

# ECE 544- Homework 2

Ahmadreza Eslaminia (Ae15)

## Questions:

Ahmadreza Eslaminia

ECE 544

HW2 ✖

AE15

1. [Max-Margin SVM]  $D = \left\{ \begin{matrix} x_1 \\ x_2 \end{matrix} \right\} \in \mathbb{R}^2$   $y \in \{-1, 1\}$

$$\min_{w, b} \frac{\|w\|_2^2}{2} \quad \forall (x^i, y^i) \in D \quad y^i (w^T x^i + b) \geq 1$$

a) we will have 4 constraints as follows:

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\begin{aligned} x^1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, y^1 = 1 &\Rightarrow y^1 (w^T x^1 + b) \geq 1 \Rightarrow w_1 + b \geq 1 \\ x^2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, y^2 = 1 &\Rightarrow y^2 (w^T x^2 + b) \geq 1 \Rightarrow w_2 + b \geq 1 \\ x^3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, y^3 = -1 &\Rightarrow y^3 (w^T x^3 + b) \geq 1 \Rightarrow -b \geq 1 \\ x^4 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, y^4 = -1 &\Rightarrow y^4 (w^T x^4 + b) \geq 1 \Rightarrow -(w_1 + w_2 + b) \geq 1 \end{aligned}$$

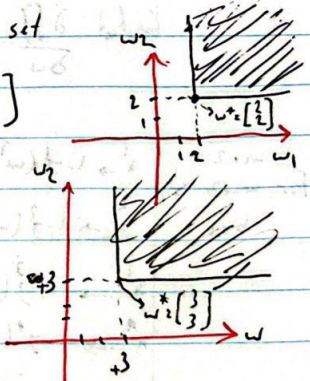
$$\left\{ \begin{aligned} w_1 + b &\geq 1 \\ w_2 + b &\geq 1 \\ b &\leq -1 \\ w_1 + w_2 &\geq 1 + b \end{aligned} \right.$$

b)  $b = 0 \Rightarrow$  does not obey  $b \leq -1$  condition  $\Rightarrow$  no feasible set

$$b = -1, \begin{cases} w_1 + b \geq 1 \Rightarrow w_1 \geq 2 \\ w_2 + b \geq 1 \Rightarrow w_2 \geq 2 \\ b = -1 \checkmark \end{cases} \Rightarrow \text{optimal } w^* = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$b = -2, \begin{cases} w_1 + b \geq 1 \Rightarrow w_1 \geq 3 \\ w_2 + b \geq 1 \Rightarrow w_2 \geq 3 \\ b = -2 \checkmark \end{cases} \Rightarrow \text{optimal } w^* = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

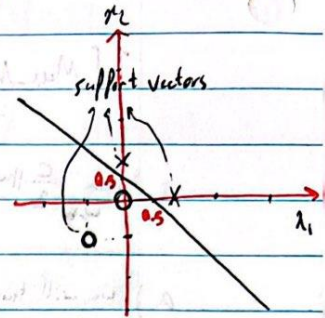
Since  $b \geq -1 \Rightarrow \begin{cases} w_1 \geq 1 - b \\ w_2 \geq 1 - b \end{cases} \Rightarrow \begin{cases} w_1 \geq 2 \\ w_2 \geq 2 \end{cases} \Rightarrow$  for this objective function  $w^* \begin{bmatrix} 2 \\ 2 \end{bmatrix}$  is the optimal choice in feasible space with  $b = -1$



7.c)

we find points  $n^1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ ,  $n^2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ ,  $n^3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$   
are support vector points.

$$\Rightarrow \begin{cases} n^1 \Rightarrow w_1 + b = 1 \\ n^2 \Rightarrow w_2 + b = 1 \\ n^3 \Rightarrow b = -1 \end{cases} \Rightarrow \begin{cases} b = 1 \\ w_1 + w_2 = 2 \end{cases} \Rightarrow w_1 n^1 + w_2 n^2 + b = 0 \\ = 2n_1 + 2n_2 = 1$$



d) they are linearly separable s.t. Eq. 7 has a feasible solution.

