

# Computational Mathematics for Learning and Data Analysis

Gemma Ragadini

AY: 2023-2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Model 1: Neural Network</b>	<b>2</b>
2.1	Implementation Overview . . . . .	3
2.2	Algorithm 1: Momentum Descent (Heavy Ball) Approach . . . .	4
2.2.1	What to expect from the algorithm . . . . .	5
2.2.2	Implementation Overview . . . . .	6
2.2.3	Considerations and Results . . . . .	7
2.3	Algorithm 2: A Deflected Subgradient Method . . . . .	8
2.3.1	What to expect from the algorithm . . . . .	9
2.3.2	Implementation Overview . . . . .	9
2.3.3	Considerations and Results . . . . .	11
2.4	Experiments . . . . .	12
<b>3</b>	<b>Model 2: L2 Linear Regression</b>	<b>16</b>
3.1	Algorithm 3: QR solver . . . . .	16
3.1.1	Implementation Overview . . . . .	17
3.1.2	Considerations and Results . . . . .	18
<b>4</b>	<b>Bibliography</b>	<b>19</b>

# 1 Introduction

As outlined in the project guidelines (ML Project 6) the goal of this project was to develop:

- **M1**: a neural network with piecewise-linear activation function and any regularization allowed;
- **M2**: a standard  $L_2$  linear regression (least squares);
- **A1**: a standard momentum descent (heavy ball) approach applied to (M1);
- **A2**: an algorithm of the class of deflected subgradient methods, applied to (M1);
- **A3**: a basic version of the direct linear least squares solver applied to (M2).

## 2 Model 1: Neural Network

A Neural Network (NN), also known as a Multi-Layer Perceptron (MLP), is a computational model designed to produce a desired output from a given input (such as for tasks like classification or prediction) after undergoing a learning phase. During this phase, its internal parameters are adjusted using a learning algorithm. In more detail, an NN consists of multiple fully-connected layers. Each layer contains a set of weights  $\mathbf{w}$  determined by the number of neurons (or computing units) and includes a bias  $b$ . A layer receives an input  $\mathbf{x}$ , either a sample if it's the first layer or the output from the preceding layer, and produces an output by combining the input with the weights (forming what is known as the *net*) and applying an *activation function* to introduce non-linearity.

$$\mathbf{net} = \mathbf{w}^T \mathbf{x} + b \quad (1)$$

$$\mathbf{o} = f_\sigma(\mathbf{net}) \quad (2)$$

As a result, it becomes feasible to formulate the complete equation of a neural network based on its architecture. For instance, consider the equation representing a neural network with one hidden layer and one output layer. Here,  $f_o$  and  $f_h$  represent the activation functions associated with the output and hidden layers respectively and the name of the function  $h$  derives from the fact that it's called the *hypothesis function*.

$$h(\mathbf{x}) = f_o(\mathbf{w}_o^T (f_h(\mathbf{w}_h^T \mathbf{x} + b_h)) + b_o) \quad (3)$$

As requested by the project assignment, the *ReLU* function (piecewise linear and non-differentiable at 0) has been used as the activation function for the hidden layer.

$$ReLU(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{otherwise} \end{cases}$$

Without the addition of an activation function, a neural network wouldn't be able to solve more than linear problems. However, with this addition, the validity of the Universal Approximation Theorem is achieved ([4], 11.6). Our model will learn by leveraging a dataset: a collection of examples with corresponding labels that describe the desired outcomes.

$$\{(\mathbf{x}_i, \mathbf{d}_i)\}_{i=1}^n \quad (4)$$

where  $n$  is the amount of samples, and the label  $\mathbf{d}_i$  can assume values according to the specific task. The Loss function chosen for this project is the *Mean Squared Error* (MSE):

$$E = \frac{1}{n} \sum_i^n (\mathbf{d}_i - h(\mathbf{x}_i))^2 \quad (5)$$

Where  $d_i$  is the label of the sample  $i$  and  $h(\mathbf{x}_i)$  the correspondent value produced by the model. The objective is to approach the minimum of the loss function by adjusting the parameters of the network, weights, and biases.

*L2 regularization* was involved, which consists in adding to the loss function the term  $\frac{\lambda}{2} \|\mathbf{w}\|$ , where  $\lambda$  is the regularization parameter, in order to simplify the model and enable better generalization, penalizing weights with a larger norm. So the error of the model will be:

$$E = \frac{1}{n} \sum_i^n (\mathbf{d}_i - h(\mathbf{x}_i))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2 \quad (6)$$

## 2.1 Implementation Overview

I chose to use the firsts 100 rows of **ML-CUP-23** dataset from Machine Learning course (regression task), which is why I opted to implement a neural network with 1 hidden layer and to fix the problem to optimize as follows: *topology* = [10, 100, 3] and  $\lambda = 10^{-5}$ . I wrote the code in Python and created classes **m1** and **m2** for the models. Moreover a **gradient** class was created, composed of matrices with dimensions suitable for containing the weights of the network, as well as the gradient, and finally the deflection. I implemented the ability to learn through minibatches, to plot the objective function against the number of epochs and to create and plot a tomography. There is a text file "**config\_NN.txt**", its purpose is to set the parameters of the model.

