

## Theoretical Assignment

- 1) **Step-size Perceptron.** Consider the modification of Perceptron algorithm with the following update rule:

$$w_{t+1} \leftarrow w_t + \eta_t y_t x_t$$

whenever  $\hat{y}_t \neq y_t$  ( $w_{t+1} \leftarrow w_t$  otherwise). Assume that data is separable with margin  $\gamma > 0$  and that  $\|x_t\| = 1$  for all  $t$ . For simplicity assume that the algorithm makes  $M$  mistakes at the first  $M$  rounds, after which it has no mistakes. For  $\eta_t = \frac{1}{\sqrt{t}}$ , show that the number of mistakes step-size Perceptron makes is at most  $\frac{4}{\gamma^2} \log\left(\frac{1}{\gamma}\right)$ . (Hint: use the fact that if  $x \leq a \log(x)$  then  $x \leq 2a \log(a)$ ). It's okay if you obtain a bound with slightly different constants, but the asymptotic dependence on  $\gamma$  should be tight.

**Solution:**

Assume  $w^*$  is a linear separator of the data and w.l.o.g  $\|w^*\| = 1$ .

(\*) Moreover, assume  $\gamma > 0$  is the margin of  $w^*$ , i.e., for every  $t$  it holds that  $y_t w^* \cdot x_t = \frac{|w^* \cdot x_t|}{\|w^*\|} = \text{dist}(x_t, \text{hyperplane}(w^*)) \geq \gamma$

After each mistake at time  $t$  it holds that:

$$\begin{aligned} w_{t+1} \cdot w^* &= (w_t + \eta_t y_t x_t) \cdot w^* = w_t \cdot w^* + \eta_t y_t x_t \cdot w^* \stackrel{(*)}{\geq} w_t \cdot w^* + \eta_t \gamma \\ &= w_t \cdot w^* + \frac{1}{\sqrt{t}} \gamma \end{aligned}$$

Therefore, after  $M$  mistakes (that made at first) we have:

$$\begin{aligned} (i) \quad w_{t+1} \cdot w^* &\geq w_M \cdot w^* + \frac{1}{\sqrt{M}} \gamma \geq \dots \geq \gamma \sum_{i=1}^M \frac{1}{\sqrt{i}} \geq \gamma \int_1^{M+1} \frac{dx}{\sqrt{x}} \\ &= \gamma(2\sqrt{M+1} - 2) \stackrel{(**)}{\geq} \gamma \sqrt{eM} \end{aligned}$$

Where (\*\*) holds for  $M$  sufficiently large:

$$\lim_{M \rightarrow \infty} \frac{\sqrt{eM}}{2\sqrt{M+1} - 2} = \frac{\sqrt{e}}{2} < 1$$

In addition, after each mistake at time  $t$  it holds that:

$$\begin{aligned} \|w_{t+1}\|^2 &= \|w_t + \eta_t y_t x_t\|^2 = \|w_t\|^2 + \underbrace{2\eta_t y_t w_t \cdot x_t}_{\substack{\text{negative term since} \\ \text{we did a mistake}}} + \|\eta_t y_t x_t\|^2 \\ &\leq \|w_t\|^2 + \eta_t^2 \|x_t\|^2 =_{\|x_t\|=1} \|w_t\|^2 + \frac{1}{t} \end{aligned}$$

Therefore, after  $M$  mistakes (that made at first) we have:

$$(ii) \quad \|w_{t+1}\|^2 \leq \|w_t\|^2 + \frac{1}{t} \leq \dots \leq \sum_{i=1}^M \frac{1}{i} \leq 1 + \int_1^M \frac{dx}{x} = 1 + \ln(M) = \ln(eM)$$

By putting all together, we get:

$$\gamma \sqrt{eM} \stackrel{(i)}{\leq} w_{t+1} \cdot w^* \leq_{C-S} \|w_{t+1}\| \|w^*\| = \|w_t\| \stackrel{(ii)}{\leq} \sqrt{\ln(eM)}$$

Hence,

$$0 < \gamma \sqrt{eM} \leq \sqrt{\ln(eM)} \Rightarrow \gamma^2 eM \leq \ln(eM) \Rightarrow eM \leq \frac{1}{\gamma^2} \ln(eM)$$

By using the hint that mentioned in the question:

$$M \leq eM \leq 2 \frac{1}{\gamma^2} \ln\left(\frac{1}{\gamma^2}\right) = \frac{4}{\gamma^2} \ln\left(\frac{1}{\gamma}\right) \blacksquare$$

2) Convex functions.

