

UNIVERSITÀ DEL SALENTO



Facoltà di Ingegneria  
Corso di Laurea Magistrale in Computer Engineering

Advanced Control Techniques Project  
**Multinomial Logistic Regression**

Professor: **Giuseppe Notarstefano**

Students: **Alfarano Gianluca,**  
**Basile Davide,**  
**Capoccia Leonardo,**  
**Larini Ludovico**

Academic year 2017/2018



# Abstract

In this report, it will be shown how to solve a Multinomial Logistic Regression problem (also known as *Softmax Regression*) using a distributed method. The sub-gradient method has been used to distribute calculations among a configurable number of agents. A portion of the dataset is given to each agent and used to minimize a cost function related to the portion of the dataset in its possession. A Python3 program has been written supported by several libraries. The most important library is MPI4py which provides implementation of the MPI specifics. MPI is a communication protocol between nodes that executes the same program in parallel. We compared the results obtained through this program with MATLAB script which executes all operation with a centralized architecture ( meaning that all operations are executed by a single agent).

Several graphs are plotted in order to see the graph of communicating agents, how the consensus is reached and how the cost diminishes to a minimum.

# Contents

<b>Introduction</b>	<b>6</b>
<b>1 Problem and its implementation</b>	<b>8</b>
1.1 Theory of the problem . . . . .	8
1.1.1 Distributed Subgradient Methods for Multi-Agent Op- timization . . . . .	8
1.1.2 Multinomial Logistic Regression . . . . .	9
1.1.3 Pseudocode . . . . .	10
1.2 Code Implementation . . . . .	10
1.2.1 Creation of the communication graph . . . . .	11
1.2.2 Loss function calculation . . . . .	12
1.2.3 Gradient calculation . . . . .	13
1.2.4 Reaching consensus . . . . .	13
<b>2 Results of simulations</b>	<b>17</b>
2.1 Minimization of the Softmax function . . . . .	17
2.2 Dataset, graph description and minimization function . . . . .	20
2.3 Performance . . . . .	20
2.3.1 Number of nodes . . . . .	20
2.3.2 Epsilon . . . . .	20
2.3.3 Learning rate . . . . .	21
2.3.4 Initial condition . . . . .	21
2.4 Four Nodes Setup . . . . .	21
2.4.1 Fixed step-size . . . . .	21
2.4.2 Diminishing step-size . . . . .	22
2.5 Seventeen Nodes Setup . . . . .	22
<b>Conclusions</b>	<b>24</b>
<b>Bibliography</b>	<b>25</b>

# List of Figures

1.1	Graph topology resulting from the $-k/2$ option. . . . .	11
2.1	Communication graph used for the simulation. . . . .	17
2.2	Value of the normalized cost function in a logarithmic scale over all the iterations. . . . .	18
2.3	Value of the normalized cost function in a linear scale. . . .	19
2.4	Consensus on the third component of $\theta_1$ . . . . .	19

# Introduction

## Motivations

In the past there was a single *Mainframe* that executed all digital computations. Years after, with the creation of the Personal Computer, more people could execute the same operations in private. Today's *Microcontrollers* allow to make smart all kind of devices. More algorithms have been created to connect these devices to distribute computations and tasks.

In general, this approach is useful because it allows cooperation between agents, to reach a common goal. In this case, the ability to split the workload in to several agents guarantees faster execution times. The Multinomial Logistic Regression generalizes the logistic regression for a multiclass problems, with more than two possible discrete outcomes, where each observation has several features.

## Contributions

This work implements a scenario in which there are some agents that estimate a cost function using their own information and those of other agents; they use a *Distributed Sub-gradient Method* to update their own estimates and, in particular, they resolve a *Multinomial Logistic Regression* problem.

The first step was to create a MATLAB script that solves a quadratic optimization problem using a gradient method. In this script it has been *simulated* a distributed approach in order to test the effectiveness of the method saw in Distributed Subgradient Methods for Multi-Agent Optimization [1]. MATLAB was only used in the prototype phase, because it doesn't allow real parallel computations. After the confirmation that the approach was working, it started the implementation of the algorithm in Python3 using the Message Passing Interface.

The quadratic form used was

$$x^T Q x + r^T x \tag{I}$$

This form was used only to test the convergence of the method.

