

# EECS 553 HW4

Lingqi Huang

September 2024

## 1 Problem 1

**Part(a):** Notice that the loss function  $J(w, b)$  can be written as

$$J(w, b) = \sum_{i=1}^n \left( \frac{1}{n} (L(y_i, w^T x_i + b) + \frac{\lambda}{2} \|w\|^2) \right)$$

Thus, we observe that

$$J_i(w, b) = \frac{1}{n} \left( L(y_i, w^T x_i + b) + \frac{\lambda}{2} \|w\|^2 \right) = \frac{1}{n} \left( \max\{0, 1 - y_i(w^T x_i + b)\} + \frac{\lambda}{2} \|w\|^2 \right)$$

Now we set  $\tilde{x}_i = [1, x_{i1}, \dots, x_{id}]^T$ ,  $\theta = [b, w^T]^T$ , we conclude that our  $J_i$  can be re-written as

$$J_i(\theta) = \frac{1}{n} \left( \max\{0, 1 - y_i \theta^T \tilde{x}_i\} + \frac{\lambda}{2} \|w\|^2 \right)$$

Now for convenience, we only discuss two cases, which are  $y_i \theta^T \tilde{x}_i < 1$ , and  $y_i \theta^T \tilde{x}_i \geq 1$ . The last case make senses because the hinge loss function is non-differentiable as point 0, and so we let the subgradient to be 0 which is the right derivative of the hinge loss if the hinge loss is 0, ie,  $1 - y_i \theta^T \tilde{x}_i = 0$ . Thus, we can easily show that

$$u_i = \nabla J_i(\theta) = \begin{cases} \frac{1}{n} \begin{pmatrix} -y_i \tilde{x}_i + \lambda \begin{bmatrix} 0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \end{pmatrix} & \text{if } y_i \theta^T \tilde{x}_i < 1 \\ \frac{\lambda}{n} \begin{bmatrix} 0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} & \text{if } y_i \theta^T \tilde{x}_i \geq 1 \end{cases}$$

by the fact that hinge loss is convex and chain rule.

**Part(b):** By the code in the attachment, we find that the estimated parameters are  $w_1 = -17.816, w_2 = -9.117, b = 12.06$ , the margin is  $\frac{1}{\|w\|} = 0.04997$ , and the minimum achieved value of the objective function is 0.4498. You can check the diagrams in the next page.

**Part(c):** By the code in the attachment, we find that the estimated parameters are  $w_1 = -5.82, w_2 = -4.41, b = 4.005$ , the margin is 0.13683, and the minimum achieved value of the objective function is 0.25827. We can see that the applying SGD algorithm could converges much faster than using subgradient method.

### Pictures in 1(b)

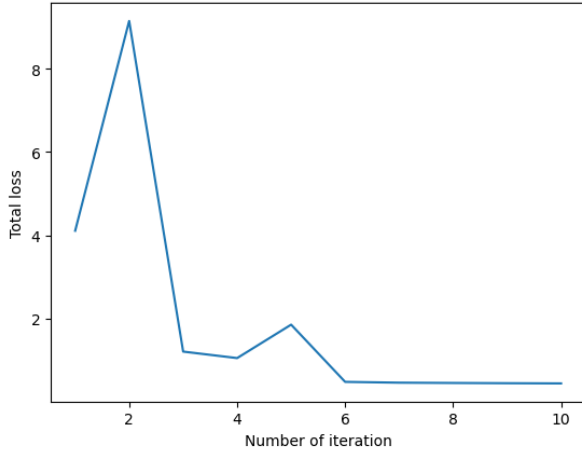


Figure 1: The loss function value VS Number of iteration using subgradient method

