

Homework 2

Problem 1 We have seen in class that for strongly convex and Lipschitz smooth functions, the convergence of the gradient descent algorithm is exponentially fast; whereas for convex and Lipschitz continuous functions, the convergence rate is $\tilde{O}(1/\sqrt{K})$. You must have wondered which factor is more important in contributing to the faster convergence in the first case, the strong convexity or the smoothness. Now let's explore this problem by considering an objective function that is convex (but not strongly convex) and β -Lipschitz smooth.

1. First prove what we left off in class that if a function \mathcal{L} is convex and β -Lipschitz smooth (i.e., $\mathcal{L}(x) - \mathcal{L}(y) \leq \langle \nabla \mathcal{L}(y), x - y \rangle + \frac{\beta}{2} \|x - y\|^2$), then

$$\mathcal{L}(x) - \mathcal{L}(y) \geq \langle \nabla \mathcal{L}(y), x - y \rangle + \frac{1}{2\beta} \|\nabla \mathcal{L}(x) - \nabla \mathcal{L}(y)\|^2.$$

2. What's the convergence rate of gradient descent over such a function \mathcal{L} that is convex (but not strongly convex) and β -Lipschitz smooth? Choose a set of step sizes h_k and prove the result using the technique we learned.

Problem 2 Consider the following stochastic gradient descent algorithm.

Algorithm 1: Stochastic gradient descent algorithm

Input: θ_0

for $k = 0, \dots, K - 1$ **do**

 Sample an i.i.d. set S_k uniformly at random from $\{1, \dots, N\}$

 Compute stochastic gradient $\tilde{\mathcal{L}}(\theta_k|S_k) = \frac{1}{|S_k|} \sum_{i \in S_k} \nabla \text{loss}(\theta_k|z_i)$

 Update $\theta_{k+1} = \theta_k - h_k \nabla \tilde{\mathcal{L}}(\theta_k|S_k)$

end

Return: θ_K

Assume that the loss functions $\text{loss}(\theta|z)$ are convex and L -Lipschitz continuous ($\|\nabla \text{loss}(\theta|z)\| \leq L$) for all data points z (No need to make assumptions about the stochastic gradient variance in this case). What's the convergence rate of the stochastic gradient descent algorithm over risk $\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \nabla \text{loss}(\theta|z_i)$? Choose a set of step sizes h_k and prove the result using the technique we learned.

Using your derivation, what's the optimal choice of the batch size $n = |S_k|$? Is the stochastic gradient descent algorithm more efficient than gradient descent in this case?

Machine Learning Assignment - 02

Problem 1

1. To prove:-

$$L(x) - L(y) \geq \langle \nabla L(y), x - y \rangle + \frac{1}{2\beta} \|\nabla L(x) - \nabla L(y)\|^2$$

Now, we know that if a function is convex:-
 $L(z) - L(y) \geq \nabla L(y)^T (z - y)$

Condition for Lipschitz smoothness:-

$$L(x) - L(z) \geq \nabla L(x)^T (x - z) - \frac{\beta}{2} \|x - z\|^2$$

$$\textcircled{1} \Rightarrow L(x) - L(z) + L(z) - L(y) \geq \nabla L(y)^T (z - y) + \nabla L(x)^T (x - z) - \frac{\beta}{2} \|x - z\|^2$$

To choose z , w.r.t z and equating
 Differentiating the right side, to 0

$$0 = \nabla L(y) - \nabla L(x) - \beta(x - z)$$

$$z = x + \frac{1}{\beta} (\nabla L(x) - \nabla L(y))$$

Substituting z in $\textcircled{1}$

$$\begin{aligned} \text{L.H.S} &\geq \nabla L(y)^T \left(x + \frac{1}{\beta} (\nabla L(x) - \nabla L(y)) - y \right) + \\ &\quad \nabla L(x)^T \left(x - x - \frac{1}{\beta} (\nabla L(x) - \nabla L(y)) \right) - \frac{\beta}{2} \|x - z\|^2 \end{aligned}$$

$$\begin{aligned}
&\geq \nabla L(y)^T (x-y) + \frac{1}{\beta} \nabla L(y)^T (\nabla L(x) - \nabla L(y)) \\
&+ \frac{1}{\beta} \nabla L(x)^T (\nabla L(x) - \nabla L(y)) \\
&- \frac{\beta}{2} \|x - y\|^2 - \frac{1}{\beta} \|\nabla L(x) - \nabla L(y)\|^2
\end{aligned}$$

