = 0.9
115     v_1 = np.zeros(shape=(n, c)) #TODO
116     v_t = v_1
117     for step in range(train_steps):
118         x_batch, y_batch = self.get_batch(x, y, batch_size)
119         theta_grad = self.grad(x_batch, y_batch, theta, alpha1, alpha2)
120         # TODO Update v_t and theta
121         v_t = gamma * v_t + learning_rate * theta_grad
122         theta = theta - v_t
123     self.theta = theta
124     return self
125 def predict(self, x):
126     x = self.preprocess(x,1)
127     x = self.normalize(x)
128     return self.h(x, self.theta).argmax(axis=1)
129 def score(self, x, y):
130     y_pred = self.predict(x)
131     return (y == y_pred).mean() * 100

```

## Опис проведених досліджень

Запустивши код із стандартними параметрами ми отримали наступний вивід:

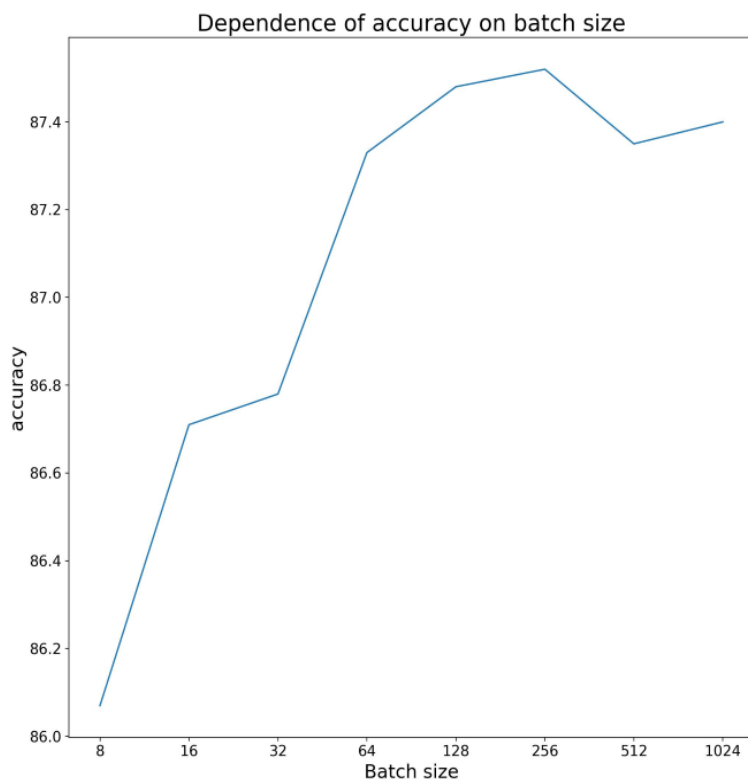
```
Processing... accuracy: 87.19000% ; time: 0 : 0 : 52
```

У рамках даної лабораторної роботи ми також виконали численні експерименти з різними значеннями таких основних параметрів моделі як коефіцієнти alpha1 та alpha2, коефіцієнт швидкості навчання (learning rate), розмір частини вибірки (batch\_size), кількість кроків навчання (train\_steps), будували також

поліноміальну логістичну регресію, аналізуючи параметр `poly_deg`. Вони підбирались так, щоб побачити межу між недостатнім навчанням та перенавчанням.

