

Problem 1. Consider Human Activity Recognition Using Smartphones dataset $\{(x_i, y_i)\}_{i \in [N]}$, with inputs defined as the accelerometer and gyroscope sensors, and outputs defined as moving (e.g., walking, running, dancing) or not (sitting or standing). Consider the logistic ridge regression loss function

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

where $f_i(\mathbf{w}) = \log (1 + \exp \{-y_i \mathbf{w}^T \mathbf{x}_i\})$. Then, address the following questions:

(a) Is f Lipschitz continuous? If so, find a small B?

Proof. We need to show that

$$\|\boldsymbol{w}\|_{2} \leq D \Rightarrow \|\nabla f(\boldsymbol{w})\|_{2} \leq B.$$

We know that for $h: \mathbb{R} \to \mathbb{R}$ and $g: \mathbb{R}^d \to \mathbb{R}$, we have

$$\nabla h\left(g\left(\boldsymbol{z}\right)\right) = h'\left(g\left(\boldsymbol{z}\right)\right) \nabla g\left(\boldsymbol{z}\right). \tag{2}$$

By substituting $h(z) = \log(z)$ and $g(\boldsymbol{w}) = 1 + \exp\{-y_i \boldsymbol{w}^T \boldsymbol{x}_i\}$ in (2), we have

$$\nabla f_i(\mathbf{w}) = h'(g(\mathbf{w})) \nabla g(\mathbf{w})$$

$$= \frac{-y_i}{1 + \exp\{y_i \mathbf{w}^T \mathbf{x}_i\}} \mathbf{x}_i.$$
(3)

Therefore, for $\|\boldsymbol{w}\|_2 \leq D$, we have

$$\|\nabla f_i(\boldsymbol{w})\|_2 = \frac{\|y_i\| \|\boldsymbol{x}_i\|_2}{1 + \exp\{y_i \boldsymbol{w}^T \boldsymbol{x}_i\}} \le \frac{\|y_i\| \|\boldsymbol{x}_i\|_2}{1 + \exp\{-D\|y_i\| \|\boldsymbol{x}_i\|_2\}} \triangleq B_{f_i}$$

For $r(\boldsymbol{w}) = \lambda \|\boldsymbol{w}\|_2^2$ and $\|\boldsymbol{w}\|_2 \leq D$, we have

$$\|\nabla r(\boldsymbol{w})\|_2 = 2\lambda \|\boldsymbol{w}\|_2 \le 2\lambda D \triangleq B_r$$

Theorem 1. Suppose that f_1, \ldots, f_n are Lipschitz continuous on I with Lipschitz constants B_1, \ldots, B_n , respectively. Then the linear combination $c_1f_1 + \cdots + c_nf_n$ is Lipschitz continuous on I with Lipschitz constant $|c_1|B_1 + \cdots + |c_n|B_n$.

According to (1), for $\|\boldsymbol{w}\|_2 \leq D$, $f(\boldsymbol{w})$ is a linear combination of Lipschitz continuous functions. Therefore, we have

$$\|\nabla f(\boldsymbol{w})\|_{2} \leq \frac{1}{N} \sum_{i \in [N]} \frac{|y_{i}| \|\boldsymbol{x}_{i}\|_{2}}{1 + \exp\{-D|y_{i}| \|\boldsymbol{x}_{i}\|_{2}\}} + 2\lambda D \triangleq B$$

Assume that $|y_i| \|\boldsymbol{x}_i\|_2 \leq C$ for all $i \in [N]$, then

$$\left\|\nabla f\left(\boldsymbol{w}\right)\right\|_{2} \leq \frac{C}{1 + \exp\left\{-DC\right\}} + 2\lambda D \triangleq B$$



(b) Is f_i smooth? If so, find a small L for f_i ? What about f?

Proof. We know that Hessian of a function is the Jacobian of its gradient. Assuming $J(\cdot)$ is the Jacobian operator, we have

$$\nabla^{2} f_{i}(\boldsymbol{w}) = J(\nabla f_{i}(\boldsymbol{x})) = \boldsymbol{x}_{i} \boldsymbol{x}_{i}^{T} \frac{y_{i}^{2} \exp\{y_{i} \boldsymbol{w}^{T} \boldsymbol{x}_{i}\}}{(1 + \exp\{y_{i} \boldsymbol{w}^{T} \boldsymbol{x}_{i}\})^{2}}$$

$$\leq \frac{\boldsymbol{x}_{i} \boldsymbol{x}_{i}^{T} y_{i}^{2}}{4}$$
(4a)