$$e) \min \frac{c}{2} \|w\|_2^2 + \sum_{i=0}^I y^i \quad \text{s.t. } y^i: \max\{0, 1 - y^i (w^T n^i + b)\}$$

$$\Rightarrow \min \frac{c}{2} \|w\|_2^2 + \sum_{i=0}^I \max\{0, 1 - y^i (w^T n^i + b)\}$$

$$\text{Gradient: } \frac{\partial f}{\partial w} = \begin{bmatrix} c w_1 \\ c w_2 \end{bmatrix} + \sum_{i=0}^I y^i \begin{bmatrix} n_1^i \\ n_2^i \end{bmatrix} \quad \text{that } 1 - y^i (w^T n^i + b) > 0 \text{ for } n \neq 0$$

$$\begin{aligned} w_1 = 2 \\ \text{For } w_2 = 2 \\ b = -1 \end{aligned} \Rightarrow \begin{cases} n^1 \rightarrow 1 - y^1 (w^T n^1 + b) = 1 - (2-1) = 0 \\ n^2 \rightarrow 1 - y^2 (w^T n^2 + b) = 1 - (2-1) = 0 \\ n^3 \rightarrow 1 - y^3 (w^T n^3 + b) = 1 - 1 = 0 \\ n^4 \rightarrow 1 - y^4 (w^T n^4 + b) = 1 + (-4-1) = -4 \end{cases} \Rightarrow \frac{\partial f}{\partial w} = c \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 2c \\ 2c \end{bmatrix}$$

Since  $c > 0 \Rightarrow \frac{\partial f}{\partial w} \neq 0$  so it is not optimal

Eff) from previous section we found the  $w = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$  is not optimal.

from the python file the solution is  $w = \begin{bmatrix} 0.667 \\ 0.667 \end{bmatrix}$ ,  $b = 0.333$ , loss: 1.778835

we can see by introducing new terms to loss the solution misclassifies one point (C=1)  
however, by choosing C closer to zero we can give more importance to the new terms so the problem can be solved.



2. [L2 SVM]  $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^{101}$   $x^{(i)} \in \mathbb{R}^d$   $y^{(i)} \in \{-1, 1\}$

$$\min \frac{1}{2} w^T w + \frac{c}{2} \sum_{i=1}^{101} (\tilde{y}^{(i)})^2 \quad \text{s.t.} \quad y^{(i)} (w^T x^{(i)} + b) \geq 1 - \tilde{y}^{(i)} \quad i \in \{1, \dots, 101\}$$

$\tilde{y}^{(i)} \geq 0$

a)  ~~$\frac{1}{2} (w^*)^T w^* + \frac{c}{2} \sum_{i=1}^{101} (\tilde{y}^{(i)})^2$~~   $\frac{c}{2} ((\tilde{y}^{(1)})^*)^2 + \frac{c}{2} ((\tilde{y}^{(2)})^*)^2$

b) since  $(\tilde{y}^{(3)})^*$  is not zero  $(\tilde{y}^{(3)})^* > 0 \Rightarrow$  loss function is as follow;

$$\frac{1}{2} (w^*)^T w^* + \frac{c}{2} ((\tilde{y}^{(1)})^*)^2 + \frac{c}{2} ((\tilde{y}^{(2)})^*)^2 + \frac{c}{2} ((\tilde{y}^{(3)})^*)^2 > \frac{c}{2} \text{ previous loss}$$

c) we show that  $\tilde{y}^i = 0$  does not violate the previous constraint so it's feasible solution

Previously we had  $(\tilde{y}^{(i)})^* < 0$  which results in constraint below:

$$y^{(i)} (w^*)^T x^{(i)} + b^* \geq 1 - (\tilde{y}^{(i)})^* > 1 \quad \textcircled{\text{I}}$$

however if we consider  $\tilde{y}^i = 0$  we have the following constraint:

$$y^{(i)} (w^*)^T x^{(i)} + b^* \geq 1 \quad \textcircled{\text{II}}$$

As you can see  $\textcircled{\text{II}}$  does not violate  $\textcircled{\text{I}}$  which means  $\tilde{y}^i = 0$  is a feasible solution

d) As we have shown in previous sections not only  $(w^*, b^*, \tilde{y}^*)$  is a feasible solution but also it results in smaller loss, so this would violate the first assumption that  $(w^*, b^*, \tilde{y}^*)$  is optimal solution. so we can conclude that

For the optimal solution  $\tilde{y}^i \geq 0$  always hold so we can remove it.