## 2.2 Algorithm 1: Momentum Descent (Heavy Ball) Approach

The main idea of gradient descent approach is to use  $\Delta \mathbf{w} = -\nabla L(\mathbf{w})$  into the descent rule

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta \Delta \mathbf{w}^i$$

The rule used to update the weights of the neural network using a momentum descent approach is:

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla L^i + \beta \nabla L^{i-1} \quad (7)$$

Where  $\mathbf{w}$  is the nn's weights matrix,  $\eta$  is the stepsize,  $\beta$  is the momentum parameter and  $\nabla L^i$  is the gradient of the Loss function: the matrix of the derivatives  $\frac{\partial L^i}{\partial w_{jk}^i}$ , the derivatives of the Loss function with respect to each component of the weights matrix. The addition of momentum in the learning rule allows considering the descent direction of the previous iteration and leads to better convergence (less zig-zag in the descent), depending on the value of  $\beta$ . At each iteration, the new gradient of the Loss function with respect to the  $(j, k)$  weight (the weight from unit  $j$  to unit  $k$ ) is calculated through error *backpropagation*, omitting  $i$  for the  $i^{th}$  iteration:

$$\frac{\partial L}{\partial w_{jk}} = \delta_k o_j \quad (8)$$

where  $\delta$  is different if the event  $j$  is a unit of output (case 1) or it's hidden layer (case 2). Let's call  $f_j$  the activation function correspondent to the layer of unit  $j$  and  $net_j = \sum_s w_{sj} o_s$ , the weighted sum computed by  $j$  before the application of  $f_j$  (index  $t$  iterates on the units of the  $k$ 's layer).

$$\begin{aligned} (1) \quad \delta_j &= (d_j - o_j) f_j'(net_j) \\ (2) \quad \delta_j &= \left( \sum_t \delta_t w_{jt} \right) f_j^i(net_j) \end{aligned} \quad (9)$$

Fundamentally, the gradient is computed layer by layer, starting from the output layer and will represent for each weight, its contribution to the error. Summarizing, the proposed approach to the neural network is:

---

### Algorithm 1: Momentum Descent Approach

---

**Input:** initial  $w$ , initial  $\nabla L$ ,  $\eta$ ,  $\beta$   
**for** every example input  $x$  **do**  
    compute  $o$  and  $L$   
    using backprop compute  $\nabla L$   
    using (7) update  $w$

---

### 2.2.1 What to expect from the algorithm

The loss function that we aim to minimize is not differentiable, as the activation function used is not differentiable (ReLU is not differentiable at  $x = 0$ ), so I choose to compute a subgradient. There are no guarantees that the chosen subgradient will provide a descent direction, and for this reason, we cannot expect great performance. However, using a small stepsize seems to work quite well, as can be seen from the tests. Additionally, the objective function is also not-convex, as a NN, but to provide evidence of its non-convexity, I plotted a tomography along a random direction, an array generated with `np.random.uniform(-1, 1)` and normalized, where it can be seen that the function is not convex along that direction and, consequently, is not convex overall (shown in Figure 1). The tomography has been created with  $\varphi_{w,d}(t) = L(x + \alpha d)$ . As we can see

