

Optimization and Algorithms

Project report

Group 42

José Neves 89683, Leonardo Duarte 89691, and Gustavo Bakker 100660

1 Part 1

2 Part 2

2.1 Task 1

In this Part, the goal is to solve the optimization problem

$$\underset{(s,r) \in \mathbf{R}^n \times \mathbf{R}}{\text{minimize}} \quad \frac{1}{K} \sum_{k=1}^K (\log(1 + \exp(s^T x_k - r)) - y_k(s^T x_k - r)). \quad (1)$$

From (1), the objective function is $f : \mathbf{R}^n \times \mathbf{R} \rightarrow \mathbf{R}$

$$f(s, r) = \frac{1}{K} \sum_{k=1}^K (\log(1 + \exp(s^T x_k - r)) - y_k(s^T x_k - r)), \quad (2)$$

which can be written as

$$f(s, r) = \sum_{k=1}^K \frac{1}{K} h_k(s, r) + \sum_{k=1}^K \frac{1}{K} l_k(s, r), \quad (3)$$

where, for $k = 1, \dots, K$, $h_k : \mathbf{R}^n \times \mathbf{R} \rightarrow \mathbf{R}$

$$h_k(s, r) = \log(1 + \exp(s^T x_k - r)), \quad (4)$$

and $l_k : \mathbf{R}^n \times \mathbf{R} \rightarrow \mathbf{R}$

$$l_k(s, r) = -y_k(s^T x_k - r). \quad (5)$$

From (5), the functions l_k can be written as

$$l_k(s, r) = \begin{bmatrix} -y_k x_k \\ y_k \end{bmatrix}^T \begin{bmatrix} s \\ r \end{bmatrix} + 0, \quad (6)$$

Therefore, the functions l_k are affine. Since l_k are affine, they are also convex. In addition, from (4), the functions h_k can be written as

$$h_k(s, r) = (t_k \circ q_k)(s, r), \quad (7)$$

where, for $k = 1, \dots, K$, $t_k : \mathbf{R} \rightarrow \mathbf{R}$

$$t_k(z) = \log(1 + \exp(z)) \quad (8)$$

and $q_k : \mathbf{R}^n \times \mathbf{R} \rightarrow \mathbf{R}$

$$q_k(z) = s^T x_k - r. \quad (9)$$

From (8), it can be observed that the functions t_k are logistic functions and, therefore, convex. The functions q_k , in the other hand, can be written from (9) as

$$q_k(s, r) = \begin{bmatrix} x_k \\ -1 \end{bmatrix}^T \begin{bmatrix} s \\ r \end{bmatrix} + 0,$$

being, therefore, affine. Since the functions q_k are affine and t_k are convex, it is known from (7) that each function h_k can be written as the composition of an affine function with a convex function, which means the functions h_k are all convex. Since all functions h_k and l_k are convex, it is also known from (2) that the objective function f is the sum with positive coefficients of convex functions, which means f is convex.

2.2 Task 2

In this Task, the gradient descent algorithm will be used to solve (1). To use this algorithm, it is necessary to know the gradient of f . From (3), one can write

$$\nabla f(s, r) = \frac{1}{K} \sum_{k=1}^K (\nabla h_k(s, r) + \nabla l_k(s, r)). \quad (10)$$

From (5), the gradient of l_k can be written as

$$\nabla l_k(s, r) = \begin{bmatrix} -y_k x_k \\ y_k \end{bmatrix} = -y_k \begin{bmatrix} x_k \\ -1 \end{bmatrix}. \quad (11)$$

