
Convex optimization and application for machine learning

Encadrants :

-Alexandre d'Aspremont

1 Introduction

This report concerns the third homework related to the course convex optimization and application to machine learning. It deals with the application of a particular *interior points* method : the barrier method. The interior point methods are recent methods coined to deal with the minimisation of objective functions subject to convex inequality constraints. Indeed, affine constraints are the only one convex, and it is well known how to take them into account in the minimization algorithm (i.e Newton's method for affine equality constraint). A barrier method consists in removing the inequality constraints with the help of a barrier function. This function has to diverge when x is approaching the unfeasible domain $\mathcal{D} = \{x | f_i(x) \leq 0 \quad \forall i \in \{1, \dots, m\}\}$. Such a procedure can be applied to problems of the form :

$$\begin{aligned} \min_x \quad & f_0(x) \\ \text{s.t.} \quad & Ax = b \\ & f_i(x) \leq 0 \quad \forall i \in \{1, \dots, m\} \end{aligned} \tag{1}$$

where f_0 and $(f_i)_i$ are convex. The barrier method establish then the following non-equivalent form :

$$\begin{aligned} \min_x \quad & t f_0(x) + \phi(x) \\ \text{s.t.} \quad & Ax = b \end{aligned} \tag{2}$$

where ϕ is the barrier playing the role of a penalisation. As stated sooner, one have to impose ϕ the following property : $\phi(x) = +\infty$ if $x \notin \mathcal{D}$. The barrier method only introduces two new parameters: t and μ . t will control how close we wish to be from the ideal solution. Indeed, if we increase t , we observe that the original objective function f_0 becomes the priority to minimize since we may have for a certain t big enough $tf_0 \gg \phi$. However, it will hold only in \mathcal{D} , yet imposing the inequality constraints. μ represents the step length we are taking every step to increase our previous calculated solution i.e at step $(n+1)$, $t^{(n+1)} = \mu t^n$. These two parameters are key, and will be crucial for this homework.

The hardest part in the method is not the running of the algorithm itself (a succession of centering step via Newton's method), but the initialisation of its parameters. There is not only one initial condition, but two : $x^{(0)}$ and $t^{(0)}$. The initial point $x^{(0)}$ has to be strictly feasible i.e $f_i(x^{(0)}) < 0 \quad \forall \quad i \in \{1, \dots, m\}$, which is not always possible. The setting of initialisation is known as phase I, and it requires its own treatment like for example another barrier method.

The main advantage of this method is that even though the problems are not strictly equivalent, it can give pretty good approximations of the real solution p^* , with a controlled computational error : $f_0(x^*(t)) - p^* \leq \frac{m}{t}$, at least in the convex case. Also, the computational cost is even controlled in certain case. It can indeed to be shown that $\#steps \approx \frac{\log(m/\epsilon t^{(0)})}{\log(\mu)}$ where ϵ is the desired accuracy. , as well as the number of Newton's method iteration in each centering step (≈ 40 iterations), which is a very important property. I will not detail further the functioning of the barrier method.

In this report, we will apply it to the LASSO (Least Absolute Shrinkage Operator and Selection Operator) objective, which is a classic least-norm square objective granted with a regularization term, in order to impose a certain regularity to the solution. Indeed, the l_1 norm term will tend to have the effect to produce sparsity in w .

2 Formulation of the problem

The LASSO optimization problem is formulated this way :

$$\min_{w \in \mathbf{R}^d} \quad \frac{1}{2} \|Xw - y\|_2^2 + \lambda \|w\|_1 \quad (\text{LASSO})$$

where $X \in \mathcal{M}_{n \times d}(\mathbf{R})$, $w \in \mathbf{R}^d$, $y \in \mathbf{R}^n$, $\lambda > 0$. $\|\cdot\|$ is the l_1 norm, defined as $\|w\|_1 = \sum_{i=1}^d |w_i|$. Let's first notice this problem is convex. The domain is trivially convex and so is the objective function, as the sum of two convex functions. Moreover, we can assure a solution exists since the objective function is coercive and it is lower semi-continuous. It is trivial to see that $\|Xw - y\|_2 \rightarrow +\infty$ and $\|w\|_1 \rightarrow +\infty$ when $\|w\|_2 \rightarrow +\infty$. This form is not appropriate to solve via the barrier method, so we will compute its dual and put it under the general form of Quadratic Program (QP) :

$$\begin{aligned} \min_{v \in \mathbf{R}^n} \quad & \frac{1}{2} v^T Q v + p^T v + r \\ \text{s.t.} \quad & A v = b \\ & G v \preceq h \end{aligned} \quad (\text{QP})$$

