

MATH-329 Nonlinear optimization Homework 1: gradient descent and convexity

Benoît Müller

Instructor: Nicolas Boumal

TAs: Quentin Rebjock and Axel Séguin

October 4, 2023

Remarks on the code

Each file named questionx.m is the main file of the question number x. The other files are data or functions named after their utility. Running the main script of a question give the output detailed in this document.

We use showMNSTImage.m and showMNSTImages_many.m as it has been given in the starting file, without changing them.

Answers to the questions

1. Using the fact that $-\log(x^a y^b) = a \log(1/x) + b \log(1/y)$, we obtain

$$\begin{aligned} f(\theta) &= -\log l(\theta) = -\log \prod_{j=1}^m \sigma(\langle x_j, x \rangle + b)^{y_j} \sigma(-\langle x_j, x \rangle - b)^{1-y_j} \\ &= -\sum_{j=1}^m \left(y_j \log(\sigma(\langle \tilde{x}_j, \theta \rangle)) + (1 - y_j) \log(\sigma(-\langle \tilde{x}_j, \theta \rangle)) \right) \\ &= \sum_{j=1}^m \left(y_j \log(1 + e^{-\langle \tilde{x}_j, \theta \rangle}) + (1 - y_j) \log(1 + e^{\langle \tilde{x}_j, \theta \rangle}) \right). \end{aligned}$$

2. We call it $\tau : z \mapsto \log(1 + e^z)$ and it is clearly twice differentiable with first derivative

$$\tau'(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}} = \sigma(z),$$

and second derivative

$$\tau''(z) = \frac{e^z}{(1 + e^z)^2} > 0.$$

Then by a classic result of calculus, having a positive second derivative, it must be convex.

3. We know that a linear combination (with positive factors) of convex functions is convex again. We then just have to show that $\log(1 + e^{\pm \langle \tilde{x}_j, \theta \rangle})$ are convex in θ , since $\{y_j, 1 - y_j\} = \{0, 1\} \subset \mathbb{R}_+$ for all $j \leq m$.

Moreover, for two twice differentiable functions f and g , $(h \circ g)'' = ((h' \circ g)g')' = ((h'' \circ g)(g')^2 + (h' \circ g)g'')$ is positive if h' , h'' and g'' are positive, that mean that $h \circ g$ is convex if g and h are convex too and h is non-decreasing.

In our case, the function τ is convex and has a positive derivative so it's increasing. $\pm\langle\tilde{x}_j, \theta\rangle$ are both linear in θ and so convex too. We can apply what we've just said and conclude that $\log(1 + e^{\pm\langle\tilde{x}_j, \theta\rangle})$ are convex in θ , implying that f is convex.

4. By construction, $f_\lambda(\theta) - \frac{\lambda}{2}\|\theta\|^2 = f(\theta)$, which is convex (Question 3.). By definition, f_λ is then λ -strongly convex.
5. By the course, λ -strong convexity implies the existence of a unique global minimum and since f_λ is continuously differentiable, it's a point where the gradient vanishes because it is in particular a local minimum. We can conclude directly thanks to Question 4.
6. We claim that, with respect to the standard euclidean structure of \mathbb{R}^m ,

$$\nabla f_\lambda(\theta) = \sum_{j=1}^m \log(1 + e^{\langle\tilde{x}_j, \theta\rangle}) - Xy + \lambda\theta = X(\tau(X^\top\theta) - y) + \lambda\theta$$

with the matrix $X = (\tilde{x}_1 \mid \dots \mid \tilde{x}_m) \in \mathbb{R}^{d \times m}$, and where we extend τ (and other functions like exp etc...) to matrix entries with $\tau(A) = (\tau(A_{ij}))_{ij}$. Let's prove it:

First, noticing that $\tau(z) = \log(1 + e^z) = z + \log(e^{-z} + 1) = z + \tau(-z)$, we rewrite f as follow :

$$\begin{aligned} f(\theta) &= \sum_{j=1}^m \left(y_j \log(1 + e^{-\langle\tilde{x}_j, \theta\rangle}) + (1 - y_j) \log(1 + e^{\langle\tilde{x}_j, \theta\rangle}) \right) \\ &= \sum_{j=1}^m \left(y_j (-\langle\tilde{x}_j, \theta\rangle + \log(1 + e^{\langle\tilde{x}_j, \theta\rangle})) + (1 - y_j) \log(1 + e^{\langle\tilde{x}_j, \theta\rangle}) \right) \\ &= \sum_{j=1}^m \left(\log(1 + e^{\langle\tilde{x}_j, \theta\rangle}) - y_j \langle\tilde{x}_j, \theta\rangle \right) = \mathbf{1}^\top \log(1 + \exp(X^\top\theta)) - (Xy)^\top\theta \\ &= \mathbf{1}^\top \tau(X^\top\theta) - (Xy)^\top\theta. \end{aligned}$$

with the vector $\mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^{d \times 1}$.

