

Лабораторна робота 2 - Лінійна та поліноміальна регресія.

Завдання: одне з безлічі завдань, яким займається сучасна фізика - це пошук матеріалу для виготовлення надпровідника, що працює за кімнатної температури. Крім теоретичних методів є і підхід з боку статистики, який передбачає аналіз бази даних матеріалів для знаходження залежності критичної температури від інших фізичних характеристик. Саме цим Ви і займетеся.

У файлі data.csv міститься весь датасет. Разом маємо 21 тисячу рядків і 169 колонок, з яких перші 167 - ознаки, колонка critical_temp містить величину, яку треба передбачити. Колонка material - містить хімічну формулу матеріалу, її можна відкинути.

Реалізуйте методи з позначкою #TODO класу PolynomialRegression:

Метод preprocess повинен виконувати таке перетворення:

$$X = [x_{i,j}]_{m \times n}$$
$$\begin{aligned} & \text{preprocess}(X) \\ = & \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,n} & x_{1,1}^2 & \dots & x_{1,n}^2 & \dots & x_{1,1}^p & \dots & x_{1,n}^p \\ 1 & x_{2,1} & \dots & x_{2,n} & x_{2,1}^2 & \dots & x_{2,n}^2 & \dots & x_{2,1}^p & \dots & x_{2,n}^p \\ \vdots & & & & & & & & & & \\ 1 & x_{m,1} & \dots & x_{m,n} & x_{m,1}^2 & \dots & x_{m,n}^2 & \dots & x_{m,1}^p & \dots & x_{m,n}^p \end{bmatrix}_{m,N} \end{aligned}$$

де p - ступінь полінома (self.poly_deg у коді). Таким чином, preprocess додає поліноміальні ознаки до X . Метод J має обчислювати оціночну функцію регресії:

$$J(\theta) = MSE(Y, h_{\theta}(X)) + \alpha_1 \sum_{i=1}^N \sum_{j=1}^k |\hat{\theta}_{i,j}| + \alpha_2 \sum_{i=1}^N \sum_{j=1}^k \hat{\theta}_{i,j}^2$$

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

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

Метод `moments` має повертати вектор-рядки μ, σ 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` и `np.std`.

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

Метод `get_batch` повинен 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 γ встановіть рівним 0.9 (можете поекспериментувати з іншими величинами), $v_1 = [0]_{N,k}$.

Код класу PolynomialRegression:

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3  import matplotlib.ticker as ticker
4  import seaborn as sns
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 PolynomialRegression:
12     def __init__(
13         self,
14         # alpha1,
15         # alpha2,
16         poly_deg,
17         # learning_rate,
18         # batch_size,
19         # train_steps
20     ):
21         # self.alpha1 = alpha1
22         # self.alpha2 = alpha2
23         self.poly_deg = poly_deg
24         # self.learning_rate = learning_rate
25         # self.batch_size = batch_size
26         # self.train_steps = train_steps
```

```

27 def preprocess(self, x): # TODO
28     PolynomialX = [x]
29     for degree in range(2, self.poly_deg + 1):
30         PolynomialX.append(x ** degree)
31     newx_poly = np.concatenate((PolynomialX[0], np.ones((x.shape[0], 1))), axis = 1)
32     newx_poly[:, :1:] = 1
33     return newx_poly
34     # previous preprocess (it works slowly)
35     """CountOfRows=len(x)
36     CountOfColumns=len(x[0])
37     NewX=np.zeros((CountOfRows,CountOfColumns*self.poly_deg+1))
38     for i in range(CountOfRows):
39         NewX[i][0]=1
40         l=1
41         for j in range(CountOfColumns):
42             NewX[i][j+1]=x[i][j]
43             k=CountOfColumns+1
44             for j in range(2,self.poly_deg+1):
45                 for l in range(CountOfColumns):
46                     NewX[i][k]=x[i][l]**j
47                     k+=1
48     return NewX"""

