

# Analysis of optimization and numerical approaches to solve the linear least square problem

Emanuele Cosenza\*, Riccardo Massidda†

Department of Computer Science  
University of Pisa

\* e.cosenza3@studenti.unipi.it, † r.massidda@studenti.unipi.it

**Abstract**—The linear least square problem can be tackled using a wide range of optimization or numerical methods. The L-BFGS method of the class of limited-memory quasi-Newton algorithms has been chosen for the former, whilst the thin QR factorization with Householder reflectors for the latter. Both these algorithms have been implemented from scratch using Python language, to finally experiment over their performances in terms of precision, stability and speed. The accordance of the implementations with the underlying theoretical models is also studied and discussed.

## INTRODUCTION

Given a dataset composed by a matrix  $\hat{X} \in \mathbb{R}^{m \times n}$  with  $m \geq n$  and a vector  $y \in \mathbb{R}^m$ , the solution of the linear least square (LLS) problem is the vector  $w \in \mathbb{R}^n$  that fits best the data assuming a linear function between  $\hat{X}$  and  $y$ . (Nocedal and Wright 2006, 50) This can be formalized as the following minimization problem:

$$w_* = \min_w \|\hat{X}w - y\|_2^2$$

The matrix  $\hat{X}$  is actually composed in the following way:

$$\hat{X} = \begin{bmatrix} X^T \\ I \end{bmatrix}$$

Where  $X \in \mathbb{R}^{n \times k}$  is a tall thin matrix, thus  $m = k + n$ . The LLS problem can be dealt both with iterative methods or with direct numerical methods. One algorithm has been chosen for each of these fields to finally discuss their experimental results.

## L-BFGS

The Limited-memory BFGS, L-BFGS, is an iterative method of the quasi-Newton limited-memory class. This method is a variation of the BFGS method, with which it shares the update rule. At the  $i + 1$ -th iteration the point is updated as follows:

$$w_{i+1} = w_i - \alpha_i H_i \nabla f_i$$

The smaller memory requirements of this variation are due to the fact that the Hessian approximation  $H_i$  is stored implicitly, and built over a fixed number of vector pairs  $\{s_j, y_j\}$  of the previous  $t$  iterations and an initial matrix  $H_i^0$ . Where

$$s_i = w_{i+1} - w_i, \quad y_i = \nabla f_{i+1} - \nabla f_i$$

$$V_i = I - \rho_i y_i s_i^T, \quad \rho_i = \frac{1}{y_i^T s_i}$$

so that  $H_i$  satisfies the following condition

$$H_i = (V_{i-1}^T \dots V_{i-t}^T) H_i^0 (V_{i-t} \dots V_{i-1})$$

$$+ \rho_{i-t} (V_{i-1}^T \dots V_{i-t}^T + 1) s_{i-t} s_{i-t}^T (V_{i-t+1} \dots V_{i-1})$$

$$+ \rho_{i-t+1} (V_{i-1}^T \dots V_{i-t}^T + 2) s_{i-t+1} s_{i-t+1}^T (V_{i-t+2} \dots V_{i-1})$$

$$+ \dots$$

$$+ \rho_{i-1} s_{i-1} s_{i-1}^T$$

Different strategies to initialize the  $H_i^0$  matrix are proposed in the literature, and so they will be tested experimentally. Finally, the step size  $\alpha_i$  is found by performing an inexact line search based on the Armijo-Wolfe conditions.

## Thin QR factorization

For the numerical counterpart, the thin QR factorization with Householder reflectors has been implemented as described in (Trefethen and Bau 1997).

By using the Householder QR factorization, the matrix  $R \in \mathbb{R}^{m \times n}$  is constructed in place of  $\hat{X}$  and the  $n$  reflection vectors  $v_1, \dots, v_n$  are stored. The reduced matrix  $\hat{R} \in \mathbb{R}^{n \times n}$  is trivially obtainable by slicing as in  $\hat{R} = R_{1:n, 1:n}$ , noting that constructing directly the reduced matrix would yield no significant advantage since we are already using  $O(mn)$  space to store the reflection vectors.