Scanned with CamScanner

P3

2. e) The Lagrangian would be as follows

(1) e

$$L(\omega, b, \gamma, \alpha, y, x^{(i)}, c, D) = \frac{1}{2} \omega^T \omega + \frac{c}{2} \sum_{i=1}^{|D|} (y^{(i)})^2 - \sum_{i=1}^{|D|} \alpha^{(i)} (y^{(i)} (\omega^T x^{(i)} + b) - 1 + y^{(i)})$$

$$\forall i \Rightarrow \alpha^{(i)} \geq 0$$

f) show the dual of the program is

$$\max_{\alpha \geq 0} -\frac{1}{2} \alpha^T (Q + I \frac{1}{c}) \alpha + \gamma^T \alpha$$

$$\text{s.t. } \gamma^T \alpha = 0$$

$$\text{where } Q_{ij} = \gamma_i \gamma_j x_i^T x_j \quad \alpha \in \mathbb{R}^{|D|} \quad \gamma \in \mathbb{R}^{|D|}$$

Take following partial derivatives:

$$\nabla_{\omega} L = 0 \Rightarrow \omega - \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} x^{(i)} = 0 \Rightarrow \omega = \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} x^{(i)}$$

$$\nabla_b L = 0 \Rightarrow \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} = 0$$

$$\nabla_{\gamma} L = 0 \Rightarrow c y^{(i)} - \alpha^{(i)} = 0 \Rightarrow y^{(i)} = \frac{\alpha^{(i)}}{c} \quad / \text{ placing above in Lagrangian:}$$

$$L = \frac{1}{2} \left( \sum_{i=1}^{|D|} (\alpha^{(i)} y^{(i)} x^{(i)})^T \right) \left( \sum_{j=1}^{|D|} \alpha^{(j)} y^{(j)} x^{(j)} \right) + \frac{c}{2} \sum_{i=1}^{|D|} \frac{(\alpha^{(i)})^2}{c^2} - \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} \left( \sum_{j=1}^{|D|} x_j^T y^{(j)} x_j \right) x^{(i)} + b + \sum_{i=1}^{|D|} \alpha^{(i)} - \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)}$$

$$= -\frac{1}{2} \alpha^T \left( \sum_{i,j=1}^{|D|} \gamma_i \gamma_j x_i^T x_j \right) \alpha - \frac{1}{2} \alpha^T I \frac{1}{c} \alpha + \gamma^T \alpha = -\frac{1}{2} \alpha^T (Q + I \frac{1}{c}) \alpha + \gamma^T \alpha$$

$$\sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} = 0 \Rightarrow \gamma^T \alpha = 0$$



Scanned with CamScanner

Pa



3. [Support Vector Machine]  $\min \frac{1}{2} \|\omega\|_2^2 \text{ s.t. } y^{(i)}(\omega^T x^{(i)} + b) \geq 1, \forall (x^{(i)}, y^{(i)}) \in D$

a)  $\Rightarrow$  Lagrangian with multipliers  $\alpha_i$ :

$$L(\omega, b, \alpha) = \frac{1}{2} \|\omega\|_2^2 + \sum_{i=1}^{|D|} \alpha^{(i)} (1 - y^{(i)} \omega^T x^{(i)} + b)$$

b) consider  $L(\omega, b, \alpha) = \frac{1}{2} \|\omega\|_2^2 + \sum_{i=1}^{|D|} \alpha^{(i)} (1 - y^{(i)} \omega^T x^{(i)} + b)$

$$\nabla_{\omega} L = 0 \Rightarrow \omega - \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} x^{(i)} = 0 \Rightarrow \omega = \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} x^{(i)}$$

$$\Rightarrow L = \frac{1}{2} \left\| \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} x^{(i)} \right\|_2^2 + \sum_{i=1}^{|D|} \alpha^{(i)} - \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} \left( \sum_{j=1}^{|D|} \alpha^{(j)} y^{(j)} x^{(j)} \right)^T x^{(i)}$$

