Class ID-19   Name: Sayeed Shafayet Chowdhury

## Problem #1(a)(i)

The problem asks us to show that in linearly separable case, the logistic regression solution makes the weights $\bar{w} \to \infty$. The $w_0$ and logistic function is defined as $\sigma(z) = \dfrac{1}{1+\exp\{-z\}}$

Now given the dataset, $\left\{(x_i, d_i)\right\}_{i=1}^{N}$, where $x_i \in \mathbb{R}^d$, $y_i \in \{0,1\}$, $\bar{w} \in \mathbb{R}^d$

the likelihood function becomes,

$$P\left(\{y_n\}_{n=}^{N} \mid \{x_n\}_{n=1}^{N}, \bar{w}\right) = \prod_{n=1}^{N} \sigma\left(\bar{w}^T x_n\right)^{y_n} \left(1 - \sigma\left(\bar{w}^T x_n\right)\right)^{1-y_n}$$

This is the function we want to maximize with respect to $\bar{w}$. So,

$$\bar{w}^* = \underset{\frac{d}{\bar{w}}}{\arg\max}\left\{\sum_{n=1}^{N} y_n \log \sigma\left(\bar{w}^T x_n\right) + (1-y_n) \log\left(1-\sigma\left(\bar{w}^T x_n\right)\right)\right\} \quad \text{——} ①$$

This is the same as minimizing the logistic regression

loss function, $J = \sum_{n=1}^{N} - \left\{y_n \log \sigma\left(\bar{w}^T x_n\right) + (1-y_n)\log\left(1-\sigma\left(\bar{w}^T x_n\right)\right)\right\}$

So, we will focus on solving ①, but it is not possible to solve ① directly by setting the gradient to 0 and solving for $\bar{w}$. So, we'll solve for ① using gradient ascent.

Now, let us examine the gradient arising from a single point of data.

$$\nabla_{\bar{w}} \left\{ y_n \log \sigma(\bar{w}^T x_n) + (1-y_n) \log (1 - \sigma(\bar{w}^T x_n)) \right\}$$

Here we note that, $1 - \sigma(z) = 1 - \dfrac{1}{1+e^{-z}} = \dfrac{1 + e^{-z} - 1}{1 + e^{-z}}$

$$= \dfrac{e^{-z}}{1+e^{-z}}$$

$$= \dfrac{1}{1+e^{z}}$$

$$= \sigma(-z) \quad \text{—} \textcircled{A}$$

and $\dfrac{d}{dz} \sigma(z) = (1 - \sigma(z)) \, \sigma(z) \quad \text{—} \textcircled{B}$

let us have a look at the gradient.

$$\nabla_{\bar{w}} \sum_{n=1}^{N} \left\{ y_n \log \sigma(\bar{w}^T x_n) + (1-y_n) \log (1 - \sigma(\bar{w}^T x_n)) \right\}$$

$$= \sum_{n=1}^{N} \left[ y_n \dfrac{1}{\sigma(\bar{w}^T x_n)} \nabla_{\bar{w}} \left\{ \sigma(\bar{w}^T x_n) \right\} + (1-y_n) \dfrac{1}{\sigma(-\bar{w}^T x_n)} \nabla_{\bar{w}} \left\{ \sigma(-\bar{w}^T x_n) \right\} \right]$$

$$\left[ \because 1 - \sigma(\bar{w}^T x_n) \atop = \sigma(-\bar{w}^T x_n) \right]$$

$$= \sum_{n=1}^{N} \left[ y_n x_n (1 - \sigma(\bar{w}^T x_n)) - (1-y_n) x_n \sigma(\bar{w}^T x_n) \right] ; \quad \left[ \text{using } \textcircled{B} \right]$$

$$\left[ \nabla_{\bar{w}} \left\{ \sigma(-\bar{w}^T x_n) \right\} \atop = -x_n \sigma(-\bar{w}^T x_n) (1 - \sigma(-\bar{w}^T x_n)) \right]$$