By using the Householder vectors it is also possible to implicitly compute  $\hat{Q}^T y$  to finally obtain  $w_*$  by back substitution over the upper-triangular system  $\hat{R}w = \hat{Q}^T y$ .

# ALGORITHMIC ANALYSIS

## Convergence of L-BFGS

Liu and Nocedal (1989) define three necessary assumptions to prove that the L-BFGS algorithm globally converges and that there exists a constant  $0 \leq r < 1$  such that

$$f(w_i) - f(w_*) \leq r^i(f(w_0) - f(w_*))$$

so that the sequence  $\{w_i\}$  converges R-linearly.

Firstly the objective function  $f$  should be twice continuously differentiable. Given the formulation of the least squares problem this is immediately true, the gradient and the Hessian of the objective function are definable as in:

$$\nabla f(w) = \hat{X}^T(\hat{X}w - y)$$

$$\nabla^2 f(w) = \hat{X}^T \hat{X}$$

Moreover the Hessian can be proven to be positive definite, as can be easily seen by rearranging it in the following way:

$$\begin{aligned} \nabla^2 f(w) &= \hat{X}^T \hat{X} \\ &= \begin{bmatrix} XI \\ I \end{bmatrix} \begin{bmatrix} X^T \\ I \end{bmatrix} \\ &= XX^T + I \end{aligned}$$

The matrix  $XX^T$  is positive semi-definite, since  $\forall z : z^T XX^T z = \|X^T z\|^2 \geq 0$ , therefore all the eigenvalues of the matrix are non-negative. Furthermore, according to the spectral theorem, since  $XX^T$  is symmetric, there exists  $U$  orthogonal matrix and  $D$  diagonal containing the eigenvalues of  $XX^T$ .

$$\begin{aligned} \nabla^2 f(x) &= XX^T + I \\ &= UDU^T + I \\ &= UDU^T + UIU^T \\ &= U(D + I)U^T \end{aligned}$$

The eigenvalues of the Hessian are contained in  $D + I$  and all of them are positive, therefore  $\nabla^2 f(w)$  is positive definite.

Being the Hessian positive definite, the objective function  $f$  is a convex function. This comes in handy for the second assumption requiring the sublevel set  $D = \{w \in \mathbb{R}^n | f(w) \leq f(w_0)\}$  to be convex. It can be easily proved that if a function is convex all of its sublevel sets are convex sets.

$$\forall x, y \in D, \lambda \in [0, 1]$$

$f$  convex

$$\begin{aligned} &\implies f(\lambda x + (1 - \lambda)y) \\ &\leq \lambda f(x) + (1 - \lambda)f(y) \\ &\leq \lambda f(w_0) + (1 - \lambda)f(w_0) \\ &= f(w_0) \\ &\implies \lambda x + (1 - \lambda)y \in D \end{aligned}$$

The third and last assumption requires the existence of two positive constants  $M_1$  and  $M_2$  such that  $\forall z \in \mathbb{R}^n, w \in D$ :

$$M_1 \|z\|^2 \leq z^T \nabla^2 f(w) z \leq M_2 \|z\|^2$$

or equivalently

$$M_1 I \preceq \nabla^2 f(w) \preceq M_2 I$$

Since  $\nabla^2 f(w)$  is positive definite the previous condition is true for  $M_1 = \lambda_{min}$  and  $M_2 = \lambda_{max}$ , where  $\lambda_{min} > 0$ .