Now using basic properties of jacobians and gradients, $\nabla f_\lambda(\theta) = D_\theta(\mathbf{1}^\top \tau(X^\top\theta)) - Xy + \lambda\theta$. The first term is computed as follow :

$$(D_\theta(\mathbf{1}^\top \tau(X^\top\theta))) = \sum_{j=1}^m D_\theta \tau(\langle\tilde{x}_j, \theta\rangle) = \sum_{j=1}^m \tilde{x}_j \tau'(\langle\tilde{x}_j, \theta\rangle) = X\sigma(X^\top\theta).$$

Together, this give us $X\sigma(X^\top\theta) - Xy + \lambda\theta = X(\sigma(X^\top\theta) - y) + \lambda\theta$

7. First we implement τ which will be needed and which has bad computing precision properties when the argument become big. The value of $\tau(z)$ approach z but before the log take the value back down, the exp take the value very high and thus, is interpreted as infinity before the final value become too big to be actually stored. Our strategy is then to factorise and take a z out to have the sum of z and the small rest : $\tau(z) = \log(1 + e^z) = \log(e^z(e^{-z} + 1)) = \log(e^z) + \log(e^{-z} + 1) = z + \log(e^{-z} + 1)$. Like this, the e^{-z} can go as small as he want (even 0), it will make the log vanish and leave the z alone. We can now separate the cases where z is big or small and use the right formula for each :

```

1 function [y]=tau(x)
2 % Input any matrix
3 % compute for each index log(1+e^x) with ~stable computation

```

```

4 index = x>0;
5 y = x;
6 y(index) = x(index) + log(exp(-x(index)) + 1);
7 y(~index) = log(exp(x(~index)) + 1);
8 end

```

Now we can implement the function that compute f_λ and its gradient with the formulas of Question 6. We put the data in the input in order not to load it each time we call the function. Moreover the function works with multiples points stored in a matrix. In order to gain time when we only want one of the two outputs we ask for a third input to precise which output to compute.

```

1 function [f, g] = logistic_regression(train,theta,output)
2 % Input:
3 %
4 %   train : the data to train on
5 %   theta : A tentative parameter vector (column) for logistic regression.
6 %           if theta is a matrix [theta(1),...,theta(N)] it will return
7 %           each output concatenated in a row.
8 %   output: 'f','g','fg', if we need f, g, or f and g respectively
9 %
10 % Output:
11 %
12 %   f: value of the logistic regression negative log-likelihood cost function
13 %   g: gradient of the cost function at theta, as a column vector
14
15
16 % train: The training data set (see Moodle)
17 %       This is a structure with fields train.X and train.y
18 %       containing, respectively, the examples (the images) and the
19 %       labels to be used for training.
20
21   X = train.X;
22   y = train.y;
23
24   lambda = 1e-4 ; % regularization parameter
25   sigma = @(x) (exp(-x) + 1).^(-1);
26   product = X' * theta;
27   if output=='f' | output=='fg'
28       f = sum(tau(product)) - y'*product + lambda/2 * sum(theta.^2);
29   else
30       f=[];
31   end
32
33   if output=='g' | output=='fg'
34       g= X * (sigma(product) - y) + lambda * theta;
35   else
36       g=[];
37   end