49 def normalize(self, x): # TODO
50     # 2-масштабування даних на основі середнього значення та стандартного відхилення:
51     # ділення різниці між змінною та середнім значенням на стандартне відхилення.
52     CountOfRows=len(x)
53     CountOfColumns=len(x[0])
54     VerySmallNumber=pow(10,-8)
55     for i in range(CountOfRows):
56         for j in range(CountOfColumns):
57             x[i][j]=(x[i][j]-self.mu[j])/(self.sigma[j]+VerySmallNumber)
58     return x
59 def moments(self, x): # TODO
60     CountOfRows=len(x)
61     CountOfColumns=len(x[0])
62     MeanDeviations=[0]
63     StandardDeviations=[1]
64     for i in range(1,CountOfColumns):
65         column=[]
66         for j in range(CountOfRows):
67             column.append(x[j][i])
68         MeanDeviations.append(np.mean(column,axis=0))
69         StandardDeviations.append(np.std(column,axis=0))
70     return [MeanDeviations,StandardDeviations]

71 def get_batch(self, x, y, batch_size): # TODO
72     RandomIndexes=np.random.randint(len(x),size=batch_size)
73     return x[RandomIndexes],y[RandomIndexes]
74     # previous get_batch 2 (it works faster)
75     """XSize=len(x)
76     RandomIndexes=np.array([i for i in range(XSize)])
77     random.shuffle(RandomIndexes)
78     return np.array([x[RandomIndexes[i]] for i in
79     range(batch_size)],np.array([y[RandomIndexes[i]] for i in range(batch_size)]))"""
80     # previous get_batch 1 (it works slowly)
81     """XSize=len(x)
82     YSize=len(y)
83     XBatch=np.zeros((self.batch_size,len(x[0])))
84     YBatch=np.zeros((self.batch_size,len(y[0])))
85     RandomIndexes=[]
86     for i in range(XSize):
87         RandomIndexes.append(i);
88     for i in range(XSize):
89         j=random.randrange(0,XSize)
90         t=RandomIndexes[i]
91         RandomIndexes[i]=RandomIndexes[j]
92         RandomIndexes[j]=t
93     for i in range(self.batch_size):
94         XBatch[i]=x[RandomIndexes[i]]
95         YBatch[i]=y[RandomIndexes[i]]
96     return [XBatch,YBatch]"""

97 def grad(self, x, y, theta, alpha1, alpha2): # TODO
98     return ((2/len(x)*-1)*x.T@(y-(x@theta)))+(alpha1*np.sign(theta)+2*alpha2*theta)
99 def J(self, x, y, theta, alpha1, alpha2): # TODO
100     prev_y=x@theta
    mse=np.square(np.subtract(y,prev_y)).mean()
    return mse+self.alpha1*sum(sum(abs(theta)))+self.alpha2*sum(sum(theta**2))

```