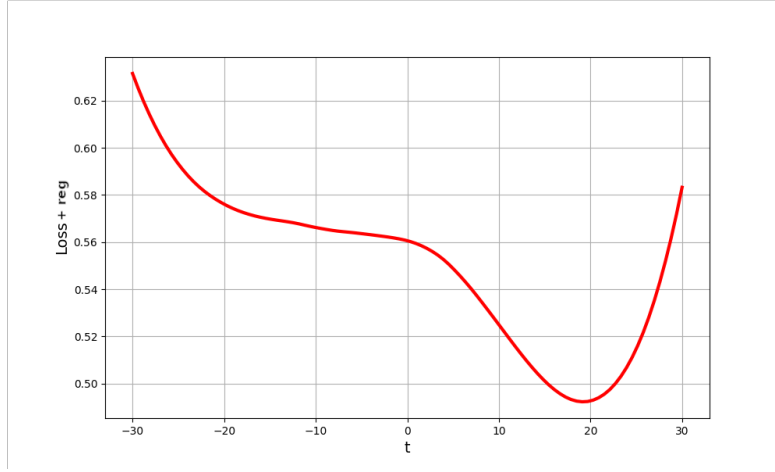


Figure 1: Tomography along a direction

from the tomography, the function is quasiconvex. To obtain this tomography, where a portion of non-convexity is evident, several attempts were necessary. This suggests that, while our objective function is not smooth, it might not be so different from a smooth function. We know that the Heavy-Ball (HB) algorithm for  $\tau$ -convex and  $L$ -smooth objective function has *linear convergence* with optimal rate  $r = \frac{\sqrt{k}-1}{\sqrt{k}+1}$ , where  $k = \frac{L}{\tau}$  ([6], Th. 3.15), choosing good parameters. Although our objective function is not differentiable, we expect it to behave similarly to a convex function since the ReLU introduces a point of non-differentiability, but we expect that regularization will mitigate this aspect. Due to L2 regularization, by the definition of  $\tau$ -convexity, we expect our regularization parameter  $\lambda$  to influence  $r$  in a manner similar to how  $\tau$  would. Due to these considerations, one would expect the convergence of the HB algorithm on the problem to be similar to linear convergence, as the problem is not very different from what one would encounter in the case of smoothness and

convexity.

### **2.2.2 Implementation Overview**

The following are the main implementation choices made for this algorithm. The algorithm is implemented into a class `a1` inheriting from class `m1`. `a1` consists of a main method `learning` wich calls the others: `backpropagation` and `update_weights`. `main1.py` calls the learning method of class `a1` if specified by command line. The parameters used are specified in a text file named "`config_a1.txt`" along with some model parameters. Those related to algorithm `a1` are: `stepsize` and `momentum`.

### 2.2.3 Considerations and Results

Using  $stepsize = 10^{-4}$  and  $momentum = 10^{-5}$  on dataset heavy ball algorithm converges on model 1, with a trend shown in the graph (Figure 2), where the value of  $\log(E_r)$  varies on 20000 epochs and  $E_r$  is the error relative to the minimum value found for the objective function (Loss added regularization) computed between algorithm a1 and a2 (0.00036).

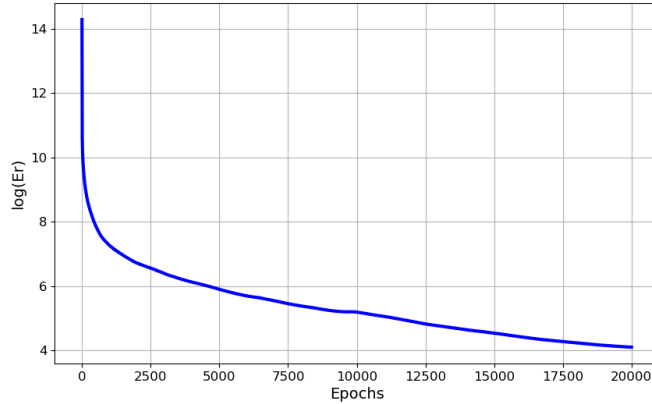


Figure 2: objective function varying with the number of epochs

Let's define *convergence rate*  $r$  of a sequence  $f^i$  to the optimal value  $f^*$ :

$$r = \lim_{i \rightarrow \infty} \frac{f^{i+1} - f^*}{(f^i - f^*)^p} \quad (10)$$

As can be seen from the Figure 2, the algorithm exhibits *sublinear convergence*, which, however, appears at the tail, to be quite similar to *linear convergence*. This seems consistent with the assumptions regarding the properties of the objective function. In Figure, the optimization curves using various algorithm parameter configurations are compared.

- *Config 1*:  $stepsize = 0.0001$ ,  $momentum = 0.00001$
- *Config 2*:  $stepsize = 0.001$ ,  $momentum = 0.0001$
- *Config 3*:  $stepsize = 0.00015$ ,  $momentum = 0.000015$

As we will see in 2.4, we are assuming that all algorithms converge to the same minimum. From the Figure 3, we can see that in each case, convergence is sublinear; there aren't significant performance differences, but the issue might be the speed of convergence.

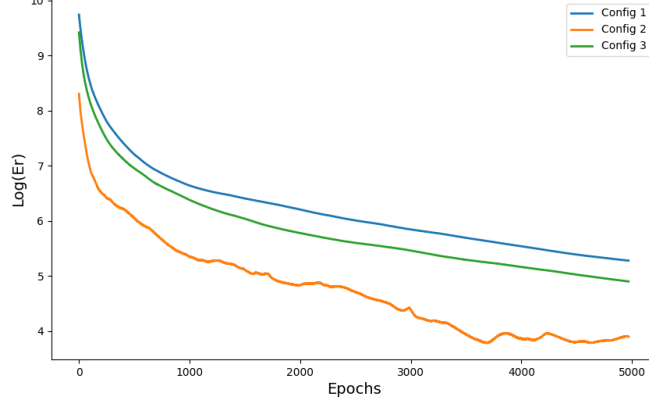


Figure 3: Comparison between different parameters configurations

### 2.3 Algorithm 2: A Deflected Subgradient Method

In this section, we address a simple algorithm from the class of subgradient methods, whose implementation is the result of studying and applying the methods discussed in [1], [2]. An algorithm from the class of deflected subgradient methods is used due to the non-differentiability of the function to minimize, as previously discussed. The arbitrary choice made regarding its value at the non-differentiable point is that it equals 1. Therefore, we will have:

$$\frac{\partial ReLu(x)}{\partial x} = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{otherwise} \end{cases} \quad (11)$$