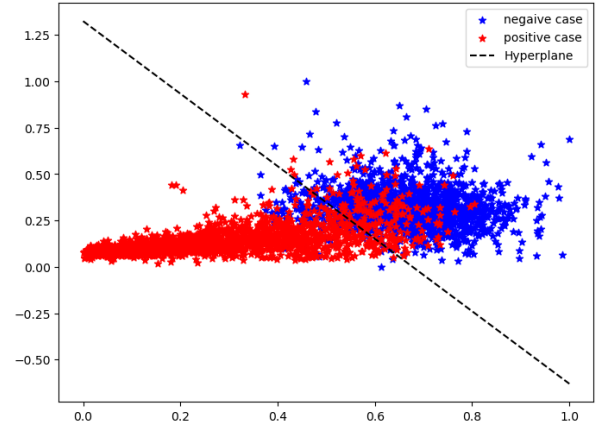


Figure 2: Visualization of data and the learned line using subgradient method

### Pictures in 1(c)

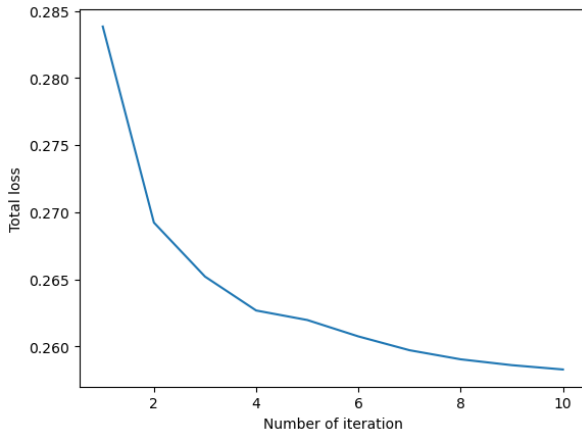


Figure 3: The loss function value VS Number of iteration using stochastic gradient descent method

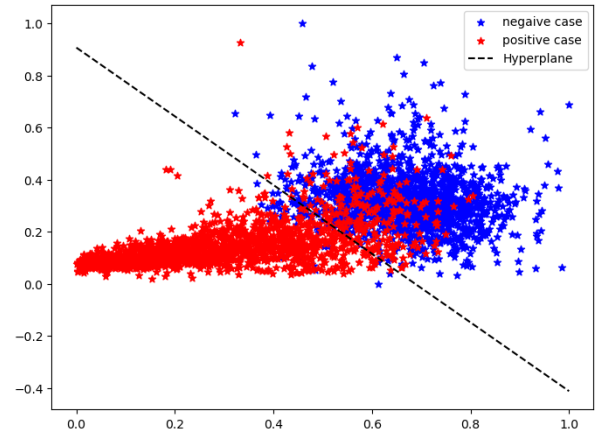


Figure 4: Visualization of data and the learned line using stochastic gradient descent method

## 2 Problem 2

**Part(a):** Notice that we can write  $g(w_j)$  at iteration  $t$  as

$$g(w_j) = \frac{1}{n} \sum_{i=1}^n (y_i - \sum_{k=1}^d w_k^{(t)} x_{ik} - b^{(t)})^2 + \lambda \|w\|_1$$

Now we define function  $h(x)$  such that

$$h(x) = \begin{cases} -1 & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Thus, we find that the sub-differential could be written as

$$\partial g(w_j) = \frac{2}{n} \sum_{i=1}^n (-y_i x_{ij} + x_{ij} \sum_{k=1}^d w_k^{(t)} x_{ik} + b^{(t)} x_{ij}) + \lambda h(w_j)$$

Notice that we can observe that

$$\frac{2}{n} \sum_{i=1}^n x_{ij} \left( \sum_{k=1}^d w_k^{(t)} x_{ik} \right) = \left( \frac{2}{n} \sum_{i=1}^n x_{ij}^2 \right) w_j + \frac{2}{n} \sum_{i=1}^n x_{ij} \left( \sum_{\substack{k=1 \\ k \neq j}}^d w_k^{(t)} x_{ik} \right)$$

where the second term could be written as the

$$\frac{2}{n} \sum_{i=1}^n x_{ij} \left( \sum_{\substack{k=1 \\ k \neq j}}^d w_k^{(t)} x_{ik} \right) = \frac{2}{n} \sum_{i=1}^n x_{ij} (w_{-j}^{(t)})^T x_{i,-j}$$

where  $w_{-j}^{(t)}, x_{i,-j}$  are the notation defined in the problem. Thus, we get that

$$\partial g(w_j) = \left( \frac{2}{n} \sum_{i=1}^n x_{ij}^2 \right) w_j - \frac{2}{n} \sum_{i=1}^n x_{ij} (y_i - (w_{-j}^{(t)})^T x_{i,-j} - b^{(t)}) + \lambda h(w_j)$$

Now if we define

$$a_j^{(t)} = \frac{2}{n} \sum_{i=1}^n x_{ij}^2, \quad c_j^{(t)} = \frac{2}{n} \sum_{i=1}^n x_{ij} (y_i - (w_{-j}^{(t)})^T x_{i,-j} - b^{(t)})$$

and combine the definition of  $h(x)$ , we finally conclude that

$$\partial g(w_j) = \begin{cases} a_j^{(t)} w_j - c_j^{(t)} - \lambda, & w_j < 0 \\ [a_j^{(t)} w_j - c_j^{(t)} - \lambda, a_j^{(t)} w_j - c_j^{(t)} + \lambda], & w_j = 0 \\ a_j^{(t)} w_j - c_j^{(t)} + \lambda, & w_j > 0 \end{cases}$$

This completes the proof.

