

DE LORE Yonatan Convex Optimization Homework 3

1) * Let us derive the dual of the (LASSO) problem which is the following:

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

where $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$, $y = (y_1, \dots, y_n) \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}^{+*}$ are known.

$$\begin{aligned} (\text{LASSO}) \Rightarrow \text{Min}_{w \in \mathbb{R}^d, z \in \mathbb{R}^n} \quad & \frac{1}{2} \|z\|^2 + \lambda \|w\|_1 \\ \text{st.} \quad & z = Xw - y \end{aligned}$$

* The Lagrangian associated to this convex problem with linear equality constraints is the following:

$$(w, z, \mu) \in \mathbb{R}^d \times \mathbb{R}^n \times \mathbb{R}^n$$

$$\rightarrow \mathcal{L}(w, z, \mu) = \frac{1}{2} \|z\|^2 + \lambda \|w\|_1 + \mu^T (Xw - y - z)$$

$$= \frac{1}{2} \|z\|^2 - \mu^T z + \lambda \|w\|_1 + \mu^T Xw - \mu^T y$$

• $z \mapsto \frac{1}{2} \|z\|^2 - \mu^T z$ is a convex function which reaches its minimum when $z = \mu = 0 \Rightarrow z = \mu$
thus $\inf_z \frac{1}{2} \|z\|^2 - \mu^T z = \frac{1}{2} \|\mu\|^2 - \|\mu\|^2 = -\frac{1}{2} \|\mu\|^2$

$$\begin{aligned} \bullet \quad \lambda \|w\|_1 + \mu^T Xw &= -\lambda \left(-\frac{1}{\lambda} \mu^T Xw - \|w\|_1 \right) \\ &= -\lambda \left(-\frac{1}{\lambda} (X^T \mu)^T w - \|w\|_1 \right) \end{aligned}$$

(dual)

$$\begin{aligned} \text{Thus, } \inf_w \lambda \|w\|_1 + \mu^T X w &= - \sup_w \left[-\lambda \left(-\frac{1}{\lambda} X^T \mu \right)^T w - \|w\|_1 \right] \\ &= - f^* \left(-\frac{1}{\lambda} X^T \mu \right) \end{aligned}$$

where $f^*(y) = \sup_x y^T x - \|x\|_1$ is the conjugate function of the norm $\|\cdot\|_1$.

We already computed this conjugate in the previous homework. But for reminder,

$$\begin{aligned} f^*(y) &= \sup_x \left(\sum_{i=1}^d y_i x_i - \sum_{i=1}^d |x_i| \right) \quad (x \in \mathbb{R}^d) \\ &= \sup_x \left(\sum_{i=1}^d \underbrace{(y_i - \text{sgn}(x_i))}_{:= \varphi_y(x)} x_i \right) \end{aligned}$$

- if there exists $j \in \{1, d\}$ such that $y_j > 1$, we have considering the sequence $x(p) = (\underbrace{+p \delta_{ij}}_{1 \leq i \leq n})$ ($p \in \mathbb{N}$).

$$\varphi_y(x(p)) = \underbrace{p(y_j - 1)}_{> 0} \xrightarrow{p \rightarrow +\infty} +\infty$$

- similarly, if there exists j such that $y_j < -1$, considering the sequence $x(p) = (\underbrace{-p \delta_{ij}}_{1 \leq i \leq n})$, we have:

$$\varphi_y(x(p)) = \underbrace{-p(y_j + 1)}_{< 0} \xrightarrow{p \rightarrow +\infty} +\infty$$

- finally, if for all $j \in \{1, d\}$, $y_j \in [-1, 1]$,

$$\begin{aligned} |\varphi_y(x)| &\leq \sum_{i=1}^d |y_i x_i| - \sum_{i=1}^d |x_i| \\ &= \sum_{i=1}^d \underbrace{|y_i|}_{\leq 1} |x_i| - \sum_{i=1}^d |x_i| \\ &\leq 0 \end{aligned}$$

3/4

As $\varphi_y(0) = 0$, we can conclude that

$$f^*(y) = \begin{cases} 0 & \text{if } |y_i| \leq 1 \quad \forall i \in [1, d] \\ +\infty & \text{otherwise} \end{cases}$$

Hence $\inf_w \lambda \|w\|_1 + \mu^T X w = \begin{cases} 0 & \text{if } |-\frac{1}{\lambda} (X^T \mu)_i| \leq 1 \\ -\infty & \text{otherwise} \end{cases} \quad \forall i \in [1, d]$