After having established that, considering that the Supervised Learning problem known as Multinomial Logistic Regression gives as output a probability value for every considered class, it was necessary to implement the softmax function.

$$g_i(\theta) := \sum_{k=1}^K 1\{y^{(i)} = e_k\} \log \left( \frac{\exp(\theta^{(k)\top} x^{(i)})}{\sum_{k=1}^K \exp(\theta^{(j)\top} x)} \right) \quad (\text{II})$$

[3].

The formula saw in 1.7 was minimized in a centralized setup using MATLAB, and the result was compared to the one of the Python distributed program. It was seen that the algorithm was converging, several tests were executed to find the optimal values for the step-size.

The present work is divided into two chapters. In Chapter 1, it is introduced the theory behind the problem, it is shown and commented the code. In Chapter 2 there are the results of simulations with some considerations.

# Chapter 1

## Problem and its implementation

### 1.1 Theory of the problem

#### 1.1.1 Distributed Subgradient Methods for Multi-Agent Optimization

In this problem there are  $m$  agents that cooperatively minimize a common additive cost. The general optimization problem is:

$$\min \sum_{i=1}^m f_i(x) \quad \text{subject to } x \in \mathbb{R}^n, \quad (1.1)$$

where  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$  is the cost function of agent  $i$ , known only by the agent  $i$ , and  $x \in \mathbb{R}^n$  is the decision variable. It is assumed that:

- the cost function is convex;
- the agents are distributed over a time-varying topology;
- the graph  $(V, E_\infty)$  is connected, where  $E_\infty$  is the set of edges  $(j, i)$  representing agent pairs communicating directly an indefinite number of times;
- there is no communication delay.

Each agent  $i$  generates and maintains an estimate of the optimal decision variable based on the gradient of its own cost function and the exchanged estimates of its neighbours at discrete times  $t_0, t_1, t_2, \dots$ . A weights matrix is associated to the communication graph, and its elements are defined as follows:

$$\begin{cases} > 0, & \text{if } (i, j) \in E(k) \\ 0, & \text{otherwise} \end{cases}$$



Each agent  $j$  that receives an information by a neighbour agent  $i$  weights this information by multiplying it by  $a_j^i(k) \in \mathbb{R}$ . Then the estimate of the optimal solution is updated according to:

$$x^i(k+1) = \sum_{j=1}^m a_j^i(k) x^j(k) - \alpha^i(k) d_i(k) \quad (1.2)$$

where  $\alpha^i(k) > 0$  is the step-size used by agent  $i$  and the vector  $d_i(k)$  is a sub-gradient of the objective function  $f_i(x)$  calculated at  $x = x^i(k)$ . [1]

### 1.1.2 Multinomial Logistic Regression

The problem to be solved is a Supervised Learning problem called Multinomial Logistic Regression, also known as Softmax Regression, and it generalizes the more common Logistic Regression. The difference between them is that in the former there are several classes to be considered, in the latter, there are only two classes (or equivalently a binary class).

The problem to be solved is to find a set of coefficients based on a given dataset to predict the belonging class for an unseen set of features, while minimizing a cost function. The dataset is composed of  $N$  labelled examples  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ . Each  $x^{(i)} \in R^{d_x}$  for  $i = 1, \dots, N$  is composed of some features which represent the value upon which is based the estimation of the belonging class, while  $y^{(i)} \in R^{d_y}$  is the belonging class for the  $i$ -th example, and can be a value in  $\{1, \dots, K\}$ .

Given a single training example  $(x^{(i)}, y^{(i)})$ , the definition of the cost function is:

$$f_i(\omega) := \left\| h_\omega(x^{(i)}) - y^{(i)} \right\|^2 \quad (1.3)$$

where the  $\omega \in R^{d_x}$  are the weights of the hypothesis function  $h_\omega$ . The overall cost function can be defined as:

$$f(\omega) := \sum_{i=1}^N f_i(\omega) \quad (1.4)$$

The problem is solved by finding the solution of the following optimization problem:

$$\omega^* := \arg \min_{\omega} f(\omega) \quad (1.5)$$

In the Multinomial Logistic Regression, a common choice for the hypothesis function is the following:

$$h_\theta = \frac{1}{\sum_{j=1}^K \exp(\theta^{(j)\top} x)} \begin{bmatrix} \exp(\theta^{(1)\top} x) \\ \vdots \\ \exp(\theta^{(K)\top} x) \end{bmatrix} \quad (1.6)$$

where the weights  $\omega = \theta = (\theta^{(1)}, \dots, \theta^{(K)}) \in R^{d_x \times K}$ .