**Part(b):** Let's first consider the case of  $c_j^{(t)} < -\lambda$ , which means  $w_j = \frac{c_j^{(t)} + \lambda}{a_j^{(t)}} < 0$ , and this is the unique  $w_j$  such that  $0 \in \partial g(w_j)$ . when  $c_j^{(t)} \in [-\lambda, \lambda]$ , we automatically have  $w_j = 0$  that only the second condition can be satisfied. For the similar reason, when  $c_j^{(t)} > \lambda$ , we have the unique  $w_j$  that satisfies the third condition, where  $w_j = \frac{c_j^{(t)} - \lambda}{a_j^{(t)}} > 0$ . Thus, we conclude that the optimal value of  $w_j$  is given by the soft-thresholding formula, that

$$w_j = \text{soft} \left( \frac{c_j^{(t)}}{a_j^{(t)}}, \frac{\lambda}{a_j^{(t)}} \right)$$

where

$$\text{soft}(\alpha, \beta) = \begin{cases} \alpha - \beta, & \alpha > \beta \\ 0, & \alpha \in [-\beta, \beta] \\ \alpha + \beta, & \alpha < -\beta \end{cases}$$

This completes the proof.

### 3 Problem 3

**Part(a):** We have verified that rows of thus sphered  $\mathbb{X}$  training data matrix have zero sample mean and unit variance. You can see the picture below as well as code in the end of this pdf.

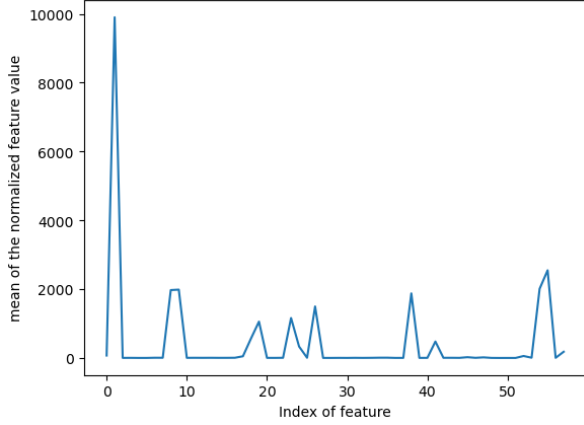


Figure 5: feature index VS Mean of feature

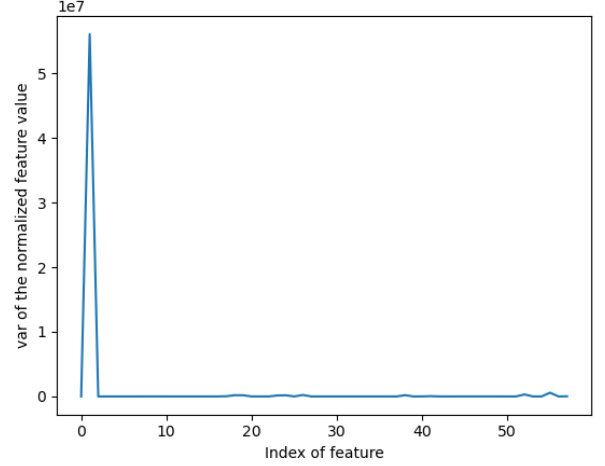


Figure 6: feature index VS Var of feature

**Part(b):** Finally we find that the Mean squared error is about 754.79, and there are 4 entries are 0 inside of the final parameter  $\omega$ . You can see the code at the end of the pdf.