To do that, let's add an artificial variable to our problem (LASSO) by setting $z = Xw - y$. We

then have two variables, and we can express it with this form :

$$\begin{aligned} \min_{(z,w) \in \mathbf{R}^n \times \mathbf{R}^d} \quad & \frac{1}{2} \|z\|_2^2 + \lambda \|w\|_1 \\ \text{s.t.} \quad & z = Xw - y \end{aligned} \quad (\text{LASSO}')$$

We introduce the lagrangian $\mathcal{L}(z, w, v) = \frac{1}{2} \|z\|_2^2 + \lambda \|w\|_1 + v^T(z - Xw + y)$ with $v \in \mathbf{R}^n$ the lagrange multiplier. By taking the $\inf_{z,w} \sup_{v \in \mathbf{R}^n}$, one find back the original problem (LASSO'). To find the dual problem, we exchange the sup and inf, and we first have to compute :

$$\inf_{(z,w) \in \mathbf{R}^{n \times d}} \mathcal{L}(z, w, v) = \inf_{(z,w) \in \mathbf{R}^{n \times d}} \frac{1}{2} \|z\|_2^2 + \lambda \|w\|_1 + v^T(z - Xw + y)$$

We can split the variables z and w , so we have next : $\mathcal{L}(z, w, v) = \frac{1}{2} \|z\|_2^2 + \lambda \|w\|_1 + v^T(z - Xw + y) + \inf_{z \in \mathbf{R}^n} v^T z + \frac{1}{2} \|z\|_2^2 + v^T z + \inf_{w \in \mathbf{R}^d} \lambda \|w\|_1 - v^T Xw$.

- $\inf_{z \in \mathbf{R}^n} v^T z + \frac{1}{2} \|z\|_2^2$ is calculated easily by setting its gradient to 0. We have immediately $z^* = -v$ and so $\inf_{z \in \mathbf{R}^n} \frac{1}{2} \|z\|_2^2 + v^T z = -\frac{\|v\|_2^2}{2}$
- $\inf_{w \in \mathbf{R}^d} \lambda \|w\|_1 - v^T Xw = -\lambda \sup_{w \in \mathbf{R}^d} \frac{v^T Xw}{\lambda} - \|w\|_1 = -f^*(\frac{X^T v}{\lambda})$ where f is the l_1 norm function and f^* its Legendre-Fenchel conjugate. It is easy to see that

$$f^*(y) = \begin{cases} 0 & \text{if } \|y\|_\infty \leq 1 \\ +\infty & \text{otherwise} \end{cases}$$

A condition here is appearing : $\frac{X^T v}{\lambda} \in \text{dom}(g) = \{y \in \mathbf{R}^n \mid \|y\|_\infty \leq 1\}$ where g is obviously $g(v) = \inf_{(z,w) \in \mathbf{R}^{n \times d}} \mathcal{L}(z, w, v)$.

Eventually, we have $\inf_{(z,w) \in \mathbf{R}^{n \times d}} \mathcal{L}(z, w, v) = v^T y - f^*(X^T w) - \frac{\|v\|_2^2}{2} = v^T y - \frac{\|v\|_2^2}{2}$ if $v \in \text{dom}(g)$

We can finally present the dual problem :

$$\begin{aligned} \max_{v \in \mathbf{R}^n} \quad & v^T y - \frac{\|v\|_2^2}{2} \\ \text{s.t.} \quad & X^T v \leq \lambda e \\ & -X^T v \leq \lambda e \end{aligned}$$

where $e_d = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbf{R}^d$. By setting $Q = I_n$, $p = -y$, $F = \begin{pmatrix} X^T \\ -X^T \end{pmatrix}$, $h = \lambda e_{2d}$, $A = 0$

and $b = 0$, we find back the (QP) formulation, except for the minus sign in front of the inf. We implement the barrier method with the log barrier

$$\phi(x) = -\sum_{i=1}^{2d} \log(\lambda - f_i^T x)$$

, with $f_i \in \mathbf{R}^n$ are the row of the F. Therefore, the point is to solve the problem:

$$\min_{v \in \mathbf{R}^n} \frac{\|v\|_2^2}{2} - v^T y + \phi(v) \quad (\text{Interior Point})$$

To solve this minimization problem, I will use the classical Newton's method, described in the next section. We can one more time state that the problem is convex, there is no equality/inequality constraint, the feasible set is (trivially) convex, and the objective function is convex as the sum of two convex function (Q is definite positive). The existence is also assured as the objective function is coercive and lower semi-continuous.

