

EP3260: Fundamentals of Machine Learning over Networks

Homework 3

Group 3

Notation

Upper-case letters with a double underline denotes matrices, e.g., $\underline{\underline{A}}$. Lower-case letters with a single underline denotes vectors, e.g., \underline{a} . The i -th element of the vector \underline{a} is denoted by either $a[i]$ or a_i , and element in the i -th row and j -th column of the matrix $\underline{\underline{A}}$ is denoted by $\underline{\underline{A}}[i, j]$. An i -th column vector of a matrix $\underline{\underline{A}}$ is denoted as either $\underline{\underline{A}}_i$ or $\underline{\underline{A}}[:, i]$.

HW3(a): Show that for convex and closed f : $\mathbf{Aw} - \mathbf{b} \in \partial g(\boldsymbol{\lambda})$, where $\partial g(\boldsymbol{\lambda})$ is the set of subgradients

The optimization problem is given by

$$\underset{\mathbf{w}}{\text{minimize}} \quad f(\mathbf{w}) \quad \text{subject to} \quad \mathbf{Aw} = \mathbf{b}. \quad (1)$$

The Lagrange dual function, utilizing the Lagrangian, can be written as

$$g(\boldsymbol{\lambda}) = \inf_{\mathbf{w}} \{L(\mathbf{w}, \boldsymbol{\lambda}) := f(\mathbf{w}) + \boldsymbol{\lambda}^T (\mathbf{Aw} - \mathbf{b})\} \quad (2)$$

The Lagrange dual problem, utilizing the conjugate function, can be expressed as

$$\underset{\boldsymbol{\lambda}}{\text{maximize}} \quad \{g(\boldsymbol{\lambda}) := -f^*(-\mathbf{A}^T \boldsymbol{\lambda}) - \boldsymbol{\lambda}^T \mathbf{b}\} \quad (3)$$

There can at least be two ways to show that $\mathbf{Aw} - \mathbf{b} \in \partial g(\boldsymbol{\lambda})$:

- Utilizing the KKT conditions, the primal feasibility condition, in particular (2), yields that $\mathbf{Aw} - \mathbf{b} \in \partial g(\boldsymbol{\lambda})$ since f is closed and convex function. In other words, taking the derivative of the Lagrangian

$$L(\mathbf{w}, \boldsymbol{\lambda}) := f(\mathbf{w}) + \boldsymbol{\lambda}^T (\mathbf{Aw} - \mathbf{b}) \quad (4)$$

with respect to the dual variable $\boldsymbol{\lambda}$ yields that

$$\mathbf{Aw} - \mathbf{b} \in \partial g(\boldsymbol{\lambda}).$$

- Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the conjugate function reads

$$f^*(\mathbf{y}) = \max_{\mathbf{x}} \mathbf{y}^T \mathbf{x} - f(\mathbf{x}) = - \left\{ \min_{\mathbf{x}} f(\mathbf{x}) - \mathbf{y}^T \mathbf{x} \right\}. \quad (5)$$

For a closed and convex function f , according to Fenchel duality theorem, any \mathbf{x} and \mathbf{y} ,

$$\mathbf{x} \in \partial f^*(\mathbf{y}) \iff \mathbf{y} \in \partial f(\mathbf{x}) \iff \mathbf{x} \in \arg \min_{\mathbf{z}} f(\mathbf{z}) - \mathbf{y}^T \mathbf{z}. \quad (6)$$

To this end, if we take the derivative of (3) with respect to $\boldsymbol{\lambda}$ and utilizing the chain rule, i.e.,

$$\partial g(\boldsymbol{\lambda}) = \mathbf{A} \partial f^*(-\mathbf{A}^T \boldsymbol{\lambda}) - \mathbf{b}. \quad (7)$$

Since $\mathbf{w} \in \partial f^*\left(\underbrace{-\mathbf{A}^T \boldsymbol{\lambda}}_{=\mathbf{y}}\right)$, then $\mathbf{Aw} - \mathbf{b} \in \partial g(\boldsymbol{\lambda})$, where

$$\mathbf{w} \in \arg \min_{\mathbf{z}} f(\mathbf{z}) - (-\mathbf{A}^T \boldsymbol{\lambda})^T \mathbf{z} = \arg \min_{\mathbf{z}} f(\mathbf{z}) + \boldsymbol{\lambda}^T \mathbf{Az}.$$

HW3(b): Analyze the convergence of dual ascent for L -smooth and μ -strongly convex f . Is the solution primal feasible?