- (a) Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  a convex function,  $A \in \mathbb{R}^{n \times n}$  and  $b \in \mathbb{R}^n$ . Show that,  $g(x) = f(Ax + b)$  is convex.

**Solution:**

Recall: For convex set  $C$  we say that  $f: C \rightarrow \mathbb{R}$  is convex if  $\forall \alpha \in [0,1], \forall x_1, x_2 \in C : f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2)$ .

First, notice that  $\mathbb{R}^n$  is convex set. Let  $x_1, x_2 \in \mathbb{R}^n$  and let  $\alpha \in [0,1]$ . It holds that:

$$\begin{aligned} g(\alpha x_1 + (1 - \alpha)x_2) &= f(A(\alpha x_1 + (1 - \alpha)x_2) + b) \\ &= f(\alpha(Ax_1 + b) + (1 - \alpha)(Ax_2 + b)) \leq_{\text{convexity of } f} \\ &\leq \alpha f(Ax_1 + b) + (1 - \alpha)f(Ax_2 + b) = \alpha g(x_1) + (1 - \alpha)g(x_2) \blacksquare \end{aligned}$$

- (b) Consider  $m$  convex functions  $f_1(x), \dots, f_m(x)$ , where  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ . Now define a new function  $g(x) = \max_i f_i(x)$ . Prove that  $g(x)$  is a convex function. (Note that from (a) and (b) you can conclude that the hinge loss over linear classifiers is convex.)

**Solution:**

Let  $x_1, x_2 \in \mathbb{R}^n$  and let  $\alpha \in [0,1]$ . It holds that:

$$\begin{aligned} g(\alpha x_1 + (1 - \alpha)x_2) &= \max_i f_i(\alpha x_1 + (1 - \alpha)x_2) \leq_{\text{convexity of each } f_i} \\ &\leq \max_i \alpha f_i(x_1) + (1 - \alpha)f_i(x_2) \leq \max_i \alpha f_i(x_1) + \max_i (1 - \alpha)f_i(x_2) \\ &= \alpha g(x_1) + (1 - \alpha)g(x_2) \blacksquare \end{aligned}$$

- (c) Let  $\ell_{\log}: \mathbb{R} \rightarrow \mathbb{R}$  be the log loss, defined by  $\ell_{\log}(z) = \log_2(1 + e^{-z})$ . Show that  $\ell_{\log}$  is convex, and conclude that the function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  defined by  $f(w) = \ell_{\log}(y\mathbf{w} \cdot \mathbf{x})$  is convex with respect to  $\mathbf{w}$ .

**Solution:**

Recall (Calculus 1): Let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be a twice differentiable function. Then,  $f$  is convex if and only if  $f''(x) \geq 0$ .

$$\begin{aligned} \ell_{\log} \text{ is differentiable twice and } \ell'_{\log}(z) &= \ln(2) \frac{-e^{-z}}{1+e^{-z}} \\ \Rightarrow \ell''_{\log} &= \ln(2) \frac{e^{-z}(1+e^{-z}) - e^{-z} \cdot e^{-z}}{(1+e^{-z})^2} = \ln(2) \frac{e^{-z}(1+e^{-z} - e^{-z})}{(1+e^{-z})^2} \\ &= \ln(2) \frac{e^{-z}}{(1+e^{-z})^2} \geq 0 \end{aligned}$$

Therefore,  $\ell_{\log}$  is convex function on  $\mathbb{R}$ . In addition, for given  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}$  we will show that  $f_{x,y}: \mathbb{R}^d \rightarrow \mathbb{R}$  which is defined by  $f(w) = \ell_{\log}(y\mathbf{w} \cdot \mathbf{x})$  is convex wrt  $w$ . Notice that  $y\mathbf{w} \cdot \mathbf{x} = \mathbf{w} \cdot (y\mathbf{x}) = A\mathbf{w}$  where  $A = (y\mathbf{x})^t$ . Therefore, by section (a)  $\ell_{\log}(A\mathbf{w}) = f(w)$  is convex function.  $\blacksquare$

3) Ranking.