```

101 def PrepareX(self,x,batch_size):
102     x = self.preprocess(x)
103     self.mu, self.sigma = self.moments(x)
104     x = self.normalize(x)
105     try:
106         assert np.allclose(x[:, 1:].mean(axis=0), 0, atol=1e-3)
107         assert np.all((np.abs(x[:, 1:].std(axis=0)) < 1e-2) | (np.abs(x[:, 1:].std(axis=0)
108             - 1) < 1e-2))
109     except AssertionError as e:
110         print('Something wrong with normalization')
111         raise e
112     x_batch, y_batch = self.get_batch(x, y, batch_size)
113     try:
114         assert x_batch.shape[0] == batch_size
115         assert y_batch.shape[0] == batch_size
116     except AssertionError as e:
117         print('Something wrong with get_batch')
118         raise e
119     return x
120
121 def fit(self, x, y, alpha1, alpha2, learning_rate, batch_size, train_steps):
122     (m, N), (_, k) = x.shape, y.shape
123     theta = np.zeros(shape=(N, k))
124     gamma = 0.9
125     v_1 = np.zeros(shape=(N, k)) # TODO
126     v_t = v_1
127     for step in range(train_steps):
128         x_batch, y_batch = self.get_batch(x, y, batch_size)
129         theta_grad = self.grad(x_batch, y_batch, theta, alpha1, alpha2)
130         # TODO Update v_t and theta
131         v_t = gamma * v_t + learning_rate * theta_grad
132         theta = theta - v_t
133     self.theta = theta
134     return self
135
136 def predict(self, x):
137     x = self.preprocess(x)
138     x = self.normalize(x)
139     return x @ self.theta
140
141 def score(self, x, y):
142     y_pred = self.predict(x)
143     return np.abs(y - y_pred).mean()

```

Опис проведених досліджень

Запустивши код із стандартними параметрами ми отримали наступний вивід:

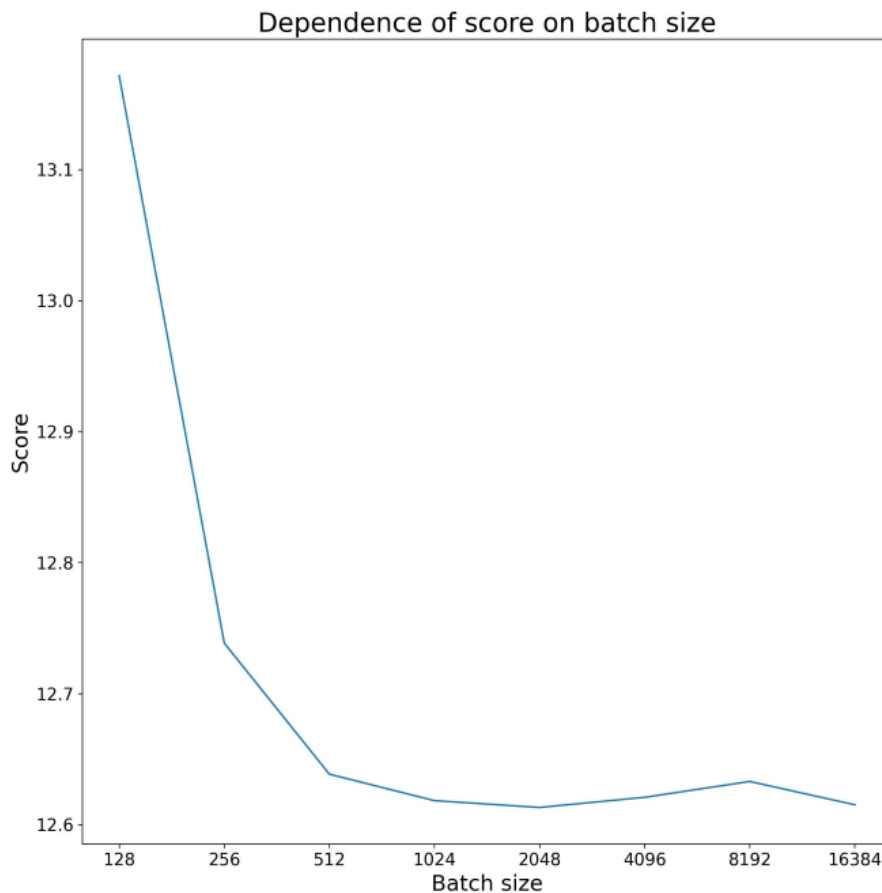
```