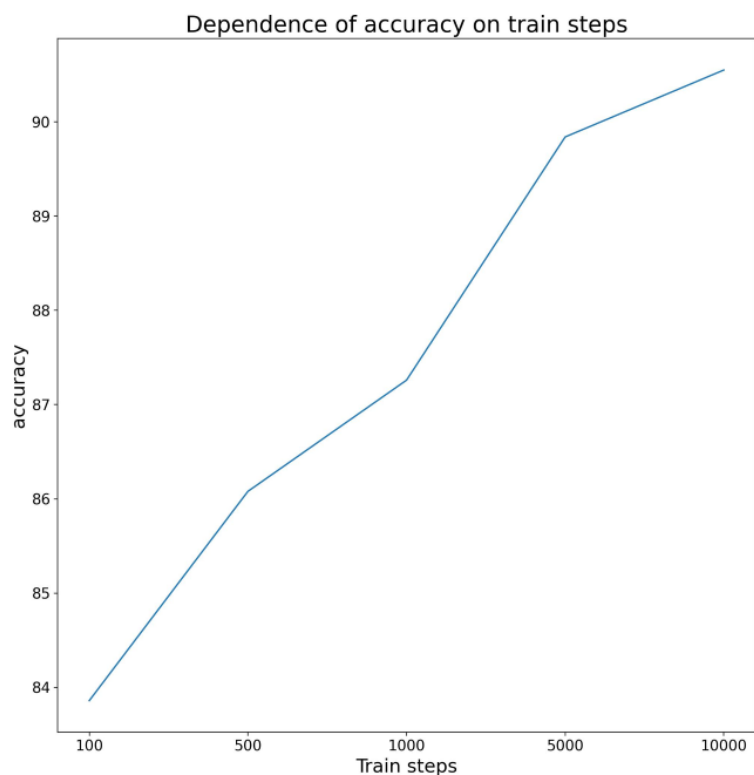
Ми створили декілька версій методів `preprocess` та `get_batch`, задіявши під час реальних випробувань їхні найбільш швидкодіючі версії. Також слід відмітити зміну конструктора та виокремлення деяких підготовчих інструкцій у процедуру `PrepareX` з метою оптимізації часових витрат на попередню підготовку даних, яка виконується лише один раз перед певним набором експериментів. Окрім тестів з різними параметрами програма також містить код для визначення часових витрат, фіксації результатів у вигляді графіків, діаграм різних типів та текстового csv-файлу, який при необхідності може бути проаналізований табличним процесором.

Тестування створеного рішення проводилось на ПК з процесором Intel Core i5-6600 без використання окремого GPU та хмарних сервісів. Тому деякі набори експериментів виконувались у послідовності за спаданням обчислювальної складності, використовуючи в кожній наступній серії експериментів оптимальні параметри з попередніх. Спочатку ми проаналізували залежність `score` від різних розмірів частини вибірки (`batch_size`):

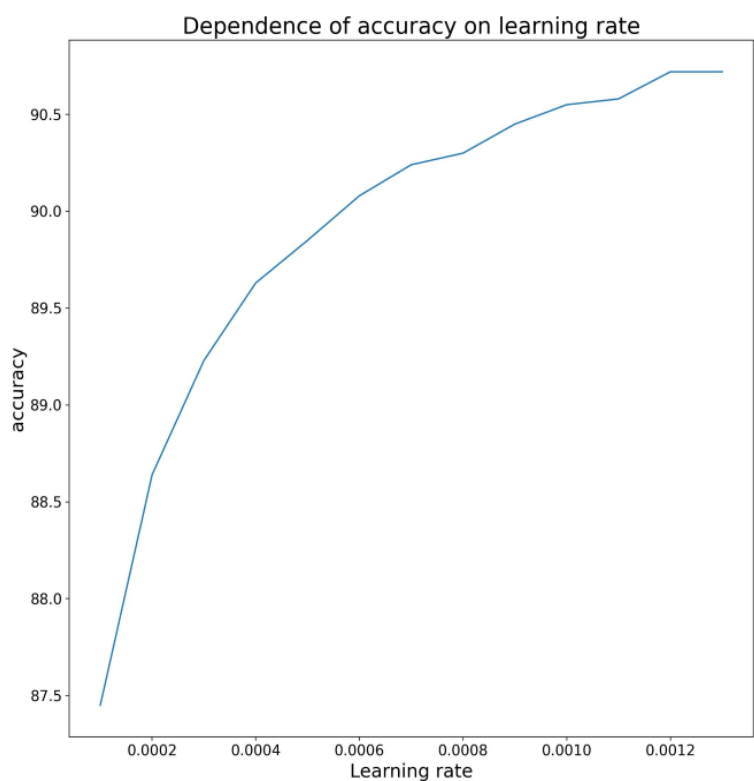


Можна зробити висновок, що найкращим значенням є 256, проте при подальших експериментах, можливо, знайдеться більш кращий розмір.