From (7), it can be written that

$$\nabla h_k(s, r) = \nabla t_k(q_k(s, r)) \nabla q_k(s, r), \quad (12)$$

where, from (9),

$$\nabla q_k(s, r) = \begin{bmatrix} x_k \\ -1 \end{bmatrix} \quad (13)$$

and, from (8),

$$\nabla t_k(z) = \frac{dt_k}{dz}(z) = \frac{\exp(z)}{1 + \exp(z)}. \quad (14)$$

Combining (9), (12), (13), and (14) leads to

$$\nabla h_k(s, r) = \begin{pmatrix} 1 - \frac{1}{\exp\left(\begin{bmatrix} x_k \\ -1 \end{bmatrix}^T \begin{bmatrix} s \\ r \end{bmatrix}\right) + 1} \end{pmatrix} \begin{bmatrix} x_k \\ -1 \end{bmatrix}. \quad (15)$$

Finally, combining (10), (11), and (15), one can write

$$\nabla f(s, r) = \frac{1}{K} \sum_{k=1}^K \begin{pmatrix} 1 - y_k - \frac{1}{\exp\left(\begin{bmatrix} x_k \\ -1 \end{bmatrix}^T \begin{bmatrix} s \\ r \end{bmatrix}\right) + 1} \end{pmatrix} \begin{bmatrix} x_k \\ -1 \end{bmatrix}. \quad (16)$$

In order to apply the gradient method, the following MATLAB script was written. This script returns the results of the method for each of the datasets given. This script implements the Gradient Descent algorithm through the MATLAB function *gradientDescent* also given below.

```

1 % GradientMethod.m
2 %% Initialization
3 clear;
4 clc;
5 NDataSets = 4;
6
7 %% Setup parameters
8 epsl = 1e-6; % stopping criterion
9 alpha_hat = 1; % initialization of alpha_k for the backtracking routine
10 gamma = 1e-4; % gamma of backtraking routine
11 beta = 0.5; % beta of backtraking routine
12 maxIt = [1e4; 1e4; 1e4; 1e5]; % maximum number of iterations
13
14 %% GD for each data set
15 for i = 1:NDataSets
16     %% Upload data
17     load(sprintf('./data%d.mat', i), 'X', 'Y'); % upload data set
18     K = length(Y);
19     n = size(X, 1);
20
21     %% Set up x0 (note that x = [s; r])
22     x0 = [-ones(n, 1); 0];
23
24     %% Setup objective function and gradient
25     h = [X; -ones(1, K)];
26     F = @(x) (1/K) * ...
27         sum(log(1+exp((h'*x)')) - Y.*(h'*x));

```

```

28     gradF = @(x) (1/K)*sum((exp((h'*x')) ./ ...
29         (1+exp((h'*x')))-Y) .*h, 2);
30
31     %% Run GD
32     fprintf("Running gradient descent for dataset %d (n = %d | K = %d).\n", ...
33         i, n, K);
34     tic
35     [xGD, ItGD, normGradGD] = gradientDescent(F, gradF, x0, eps1, ...
36         alpha_hat, gamma, beta, maxIt(i));
37     elapsedTimeGD = toc;
38     if ~isnan(xGD)
39         fprintf("Gradient descent for dataset %d"+...
40             " converged in %d iterations.\n", i, ItGD);
41         fprintf("Elapsed time is %f seconds.\n", elapsedTimeGD);
42         if i ≤ 2
43             fprintf("s = [%g; %g] | r = %g.\n", xGD(1), xGD(2), xGD(3));
44         end
45     else
46         fprintf("Gradient descent for dataset %d "+...
47             "exceeded the maximum number of iterations.\n", i);
48         fprintf("Elapsed time is %f seconds.\n", elapsedTimeGD);
49     end
50     save(sprintf("./DATA/GradientDescent/GDsolDataset%d.mat", i), ...
51         'xGD', 'ItGD', 'normGradGD', 'elapsedTimeGD');
52
53     %% Plot result
54     plotResults = false;
55     if plotResults
56         if i ≤ 2
57             figure('units', 'normalized', 'outerposition', [0 0 1 1]);
58             set(gca, 'FontSize', 35);
59             hold on;
60             ax = gca;
61             ax.XGrid = 'on';
62             ax.YGrid = 'on';
63             axis equal;
64             for k = 1:K
65                 if Y(k)
66                     scatter(X(1,k), X(2,k), 200, 'o', 'b', 'LineWidth', 3);
67                 else
68                     scatter(X(1,k), X(2,k), 200, 'o', 'r', 'LineWidth', 3);
69                 end
70             end
71             %ylim([-4 8]);
72             %xlim([-4 8]);
73             title(sprintf("Dataset %d", i));
74             ylabel('$x_2$', 'Interpreter', 'latex');
75             xlabel('$x_1$', 'Interpreter', 'latex');
76             x1 = (min(X(1,:)):(max(X(1,:))-min(X(1,:)))/100:max(X(1,:)));
77             plot(x1, (xGD(3)-xGD(1)*x1)/xGD(2), '-g', 'LineWidth', 4);
78             saveas(gcf, sprintf("./DATA/GradientDescent/GDsolDataset%d.fig", i));