(a) Consider the hinge loss function for the ranking objective problem:

$$\ell(h_{\mathbf{w}}(\bar{\mathbf{x}}), y) = \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, 1 - \text{sgn}(y_j - y_r) \mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)\}$$

We shall show that  $\ell$  is convex wrt  $\mathbf{w}$ . Let  $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$  and let  $\alpha \in [0, 1]$ . It holds that,

$$\begin{aligned} \ell(h_{\alpha \mathbf{w}_1 + (1-\alpha) \mathbf{w}_2}(\bar{\mathbf{x}}), y) &= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, 1 - \text{sgn}(y_j - y_r) (\alpha \mathbf{w}_1 + (1-\alpha) \mathbf{w}_2) \cdot (\mathbf{x}_j - \mathbf{x}_r)\} \\ &= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, \alpha (1 - \text{sgn}(y_j - y_r) \mathbf{w}_1 \cdot (\mathbf{x}_j - \mathbf{x}_r)) \\ &\quad + (1-\alpha) (1 - \text{sgn}(y_j - y_r) \mathbf{w}_2 \cdot (\mathbf{x}_j - \mathbf{x}_r))\} \\ &\leq_{(*)} \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, \alpha (1 - \text{sgn}(y_j - y_r) \mathbf{w}_1 \cdot (\mathbf{x}_j - \mathbf{x}_r))\} \\ &\quad + \max\{0, (1-\alpha) (1 - \text{sgn}(y_j - y_r) \mathbf{w}_2 \cdot (\mathbf{x}_j - \mathbf{x}_r))\} \\ &= \alpha \ell(h_{\mathbf{w}_1}(\bar{\mathbf{x}}), y) + (1-\alpha) \ell(h_{\mathbf{w}_2}(\bar{\mathbf{x}}), y) \end{aligned}$$

Where (\*) is true  $\alpha$  and  $1-\alpha$  are non-negative numbers and we sum all potential non-negative terms in the max function. Therefore,  $\ell$  is convex wrt  $\mathbf{w}$ . ■

(b) Consider the Kendall-Tau loss function for the ranking objective problem:

$$\Delta(y, y') = \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbf{1}\{\text{sgn}(y'_j - y'_r) \neq \text{sgn}(y_j - y_r)\}$$

We shall show that the hinge loss upper-bounds the Kendall-Tau loss for linear classifiers:  $h_{\mathbf{w}}((x_1, \dots, x_k)) = (\mathbf{w} \cdot x_1, \dots, \mathbf{w} \cdot x_k)$ .

Let  $\mathbf{w} \in \mathbb{R}^d, \bar{\mathbf{x}} \in \mathcal{X}^k, \mathbf{y} \in \mathbb{R}^k$ . It holds that,

$$\begin{aligned} \Delta(h_{\mathbf{w}}(\bar{\mathbf{x}}), \mathbf{y}) &= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbf{1}\{\text{sgn}(h_{\mathbf{w}}(\bar{\mathbf{x}})_j - h_{\mathbf{w}}(\bar{\mathbf{x}})_r) \neq \text{sgn}(y_j - y_r)\} \\ &= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbf{1}\{\text{sgn}(\mathbf{w} \cdot \mathbf{x}_j - \mathbf{w} \cdot \mathbf{x}_r) \neq \text{sgn}(y_j - y_r)\} \\ &= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbf{1}\{\text{sgn}(\mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)) \neq \text{sgn}(y_j - y_r)\} \\ &\leq_{(*)} \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbf{1}\{\text{sgn}(\mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)) \text{sgn}(y_j - y_r) \in \{0, -1\}\} \\ &\leq_{(**)} \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, 1 - \text{sgn}(y_j - y_r) \mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)\} = \ell(h_{\mathbf{w}}(\bar{\mathbf{x}}), y) \end{aligned}$$

Where (\*) is true since we took also cases which both sides are equal to zero. Part

(\*\*) is true since if  $\text{sgn}(\mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)) \text{sgn}(y_j - y_r) = -1$  then  $-\text{sgn}(y_j - y_r) \mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r) \geq 0$  and therefore  $\max\{0, 1 - \text{sgn}(y_j - y_r) \mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)\} \geq 1$ . In the

other cases, we don't decrease the sum by taking 0 at the max function. ■

- (c) Assume the data is linearly separable with margin  $\gamma > 0$ . We shall show now that minimizing the hinge loss will result a ranking function which minimizes the Kendall-Tau loss. Since the hinge loss upper-bounds the Kendall-Tau loss for linear classifiers, it suffices to show that the hinge loss can get arbitrarily close to Kendall-Tau loss. By the separability assumption, denote  $\mathbf{w}^* \in \mathbb{R}^d$  the true classifier such that  $\text{sgn}(y_j^i - y_r^i) \mathbf{w}^* \cdot (\mathbf{x}_j^i - \mathbf{x}_r^i) \geq \gamma$  for all  $1 \leq i \leq n$  and  $1 \leq j < r \leq k$  (\*).