* Finally we can conclude that:

$$\inf_{w, \lambda} \mathcal{L}(w, \lambda, \mu) = \begin{cases} -\frac{1}{2} \|\mu\|^2 - \mu^T y & \text{if } |(X^T \mu)_i| \leq \lambda \\ -\infty & \text{otherwise} \end{cases} \quad \forall i \in [1, d]$$

The dual of the (LASSO) problem is therefore:

$$\begin{aligned} \text{Max}_{\mu \in \mathbb{R}^m} \quad & -\frac{1}{2} \|\mu\|^2 - \mu^T y \\ \text{st.} \quad & |(X^T \mu)_i| \leq \lambda \quad \forall i \in [1, d] \end{aligned}$$

$$\Leftrightarrow \begin{aligned} \text{Min}_{\mu \in \mathbb{R}^m} \quad & \frac{1}{2} \|\mu\|^2 + y^T \mu \\ \text{st.} \quad & (X^T \mu)_i \leq \lambda \quad \forall i \in [1, d] \\ & -(X^T \mu)_i \leq \lambda \quad \forall i \in [1, d] \end{aligned}$$

$$\Leftrightarrow \begin{aligned} \text{Min}_{v \in \mathbb{R}^m} \quad & v^T Q v + p^T v \quad (\text{QP}) \\ \text{st.} \quad & A v \leq b \end{aligned}$$

where

$$\begin{aligned} Q &= \frac{1}{2} I_m & \in \mathbb{R}^{m \times m} \\ p &= y & \in \mathbb{R}^m \\ A &= (X^T, -X^T) & \in \mathbb{R}^{2d \times m} \\ b &= \lambda \mathbf{1}_{2d} & \in \mathbb{R}^{2d} \end{aligned}$$

where we denote:

I_p the identity matrix of dimension p
 $\mathbf{1}_p$ the vector of ones of size p
 (U, V) the concatenation of the columns of U and the columns of V

- 2) At each centering step of the barrier method, we are solving, for a certain t , the following unconstrained minimization problem:

$$\min_{v \in \mathbb{R}^m} t(v^T Q v + r^T v) - \sum_{i=1}^{2d} \log(b_i - (Av)_i)$$

As we will use the Newton method to solve it, we need to compute the gradient and the hessian of the objective function we will denote f_t .

$$f_t(v) = t(v^T Q v + r^T v) + \varphi(v) \quad \text{where} \quad \varphi(v) = - \sum_{i=1}^{2d} \log(b_i - (Av)_i)$$

$$\varphi(v) = - \sum_{i=1}^{2d} \log(b_i - \sum_{j=1}^m A_{ij} v_j)$$

$$\begin{aligned} * \text{ Then, } \frac{\partial \varphi}{\partial v_k} &= - \sum_{i=1}^{2d} \frac{-A_{ik}}{b_i - \sum_j A_{ij} v_j} = \sum_{i=1}^{2d} \frac{A_{ik}}{b_i - (Av)_i} \\ &= \sum_{i=1}^{2d} (A^T)_{ki} \frac{1}{b_i - (Av)_i} \end{aligned}$$

Thus

$$\nabla f_t = t(2Qv + r) + A^T \text{inv}(b - Av)$$

where $\text{inv}(x)$ is the vector x to which was applied the inverse function component-wise.

* We also get:

$$\begin{aligned} \frac{\partial^2 \varphi}{\partial v_k \partial v_l} &= \sum_{i=1}^{2d} A_{ik} \left(- \frac{-A_{il}}{(b_i - (Av)_i)^2} \right) \\ &= \sum_{i=1}^{2d} (A^T)_{ki} \frac{1}{(b_i - (Av)_i)^2} A_{il} \end{aligned}$$

Thus we have:

$$\nabla_t^2 f = 2tQ + A^T \text{Diag} \left(\left(\frac{1}{(b_i - (Av)_i)^2} \right)_{1 \leq i \leq 2d} \right) A$$

where $\text{Diag}(v)$ denotes the diagonal matrix whose diagonal is the vector v

Each Newton iteration will be therefore performed as:

$$v \leftarrow v - \alpha(v, f_t) (\nabla^2 f_t(v))^{-1} \nabla f_t(v)$$

where $\alpha(v, f_t)$ is the stepsize given by the backtracking line search

Experiments and comments

Below, I show results obtained for the following experiment :

I generated a random matrix X of size $(N=50, D=500)$ (values between -1 and 1), a vector of weights w of size D with 20 % of non-zero random entries (values between -1 and 1), and the vector of observations Y of size N by computing Xw .