It can be shown that the previous problem is equivalent to the following:

$$\theta^* = \arg \min_{\theta} - \sum_{i=1}^N g_i(\theta) \quad (1.7)$$

with

$$g_i(\theta) := \sum_{k=1}^K \mathbf{1}\{y^{(i)} = e_k\} \log \left( \frac{\exp(\theta^{(k)\top} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)\top} x)} \right) \quad (1.8)$$

with  $\mathbf{1}\{\cdot\}$  being the *indicator function*. [2]

### 1.1.3 Pseudocode

---

#### Algorithm 1

---

- 1: *Stop Rules:*
  - 2:  $\|\theta_{k+1} - \theta_k\| \leq \varepsilon \quad \varepsilon \text{ fixed}$
  - 3: Number of maximum iterations reached
  - 4: **Start:**
  - 5: Fix initial conditions for each node  $\theta_i(0) = [0 \quad \dots \quad 0]^T$
  - 6: Define the Adjacency Matrix, Weights Matrix,  $\alpha^i = \alpha$  constant for each iteration
  - 7: **while** No stop rule is true, each node  $i$  does: **do**
  - 8:     calculate  $\nabla f_i$
  - 9:     **for** each neighbour  $j$  **do**
  - 10:          $\theta_i(k+1) = \theta_i(k+1) + a_j^i(k)\theta^j(k)$
  - 11:     **end for**
  - 12:      $\theta_i(k+1) = \theta_i(k+1) - \alpha \nabla f_i$
  - 13: **end while**
  - 14: **Result:**
  - 15: Each node  $i$  should converge to  $\theta^*$
  - 16: The minimum of function is  $\sum_{i=1}^m f_i(x^*)$
- 

## 1.2 Code Implementation

In this section it will be explained the Python implementation that solves the given problem. The solution has been implemented using the following libraries:

- *numpy* which computes all vector and matrix operations such as transposition, product, division, summation

- *networkx* which creates and manages the adjacency matrices for all the agents
- *matplotlib* that plots the results of the tests and the communication graph.

### 1.2.1 Creation of the communication graph

The communication graph is created by the *createAdjM* function. It creates a graph that is strongly connected and that has a cycle-like topology. In particular, from the command line is possible to pass a number (through the *-k* option) to set the number of connection that each agent has. Suppose that the user set the option *-k 2*, then the resulting graph will have a topology as in Fig. 1.1.

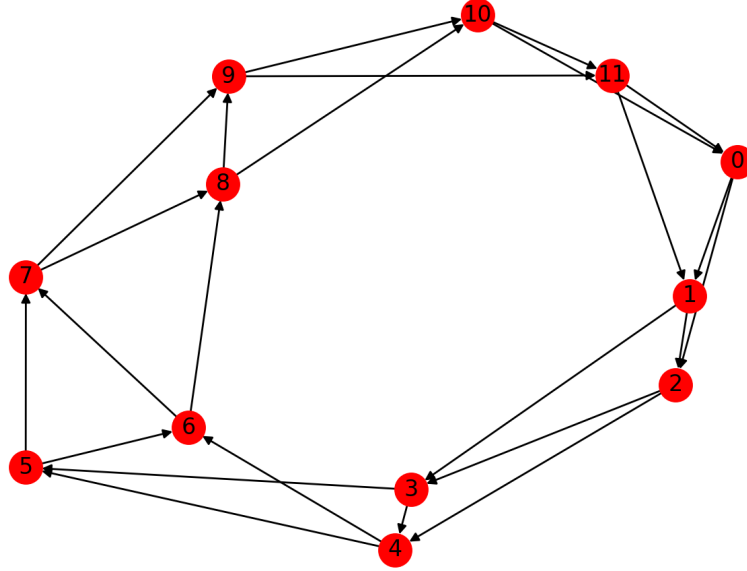


Figure 1.1: Graph topology resulting from the *-k 2* option.

It can be seen that each agent  $i$  has as out-neighbours the agents  $i+1$  and  $i+2$  and as in-neighbours the agents  $i-1$  and  $i-2$ .

This can be formalized as follows:

$$N_i^{\text{OUT}} = \{i+1, i+2, \dots, i+k\} \quad (1.9)$$

$$N_i^{\text{IN}} = \{i-1, i-2, \dots, i-k\} \quad (1.10)$$

It's important to treat the numbers of the agents as if they were ordered with an agent following another one. This means that to the last agent  $N$  follows the agent 0.