Після цього, вже з використанням кращого `batch_size`, був проаналізований параметр `train_steps` (кількість кроків навчання). Звісно, чим він більший, тим точнішою є модель, проте визначення можливої межі початку перенавчання потребує значних обчислювальних потужностей або великих часових витрат:



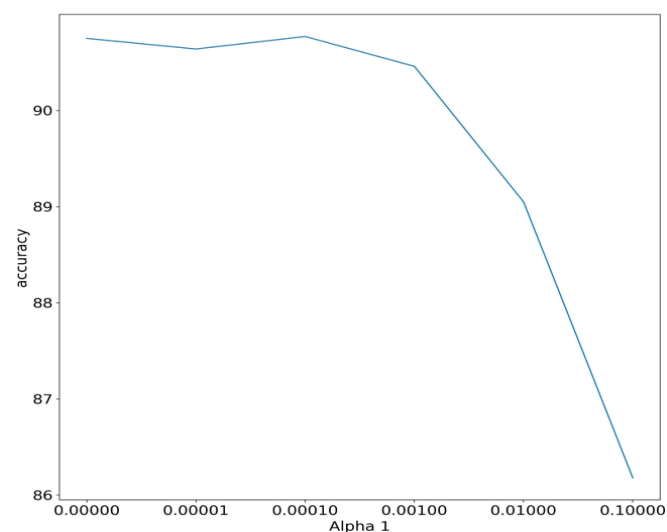
Після цього були досліджені різні значення `learning_rate` (коефіцієнт швидкості навчання):



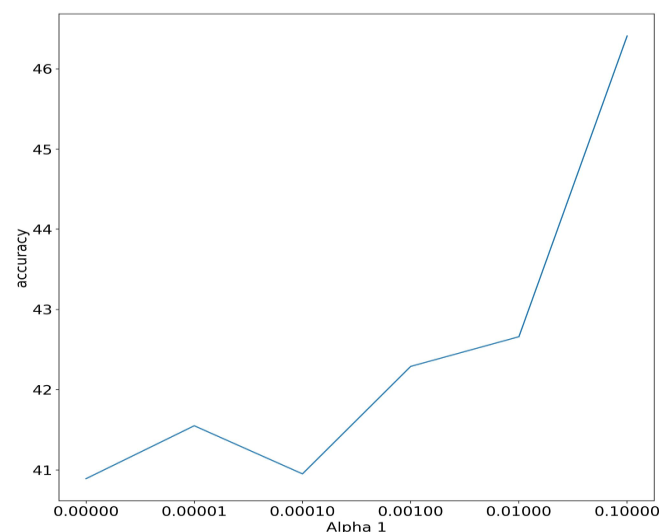
Найкращим значенням із досліджених виявилось 0.0012. Можна стверджувати, що його збільшення призводить до покращення моделі, проте визначення межі, де починається стрімке перенавчання, потребує більшої кількості додаткових експериментів.