Now, $1 - \sigma(-\bar{w}^T x_n) = \sigma(\bar{w}^T x_n)$

$$= \sum_{n=1}^{N} \left[ y_n x_n - y_n x_n \sigma\left(\overline{w}^T x_n\right) - x_n \sigma\left(\overline{w}^T x_n\right) + x_n y_n \sigma\left(\overline{w}^T x_n\right) \right]$$

$$= \sum_{n=1}^{N} \left[ y_n x_n - x_n \sigma\left(\overline{w}^T x_n\right) \right] = \sum_{n=1}^{N} \left[ x_n \left(y_n - \sigma\left(\overline{w}^T x_n\right)\right) \right]$$

Then using the gradient ascent update, we get-

$$\overline{w}^{(t+1)} \leftarrow \overline{w}^{(t)} + \alpha \left[ \sum_{n=1}^{N} \left( x_n \left(y_n - \sigma\left(\left(\overline{w}^{(t)}\right)^T x_n\right)\right)\right) \right] \quad —②$$

Now, let us consider the linear separability condition, it creates a problem for the unregularized logistic regression problem (the case considered here). We note that, decision boundary in a linear classifier is independent of the scale of the parameters. This can be seen as the decision boundary is the set $\{\overline{x} : \overline{w}^T \overline{x} = 0\}$ and this set doesn't change if $\overline{w}$ is multiplied by some constant. For a given decision boundary, however, the scale does $\ell$ affect the likelihood in logistic regression by causing the logistic function to become more steep. We can see this by thinking

about the derivative of $\sigma(z)$ evaluated at $z=0$ versus the derivative of $\sigma(10z)$ at $z=0$. The derivatives are $\sigma(z)(1-\sigma(z))$ and $10\sigma(z)(1-\sigma(z))$, respectively, and so the linear regime in the middle is ten times steeper when input is scaled by a factor of ten.

If we currently have a decision boundary such that the data are all correctly classified, then increasing the scale of the weights pushes the predictions further towards their correct answers. Say, we have a set of weights $\hat{w}$ with unit norm, i.e. $\|\hat{w}\|=1$. Then we construct a logistic regression classifier with weights $\bar{w} = c\hat{w}$ and seek only to fit the constant $c>0$ to the data. Here, changing $c$ does not move the decision boundary for the classifier. Let $\mathbf{x}_n$ the $n$th example and examine the derivative of its log likeli-hood with respect to $c$:

$$\frac{\partial}{\partial c}\left\{ y_n \log \sigma(c\hat{w}^T x_n) + (1-y_n)\log(1-\sigma(c\hat{w}^T x_n)) \right\}$$
$$= \hat{w}^T x_n \left(y_n - \sigma(c\hat{w}^T x_n)\right)$$

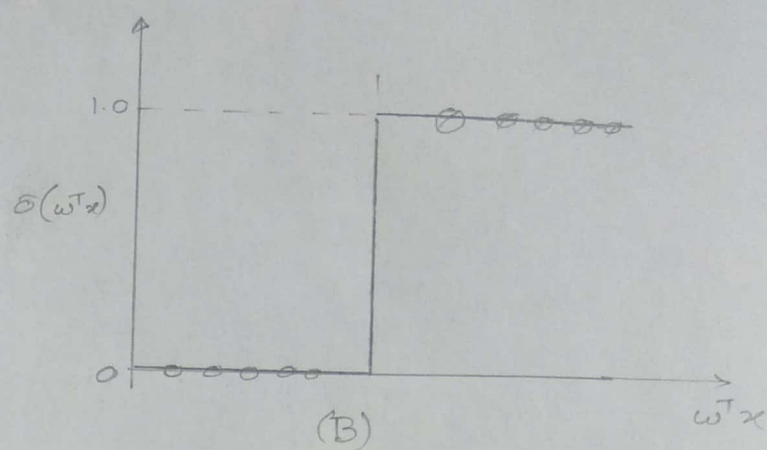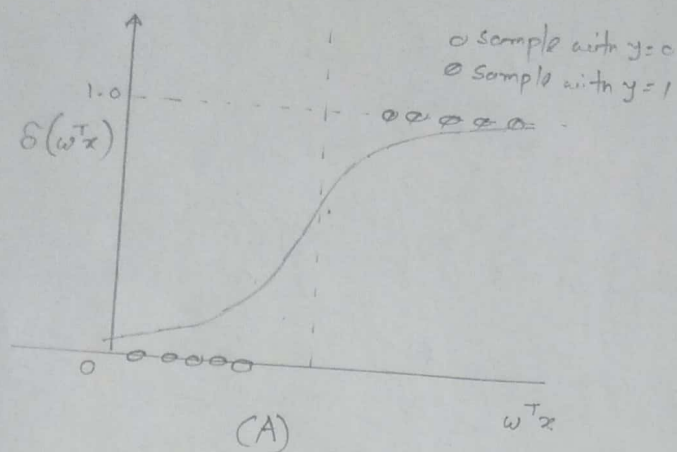If $y_n = 0$ then $(y_n - \delta(c\hat{w}^T x_n)) < 0$ and if $y_n = 1$ then $(y_n - \delta(c\hat{w}^T x_n)) > 0$. We also note that due to the fixed decision boundary, if $y_n = 0$ is classified correctly then $\hat{w}^T x_n < 0$ and is positive otherwise. Similarly if $y_n = 1$ is classified correctly, then $\hat{w}^T x_n > 0$ and is negative otherwise. Thus the derivative of the log likelihood w.r.t. $c$ is always positive for an example that $\hat{w}$ classifies correctly. Since in our case, the data are linearly separable, there exists a $\hat{w}$ such that all of the data have log likelihoods with positive derivaties with respect to $c$. In that situation, gradient ascent on $c$ would cause it to grow without bound, this essentially drives the sigmoid function to be sharper and sharper until it becomes a theaviside step function. So, to predict all the samples correctly wit certainty, the gradient ascent on maximizing