In the *backpropagation* process as in *Algorithm 1*, the subgradient is computed. However, an additional deflection is introduced, calculated at iteration  $i$  as follows:

$$\mathbf{d}_i = \|\alpha_i \mathbf{g}_i + (1 - \alpha_i) \mathbf{g}_{i-1}\|_2^2 \quad (12)$$

Where  $\mathbf{g}_i$  is the subgradient, and  $\alpha_i$  is

$$\alpha_i = \operatorname{argmin}_{\alpha_j \in [0,1]} \|\alpha_i \mathbf{g}_i + (1 - \alpha_i) \mathbf{g}_{i-1}\|_2^2 \quad (13)$$

When, omitting  $i$  and defining  $\mathbf{v} = \mathbf{g} - \mathbf{d}$ , it is computed as follows:

$$\|\alpha \mathbf{v} + \mathbf{d}\|_2^2 = \alpha^2 \|\mathbf{v}\|_2^2 + 2\alpha \langle \mathbf{v}, \mathbf{d} \rangle + \|\mathbf{d}\|_2^2$$

So, we can find the minimum as:

$$\alpha = -\frac{\langle \mathbf{v}, \mathbf{d} \rangle}{\|\mathbf{v}\|_2^2} \text{ if } \mathbf{g} \neq \mathbf{d}$$

Otherwise, we can chose a value for  $\alpha$  in  $[0, 1]$ ; for the implementation I choose 0.5. If the minimum is not into the wished range, we can check values at



the extremes. Moreover, *stepsize* has to be different for each iteration due to the non differentiability of the function, so I choose a stepsize rule similar to *Diminishing-Square Summable stepsize (DSS)*  $\eta$ , as follows:

$$\eta = \frac{\beta}{i} \quad (14)$$

By defining  $j$  as a counter starting from  $j = 1$  for the number of consecutive non-improving iterations, and fixing  $k$  chosen a priori,  $j$  will be reset if its value exceeds  $k$ , and  $i$  is the counter to be incremented in that case. At the end the rule for updating weights at iteration  $i$  will be:

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \eta_i \mathbf{d}_i \quad (15)$$

---

**Algorithm 2:** Deflected Subgradient Approach

---

**Input:** initial  $w$ , loss function  $L$ , initial  $d$ ,  $\beta$   
**for** every example input  $x$  **do**  
    compute  $o = wx$  and  $L$   
    using backprop compute  $g = \nabla L$   
    using (13) compute  $\alpha$   
    using (12) compute  $d$   
    using (14) compute  $\eta$   
    using (15) update  $w$

---

### 2.3.1 What to expect from the algorithm

Due to the non-differentiability and non-convexity of my objective function, in theory I would expect a subgradient method to perform better than the heavy ball approach. This is certainly true in general for non-differentiable problems, but we are in a limit case due to the L2 regularization and the almost-convexity of the objective function. Anyway, with the right parameters, the algorithm is expected to converge to a minimum, which may be local or global. As in [1] a subgradient method requires  $\Theta(\frac{1}{\epsilon^2})$  iterations to solve up to absolute error  $\epsilon$ , which means that it is not practical for attaining any more than a modest accuracy. Besides, the complexity is independent of the size  $n$  of the problem. Therefore, SM may be promising for very-large-scale problems where a high accuracy is not necessary, whereas a short running time is a primary concern.

### 2.3.2 Implementation Overview

The following are the main implementation choices made for this algorithm. The algorithm is implemented into a class `a2` inheriting from class `m1`. `a1` consists of a main method `learning` which calls the others: `compute_alpha`,

`compute_deflection`, `compute_stepsize`, `update_weights` and `backpropagation`.  
`main1.py` calls the learning method of class `a2` if specified by command line.  
The parameters used by the algorithm ( $k$  and  $\beta$ ) are specified in a text file  
named "`config_a2.txt`" along with some model parameters.