```

```

79         close(gcf);
80         hold off;
81     end
82
83     figure('units','normalized','outerposition',[0 0 1 1]);
84     plot(0:ItGD,normGradGD,'LineWidth',3);
85     hold on;
86     set(gca,'FontSize',35);
87     ax = gca;
88     ax.XGrid = 'on';
89     ax.YGrid = 'on';
90     title(sprintf("Gradient method | Dataset %d",i));
91     ylabel('$||\Delta f(s_k,r_k)||$', 'Interpreter','latex');
92     xlabel('$k$', 'Interpreter','latex');
93     set(gca, 'YScale', 'log');
94     saveas(gcf,sprintf("./DATA/GradientDescent/GDNormGradDataset%d.fig",i));
95     close(gcf);
96     hold off;
97     end
98
99 end

```

```

1  % gradientDescent.m
2  function [xk,k,normGk] = gradientDescent(F,gradF,x0,epsl,...
3      alpha_hat,gamma,beta,maxIt)
4      %% Description
5      % Inputs: 1. F: objective function (as a function handle)
6      %          2. gradF: gradient of teh objective function (as a function
7      %          handle)
8      %          3. x0: initialization
9      %          4. epsl: stopping criterion
10     %          5. alpha_hat: initialization of alpha_k for the backtracking
11     %          routine
12     %          6. gamma: gamma of backtraking routine
13     %          7. beta: beta of backtraking routine
14     %          8. maxIt: maximum number of iterations
15     % Outputs: 1. x: output of the gradient descent method (returns NaN if
16     %          stopping criterion not met after the maximum number of
17     %          iterations chosen
18     %          2. k: number of iterations required for convergence if a
19     %          solution was found
20     %          3. normGk: norm of the gradient of the objective function
21     %% Gradient descent routine
22     k = 0;
23     xk = x0;
24     normGk = zeros(maxIt,1);
25     while k < maxIt
26         gk = gradF(xk); % Compute gradient at xk
27         normGk(k+1) = norm(gk);
28         if normGk(k+1) < epsl % Stopping criterion

```

```

29         break;
30     end
31     % ----- backtracking routine -----
32     alpha_k = alpha_hat;
33     % It is guaranteed that there is convergence, no maximum number of
34     % iterations needed (obviously for beta < 0)
35     while true
36         % check if F(alpha_k) < phi(0)+gamma*phi_dot(0)+alpha_k
37         if F(xk-alpha_k*gk) < F(xk)-gamma*alpha_k*(gk'*gk)
38             break; % alpha_k found
39         else
40             alpha_k = beta*alpha_k; % Update alpha_k
41         end
42     end
43     xk = xk - alpha_k*gk; % update xk
44     % ----- End backtracking routine -----
45     k = k + 1; % Increment iteration count
46 end
47 if k == maxIt
48     % No solution found within the maximum number of iterations
49     xk = NaN;
50 else
51     normGk = normGk(1:k+1);
52 end
53 end

```

For Task 2, the results obtained were $s = (1.3495, 1.0540)$ and $r = 4.8815$. The dataset 1 and the line defined by $\{x \in \mathbf{R}^2 : s^T x = r\}$ are represented in Fig. 1. In Fig. 2, the norm of the gradient along iterations is represented.