### 4 Problem 4

**Part(a):** We denote  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ , and so we notice that

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^3 = \left( \sum_{i=1}^d u_i v_i \right)^3 = \sum_{k_1+k_2+\dots+k_d=3} \frac{3!}{k_1!k_2!\dots k_d!} \prod_{i=1}^d (u_i v_i)^{k_i}$$

Notice that the last term could be written as an inner product, where the vector of  $\Phi(\mathbf{u}), \Phi(\mathbf{v})$  has dimension of  $\binom{d+2}{d-1}$  by 1. Now we can find a way to order the element of  $\Phi(\mathbf{u})$  (same for  $\Phi(\mathbf{v})$ ), that is we first consider all cases of  $k_1 = 3, k_2 = 3, \dots, k_d = 3$ , and then all cases of  $k_2 = 2, k_2 = 2, \dots, k_d = 2$ . Thus, we can find a order function  $\phi(k_1, k_2, \dots, k_d)$  ( $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{\binom{d+2}{d-1}}$ ) mapping to the position of the inner product vector  $\Phi(\mathbf{u})$  where  $\forall 1 \leq i \leq d$ ,  $0 \leq k_i \leq 3$ , and  $\sum_{i=1}^d k_i = 3$ . Thus, given the order function  $\phi$ , we can  $k(\mathbf{u}, \mathbf{v})$  as

$$k(\mathbf{u}, \mathbf{v}) = \left\langle \begin{bmatrix} \vdots \\ \sqrt{\frac{3!}{k_1!k_2!\dots k_d!}} \prod_{i=1}^d u_i^{k_i} \\ \vdots \end{bmatrix}_{\phi}, \begin{bmatrix} \vdots \\ \sqrt{\frac{3!}{k_1!k_2!\dots k_d!}} \prod_{i=1}^d v_i^{k_i} \\ \vdots \end{bmatrix}_{\phi} \right\rangle$$

This completes the proof.

**Part(b):** Suppose that  $k$  is an inner product kernel, then  $\exists \Phi, \forall \mathbf{x}_i, \mathbf{x}_j, k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$  by definition of inner product, so  $k$  is symmetric. Now notice that the matrix of  $k$  with dimension of  $n \times n$  can be rewritten as

$$k = \begin{bmatrix} \Phi(\mathbf{x}_1)^T \\ \Phi(\mathbf{x}_2)^T \\ \vdots \\ \Phi(\mathbf{x}_n)^T \end{bmatrix} \cdot [\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_n)]$$

Now denote

$$\varphi = \begin{bmatrix} \Phi(\mathbf{x}_1)^T \\ \Phi(\mathbf{x}_2)^T \\ \vdots \\ \Phi(\mathbf{x}_n)^T \end{bmatrix}$$

Then for any vector  $z \in \mathbb{R}^{n \times 1}$ , we must have

$$z^T \cdot k \cdot z = z^T \varphi \varphi^T z = \langle \varphi^T z, \varphi^T z \rangle \geq 0$$

Then by definition of semi-positive matrix, we conclude that  $k$  is a symmetric, positive definite kernel. This completes the proof.

**Part(c):** By lecture notes, we know that

$$\hat{b} = \bar{y} - \hat{w}^T \bar{\mathbf{x}}$$

where

$$\hat{w}^T = \tilde{y}^T (\tilde{G} + n\lambda I)^{-1} \tilde{\mathbf{X}}$$

$$\tilde{y} = \begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{bmatrix}$$

$$\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} \\ \mathbf{x}_2 - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}_d - \bar{\mathbf{x}} \end{bmatrix}$$

And  $\hat{G}$  is the matrix that presented in the lecture slide, which is a SPD kernel. Thus, we can rewrite  $\hat{b}$  as

$$\hat{b} = \bar{y} - \hat{y}^T (\hat{G} + n\lambda I)^{-1} \tilde{\mathbf{X}} \bar{\mathbf{x}}$$

Now we can write  $\tilde{\mathbf{X}} \bar{\mathbf{x}}$  as  $\tilde{g}(\tilde{\mathbf{x}})$  where

$$\tilde{g}(\tilde{\mathbf{x}}) = \tilde{\mathbf{X}} \bar{\mathbf{x}} = \begin{bmatrix} \langle \tilde{\mathbf{x}}_1, \bar{\mathbf{x}} \rangle \\ \vdots \\ \langle \tilde{\mathbf{x}}_n, \bar{\mathbf{x}} \rangle \end{bmatrix}$$

And we have that

$$\langle \tilde{\mathbf{x}}_i, \bar{\mathbf{x}} \rangle = \frac{1}{n} \sum_{j=1}^n \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

Thus, we have shown that the offset  $b$  can also be evaluated using the kernel.