```

8. Here is the code for it and the plot in Figure 1 :

```

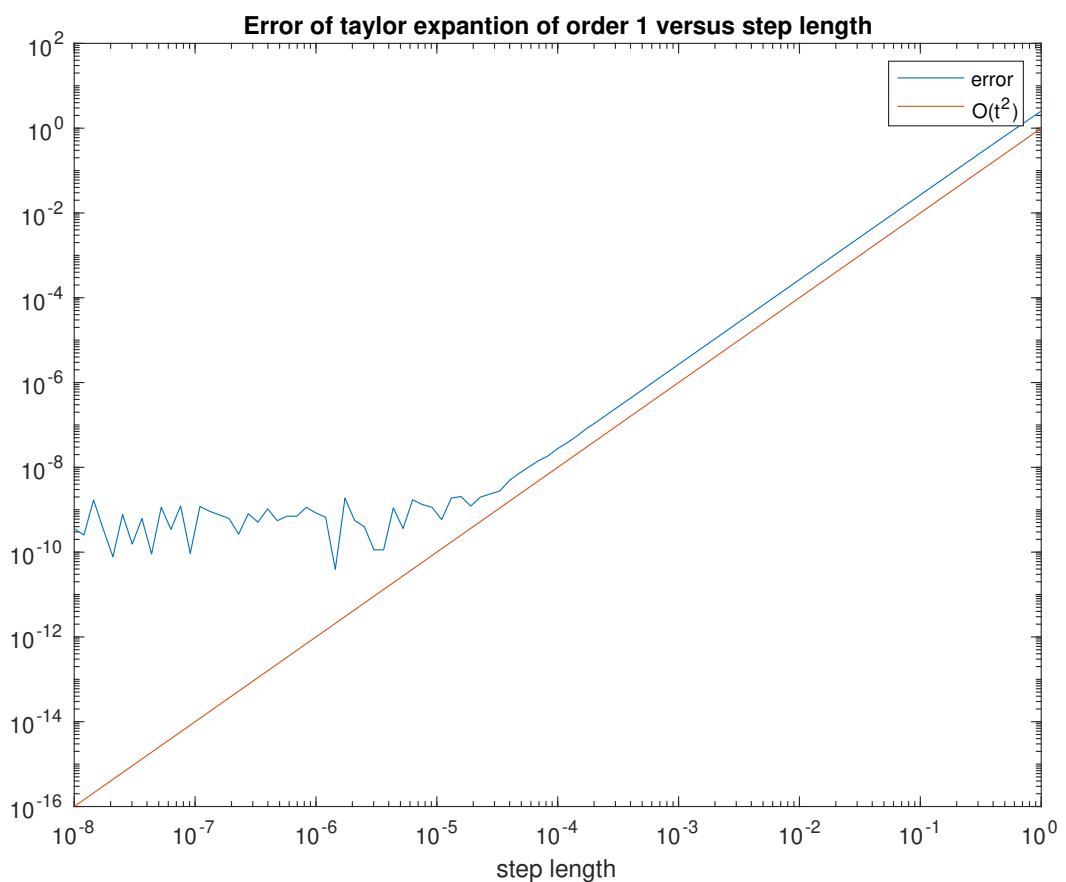
1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2 % Question 8
3 % Test of logistic_regression and its gradient by...
4 % taylor expansion of order 2...
5 % in a random point and direction.
6 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
7

```

```

8 load('mnist_train_test.mat');
9 theta = rand(28*28+1, 1); % random start
10 v = rand(size(theta));
11 v = v / norm(v); % random direction
12 t = logspace(-8, 0, 101); % length of step
13
14 [f g] = logistic_regression(train, theta, 'fg');
15 [ft, ~] = logistic_regression(train, t.*v + theta, 'f');
16 error = abs(ft - f - (g'*v)*t);
17
18 loglog(t, error, t, t.^2)
19 title('Error of Taylor expansion of order 1 versus step length')
20 xlabel('step length')
21 legend('error', 'O(t^2)')

```



The basic idea of this is to compute a Taylor expansion of order two that involve both the gradient and different values of the function, and to see if this work. In some sense it could be a method to approximate the gradient and it's called Finite difference method. As it is a order two Taylor polynomial, the error is by definition of order two, which is what we obtain in the plot and confirm the coherence of our code. We notice that for too small values, the curve tend to be irregular and go away of the line, showing the limits and the round-off errors of the computer.