2.3 Task 3

In this Task, the code used was the one presented in the previous section. The results can be observed in Fig. 3 and 4 and the values obtained for s and r were $s = (0.7402, 2.3577)$ and $r = 4.5553$.

2.4 Task 4

In this Task, the gradient method was applied to two different datasets. However, in these datasets the points are no longer two-dimensional. In the dataset 3, $x_k \in \mathbf{R}^{30}$, for $k = 1, \dots, 500$, and in the dataset 4, $x_k \in \mathbf{R}^{100}$, for $k = 1, \dots, 8000$, which means it is no longer possible to represent the datasets. Therefore, only the global minimizers, s and r , and the evolution of the norm of the gradient along iterations will be presented.

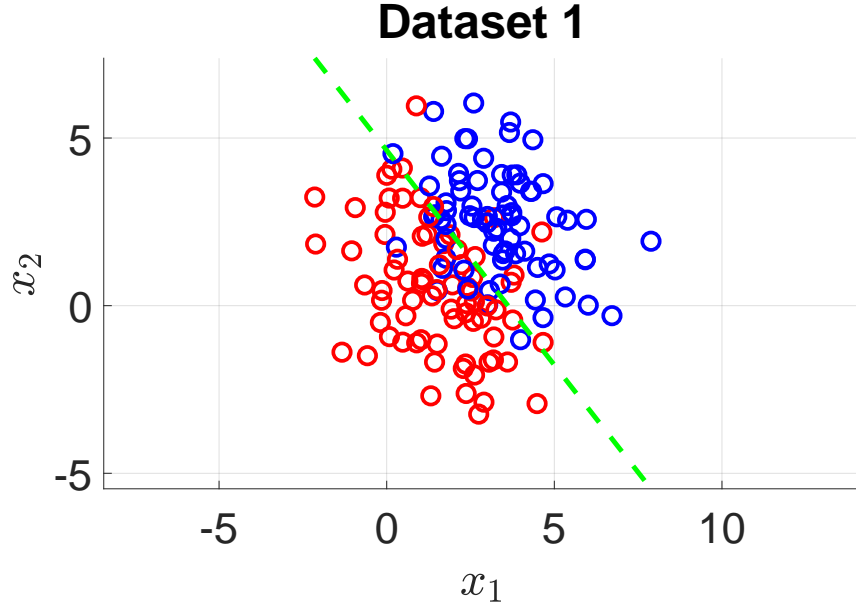


Figure 1: Dataset 1 and the corresponding line defined by $\{x \in \mathbf{R}^2 : s^T x = r\}$.