the log conditional likelihood (LCL) sends the $\vec{w}$ to infinity. Again, this can be seen as, for weight parameters $w_j$, $\frac{\partial}{\partial w_j} LCL = \sum_i (y_i - p_i) x_{ij}$, here we consider only positive samples ($y_i = 1$), $x_{ij}$ is the $j^{th}$ feature of $i^{th}$ sample.

Now this derivative will always be positive, as long as the predicted probability $p_i$ is not perfectly one for all the positive samples. So, the logistic loss keeps on decreasing with iteration as $\vec{w}$ keeps increasing and eventually $\vec{w}$ goes towards infinity.

This fact can also be viewed as — in order to maximize the likelihood, we need to minimize $\log(1 + e^{-y_i w^T x})$, so we want the exponent to be as negative as possible.

Now, the exponent $y_i \cdot \|\vec{w}\|_2 \cdot \|x_i\|_2 \cos\theta$ to make as large as possible, the algorithm simply increases $\|\vec{w}\|_2$ without bound.

Legend for graph (A): o Sample with y=0, ⊘ Sample with y=1. Graphs: (A) shows $\sigma(w^T x)$ vs $w^T x$ (sigmoid); (B) shows $\sigma(w^T x)$ vs $w^T x$ (step function).

We can also visualize this graphically. Although in (A), we can correctly classify the samples, still samples with $y=0$ have non-zero probability, while samples with $y=1$ have $p<1$. So the algorithm while trying to minimize the logistic loss drives $\bar{w}$ towards a position where all $x$'s with $y_i=0$ have $p_i=0$ and all $x$'s with $y_i=1$ have $p_i=1$ (as shown in (B)). This essentially drives the $\bar{w} \to \infty$, as the slope of the sigmoid $\to \infty$ at the middle.

Since with $\bar{w} \to \infty$, $h(x) = \dfrac{1}{1+e^{-w^T x}}$, $\lim\limits_{w \to \infty} h(x) = \text{sign}(w^T x)$

So $h(x)$ becomes a step function and to make it a step function, $w \to \infty$.

Now, from ②, $\bar{w}^{(k+1)} \leftarrow \bar{w}^{(k)} + \alpha\left[\sum_{n=1}^{N} \bar{x}_n \left(y_n - \sigma\left(\bar{w}^{T(k)} x_n\right)\right)\right]$

if we allowed the algo to run forever (only let it stop when $\|\bar{w}^{(k+1)} - \bar{w}^{(k)}\|_2 = 0$), the algo would not stop until the gradient $= 0$, meaning $y_n = \sigma\left(\bar{w}^{T(k)} x_n\right)$ for all $n$.

