

MLoN Homework2

Group 1

1 Problem 1

Consider the logistic ridge regression loss function:

$$\text{minimize}_{\mathbf{w}} f(\mathbf{w}) = \frac{1}{N} \sum_{i \in [N]} f_i(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \quad (1)$$

where $f_i(\mathbf{w}) = \log(1 + \exp\{-y_i \mathbf{w}^T \mathbf{x}_i\})$

1.1

To prove the Lipschitz continuity, if we have $\|\mathbf{w}\|_2 \leq D$, $\lambda \geq 0$

$$\begin{aligned} \nabla f(\mathbf{w}) &= 2\lambda \mathbf{w} - \frac{1}{N} \sum_{i \in [N]} \frac{y_i \mathbf{x}_i}{1 + e^{y_i \mathbf{w}^T \mathbf{x}_i}} \\ \|\nabla f(\mathbf{w})\|_2 &\leq 2\lambda \|\mathbf{w}\|_2 + \frac{1}{N} \sum_{i \in [N]} \frac{\|y_i \mathbf{x}_i\|_2}{|1 + e^{y_i \mathbf{w}^T \mathbf{x}_i}|} \\ &\leq 2\lambda D + \frac{1}{N} \sum_{i \in [N]} \|y_i \mathbf{x}_i\|_2 \end{aligned} \quad (2)$$

Assume that $\|y_i \mathbf{x}_i\|_2 \leq E$, then

$$\|\nabla f(\mathbf{w})\|_2 \leq 2\lambda D + E \quad (3)$$

So $f(\mathbf{w})$ is Lipschitz continuous.

To find a small B, we apply the whole human Activity Recognition Using Smartphones dataset, where \mathbf{x}_i is normalized to range $[-1, 1]$ and y_i has range $[0, 5]$ (standing, sitting, laying, walking, walking_downstairs, walking_upstairs):

$$E = \frac{1}{N} \sum_{i \in [N]} \|y_i \mathbf{x}_i\|_2 = 36.078 \quad (4)$$

Then we have:

$$B = 2\lambda D + E = 2\lambda D + 36.078 \quad (5)$$

1.2

To prove the smoothness of $f_i(\mathbf{w})$, we only need to prove that $f_i(\mathbf{w})$ is differentiable (obviously) and $\nabla f_i(\mathbf{w})$ is Lipschitz continuous:

$$\begin{aligned}\|\nabla^2 f_i(\mathbf{w})\|_2 &= \left\| \frac{(y_i \mathbf{x}_i)^2}{e^{y_i \mathbf{w}^T \mathbf{x}_i} + 2 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}} \right\|_2 \\ &\leq \frac{\|y_i \mathbf{x}_i\|_2^2}{4} \\ &\leq \frac{25}{4}\end{aligned}\tag{6}$$

Thus we can conclude that $\nabla f_i(\mathbf{w})$ is Lipschitz continuous, then $f_i(\mathbf{w})$ is L-smoothness and L is 6.25

As for $f(\mathbf{w})$, to apply the whole human Activity Recognition Using Smartphones dataset,

$$\begin{aligned}\|\nabla^2 f(\mathbf{w})\|_2 &= \left\| 2\lambda + \frac{1}{N} \sum_{i \in [N]} \frac{(y_i \mathbf{x}_i)^2}{e^{y_i \mathbf{w}^T \mathbf{x}_i} + 2 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}} \right\|_2 \\ &\leq 2\lambda + \frac{1}{4N} \sum_{i \in [N]} \|y_i \mathbf{x}_i\|_2^2 \\ &\leq 2\lambda + 465.870\end{aligned}\tag{7}$$

Thus $f(\mathbf{w})$ is L-smoothness with $L = 2\lambda + 465.870$

1.3

To prove the strongly-convexity of $f(\mathbf{w})$,

$$\nabla^2 f(\mathbf{w}) = 2\lambda + \frac{1}{N} \sum_{i \in [N]} \frac{(y_i \mathbf{x}_i)^2}{e^{y_i \mathbf{w}^T \mathbf{x}_i} + 2 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}} \geq 2\lambda\tag{8}$$

Thus $f(\mathbf{w})$ is μ -strongly convex with $\mu = 2\lambda$

2 Problem 2

Problem Let us assume that there exist scalars $c_0 \geq c > 0$ such that for all $k \in N$

$$\nabla f(\mathbf{w}_k)^T \mathbb{E}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)] \geq c \|\nabla f(\mathbf{w}_k)\|_2^2 \quad (1a)$$

$$\|\mathbb{E}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)]\|_2 \leq c_0 \|\nabla f(\mathbf{w}_k)\|_2 \quad (1b)$$