In the convergence proof the  $M_2$  constant is used to upper bound the trace of the next Hessian substitute  $H_{i+1}$ , implying an upper bound for the largest eigenvalue in the sequence of Hessian substitutes.

$$tr(H_{i+1}) \leq tr(H_i^0) + tM_2 \leq M_3$$

On the other hand the  $M_1$  constant is used, to lower bound the determinant of  $H_{i+1}$ , implying a lower bound for the smallest eigenvalue in the sequence of Hessian substitutes.

$$det(H_{i+1}) \geq det(H_i^0) + \left(\frac{M_1}{M_3}\right)^t \geq M_4$$

These two assertions are used to prove the existence of constant  $\delta > 0$  such that

$$\forall i : \cos \theta_i = \frac{s_i^T H_i s_i}{\|s_i\| \|H_i s_i\|} \geq \delta$$

where  $\theta_i$  is the angle between the chosen direction and  $-\nabla f(w_i)$ . If the constant  $M_1$  was to be equal to zero, it would not be enough to prove the existence of  $\delta > 0$  for each step, possibly having directions orthogonal to steepest one. As already pointed out, given that the Hessian is positive definite, its eigenvalues and consequently  $M_1$  are positive.

Other than the three discussed assumptions, the theorem requires for the sequence of initializers  $\{\|H_i^0\|\}$  to be bounded. This obviously depends on the initialization technique used to generate  $H_i^0$ , various techniques are suggested in the literature such as  $H_k^0 = \gamma_k I$  or  $H_k^0 = \gamma_k H_0$  where

$$\gamma_k = \frac{s_{k-1}^T y_{k-1}}{\|y_{k-1}\|^2}$$

Other initialization techniques may possibly be tested and evaluated experimentally.

#### Armijo-Wolfe line search

The convergence proof requires the algorithm to perform a line search respectful of the Armijo-Wolfe conditions. By defining the function  $\phi$ , used to evaluate the value of  $f$  at a certain step-size  $\alpha$ , the conditions can be defined as follows.

$$\phi(\alpha) = f(w_i + \alpha d_i)$$

$$\phi(\alpha) \leq \phi(0) + \alpha \rho \phi'(0) \quad (\text{A})$$

$$\phi'(\alpha) \geq \sigma \phi'(0) \quad (\text{W})$$

Given the quadratic nature of the least squares problem, it is possible to compute the exact optimal step-size  $\bar{\alpha}$  by solving  $\phi'(\alpha) = 0$ . To simplify the following discussion the objective function  $f(w)$  is described in the form  $\frac{1}{2}w^T Q w + q^T w + c$  where

$$Q = \hat{X}\hat{X}^T, q^T = -y^T \hat{X}, c = -\frac{1}{2}\|y\|^2$$

Therefore the function  $\phi'(\alpha)$  can be defined as follows

$$\phi'(\alpha) = \frac{\partial f(w + \alpha d)}{\partial \alpha} = \nabla f(w)^T d + \alpha d^T Q d$$

where  $d$  is the descent direction computed by the L-BFGS algorithm. The optimal step-size is then computable as

$$\bar{\alpha} = -\frac{\nabla f(w)^T d}{d^T Q d}$$

It can be proven that the step size  $\bar{\alpha}$  satisfies both (A) for any  $\rho \leq \frac{1}{2}$  and (W) for any positive  $\sigma$ . The core steps of the proofs are hereby reported, it should be remarked that these depend on the fact that  $d$  is a descent direction, and so that  $f(w)^T d < 0$ .

#### Armijo satisfaction proof:

$$\phi(\bar{\alpha}) \leq \phi(0) + \rho \bar{\alpha} \phi'(0)$$

$$f(w + \bar{\alpha} d) - f(w) \leq \rho \bar{\alpha} \nabla f(w)^T d$$

$$\bar{\alpha}(w^T Q + q^T) d + \frac{\bar{\alpha}^2}{2} d^T Q d \leq \rho \bar{\alpha} \nabla f(w)^T d$$

$$\nabla f(w)^T d + \frac{\bar{\alpha}}{2} d^T Q d \leq \rho \nabla f(w)^T d$$

$$\nabla f(w)^T d - \frac{1}{2} \frac{\nabla f(w)^T d}{d^T Q d} d^T Q d \leq \rho \nabla f(w)^T d$$

$$\rho \leq \frac{1}{2}$$

#### Wolfe satisfaction proof:

$$\phi'(\alpha) \geq \sigma \phi'(0)$$

$$\nabla f(w)^T d + \bar{\alpha} d^T Q d \geq \sigma \nabla f(w)^T d$$

$$\nabla f(w)^T d - \nabla f(w)^T d \geq \sigma \nabla f(w)^T d$$

$$\sigma \geq 0$$

#### Analysis of standard and modified thin QR

Ignoring constants, we know from theory that the standard QR factorization algorithm applied on the matrix  $\hat{X}$  yields a time complexity of  $O(mn^2)$ . Actually, given that we are generally dealing with a very tall and thin matrix  $X$ , the resulting  $\hat{X}$  is going to be squarish ( $m \approx n$ ). This means that we can consider the complexity to be cubic in  $n$ .