Now, say $y_n = 1$, then $\sigma\left(\bar{w}^{T(k)} x_n\right)$ must be 1 for algo to stop but $\sigma\left(\bar{w}^{T(k)} x_n\right) = \dfrac{1}{1+e^{-\bar{w}^{(k)} x_n}}$, to make it exactly $=1$,

$e^{-\bar{w}^{T(k)} x_n}$ must be $= 0 \Rightarrow \bar{w}^{(k)}$ must be $\rightarrow \infty$. So, the algo would keep chasing this criterion. However to prove that it would not converge in a finite no. of steps, meaning the $\nabla_{\bar{w}} J$ will not become 0, we need to investigate the convexity of $J$.

$$J(\theta) = -\sum_{n=1}^{N} \left\{ y_n \log h_\theta(x_n) + (1-y_n) \log (1 - h_\theta(x_n)) \right\}$$

where $h_\theta(x)$ is the logistic function defined as,

$$h_\theta(x_n) = \frac{1}{1 + e^{-\theta^T x_n}}$$

Now, $J(\theta)$ can be written as,

$$J(\theta) = \sum_{n=1}^{N} - \left\{ y_n \log \left( \frac{h_\theta(x_n)}{1 - h_\theta(x_n)} \right) + \log (1 - h_\theta(x_n)) \right\}$$

Now, $\log \left( \frac{h_\theta(x_n)}{1 - h_\theta(x_n)} \right) = \log \left( \frac{\frac{1}{1 + e^{-\theta^T x_n}}}{\frac{e^{-\theta^T x_n}}{1 + e^{-\theta^T x_n}}} \right) = \log \left( e^{\theta^T x_n} \right)$

$$= \theta^T x_n$$

So, $J(\theta) = \sum_{n=1}^{N} - \left\{ y_n \theta^T x_n + \log (1 - h_\theta(x_n)) \right\}$ ———①

the 1st term is linear in $\theta$, so it is convex

let's take gradient of the $2^{nd}$ term,

$$\nabla_\theta \left[ -\log \left(1 - h_\theta(x)\right)\right] = -\nabla_\theta \left[\log \left(1 - \frac{1}{1+e^{-\theta^T x}}\right)\right]$$

$$= -\nabla_\theta \left[\log \frac{e^{-\theta^T x}}{1+e^{-\theta^T x}}\right]$$

$$= -\nabla_\theta \left[\log e^{-\theta^T x} - \log \left(1+e^{-\theta^T x}\right)\right]$$

$$= -\nabla_\theta \left[-\theta^T x - \log \left(1+e^{-\theta^T x}\right)\right]$$

$$= x + \nabla_\theta \left[\log \left(1+e^{-\theta^T x}\right)\right]$$

$$= x + \left(\frac{-e^{-\theta^T x}}{1+e^{-\theta^T x}}\right) x = x\left(1 - \frac{e^{-\theta^T x}}{1+e^{-\theta^T x}}\right)$$

$$= h_\theta(x) \, x$$

the $1^{st}$ term of ① was $-y_n \theta^T x_n$

gradient of $1^{st}$ term is $\nabla_\theta\left(-y_n \theta^T x_n\right) = -y_n x_n$

So, $\nabla_\theta J(\theta) = \sum_{n=1}^{N} \left\{ -y_n x_n + h_\theta(x_n^\circ) \, x_n\right\}$

## Proof that $J(\theta)$ is convex

From ①, $J(\theta) = \sum\limits_{n=}^{N} - \left\{ y_n \theta^T x_n + \log\left(1 - h_\theta(x_n)\right)\right\}$

the $1^{st}$ term is linear in $\theta$, so it is convex.