Після цього експериментальним шляхом було встановлено, що коефіцієнт поліноміальності (poly\_deg) більший ніж 5 призводить до перенавчання, тобто до збільшення score. Ми намагались підібрати його у різних комбінаціях разом із параметрами alpha1 та alpha2, отримавши схожі один на одного графіки та діаграми:

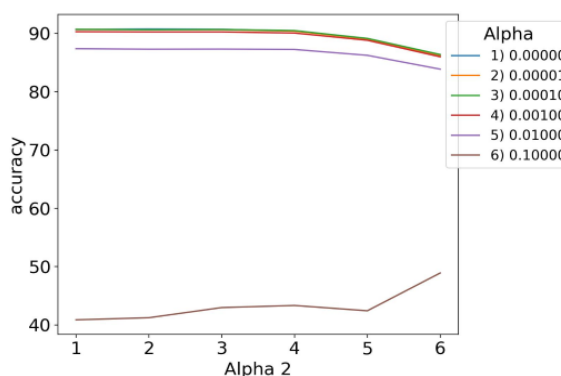
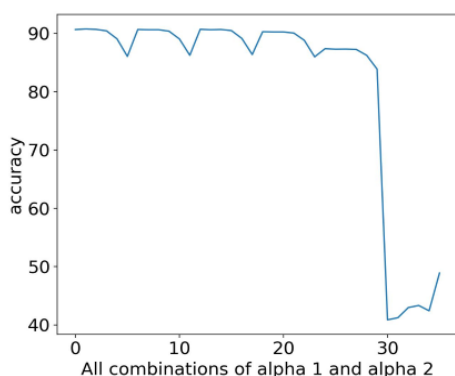
Polydeg: 1; alpha 1: 0.00010



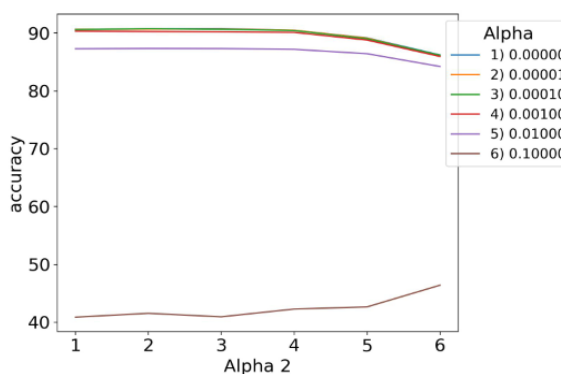
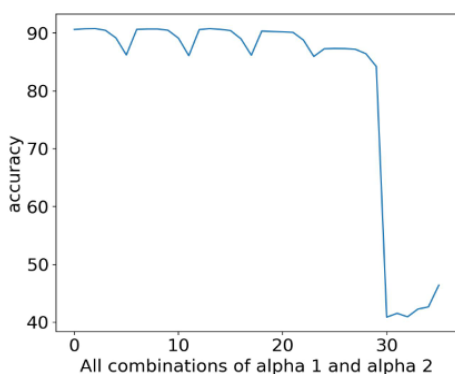