$$\leq \sigma_{\max} \left(\frac{\boldsymbol{x}_i \boldsymbol{x}_i^T \boldsymbol{y}_i^2}{4} \right) I$$
(4b)

$$= \frac{y_i^2 \sigma_{\text{max}} \left(\boldsymbol{x}_i^T \boldsymbol{x}_i \right)}{4} I \tag{4c}$$

$$= \frac{y_i^2 \|\mathbf{x}_i\|_2^2}{4} I, \tag{4d}$$

where in (4a) we used the inequality $\frac{\exp\{x\}}{(1+\exp\{x\})^2} \leq \frac{1}{4}$. In (4b), $\sigma_{\max}(X)$ is the largest eigenvalue of a (symmetric positive semidefinite) matrix X. In (4c), we used the facts that $\operatorname{eig}(AB) = \operatorname{eig}(BA)$. Hence, f_i is L_i -smooth with

$$L_i \le \frac{y_i^2 \|\boldsymbol{x}_i\|_2^2}{4}.$$

Similarly for f we have

$$\nabla^{2} f(\boldsymbol{w}) = J(\nabla f(\boldsymbol{x})) = \frac{1}{N} \sum_{i \in [N]} \boldsymbol{x}_{i} \boldsymbol{x}_{i}^{T} \frac{y_{i}^{2} \exp\left\{y_{i} \boldsymbol{w}^{T} \boldsymbol{x}_{i}\right\}}{\left(1 + \exp\left\{y_{i} \boldsymbol{w}^{T} \boldsymbol{x}_{i}\right\}\right)^{2}} + 2\lambda I$$
 (5a)

$$\leq \frac{1}{4N} \sum_{i \in [N]} \boldsymbol{x}_i \boldsymbol{x}_i^T y_i^2 + 2\lambda I \tag{5b}$$

$$= \frac{1}{4N}AA^T + 2\lambda I \tag{5c}$$

$$\leq \sigma_{\text{max}} \left(\frac{1}{4N} A A^T + 2\lambda I \right) I$$
(5d)

$$= \frac{1}{4N} \left(\sigma_{\text{max}} \left(A^T A \right) + 2\lambda \right) I, \tag{5e}$$

where $A \triangleq [y_1 x_1, y_2 x_2, \dots, y_N x_N]$. In (5e), we used the facts that $\operatorname{eig}(X + \lambda I) = \operatorname{eig}(X) + \lambda$ and $\operatorname{eig}(AB) = \operatorname{eig}(BA)$. Hence, f is L-smooth with

$$L \le \frac{1}{4N} \sigma_{\max} \left(A^T A \right) + 2\lambda.$$

(c) Is f strongly convex? If so, find a high μ ?

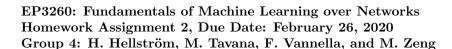
Proof. From (5a) and for all $\boldsymbol{v} \in \mathbb{R}^d$, we have

$$\mathbf{v}^{T} \left(\nabla^{2} f\left(\mathbf{w}\right) - 2\lambda I \right) \mathbf{v} = \frac{1}{N} \sum_{i \in [N]} \frac{y_{i}^{2} \exp\left\{ y_{i} \mathbf{w}^{T} \mathbf{x}_{i} \right\}}{\left(1 + \exp\left\{ y_{i} \mathbf{w}^{T} \mathbf{x}_{i} \right\} \right)^{2}} \mathbf{v}^{T} \mathbf{x}_{i} \mathbf{x}_{i}^{T} \mathbf{v}$$

$$= \frac{1}{N} \sum_{i \in [N]} \frac{y_{i}^{2} \exp\left\{ y_{i} \mathbf{w}^{T} \mathbf{x}_{i} \right\}}{\left(1 + \exp\left\{ y_{i} \mathbf{w}^{T} \mathbf{x}_{i} \right\} \right)^{2}} \left\| \mathbf{v}^{T} \mathbf{x}_{i} \right\|_{2}^{2}$$

$$(6)$$

$$\geq 0. \tag{7}$$





Therefore, $\nabla^{2} f(\boldsymbol{w}) - 2\lambda I \succeq 0$ or equivalently $\nabla^{2} f(\boldsymbol{w}) \succeq 2\lambda I$. Hence, f is strongly convex with $\mu = 2\lambda$. In summary we have

$$\underbrace{2\lambda I}_{\mu I} \preceq \nabla^{2} f\left(\boldsymbol{w}\right) \preceq \underbrace{\left(\frac{1}{4N}\sigma_{\max}\left(A^{T}A\right) + 2\lambda\right)I}_{LI},$$

where
$$A \triangleq [y_1 \boldsymbol{x}_1, y_2 \boldsymbol{x}_2, \dots, y_N \boldsymbol{x}_N].$$