9. We compute the hessian of f_λ :

$$\begin{aligned}\nabla^2 f_\lambda(\theta) &= D(\nabla f_\lambda)(\theta) = X D_\theta(\sigma(X^\top \theta)) + \lambda I = X(D\sigma)(X^\top \theta)X^\top + \lambda I = X \text{diag}(\sigma'(X^\top \theta))X^\top + \lambda I \\ &= X \text{diag}(\sigma'(X^\top \theta))X^\top + \lambda I.\end{aligned}$$

Its spectral norm can be bounded as follow, using the under-linearity and under-multiplicity :

$$\|\nabla^2 f_\lambda(\theta)\| \leq \|X\| \|\text{diag}(\sigma'(X^\top \theta))\| \|X^\top\| + \lambda \leq \|X\|^2 \sup |s'| + \lambda.$$

We see that the function $|\sigma'| = \sigma'$ has a max in 0 by looking at its derivative that change of sign in 0, because

$$\sigma'(z) = \frac{e^z}{(1+e^z)^2} \text{ and } \sigma''(z) = \frac{e^z(1-e^z)}{(1+e^z)^3}.$$

The bound is then $\sup |\sigma'| = \max \sigma' = \sigma'(0) = 1/4$ Finally, the spectral norm of the hessian is smaller than $L := 1/4\|X\|^2 + \lambda$. By a result of the course, f_λ being twice continuously differentiable, is L -strongly convex.

10. By running `1/4 * norm(X)^2` in MATLAB, we obtain $1/4\|X\|^2 \approx 198\,570$

11. A classic implementation of the gradient descent :

```

1 function [theta, grad, time]= my_optimizer(theta0, maxtime)
2
3     % optimize logistic regression by const step GD
4     %
5     % INPUT  theta0 : initial point
6     %         maxtime : maximal time of research
7     %
8     % OUTPUT theta : last point of search
9     %         grad  : vector of successive norm of gradient
10    %         time  : vector of successive time
11
12    % Answer to Question 11.
13
14    load('mnist_train_test.mat');
15    f = @(theta) logistic_regression(train, theta, 'g');
16    L = norm(train.X)^2 / 4 + 1e-4;
17    alpha = 1/L;
18
19    theta = theta0;
20    grad = [];
21    time = [];
22    tic
23    while toc < maxtime
24        time = [ time toc ];
25        [~, G] = f(theta);
26        grad = [grad norm(G)];
27        theta = theta - alpha * G;
28    end
29 end

```

12. The code and the plot of the norm of the gradient $\|\nabla f_\lambda(\theta_k)\|$ as a function of k (log-scale on the vertical axis):

```

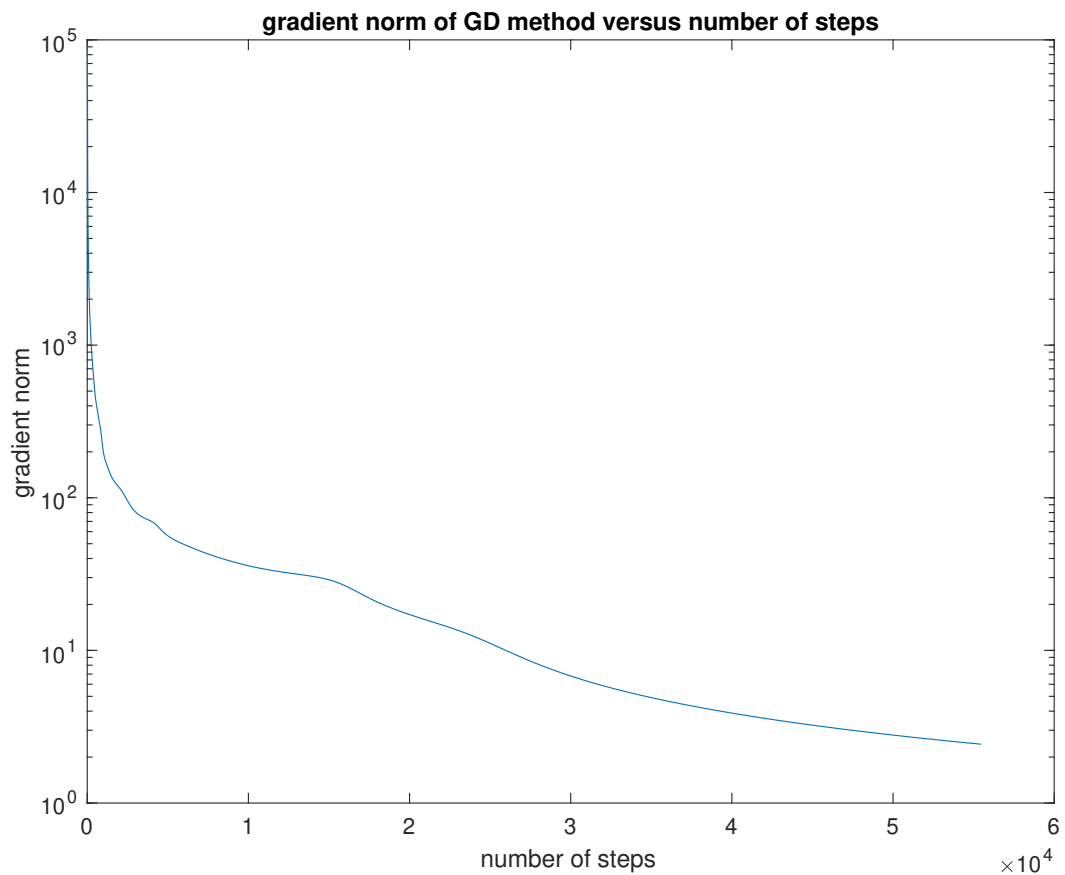
1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

