

## Problem 1.

(a) Fundamental theorem of calculus:

$$F(\vec{y}) = F(\vec{x}) + \int_0^1 (\vec{y} - \vec{x})^T \nabla F(\vec{x} + t(\vec{y} - \vec{x})) dt$$

Lipschitz Continuous:

$$\|\nabla F(\vec{y}) - \nabla F(\vec{x})\|_2 \leq L \|\vec{y} - \vec{x}\|_2$$

$$\begin{aligned} F(\vec{y}) &= F(\vec{x}) + \int_0^1 (\vec{y} - \vec{x})^T \nabla F(\vec{x} + t(\vec{y} - \vec{x})) dt \\ &= F(\vec{x}) + \int_0^1 \nabla F(\vec{x} + t(\vec{y} - \vec{x}))^T (\vec{y} - \vec{x}) dt \\ &= F(\vec{x}) + \nabla F(\vec{x})^T (\vec{y} - \vec{x}) + \int_0^1 [\nabla F(\vec{x} + t(\vec{y} - \vec{x})) - \nabla F(\vec{x})]^T (\vec{y} - \vec{x}) dt \end{aligned}$$

By the definition of Lipschitz continuous, we have

$$\|\nabla F(\vec{x} + t(\vec{y} - \vec{x})) - \nabla F(\vec{x})\|_2 \leq L \|t(\vec{y} - \vec{x})\|_2$$

Then

$$F(\vec{y}) \leq F(\vec{x}) + \nabla F(\vec{x})^T (\vec{y} - \vec{x}) + \int_0^1 L \|t(\vec{y} - \vec{x})\|_2 \|\vec{y} - \vec{x}\|_2 dt$$

$$F(\vec{y}) \leq F(\vec{x}) + \nabla F(\vec{x})^T (\vec{y} - \vec{x}) + \frac{1}{2} L \|\vec{y} - \vec{x}\|_2^2$$

(b) Lemma (sequence for stochastic gradient):  $\vec{W}_{k+1} = \vec{W}_k - \alpha_k \hat{g}_k$  satisfies in expectation  $E[F(\vec{W}_{k+1})] - F(\vec{W}_k) \leq -\alpha_k [\nabla F(\vec{W}_k)]^T$ 

$$E[F(\vec{W}_{k+1})] - F(\vec{W}_k) \leq -\alpha_k [\nabla F(\vec{W}_k)]^T E[\hat{g}_k] + \frac{\alpha_k^2 L}{2} E[\|\hat{g}_k\|_2^2]$$

Want to minimize R.H.S

$$\text{let } f(\alpha_k) = -\alpha_k [\nabla F(\vec{W}_k)]^T E[\hat{g}_k] + \frac{\alpha_k^2 L}{2} E[\|\hat{g}_k\|_2^2]$$

$$\frac{df(\alpha_k)}{d\alpha_k} = 0 = -[\nabla F(\vec{W}_k)]^T E[\hat{g}_k] + \alpha_k L E[\|\hat{g}_k\|_2^2]$$

By definition of variance,  $\text{Var}[\|\hat{g}_k\|_2^2] = E[\|\hat{g}_k\|_2^2] - E[\|\hat{g}_k\|_2]^2$ 

So we have

$$\alpha_k = \frac{[\nabla F(\vec{W}_k)]^T E[\hat{g}_k]}{L (\text{Var}[\|\hat{g}_k\|_2^2] + E[\|\hat{g}_k\|_2]^2)} = \frac{\|\hat{g}_k\|_2^2}{L (M + E[\|\hat{g}_k\|_2^2])}, \quad \text{var}(\|\hat{g}_k\|_2) \leq M$$

if we consider an unbiased gradient estimator, and assume that the variance is bounded, we have  $E[\hat{g}_k] = g_k$ ,  $\text{var}(\|\hat{g}_k\|_2^2) \leq M$ If ①  $\alpha_k = \frac{1}{L} = \frac{\|\hat{g}_k\|_2^2}{L(M + E[\|\hat{g}_k\|_2^2])}$ , we know that  $M$  (the upper bound of the variance) must be 0.②  $\alpha_k = \frac{1}{2L} = \frac{\|\hat{g}_k\|_2^2}{L(M + E[\|\hat{g}_k\|_2^2])}$ , we have  $M = \|\hat{g}_k\|_2^2$  for  $\alpha_k$  to be optimal③  $\alpha_k = \frac{1}{10L} = \frac{\|\hat{g}_k\|_2^2}{L(M + E[\|\hat{g}_k\|_2^2])}$ , we have  $M = 9\|\hat{g}_k\|_2^2$  for  $\alpha_k$  to be optimal  
In general, if  $\alpha_k = \frac{1}{cL}$ , the variance is bounded to be  $\leq (c-1)\|\hat{g}_k\|_2^2$  for  $\alpha_k$  to be optimal.