### 1.2.2 Loss function calculation

The loss function (described in Eq. 1.7) is implemented in the Python function called *loss\_softmax*. It has three parameters:

- *all\_theta*: the matrix with the current estimation of  $\theta^*$
- *category\_count*: the number of categories in the dataset. In the IRIS dataset is 3
- *personal\_dataset*: represents the dataset accessible to the considered agent

The only addition to the function is the calculation of the *const\_to\_subtract*, a constant used to prevent the overflow. The property used to justify the use of this method is the following:

$$\frac{e^a}{e^b} = \frac{e^{a-c}}{e^{b-c}} \quad (1.11)$$

where  $c$  represent the constant to subtract and it is calculated as

$$c = \max \left\{ \theta^{1^\top} x, \dots, \theta^{K^\top} x \right\} \quad (1.12)$$

```

1  def loss_softmax(all_theta, category_count, personal_dataset):
2      the_sum = 0
3      for index in range(0, len(personal_dataset)):
4          denominator = 0
5          const_to_subtract = find_const_to_subtract(all_theta,
6                                                       personal_dataset[index][0:4])
7
8          for theta in all_theta:
9              denominator = denominator + np.exp(np.dot(theta,
10                                                           personal_dataset[index][0:4]) - const_to_subtract)
11
12         for category in range(0, category_count):
13             if category == personal_dataset[index][4]:
14                 _exp = np.exp(np.dot(all_theta[category],
15                                       personal_dataset[index][:4]) - const_to_subtract)
16                 _log = np.log(np.divide(_exp, denominator))
17                 the_sum = the_sum - _log
18
19         return the_sum

```

This function could be omitted because it is not needed for solving the problem. It is used only for debug purpose as it is needed to draw the graphics of the overall loss function and to demonstrate that the cost gets minimized.

### 1.2.3 Gradient calculation

Given the nature of the method used to solve the problem, this is the most important piece of software. The gradient of the cost function has been calculated to be:

$$\nabla_{\theta^{(k)}} J(\theta) = - \sum_{i=1}^m \left[ x^{(i)} \left( 1 \{y^{(i)} = k\} - \frac{\exp(\theta^{(k)\top} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)\top} x^{(i)})} \right) \right] \quad (1.13)$$

The Python function *gradient\_softmax* calculate the value of the equation 1.13, using the same method saw in the previous section to prevent overflow.

```

1  def gradient_softmax(all_theta , category_count , dimensions ,
2      personal_dataset , normalized):
3
4      thetas = np.zeros(dimensions)
5
6      if normalized:
7          m = len(personal_dataset)
8      else:
9          m = 1
10
11     for index in range(0, len(personal_dataset)):
12         denominator = 0
13         const_to_subtract = find_const_to_subtract(all_theta ,
14             personal_dataset[index][0:4])
15
16         for theta in all_theta:
17             denominator = denominator + np.exp(np.dot(theta ,
18                 personal_dataset[index][:4]) - const_to_subtract)
19
20         for category in range(0, category_count):
21             coeff = 0
22             if category == personal_dataset[index][4]:
23                 coeff = 1
24             _exp = np.exp(np.dot(all_theta[category] , personal_dataset
25                 [index][:4]) - const_to_subtract)
26             coeff = coeff - np.divide(_exp , denominator)
27             thetas[category] = thetas[category] - ((1 / m) * np.
28                 multiply(personal_dataset[index][:4] , coeff))
29
30     return thetas

```

### 1.2.4 Reaching consensus

In this section in will be explained how the consensus is reached.

The distributed sub-gradient method is an iterative algorithm, so it make use of a cycle in which the calculations are made. In the Python script there is a constant, called *MAX\_ITERATIONS*, that represents the limit of iterations that can be made.

The first thing done in the cycle is the calculation of the step-size, basing the computation on the type chosen from the command line argument.

```

1  if alpha_type == "diminishing":
2      alpha = psi_coefficient * (1 / tt) ** alpha_coefficient
3  else:
4      alpha = alpha_coefficient

```

If a diminishing step-size is chosen, then the formula implemented is

$$\alpha_t = \psi \left( \frac{1}{t} \right)^\gamma \quad (1.14)$$

otherwise, it is used a constant step-size defined as follows:

$$\alpha_t = \gamma \quad (1.15)$$

with  $t$  indicating the current iteration.