Processing... Score: 12.63077 ; time: 0 : 0 : 4
```

У рамках даної лабораторної роботи ми також виконали численні експерименти з різними значеннями таких основних параметрів моделі як ступінь поліноміальності (poly_deg), коефіцієнти alpha1 та alpha2, коефіцієнт швидкості навчання (learning rate), розмір частини вибірки (batch_size), кількість кроків навчання (train_steps). 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,

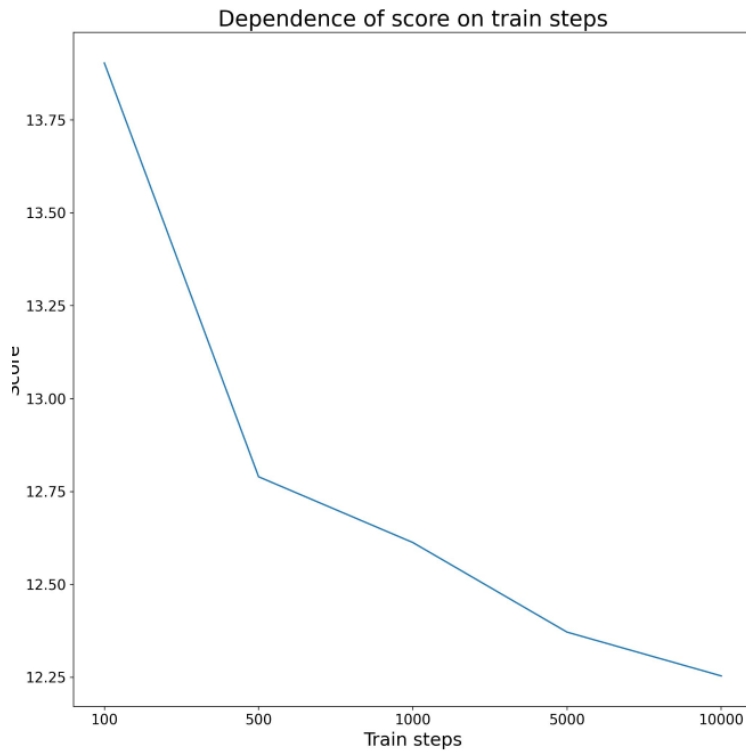
фіксації результатів у вигляді графіків, діаграм різних типів та текстового csv-файлу, який при необхідності може бути проаналізований табличним процесором.

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):

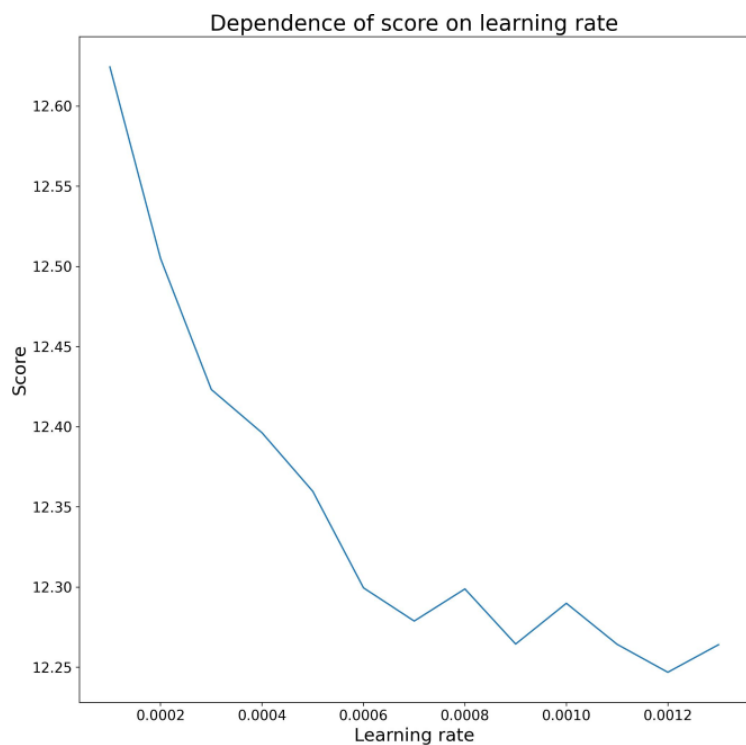


Можна зробити висновок, що найкращим значенням є 2048, проте при подальших експериментах, можливо, знайдеться більш кращий розмір.

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:



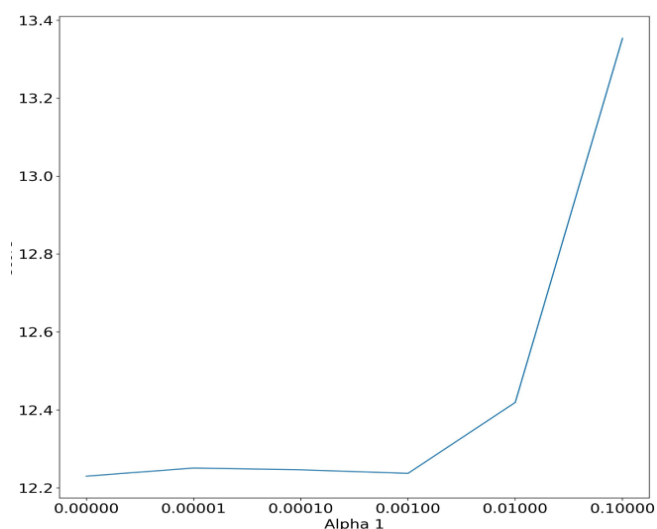
Після цього були досліджені різні значення `learning_rate` (коефіцієнт швидкості навчання):



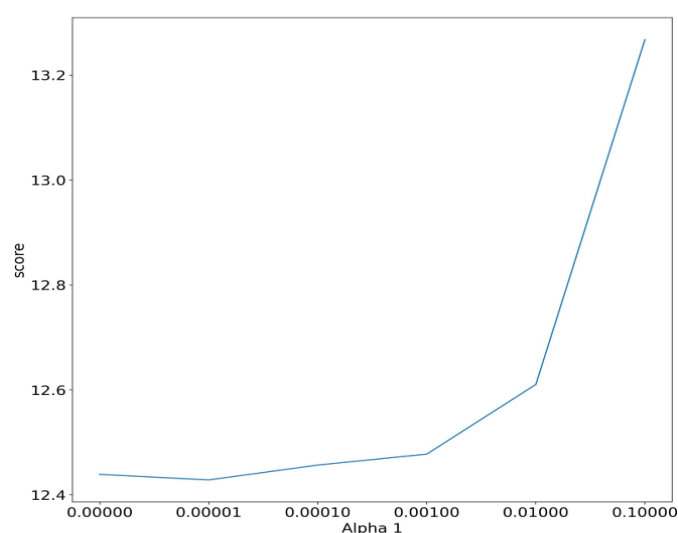
Найкращим значенням із досліджених виявилось 0.0012. Можна стверджувати, що його збільшення призводить до покращення моделі, проте визначення межі, де починається стрімке перенавчання, потребує більшої кількості додаткових експериментів.

Після цього експериментальним шляхом було встановлено, що коефіцієнт поліноміальності (poly_deg) більший ніж 6 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:

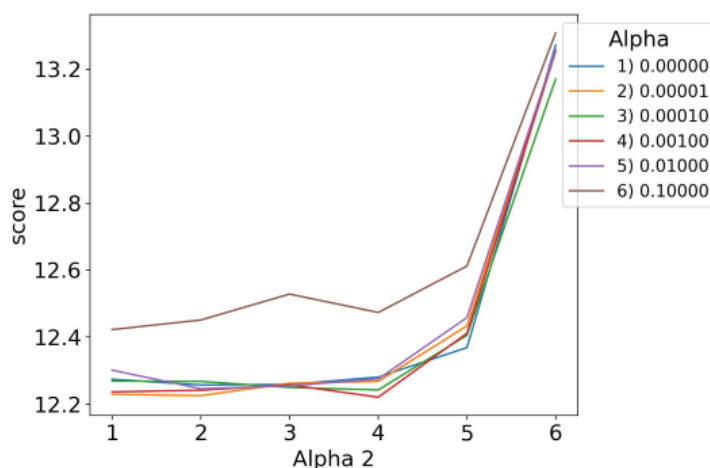
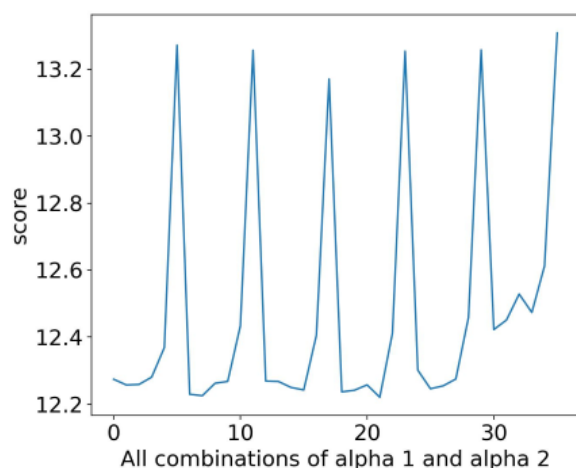
Polydeg: 1; alpha 1: 0.00100



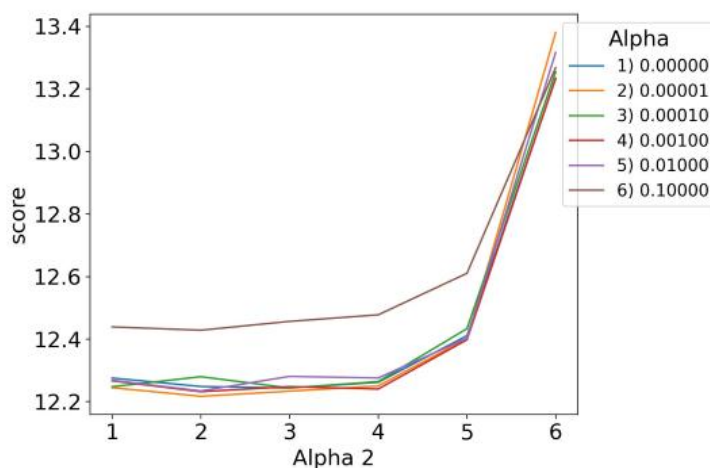
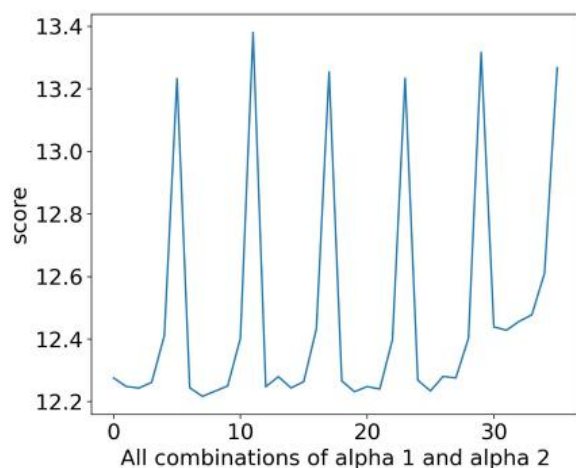
Polydeg: 8; alpha 1: 0.10000

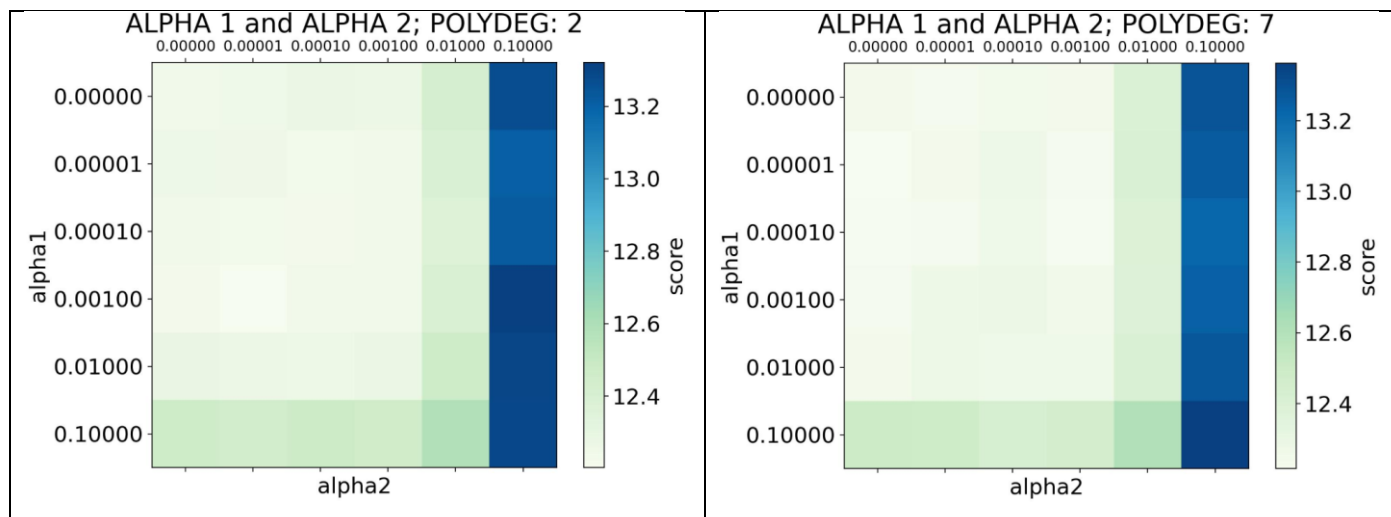


Poly deg: 3