The dual ascent scheme for the problem (1) can be summarized as

$$\mathbf{w}^{(k+1)} = \arg \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\lambda}^{(k)}) \in \arg \min_{\mathbf{w}} f(\mathbf{w}) + (\boldsymbol{\lambda}^{(k)})^T \mathbf{A} \mathbf{w} \quad (8)$$

$$\boldsymbol{\lambda}^{(k+1)} = \boldsymbol{\lambda}^{(k)} + \alpha^{(k)} (\mathbf{A} \mathbf{w}^{(k+1)} - \mathbf{b}) , \quad (9)$$

where $L(\mathbf{w}, \boldsymbol{\lambda}^{(k)})$ is defined in (4).

For brevity, here onwards, we redefine $\mathbf{x} := \mathbf{w}$ and $\mathbf{u} := \boldsymbol{\lambda}$.

Some preliminaries are presented in order to support the convergence of dual ascent scheme.

Since f is (strongly) convex, then according to Fenchel duality theorem, any \mathbf{x} and \mathbf{y} ,

$$\mathbf{x} \in \partial f^*(\mathbf{y}) \iff \mathbf{y} \in \partial f(\mathbf{x}) \iff \mathbf{x} \in \arg \min_{\mathbf{z}} f(\mathbf{z}) - \mathbf{y}^T \mathbf{z} . \quad (10)$$

In addition, f is strongly convex, i.e., $\partial f^*(\mathbf{y}) = \{\nabla f^*(\mathbf{y})\}$, then

$$\mathbf{x} = \nabla f^*(\mathbf{y}) = \arg \min_{\mathbf{z}} \{f(\mathbf{z}) - \mathbf{y}^T \mathbf{z}\} . \quad (11)$$

On the other hand, if a function, say $g(\mathbf{x})$, is also μ -strongly convex (where $\mu > 0$),

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 . \quad (12)$$

If \mathbf{x} is a minimizer, then $\nabla g(\mathbf{x}) = 0$ yields

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 . \quad (13)$$

To this end, we can define $g_{\mathbf{u}}(\mathbf{x}) := f(\mathbf{x}) - \mathbf{u}^T \mathbf{x}$ such that the minimizer $\mathbf{x}_{\mathbf{u}} = \nabla f^*(\mathbf{u})$, see (11), then for any \mathbf{y}

$$\underbrace{g(\mathbf{y})}_{=f(\mathbf{y})-\mathbf{u}^T \mathbf{y}} \geq \underbrace{g(\mathbf{x})}_{=f(\mathbf{x}_{\mathbf{u}})-\mathbf{u}^T \mathbf{x}_{\mathbf{u}}} + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad (14)$$

$$f(\mathbf{y}) - \mathbf{u}^T \mathbf{y} \geq f(\mathbf{x}_{\mathbf{u}}) - \mathbf{u}^T \mathbf{x}_{\mathbf{u}} + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 . \quad (15)$$

Similarly, we define $g_{\mathbf{v}}(\mathbf{x}) := f(\mathbf{x}) - \mathbf{v}^T \mathbf{x}$ such that the minimizer $\mathbf{x}_{\mathbf{v}} = \nabla f^*(\mathbf{v})$, then

$$f(\mathbf{y}) - \mathbf{v}^T \mathbf{y} \geq f(\mathbf{x}_{\mathbf{v}}) - \mathbf{v}^T \mathbf{x}_{\mathbf{v}} + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 . \quad (16)$$

Now, plugin $\mathbf{y} = \mathbf{x}_{\mathbf{v}}$ in (15) and $\mathbf{y} = \mathbf{x}_{\mathbf{u}}$ in (16), and thereby add them such that,

$$(\mathbf{u} - \mathbf{v})^T (\mathbf{x}_{\mathbf{u}} - \mathbf{x}_{\mathbf{v}}) \geq \mu \|\mathbf{x}_{\mathbf{u}} - \mathbf{x}_{\mathbf{v}}\|_2^2 . \quad (17)$$