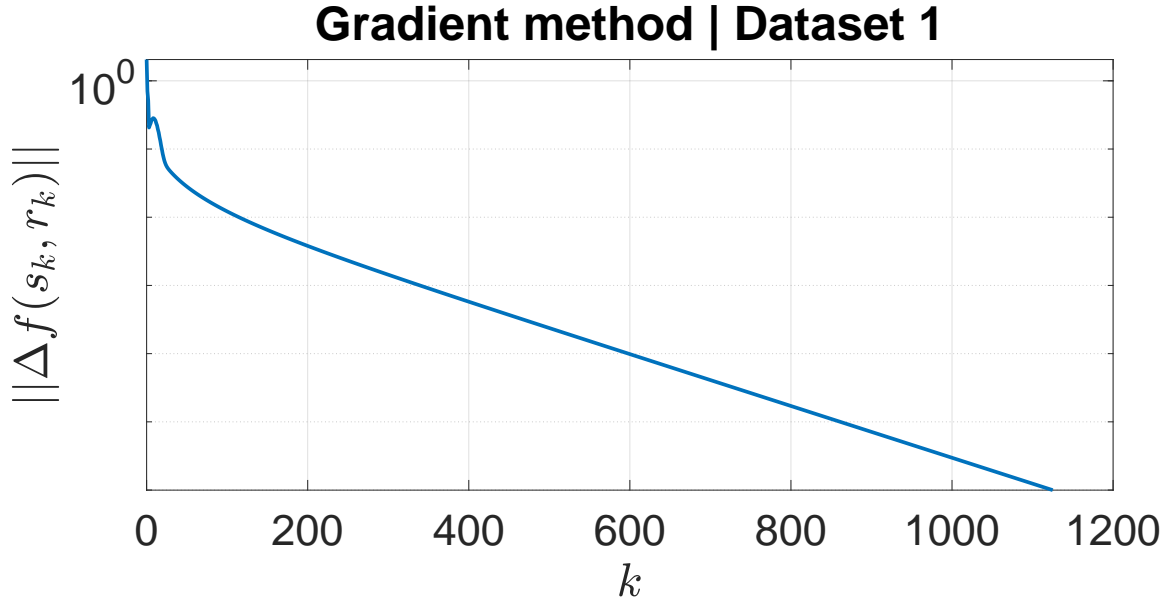


Figure 2: Norm of the gradient along iterations for the dataset 1.

For dataset 3,

$$s = [-1.3082, 1.4078, 0.8049, -1.0024, 0.5548, -0.5489, -1.1997, 0.0792, -1.8279, -0.1484, \\ 1.9241, -0.3586, -0.2900, 0.1925, 1.0614, 0.2107, -0.0929, 1.0476, -1.1248, -1.3311, \\ 0.7661, -0.2729, -0.5349, 0.9996, -0.4191, -0.3133, 0.4075, -0.1965, -0.7379, \\ -0.9814],$$

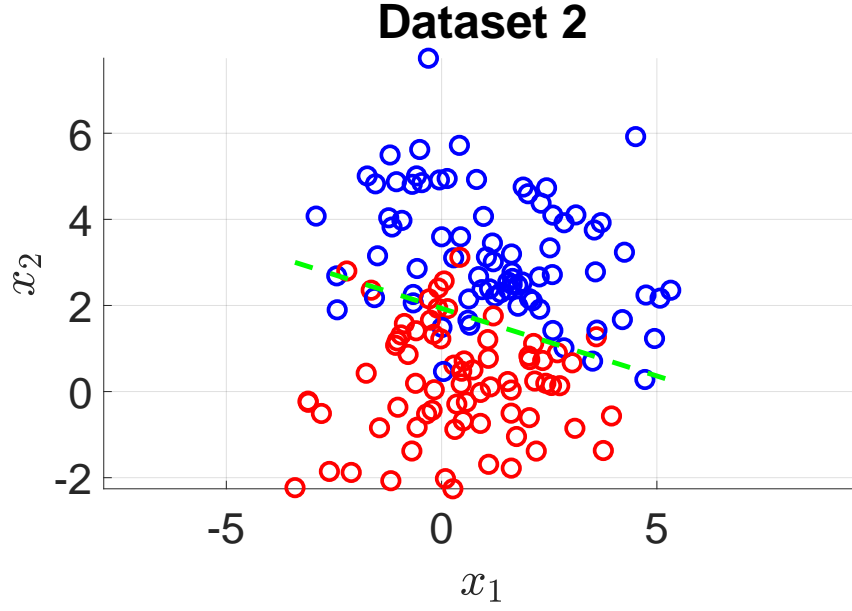


Figure 3: Dataset 2 and the corresponding line defined by $\{x \in \mathbf{R}^2 : s^T x = r\}$.