### 2.3.3 Considerations and Results

We know that in cases of non-differentiability, fixed stepsizes cannot be used. For this reason, I have chosen a rule similar to the *diminishing/square summable stepsize*, which maintains simplicity but enhances slightly faster convergence. The use of a convex combination in (12) is crucial, because it ensures that  $d_i$  is always an approximate subgradient of the loss function. The adding of deflection allows to face the "zig-zagging" behaviour whereby  $g_{i+1} \approx -g_i$ , so that two "reasonably long" steps combined make an "unfeasibly short" one. In Figure 4 we can see the value of objective function varying with the number of epochs, using the the 7-th configuration in Table 5. In this case as well, the

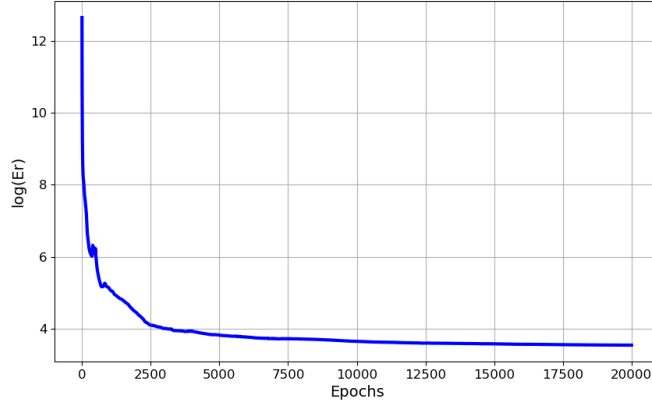


Figure 4: objective function varying with the number of epochs

convergence is sublinear, more markedly than in Figure 2. This is consistent with the fact that since the objective function is almost convex, it is expected that a1 performs better than a2. This is because a2 is designed for scenarios where the lack of differentiability is significant enough to cause problems for algorithms like a1 that assume smoothness.

- *Config1*:  $\beta = 0.003$ ,  $k = 10000$
- *Config2*:  $\beta = 0.005$ ,  $k = 15000$
- *Config3*:  $\beta = 0.001$ ,  $k = 12000$

As we will see in 2.4, we are assuming that all algorithms converge to the same minimum. From the Figure 5, we can see that in each case, convergence is sublinear; there aren't significant performance differences, but the issue might be the speed of convergence.

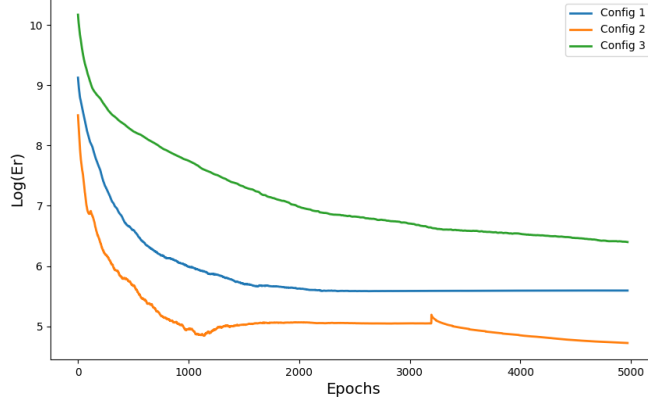


Figure 5: omparison between different parameters configurations

## 2.4 Experiments

Naming  $f$  the objective function MSE with regularization, in this section, we attempt to present and compare the results of the two algorithms discussed so far. The notation  $f^*$  specifies the minimum value reached by a specific configuration of the problem. In Table 1 I compared the performance of  $a1$  while keeping momentum fixed, as stepsize varied, excluding extreme cases that showed poor performance. There do not seem to be significant differences in performance. Within the range of values where everything "works well". The values in the middle of this tested range appear to be the best.

stepsize	momentum	$f^*$
0.001	0.0001	0.02997
0.0013	0.0001	0.02191
0.0015	0.0001	0.02409
0.0017	0.0001	0.02123
0.002	0.0001	0.02306

Table 1: Comparison between stepsize values on 5000 epochs

I did the same for the momentum in  $a1$ , keeping the step size fixed as shown in Table 2. Again, the oscillation in convergence speed is not very high. However, even in this case, the best values appear to be in the middle of the tested range.

In Table 3 I compared the performance of  $a2$  while keeping  $\beta$  fixed, as  $k$  varied, excluding extreme cases that showed poor performance.

stepsize	momentum	$f^*$
0.001	0.00001	0.03317
0.001	0.0001	0.02997
0.001	0.00015	0.01389
0.001	0.0002	0.03039
0.001	0.0003	0.03269

Table 2: Comparison between momentum values on 5000 epochs

k	$\beta$	$f^*$
$10^2$	0.003	6.45062
$10^3$	0.003	0.12248
$10^4$	0.003	0.03519
$10^5$	0.003	0.04381
$10^6$	0.003	0.04785
$10^7$	0.003	0.05791
$10^8$	0.003	0.03255
$10^9$	0.003	0.05654

Table 3: Comparison between k values on 5000 epochs

Since the performance seemed to increase with the orders of magnitude of  $k$ , but then remained approximately constant beyond  $10^4$ , I decided to test  $k$  around the order of  $10^4$ . Fixing  $k$ , I did the same for  $\beta$  in  $a2$ , varying it in Table 4.

