



**Problem 1.** Consider the optimization problem on slide 11 of Lecture 6. Show that for convex and closed  $f : A\mathbf{w} - \mathbf{b} \in \partial g(\boldsymbol{\lambda})$ , where  $\partial$  is the set of subgradients.

*Proof.*

**Def. 1 (Subgradient).** We say a vector  $\mathbf{c} \in \mathbb{R}^n$  is a subgradient of  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  at  $\mathbf{x} \in \text{dom } g$  if for all  $\mathbf{z} \in \text{dom } g$ ,

$$g(\mathbf{z}) \geq g(\mathbf{x}) + \mathbf{c}^T (\mathbf{z} - \mathbf{x}). \quad (1)$$

We know that  $g(\boldsymbol{\lambda}) := \inf_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}) = f(\mathbf{w}) + \boldsymbol{\lambda}^T (A\mathbf{w} - \mathbf{b})$ . Therefore,

$$g(\boldsymbol{\lambda}_1) = f(\mathbf{w}) + \boldsymbol{\lambda}_1^T (A\mathbf{w} - \mathbf{b}) \quad (2a)$$

$$g(\boldsymbol{\lambda}_2) = f(\mathbf{w}) + \boldsymbol{\lambda}_2^T (A\mathbf{w} - \mathbf{b}). \quad (2b)$$

Hence,

$$g(\boldsymbol{\lambda}_2) = g(\boldsymbol{\lambda}_1) + (A\mathbf{w} - \mathbf{b})^T (\boldsymbol{\lambda}_2 - \boldsymbol{\lambda}_1), \quad (3)$$

here the equality holds, which means that  $(A\mathbf{w} - \mathbf{b}) \in \partial g(\boldsymbol{\lambda})$ . □

**Problem 2.** Consider the dual ascent algorithm on slide 11 of Lecture 6. Analyze the convergence of dual ascent for  $L$ -smooth and  $\mu$ -strongly convex  $f$ . Is the solution primal feasible?

### Problem Description

Consider

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && z = f(\mathbf{w}) \\ & \text{subject to} && \mathbf{A}\mathbf{w} = \mathbf{b} \end{aligned} \quad (4)$$

Let the Lagrange dual function be

$$g(\boldsymbol{\lambda}) = \inf_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\lambda}) := f(\mathbf{w}) + \boldsymbol{\lambda}^T(\mathbf{A}\mathbf{w} - \mathbf{b}). \quad (5)$$

A lower bound for problem (4) can then be found by the Dual Ascent Algorithm in which each iteration consists of two steps. Step one is to update the primal variable according to

$$\mathbf{w}_{k+1} \in \underset{\mathbf{w}}{\text{argmin}} L(\mathbf{w}, \boldsymbol{\lambda}), \quad (6)$$

and step two is to update the dual variable according to

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k(\mathbf{A}\mathbf{w}_{k+1} - \mathbf{b}). \quad (7)$$

The problem is to analyze the convergence of dual ascent for  $L$ -smooth and  $\mu$ -strongly convex  $f$ , and then to check if the solution is primal feasible.

### Convergence Analysis - Primal Variable Update

The function  $L(\mathbf{w}, \boldsymbol{\lambda})$  consists of the sum of two functions  $f$  and  $\boldsymbol{\lambda}^T(\mathbf{A}\mathbf{w} - \mathbf{b})$ .  $f$  is  $\mu$ -strongly convex and  $L$ -smooth per the problem description, and  $\boldsymbol{\lambda}^T(\mathbf{A}\mathbf{w} - \mathbf{b})$  is convex since it's an affine function in  $\mathbf{w}$ . As we showed in HW1 [1], the sum of a convex function and a  $\mu$ -strongly convex function is still  $\mu$ -strongly convex, so  $L(\mathbf{w}, \boldsymbol{\lambda})$  is  $\mu$ -strongly convex.

Similarly, we can prove that  $L(\mathbf{w}, \boldsymbol{\lambda})$  is  $L$ -smooth. First, since  $\boldsymbol{\lambda}^T(\mathbf{A}\mathbf{w} - \mathbf{b})$  is affine, it can be upper-bounded by a quadratic with no curvature, in other words it is Lipschitz-smooth with Lipschitz-constant zero. Let's define  $g(\mathbf{x}) := \boldsymbol{\lambda}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$ , and call it  $L_2$ -smooth. We can then use a condition for smoothness from HW1 [1] on both  $g(\mathbf{x})$  and  $f(\mathbf{x})$  to get:

$$f(\mathbf{x}_2) \leq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T(\mathbf{x}_2 - \mathbf{x}_1) + \frac{L}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \quad (8)$$

$$g(\mathbf{x}_2) \leq g(\mathbf{x}_1) + \nabla g(\mathbf{x}_1)^T(\mathbf{x}_2 - \mathbf{x}_1) + \frac{L_2}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \quad (9)$$