After this, the state of each agent is updated using the current state and the states of the in-neighbours. The new calculated state is then sent to all the out-neighbours. This is done with the following piece of code:

```

1  # Update with my previous state
2  u_i = np.multiply(XX[tt - 1], weight)
3
4  # Send the state to neighbors
5  for node in graph.successors(rank):
6      world.send(XX[tt - 1], dest=node)
7
8  # Update with state of all nodes before me
9  for node in graph.predecessors(rank):
10     u_i = u_i + world.recv(source=node) * weight

```

The variable used here are:

- $u_i$ : represents the current estimate of the state of the agents
- *world*: a variable used for communicating between nodes. It is made available by the MPI environment
- *graph*: a variable representing the communication graph. It is created using the functions in the *networkx* library. It has two main functions that can be called: *successors(node)* that returns a list with the out-neighbours of the agent *node* and *predecessors(node)* that returns the in-neighbours
- *rank*: represent the identifier of each agent. It is made available by the MPI environment and it is an integer number, that goes from 0 to *numberofagents* - 1

After this, the value of the gradient is calculated, using the function name passed as command line argument:

```

1  # Variable for storing the gradient value
2  gradient = 0
3
4  if function_name == "softmax":
5      gradient = func.gradient_softmax(XX[tt - 1], category_n,
6                                       dimensions, personal_dataset, normalized)
7
8  elif function_name == "quadratic":
9      gradient = func.gradient_quadratic(XX[tt - 1], category_n,
10                                       dimensions, personal_dataset, Q, r)
11
12 elif function_name == "exponential":
13     gradient = func.gradient_exponential(XX[tt - 1], category_n,
14                                         dimensions, personal_dataset)
15
16 # Multiply the gradient value by alpha
17 grad = np.multiply(alpha, gradient)

```

In particular, a variable is firstly defined to store the value of the gradient that will be calculated. After computing it, the gradient is multiplied by the step-size and the result is stored into the *grad* variable.

After this phase, the state is updated using the gradient descent step:

```

1  for i in range(0, dimensions[0]):
2      u_i[i] = np.subtract(u_i[i], grad[i])
3
4  # Store my new state
5  XX[tt] = u_i

```

Since the *u\_i* and *grad* variable are matrices, then each component of the gradient (multiplied by the step-size) is subtracted from the current estimate of the state. The result is then stored into an array *XX*, which keeps in memory all the states at each iteration.

Consequently, the loss value is calculated at each iteration using the right function. The value is then stored into an array that keeps in memory all the losses at each iteration:

```

1  if function_name == "softmax":
2      losses[tt] = func.loss_softmax(XX[tt], category_n,
3                                     personal_dataset)
4
5  elif function_name == "quadratic":
6      losses[tt] = func.loss_quadratic(XX[tt], category_n,
7                                       dimensions, personal_dataset, Q, r)
8
9  elif function_name == "exponential":
10     losses[tt] = func.loss_exponential(XX[tt - 1], category_n,
11                                         dimensions, personal_dataset)

```

After this, the stop conditions are checked:

```

1      # Checking epsilon reached condition
2      if np.linalg.norm(np.subtract(XX[tt], XX[tt - 1])) < epsilon:
3          buff = True
4
5      # Rank 0 get all epsilon and check if all reached it
6      buffer = world.gather(buff, root=0)
7
8      # If true it set epsilon reached
9      if rank == 0:
10         if False not in buffer:
11             epsilon_reached = True
12
13     # Send epsilon reached to all agents
14     epsilon_reached = world.bcast(epsilon_reached, root=0)
15
16     # Check if all agent have reached epsilon condition and then exit
17     # from loop
18     if epsilon_reached:
19         if rank == 0:
20             print("Exiting_at_iteration_", tt, "/", MAX_ITERATIONS, "
                Condition_on_epsilon_reached")
                break

```

This is done in different steps:

1. The following condition is checked:

$$\|\theta_{t+1} - \theta_t\| \leq \epsilon \quad (1.16)$$

If it is met, the variable *buff* is set to true.

2. The agent *0* collects all the values of *buff* of each agent in the *buffer* variable.
3. Agent *0* check if all the conditions in the *buffer* variable are set to true. If they are, then the variable *epsilon\_reached* is set to true.
4. This variable is sent to all the other agents using the *bcast* function, that broadcast the value of the variable to all the nodes in the graph.
5. As last step, all the agents exit from the main cycle if all of them have met the stop condition.