Note that for every  $c > 0$ ,  $c\mathbf{w}^*$  yields the same classifier as  $\mathbf{w}^*$ . Therefore:

$$\begin{aligned} \sum_{i=1}^n \Delta(h_{c\mathbf{w}^*}(\bar{\mathbf{x}}^i), \mathbf{y}^i) &\leq_{U-B} \sum_{i=1}^n \ell(h_{c\mathbf{w}^*}(\bar{\mathbf{x}}^i), \mathbf{y}^i) \\ &= \sum_{i=1}^n \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, 1 - \text{sgn}(y_j^i - y_r^i) c\mathbf{w}^* \cdot (\mathbf{x}_j^i - \mathbf{x}_r^i)\} \\ &\leq_{(*)} \sum_{i=1}^n \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, 1 - c\gamma\} = n \cdot \max\{0, 1 - c\gamma\} \end{aligned}$$

Taking  $c \rightarrow \frac{1}{\gamma}$  and we get that  $\sum_{i=1}^n \Delta(h_{c\mathbf{w}^*}(\bar{\mathbf{x}}^i), \mathbf{y}^i) \leq n \cdot \max\{0, 1 - c\gamma\} \rightarrow 0$

Hence, we minimize the K-T loss. ■

#### 4) Gradient Decent on Smooth Functions.

$f$  is  $\beta$ -smooth function, i.e.,  $\forall y, x \in \mathbb{R}^n$ :

$$f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{\beta}{2} \|x - y\|^2$$

Substitute  $y = x_{t+1}$  and  $x = x_t$ , and the iteration rule  $x_t - x_{t+1} = \eta \nabla f(x_t)$ . Then, for every  $t \geq 0$ :

$$\begin{aligned} 0 &\leq \eta \|\nabla f(x_t)\|^2 = \nabla f(x_t) \cdot (\eta \nabla f(x_t)) = \nabla f(x_t) \cdot (x_t - x_{t+1}) \\ &\leq f(x_t) - f(x_{t+1}) + \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \end{aligned}$$

Consider the infinite sum of both non-negative sides, and we get:

$$\begin{aligned} \sum_{t=0}^T \|\nabla f(x_t)\|^2 &\leq \frac{1}{\eta} \sum_{t=0}^T f(x_t) - f(x_{t+1}) + \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \\ &= \frac{1}{\eta} (f(x_0) - f(x_{T+1})) + \frac{\beta}{2\eta} \sum_{t=0}^{\infty} \|\eta \nabla f(x_t)\|^2 \end{aligned}$$

So,

$$\left(1 - \frac{\beta\eta}{2}\right) \sum_{t=0}^T \|\nabla f(x_t)\|^2 \leq \frac{1}{\eta} (f(x_0) - f(x_{T+1}))$$

Since  $\eta < \frac{2}{\beta}$ ,  $1 - \frac{\beta\eta}{2} > 0$ . Recall that  $f$  is non-negative function so  $f(x_{T+1}) \geq 0$ .

Therefore,

$$\sum_{t=0}^T \|\nabla f(x_t)\|^2 \leq \frac{1}{1 - \frac{\beta\eta}{2}} \frac{1}{\eta} (f(x_0) - f(x_{T+1})) \leq \frac{1}{1 - \frac{\beta\eta}{2}} \frac{1}{\eta} f(x_0)$$

Taking  $T \rightarrow \infty$ :