$$\begin{aligned}
&\geq \nabla L(y)^T (x-y) + \frac{1}{\beta} (\nabla L(x) - \nabla L(y))^T (\nabla L(y) - \nabla L(x)) \\
&- \frac{1}{2\beta} \|\nabla L(x) - \nabla L(y)\|^2
\end{aligned}$$

$$\begin{aligned}
&\geq \nabla L(y)^T (x-y) + \frac{1}{\beta} \|\nabla L(x) - \nabla L(y)\|^2 - \\
&\quad \frac{1}{2\beta} \|\nabla L(x) - \nabla L(y)\|^2
\end{aligned}$$

$$\geq \nabla L(y)^T (x-y) + \frac{1}{2\beta} \|\nabla L(x) - \nabla L(y)\|^2$$

Hence Proved

Problem 1from β -Lipschitz Smooth eqⁿ:-

$$2. L(x) - L(y) \leq \langle \nabla L(y), (x-y) \rangle + \frac{\beta}{2} \|x-y\|^2$$

$$\Rightarrow L(x_{t+1}) - L(x_t) \leq \langle \nabla L(x_t), (x_{t+1} - x_t) \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$

L (1)

Differentiating R.H.S w.r.t x_{t+1} & equate to 0

$$\Rightarrow 0 = \nabla L(x_t) + \beta \|x_{t+1} - x_t\|$$

e.g. D eqn:-

$$\Rightarrow x_{t+1} = x_t - \frac{1}{\beta} \nabla L(x_t) \Rightarrow x_{t+1} - x_t = -\frac{1}{\beta} \nabla L(x_t)$$

L (2)

So we will make most progress when

$$h_t = \frac{1}{\beta}$$

$$\Rightarrow L(x_{t+1}) - L(x_t) \leq \langle \nabla L(x_t), (x_{t+1} - x_t) \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$

L (3)

Now, substituting equation (2) in (3)

$$L(x_t) - L(x_{t+1}) \geq \frac{1}{2\beta} \|\nabla L(x_t)\|^2 \quad - (4)$$

$$(L(x_t) - L(x^*))^2 \leq (\nabla L(x_t)^T (x_t - x^*))^2 \quad \left[\begin{array}{l} \text{Squaring} \\ \text{convexity} \\ \text{eqn} \end{array} \right]$$

$$(L(x_t) - L(x^*))^2 \leq \|\nabla L(x_t)\|^2 \|x_t - x^*\|^2$$

$$\frac{(L(x_t) - L(x^*))^2}{\|x_t - x^*\|^2} \leq \|\nabla L(x_t)\|^2$$

Substituting (4) in above eqn

$$\frac{(L(x_t) - L(x^*))^2}{\|x_t - x^*\|^2} \leq 2\beta [L(x_t) - L(x_{t+1})]$$

Because $\|x_t - x^*\| \leq \|x_0 - x^*\|$

$$\frac{(L(x_t) - L(x^*))^2}{\|x_0 - x^*\|^2} \leq 2\beta [L(x_t) - L(x_{t+1})]$$

Put $a_t = L(x_t) - L(x^*)$

$a_{t+1} = L(x_{t+1}) - L(x^*)$

$$\frac{(a_t)^2}{\|x_0 - x^*\|^2} \leq 2\beta (a_t - a_{t+1})$$

Divide the whole inequality by $a_t a_{t+1}$

$$\frac{a_t}{a_{t+1}} \frac{1}{\|x_0 - x^*\|^2} \leq 2\beta \left(\frac{1}{a_{t+1}} - \frac{1}{a_t} \right)$$

$$-\frac{1}{2\beta} \frac{a_t}{a_{t+1}} \frac{1}{\|x_0 - x^*\|^2} \leq -\frac{1}{a_{t+1}} + \frac{1}{a_t}$$

