

EP3260: Machine Learning Over Networks

Homework 2

Stefanos Antaris^{*1}, Amaru Cuba Gyllensten^{†1,2},
Martin Isaksson^{‡1,2,3}, Sarit Khirirat^{§1}, and Klas Segeljak^{¶1,2}

¹*KTH Royal Institute of Technology*

²*RISE AI*

³*Ericsson Research*

February, 2019

Contents

1 Homework assignment	1
1.1 Human Activity Recognition Using Smartphones	1
1.2 find α and β	4
1.3 Theorem 5	4

1 Homework assignment

1.1 Human Activity Recognition Using Smartphones

Problem 1.1.1. *Consider logistic ridge regression:*

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

where

$$f_i(\mathbf{w}) = \log(1 + \exp\{-y_i \mathbf{w}^\top \mathbf{x}_i\}) \quad (2)$$

^{*}antaris@kth.se

[†]amaru.cuba.gyllensten@ri.se

[‡]martisak@kth.se

[§]sarit@kth.se

[¶]klasseg@kth.se

1. Is f Lipschitz continuous? If so, find a small B ?
2. Is f_i smooth? If so, find a small L for f_i ? What about f ?
3. Is f strongly convex? If so, find a high μ ?

Proof. **Statement 1.**

Lipschitz continuity (bounded gradients)

$$\begin{aligned}\|\mathbf{w}\|_2 \leq D &\Rightarrow \|\nabla f(\mathbf{w})\|_2 \leq B \\ \|\mathbf{w}_1\|_2, \|\mathbf{w}_2\|_2 \leq D &\Rightarrow |f(\mathbf{w}_2) - f(\mathbf{w}_1)| \leq B\|\mathbf{w}_2 - \mathbf{w}_1\|_2\end{aligned}$$

Let $\sigma(a)$ be the Sigmoid function

$$\sigma(a) := \frac{1}{1 + \exp(-a)} = \frac{\exp(a)}{\exp(a) + 1} \quad (3)$$

Since

$$\begin{aligned}\nabla f_i(\mathbf{w}) &= -y_i \mathbf{x}_i / (1 + \exp(y_i \mathbf{w}^\top \mathbf{x}_i)), \\ &= -\sigma(y_i \mathbf{w}^\top \mathbf{x}_i) y_i \mathbf{x}_i\end{aligned}$$

ok

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

By the triangle inequality,

$$\|\nabla f(w)\|_2 \leq 2\lambda \|\mathbf{w}\|_2 + \frac{1}{N} \sum_{i \in [N]} \frac{1}{|1 + e^{y_i \mathbf{w}^\top \mathbf{x}_i}|} \|y_i \mathbf{x}_i\|_2.$$

Assume that $y_i \leq C < \infty$ and $\|y_i \mathbf{x}_i\| \leq \tilde{C} < \infty$. Then,

$$\|\nabla f(w)\|_2 \leq 2\lambda D + \tilde{C}.$$

Using the dataset Human Activity Recognition Using Smartphones and finding a w that minimizes $f(\mathbf{w})$ we find that $\|f(\mathbf{w})\| \lesssim 410$ ■

Proof. **Statement 2.** The Hessian information of f_i is

$$\nabla^2 f_i(w) = \frac{\exp(y_i w^T x_i)}{(1 + \exp(y_i w^T x_i))^2} x_i x_i^T \leq \frac{1}{4} x_i x_i^T \quad \text{ok}$$

Therefore, the matrix norm of $\nabla^2 f_i(w)$ is

$$\|\nabla^2 f_i(w)\| \leq L,$$

where

$$L \leq \sigma_{\max} \left(\frac{1}{4} x_i x_i^T \right). \quad \text{It can be proved that } \mathbf{xx}' \leq (\text{norm}(\mathbf{x})^2) \cdot \mathbf{I}$$

Here, $\sigma_{\max}(A)$ is the largest eigenvalue of a positive semidefinite matrix A .

Next, we can easily compute the Hessian information of f as follows:

$$\nabla^2 f(w) = \frac{1}{N} \sum_{i=1}^N \frac{\exp(y_i w^T x_i)}{(1 + \exp(y_i w^T x_i))^2} x_i x_i^T + 2\lambda I \leq \frac{1}{4N} \sum_{i=1}^N x_i x_i^T + 2\lambda I.$$

From the definition of the matrix norm, and by the triangle inequality,

$$\|\nabla^2 f(w)\| \leq L$$

where

$$L \leq \sigma_{\max} \left(\frac{1}{4N} \sum_{i=1}^N x_i x_i^T + 2\lambda I \right). \quad \text{Similarly, we know that } f \text{ is } \{\sum(\mathbf{x}'\mathbf{x})/4N + 2\lambda\} \text{-smooth}$$

■

Proof. **Statement 3.**

Define any $v \in \mathbb{R}^d$. Then, it is obvious that

$$v^T x_i x_i^T v = (x_i^T v)^T (x_i^T v) = \|x_i^T v\|^2 \geq 0,$$

for each $x_i \in \mathbb{R}^d$ and $i \in [N]$. This implies that $x_i x_i^T$ is a positive semi-definite matrix. Therefore,

$$\frac{1}{N} \sum_{i=1}^N v^T x_i x_i^T v = \frac{1}{N} \sum_{i=1}^N \|x_i^T v\|^2 \geq 0.$$