k	$\beta$	$f^*$
$10^4$	0.005	0.02096
$10^4$	0.004	0.03684
$10^4$	0.003	0.02140
$10^4$	0.002	0.08803
$10^4$	0.001	0.29254

Table 4: Comparison between  $\beta$  values on 5000 epochs

In Table 4 it can be observed that there doesn't seem to be a significant difference in performance apart from  $\beta = 0.001$ . However, it is also noted that these values of the objective function, which differ slightly, are interchangeable when the starting point is changed randomly. For this reason, in Table 5, some of these are tested. As mentioned the objective function is not convex. For this reason, using the same notation of precious sections, in Table 5 I try to understand if the methods used were approaching the same optimal value or different ones, starting from the same initial point, created with `np.random.uniform`.

At first glance, we are not sure that the values found represent the same minimum as we would have expected. However, to evaluate this more

Optimizer	Stepsize	Momentum	$\beta$	k	$f^*$
A1	0.0001	0.00001	-	-	0.03457
A1	0.001	0.0001	-	-	0.00672
A1	0.001	0.0005	-	-	0.00285
A1	0.00015	0.000015	-	-	0.01814
A2	DSS	-	0.003	10000	0.01561
A2	DSS	-	0.005	15000	0.01897
A2	DSS	-	0.004	12000	0.01180
A2	DSS	-	0.003	15000	0.01647

Table 5: Comparison of minima on 15000 iterations

accurately, this ratio was calculated:

$$\frac{\|w_i - \text{mean}_j w_j\|}{\|\max_j w_j\|}$$

And the obtained relative distances related to the data in the Table 5 are:

$$[0.45, 0.24, 0.16, 0.27, 0.29, 0.59, 0.30, 0.29]$$

We were interested in understanding whether the minimum reached by the two algorithms and some of their configurations was the same. We could reasonably assume so, given the shape of the curve we envision (quite convex). However, the convergence to the optimal value of the two algorithms over 15000 epochs does not seem to confirm this with certainty, especially for **a2**. To have reasonable confirmation of convergence to a single minimum, we could evaluate the minimum value reached in 15000 epochs, using the second parameters configurations from Table 5 and for each execution a different random initial point (`np.random.uniform` in  $[-2, 2]$  ).

ex 1	0.00066
ex 2	0.00039
ex 3	0.00094
ex 4	0.00081
ex 5	0.00087
ex 6	0.00049
ex 7	0.00071
ex 8	0.00066

Table 6: Comparison of minima with random initial point using second configuration parameter in A1

To assess whether the obtained values corresponded to the same minimum, I calculated the distances between the weight arrays, as done for the data in Table 5:

[0.0243, 0.033, 0.087, 0.069, 0.013, 0.097, 0.068, 0.052]

I have also calculated in the same way the gaps between the correspondent values obtained from the objective function:

[0.04, 0.43, 0.35, 0.17, 0.25, 0.29, 0.027, 0.045]

We can conclude in this case that, despite the different initial point, the algorithm seems to converge to the same minimum value. Regarding the performance of the two algorithms on the problem, they both exhibit sublinear convergence, the convergence of **a2** is faster, but it appears to stop at a minimum value that is much higher than the convergence of **a1**. Yet, it appears that **a1** requires more iterations to converge, but it seems to continue decreasing in the tail, whereas **a2** continues to decrease, but very slowly.

### 3 Model 2: L2 Linear Regression

The second model aimed to solve a linear regression problem using *least squares*. So, my objective was to find  $\bar{\mathbf{x}} \in \mathbb{R}^{n \times k}$  s.t.

$$\bar{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^{n \times k}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$$

Where  $A \in \mathbb{R}^{m \times n}$  and  $\mathbf{y} \in \mathbb{R}^{m \times k}$ . We can think of it as finding the vector closest to  $\mathbf{y}$  in the hyperplane  $\operatorname{Im}(A)$ .

#### 3.1 Algorithm 3: QR solver

Among the direct linear least squares solvers, I chose the *QR factorization method* to solve least square problems.

**Theorem.** For every  $A \in \mathbb{R}^{m \times n}$ , there exists  $Q \in \mathbb{R}^{m \times m}$  orthogonal,  $R$  upper triangular s.t.  $A = QR$ .

*Proof in Theorem 7.1 of [3].*

To compute the QR decomposition, I used the *Householder Reflectors* method.