Because $\frac{a_t}{a_{t+1}} = \frac{L(x_t) - L(x^*)}{L(x_{t+1}) - L(x^*)} > 1$

Multiply R.H.S by $\frac{a_{k+1}}{a_k} < 1$ - inequality remains same

$$\sum_{t=0}^{T-1} \frac{1}{a_{t+1}} - \frac{1}{a_t} \geq \sum_{t=0}^{T-1} \frac{1}{2\beta \|x_0 - x^*\|^2}$$

$$\frac{1}{a_T} - \frac{1}{a_0} \geq (T) \frac{1}{2\beta \|x_0 - x^*\|^2}$$

Ignoring $\frac{1}{a_0}$ on L.H.S

$$\frac{1}{a_T} \geq \frac{T}{2\beta \|x_0 - x^*\|^2}$$

$$\boxed{L(x_T) - L(x^*) \leq \frac{2\beta \|x_0 - x^*\|^2}{T}}$$

∇ converges $O\left(\frac{1}{T}\right)$ with step size $\frac{1}{\beta}$

Problem 2

$\mathcal{L}(\theta|z) \sim$ convex L -Lipschitz continuous for all point z

$\mathbb{E}_{z \sim 0} \|\nabla \mathcal{L}(\theta_k|z_k)\|^2 \leq b^2$ for some constant b^2 variance bounded - (1)

$$\|\theta_{k+1} - \theta^*\|^2 = \|\theta_k - \theta^*\|^2 - 2h_k \langle \nabla \mathcal{L}(\theta_k), \theta_k - \theta^* \rangle + h_k^2 \|\nabla \mathcal{L}(\theta_k)\|^2$$

$$\langle \nabla \mathcal{L}(\theta_k), \theta_k - \theta^* \rangle = \frac{1}{|S_k|} \sum_{i \in S_k} \langle \nabla \text{loss}(\theta_k|z_i), \theta_k - \theta^* \rangle$$

$$= \frac{1}{|S_k|} \sum_{i \in S_k} \langle \nabla \text{loss}(\theta_k|z_i), \theta_k - \theta^* \rangle$$

$$\Rightarrow \sum_{i \in Z_i} \langle \nabla \text{loss}(\theta_k|z_i), \theta_k - \theta^* \rangle \leq \sum_{i \in Z} \text{loss}(\theta_k|z_i) - \sum_{i \in Z_i} \text{loss}(\theta^*|z_i)$$

$$\Rightarrow \frac{1}{|S_k|} \langle \nabla \mathcal{L}(\theta_k), \theta_k - \theta^* \rangle \leq \frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta_k|z_i) - \frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta^*|z_i)$$

$$\Rightarrow \|\nabla \mathcal{L}(\theta_k)\|^2 = \left\| \frac{1}{|S_k|} \sum_{i \in S_k} \nabla \text{loss}(\theta_k|z_i) \right\|^2 \leq \frac{1}{|S_k|^2} \left(\sum_{i \in S_k} \|\nabla \text{loss}(\theta_k|z_i)\| \right)^2$$

\Rightarrow Since $\text{loss}(\theta|z_i)$ is convex and L -Lipschitz continuous

$$\|\nabla \text{loss}(\theta|z)\| \leq L$$

$$\|\nabla \mathcal{L}(\theta_k|z_i)\|^2 \leq \frac{1}{|S_k|^2} \cdot |S_k|^2 \cdot L^2 = L^2$$

$$\Rightarrow \|\theta_{k+1} - \theta^*\|^2 \leq \|\theta_k - \theta^*\|^2 - 2h_k \left[\frac{1}{S_k} (\text{loss}(\theta_k | z_i) - \text{loss}(\theta^* | z_i)) \right] + (h_k L)^2$$

$$\Rightarrow \|\theta_{k+1} - \theta^*\|^2 - \|\theta_k - \theta^*\|^2 - (h_k L)^2 \leq -2h_k \left[\frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta_k | z_i) - \frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta^* | z_i) \right]$$

$$\Rightarrow 2h_k \left[\frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta_k | z_i) - \frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta^* | z_i) \right] \leq \|\theta_k - \theta^*\|^2 - \|\theta_{k+1} - \theta^*\|^2 + (h_k L)^2$$

✓ Take expectation

$$\Rightarrow -2h_k \left[\mathbb{E} \left[\frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta_k | z_i) \right] - \mathbb{E} \left[\frac{1}{|S_k|} \sum_{i \in Z_i} \text{loss}(\theta^* | z_i) \right] \right] \leq \mathbb{E} \|\theta_k - \theta^*\|^2 - \mathbb{E} \|\theta_{k+1} - \theta^*\|^2 + (h_k L)^2$$

$$\Rightarrow 2h_k [L(\theta_k) - L(\theta^*)] \leq \mathbb{E} \|\theta_k - \theta^*\|^2 - \mathbb{E} \|\theta_{k+1} - \theta^*\|^2 + (h_k L)^2$$