We have shown above that $\nabla_\theta\left(-\log\left(1 - h_\theta(x)\right)\right) = h_\theta(x) x$

the hessian of this $2^{nd}$ term is,

$$\nabla_\theta^2\left[-\log\left(1 - h_\theta(x)\right)\right] = \nabla_\theta\left[h_\theta(x)\, x\right]$$

$$= \nabla_\theta\left[\left(\frac{1}{1 + e^{-\theta^T x}}\right) x\right]$$

$$= \frac{1}{\left(1 + e^{-\theta^T x}\right)^2}\left(-e^{-\theta^T x}\right) x x^T$$

$$= \left(\frac{1}{1 + e^{-\theta^T x}}\right)\left(1 - \frac{1}{1 + e^{-\theta^T x}}\right) x x^T$$

$$= h_\theta(x)\left[1 - h_\theta(x)\right] x x^T$$

Now, for any $v \in \mathbb{R}^d$, we have,

$$v^T \nabla_\theta^2\left[-\log\left(1 - h_\theta(x)\right)\right] v = v^T\left[h_\theta(x)\left[1 - h_\theta(x)\right] x x^T\right] v$$

$$= h_\theta(x)\left[1 - h_\theta(x)\right] \|v^T x\|^2$$

we have $h_\theta(x) \in [0,1]$, so $[1 - h_\theta(x)] \in [0,1]$, so

$$r^T \nabla_\theta^2 \left[ -\log\left(1 - h_\theta(x)\right)\right] r = h_\theta(x)\left[1 - h_\theta(x)\right] \|w^T x\|^2 \geq 0$$

So the Hessian of $-\log\left(1 - h_\theta(x)\right)$ is positive semi-definite, so $-\log\left(1 - h_\theta(x)\right)$ is convex in $\theta$.

So, both the 1st and 2nd terms in $J(\theta)$ are convex, so, $J(\theta)$ is convex in $\theta$. Since $J(\theta)$ is convex, it would not get stuck at some other critical point (i.e. saddle point), rather there is a global minima with $w$ of infinite magnitude and the algo would keep chasing this global minima. Again, with proper $\alpha$, the algo would not converge. So, gradient descent would just push the loss down after every descent step. Since if $y_n = 1$, the gradient is +ve, so $w^{(k+1)}$ is $> w^{(k)}$ from ②. So, $\sigma\left(w^{T(k+1)} x_n\right)$ gets closer to 1 compared to $\sigma\left(w^{T(k)} x_n\right)$. So, the loss keeps going down with every descent step toward the unique global minima, without worrying about the descent being stuck at some finite critical point.

So, the algo would not converge in finite no. of steps. (Proved)

(ii) If we restrict $\|\bar{\omega}\|_2 \leq c_1$ and $|\omega_0| < c_2$ for some $c_1, c_2 > 0$, then the algo can not keep chasing the $\bar{\omega} \to \infty$, rather it is bounded by $\|\bar{\omega}\|_2 \leq c_1$, so it is like a forced convergence condition.

This non-convergence issue can be dealt with using $L^2$ regularization. The gradient of the resulting objective is

$$\nabla_{\bar{\omega}} \left\{ \sum_{n=1}^{N} y_n \log \varepsilon(\bar{\omega}^T x_n) + (1-y_n) \log (1 - \varepsilon(\bar{\omega}^T x_n)) - \frac{\lambda}{2} \bar{\omega}^T \bar{\omega} \right\}$$

$$= \sum_{n=1}^{N} x_n (y_n - \varepsilon(\bar{\omega}^T x_n)) - \lambda \bar{\omega}$$

We can adjust $\lambda$, scale $\lambda$ down by a factor of $N$. So the rule is:

$$\bar{\omega}^{(k+1)} \leftarrow \bar{\omega}^{(k)} + \alpha \left( x_n (y_n - \varepsilon((\bar{\omega}^{(k)})^T x_n)) - \frac{\lambda}{N} \bar{\omega}^{(k)} \right)$$

which can be rewritten as:

$$\bar{\omega}^{(k+1)} \leftarrow \left(1 - \alpha \frac{\lambda}{N}\right) \bar{\omega}^{(k)} + \alpha x_n (y_n - \varepsilon((\bar{\omega}^{(k)})^T x_n))$$

This is why $L^2$ is often referred as "weight decay".

