

Elwyn CW2
104973012

M46 HW2
1.2) $y = \vec{w}^T \vec{x} + b$
 $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow y \geq 0$
 $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow y \geq 0$
 $x = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \rightarrow y \geq 0$
 $x = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \rightarrow y \geq 0$



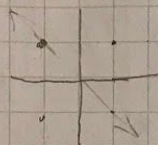
$y = [w_1, w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b$
 choose $b=0, w_1=1, w_2=+1$

This is not unique:

$b=-1, w_1=1, w_2=2 \cdot y = [1, 2]^T \vec{x} + 1$

b.) XOR

$x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow y \neq 0$
 $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow y \neq 0$
 $x = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \rightarrow y \neq 0$
 $x = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \rightarrow y \neq 0$



$y = [w_1, w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b$
 $w_1 x_1 + w_2 x_2 + b > 0$
 $w_1 = +1, w_2 = -1, b = 0$

there is no such linear model

(can't separate w/ a single line)

2.) $\frac{1}{2w} \frac{1}{1+e^{-w^T x}} = \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$

$\beta = w^T x \Rightarrow \sigma(w^T x)(1-\sigma(w^T x))$

$y_n \log h(w, x_n) + (1-y_n) \log (1-h(w, x_n))$
 $y_n \log \left(\frac{1}{1+e^{-w^T x}} \right) + (1-y_n) \log \left(\frac{1}{1+e^{w^T x}} \right)$
 $= y_n \log \left(\frac{1}{1+e^{-w^T x}} \right) + (1-y_n) \log \left(\frac{1}{1+e^{w^T x}} \right)$
 $= y_n \log \left(\frac{1}{1+e^{-w^T x}} \right) + (1-y_n) \log \left(\frac{1}{1+e^{w^T x}} \right)$
 $= y_n \log \left(\frac{1}{1+e^{-w^T x}} \right) + (1-y_n) \log \left(\frac{1}{1+e^{w^T x}} \right)$

$= y_n - y_n \frac{w^T x}{1+e^{-w^T x}} + (-\sigma(w^T x) + y_n \frac{w^T x}{1+e^{w^T x}})$
 $= y_n - \sigma(w^T x)$

$\frac{\partial J}{\partial w_j} = y_n - \frac{1}{1+e^{w^T x}}$

3a) $\frac{\partial J}{\partial w_1} = \sigma_n(2(w_0 + w_1 x_{n,1} - y_n))(x_{n,1})$

$\frac{\partial J}{\partial w_0} = \sigma_n(2(w_0 + w_1 x_{n,1} - y_n))(1)$

$\nabla J = (\sigma_n(2(w_0 + w_1 x_{n,1} - y_n))(x_{n,1}), \sigma_n(2(w_0 + w_1 x_{n,1} - y_n)))(1)$
 $= \sum_n \sigma_n(2(w_0 + w_1 x_{n,1} - y_n)) \begin{bmatrix} x_{n,1} \\ 1 \end{bmatrix}$

$$3b) \sum_{n=1}^N \alpha_n (w_0 + w_1 x_{n,1} - y_n) [x_{n,1}] = [0]$$

$$\sum_{n=1}^N \alpha_n (w_0 + w_1 x_{n,1} - y_n) = 0$$

$$\textcircled{1} \quad \begin{aligned} \sum \alpha_n w_0 + \sum \alpha_n w_1 x_{n,1} - \sum \alpha_n y_n &= 0 \\ \sum \alpha_n w_0 + \alpha_n w_1 x_{n,1} &= \sum \alpha_n y_n \end{aligned}$$

$$\sum \alpha_n (w_0 + w_1 x_{n,1} - y_n) x_{n,1} = 0$$

$$\sum \alpha_n x_{n,1} w_0 + \sum \alpha_n x_{n,1} x_{n,1} w_1 - \sum \alpha_n x_{n,1} y_n = 0$$

$$\textcircled{2} \quad \sum \alpha_n x_{n,1} w_0 + \sum \alpha_n x_{n,1} x_{n,1} w_1 = \sum \alpha_n x_{n,1} y_n$$

matrix:

$$\begin{bmatrix} \sum_{n=1}^N \alpha_n & \sum_{n=1}^N \alpha_n x_{n,1} \\ \sum_{n=1}^N \alpha_n x_{n,1} & \sum_{n=1}^N \alpha_n x_{n,1} x_{n,1} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \sum_{n=1}^N \alpha_n y_n \\ \sum_{n=1}^N \alpha_n y_n x_{n,1} \end{bmatrix}$$

solve using

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \sum_{n=1}^N \alpha_n & \sum_{n=1}^N \alpha_n x_{n,1} \\ \sum_{n=1}^N \alpha_n x_{n,1} & \sum_{n=1}^N \alpha_n x_{n,1} x_{n,1} \end{bmatrix}^{-1} \begin{bmatrix} \sum_{n=1}^N \alpha_n y_n \\ \sum_{n=1}^N \alpha_n y_n x_{n,1} \end{bmatrix}$$

given:

4a) D is linearly separable
Then: $y_i = \begin{cases} 1 & \text{if } w^T x_i + b \geq 0 \\ -1 & \text{if } w^T x_i + b < 0 \end{cases}$

Define subplane $w^T x_i + b$

where $y_i \in \{-1, 1\}$

then, if $b=0$, the hyper plane lies on the data point (x_i, y_i)