Since $\exp(-x)/(1 + \exp(-x))$ is always a positive number, and $(1/N) \sum_{i=1}^N x_i x_i^T$ is the positive-semidefinite matrix, the Hessian information on the form

$$\nabla^2 f(w) = \frac{1}{N} \sum_{i=1}^N \frac{\exp(y_i w^T x_i)}{(1 + \exp(y_i w^T x_i))^2} x_i x_i^T + 2\lambda I$$

implies

$$\begin{aligned} v^T \nabla^2 f(w) v &= \frac{1}{N} \sum_{i=1}^N \frac{\exp(y_i w^T x_i)}{(1 + \exp(y_i w^T x_i))^2} v^T x_i x_i^T v + 2\lambda v^T v \\ &\geq v^T (2\lambda \cdot I) v. \end{aligned}$$

ok

Thus, $\nabla^2 f(w)$ has the smallest eigenvalue with 2λ , i.e. $\nabla^2 f(w) \geq \mu I$ with $\mu = 2\lambda$. This means that f is strongly convex with $\mu = 2\lambda$. ■

1.2 find α and β

Problem 1.2.1.

$$\mathbb{E}_{\zeta_k} \left[\|g(\mathbf{w}_k; \zeta_k)\|_2^2 \right] \leq \alpha + \beta \|\nabla f(\mathbf{w}_k)\|_2^2 \quad (4)$$

Proof. Notice that $\mathbb{E}_{\zeta_k} g(\mathbf{w}_k; \zeta_k) = \nabla f(w_k)$ and from the definition of the Euclidean norm,

$$\begin{aligned} \mathbb{E}_{\zeta_k} \left[\|g(\mathbf{w}_k; \zeta_k)\|_2^2 \right] &= \mathbb{E}_{\zeta_k} \left[\|g(\mathbf{w}_k; \zeta_k) - \nabla f(w_k) + \nabla f(w_k)\|_2^2 \right] \\ &= \mathbb{E}_{\zeta_k} \left[\|g(\mathbf{w}_k; \zeta_k) - \nabla f(w_k)\|^2 - \|\nabla f(w_k)\|^2 \right. \\ &\quad \left. + 2 \mathbb{E}_{\zeta_k} \langle g(\mathbf{w}_k; \zeta_k), \nabla f(w_k) \rangle \right] \\ &\leq \mathbb{E}_{\zeta_k} \left[\|g(\mathbf{w}_k; \zeta_k) - \nabla f(w_k)\|^2 - \|\nabla f(w_k)\|^2 \right. \\ &\quad \left. + 2 \left\| \mathbb{E}_{\zeta_k} g(\mathbf{w}_k; \zeta_k) \right\| \|\nabla f(w_k)\| \right]. \end{aligned}$$

Since

$$\begin{aligned} \mathbb{E}_{\zeta_k} \left[\|g(\mathbf{w}_k; \zeta_k) - \nabla f(w_k)\|^2 - \|\nabla f(w_k)\|^2 \right] &= \text{Var}_{\zeta_k} g(w_k; \zeta_k) \leq M + M_V \|\nabla f(w_k)\|^2 \\ \mathbb{E}_{\zeta_k} \|g(\mathbf{w}_k; \zeta_k)\| &\leq c_0 \|\nabla f(w_k)\| \end{aligned}$$

we have:

$$\mathbb{E}_{\zeta_k} \left[\|g(\mathbf{w}_k; \zeta_k)\|_2^2 \right] \leq M + (M_V - 1 + 2c_0) \|\nabla f(w_k)\|^2. \quad \text{ok}$$

Since $c_0 = 1$, we have: $\alpha = M$ and $\beta = M_V + 1$. ■

1.3 Theorem 5

Theorem 1.3.1. *With square summable but not summable step-size, we have for any $K \in \mathbb{N}$*

$$\mathbb{E} \left[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 \right] < \infty \quad (5)$$

and therefore

$$\lim_{K \rightarrow \infty} \mathbb{E} \left[\frac{1}{\sum_{k \in [K]} \alpha_k} \sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 \right] = 0 \quad (6)$$

The expected gradient norm cannot stay bounded away from zero.

Proof. Since $\sum_{k \in [K]} \alpha_k^2 < \infty$ and $\sum_{k \in [K]} \alpha_k = \infty$, we have:

$$\mathbb{E} \left[\frac{1}{\sum_{k \in [K]} \alpha_k} \sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 \right] = \frac{1}{\sum_{k \in [K]} \alpha_k} \mathbb{E} \left[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 \right].$$

Since $\sum_{k \in [K]} \alpha_k = \infty$ and $\mathbb{E} \left[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 \right] = C < \infty$ for a finite constant C ,

Why? Prove that

$$\begin{aligned} \lim_{K \rightarrow \infty} \mathbb{E} \left[\frac{1}{\sum_{k \in [K]} \alpha_k} \sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 \right] &= \lim_{K \rightarrow \infty} \frac{1}{\sum_{k \in [K]} \alpha_k} \mathbb{E} \left[\sum_{k \in [K]} \alpha_k \|\nabla f(\mathbf{w}_k)\|_2^2 \right] \\ &= C \lim_{K \rightarrow \infty} \frac{1}{\sum_{k \in [K]} \alpha_k} = 0. \end{aligned}$$

■