

UNIVERSITÀ DEL SALENTO



Facoltà di Ingegneria
Corso di Laurea Magistrale in Computer Engineering

Advanced Control Techniques Project
**Multinomial Logistic Regression for a
Supervised Learning problem**

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 (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.

Contents

Introduction	6
1 Chapter 1 Problem and its implementation	7
1.1 Theory of the problem	7
1.1.1 Distributed Subgradient Methods for Multi-Agent Op- timization	7
1.1.2 Multinomial Logistic Regression	8
1.1.3 Pseudocode	9
1.2 Code Implementation	9
1.2.1 working.py	11
2 Chapter 2 Results of simulations	14
2.1 Dataset, graph description and minimization function	14
2.2 Performance	14
2.2.1 Number of nodes	15
2.2.2 Epsilon	15
2.2.3 Learning rate	15
2.2.4 Fixed step-size	15
2.2.5 Diminishing step-size	16
Conclusions	17
Bibliography	18

List of Figures

Introduction

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 an infinity of devices. More algorithms have been created to connect these devices to distribute.

This work implements a scenario in which there are some agents that estimate a cost function using their own information and those of the other agents; they use a *Distributed Sub-gradient Method* to update their own estimate and, in particular, they resolve a *Multinomial Logistic Regression*. After some test in *MATLAB*, it is used *MPI* implemented with *Python*.

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

Chapter 1

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:

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

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

- the cost function is convex;
- the agents are distributed over a time-varying topology;
- the graph (V, E_∞) is connected, where E_∞ is the set of edges (j, i) representing agent pairs communicating directly an indefinite number of times;
- there isn't communication delay.

Every agent i generates and maintains estimates of the optimal decision vector based on information concerning its own cost function and exchanges this estimate with its directly neighbors at discrete times t_0, t_1, t_2, \dots . Moreover, each agent i has a vector of weights $a^i(k) \in \mathbb{R}^m$ at any time t_k ; the scalar $a_i^j(k)$ is zero if the agent i doesn't directly communicate with j , else it is the weight assigned from the agent i to the information x^j obtained from

j during the time interval (t_k, t_{k+1}) . The estimates are updated according to the update rule:

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

where $\alpha^i(k) > 0$ is the stepsize used by agent i and the vector $d_i(k)$ is a subgradient of agent i objective function $f_i(x)$ 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 we base 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 values 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_ω . The overall cost function can be defined as:

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

We solve the problem 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 $\omega = \theta = (\theta^{(1)}, \dots, \theta^{(K)}) \in R^{d_x}$.

Using this function, the solution of the problem is given by finding:

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

with

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

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

1.1.3 Pseudocode

Algorithm 1

- 1: *Stop Rules:*
 - 2: $\|\theta_{k+1} - \theta_k\| \leq \varepsilon$ ε 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 ∇f_i
 - 9: **for** each neighbor 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 θ^*
 - 16: The minimum of function is $\sum_{i=1}^m f_i(x^*)$
-

1.2 Code Implementation

Here we are defining and explaining what has been implemented in Python to solve the problem. The solution has been implemented using: *numpy* library which computes all vector and matrix operation such as transposition, product, division, summation; *networkx* which creates and manage adjacency matrices for all agent; *matplotlib* that create and plot result which we have caught during test and training phase. At a first time the

code check if any parameter are send to every agent, then start to create environment, such as get the number of agent, create adjacency matrix and initialize local variable where it will save collected data. The weights depends about how many in-neighbors each agent has, because our adjacency matrix is setted with k in-neighbors that can be setted when program is launched.

The adjacency matrix is created by *createAdjM* function, that make a communication graph using *world_d* as number of mpi agents, *n_edges* setted by default as 1 (if it is not given or bigger than number of agent), *phi* that is simply a phase inserted into graph which just shift the connection, by default is 0. The first *for* cycle sets the number of agents that the system has, then edges are computed as follows:

- The sum $j + k + \text{phi} + 1$ means the position of edge in the graph where item j is the current row, k is the count of agent that the agent j will send a message, phi is the phase and just shift the agent that j will send a message and the "+1"
- The instruction *if e >= world_d* check if the sum is over the number of agents and if is true just subtracts this one
- Instead *if e == j* removes the ipotetical unnecessary self loop.