The last things done in the main cycle are to print the iterations to the screen. This is done by the following piece of code:

```

1      if tt in range(0, MAX_ITERATIONS, 100):
2          if rank == 0:
3              print("Iteration_", tt, "/", MAX_ITERATIONS)

```



## Chapter 2

# Results of simulations

To verify the correctness of the software, a couple of simulations were run, divided into two parts:

- Minimization of the Softmax function to prove the convergence to the optimal value
- Application of the minimization algorithm to find the best coefficients to do the best predictions on the test dataset

### 2.1 Minimization of the Softmax function

The simulation has been carried out using a random directed graph made up of 17 agents, as shown in figure 2.1.

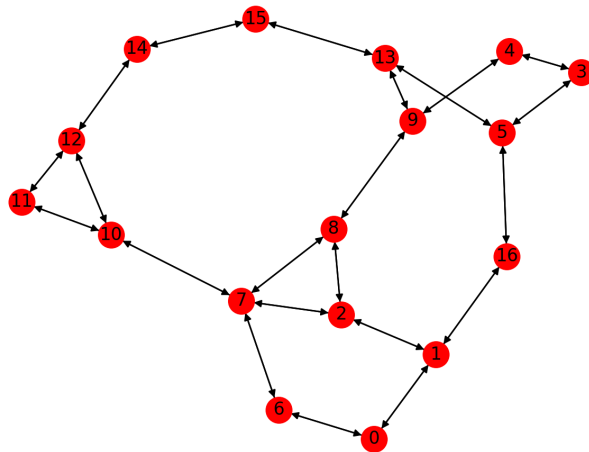


Figure 2.1: Communication graph used for the simulation.

For completing this simulation, a diminishing step-size has been used, with the following formulation:

$$\alpha^k = 0.01 \left( \frac{1}{k} \right)^{0.4} \quad (2.1)$$

with  $k$  indicating the current iteration.

The problem to be solved is the one formulated in 1.7. Let's suppose to indicate the whole dataset as  $S$ . Each agent  $i$  can access to a limited set  $s_i$  of examples contained in  $S$ . To each agent is given a set  $s_i$ , chosen as:

$$S = \bigcup_{i=1}^{17} s_i \quad (2.2)$$

and

$$\bigcap_{i=1}^{17} s_i = \emptyset \quad (2.3)$$

In the simulation, the training set  $S$  is composed by 120 examples, split into sets  $s_i$  of 6 or 7 examples, chosen in order to equation 2.2 and 2.3 hold.

The iterations run are 15000, sufficient enough to see the function approaching the convergence, but not sufficient enough to reach the minimum value of the function, calculated with MATLAB to be 8.945. The figure 2.2 shows in a logarithmic scale how the algorithm tends to minimize the cost value at each iteration.

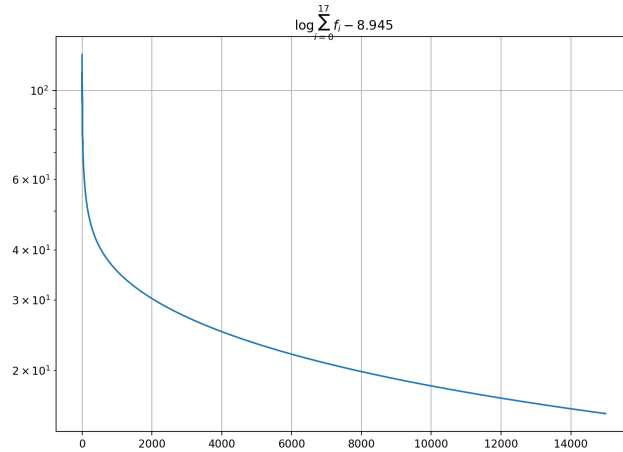


Figure 2.2: Value of the normalized cost function in a logarithmic scale over all the iterations.

The value of the cost function is shown in figure 2.3 in a linear scale too, so it can be seen how the algorithm approaches the minimum value  $\theta^*$  in order to have the minimum of the problem 1.7.