Polydeg: 8; alpha 1: 0.10000

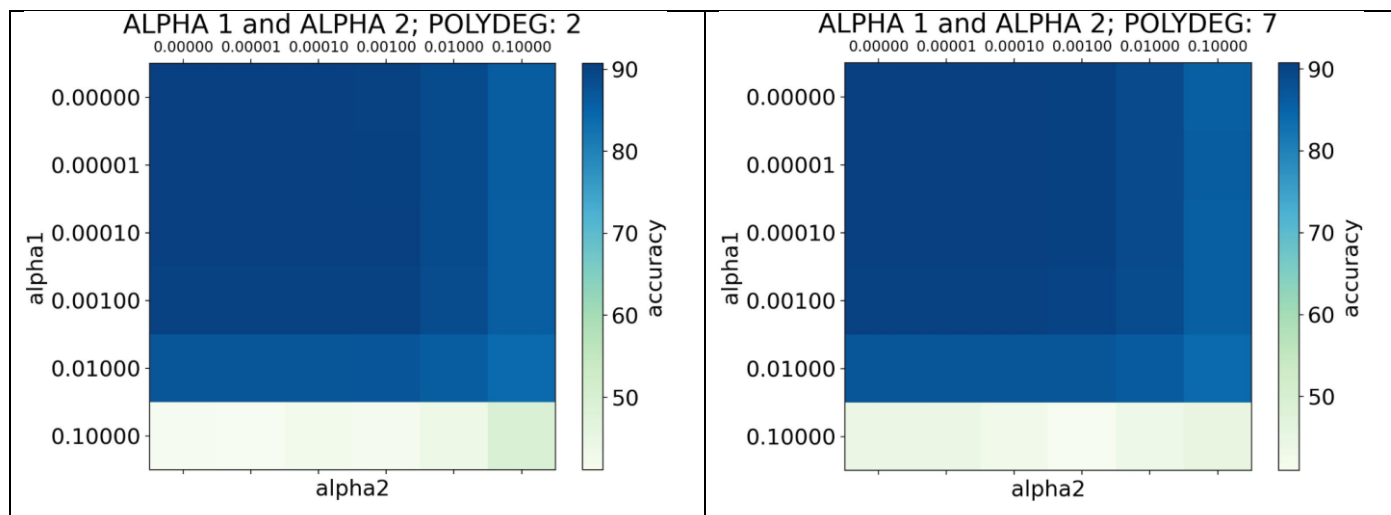


Poly deg: 3



Poly deg: 8



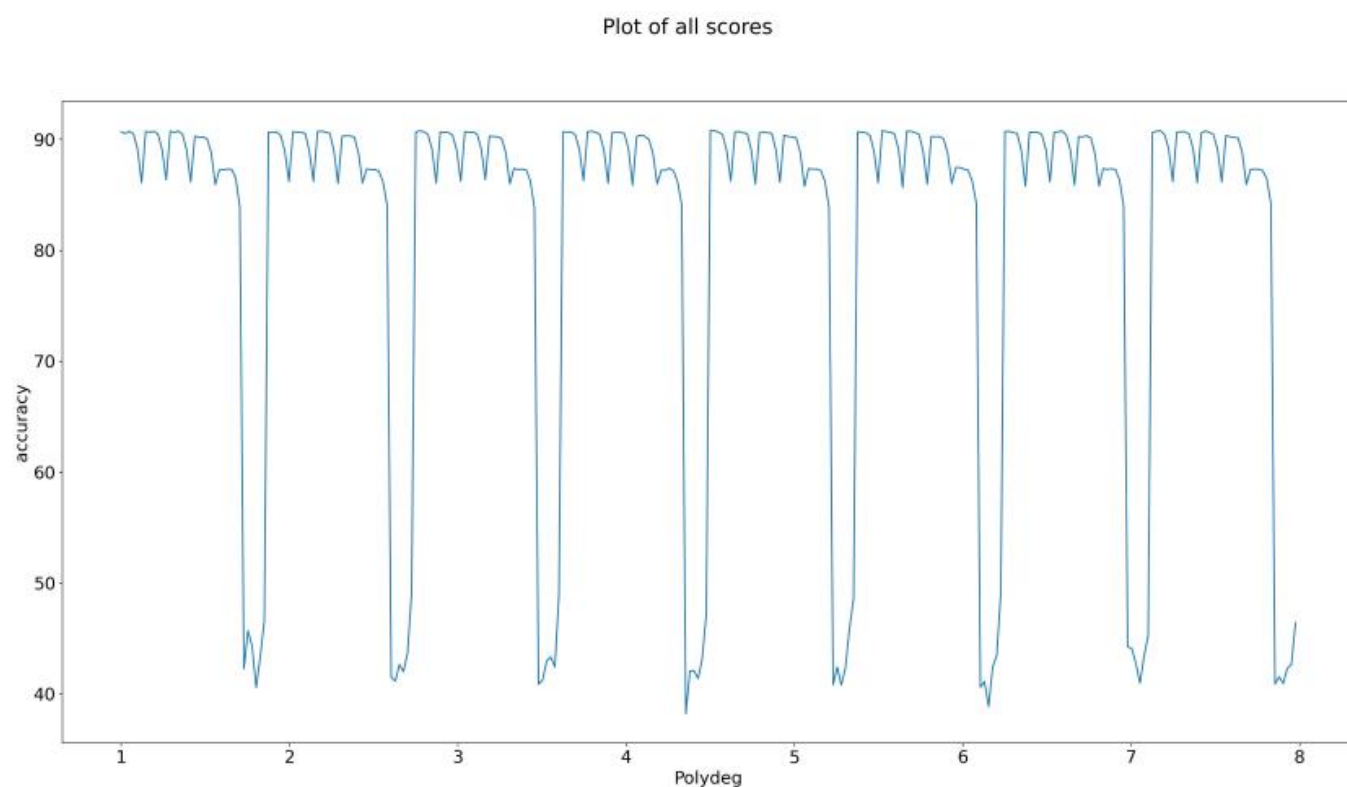


Можна зробити висновок, що оптимальні значення  $\alpha_1$  та  $\alpha_2$  знаходяться в межах до 0.01, вище – стрімке зменшення якості моделі, проте найкращим значенням  $\alpha_1$  та  $\alpha_2$  виявився 0,  $\text{poly\_deg} = 5$ . Це свідчить про негативний вплив коефіцієнтів регуляризації на якість моделі, а поліноміальність покращує точність лише в певній мірі.