Problem 2. Let us assume that there exist scalars $c_0 \geq c > 0$ such that for all $k \in \mathbb{N}$

$$\nabla f(\boldsymbol{w}_k)^T \mathbb{E}_{\zeta_k} \left[g(\boldsymbol{w}_k; \zeta_k) \right] \ge c \left\| \nabla f(\boldsymbol{w}_k) \right\|_2^2, \tag{8a}$$

$$\|\mathbb{E}_{\zeta_k} \left[g\left(\boldsymbol{w}_k; \zeta_k \right) \right] \|_2 \le c_0 \|\nabla f\left(\boldsymbol{w}_k \right) \|_2. \tag{8b}$$

Furthermore, let us assume that there exist scalars $M \geq 0$ and $M_V \geq 0$ such that for all $k \in \mathbb{N}$

$$\operatorname{Var}_{\zeta_k}\left[g\left(\boldsymbol{w}_k;\zeta_k\right)\right] \le M + M_V \left\|\nabla f\left(\boldsymbol{w}_k\right)\right\|_2^2. \tag{9}$$

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

$$\mathbb{E}_{\zeta_k} \left[\|g\left(\boldsymbol{w}_k; \zeta_k\right)\|_2^2 \right] \le \alpha + \beta \|\nabla f\left(\boldsymbol{w}_k\right)\|_2^2.$$
(10)

Proof. By assumption 9 we know that

$$\operatorname{Var}_{\zeta_{k}}\left[g\left(\boldsymbol{w}_{k};\zeta_{k}\right)\right] = \mathbb{E}_{\zeta_{k}}\left[\left\|g\left(\boldsymbol{w}_{k};\zeta_{k}\right)\right\|_{2}^{2}\right] - \left\|\mathbb{E}_{\zeta_{k}}\left[g\left(\boldsymbol{w}_{k};\zeta_{k}\right)\right]\right\|_{2}^{2}$$

$$\leq M + M_{V}\left\|\nabla f\left(\boldsymbol{w}_{k}\right)\right\|_{2}^{2}.$$
(11)

This is equivalent to

$$\mathbb{E}_{\zeta_{k}}\left[\left\|g\left(\boldsymbol{w}_{k};\zeta_{k}\right)\right\|_{2}^{2}\right] \leq M + M_{V}\left\|\nabla f\left(\boldsymbol{w}_{k}\right)\right\|_{2}^{2} + \left\|\mathbb{E}_{\zeta_{k}}\left[g\left(\boldsymbol{w}_{k};\zeta_{k}\right)\right]\right\|_{2}^{2}.$$
(12)

Also, by assumption 8b we have

$$\mathbb{E}_{\zeta_{k}} \left[\|g\left(\boldsymbol{w}_{k}; \zeta_{k}\right)\|_{2}^{2} \right] \leq M + M_{V} \|\nabla f\left(\boldsymbol{w}_{k}\right)\|_{2}^{2} + c_{0}^{2} \|\nabla f\left(\boldsymbol{w}_{k}\right)\|_{2}^{2}$$

$$= M + \left(M_{V} + c_{0}^{2}\right) \|\nabla f\left(\boldsymbol{w}_{k}\right)\|_{2}^{2}$$
(13)

By comparing Equation 13 to Equation 10 we observe that $\alpha = M$, and $\beta = (M_V + c_0^2)$.



Problem 3. For the SGD with non-convex objective function, prove that with square summable but not summable step-size, we have for any $K \in \mathbb{N}$

$$\mathbb{E}\left[\sum_{k\in[K]}\alpha_k \|\nabla f\left(\boldsymbol{w}_k\right)\|_2^2\right] < \infty, \tag{14}$$

and therefore

$$\mathbb{E}\left[\frac{1}{\sum_{k \in [K]} \alpha_k} \sum_{k \in [K]} \alpha_k \left\|\nabla f\left(\boldsymbol{w}_k\right)\right\|_2^2\right] \xrightarrow{K \to \infty} 0.$$
 (15)

Proof. Generic SG algorithm on L-smooth function satisfies

$$\mathbb{E}\left[f(\boldsymbol{w}_{k+1})\right] - f(\boldsymbol{w}_{k}) \le -(c - 0.5\alpha_{k}LM_{G})\alpha_{k}\|\nabla f(\boldsymbol{w}_{k})\|_{2}^{2} + 0.5\alpha_{k}^{2}LM. \tag{16}$$