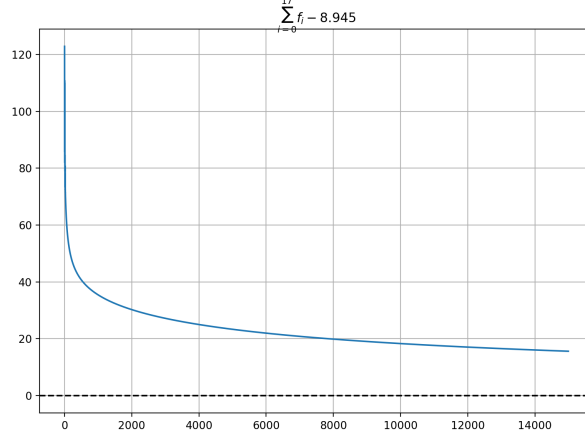
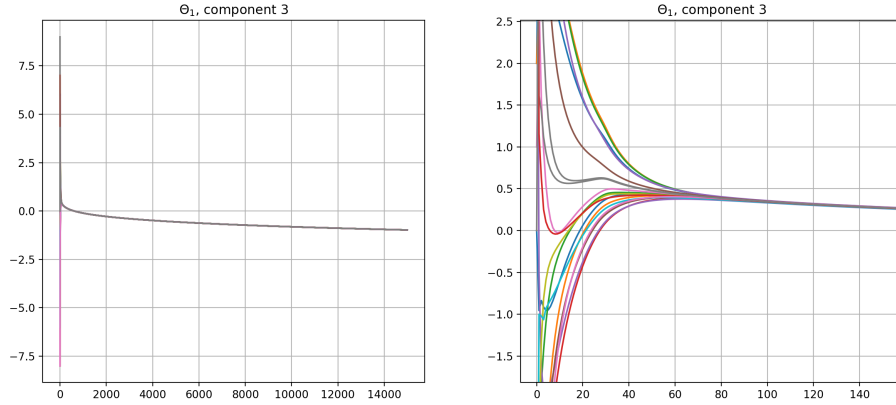


Figure 2.3: Value of the normalized cost function in a linear scale.

As can be seen in figure 2.4, all the agents try to reach the consensus for each component of the local minimum of the problem.

To better see the consensus, in this simulation the initial condition have been chosen to be random integer number between  $-10$  and  $+10$ . With this choice is more visible the trend of the agents that try to reach consensus.



(a) Overall trend over all the iterations      (b) Detail for the first 140 iterations

Figure 2.4: Consensus on the third component of  $\theta_1$

## 2.2 Dataset, graph description and minimization function

This dataset is composed of 150 instances, 120 of them used for training while the remaining 30 used for testing purpose. The instances contain 3 classes, each representing a type of Iris flower. Every instance has 4 features: sepal and petal length, sepal and petal width, expressed in *cm*. We tried different types of graphs, all of them chosen to be strongly connected. The results discussed in this chapter, if not differently noted, refer to cyclic graphs with a variable number of nodes. The program can minimize all kinds of loss functions. As shown in Chapter 1.2, the quadratic and exponential functions can also be optimized, even if they don't guarantee useful results to apply to our problem. The minimization function used in this chapter is the one described in Chapter 1.1.2.

## 2.3 Performance

There are some key factors that influence *the computational time*, *the numbers of iterations necessary* and *the accuracy of the results*. They are, but not limited to, number of nodes, type of graph, step-size type and value, etc.

### 2.3.1 Number of nodes

The Python program, thanks to the MPI platform, is capable of running on an arbitrary number of nodes. It was tested on as little as 2 nodes to as many as 60 nodes, which means that every node was processing the data of 2 instances (120 instances divided into 60 nodes). The best performances are obtained when the number of nodes corresponds to the number of physical cores of the machine where it runs. When the number of nodes exceeds greatly the number of physical cores, the resources are oversubscribed. In this case, the performances degrade notably as the nodes compete for cache and memory and the processors' schedulers are put in a difficult situation. On a 4-core test machine, a computation with 5000 iterations and 30 nodes is done in 5 minutes. The same machine can do the same number of iterations, but with 60 nodes, in 15 minutes. Therefore, to show the scalability of the software, the following tests will be shown on a 4 nodes set-up and a 17 nodes setup (the same of figure 2.1) .

### 2.3.2 Epsilon

This is a small constant used as stop condition. If the norm of the result of the current calculation differs less than  $\epsilon$  (chosen arbitrary small) from

the previous, the algorithm is stopped. This is used as convergence stop condition. In fact, given the nature of the method used to minimize the cost function, the converge to the optimal value can be reached in a huge amount of time. To prevent this from happening, it has been used this condition.