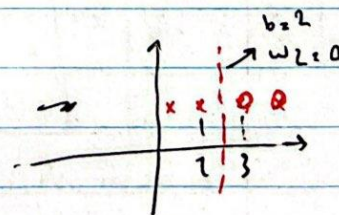
$$= \sum_{i=1}^{|D|} \alpha^{(i)} - \frac{1}{2} \left\| \sum_{i=1}^{|D|} \alpha^{(i)} y^{(i)} x^{(i)} \right\|_2^2$$

s.t.  $\alpha_i \geq 0 \quad \forall i \in D$

c)  $k(x, z) = (x^T z)^2 + 1 = \phi(x)^T \phi(z) \xRightarrow{\text{is 2D}} \phi(x) = \begin{bmatrix} x^2 \\ 1 \end{bmatrix}$

$i$	$x^{(i)}$	$\phi(x^{(i)})$	$y_i$
1	$\frac{1}{2}$	$\begin{bmatrix} \frac{1}{4} \\ 1 \end{bmatrix}$	+1
2	-1	$\begin{bmatrix} 1 \\ 1 \end{bmatrix}$	+1
3	$\sqrt{3}$	$\begin{bmatrix} 3 \\ 1 \end{bmatrix}$	-1
4	2	$\begin{bmatrix} 4 \\ 1 \end{bmatrix}$	-1

2, 3 support vectors  $\rightarrow$  midpoint is  $\frac{1}{2}$



$$\phi_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \omega^T \phi + b = 1$$

$$\phi_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \omega^T \phi + b = 0$$

$$\phi_3 = \begin{bmatrix} 3 \\ 1 \end{bmatrix} \Rightarrow \omega^T \phi + b = -1$$

$$\omega^T [\omega^1, \omega^2]$$

$$\Rightarrow \boxed{\omega_1 = -1, b = 2}$$

3.d) points at ②, ③ are support vectors.

we have:  $w^* = \sum_{i=1}^N \alpha^{(i)} y^{(i)} z^{(i)}$   $z = \phi(h)$

$\alpha^{(1)} = \alpha^{(4)} = 0$

$\Rightarrow \begin{bmatrix} -1 \\ 0 \end{bmatrix} = \alpha_2 \times 1 \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \alpha_3 \times -1 \times \begin{bmatrix} 3 \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha_2 - 3\alpha_3 \\ \alpha_2 - \alpha_3 \end{bmatrix} \Rightarrow \alpha_2 = \alpha_3 = \frac{1}{2}$   
 $\alpha^{(1)}, \alpha^{(4)} = 0$  (not support vector)



4. [Multi logistic regression]  $\min_{(w,b)} \sum_{(x,y) \in D} p(y|x)$  where  $p(y|x) = \frac{e^{w_y^T [x; 1]}}{\sum_{j \in \{0,1,2\}} e^{w_j^T [x; 1]}}$

2 weight + 1 bias

a)  $w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$  where  $w_i = \begin{bmatrix} w_i^{(2)} \\ w_i^{(1)} \\ w_i^{(0)} \end{bmatrix}$   $w_1 = \begin{bmatrix} w_1^{(2)} \\ w_1^{(1)} \\ w_1^{(0)} \end{bmatrix}$   $w_2 = \begin{bmatrix} w_2^{(2)} \\ w_2^{(1)} \\ w_2^{(0)} \end{bmatrix}$

we have one  $w_j$  for each classes. so since we have 3 parameter at each class

and 3 classes  $\Rightarrow 3 \times 3 = 9$  Parameter  $\Rightarrow \dim w = 9 \times 1$

b)  $w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$   $\Psi(x,y) = \begin{bmatrix} \phi(x) & \delta(y=0) \\ \phi(x) & \delta(y=1) \\ \phi(x) & \delta(y=2) \end{bmatrix} = \begin{bmatrix} [x; 1] & \delta(y=0) \\ [x; 1] & \delta(y=1) \\ [x; 1] & \delta(y=2) \end{bmatrix}$