Taking the sum of (8) and (9) yields:

$$L(\mathbf{w}, \boldsymbol{\lambda}) = g(\mathbf{x}_2) + f(\mathbf{x}_2) \leq g(\mathbf{x}_1) + f(\mathbf{x}_1) + (\nabla g(\mathbf{x}_1) + \nabla f(\mathbf{x}_1))^T(\mathbf{x}_2 - \mathbf{x}_1) + \frac{L_2 + L}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \quad (10)$$

which is equivalent to saying that  $L(\mathbf{w}, \boldsymbol{\lambda})$  is  $L_2 + L$ -smooth, but  $L_2 = 0$  so it is  $L$ -smooth. In total, we now know that, just like  $f$ ,  $L(\mathbf{w}, \boldsymbol{\lambda})$  is a  $\mu$ -strongly convex and  $L$ -smooth function. If we use gradient descent with constant step size to find the minimum of  $L(\mathbf{w}, \boldsymbol{\lambda})$ , it will converge in  $\mathcal{O}(\log(\frac{1}{\epsilon}))$  operations [2] (slide 19).

### Convergence Analysis - Dual Variable Update

Since  $\boldsymbol{\lambda}^T(\mathbf{A}\mathbf{w} - \mathbf{b})$  is in the set of subgradients for the lagrange dual function (as proven in the previous problem), the dual variable update is being performed with gradient ascent. So once again, if we can establish

strong convexity and Lipschitz-smoothness of the lagrange dual function  $g(\boldsymbol{\lambda})$ , the algorithm will converge in  $\mathcal{O}(\log(\frac{1}{\epsilon}))$  operations.

$g(\boldsymbol{\lambda})$  consists of the sum of two functions,  $f^*$  and  $\boldsymbol{\lambda}^T \mathbf{b}$ . As explained in the previous section, the sum of a strongly convex and Lipschitz-smooth function with an affine function is also strongly convex and Lipschitz-smooth. Since  $\boldsymbol{\lambda}^T \mathbf{b}$  is affine in  $\boldsymbol{\lambda}$  it is sufficient to show that  $f^*$  is strongly convex and Lipschitz-smooth.

For all  $\mu$ -strongly convex functions  $f$ , the convex conjugate  $f^*$  is  $\frac{1}{\mu}$ -smooth [5] (slide 19). Also, due to the biconjugate property [5] (slide 13), the opposite holds as well. Using these two properties we can establish that  $f^*$  of the  $f$  in our problem is  $\frac{1}{L}$ -strongly convex and  $\frac{1}{\mu}$ -smooth. Thus, we have established that  $g(\boldsymbol{\lambda})$  is strongly convex and Lipschitz smooth. Which means that the dual ascent algorithm will converge in  $\mathcal{O}(\log(\frac{1}{\epsilon}))$  operations.

### **Primal feasibility**

Since the primal problem is a convex programming problem with no inequality constraints and only linear equality constraints, we know that Slater's condition holds and thus that we have strong duality [4] (theorem 1). Since we have strong duality, the solution is not only primal feasible but also primal optimal.

**Problem 3.** Consider the optimization problem (P2) on slide 21 of Lecture 6. Extend the dual decomposition of Slide 6–12 to solve (P2). Compare it to the primal method (analytically or numerically) in terms of total communication cost and convergence rate on a random geometric communication graph.

*Proof.* Let us consider the problem

$$\begin{aligned} & \text{minimize } \frac{1}{N} \sum_{i \in [N]} f_i(w_i) \\ & \text{s.t. } w_i = w_j, \quad \text{for all } j \in \mathcal{N}_i \end{aligned} \quad (11)$$

The Lagrangian function in the dual formulation will be

$$\mathcal{L}(\lambda) = \frac{1}{N} \sum_{i \in [N]} f_i(w_i) - \sum_{i \in [N]} \sum_{j \in \mathcal{N}_i} \lambda_{i,j}^T (w_i - w_j) \quad (12)$$

where we denote the dual variable  $\lambda_{i,j} \forall i, j \in \mathcal{V}$ , and  $\lambda := [\lambda_{11}^T, \dots, \lambda_{1|\mathcal{N}_1|}^T, \lambda_{21}^T, \dots, \lambda_{2|\mathcal{N}_2|}^T, \dots, \lambda_{N|\mathcal{N}_N|}^T]^T$ . The corresponding dual problem with respect to the primal variables  $w_i, \forall i \in [N]$ :

$$\mathcal{D}(\lambda) = \inf_{w_i: i \in [N]} L(\lambda) = \inf_{w_i: i \in [N]} \frac{1}{N} \sum_{i \in [N]} f_i(w_i) - \sum_{i \in [N]} \sum_{j \in \mathcal{N}_i} \lambda_{i,j}^T (w_i - w_j). \quad (13)$$