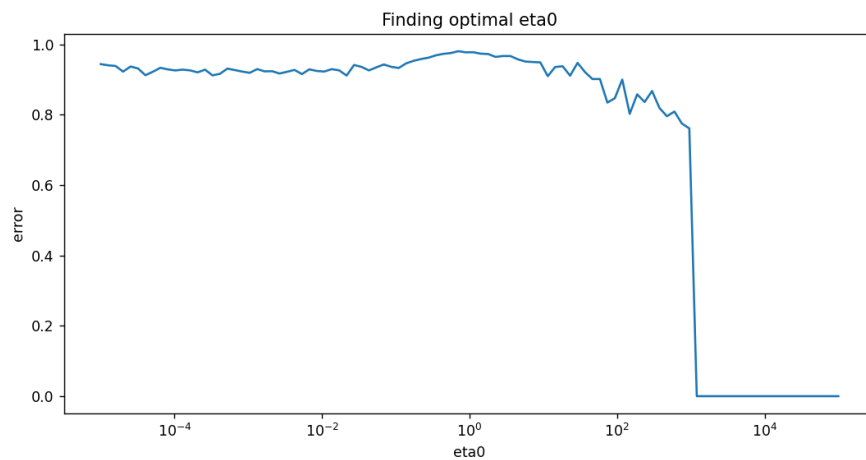
$$\sum_{t=0}^{\infty} \|\nabla f(x_t)\|^2 \leq \frac{1}{1 - \frac{\beta\eta}{2}} \frac{1}{\eta} f(x_0) < \infty$$

Therefore,  $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\|^2 = 0$  and also  $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$ . ■

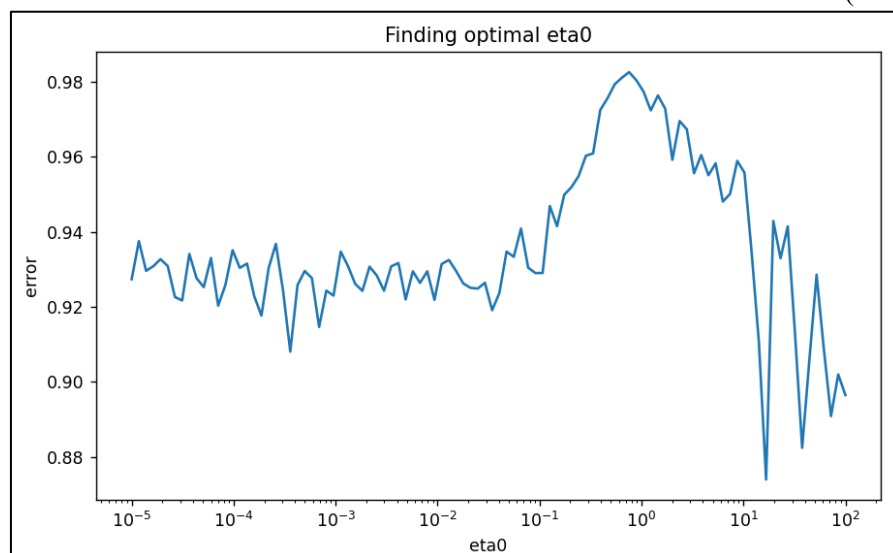
## Programming Assignment

### (1) SGD for Hinge loss

(a) The average accuracy on the validation set as a function of  $\eta_0$ :

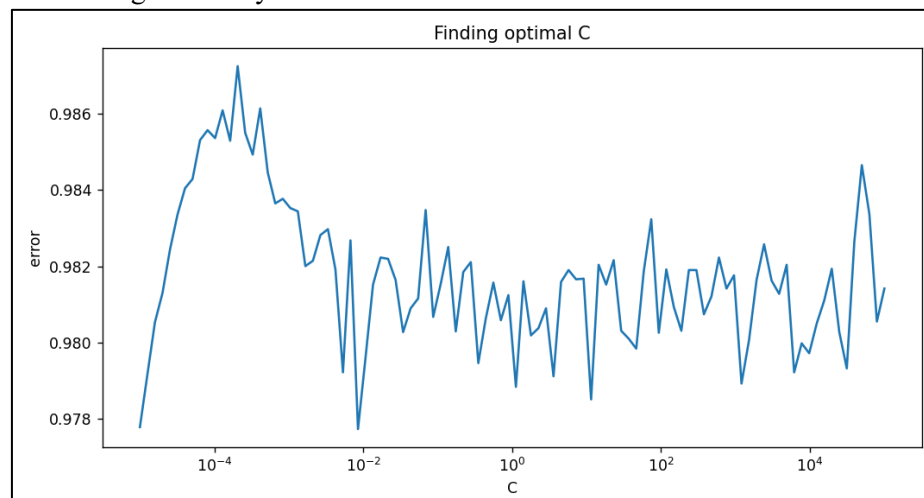


(First run)