Now, we invoke Cauchy-Schwartz inequality on the left hand side and rearrange such that

$$\left\| \underbrace{\mathbf{x}_{\mathbf{u}}}_{=\nabla f^*(\mathbf{u})} - \underbrace{\mathbf{x}_{\mathbf{v}}}_{=\nabla f^*(\mathbf{v})} \right\|_2 \geq \frac{1}{\mu} \|(\mathbf{u} - \mathbf{v})\|_2 \quad (18)$$

$$\|\nabla f^*(\mathbf{u}) - \nabla f^*(\mathbf{v})\|_2 \geq \frac{1}{\mu} \|(\mathbf{u} - \mathbf{v})\|_2 . \quad (19)$$

Thus, if f is μ -strongly convex, then ∇f^* is Lipschitz with parameter $\frac{1}{\mu}$.

To this end, applying what we know about the convergence rate of (primal) gradient descent.

- If f is μ -strongly convex, then the dual ascent with constant step size $\alpha^{(k)} \leq \mu$ converges at sublinear rate $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ same as the primal gradient descent.
- If f is μ -strongly convex and also L -smooth, i.e., ∇f is Lipschitz, then the primal gradient descent converges at linear rate $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$.
- Since we have shown above that if f is μ -strongly convex, then ∇f^* is Lipschitz with parameter $\frac{1}{\mu}$, where the converse is also true, i.e., if ∇f^* is Lipschitz with parameter $\frac{1}{\mu}$ then f is μ -strongly convex. For closed and convex f , the dual of the dual is primal, i.e., $f^{**} = f$. So, if f has Lipschitz gradient and it is strongly convex, then it is also true for f^* . Thus, dual ascent also converges at the linear rate of $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$.

Problem C

Problem Description Extend the dual decomposition of the following problem

$$\min \frac{1}{N} \sum_{i \in N} f_i(w_i) \quad (20)$$

$$\text{such that } w_i = w_j \quad \forall i \in \mathcal{N}_i \quad (21)$$

To extend the dual decomposition of the above we first convert the equality constraint to the $Aw = b$. Note that here $A = \text{lap}(G)$ such that $Aw = 0$ is a correct constraint, where the laplacian of a graph is defined as the following

$$\text{lap}(G) = D - \text{Adj}(G),$$

where D is a diagonal matrix with each entry $D_{ii} = \text{deg}(i)$ and $\text{deg}(i)$ is the degree of node of node i . Further, $\text{Adj}(G)$ is the adjacency matrix of the network G . For example for a triangle graph between three nodes 1, 2, 3, $\text{lap}(G)$ looks like the following

$$\text{lap}(\Delta) = \begin{bmatrix} -2 & 1 & 1 \\ 1 & -2 & 1 \\ 1 & 1 & -2 \end{bmatrix}$$

Hence the problem becomes

$$\min \frac{1}{N} \sum_{i \in N} f_i(w_i) \quad (22)$$

$$\text{such that } L_i w = 0 \quad \forall i \in \mathcal{N}_i \quad (23)$$

Where L_i is the i^{th} column vector of $\text{lap}(G)$

Dual Decomposition Here the dual of the problem is

$$\mathcal{L}(w, \lambda) = \sum_{i \in N} (f_i(w_i) + \lambda^T A_i w) - \frac{1}{N} \lambda^T b$$

Using the parameters from our problem this becomes

$$\mathcal{L}(w, \lambda) = \sum_{i \in N} (f_i(w_i) + \lambda^T L_i w)$$

Since the langrangian as described above is separable thus we can proceed with the usual Lagrangian separation on different nodes. Firstly we give the primal step.

Primal update for optimization In the primal step, each node minimizes $L_i(w_i, \lambda_k)$, that is for every node we need to update w_i^k which is given by the following

$$w_i^{k+1} \in \operatorname{argmin}_{w_i} L_i(w_i, \lambda_k)$$

for $i \in \{1, \dots, N\}$.