By recursively adding them up over $k \in [K]$, we obtain

$$f_{\inf} - f(\boldsymbol{w}_1) \le \mathbb{E}\left[f(\boldsymbol{w}_{k+1})\right] - f(\boldsymbol{w}_1) \tag{17a}$$

$$\leq \sum_{k=1}^{K} \left(-\left(c - 0.5\alpha_{k} L M_{G} \right) \alpha_{k} \left\| \nabla f\left(\boldsymbol{w}_{k} \right) \right\|_{2}^{2} + 0.5\alpha_{k}^{2} L M \right)$$
(17b)

$$\leq -c \sum_{k=1}^{K} \alpha_{k} \|\nabla f(\boldsymbol{w}_{k})\|_{2}^{2} + 0.5LM_{G} \sum_{k=1}^{K} \alpha_{k}^{2} \|\nabla f(\boldsymbol{w}_{k})\|_{2}^{2} + 0.5LM \sum_{k=1}^{K} \alpha_{k}^{2}$$
(17c)

Taking expectation over the above equation and then re-arrange the terms, we have

$$\mathbb{E}\left[\sum_{k\in[K]}\alpha_{k}\left\|\nabla f\left(\boldsymbol{w}_{k}\right)\right\|_{2}^{2}\right] \leq \frac{f(\boldsymbol{w}_{1}) - f_{\inf}}{c} + 0.5LM_{G}\mathbb{E}\left[\sum_{k=1}^{K}\alpha_{k}^{2}\left\|\nabla f\left(\boldsymbol{w}_{k}\right)\right\|_{2}^{2}\right] + 0.5LM\mathbb{E}\left[\sum_{k=1}^{K}\alpha_{k}^{2}\right]$$
(18)

Note that a L-smooth function is also Lipschitz continuous, and thus we have

$$\|\nabla f\left(\boldsymbol{w}_{k}\right)\|_{2} \leq B\tag{19}$$

if $\|\nabla \boldsymbol{w}_k\|_2 \leq D$. Therefore, we have

$$\mathbb{E}\left[\sum_{k\in[K]}\alpha_k \left\|\nabla f\left(\boldsymbol{w}_k\right)\right\|_2^2\right] \leq \frac{f(\boldsymbol{w}_1) - f_{\inf}}{c} + 0.5LM_G B^2 \mathbb{E}\left[\sum_{k=1}^K \alpha_k^2\right] + 0.5LM \mathbb{E}\left[\sum_{k=1}^K \alpha_k^2\right]$$
(20)

Since the square is summable, and thus, we have the three terms on the right side bounded. Therefore, the left side is also bounded, namely

$$\|\nabla f\left(\boldsymbol{w}_{k}\right)\|_{2} \leq B \tag{21}$$

if $\|\nabla \boldsymbol{w}_k\|_2 \leq D$. Therefore, we have

$$\mathbb{E}\left[\sum_{k\in[K]}\alpha_k \left\|\nabla f\left(\boldsymbol{w}_k\right)\right\|_2^2\right] \leq \infty.$$
(22)

Meanwhile, we have

$$\mathbb{E}\left[\frac{1}{\sum_{k\in[K]}\alpha_{k}}\sum_{k\in[K]}\alpha_{k}\left\|\nabla f\left(\boldsymbol{w}_{k}\right)\right\|_{2}^{2}\right] \leq \frac{f(\boldsymbol{w}_{1}) - f_{\inf}}{c\sum_{k\in[K]}\alpha_{k}} + 0.5LM_{G}B^{2}\mathbb{E}\left[\frac{\sum_{k=1}^{K}\alpha_{k}^{2}}{\sum_{k\in[K]}\alpha_{k}}\right] + 0.5LM\mathbb{E}\left[\frac{\sum_{k=1}^{K}\alpha_{k}^{2}}{\sum_{k\in[K]}\alpha_{k}}\right]$$

$$(23)$$





Since α_k is not summable step-size, it is clear that the three terms on the right side approach zero. Therefore, we have

$$\mathbb{E}\left[\frac{1}{\sum_{k\in[K]}\alpha_k}\sum_{k\in[K]}\alpha_k \left\|\nabla f\left(\boldsymbol{w}_k\right)\right\|_2^2\right] \xrightarrow{K\to\infty} 0.$$
(24)