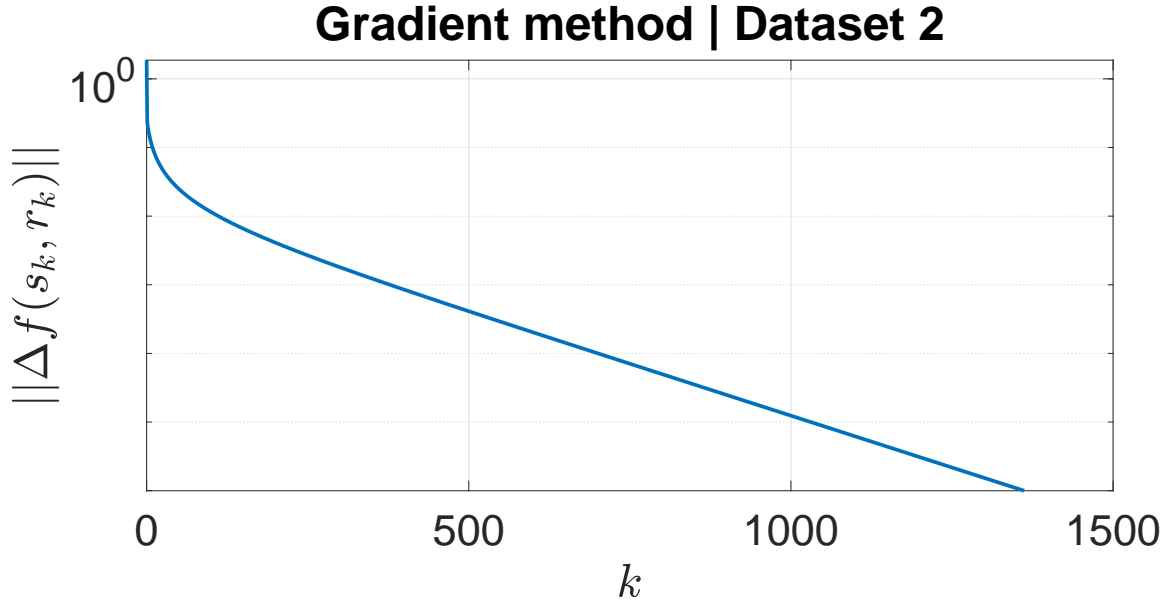


Figure 4: Norm of the gradient along iterations for the dataset 2.

$r = 4.7984$, and the evolution of the norm of the gradient along iterations is represented in

Fig. 5. For dataset 4,

$s = [0.1098, -0.6423, 0.1019, 1.2428, -1.6431, 1.0244, 0.0512, 0.8271, 0.3136, 0.7449, -0.5858,$
 $0.6267, 1.3611, 0.1534, 2.3234, -0.0840, -0.9489, 2.4699, -0.8678, -1.6516, 0.6460,$
 $-0.4779, 1.6397, 0.9034, -1.2293, -0.7587, -0.4887, 1.0306, 0.0888, -1.0917, -1.2717,$
 $-2.0333, -0.2505, -0.3518, -0.3486, -2.5610, -0.3132, -0.4902, 0.7258, 0.5774,$
 $-1.0528, 0.6400, 0.3759, -0.1547, 0.0298, 0.9547, -0.2863, 0.6364, 0.7859, 0.7584,$
 $0.2880, 0.1648, 0.6776, 2.0550, 1.0996, 0.5261, -0.5770, 1.1454, -0.5617, 0.0065, 0.4768,$
 $-2.3677, -1.1561, -2.6619, 0.0622, 0.1037, -0.6237, 0.1913, 0.6672, -1.0493, -0.3240,$
 $-0.3207, -1.0904, -0.8293, -0.3104, -0.4879, -0.1060, -0.1646, 2.2683, -1.2380,$
 $-0.8575, -2.4781, -0.4158, 0.1660, 0.7931, 0.3685, -0.0524, -0.9997, -0.5732, 0.3971,$
 $1.1911, 1.8318, -1.7287, 0.2329, -1.1921, 1.6558, 0.4612, -0.6431, 0.8295, 0.2975],$

$r = 7.6701$, and the evolution of the norm of the gradient along iterations is represented in Fig. 6.