$w^T \Psi(x,y) = w_0^T [x^{(1)}; 1] \delta(y=0) + w_1^T [x^{(1)}; 1] \delta(y=1) + w_2^T [x^{(1)}; 1] \delta(y=2)$   
 $= w_j^T [x^{(1)}; 1] \quad \forall y^{(1)} \text{ as } i \in \{0,1,2\}$

c) we had  $F$  as following in past:

$F(w, x, y) = \begin{bmatrix} w_0^T [x; 1] \\ w_1^T [x; 1] \\ w_2^T [x; 1] \end{bmatrix} = \begin{bmatrix} w_0^{(1)} x_2 + w_0^{(1)} x_1 + b_0 \\ w_1^{(1)} x_2 + w_1^{(1)} x_1 + b_1 \\ w_2^{(1)} x_2 + w_2^{(1)} x_1 + b_2 \end{bmatrix} = [W X + b]$

where  $W = \begin{bmatrix} w_0^{(1)} & w_0^{(1)} \\ w_1^{(1)} & w_1^{(1)} \\ w_2^{(1)} & w_2^{(1)} \end{bmatrix}$   $3 \times 2$

$b = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}$   $3 \times 1$

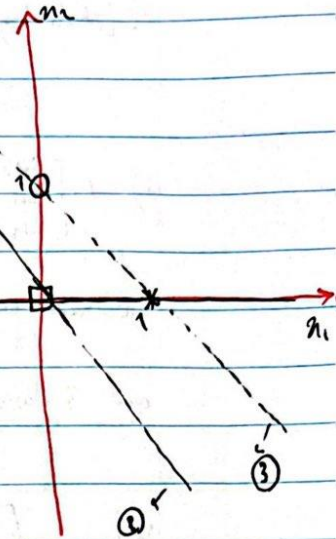
in this case row of  $W$  are the  $w$  vector with ignoring biases. That's how it relate to previous  $w$

$$4.d) W = \begin{bmatrix} 3 & 0.5 \\ 0 & 1 \\ -1.5 & -1.5 \end{bmatrix}$$

$$b = \begin{bmatrix} 0 \\ 0 \\ 1.5 \end{bmatrix}$$

$$\begin{aligned} * & \leftarrow x^1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, y^1 = 0 \\ \bigcirc & \leftarrow x^2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, y^2 = 1 \\ \square & \leftarrow x^3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, y^3 = 0.2 \end{aligned}$$

$$[Wx + b] = \begin{cases} 3x_1 + 0.5x_2 = 0 & (1) \\ x_2 = 0 & (2) \\ -1.5x_1 - 1.5x_2 + 1.5 = 0 & (3) \end{cases}$$



$$Wx^{(1)} + b = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{class 1} \checkmark$$

$$Wx^{(2)} + b = \begin{bmatrix} 0.5 \\ 1 \\ 0 \end{bmatrix} \rightarrow \text{class 2} \checkmark$$

$$Wx^{(3)} + b = \begin{bmatrix} 0 \\ 0 \\ 1.5 \end{bmatrix} \rightarrow \text{class 3} \checkmark$$

As all numbers are positive we can see the higher magnitude will slow the class because  $e^x$  is not going to change that.

