
CEE 690/ECE 590: Solutions of Homework 2

Wei Wen
Duke University
wei.wen@duke.edu

Abstract

This includes solutions for homework 2.

1 Problem 1

(a)

Please refer to Appendix B in [1] (page 83).

(b)

Refer to slide 30 and 31 in lecture 5,

$$\mathbb{E}[F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) \leq -\alpha_k [\nabla F(\mathbf{w}_k)]^T \mathbb{E}[\tilde{\mathbf{g}}_k] + \frac{\alpha_k^2 L}{2} \mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2]. \quad (1)$$

The goal is to maximize the decrease of loss function per iteration, *i.e.*, minimizing the left part of Inequality 1. Therefore, we want to minimize the the upper bound which is the right side of Inequality 1.

We define the right side of Inequality 1 as

$$D(\alpha_k) \triangleq -\alpha_k [\nabla F(\mathbf{w}_k)]^T \mathbb{E}[\tilde{\mathbf{g}}_k] + \frac{\alpha_k^2 L}{2} \mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2]. \quad (2)$$

In general, given loss function and dataset distribution, the optimal α_k satisfies

$$\frac{\partial D(\alpha_k)}{\partial \alpha_k} = -[\nabla F(\mathbf{w}_k)]^T \mathbb{E}[\tilde{\mathbf{g}}_k] + \alpha_k L \cdot \mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2] = 0, \quad (3)$$

that is,

$$\alpha_k = \frac{[\nabla F(\mathbf{w}_k)]^T \mathbb{E}[\tilde{\mathbf{g}}_k]}{L \cdot \mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2]} \quad (4)$$

Substituting $\alpha_k = \frac{1}{cL}$ to Equation (4), we have

$$\mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2] = c [\nabla F(\mathbf{w}_k)]^T \mathbb{E}[\tilde{\mathbf{g}}_k], \quad (5)$$

which is the condition/regime where the optimal learning rate is $\alpha_k = \frac{1}{cL}$. Suppose we have unbiased gradient, *i.e.*,

$$\mathbb{E}[\tilde{\mathbf{g}}_k] = \nabla F(\mathbf{w}_k) \triangleq \mathbf{g}_k, \quad (6)$$

then,

$$\mathbb{E}[\|\tilde{\mathbf{g}}_k\|_2^2] = \text{Var}(\|\tilde{\mathbf{g}}_k\|_2) + \|\mathbf{g}_k\|_2^2 = c \|\mathbf{g}_k\|_2^2, \quad (7)$$

that is,

$$\text{Var}(\|\tilde{\mathbf{g}}_k\|_2) = (c - 1) \cdot \|\mathbf{g}_k\|_2^2. \quad (8)$$

where $c \geq 1$.

In case that variance is zero, *i.e.*, $Var(\|\tilde{\mathbf{g}}_k\|_2) = 0$, then $\alpha_k = \frac{1}{L}$. Otherwise, when c is larger, *i.e.*, a larger variance in the gradient estimation, the optimal learning rate α_k will be smaller.

When the variance is bounded as $Var(\|\tilde{\mathbf{g}}_k\|_2) \leq M$,

$$(c - 1) \cdot \|\mathbf{g}_k\|_2^2 \leq M, \quad (9)$$

so, the optimal learning rate is

$$\alpha_k = \frac{1}{cL} \geq \frac{\|\mathbf{g}_k\|_2^2}{L(M + \|\mathbf{g}_k\|_2^2)}. \quad (10)$$

2 Problem 2

Code and analyses will be pushed here.

References

- [1] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018.