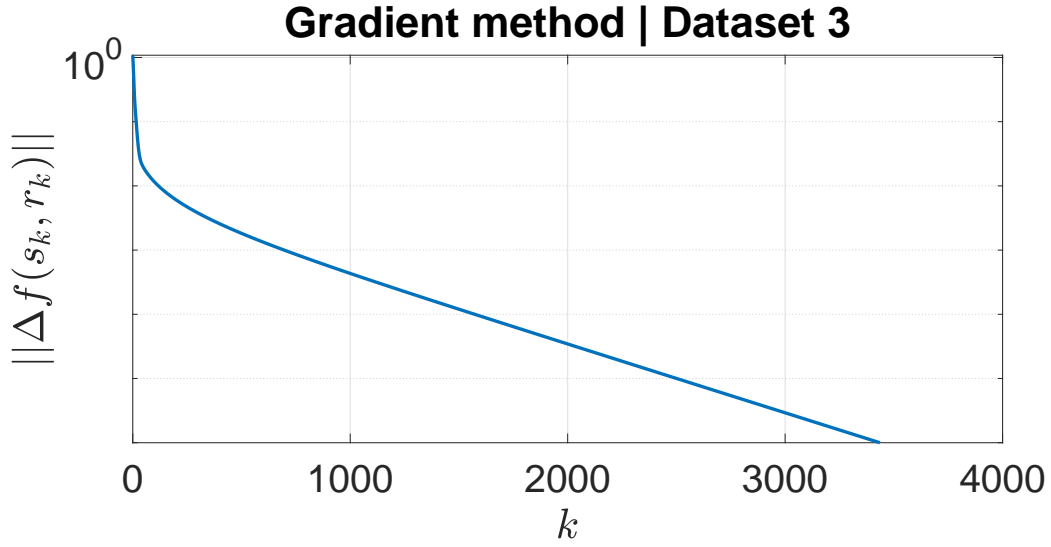


Figure 5: Norm of the gradient along iterations for the dataset 3.

2.5 Task 5

Letting $\phi : \mathbf{R} \rightarrow \mathbf{R}$ be a twice differentiable function and supposing $p : \mathbf{R}^3 \rightarrow \mathbf{R}$ is given by