So the model correctly classifies all the data points

$$e) W = \begin{bmatrix} 8.7387 & -1.6501 \\ -1.9086 & 9.0514 \\ -7.2665 & -6.9575 \end{bmatrix}$$

$$b = \begin{bmatrix} -2.5121 \\ -2.5161 \\ 5.3407 \end{bmatrix}$$

$$P(y|x) = \begin{bmatrix} 9.9981 \times 10^{-1}, 1.7380 \times 10^{-5}, 6.9603 \times 10^{-9} \\ 3.0772 \times 10^{-5}, 9.9987 \times 10^{-1}, 9.4992 \times 10^{-9} \\ 1.6024 \times 10^{-9}, 1.1721 \times 10^{-9}, 9.9835 \times 10^{-1} \end{bmatrix}$$



## 5. [Multi class classification via NN]

a) since it has to be atleast one class  $\sum f(z_i) = 1 \Rightarrow \text{Softmax}$  and  $0 \leq f(z_i) \leq 1$

$$f(z_i) = \frac{e^{z_i}}{\sum_{i=1}^3 e^{z_i}}$$

$$b) G(z) = \left( \frac{z_1}{z_1 + z_2 + z_3}, \frac{z_2}{z_1 + z_2 + z_3}, \frac{z_3}{z_1 + z_2 + z_3} \right)$$

Since  $f$  uses exponential function we know that  $f$  is translation invariant.

$$\text{about } G, G(z^{(1)}) = \left( \frac{0.1}{0.4}, \frac{0.1}{0.4}, \frac{0.2}{0.2} \right) = \left( \frac{1}{4}, \frac{1}{4}, \frac{1}{2} \right)$$

$$G(z^{(1)}) = \left( \frac{1.01}{3.04}, \frac{1.01}{3.04}, \frac{1.04}{3.04} \right) = \left( \frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right)$$

}  $G$  is not translation invariant

we can understand a translation invariant function (such as  $f$ ) is more desirable than a non-invariant function (such as  $G$ ) since the nontranslation one if the numbers, yet large in comparison to their difference they are going to give same output for all the inputs.

$$c) \mathcal{L}(y, f(z)) = -\ln \left( \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}} \right) = -z_3 + \ln(e^{z_1} + e^{z_2} + e^{z_3})$$

$$\frac{\partial \mathcal{L}(y, f(z))}{\partial z_3} = -1 + \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\mathcal{L}(y, G(z)) = -\ln \left( \frac{z_3}{z_1 + z_2 + z_3} \right) = -\ln z_3 + \ln(z_1 + z_2 + z_3)$$

} As you can see  
At  $z = (1, 10^{-5}, 10^{-5})$   
gradient of  $\mathcal{L}(y, G(z))$   
will explode and will

$$\frac{\partial \mathcal{L}(y, G(z))}{\partial z_3} = -\frac{1}{z_3} + \frac{1}{z_1 + z_2 + z_3}$$

result to optimization issues

P2

5.d) one-vs.-rest is going to be slower. Since in Multinomial classification Although the last layer has more parameters ( $3x$ ), the whole computational cost (forward) for 3 of one-vs.-rest  $ANN$  which is needed to completely assign a correct class to input is higher than multi-class classification.



### Codes A3-SVM:

```
C: > DriveA > UIUCcourses > Fall 2023 > ECE 544 pattern recognition > HWS > hw2 > homework2 > A3_SVM.

1  import torch
2  import torch.optim as optim
3  import torch.nn as nn
4  import matplotlib.pyplot as plt
5
6  torch.manual_seed(1)
7  X = torch.tensor([[1, 0, 0, -1],[0, 1, 0, -1]], dtype=torch.float32)
8  y = torch.tensor([1, 1, -1, -1], dtype=torch.float32)
9  # initialize w and b
10 w = torch.tensor([0.0, 0.0], requires_grad=True)
11 b = torch.tensor([0.0], requires_grad=True)
12 alpha = 0.001
13 C = 1
14
15 optimizer = optim.SGD([w,b], lr=alpha, weight_decay=0)
16 optimizer.zero_grad()
17
18 grads = []
19 losses = []
20
21 for iter in range(10000):
22     if iter==0:
23         print('1-y(w^T*x+b): {}'.format(1 - y*(torch.matmul(X.T, w)+b)))
24         #####
25         ## Complete this line which is our cost function
26         ## Dimensions: loss (scalar)
27         #####
```

```

24 #####
25 ## Complete this line which is our cost function
26 ## Dimensions: loss (scalar)
27 #####
28
29 New_term = 1 - y*(torch.matmul(X.T, w)+b)
30 loss = C/2 * torch.matmul(w, w) + torch.sum(torch.clamp(New_term, min=0))
31 loss.backward()
32 gn = torch.norm(w.grad)**2 + torch.norm(b.grad)**2
33 print("Iter: %d; Loss: %f; ||Grad||: %f" % (iter, loss, gn))
34 optimizer.step()
35 optimizer.zero_grad()
36
37 losses.append(loss.item())
38 grads.append(gn)
39
40 print('w: {}'.format(w.data))
41 print('b: {}'.format(b.data))
42
43 plt.figure()
44 plt.plot(losses, label='Loss')
45 plt.title('Loss')
46 plt.show()
47 plt.figure()
48 plt.plot(grads, label='Grad')
49 plt.title('||Grad||')
50 plt.show()