✓ Take expectation of θ_k

$$\Rightarrow 2h_k \mathbb{E}_{\theta_k} [L(\theta_k) - L(\theta^*)] \leq \mathbb{E}_{\theta_k} \|\theta_k - \theta^*\|^2 - \mathbb{E}_{\theta_k} \|\theta_{k+1} - \theta^*\|^2 + (h_k L)^2$$

$$\Rightarrow 2 \sum_{k=0}^{k+1} h_k \mathbb{E} [L(\theta_k) - L(\theta^*)] \leq \mathbb{E} \|\theta_k - \theta^*\|^2 + L^2 \sum_{k=0}^{k-1} h_k^2$$

$$\Rightarrow \sum_{k=0}^{k+1} h_k \frac{\mathbb{E} [L(\theta_k)] - \mathbb{E} [L(\theta^*)]}{\sum_{i=0}^{k-1} h_i} \leq \frac{1}{2} \frac{\mathbb{E} \|\theta_k - \theta^*\|^2}{\sum_{i=0}^{k-1} h_i} + \frac{1}{2} \frac{\mathbb{E} \|\theta_k - \theta^*\|^2}{\sum_{i=0}^{k-1} h_i} + \frac{L^2 \sum_{k=0}^{k-1} h_k^2}{\sum_{i=0}^{k-1} h_i}$$

Since $\mathbb{E} \|\theta_k - \theta^*\|^2 \sum_{i=0}^{k-1} h_i > 0$

$$\sum_{k=0}^{k+1} h_k \frac{\mathbb{E} [L(\theta_k)] - \mathbb{E} [L(\theta^*)]}{\sum_{i=0}^{k-1} h_i} \leq \frac{\mathbb{E} \|\theta_k - \theta^*\|^2}{\sum_{i=0}^{k-1} h_i} + \frac{L^2 \sum_{k=0}^{k-1} h_k^2}{\sum_{i=0}^{k-1} h_i}$$

Let $h_k = \sqrt{\mathbb{E} \|\theta_k - \theta^*\|^2} / L \sqrt{k+1}$

$$\Rightarrow \frac{E\|\theta_0 - \theta^*\|^2}{2 \sum_{i=0}^{k-1} h_i} + \frac{L^2}{2} \frac{\sum_{k=0}^{k-1} h_k}{\sum_{i=0}^{k-1} h_i} \leq \frac{L}{2} \sqrt{E\|\theta_0 - \theta^*\|^2} \left(\frac{1 + \sum_{k=1}^k \frac{1}{k}}{\sum_{i=1}^k \frac{1}{\sqrt{k+1}}} \right)$$

$$\Rightarrow \leq \frac{L}{2} \sqrt{E\|\theta_0 - \theta^*\|^2} \left(\frac{i+1 + \int_1^k \frac{1}{u} du}{\int_1^k \frac{1}{u} du} \right)$$

$$\leq \frac{L}{2} \sqrt{E\|\theta_0 - \theta^*\|^2} \left(\frac{1 + \frac{1}{k} \log(k)}{(\sqrt{k} - 1)} \right)$$

Since

$$E[L(\theta) - L(\theta^*)] = E\left[L\left(\frac{\sum_{k=0}^{k-1} h_k \theta_k}{\sum_{i=0}^{k-1} h_i}\right) - L(\theta^*)\right]$$

By convex property \rightarrow

$$\leq \sum_{k=0}^{k-1} \frac{h_k}{\sum_{i=0}^{k-1} h_i} \left[E[L(\theta_k)] - E[L(\theta^*)] \right]$$

$$\leq \frac{L}{2} \sqrt{E\|\theta_0 - \theta^*\|^2} \left(\frac{1 + \frac{1}{k} \log(k)}{\sqrt{k} - 1} \right)$$

So convergence rate $\sim \tilde{O}\left(\frac{1}{\sqrt{k}}\right)$

(ii) Since $\mathbb{E}[L(\theta) - L(\theta^*)] \sim \tilde{O}\left(\frac{1}{\sqrt{k}}\right)$, the optimal

choice of iid $n = |S_k|$ that optimal choice is $n=1$ with same optimal choice for gradient descent - convergence rate of GD

$\sim \tilde{O}\left(\frac{1}{\sqrt{k}}\right)$ but for each iteration, GD have

gradient more than SGD \Rightarrow (SGD is more efficient)