(iii) This linear separability does not cause nonconvergence for other linear classifiers that we have studied. Since for linear regression classifier, if data is linearly separable, the algorithm tries to find the line with

$\vec{w}^T x + w_o^*$ such that all data are separated perfectly, but once that goal is achieved, it stops. There is no objective to chase after once the perfect classification is achieved. So nonconvergence is not a problem. Similarly, for perceptron, the goal is to just classify all the samples perfectly, once it is achieved, the algorithm stops. Again, for the SVM, since it is a separable case, with the maximum possible margin, SVM tries to find the decision boundary and once it is able to find it stops.

(b)(i) Since $x_j$ is misclassified by $w^{(k)}$,

$$y_j \, w^{T(k)} x_j < 0$$

Now,
$$y_j \left(w^{(k+1)}\right)^T x_j = y_j \left(w^{(k)} + y_j x_j\right)^T x_j$$
$$= y_j \left(w^{(k)}\right)^T x_j + \|y_j x_j\|_2^2$$

So, $y_j \left(w^{k+1}\right)^T x_j > y_j \left(w^{(k)}\right)^T x_j$ ———(A)

Since $y_j w^{T(k)} x_j < 0$ but we want to make $y_j \left(w^{(k)}\right)^T x_j > 0$,

So, we need to increase $y_j \left(w^i\right)^T x_j$, as $i$ increases. So, according to (A), the move from $w^{(k)}$ to $w^{(k+1)}$ is in right direction.

(Showed)

(ii) 1. Here, $\rho = \min_j y_j (w^*)^T x_j$, since every $x_j$ is correctly classified by $w^*$, we have for all $n = 1$ to $N$, $y_n (w^*)^T x_n > 0$

So, $\rho > 0$, now, $(w^{(k)})^T w^* = \left[ (w^{(k-1)})^T + y_j x_j^T \right] w^*$

$= (w^{(k-1)})^T w^* + y_j (w^*)^T x_j$

$\geq (w^{(k-1)})^T w^* + \rho$ ——©

Now, $\overset{to}{\underset{\Lambda}{}}$ prove that $(w^{(k)})^T w^* \geq k\rho$, let us go by induction. If $k=0$, we get, $0 . w^* \geq 0$. If the thesis is true for $(k-1)$, let us prove it for $k$. Using ©, we have,

$(w^{(k)})^T w^* \geq (w^{(k-1)})^T w^* + \rho \geq (k-1)\rho + \rho = k\rho$

2. $\| w^{(k)} \|^2 = \| w^{(k-1)} + y_j x_j \|^2 = \| w^{(k-1)} \|^2 + \| y_j x_j \|^2 + 2 y_j (w^{(k-1)})^T x_j$

$\leq \| w^{(k-1)} \|^2 + \| y_j x_j \|^2$

$\leq \| w^{(k-1)} \|^2 + \| x_j \|^2 \; ; \; [\because \| y_j \|^2 = 1]$

Since $x_j$ is misclassified by $w^{(k-1)}$. So $2 y_j (w^{(k-1)})^T x_j < 0$

Now, let us prove by induction that $\| w^{(k)} \|^2 \leq k R^2$ where $R = \max_j \| x_j \|_2$. If $k=0$, we have $0 \leq 0 . R^2$. If the claim is true for $(k-1)$, let us prove it for $k$. So,

$\| w^{(k)} \|^2 \leq \| w^{(k-1)} \|^2 + \| x_j \|^2 \leq (k-1) R^2 + R^2 = k R^2$

(Showed)

3. Using 1 & 2, we get, $\dfrac{(w^{(k)})^T w^*}{\| w^{(k)} \|} > \dfrac{k\rho}{\sqrt{k} R} = \sqrt{k} \dfrac{\rho}{R}$ ——③

Since $\| w^{(k)} \| < \sqrt{k} R$ (from (2)) $\Rightarrow \dfrac{1}{\| w^{(k)} \|} > \dfrac{1}{\sqrt{k} R}$

So from ③, $k \leq \frac{R^2}{\rho^2} \frac{\left( (w^{(k)})^T w^* \right)^2}{\|w^{(k)}\|^2} = \frac{R^2}{\rho^2} \frac{\left( (w^{(k)})^T w^* \right)^2}{\|w^{(k)}\|^2 \|w^*\|^2} \|w^*\|^2$