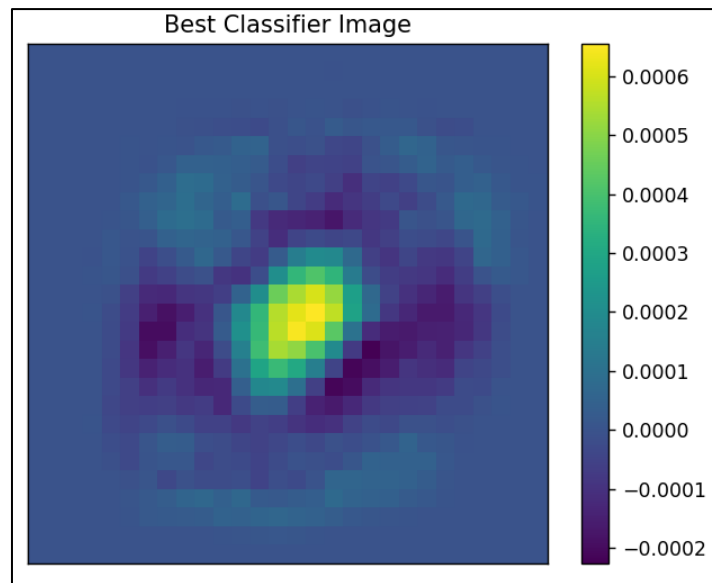
The best  $\eta_0$ : 0.7564

(b) The average accuracy on the validation set as a function of  $C$ :



The best  $C$ : 2.0565e-4

- (c) An image of the trained classifier over the best  $\eta_0, C$  that were obtained earlier, with 20000 SGD iterations:



We can interpret the image by matching the bright areas (the green and yellow ones) with the digit 8 and the dark areas (the dark blue and black ones) with the digit 0. The classifier gives the most positive weights to digits with presence in the center of the picture (in our case this is the middle of the digit 8); and gives the most negative weights to digits with presence in the mid-sides of the image's center (in our case this is the sides of the digit 0). Moreover, we can see that the classifier almost ignores the positions of top and bottom parts of 8 and 0, as we can expect since it looks same at these areas.

- (d) The accuracy of the best classifier on the test set: 0.9928.

(2) SGD for log-loss

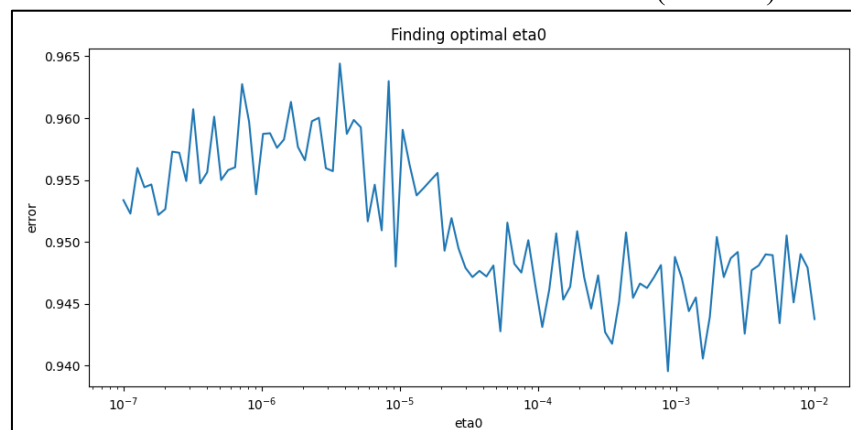
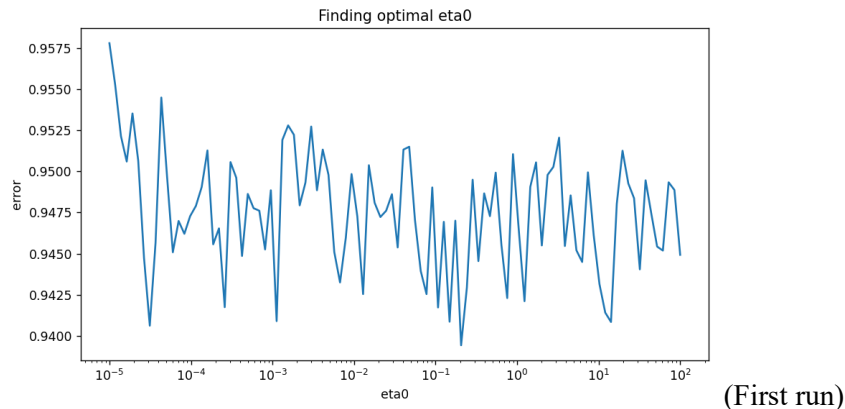
Consider the log-loss function:  $\ell_{\log}(\mathbf{w}, \mathbf{x}, y) = \log(1 + e^{-y\mathbf{w} \cdot \mathbf{x}})$ . First, we shall derive the gradient updates of SGD for this objective function. The derivative of  $\ell$  wrt  $\mathbf{w}$  is:

$$\nabla \ell = \frac{-ye^{-y\mathbf{w} \cdot \mathbf{x}}}{1 + e^{-y\mathbf{w} \cdot \mathbf{x}}} \mathbf{x}$$