```

```

2 % Question 12
3 % gradient norm of GD method versus number of steps
4 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
5
6 theta0 = rand(785,1); % random start
7 maxtime = 5*60;      % maximal time
8
9 [theta, grad, ~] = my_optimizer(theta0, maxtime);
10 semilogy( grad )
11 title('gradient norm of GD method versus number of steps')
12 xlabel('number of steps')
13 ylabel('gradient norm')

```



13. We have shown that f_λ is twice differentiable with Lipschitz continuous gradient, is strongly convex, and has a unique critical point with positive hessian, so by a theorem of the course, the method converge to this point. Moreover the evaluations of the function converge at least linearly to the minimum. the norm of the gradient converge at least linearly to zero with rate $r = \sqrt{1 - 1/k}$ where k is the condition number of the hessian of the minimizer.

In our plot the method clearly converge, very fast at first, but eventually go very slower which motivate the use of line search methods.

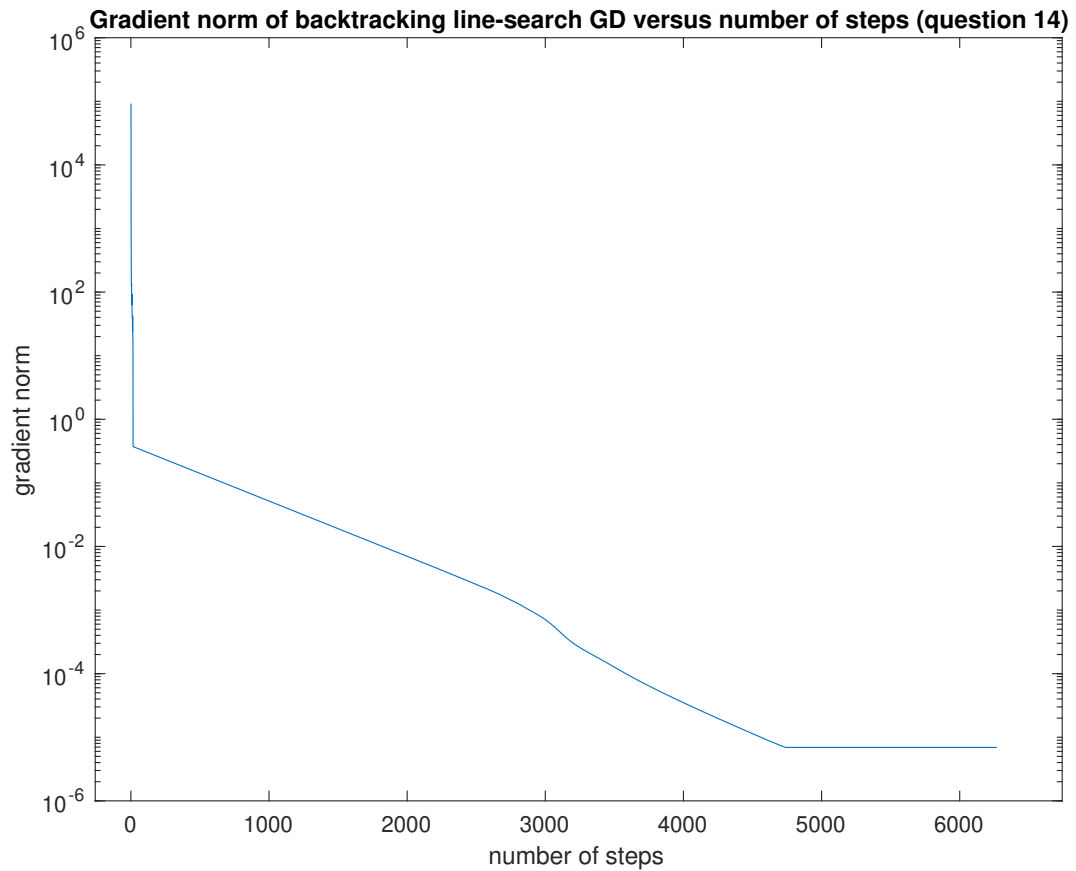
14. We implement the line search method. The parameters have been chosen like this : $\alpha = 20$ $\rho = 0.5$, $c = 10^{-4}$. The choice for c and ρ are usual values. We plotted the chosen values of $\rho^k \alpha$ for each step of the algorithm when we put a big initial α and a ρ close to 1 (for precision), so we can see the most frequent values chosen. All of this made us choose $\alpha = 20$, and keep $\rho = 0.5$ for short computational time.

Here the code and the plot :

```

1 function [theta, grad, time]= my_optimizer2(theta0, maxtime)
2
3 % Optimize logistic_regression with...
4 % backtracking line-search GD method.
5 %
6 % INPUT  theta0 : first point of research
7 %        maxtime : maximal time of research
8 %
9 % OUTPUT theta : last point of research
10 %        grad  : vector of the successive norms of gradient
11 %        time   : vector of successives times
12
13 % Answer to Question 14.
14
15 load('mnist_train_test.mat');
16 fg= @(theta) logistic_regression(train,theta,'fg');
17 f= @(theta) logistic_regression(train,theta,'f');
18 g= @(theta) logistic_regression(train,theta,'g');
19
20 alpha0=20;
21 rho=0.5; % in [0.5,0.8]
22 c=1e-4;
23
24 theta=theta0;
25 grad=[];
26 time=[];
27 tic
28 while toc<maxtime
29     time=[ time toc];
30     [F, G]=fg(theta);
31     grad=[grad norm(G)];
32     alpha=alpha0;
33     [Fnext, ~]=f(theta - alpha*G);
34     while Fnext > F - c * alpha * G'*G
35         alpha = alpha * rho;
36         [Fnext, ~]=f(theta - alpha*G);
37     end
38     theta = theta - alpha*G;
39 end
40 end

```



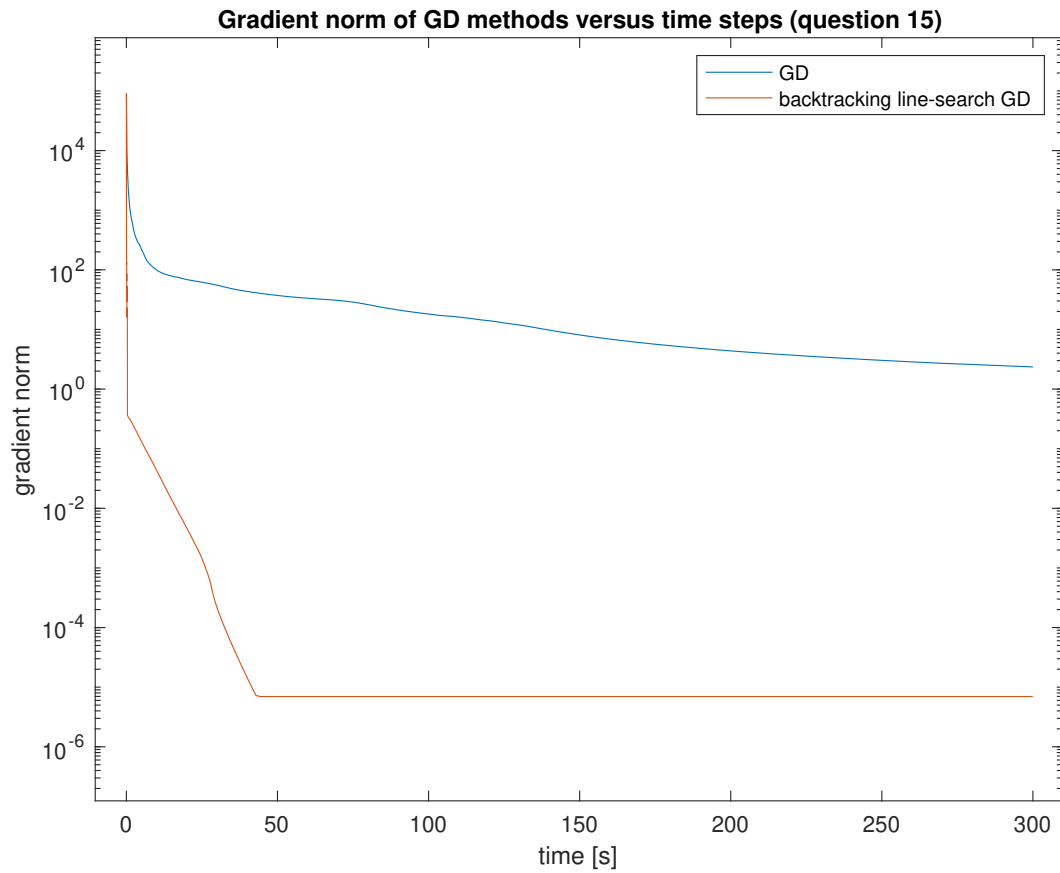
The resulting theta obtained at the end of this program is stored in a variable file (final.theta.mat) and will be used for the algorithm.