From now on, we show a way to bring down the time complexity of the algorithm from  $O(mn^2)$  to  $O(kn^2)$ , with  $k = m - n$ . The resulting modified QR factorization algorithm will become useful when  $k$  is much smaller than  $m$ , as in our case.

In the standard algorithm, at each step we focus on a single column of the input matrix, constructing a householder vector to zero out all the entries below the diagonal. Following the geometric reasoning in (Trefethen and Bau 1997), this brings the algorithm to depend on  $m$ . While this cannot be avoided in general, in our particular case we can be a little bit smarter.

Since the block matrix  $\hat{X}$  contains the identity as its lower block, at each step of the algorithm we can just focus on zeroing out the  $k = m - n$  entries below the diagonal up to the 1s of the identity. In the modified algorithm then, to obtain  $R$ , the matrix  $\hat{X}$  is multiplied on the left side by a sequence of matrices  $L_i \in \mathbb{R}^{m \times m}$  of the form ( $i = 1, \dots, n$ ):

$$L_i = \begin{bmatrix} I_{i-1} & 0 & 0 \\ 0 & H_i & 0 \\ 0 & 0 & I_{n-i} \end{bmatrix}$$

where  $H_i \in \mathbb{R}^{(k+1) \times (k+1)}$  are all Householder reflectors that zero out the  $k$  entries in the  $i$ -th column of the matrix which is being multiplied by  $L_i$ . The resulting Householder vectors will all have constant dimension  $k + 1$  and storing them will require  $\theta(kn)$  space instead of  $O(mn)$ .

To derive the time complexity of this phase we can reason as follows. The right side matrix can be divided in three blocks as in  $\begin{bmatrix} A \\ B \\ C \end{bmatrix}$ . When we multiply this matrix by  $L_i$  the only relevant operation is the matrix multiplication  $H_i B$ , which costs  $O(kn)$ . Since the total number of multiplications is  $n$ , the total complexity is  $O(kn^2)$ .

Since each  $L_i$  is orthogonal and symmetric it is then possible to reconstruct  $Q$  and the reduced  $\hat{Q}$  in the following way:

$$Q = L_1 L_2 \dots L_n$$

$$\hat{Q} = L_1 L_2 \dots L_n \begin{bmatrix} I_n \\ 0 \end{bmatrix}$$

Applying again the reasoning above, the time complexity of these reconstructions is  $O(kn^2)$ . It follows that the overall time complexity of the modified QR factorization is  $O(kn^2)$ .

The least squares problem is then solved through back substitution over the upper-triangular system  $\hat{R}w = \hat{Q}^T y$ . Since the costs for the  $\hat{Q}^T y$  product and the back substitution are dominated by the factorization cost, the overall time complexity for solving the least squares problem through QR factorization is  $O(mn^2)$  when using the standard algorithm and  $O(kn^2)$  when using the modified one.

#### *Stability and accuracy of the QR algorithm*

As stated in (Trefethen and Bau 1997, 140), the algorithm obtained by combining the standard QR algorithm, the  $\hat{Q}^T y$  product and back substitution is backward stable in the context of least squares problems.

We claim that the QR factorization step remains backward stable if we consider the modified version described in this report. Without going into details with an extended proof, this can be explained by saying that at each step of the algorithm we apply a transformation  $L_i$  doing a smaller number of operations than those of the standard algorithm. Then, since we know that each step of the standard QR factorization is backward stable, this must be true also in the modified version of the algorithm.