Now using Cauchy-Schwarz inequality, $\frac{\left( (w^{(k)})^T w^* \right)^2}{\|w^{(k)}\|^2 \|w^*\|^2} \leq 1$

So we get, $k \leq \frac{R^2 \|w^*\|^2}{\rho^2}$ 

(Showed)

(c) Here the primal problem of soft margin SVM is

$$\underset{\theta, \varepsilon}{\arg\min} \ \frac{1}{2}\left( \|\omega\|_2^2 + c\|\varepsilon\|_2^2 \right)$$

subject to $y_j \left( w^T x_j + w_o \right) \geq 1 - \varepsilon_j$, $\varepsilon_j \geq 0$, $j = 1, ..., N$

(i) the Lagrangian of the problem is.

$$\mathcal{L}\left( w, w_o, \varepsilon, \lambda \right) = \frac{1}{2}\left( \|\omega\|^2 + c\|\varepsilon\|^2 \right) + \sum_{j=1}^{N} \lambda_j \left[ 1 - \varepsilon_j - y_j \left( w^T x_j + w_o \right) \right]$$

let us minimize over $\left( w, w_c, \varepsilon \right)$:

$$\nabla_w \mathcal{L}\left( w, w_c, \varepsilon, \lambda \right) = w^* - \sum_{j=1}^{N} \lambda_j y_j x_j = 0$$

$$\Rightarrow w^* = \sum_{j=1}^{N} \lambda_j y_j x_j$$

$$\nabla_{w_o} \mathcal{L}\left( w, w_c, \varepsilon, \lambda \right) = 0 \Rightarrow \sum_{j=1}^{N} \lambda_j y_j = 0$$

$$\nabla_{\varepsilon} \mathcal{L}\left( w, w_c, \varepsilon, \lambda \right) = 0 \Rightarrow c \varepsilon_j^* = \lambda_j \ , \ \text{for} \ \forall j$$

Now, since for all $j$, $\lambda_j \geq 0$, so $\varepsilon_j = \frac{1}{c}\lambda_j$ (with $c > 0$).

So, the optimization solution is automatically

$\varepsilon_j \geq 0, \forall j$, so the constraint $\varepsilon_j > 0, \forall j$ can be

removed without affecting the solution.

(ii) We proved these in (i) taking derivative of

the Lagrangian w.r.t. $w$, $w_o$ & $\varepsilon$ respectively & setting $= 0$.

(iii) So, the Lagrangian function becomes,

$$\mathcal{L}(w^*, w_o^*, \varepsilon^*, \lambda) = \frac{1}{2}\left(\|\sum_{j=1}^{N} \lambda_j \cdot y_j \cdot x_j\|^2 + \frac{\sum \lambda_j^2}{c}\right)$$
$$+ \sum_{j=1}^{N} \lambda_j \left[1 - \varepsilon_j^* - y_j \left(w^{*T} x_j + w_o^*\right)\right] \quad\text{————(A)}$$

the 2nd term becomes,

Now, $\sum_{j=1}^{N} \lambda_j - \sum_{j=1}^{N} \lambda_j \varepsilon_j^* - \sum_{j=1, J}^{N}\sum_{K=1}^{N}\lambda_j y_j \sum_{j=1}^{N} \lambda_j y_j \left[\sum_{k=1}^{N}\lambda_k y_k x_k\right]^T x_j$

$$- \left(\sum_{j=1}^{N} \lambda_j y_j\right) w_o^*$$

$$= \sum_{j=1}^{N} \lambda_j - \sum_{j=1}^{N} \frac{\lambda_j \cdot \lambda_j}{c} - \sum_{j=1}^{N}\sum_{k=1}^{N} \lambda_j \lambda_k y_j y_k x_j^T x_k$$

the 1st term of (A) is, $\frac{1}{2}\sum_{j=1}^{N}\sum_{k=1}^{N}\lambda_j\lambda_k y_j y_k x_j^T x_k + \frac{1}{2}\sum\frac{\lambda_j^2}{c}$