15. The code and the plot for both methods with respect to the time :

```

1  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2  % Question 15
3  % gradient norm of GD methods versus time steps
4  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
5
6  theta0 = rand(785,1);
7  maxtime = 5*60;
8
9  [~, grad, time] = my_optimizer(theta0, maxtime);
10 [~, grad2, time2] = my_optimizer2(theta0, maxtime);
11
12 semilogy(time, grad, time2, grad2)
13 title('Gradient norm of GD methods versus time steps (question 15)')
14 xlabel('time [s]')
15 ylabel('gradient norm')
16 legend('GD', 'backtracking line-search GD')

```

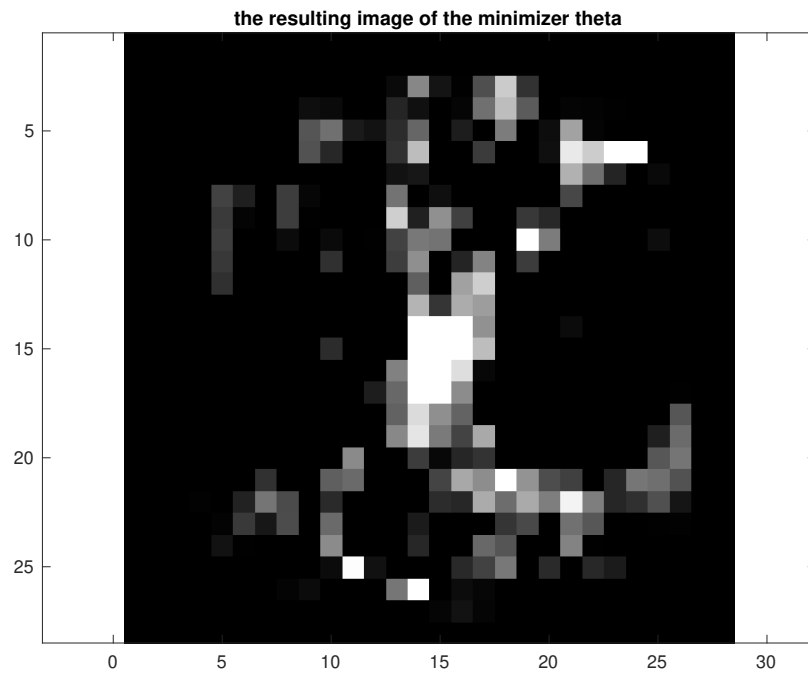



16. The code to show the resulting θ and the value of b:

```

1  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2  % Question 16
3  % Display of the image of the minimizer theta
4  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
5
6  load('final_theta.mat')
7  b=final_theta(end)
8  showMNISTImage( final_theta )
9  title('the resulting image of the minimizer theta')