Considering that both versions of the QR algorithm are backward stable, the accuracy of the algorithms will depend mostly on the conditioning of the least squares problem at hand. In fact, following from the definition of backward stability, the algorithms will both find exact solutions to slightly perturbed problems, with perturbations of the order of machine precision. This implies that if the conditioning of the problem is high the real solutions to the perturbed problems are inevitably going to be inaccurate. If  $w_*$  is the exact solution to the least squares problem and  $\tilde{w}_*$  is the solution found with one of the QR based algorithms outlined above, the accuracy of the algorithms will therefore follow the general upper bounds of relative errors found in (Trefethen and Bau 1997, 131): one relative to perturbations of the matrix  $\hat{X}$ ,

$$\frac{\|\tilde{w}_* - w_*\|}{\|w_*\|} \leq (\kappa(\hat{X}) + \kappa(\hat{X})^2 \tan \theta) \frac{\|\delta \hat{X}\|}{\|\hat{X}\|}$$

and one relative to perturbations of the vector  $y$ ,

$$\frac{\|\tilde{w}_* - w_*\|}{\|w_*\|} \leq \left( \frac{\kappa(\hat{X})}{\cos \theta} \right) \frac{\|\delta y\|}{\|y\|}$$

where  $\theta$  is the angle such that  $\cos \theta = \frac{\|\hat{X}w_*\|}{\|y\|}$ . From these upper bounds we can expect that the algorithm will be more accurate when the angle  $\theta$  is near 0 and less accurate when it is near  $\frac{\pi}{2}$ , reminding that in our context

the value of  $\theta$  will depend on the value of the random vector  $y$ .

#### INPUT DATA

Because of the bottom  $I$  block, the input matrix  $\hat{X}$  has  $n$  linearly independent rows and consequently full column rank, therefore a unique solution is expected from the linear least squares problem regardless of the possible values  $y$ . Nonetheless, as discussed in the previous section, the conditioning of the problem is dependent on  $y$ , precisely on the angle  $\theta$  between the image of  $\hat{X}$  and  $y$ .

The problem of extracting a random vector is tackled in two different ways, choosing one or the other depending on the specific experiment. The first way is to extract  $m$  random variables using a normal distribution, constructing so the  $y$  vector by associating each component to a random variable. This is done by using the Numpy method `random.rand`.

When more control over the conditioning of the problem is instead required, the vector  $y$  is extracted considering the  $\theta$  angle as a constraint. Firstly a vector  $w \in \mathbb{R}^n$  is extracted using multiple normal random variables as previously discussed. Then a vector  $v$  perpendicular to the image is found, this can be done by selecting one of the rows of  $Q_2$ , from the  $QR$  factorization of  $\hat{X}$ . The following condition must then be enforced for the  $v$  vector.

$$\|v\| = \|\hat{X}w\| \tan \theta$$

If the previous condition holds the desired vector  $y$  it is obtainable by summing  $\hat{X}w + v$ .

#### EXPERIMENTAL RESULTS

Other than LBFGS also the Newton method and the BFGS have been implemented. We compare these three with both versions of the QR factorization and the standard least squares problem solver from `numpy`.

Plot  $\theta$  against  $\kappa$  for both  $y$  and  $\hat{X}$ . (1 plot)

Plot  $\theta$  against time, steps and residual for each method. (3 plots)

Plot  $t$  against time, steps and residual for LBFGS. (3 plots)

Table different initializations for LBFGS and BFGS with number of steps and perturbation. (2 tables)

#### CONCLUSIONS

#### BIBLIOGRAPHY

Liu, Dong C., and Jorge Nocedal. 1989. "On the Limited Memory BFGS Method for Large Scale Optimization." *Mathematical Programming* 45 (1-3): 503–28. <https://doi.org/10.1007/BF01589116>.

Nocedal, Jorge, and Stephen J. Wright. 2006. *Numerical Optimization*. 2nd ed. Springer Series in Operations Research. New York: Springer.

Trefethen, Lloyd N., and David Bau. 1997. *Numerical Linear Algebra*. Philadelphia: Society for Industrial; Applied Mathematics.