Furthermore, let us assume that there exist scalars $M \geq 0$ and $M_v \geq 0$ such that for all $k \in N$:

$$\text{Var}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)] \leq M + M_v \|\nabla f(\mathbf{w}_k)\|_2^2 \quad (2)$$

For the convergence proof of SGD with an L-smooth convex objective function (see slides), prove that

$$\mathbb{E}_{\zeta_k}[\|g(\mathbf{w}_k; \zeta_k)\|_2^2] \leq \alpha + \beta \|\nabla f(\mathbf{w}_k)\|_2^2$$

Proof Start with definition of variance:

$$\begin{aligned} \text{Var}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)] &= \mathbb{E}_{\zeta_k}[\langle g(\mathbf{w}_k; \zeta_k), g(\mathbf{w}_k; \zeta_k) \rangle] - \langle \mathbb{E}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)], \mathbb{E}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)] \rangle \\ &\quad (\text{Second - moment boundary}) \\ &\leq M + M_v \|\nabla f(\mathbf{w}_k)\|_2^2 \\ &\Rightarrow \\ \mathbb{E}_{\zeta_k}[\langle g(\mathbf{w}_k; \zeta_k), g(\mathbf{w}_k; \zeta_k) \rangle] &\leq M + M_v \|\nabla f(\mathbf{w}_k)\|_2^2 + \langle \mathbb{E}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)], \mathbb{E}_{\zeta_k}[g(\mathbf{w}_k; \zeta_k)] \rangle \\ &\quad (\text{First - moment boundary}) \\ &\leq M + M_v \|\nabla f(\mathbf{w}_k)\|_2^2 + c_0^2 \|\nabla f(\mathbf{w}_k)\|_2^2 \\ &= M + (M_v + c_0^2) \|\nabla f(\mathbf{w}_k)\|_2^2 \end{aligned}$$

Thus, proved with $\alpha = M, \beta = M_v + c_0^2$.

3 Problem 3

Problem For the SGD with non-convex objective function, prove that with square summable but not summable step-size, we have for any $K \in \mathbb{N}$

$$\mathbb{E}[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2] < \infty \quad (9)$$

and therefore:

$$\lim_{k \rightarrow \infty} \mathbb{E}[\frac{\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2}{\sum_{k \in [K]} \alpha_k}] = 0 \quad (10)$$

Proof Using the inequality (4) and (5) in the slice of lecture 3, we have:

$$\begin{aligned} & \mathbb{E}[f(\omega_{k+1})] - f(\omega_k) \\ & \leq -(c - \frac{1}{2}\alpha_k LM_G)\alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 + \frac{1}{2}\alpha_k^2 LM \\ & \leq -(c - \frac{c}{2LM_G}LM_G)\alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 + \frac{1}{2}\alpha_k^2 LM \\ & \leq -\frac{c}{2}\alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 + \frac{1}{2}\alpha_k^2 LM \end{aligned} \quad (11)$$

Recursively $\forall k \in [K]$, take total expectation for both side and do summation, we have:

$$\begin{aligned} & \mathbb{E}[f(\omega_{k+1})] - \mathbb{E}[f(\omega_1)] \\ & \leq -\frac{c}{2}\mathbb{E}[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2] + \frac{1}{2}LM \sum_{k \in [K]} \alpha_k^2 \end{aligned} \quad (12)$$

Rearrange (12), we have:

$$\begin{aligned} & \mathbb{E}[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2] \\ & \leq -\frac{2\mathbb{E}[f(\omega_{k+1})] - \mathbb{E}[f(\omega_1)]}{c} + \frac{LM \sum_{k \in [K]} \alpha_k^2}{c} \end{aligned} \quad (13)$$

The step size is square summable $\sum_{k \in [K]} \alpha_k^2 < \infty$, the right hand size is smaller than ∞ , so we have:

$$\mathbb{E}[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2] < \infty \quad (14)$$

The step size is not summable $\sum_{k \in [K]} \alpha_k = \infty$, so we have:

$$\lim_{k \rightarrow \infty} \mathbb{E}[\frac{\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2}{\sum_{k \in [K]} \alpha_k}] = 0 \quad (15)$$