We introduce the sub-problem:

$$\eta_i(\lambda) = \inf_{w_i} f_i(w_i) + \sum_{j \in \mathcal{N}_i} \lambda_{i,j}^T w_i \quad (14)$$

and we note that each sub-problem  $\eta_i(\lambda)$  can be solved locally at node  $i$  and independently of all other node  $j \neq i$ . The whole dual problem in Equation 13 can be written as the sum of these sub-problems

$$\mathcal{D}(\lambda) = \sum_{i \in [N]} \eta_i(\lambda) \quad (15)$$

We know that  $\mathcal{D}(\lambda)$  is a lower bound of the solution in the primal space. We are interested in the tightest lower bound that is given by the solution of the the following of the dual problem

$$d^* = \sup_{\lambda} \mathcal{D}(\lambda) = \sup_{\lambda} \sum_{i=1}^N \eta_i(\lambda). \quad (16)$$

The dual problem can be solve with sub-gradient method, by computing:

$$\lambda_{i,j}^{k+1} = \lambda_{i,j}^k + \alpha_k g(\mathcal{D}(\lambda^k)) \quad (17)$$

where  $g(\cdot)$  denotes the sub-gradient operator with respect to the dual variable  $\lambda$ . Since the subgradient can be computed and is equal to  $g(\mathcal{D}(\lambda)) = (w_{i,k+1} - w_{j,k+1})$ , where  $w_{i,k+1}$  and is solution of the minimization problem  $\eta_i(\lambda)$  evaluated at  $\lambda^k$

$$\lambda_{i,j}^{k+1} = \lambda_{i,j}^k + \alpha_k (w_{i,k+1} - w_{j,k+1}) \quad (18)$$

Dual decomposition algorithm  
 for  $k = 1, 2, \dots$ , do  
 for  $i \in [N]$  do  
     Master node broadcasts  $\lambda_{i,j}$  for all  $j \in \mathcal{N}_i$   
     Solve in parallel dual subproblems  $\eta_i$  finding  $(w_{i,k+1}) \in \operatorname{arginf}_{w_i} \eta_i(\lambda)$   
     Master node collects  $w_{1,k+1}$  and computes  $\lambda_{i,j}^{k+1} = \lambda_{i,j}^k + \alpha_k (w_{i,k+1} - w_{j,k+1})$

For what concerns the communication cost we have two contributions: the first concerns the distribution of the primal variables  $w_j$  such that  $j \in \mathcal{N}_i$  and the second is about accumulating the dual variables  $\lambda_{ij}$ . The cost will be

$$\sum_{i=1}^N (|\mathcal{N}_{ij}| - 1) \cdot (|\lambda_{ij}| + |w_{i,k+1}|) \quad (19)$$

where we subtract 1 to the neighbour's cardinality because we exclude the cost of the node  $i$  that does not to be communicated. Compared to the transmission cost of the dual, where only the transmission of the  $a_{i,j}$  weights for every node  $i \in [n]$  is required, we can see that the dual solution requires approximately double of the information cost with respect to the primal method.  $\square$



## References

- [1] Machine Learning over Networks, homework 1. <https://github.com/hshokrig/EP3260-MLoNs-2020/blob/master/Assignments/Homework1.pdf> Visited 2020-03-10
- [2] Machine Learning over Networks, lecture 2. [https://github.com/hshokrig/EP3260-MLoNs-2020/blob/master/Lectures%20\(2020\)/Lecture2.pdf](https://github.com/hshokrig/EP3260-MLoNs-2020/blob/master/Lectures%20(2020)/Lecture2.pdf) Visited 2020-03-10
- [3] Boyd, Stephen, Stephen P. Boyd, and Lieven Vandenbergh. Convex optimization. Cambridge university press, 2004.
- [4] Laurent El Ghaoui. Convex optimization, lecture 8. <https://people.eecs.berkeley.edu/~elghaoui/Teaching/EE227A/lecture8.pdf> Visited 2020-03-10
- [5] Lieven Vandenbergh, Optimization Methods for Large-Scale Systems, lecture 5 <http://www.seas.ucla.edu/~vandenbe/236C/lectures/conj.pdf> Visited 2020-03-10