**Lemma.**  $\forall v$  in  $\mathbb{R}^m$  the matrix  $H = I - \frac{2}{\|v\|^2}vv^T$  is orthogonal and symmetric.

*Proof.* (symmetry)

$$H^T = I^T - \frac{2}{\|v\|^2}(vv^T)^T = I - \frac{2}{\|v\|^2}vv^T = H$$

(orthogonality)

$$\begin{aligned} H^T H &= H^2 = (I - \frac{2}{\|v\|^2}vv^T)(I - \frac{2}{\|v\|^2}vv^T) = I - \frac{2}{\|v\|^2}vv^T - \frac{2}{\|v\|^2}vv^T + \\ &+ \frac{4}{\|v\|^4}vv^T vv^T = I - \frac{4}{\|v\|^2}vv^T + \frac{4}{\|v\|^2}vv^T = I \end{aligned}$$

□

**Lemma.** Let  $x, y \in \mathbb{R}^m$  s.t.  $\|x\| = \|y\|$ . If  $v = x - y \Rightarrow H = I - \frac{2}{\|v\|^2}vv^T$  is s.t.  $Hx = y$ .

*Proof.*  $(I - \frac{2}{\|v\|^2}vv^T)x = y \Leftrightarrow (I - 2\frac{vv^T}{\|x\|^2 - 2x^T y})x = y \Leftrightarrow x - \frac{vv^T x}{\|x\|^2 - x^T y} = y$   
 $x - \frac{v(\|x\|^2 - y^T x)}{\|x\|^2 - x^T y} = y \Leftrightarrow x - (x - y) = y$  □

The vector  $y = Hx$  can be computed from  $v$  and  $x$  with  $O(m)$  operations. In our case, with  $A$  being the matrix of inputs and  $\mathbf{y}$  being that of outputs, it is assumed that  $m \gg n$  (and indeed, this is the case in the example I will use in the implementation), so we will have a thin matrix. Using the QR factorization, the problem will be:



$$\|A\mathbf{x} - \mathbf{y}\|_2 = \|QR\mathbf{x} - \mathbf{y}\|_2 = \|Q^T(QR\mathbf{x} - \mathbf{y})\|_2 =$$

Where the last equality stems from the fact that  $Q$  is orthogonal than it preserves the norm.

$$= \|R\mathbf{x} - Q^T\mathbf{y}\|_2 = \left\| \begin{bmatrix} R_0\mathbf{x} \\ 0 \end{bmatrix} - \begin{bmatrix} Q_0^T\mathbf{y} \\ Q_c^T\mathbf{y} \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} R_0\mathbf{x} - Q_0^T\mathbf{y} \\ -Q_c^T\mathbf{y} \end{bmatrix} \right\|_2$$

It is possible to notice that  $-Q_c^T\mathbf{y}$  does not depend on  $\mathbf{x}$ . Therefore, the objective can be reduced to minimizing the upper part.  $R_0\mathbf{x} - Q_0^T\mathbf{y}$  can be zero if and only if  $R_0$  is invertible.

**Theorem.**  $\forall$  matrix  $A$ , the matrix  $R_0$  given by its  $QR$  factorization is invertible iff  $A$  has full column rank.

*Proof.*  $A$  has full column rank  $\iff A^T A$  is invertible, but

$$A^T A = (QR)^T QR = R^T Q^T QR = R^T R = \begin{bmatrix} R_0^T & 0 \end{bmatrix} \begin{bmatrix} R_0 \\ 0 \end{bmatrix} = R_0^T R_0$$

Finally, the product of two square matrices is invertible iff they are both invertible.  $\square$

Then it is proved:

**Lemma.** Let  $A \in \mathbb{R}^{m \times n}$  be a tall thin matrix with full column rank. Then the solution of the LS problem  $\min \|A\mathbf{x} - \mathbf{y}\|_2$  is obtained with  $\mathbf{x} = R_0^{-1}(Q_0^T\mathbf{y})$  and the value of the minimum is  $\|Q_c^T\mathbf{y}\|_2$ .

### 3.1.1 Implementation Overview

To address the problem presented in the dataset `ML-CUP23-TR.csv`, which contains approximately 1000 input data points, each with 10 features, I have chosen to use the *thin QR* factorization (also known as reduced QR factorization). This method is particularly suitable for rectangular matrices where the number of rows (1000) is significantly greater than the number of columns (10), allowing for a more efficient decomposition  $A = QR$  compared to the full QR factorization. The thin QR factorization yields a matrix  $Q$  of dimensions  $mn$  (in this case  $1000 \times 10$ ) with orthonormal columns, and a matrix  $R$  of dimensions  $nn$  ( $10 \times 10$ ) that is upper triangular. Additionally, the *Huge optimization* was applied by leveraging the structure of the Householder matrix  $H$ . This approach avoids explicitly computing the product between matrices, effectively reducing the computational cost from cubic to quadratic (computing  $R \leftarrow R - 2uu^T R$  instead of  $R \leftarrow H$  and the same respectively for  $Q$ ). It is computationally much more convenient in thin cases. The algorithm computationally costs ( $O(mn^2)$  from lecture 10 in [3]). The matrix corresponding to the dataset used has full column rank, So I was able to calculate  $\mathbf{x}$  as described previously. I created a class `linear regression` that implements *model 2* for linear regression with LS, and class `a3` that implements an algorithm for LS problem via QR factorization.