Loss softmax function, *loss_softmax* gets as parameter the current state (*all_theta*), the cardinality of iris set (*category_count*), the agent dataset (*personal_dataset*) and a *CONSTANT_TO_SUBTRACT*, a constant to prevent overflow.

The algorithm wants to calculate the $\sum_{i=0}^{\text{number_of_agents}} f_i(x)$ and this function calculate the agent i $f(x)$. This can be done as follows: for every row in the personal dataset, the algorithm computes the denominator that is the sum over all exponential. Then it computes, category by category, row by row, the logarithm of ratio between the sum previously calculated and the exponential of the row. After that, there's the summation of all calculated elements.

```

1  def loss_softmax(all_theta , category_count , personal_dataset ,
2      CONSTANT_TO_SUBTRACT):
3      the_sum = 0
4      for index in range(0, len(personal_dataset)):
5          denominator = 0
6          for theta in all_theta:
7              denominator = denominator + np.exp(np.dot(theta ,
8                  personal_dataset[index][0:4]) - CONSTANT_TO_SUBTRACT)
9
10         for category in range(0, category_count):
11

```

```

12         if category == personal_dataset[index][4]:
13             _exp = np.exp(np.dot(all_theta[category],
14                                   personal_dataset[index][:4]) -
15                               CONSTANT.TO.SUBTRACT)
16             _log = np.log(np.divide(_exp, denominator))
17             the_sum = the_sum - _log
18
19     return the_sum

```

The gradient method implemented here, called *gradient_softmax* This is the gradient of softmax equation, that has the same concept as the function mentioned before: calculate in a first time se sum of all exponential theta and then sum for each coefficient, subtract from one that coefficient and finally subtract the dataset normalized with that to the respective theta.

```

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

```

1.2.1 working.py

In this file is implemented the main algorithm where the consensus and the minimization of loss function is calculated.

The complete iris training set is loaded, and environment variables are setted, as the number of agents, the name of *agent_i*. Then the dataset is splitted equal parts to all agents and each one gets its one and print how many row it has.

Every agent creates the same communication directed graph with number of agents and number of in connection. Then state and loss variable are

created, setted to 0 of *MAX_ITERATIONS* size. In order to get wighted message, all in-neighbors are found and the variable *weight* is setted as the average. If the function name inserted is "quadratic" Q and r variables are created randomly. *epsilon_reached* and *buff* are epsilon checker variable that say to agent when they have to exit from the loop (if epsilon is reached).

This *for* cycle calculates the consensus. For every *iters*, is calculated the new "diminishing" alpha; message(s) are sent and received, then local variable are weighted. In order to solve consensus to desiderated function, the three implemented ones are inserted into an *if* clause, and the right function is called. After that, the calculated gradient is multiplied by alpha and then new state is calculated. Then *lossfunction* is called and if $\|XX[tt] - XX[tt - 1]\| \leq \epsilon$ then buff is true and rank 0 check if all agent has reched the epsilon condition; if true it sends a broadcast message with a *true* value to say to all that the cycle is done.

```

1  for tt in range(1, MAXITERATIONS - 1):
2
3      if alpha_type == "diminishing":
4          alpha = psi_coefficient * (1 / tt) ** alpha_coefficient
5      else:
6          alpha = alpha_coefficient
7
8      # Update with my previous state
9      u_i = np.multiply(XX[tt - 1], weight)
10
11     # Send the state to neighbors
12     for node in adj.successors(rank):
13         world.send(XX[tt - 1], dest=node)
14
15     # Update with state of all nodes before me
16     for node in adj.predecessors(rank):
17         u_i = u_i + world.recv(source=node) * weight
18
19     # Go in the opposite direction with respect to the gradient
20     gradient = 0
21
22     if function_name == "softmax":
23         gradient = func.gradient_softmax(XX[tt - 1], category_n,
24                                         dimensions, personal_dataset, CONSTANT.TO.SUBTRACT)
25
26     elif function_name == "quadratic":
27         gradient = func.gradient_quadratic(XX[tt - 1], category_n,
28                                         dimensions, personal_dataset, Q, r)
29
30     elif function_name == "exponential":
31         gradient = func.gradient_exponential(XX[tt - 1], category_n,
32                                         dimensions, personal_dataset, CONSTANT.TO.SUBTRACT)
33
34     #print(gradient)
35
36     grad = np.multiply(alpha, gradient)
37
38     for i in range(0, dimensions[0]):
39         u_i[i] = np.subtract(u_i[i], grad[i])
40
41     # Store my new state
42     XX[tt] = u_i