```

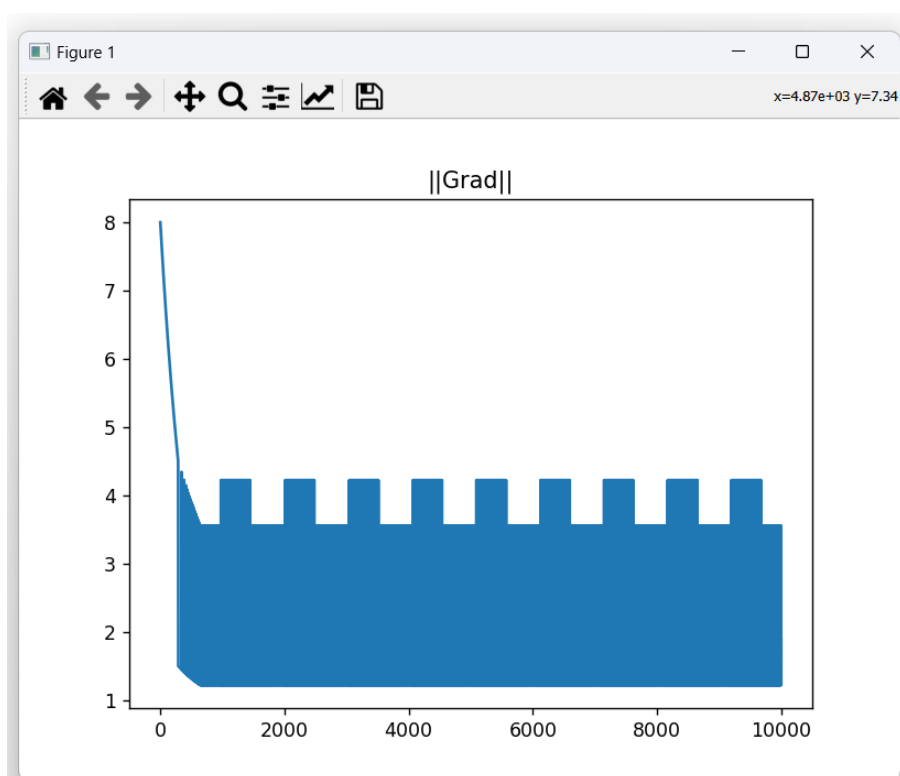
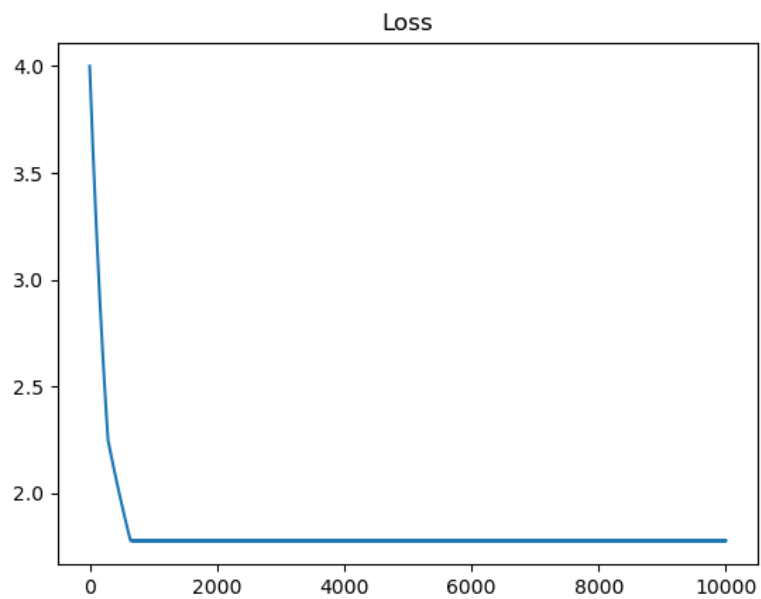
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL

```

Iter: 9997; Loss: 1.778837; ||Grad||: 1.891705
w: tensor([0.6671, 0.6671])
b: tensor([0.3320])
Iter: 9998; Loss: 1.778868; ||Grad||: 1.221706
w: tensor([0.6674, 0.6674])
b: tensor([0.3330])
Iter: 9999; Loss: 1.778391; ||Grad||: 1.890811
w: tensor([0.6667, 0.6667])
b: tensor([0.3320])