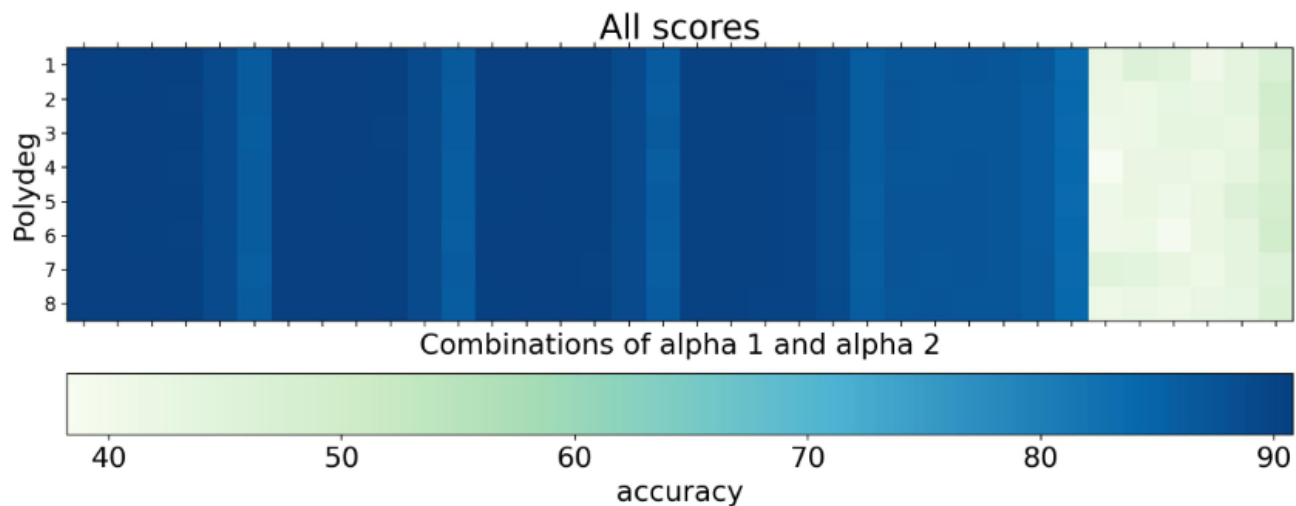
Найгірше значення accuracy 38.21000% було отримано при таких параметрах:

poly_deg	alpha1	alpha2	learning_rate	batch_size	train_steps
4	0.10000	0.00000	0.0010	8	100

Наведемо глобальний графік та діаграму типу `matshow` для всіх комбінацій коефіцієнтів (результат пошуку параметрів  $\alpha_1$  та  $\alpha_2$  разом з  $\text{poly\_deg}$ ):







Після завершення роботи програми ми отримуємо наступний вивід на екран:

```

286 / 288 ; alpha1: 0.10000 ; alpha2: 0.00100... Time: 0 : 0 : 20 ; accuracy: 42.2900
287 / 288 ; alpha1: 0.10000 ; alpha2: 0.01000... Time: 0 : 0 : 20 ; accuracy: 42.6600
288 / 288 ; alpha1: 0.10000 ; alpha2: 0.10000... Time: 0 : 0 : 20 ; accuracy: 46.4100
-----
Max accuracy: 90.82000% ; polydeg: 5 ; alpha 1: 0.00000 ; alpha 2: 0.00000
learning rate: 0.0012 ; batch size: 256 ; train steps: 10000
-----
Min accuracy: 38.21000% ; polydeg: 4 ; alpha 1: 0.10000 ; alpha 2: 0.00000
learning rate: 0.0010 ; batch size: 8 ; train steps: 100
-----
Best model... Time: 0 : 1 : 16 ; Test accuracy: 90.65000%
Total time of all calculations: 1 : 53 : 24
Drawing plots... Time: 0 : 0 : 12
  
```

Візуалізація матриці помилок для кращої моделі:

