

Optimization and Algorithms

Part 2 of the Project

João Xavier and Cláudia Soares
Instituto Superior Técnico
October 2020

Contents

1	Introduction to logistic regression	1
2	Gradient method	3
3	Newton method	6

1 Introduction to logistic regression

Automatic prediction. Consider the dataset $D = \{(x_k, y_k) \in \mathbf{R}^2 \times \{0, 1\} : k = 1, \dots, K\}$ pictured in Figure 1. In this dataset, each element is a pair of the form (x_k, y_k) , where $x_k \in \mathbf{R}^2$ is a vector with two entries and $y_k \in \{0, 1\}$ is a binary label, encoded in Figure 1 with the color red for $y_k = 0$ and blue for $y_k = 1$.

This kind of dataset arises in many real-life setups. For example, it could be generated by a doctor as follows. The doctor starts by taking two measurements (say, weight and blood pressure) from an individual 1, thereby filling the two entries of the vector $x_1 \in \mathbf{R}^2$; then, based on those measurements, the doctor decides whether individual 1 suffers from a certain illness, thereby filling the binary label $y_1 \in \{0, 1\}$ (for example, setting $y_1 = 1$ if the individual has the illness, and $y_1 = 0$, otherwise). As the doctor repeats this procedure for individuals $2, 3, \dots, K$, the dataset D is generated.

In applications like these, an interesting problem is automatic prediction, that is, prediction without the intervention of the doctor: given a *new* individual with measurements $x \in \mathbf{R}^2$, what should be his binary label $y \in \{0, 1\}$?

Probabilistic model. One way to address the prediction problem is to model how the vector of measurements x impacts the behaviour of the binary label y . A popular such model is logistic regression.

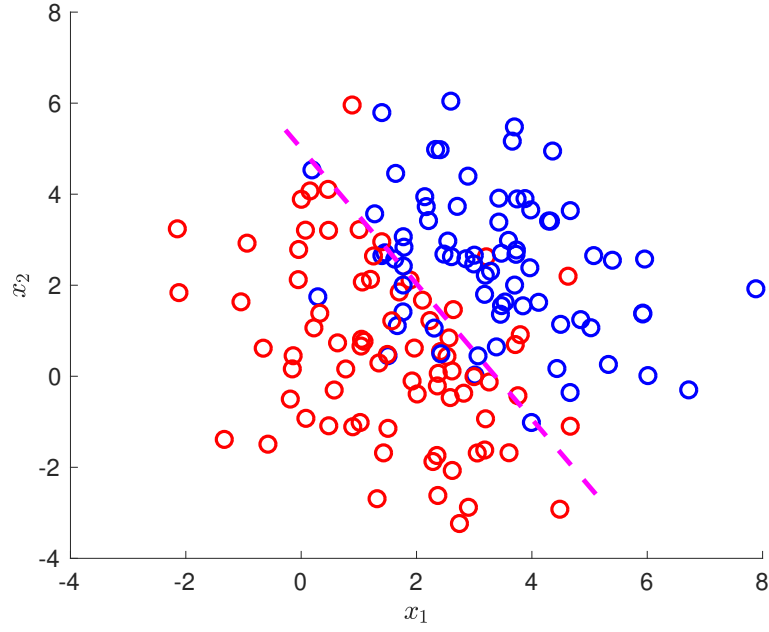


Figure 1: A dataset $D = \{(x_k, y_k) \in \mathbf{R}^2 \times \{0, 1\} : k = 1, \dots, 150\}$ where each element $x_k \in \mathbf{R}^2$ is colored blue if the corresponding binary label y_k is 1, and colored red if y_k is 0. As we move to the right, orthogonally to the magenta line, the blue label becomes more probable; conversely, as we move to the left, orthogonally to the magenta line, the red label becomes more probable. On the magenta line, both labels are equally probable.

In logistic regression, a given x creates a probability distribution for the label y , thus making the label y either more likely to be 0 or 1. Specifically, for a given x , logistic regression assigns the probability of y being 1 to be equal to $\exp(s^T x - r)/(1 + \exp(s^T x - r))$ and the probability of y being 0 to be equal to $1/(1 + \exp(s^T x - r))$. Here, $s \in \mathbf{R}^n$ and $r \in \mathbf{R}$ are the model parameters.

Roughly speaking, logistic regression says that when a vector of measurements x verifies $s^T x - r > 0$, its label y is more likely to be 1; in fact, the more positive is $s^T x - r$, the more likely the label is to be 1 (because the probability $\exp(s^T x - r)/(1 + \exp(s^T x - r))$ converges to 1 as $s^T x - r$ grows). Conversely, if $s^T x < r$, its label y is more likely to be 0. Finally, for vectors x satisfying $s^T x = r$, the label y is equally likely to be 0 or 1. In two dimensions, the equation $s^T x = r$ defines a line; in three dimensions, it defines a plane, and so on. Figure 1 also shows the line $s^T x = r$ in magenta.

To conclude, once we know a value for (s, r) matched to our dataset, we can “solve” the prediction problem easily: we simply output, for given new x , the probabilities of its corresponding y being 1 or 0 as $\exp(s^T x - r)/(1 + \exp(s^T x - r))$ and $1/(1 + \exp(s^T x - r))$, respectively.

But how do we find an (s, r) matched to our dataset?

The optimization problem. An (s, r) matched to our dataset can be found by a famous estimation principle in probability—the maximum likelihood (ML) principle. In ML, we search for the (s, r) that maximizes the likelihood of the parameters (s, r) given the dataset D . For the generic logistic regression model in dimension n (that is, in the dataset D , each vector x_k has n measurements), ML leads to the following optimization problem:

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

By solving problem (1), we find the (s, r) best matched to our dataset.

Task 1. Show that the objective function f in problem (1) is convex.

2 Gradient method

The dataset for Figure 1 is in the MATLAB file `data1.mat`, which contains a matrix \mathbf{X} of size $n \times K$ (where $n = 2$ and $K = 150$) and a vector \mathbf{Y} of length K : column k of the matrix

\mathbf{X} is x_k , and entry k of the vector \mathbf{Y} is y_k .

Task 2. Solve problem (1) for the dataset `data1.mat` using the gradient method given in slide 48 of the set of slides “Part 2: unconstrained optimization” (available in the course webpage):

- Regarding slide 48, use the initialization $s_0 = (-1, -1)$ and $r_0 = 0$, and set $\epsilon = 10^{-6}$ for the stopping criterion;
- For the backtracking subroutine, use the values of $\hat{\alpha}$, γ , and β given in slide 47.

Report the values of s and r that you obtain; plot the line $\{x \in \mathbf{R}^2: s^T x = r\}$ along with the dataset; and plot the norm of the gradient of the cost function along iterations.

To check your code, we now give the answers for task 1 (which you should reproduce). The gradient method finds $s = (1.3495, 1.0540)$ and $r = 4.8815$. Figure 2 plots the line $\{x \in \mathbf{R}^2: s^T x = r\}$ along with the dataset `data1.mat`, and Figure 3 plots the norm of the gradient along iterations.

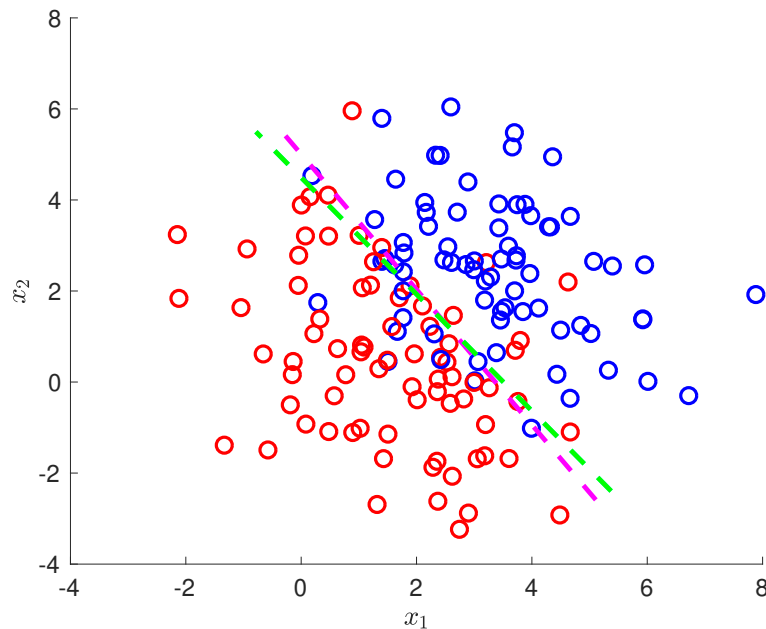


Figure 2: The dataset `data1.mat` with the green line $\{x \in \mathbf{R}^2: s^T x = r\}$ superimposed, where (s, r) is the solution of problem (1) found by the gradient method (Task 2).

Task 3. Redo task 2 for the dataset `data2.mat`.

The dataset `data2.mat` is given in Figure 4.

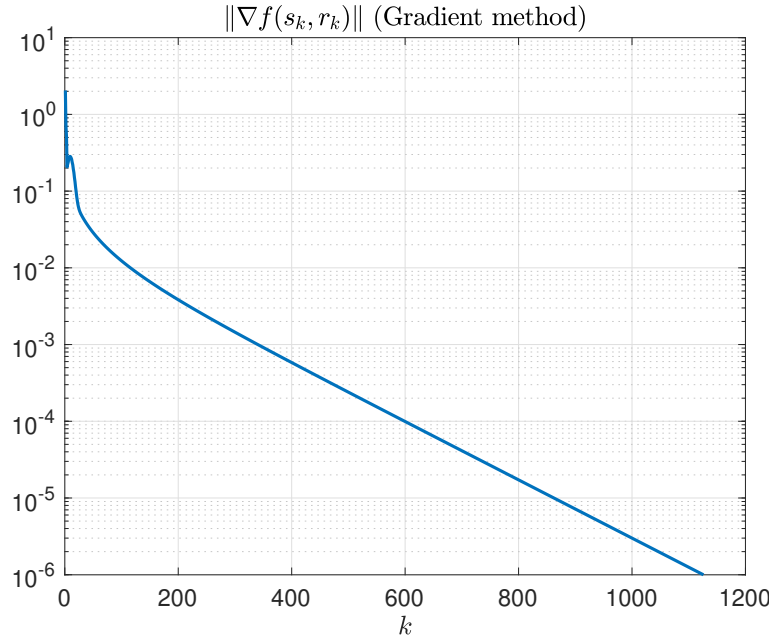


Figure 3: Norm of the gradient along the iterations of the gradient method, when the gradient method is applied to problem (1) with dataset `data1.mat` (Task 2).

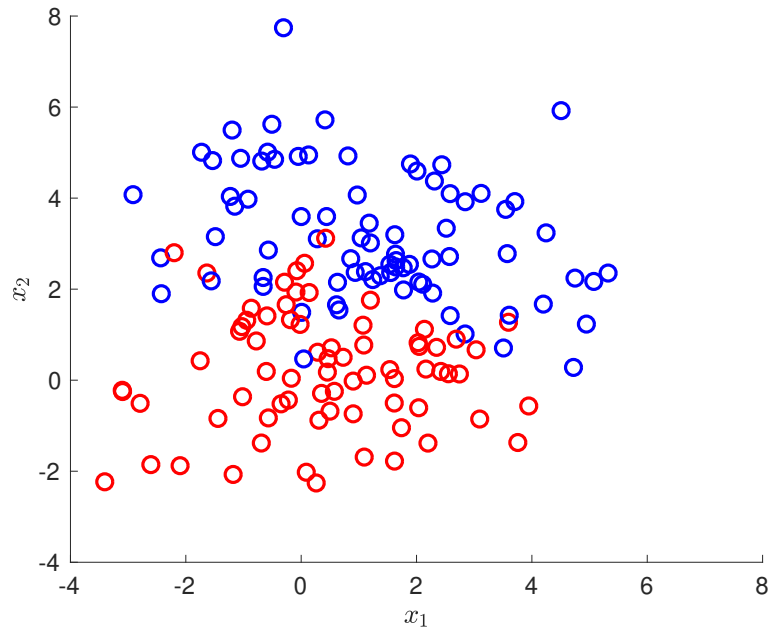


Figure 4: The dataset `data2.mat` for task 3.

Task 4. Use the gradient method to solve problem (1) for the datasets `data3.mat` (for which, $n = 30$ and $K = 500$) and `dataset4.m` (for which, $n = 100$ and $K = 8000$):

- Regarding slide 48, use the initialization $s_0 = (-1, -1, \dots, -1)$ and $r_0 = 0$, and set $\epsilon = 10^{-6}$ for the stopping criterion;
- For the backtracking subroutine, use the values of $\hat{\alpha}$, γ , and β given in slide 47.

3 Newton method

We now solve problem (1) with the Newton method given in slide 81 of the set of slides “Part 2: unconstrained optimization.”

In each iteration of Newton’s method, the main computation is solving a linear system of the form $Hd = -g$, where d is the Newton direction, and H and g are the Hessian and the gradient of the cost function at the current iterate. For an efficient implementation of Newton’s method in MATLAB, both the gradient and Hessian should be expressed through matrix operations, to avoid time-consuming `for` loops (in fact, this also applies for the gradient method in the previous section: the gradient should be computed preferably without `for` loops). The next task helps finding compact matrix expressions for both the gradient

and Hessian.

Task 5. Let $\phi : \mathbf{R} \rightarrow \mathbf{R}$ be a twice-differentiable function. Suppose the function $p : \mathbf{R}^3 \rightarrow \mathbf{R}$ is given by

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

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

- (a) Show that the gradient of p at x is given by $\nabla p(x) = Av$, where $A = [a_1 \ a_2 \ \cdots \ a_K]$ and

$$v = \begin{bmatrix} \dot{\phi}(a_1^T x) \\ \dot{\phi}(a_2^T x) \\ \vdots \\ \dot{\phi}(a_K^T x) \end{bmatrix}.$$

(Note that $\dot{\phi}$ is the derivative of ϕ .)