Then we shift the hyper plane

by ϵ so that $w^T x_i + \epsilon$ defines our new hyper plane

and $y_i (w^T x_i + \epsilon) = 1 - \delta = \epsilon$

now we normalize by dividing by ϵ

$$\frac{y_i (w^T x_i + \epsilon)}{\epsilon} = 1$$

$$\text{and } y_i \left(\frac{w^T x_i}{\epsilon} + 1 \right) = 1 \geq 1 - \delta \quad w/ \delta = 0$$

thus, we have an optimal solution w/ $\delta = 0$

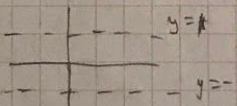
4b) Given: there is an optimal solution w/ $\delta = 0$

Then we know there is some \vec{w} and b that

$y_i (w^T x_i + b) \geq 1$, then we know that for any quantity $(w^T x_i + b) \leq 0$ (negative), y_i is also negative and for any $(w^T x_i + b) > 0$, y_i is also positive.

Thus, D is linearly separable since y_i can be written

$$y_i = \begin{cases} 1 & \text{if } w^T x_i + b \geq 0 \\ -1 & \text{if } w^T x_i + b < 0 \end{cases}$$



$$\text{taking } \vec{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \epsilon = -$$

$$x_i = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \rightarrow y_i \left(\frac{1 \cdot 1 + 1 \cdot (-1)}{0.1} \right) = 1$$

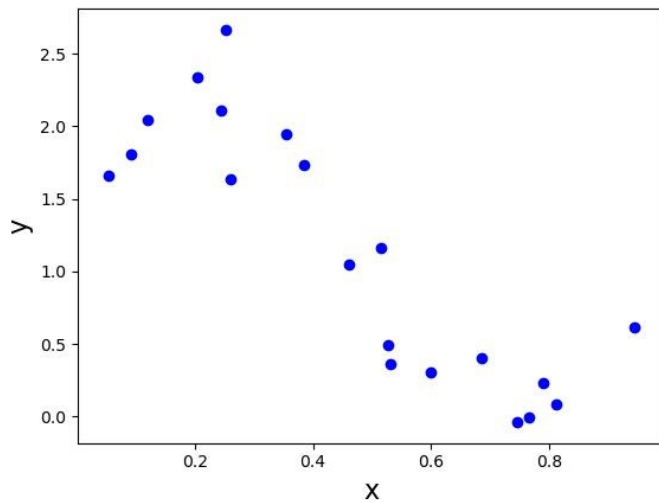
4c.) If we have a hyperplane that satisfies condition
a w/ $\delta > 0$, then we can say that there also exists
a hyper plane w/ $\delta = 0$, through normalization
and thus that the data is linearly separable.

4d.) optimal solution independent of D is
 $\vec{w} = \vec{0}$, $\theta = 0$, $\delta = 0$.

This does not give us a hyper plane at all, so this is
not a good formulation.

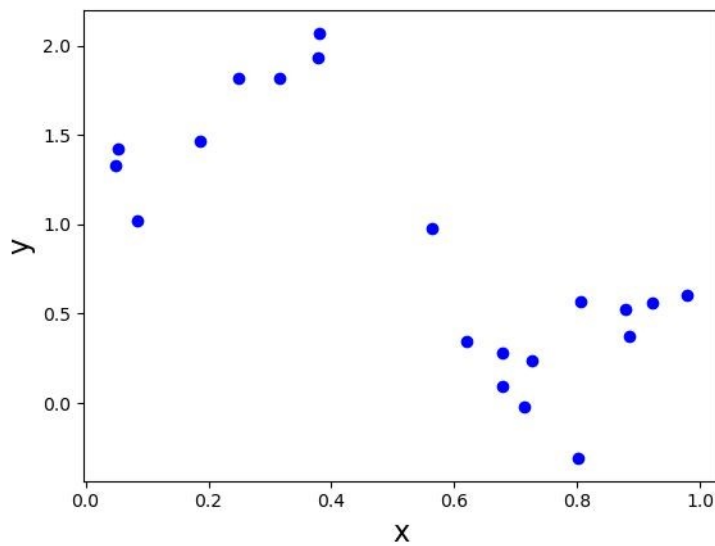
4e.) Given D has only 2 points, we know it is linearly
separable. possible optimal solutions is $w_1 + w_2 \dots w_{n/2} + 10$

5a.) Train data plot:



We can see that the train data looks somewhat linear, so linear regression seems like it would be fairly effective in predicting the data.

Test data plot:



This data looks somewhat linear, although less so than the train data plot. It would be fair to say that linear regression would be less effective in predicting the data.

5d.)

Coefficient	Iterations to Convergence	Cost	Coefficients
0.01	765	3.9125	(2.4464, -2.8163)
0.001	7021	3.9125	(2.4464, -2.8163)
0.0001	10000 (does not converge)	4.0864	(2.27045, -2.46065)
0.0407	10000 (does not converge)	2.7109e39	(-9.4047e39, -4.6523e18)

5e.) The closed form solution is (2.44640709, -2.81635359). The closed form solution runs almost 3x as fast as the gradient descent. The coefficients are nearly the same for when gradient descent converges. The costs are the same since they both

5f.) It takes the algorithm 1679 iterations to converge with the proposed learning rate.

5h.) We might prefer RMSE as a metric over the cost because while an overfit model may seem accurate on the training data and the cost will continuously decrease, the RMSE will reach a minimum then as the model becomes more overfit, increase again. Thus minimizing RMSE will allow one to minimize cost while not overfitting the model..