```





#### A4-Multiclass:

```
import torch
import torch.optim as optim
import torch.nn as nn

torch.manual_seed(1)
alpha = 1
C = 0

#####
## encode the dataset to fit the one specified in HW4.pdf (note that bias
## is part of the network now)
## Dimensions: X (2x3); y (3)
#####
X = torch.Tensor([[1, 0, 0],[0, 1, 0]])
y = torch.LongTensor([0, 1, 2])

class ShallowNet(nn.Module):
    def __init__(self):
        super(ShallowNet, self).__init__()
        self.fc1 = nn.Linear(2,3, bias=True)

    def forward(self, X):
        return self.fc1(X)

net = ShallowNet()
print(net)

print(net(torch.transpose(X,0,1)).squeeze())

optimizer = optim.SGD(net.parameters(), lr=alpha, weight_decay=C)
optimizer.zero_grad()

criterion = nn.CrossEntropyLoss()

for iter in range(10000):
    netOutput = net(torch.transpose(X,0,1))
```



C:\> DriveA > UIUCcourses > Fall 2023 > ECE 544 pattern recognition > HWS > hw2 > homework2 > A4\_Multiclass.py

```
31 optimizer = optim.SGD(net.parameters(), lr=alpha, weight_decay=C)
32 optimizer.zero_grad()
33
34 criterion = nn.CrossEntropyLoss()
35
36 for iter in range(10000):
37     netOutput = net(torch.transpose(X,0,1))
38
39     #####
40     ## provide the arguments for the criterion function
41     ## Dimensions: loss (scalar)
42     #####
43     loss = criterion( netOutput , y )
44
45     loss.backward()
46     gn = 0
47     for f in net.parameters():
48         gn = gn + torch.norm(f.grad)
49     print("Loss: %f; ||g||: %f" % (loss, gn))
50
51     #####
52     ## Use two functions within the optimizer instance to perform the update step
53     #####
54     optimizer.step()
55     optimizer.zero_grad()
56
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL

```
Parameter containing:
tensor([[ 8.7387, -1.6501],
        [-1.9086,  9.0514],
        [-7.2685, -6.9575]], requires_grad=True)
Parameter containing:
tensor([-2.5121, -2.5161,  5.3407], requires_grad=True)
tensor([[9.9981e-01, 1.7380e-05, 6.9603e-04],
        [3.0772e-05, 9.9987e-01, 9.4992e-04],
        [1.6024e-04, 1.1721e-04, 9.9835e-01]], grad_fn=<SoftmaxBackward0>)
PS C:\DriveA\UIUCcourses\Fall 2023\ECE 544 pattern recognition\HWS\hw2\homework2> |
```