- (b) Show that the Hessian of p at x is given by $\nabla^2 p(x) = ADA^T$, where D is the diagonal matrix

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

(Note that $\ddot{\phi}$ is the second derivative of ϕ .)

Task 6. Solve problem (1) for the datasets `data1.mat`, `data2.mat`, `data3.mat`, and `data4.mat` using the Newton method:

- Regarding slide 81, use the initialization $s_0 = (-1, -1, \dots, -1)$ and $r_0 = 0$, and set $\epsilon = 10^{-6}$ for the stopping criterion;
- For the backtracking subroutine, use the values of $\hat{\alpha}$, γ , and β given in slide 47.

For each dataset, plot the norm of the gradient of the cost function along iterations, and the values of the stepsizes (α_k) determined by the backtracking subroutine.

So you can check your code, we give the answers for the dataset `data1.mat` (which you should reproduce): Figure 5 plots the norm of the gradient along iterations, and Figure 6 plots the stepsizes found by the backtracking subroutine.

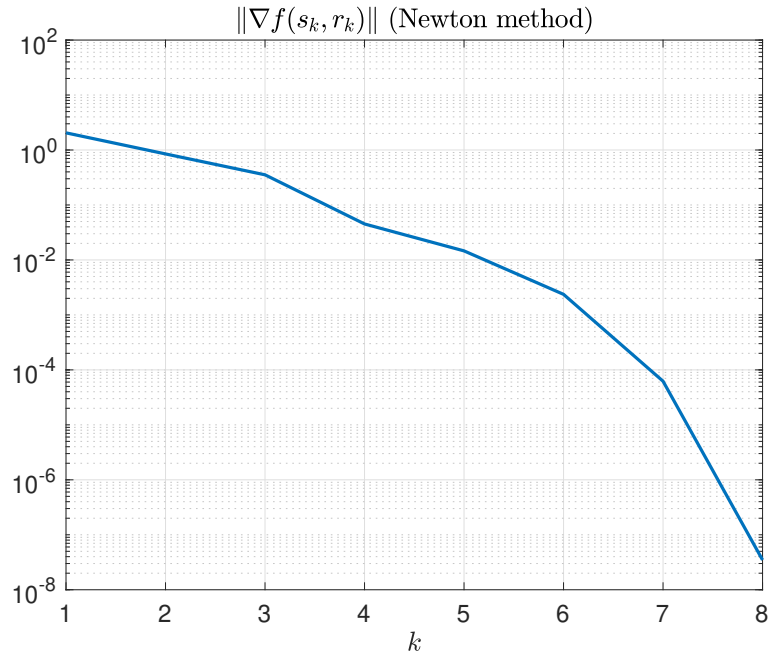


Figure 5: Norm of the gradient along the iterations of the Newton method, when the Newton method is applied to problem (1) with dataset `data1.mat` (Task 6).

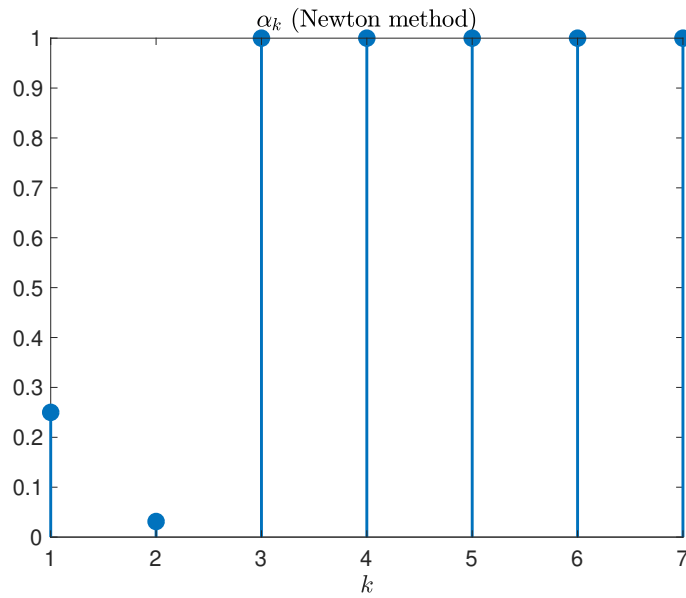


Figure 6: Values of the stepsizes along the iterations of the Newton method, when the Newton method is applied to problem (1) with dataset `data1.mat` (Task 6).

Task 7. Comment on all the results you obtained. In particular, compare the relative behaviour of the gradient and Newton method on problem (1).