3 Newton's method

The point of Newton's method is to resolve the following equation $F(x) = 0$, where $F : \mathbf{R}^n \rightarrow \mathbf{R}^m$. The Newton's method is part of a more general class of method : the gradient descent methods. First order conditions for (global in this case) minimum x^* is the well known Peano-Kantorovitch condition

$$\nabla f_t(x^*) \in (T_{x^*}X)^+$$

or rather $\nabla f_t(x^*) = 0$ in this unconstrained case. Therefore, we are going to apply Newton's method on $F = \nabla f$. So, starting from a certain point $x^{(n)}$, we want to compute a new point $x^{(n+1)}$ verifying the property : $f(x^{(n+1)}) \leq f(x^{(n)})$. The Newton's method propose the following point $x^{(n+1)} = x^{(n)} - \alpha^{(n)}(\nabla^2 f(x))^{-1} \nabla f(x)$, more precisely the descent direction $\Delta x = -(\nabla^2 f(x))^{-1} \nabla f(x)$. This method has very good convergence properties. Indeed, it is theoretically proved that the method has a q-quadratic rate of convergence, at least locally. Concerning the step size $\alpha^{(k)}$, it takes a different value at every step. Intuitively, it cannot be too big, in order to not jump on too many minimum, risking thus to miss the global one or creating a oscillatory behavior. And not too small either, one do not want to remain trapped within a local minimum. There is no local minimum in this case because of convexity, but if it is too slow, the desired precision will require too many iterations. In this work, I will use a backtracking line search technique known as Armijo's rule. Other rules exist such as Goldstein rule, Wolfe rule etc... It consists in choosing a coefficient $\alpha_k^{(n)}$ verifying

$$f(x_k + \alpha^{(k)} \Delta x_k) \leq f(x^{(k)}) - \omega_1 \alpha^{(k)} \nabla f(x)^T (\nabla^2 f(x))^{-1} \nabla f(x)$$

with $\omega_1 \in]0, 1[$. To find it, we proceed iteratively and start with $\alpha_0^{(k)} = 1$. At step i, if $\alpha_i^{(k)}$ doesn't verify this equation, we try with $\alpha_{i+1}^{(k)} = \tau \alpha_i^{(k)}$, with $\tau \in]0, 1/2[$. If it is not yet satisfied, we repeat the procedure. In my code, I have chosen the following parameters : $\omega = 10^{-4}$ and $\tau = 0.2$. I used another trick and add a counter, to prevent α from becoming too small, as I discussed above. Goldstein's rule is supposed to prevent that, but I didn't have enough time to implement it. Therefore, the maximum of k is $k_{max} = 6$.

An interesting point would be to say we can also apply the backtracking linear search on ∇f and not f . Indeed, setting ∇f to 0 could also be seen as minimizing the criteria $\|\nabla f\|_2^2$. Replacing up there f by $\|\nabla f\|_2^2$, and using the useful property of $\nabla^2 f \Delta x = -\nabla f$, the backtracking line search takes the nice form

$$\|\nabla f(x_k + \alpha^{(k)} \Delta x_k)\|_2^2 \leq \|\nabla f(x_k)\|_2^2 (1 - \omega_1 \alpha^{(k)})$$

Let's now root all this to our (QP) problem. We set here $f_t(x) = t(\frac{\|x\|_2^2}{2} - x^T y) + \phi(x) = t(\frac{1}{2}x^T Qx + p^T x) + \phi(x)$. We here have immediately :

$$\nabla f_t(x) = t(Qx + p) + \sum_{i=1}^{2d} \frac{f_i}{\lambda - f_i^T x}$$

and

$$\nabla^2 f_t(x) = tQ + \sum_{i=1}^{2d} \frac{f_i f_i^T}{(f_i^T x - \lambda)^2}$$

As advised by the guideline, μ took the values $\{2, 15, 50, 100\}$. I have tried the two backtracking strategies and both gives results. However, the gradient approach seems far superior, and do not need any kind of counter. Eventually, even though the backtracking strategy can appear as the simplest step in the algorithm, it is actually the more subtle one, because a fine tuning of the parameter can be necessary.

To stop Newton's descent, we'll use the following criteria :

$$\lambda^2(x) = \Delta x \nabla^2 f(x) \Delta x \leq 2\epsilon$$

where ϵ is an arbitrary constant. It tells us how much one can descent towards this direction. It is also a local estimator of how close we are from the minimum :

$$f(x) - p^* = \frac{1}{2}\lambda^2(x)$$

It was not necessary to implement a phase I method with another barrier method. Indeed, there is no equality constraint, and $v_O = O_n$ is an obvious strictly feasible point. We will start from this point.

4 Results

First thing I noticed was that the gradient-based backtracking strategy was more efficient and allowed much smaller ϵ . Therefore, I will show the results for both approach.

Without any contest possible, the second approach for backtracking is much more efficient, and is also much faster (less newton's method iteration to compute).

We observe also that in the classic case, the higher μ is not the better. On the contrary to the gradient-based approach, the higher μ performs the better, with around 10 iterations.