Dual update for optimization The dual update need to gather all the w_i within a connected component and update the dual as follows

$$\lambda_C^{k+1} = \lambda_C^k + \alpha_C^k \sum_{i \in [C]} A_i w_i^k$$

where C is a connected component. Recall that a connected component in an undirected graph is a set of nodes which has a path to each other. Note that here, for each component we can optimize the problem independently.

Convergence cost in Random Geometric Graph

A random geometric graph $G(n, r)$ is a graph of n nodes in $[0, 1]^2$ and two nodes are connected with an edge if their distance is at most r_n . It is well know fact that when $r_n^2 = (1 + \epsilon) \log n / n$ (see [9]) then the network is connected.

Recall that the diameter $\operatorname{diam}(G)$ of a graph is the maximum eccentricity of any vertex in the graph. That is, $\operatorname{diam}(G)$ is the greatest distance between any pair of vertices or, alternatively, $d = \max_{v \in V} \epsilon(v) d = \max_{v \in V} \epsilon(v) = \max_{v \in V} \max_{u \in V} d(v, u)$, where $d(u, v)$ is the distance between any two nodes u and v .

When a random graph is connected, then we need all the nodes in the graph to agree on the same w_i . This is similar to the *average consensus building problem*. From the lecture, we know that the convergence rate is $O(n \log n \log(\frac{1}{\epsilon}))$.

The interesting case is when $G(n, r)$ is disconnected. Here each update of the dual variable cost the $\operatorname{diam}(C)$, where C is component. The idea is as follows, every node updates its w_i and sends to a leader node of a C along a tree T edges. This tree T could be a BFS tree and the cost is equal to the $\operatorname{Depth}(T)$. This causes the leader to receive the sum in $\operatorname{Depth}(T) = O(\operatorname{Diam}(G))$ rounds all the values of w_i within the component. Note, that only aggregation of values are required and not all different values. Thus this requires $\operatorname{Depth}(T)$ rounds for one dual update. Here, we need to consider the largest diameter of a component, by [10], this is $\Theta(n^{1/2}/r)$. Hence the number of iterations required is $\Theta(n^{1/2}/r \log n \log(\frac{1}{\epsilon}))$

Comparison with Primal Method

In the primal method, we are required to compute the gradient that is

$$w_{i,k+1} = a_{ii} w_{i,k} - \alpha_k g_i(w_i, k) + \sum_{i \in \mathcal{N}_i \setminus \{i\}} a_{ij} w_{jk} \quad (24)$$

$$= \sum_{i \in \mathcal{N}_i} a_{ij} w_{jk} - \alpha_k g_i(w_i, k) \quad (25)$$

Here, instead of aggregation, each node needs to know the w_j for each node in the neighborhood and for every connected component C it will cost $|C|$ messages. Thus, this is more costly from the earlier described method.

References

- [1] Bubeck, Sbastien, "Convex optimization: Algorithms and complexity," Foundations and Trends in Machine Learning, vol. 8, no.3–4, pp. 231–357, 2015.
- [2] L. Bottou, F. Curtis, J. Norcedal, "Optimization Methods for Large-Scale Machine Learning," SIAM Rev., 60, no. 2, pp. 223-311, 2018.
- [3] T. Hastie, R. Tibshirani, and J. Friedman, "The Elements of Statistical Learning: Data Mining, Inference and Prediction," Second edition, Springer, 2009.
- [4] C. D. Meyer, "Matrix Analysis and Applied Linear Algebra," SIAM, 2000.
- [5] J. R. Magnus and H. Neudecker, "Matrix differential calculus with applications in statistics and econometrics," 2nd Edition, Wiley, UK, 1999.
- [6] K. P. Murphy, "Machine learning - a probabilistic perspective," MIT Press, 2012.
- [7] M. Schmidt, N. Le Roux, F. Bach, "Minimizing finite sums with the stochastic average gradient," Mathematical Programming, 2017.
- [8] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," NIPS, 2013.
- [9] Penrose, Mathew. Random geometric graphs. Vol. 5. Oxford university press, 2003.
- [10] Friedrich, Tobias, Thomas Sauerwald, and Alexandre Stauffer. "Diameter and broadcast time of random geometric graphs in arbitrary dimensions." Algorithmica 67.1 (2013): 65-88.