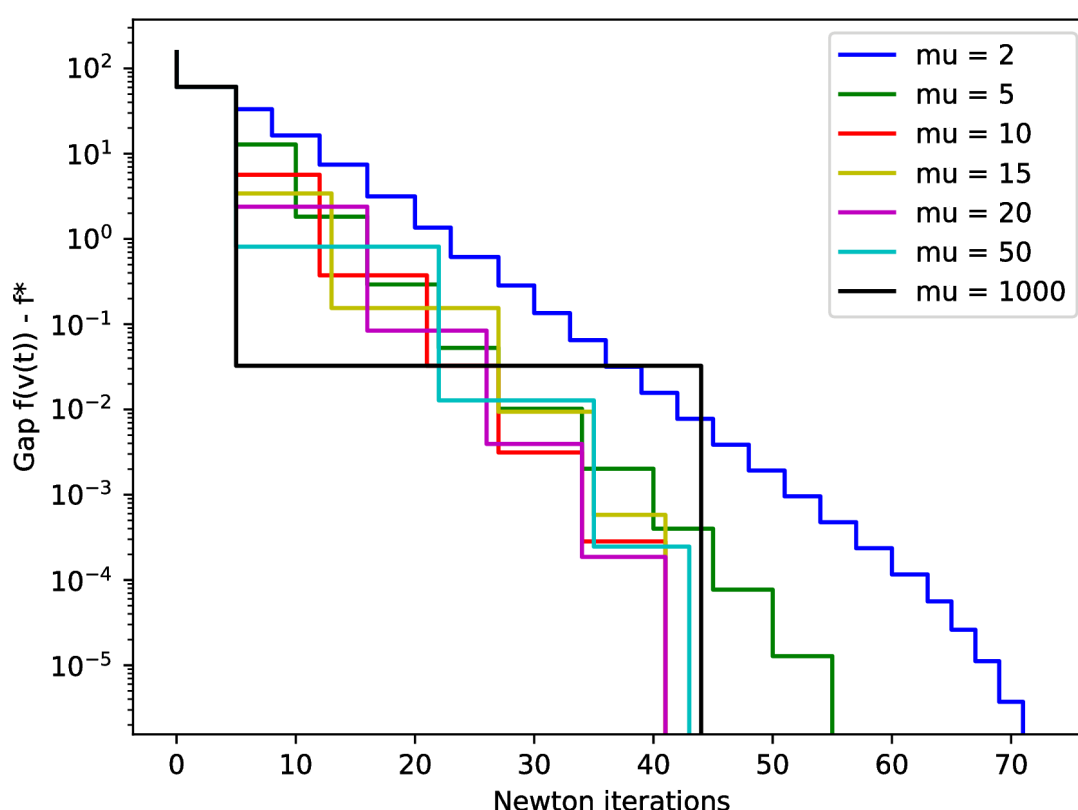
I launched the barrier method to solve (QP) (for $\lambda=10$) for various values of μ , on the same data (μ being the multiplicative factor by which we are increasing t after each centering step until convergence, t being the parameter of the barrier method introduced in page 4).

Precision criterion was set to 10^{-4} . Backtracking line search parameters were set to $\alpha=0.25$ and $\beta=0.5$, using the notations of the course.

On the figure below, I am plotting the gap $f(v(t)) - f^*$, where f is the objective function of (QP) ($f = v^T Q v + p^T v$), where $v(t)$ is the estimate of the dual solution at the end of the centering step of parameter t , and where f^* is estimated by $f(v_{\text{end}})$ with v_{end} the last estimate (for the precision criterion stated above). The length of the steps correspond to the number of Newton iterations needed per centering step.

We can see first that **the larger μ , the larger the number of Newton iterations needed per centering step but the smaller the number of outer iterations (centering steps) required for the precision criterion.**

In addition, we see that **the total number of Newton iterations (ie. for all centering steps) is the smallest for μ between 10 and 20**, range which can be seen as achieving a compromise between the number of inner Newton iterations and the number of centering steps.



However I did not manage to interpret consistently the final subquestion relative to the impact on w .

From the strong duality (convex problem, linear equality constraints), we have that $z^*=v^*$ (cf. question 1), where the dual variable v is first called μ) or equally said that $Xw^*=v^*+y$. Hence one could recover an estimate of the primal solution w^* from the estimates of the dual solution v^* using the pseudo inverse of X .

However, how to interpret the difference between the true w (the one we used to generate the data) and the ones estimated by the lasso, since the λ parameter of the lasso, ie. the weight of the L1-regularization, is not necessarily adequate with the number of non-zero entries we chose when generating w ? Therefore, I am not sure about which information we could get (in addition to the first plot) from such an error curb, regarding the choice of μ .