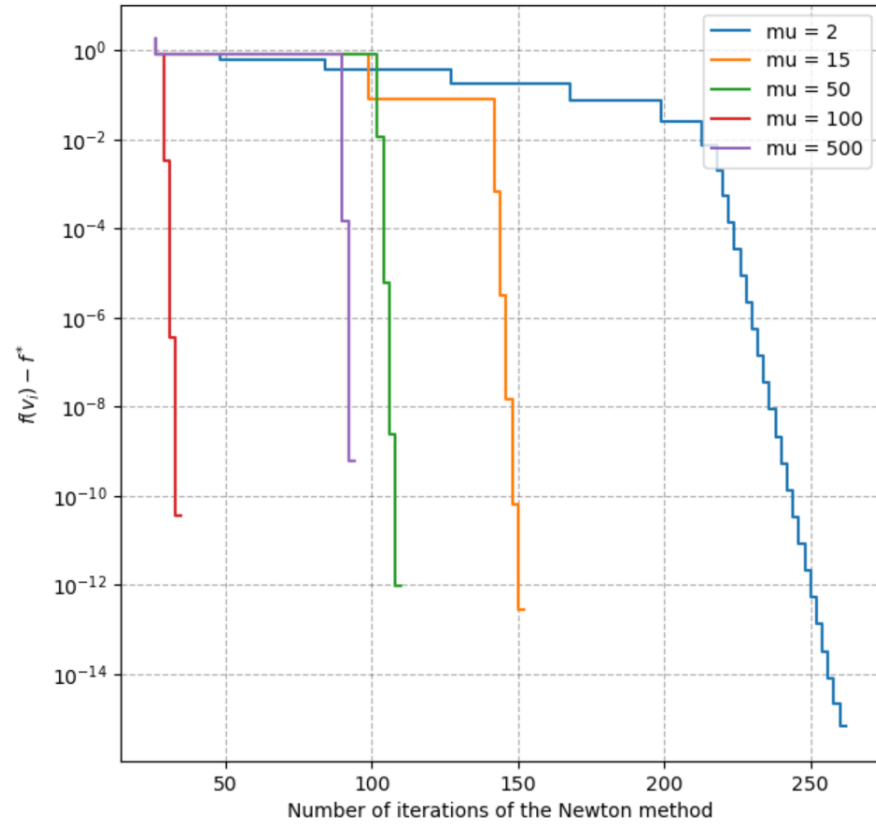


Figure 1: Interior point methods for different μ with $\epsilon = 10^{-8}$ and with classic backtracking

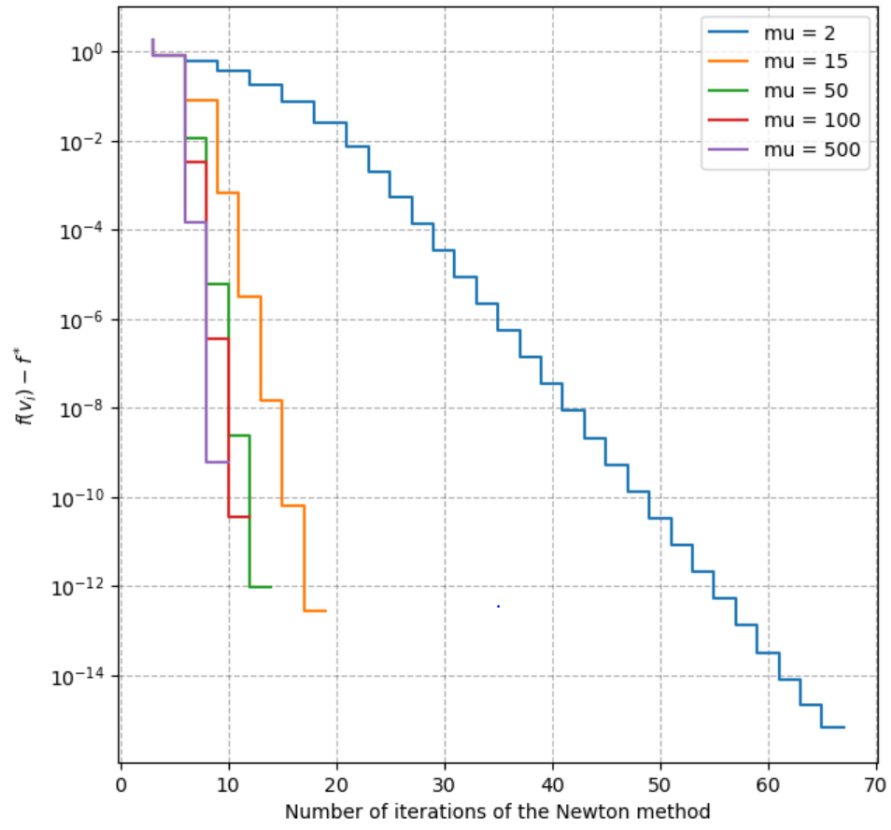


Figure 2: Interior point methods for different μ with $\epsilon = 10^{-8}$ and with gradient-based backtracking