---

**Algorithm 3:** QR thin solver for LS

---

**Input:**  $A, x$   
 $(m, n) \leftarrow \text{shape}(A)$   
 $R \leftarrow A$   
 $Q_0 \leftarrow \text{eye}(m, n)$   
**for**  $i=0:n$  **do**  
    update  $R$   
    update  $Q$   
compute  $x = Q_0^T y$   
solve  $R_0 x = Q_0^T y$   
**Output:**  $x$

---

### 3.1.2 Considerations and Results

Since QR decomposition is *backward stable* (the computed  $Q, R$  are the exact result of  $qr(A + \Delta A)$  where  $\|\Delta A\| \leq O(u) \|A\|$ ) and householder transformations are too, the total algorithm to solve LS problem with QR factorization is backward stable and we can say it "a priori"; so Algorithm 3 only have intrinsic representation errors. Since the problem is backward stable, I know that  $\frac{\|\tilde{y} - y\|}{\|y\|} \leq k_{rel}(f, x) \cdot O(u)$ . I calculated  $k_{rel}(f, x) = 7.912$ , so we can affirm that the error in the solution is  $\leq 7.912 \cdot 10^{-16}$ . Moreover how much my QR factorization differs from the original matrix  $A$ , in terms of relative error,  $\frac{\|A - QR\|}{\|A\|} = 1.09 \cdot 10^{-15}$ , and this is very close to floating order. Moreover I computed relative residuals  $\frac{\|Ax - y\|}{\|y\|} = 0.24325720795631878$  and compared with residuals from `np.linalg.lstsq`, these two values are the same up to the 18th decimal place. In order to verify if this implementation of thin QR decomposition using Householder reflections showed a linear growth in execution time as a function of the number of rows  $m$  of the matrix, as theoretically expected, Figure 6 depicts the graph of this growth. It is observed that, apart from some fluctuations, the execution time growth is regular and follows an increasing trend consistent with the expected linear growth.

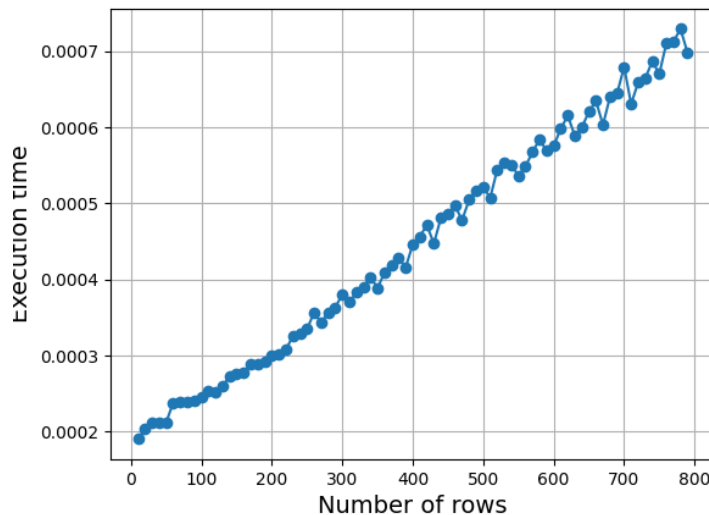


Figure 6: Thin QR Factorization Execution Time vs. Number of Rows

## 4 Bibliography

### References

- [1] Antonio Frangioni, Bernard Gendron, Enrico Gorgone. *On the computational efficiency of subgradient methods: a case study with Lagrangian bounds*, Math. Prog. Comp. (2017) 9:573–604 DOI 10.1007/s12532-017-0120-7.
- [2] Giacomo d’Antonio, Antonio Frangioni *Convergence Analysis of Deflected Conditional Approximate Subgradient Methods*, SIAM J. OPTIM. Vol. 20, No. 1, pp. 357–386
- [3] Lloyd N. Trefethen, David Bau *III Numerical linear algebra 1997* siam
- [4] Ethem Alpaydın *Introduction to Machine Learning, Second Edition* The MIT Press Cambridge, Massachusetts London, England
- [5] Euhanna Ghadimi, Hamid Reza Feyzmahdavian, and Mikael Johansson *Global convergence of the Heavy-ball method for convex optimization* 2015 European Control Conference (ECC) July 15-17, 2015. Linz, Austria
- [6] Sébastien Bubeck *Convex Optimization: Algorithms and Complexity* Foundations and Trends in Machine Learning Vol. 8, No. 3-4 (2015) 231–358 c 2015 S. Bubeck