```



$b = -1.1452$

17. the final code that report the results :

```

1  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2  % Question 17
3  % Display the images that our theta got wrong
4  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
5
6  load('mnist.train-test.mat');
7  load('final.theta');
8
9  fprintf('\n\n-----RESULTS ON THE TEST DATA :-----\n')
10 results(test,final.theta);
11 title('mispredicted images of test')
12
13 fprintf('\n\n-----RESULTS ON THE TRAIN DATA :----- \n')
14 results(train,final.theta);
15 title('mispredicted images of train')
16
17 function results(data,theta)
18     X = data.X ;
19     y = data.y' ;
20     sigma = @(x) (1+exp(-x)).^(-1);
21     z = sigma(theta' * X); % probability
22     prediction = z > 0.5;
23     error = prediction ~= y;
24     m=length(y);
25     % display on the command window :
26
27     text='On the %d image(s), %d have been mispredicted, giving %.2f%% of accuracy.\n';
28     fprintf( text,m, sum(error),100*(m-sum(error))/m );

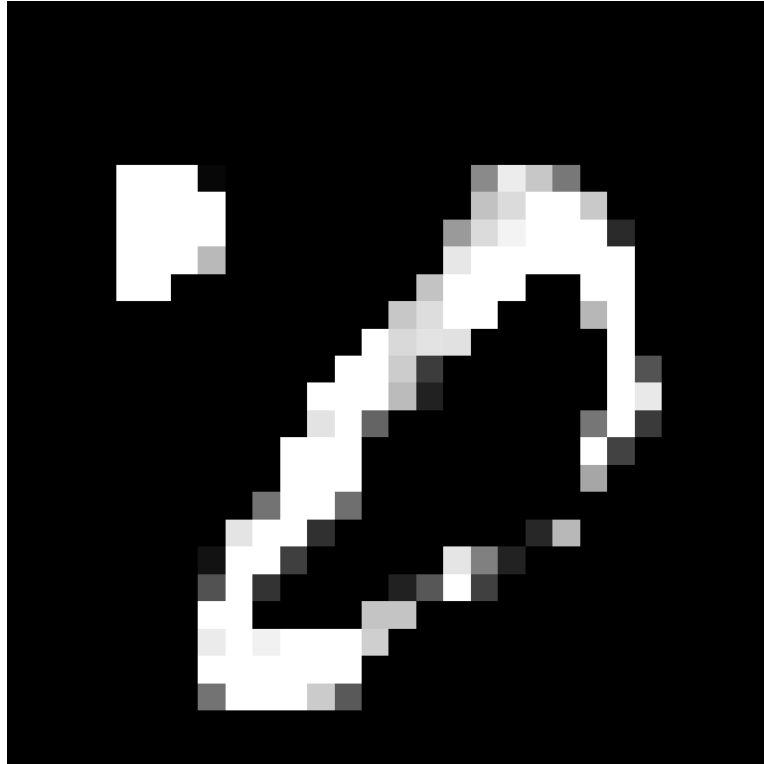
```

```
29
30     text= 'The average difference between the probability and the resulting';
31     text=[text , ' prediction is ~%.1d.\n'];
32     text= [text 'The %d mispredicted images have probability : \n'];
33     fprintf(text, sum(abs(z-prediction))/m, sum(error))
34     fails_probabilities=z(error)
35
36     showMNISTImages_many( X(:,error) )
37 end
```

With the following output in the command window :

```
1          -----RESULTS ON THE TEST DATA :-----
2 On the 2115 image(s), 1 have been mispredicted, giving 99.95% of accuracy.
3 The average difference between the probability and the resulting prediction is ~1.5e-04.
4 The 1 mispredicted images have probability :
5
6 fails_probabilities =
7
8     0.9984
9
10
11 -----RESULTS ON THE TRAIN DATA :-----
12 On the 12665 image(s), 0 have been mispredicted, giving 100.00% of accuracy.
13 The average difference between the probability and the resulting prediction is ~8.9e-08.
14 The 0 mispredicted images have probability :
15
16 fails_probabilities =
17
18     1 0 empty double row vector
```

In conclusion, the algorithm has been infallible on the train data, and has made only one mistake on the test data. Here is the mispredicted image :

mispredicted images of test

We conclude that our algorithm is coherent and allows to actually read images with a certain reliability as it worked well enough on the completely new data, with respect to our overarching goal.