So, from (A), we have,

the convex dual problem is,

$$\max_{\lambda \geq 0} \left\{ \sum_{j=1}^{N} \lambda_j - \frac{1}{2} \sum_{j=1}^{N} \sum_{k=1}^{N} \lambda_j \lambda_k y_j y_k x_j^T x_k - \frac{1}{2} \sum_{j=1}^{N} \frac{\lambda_j^2}{C} \right\}$$

subject to $\sum_{j=1}^{N} \lambda_j y_j = 0$

(d) (i) Here the problem is,

$$\min_{w, w_o} \frac{1}{2} \|w\|_2^2 \text{ subject to } y_j(w^T x_j + w_o) \geq \gamma \quad j=1,\ldots N$$

So the Lagrangian becomes.

$$\mathcal{L}(w, w_o, \lambda) = \frac{1}{2} \|w\|_2^2 + \sum_{j=1}^{N} \lambda_j \left[ \gamma - y_j(w^T x_j + w_o) \right]$$

So, $\nabla_w \mathcal{L} = 0 \Rightarrow w^* = \gamma \sum_j \lambda_j y_j x_j$ , $\nabla_{w_o} \mathcal{L} = 0 \Rightarrow \sum_j \lambda_j y_j = 0$

Now, if we proceed as shown before (and for the usual Hard margin-SVM), we can get the solution from the optimization as:

$$\frac{w^*}{\gamma} = \gamma \sum_{j \in V} \lambda_j^* y_j x_j, \quad V \text{ is set of support vectors:}$$
$$\lambda_j > 0$$

Now, $w_o^* = -\frac{(x^+ + x^-)^T w^*}{2} \cdot \gamma = \gamma w_o^* \left[ \text{where } w_o^* = -\frac{(x^+ + x^-)^T w^*}{2} \right.$
$\left. \text{ is the optimal } w_o \right.$

Since $\frac{w^*}{\gamma} = \gamma \sum_{j \in V} \lambda_j^* y_j x_j = \gamma w^*$ for $y_j(w^T x_j + w_o) \geq 1]$

$\underbrace{\phantom{XXXXXXXXX}}_{w^* \text{ if constraint}}$
$\text{was } y_j(w^T x_j + w_o) \geq 1$

So, with $y_j(w^Tx_j + w_o) \geq \gamma$, the decision boundary is,

$$\gamma w^{*T}x + \hat{w}\gamma w_o^* = 0$$

$$\Rightarrow \gamma(\hat{w}^{*T}x + w_o^*) = 0 \Rightarrow w^{*T}x + w_o^* = 0$$

So, a scalar multiple of $\bar{\theta}$ does not change the decision boundary.

(ii) Here, let's assume $w^Tx_1 + w_o = 1$, $w^Tx_2 + w_o = -1$

So, $\alpha(w, w_o, r_1, r_2) = \frac{1}{2}\|w\|_2^2 + r_1[1 - (w^Tx_1 + w_o)] + r_2[1 + w^Tx_2 + w_o]$

$\nabla_w \alpha = 0 \Rightarrow w = r_1 x_1 - r_2 x_2$ ; $\nabla_{w_o}\alpha = 0 \Rightarrow r_1 = r_2$

$\Rightarrow w = r_1(x_1 - x_2)$, so let's take $r_1 = r_2 = r$

So, the lagrangian becomes,

$\frac{1}{2}\|w\|^2 + 2r + (-r)w^T(x_1 - x_2) = \frac{1}{2}\|w\|^2 + 2r - \overbrace{w^T w}^{\|w\|^2}$

$= -\frac{1}{2}\|w\|^2 + 2r$

$= -\frac{1}{2}r^2\|x_1 - x_2\|^2 + 2r$

So, the optimization becomes,

$$\max \quad -\frac{1}{2}r^2\|x_1 - x_2\|^2 + 2r$$
$$r \geq 0$$

taking 1st derivative, we get, $-r\|x_1 - x_2\|^2 + 2 = 0$

$$\Rightarrow r = \frac{2}{\|x_1 - x_2\|^2}$$

$$\therefore w^* = \frac{2(x_1 - x_2)}{\|x_1 - x_2\|^2} \quad , \quad w_o^* = \frac{-(x_1 + x_2)^T(x_1 - x_2)}{\|x_1 - x_2\|^2}$$

So, even with 2 data points (one from each class) the decision boundary is uniquely determined.

(Showed)