$$p(x) = \sum_{k=1}^K \phi(a_k^T x), \quad (17)$$

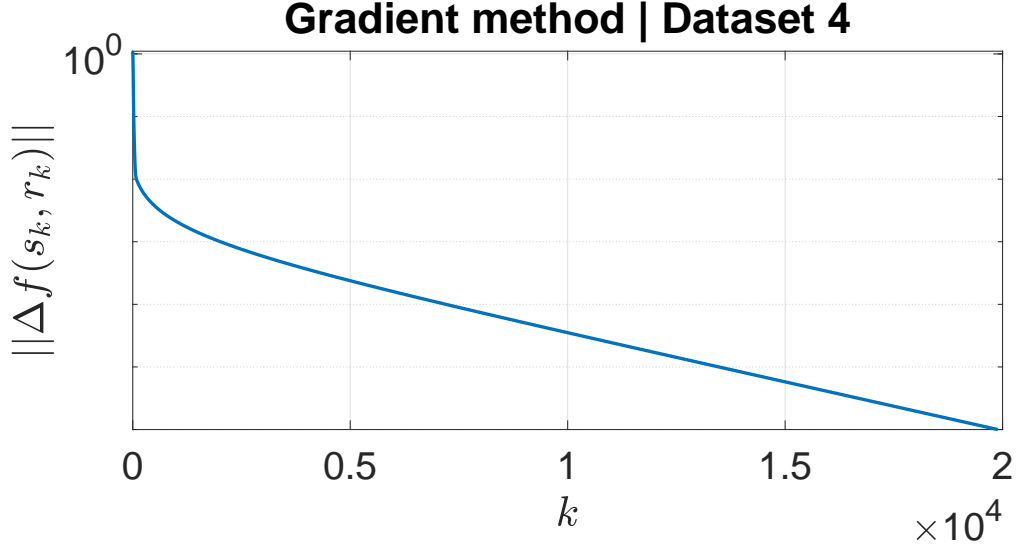


Figure 6: Norm of the gradient along iterations for the dataset 4.

where $a_k \in \mathbf{R}^3$ for $k = 1, \dots, K$, one can write

$$\nabla p(x) = \sum_{k=1}^K \dot{\phi}(a_k^T x) a_k = Av, \quad (18)$$

where $A = [a_1 \ a_2 \ \dots \ a_K]$ and $v = [\phi(a_1^T x) \ \phi(a_2^T x) \ \dots \ \phi(a_K^T x)]^T$. In order to write the Hessian of p at x , one can define, for $k = 1, \dots, K$, the functions $u_k : \mathbf{R} \rightarrow \mathbf{R}^3$

$$u_k(z) = za_k, \quad (19)$$

which, from (18), lead to

$$\nabla p(x) = \sum_{k=1}^K u_k \left[\dot{\phi}(a_k^T x) \right]. \quad (20)$$

From (19) and (20), it is possible to write

$$\nabla^2 p(x) = \sum_{k=1}^K Du_k \left[\dot{\phi}(a_k^T x) \right] D \left[\dot{\phi}(a_k^T x) \right] D(a_k^T x) = \sum_{k=1}^K a_k \ddot{\phi}(a_k^T x) a_k^T = ADA, \quad (21)$$

where D is the diagonal matrix

$$D = \begin{bmatrix} \ddot{\phi}(a_1^T x) & & & \\ & \ddot{\phi}(a_2^T x) & & \\ & & \dots & \\ & & & \ddot{\phi}(a_K^T x) \end{bmatrix} \quad (22)$$

2.6 Task 6

In this Task, the Newton method will be used to solve (1). To be able to use this method, it is necessary to know the gradient, which is given by (16), and the Hessian of the objective function. In order to write the Hessian, one starts by writing, from (2), (4), and (5),

$$f(x) = \sum_{k=1}^K \left(\phi(a_k^T x) + \frac{1}{K} \begin{bmatrix} -y_k x_k \\ y_k \end{bmatrix}^T x \right), \quad (23)$$

where $x = [s \quad r]^T$, $a_k = [x_k \quad -1]^T$ and $\phi : \mathbf{R} \rightarrow \mathbf{R}$

$$\phi(z) = \frac{1}{K} \log(1 + \exp(z)) \quad (24)$$

with second-derivative

$$\ddot{\phi}(z) = \frac{\exp(z)}{K [1 + \exp(z)]^2}. \quad (25)$$

Taking into consideration that (23) is written in the form of (17) except for the sum of affine terms whose second-derivative is null, it may be written, from (21), (22), and (25),

$$\nabla^2 f(x) = \begin{bmatrix} a_1 & a_2 & \dots & a_K \end{bmatrix} \begin{bmatrix} \ddot{\phi}(a_1^T x) & & & \\ & \ddot{\phi}(a_2^T x) & & \\ & & \dots & \\ & & & \ddot{\phi}(a_K^T x) \end{bmatrix} \begin{bmatrix} a_1 & a_2 & \dots & a_K \end{bmatrix}. \quad (26)$$

Taking into consideration (16) and (26), the Newton method was implemented according to the script below. In it, the Newton algorithm is implemented through the MATLAB function *newtonAlgorithm* also presented below.

2.7 Task 7

3 Part 3