```

```

40
41     if function_name == "softmax":
42         losses[tt] = func.loss_softmax(XX[tt], category_n,
43                                         personal_dataset, CONSTANT_TO_SUBTRACT)
44
45     elif function_name == "quadratic":
46         losses[tt] = func.loss_quadratic(XX[tt], category_n,
47                                         dimensions, personal_dataset, Q, r)
48
49     elif function_name == "exponential":
50         losses[tt] = func.loss_exponential(XX[tt - 1], category_n,
51                                         dimensions, personal_dataset, CONSTANT_TO_SUBTRACT)
52
53     # Checking epsilon reached condition
54     if np.linalg.norm(np.subtract(XX[tt], XX[tt - 1])) < epsilon:
55         buff = True
56
57     # Rank 0 get all epsilon and check if all reached it
58     buffer = world.gather(buff, root=0)
59
60     # If true it set epsilon reached
61     if rank == 0:
62         if False not in buffer:
63             epsilon_reached = True
64
65     # Send epsilon reached to all agents
66     epsilon_reached = world.bcast(epsilon_reached, root=0)
67
68     # Check if all agent have reached epsilon condition and then exit
69     # from loop
70     if epsilon_reached:
71         if rank == 0:
72             print("Exiting_at_iteration_", tt, "/", MAX_ITERATIONS, "
73                   Condition_on_epsilon_reached")
74             sys.stdout.flush()
75
76         break
77
78     if tt in range(0, MAX_ITERATIONS, 100):
79         if rank == 0:
80             print("Iteration_", tt, "/", MAX_ITERATIONS)
81             sys.stdout.flush()
82
83     ITERATION_DONE = tt

```

When consensus is reached, the theta values are printed to user, then all data are sent to rank 0 for centralized calculations and plot result.

Chapter 2

Chapter 2 Results of simulations

The software described in the previous chapter was used to solve a Multinomial Logistic Regression problem. Specifically, we classified the data of the Iris Dataset.

2.1 Dataset, graph description and minimization function

This dataset is composed of 150 instances, 120 used for training, the rest for tests. The instances contain 3 classes, each representing a type of Iris flower. Every instance has 4 features, sepal length, sepal width, petal length, petal width expressed in *cm*. We tried different types of graphs. These graphs are all strongly connected and the weight matrices for the nodes are doubly stochastic, as per the assumptions of convergence of the algorithm described in Chapter 1. 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 used. These last 2 functions don't guarantee useful results and/or convergence. The minimization function used in this chapter is the one described in Chapter 1.1.2, softmax.

2.2 Performance

There are some key factors that influence *the computational time, the numbers of iterations necessary and the accuracy of the results*.

2.2.1 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, the following tests will be shown on a 4 nodes setup.

2.2.2 Epsilon

This is a small constant used as stop condition. If the result of the current calculation differs less than epsilon from the previous, the algorithm is stopped.

2.2.3 Learning rate

The step-size alpha plays a big role in the speed of convergence of the algorithm. There are 3 kinds of step-size. Fixed, diminishing and adaptive (Armijo). For reasons not discussed in this paper, it's not possible to use Armijo rule in a distributed problem.

2.2.4 Fixed step-size

We will first deal with a simpler fixed step-size. With an epsilon equal to 0.001 (10^{-3}):

Value of fixed step-size	Iteration required	Execution time in s	Wrong guesses o.30
0.5	overflow	alpha too big	-
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 identified, it allows reaching good performance and accuracy. In fact, in only 3.1 seconds, we can make predictions with no errors, using our 30 instances test dataset.

2.2.5 Diminishing step-size

The diminishing step-size implemented in the code is in this form: TODO: SISTEMARE

$$\alpha = \text{const} \left(\frac{1}{tt} \right) \exp \quad (2.1)$$

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

psi.coeff	alpha_exp coefficient	Iteration required	Execution time in s	Wrong guesses o. 30
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 0.1 and 0.5 as multiplicative constant and exponent respectively. 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 setup where a bigger number of calculations and bigger dataset are involved, one may appreciate the advantages of this faster approach.

TODO: GRAFICI

Conclusions

In this work it has been resolved a Multinomial Logistic Regression problem using MPI in Python. Each agent used a Distributed Sub-gradient method to update its own estimate of optimal 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.