Start with  $\mathbf{w}_0$ . In each update step, we randomly choose  $i \in \{1, \dots, n\}$  and denote the loss in  $x_i$  with classifier weight  $w$  as  $f_i(w) = \log(1 + e^{-y_i \mathbf{w} \cdot \mathbf{x}_i})$ . Then we get:

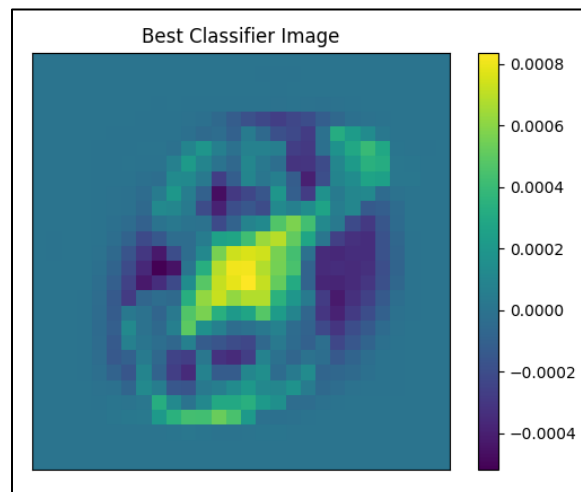
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f_i(\mathbf{w}_t) = \mathbf{w}_t + \eta_t \frac{y_i e^{-y_i \mathbf{w}_t \cdot \mathbf{x}_i}}{1 + e^{-y_i \mathbf{w}_t \cdot \mathbf{x}_i}} \mathbf{x}_i$$

(a) The average accuracy on the validation set as a function of  $\eta_0$ :



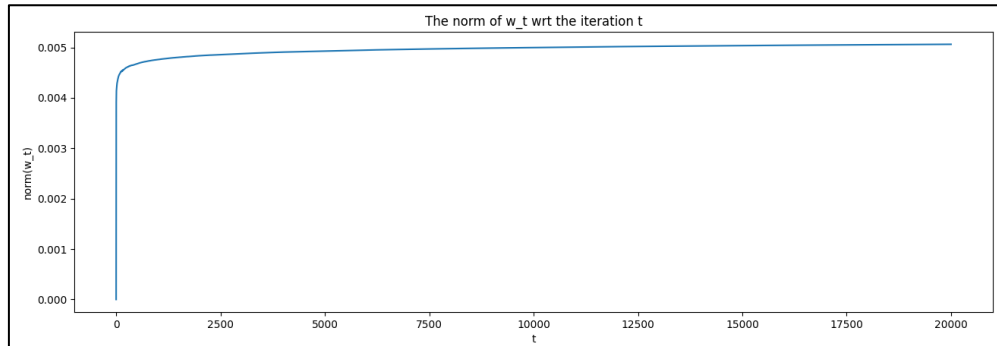
The best  $\eta_0$ : 3.274e-06

(b) An image of the trained classifier over the best that was obtained earlier, with 20000 SGD iterations:



The accuracy of the best classifier on the test set: 0.9851

(c) The norm of the classifier  $w_t$  as the iterations  $t$  growth:



As we can see, the norm of  $w_t$  get closer very fast to its optimum value (after  $\sim 130$  iterations it reaches a norm of 0.0045, after 1000 iterations 0.00475 and after 20000 iterations it reaches a norm of 0.00505). We can conclude that as  $t$  is growth,  $w_t$  does not changes a lot. Recall that the SGD choose the parameter  $\eta_t = \frac{\eta_0}{t}$  and the changes of the vector  $w$  is  $w_{t+1} = w_t - \eta_t \nabla f_i(w_t)$ . Therefore, it makes sense that  $w_t$  does not change a lot as  $t$  is big. The interesting phenomena is that indeed we obtain a very accurate result after few iterations.