### 2.3.3 Learning rate

The step-size plays a big role in the speed of convergence of the algorithm. There are 3 kinds of step-size:

1. fixed step-size: it has always the same value, which means  $\alpha^{k+1} = \alpha^k$
2. diminishing step-size: chosen such that  $\alpha^{k+1} < \alpha^k$
3. adaptive step-sizes: chosen with respect to Armijo or (limited) minimization rule.

For reasons not discussed in this paper, it's not possible to use the last type of rules in a distributed problem, so the first two types have been used in these simulations.

### 2.3.4 Initial condition

The agents start from an initial condition  $x_0$ , which represents a guess of the optimal value. The program was tested with random initial conditions for all agents and the consensus was still reached. To guarantee the repeatability of the experiments, the initial condition in the tests is always set to a matrix full of ones.

## 2.4 Four Nodes Setup

### 2.4.1 Fixed step-size

We will first deal with a simple fixed step-size such as  $\alpha^k = \alpha$ , with  $\epsilon = 10^{-3}$  equal to:

$\alpha$	Iteration required	Execution time in s	Wrong guesses
0.5	overflow occurred	-	-
0.1	>10000	>15	2
0.05	>10000	>15	1
0.01	1968	2.8	1
0.005	1558	2.3	1
0.001	2114	3.1	0
0.0005	1619	2.3	2
0.0001	527	0.7	14

These results show a general truth about the step-size. If it is too little, the learning process proceeds in a very slow way and it requires a huge amount of iterations. If the learning rate is too high, the gradient descent will most probably overshoot the minimum and it will not converge. Through trial and error, the step-size 0.001 was chosen. It allows to reach good performances and the best accuracy: only 3.1 seconds are required to obtain 0 errors over the test dataset.

### 2.4.2 Diminishing step-size

The diminishing step-size implemented is chosen to be:

$$\alpha = \psi \left( \frac{1}{k} \right)^\gamma \quad (2.4)$$

Again, several tests were run by tweaking the multiplying constant and the exponent.

$\psi$	$\gamma$	Iteration required	Execution time in s	Wrong guesses
1	0.01	overflow	-	-
1	0.1	overflow	-	-
0.1	0.01	>10000	>15	1
0.1	0.5	349	0.5	0
0.1	0.1	>10000	>16	1
0.01	0.01	1852	2.7	1
0.01	0.1	1259	1.8	1

The best result is obtained in only 349 iterations, with  $\psi = 0.1$  and  $\gamma = 0.5$ . This is done in roughly half of a second, obtaining no error. The accuracy is worst then the previous result, but this is obtained in a fraction of the time needed to obtain 0 error with a fixed step-size. A difference of 2 or 3 seconds may not seem important using the Iris Dataset, with a small amount of calculations. In a set-up where a bigger number of calculations and a bigger dataset are involved, one may appreciate the advantages of this faster approach.

## 2.5 Seventeen Nodes Setup

The following results are the ones for the set-up shown in fig. 2.1.

The computation times are greatly larger than the ones shown earlier. This is because the tests involving 17 agents are carried out on a quad-core machine, which means that MPI oversubscribe the resources and the CPU scheduler are put at hard work.

The optimal step-size is different from the ideal step-size of the 4 nodes

set-up. Only through several tests and experience is possible to find a good step-size to reduce the errors in the best way.

$\psi$	$\gamma$	Iteration required	Execution time in s	Wrong guesses
0.1	0.1	>10000	>440	2
0.5	0.1	>1000	>440	2
0.9	0.1	>10000	>440	2
0.01	0.1	1894	83	2
0.1	0.5	698	30	2
0.5	0.5	1017	44	1
0.9	0.5	2409	107	0
0.01	0.5	318	13	14
0.1	0.01	>10000	>440	1
0.5	0.01	>10000	>440	2
0.9	0.01	>10000	>440	2
0.01	0.01	3717	164	2

TODO: AGGIUNGERE GRAFICI

# Conclusions

In this work it has been resolved a Multinomial Logistic Regression problem using MPI in Python. Each agent use a Distributed Sub-gradient method to update its own estimate of the optimal solution. At the end of the iterations, all the agents converge to the same solution.



# Bibliography

- [1] A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 2009.
- [2] Stanford. Softmax regression. Available at <http://ufldl.stanford.edu/tutorial/supervised/SoftmaxRegression/>.
- [3] Wikipedia. Multinomial logistic regression. Available at [https://en.wikipedia.org/wiki/Multinomial\\_logistic\\_regression](https://en.wikipedia.org/wiki/Multinomial_logistic_regression).