Poly deg: 8





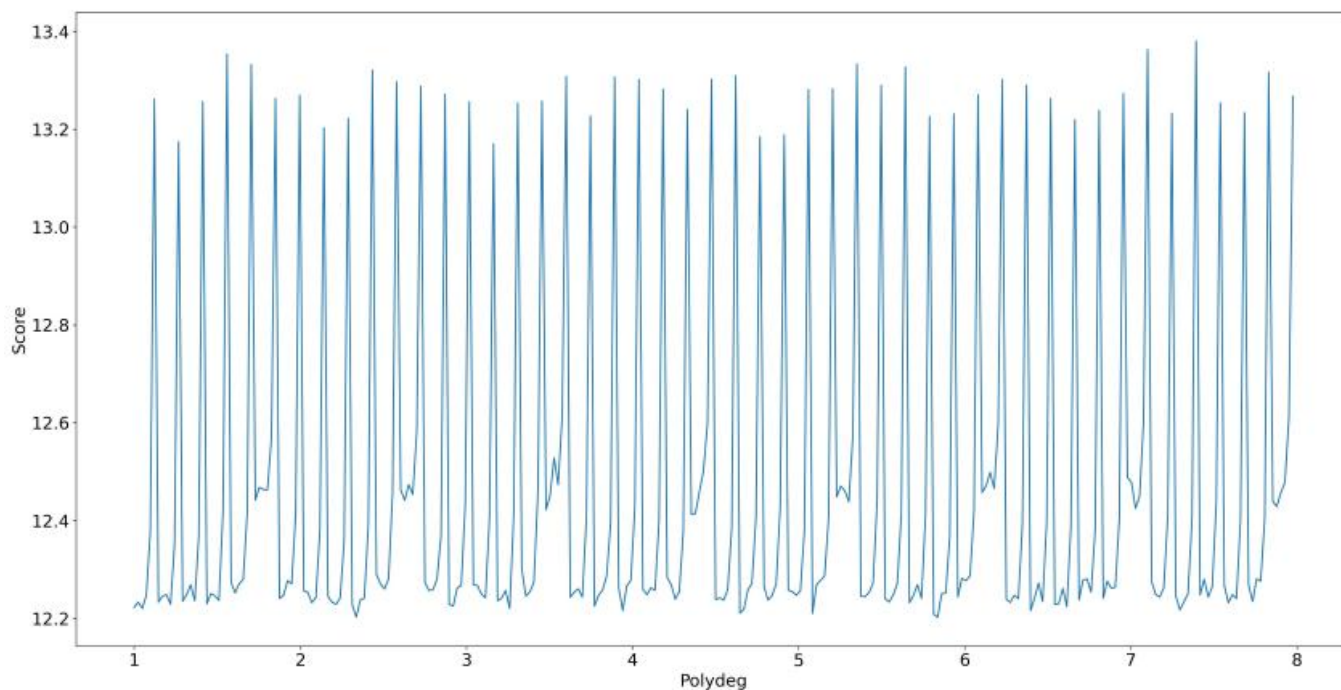
Можна зробити висновок, що оптимальні значення α_1 та α_2 знаходяться в межах до 0.01, вище – стрімке зменшення якості моделі. Найкращим α_1 виявилось 0.001, α_2 – 0.00001, poly_deg – 6.

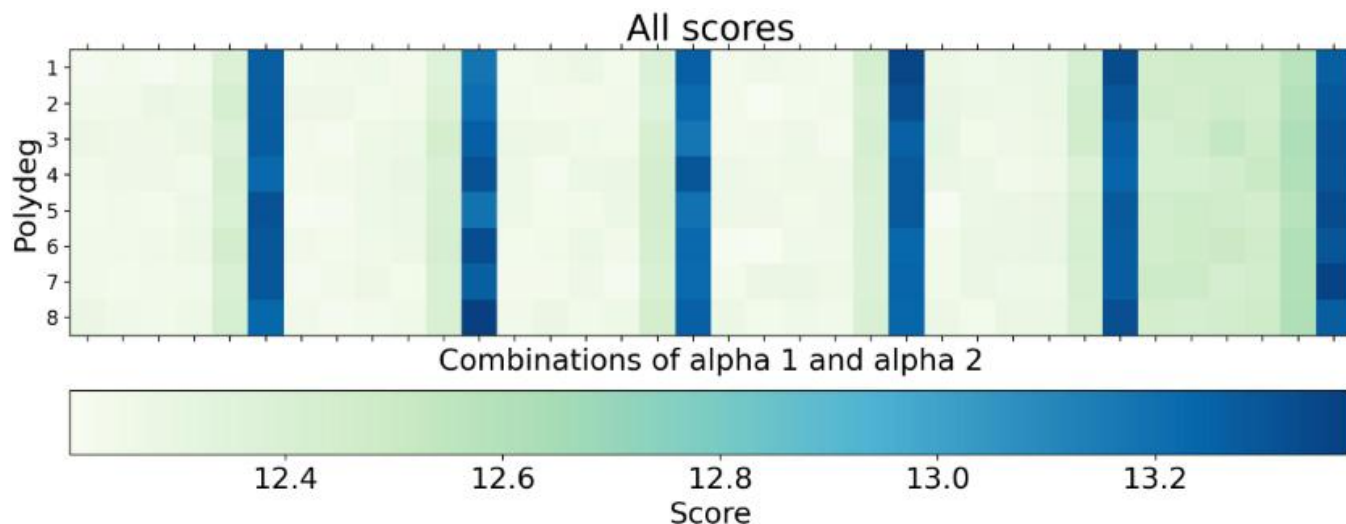
Найгіршим значенням score є 13.90356, отримане при таких параметрах:

| poly_deg | α_1 | α_2 | learning_rate | batch_size | train_steps |
|--------------------|------------|------------|-------------------------|----------------------|-----------------------|
| 1 | 0.00000 | 0.00000 | 0.001 | 128 | 100 |

Наведемо глобальний графік та діаграму типу `matshow` для всіх комбінацій коефіцієнтів (результат пошуку параметрів α_1 та α_2 разом з poly_deg):

Plot of all scores





Можна зробити висновок, що коефіцієнти α_1 та α_2 , які використовуються при обчисленні градієнту, є досить впливовими параметрами.

Після завершення роботи програми ми отримуємо наступний вивід на екран:

```
287 / 288 ; alpha1: 0.10000 ; alpha2: 0.01000... Time: 0 : 0 : 21 ; score: 12.61061
288 / 288 ; alpha1: 0.10000 ; alpha2: 0.10000... Time: 0 : 0 : 21 ; score: 13.26785
-----
Min score: 12.20193 ; polydeg: 6; alpha 1: 0.00100 ; alpha 2: 0.00001
learning rate: 0.0012000000000000003 ; batch size: 2048 ; train steps: 10000
-----
Max score: 13.90356 ; polydeg: 1; alpha 1: 0.00000 ; alpha 2: 0.00000
learning rate: 0.001 ; batch size: 128 ; train steps: 100
-----
Calculation the dependence of the predicted critical temperature on the true one...
Time: 0 : 0 : 25 ; Test MAE: 12.24457
Total time of all calculations: 1 : 49 : 45
Drawing plots... Time: 0 : 0 : 12
```

Візуалізація залежності передбаченої критичної температури від істинної (обчислення виконані з параметрами кращої моделі):

