

zID : z5224151

Name: ZANNING WANG

Q1:

(a) According Naive Bayes,
we have $P(+|X^*) = \frac{P(X^*|+)}{P(X^*)} \cdot P(+)$

Obviously, $P(+)=P(-)=\frac{4}{8}=\frac{1}{2}$ and $X^*=(1, 0, 0, 2)$

$\therefore \Theta = (\frac{1}{3}, \frac{3}{3}, \frac{9}{3}, \frac{6}{3}) = (0.2174, 0.1304, 0.3913, 0.2609)$

$P(X^*|+) = 3! \times \frac{(0.2174)^1}{1!} \times \frac{(0.2609)^2}{2!} = 0.0443$

$P(X^*)$ can be calculated by the $P(X^*(x_1, x_2, x_3, x_4))$ given by.
therefore. $P(X^*) = \frac{(1+2)!}{1! \times 2!} \times \left(\frac{16}{39}\right)^1 \times \left(\frac{8}{39}\right)^2 = 0.0517$

$P(+|X^*) = \frac{0.0443 \times 0.5}{0.0517} = 0.4284$

(b) $P(-|X^*) = \frac{P(X^*|-) \cdot P(-)}{P(X^*)} \quad X^*=(1, 0, 0, 2)$

According to the figure. $\Theta = (1, 3, 0, 2)$

By using smooth, add 1 to each column.

Smooth $\tilde{\Theta} = (\frac{1^2}{20}, \frac{4}{20}, \frac{1}{20}, \frac{3}{20})$

$P(X^*|-) = 3! \times \frac{(0.6)^1}{1!} \times \frac{(0.15)^2}{2!} = 0.0405$

$P(X^*) = \frac{(1+2)!}{1! \times 2!} \times \left(\frac{18}{47}\right)^1 \times \left(\frac{10}{47}\right)^2 = 0.0519$

$P(-|X^*) = \frac{0.0405 \times 0.5}{0.0519} = 0.3902$

(c) Similar as last two questions, but we should translate.
 x to bit vector: $(1, 0, 0, 1)$

$$P(+|x_*) = \frac{P(x_*|+)}{P(x_*|-)} \quad P(+) = P(-) = \frac{4}{8} = \frac{1}{2}$$

$$P(x_*|+) = \frac{1}{2} \times (1 - \frac{1}{4}) \times (1 - \frac{3}{4}) \times \frac{1}{2} = 0.0468$$

$$P(x_*|-) = \frac{1}{2} \times (1 - \frac{1}{4}) \times (1 - \frac{3}{4}) \times \frac{1}{2} = 0.1464$$

$$P(+|x_*) = \frac{0.0468 \times 0.5}{0.1464} = \underline{\underline{0.1598}}$$

(d) $P(x_*|-) = \frac{4}{6} \times (1 - \frac{2}{6}) \times (1 - \frac{1}{6}) \times \frac{3}{6} = 0.1851$

In order to smooth, we should add $(0, 0, 0)$ and $(1, 1, 1)$ to both
 \oplus, \ominus classes.

therefore:

$$P(x_*|-) = \frac{7}{12} \times (1 - \frac{4}{7}) \times (1 - \frac{4}{7}) \times \frac{6}{12} = 0.1296$$

$$P(+|x_*) = \frac{0.1851 \times 0.5}{0.1296} = \underline{\underline{0.7141}}$$

Q2:

(a)

(a) Because $\hat{y}_i = w_0 + w_1 x_i$ $\hat{y}_i = w^T x_i$
therefore: $w^T x = w_0 + w_1 x_i$

$$L(y_i, \hat{y}_i) = \left[\sqrt{\frac{1}{c^2} (y_i - w_0 - w_1 x_i)^2 + 1} - 1 \right]$$

$$\begin{aligned}\frac{\partial L(y_i, \hat{y}_i)}{\partial w_0} &= \frac{\partial}{\partial w_0} \sqrt{\frac{1}{c^2} (y_i - w_0 - w_1 x_i)^2 + 1} \\ &= \frac{1}{2} \left[\sqrt{\frac{1}{c^2} (y_i - w_0 - w_1 x_i)^2 + 1} \right]^{-\frac{1}{2}} \times \frac{1}{c^2} \times 2 \cdot (y_i - w_0 - w_1 x_i) \\ &= \frac{-2(y_i - w_0 - w_1 x_i)}{c^2 \left[\sqrt{\frac{1}{c^2} (y_i - w_0 - w_1 x_i)^2 + 1} \right]^{\frac{1}{2}}}\end{aligned}$$

$$\begin{aligned}\frac{\partial L(y_i, \hat{y}_i)}{\partial w_1} &= \frac{\partial}{\partial w_1} \sqrt{\frac{1}{c^2} (y_i - w_0 - w_1 x_i)^2 + 1} \\ &= \frac{1}{2} \left[\sqrt{\frac{1}{c^2} (y_i - w_0 - w_1 x_i)^2 + 1} \right]^{-\frac{1}{2}} \times \frac{2}{c^2} \cdot (y_i - w_0 - w_1 x_i) \cdot (-x_i) \\ &= \frac{-2(y_i - w_0 - w_1 x_i) \cdot (-x_i)}{c^2 \left[\sqrt{\frac{1}{c^2} (y_i - w_0 - w_1 x_i)^2 + 1} \right]^{\frac{1}{2}}}\end{aligned}$$

(b)

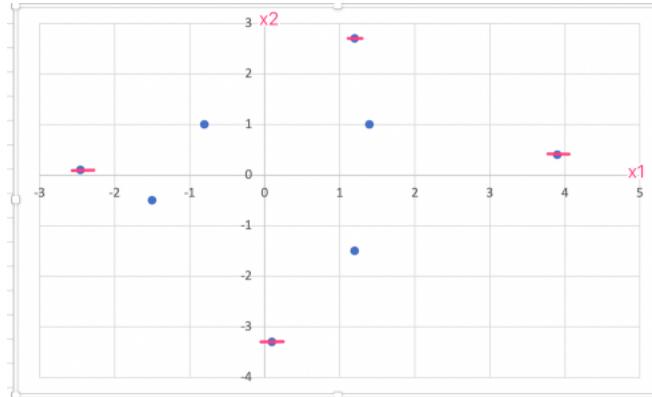
b) while $\frac{\partial L(y_i, \hat{y}_i)}{\partial w_0} \neq 0$ and $\frac{\partial L(y_i, \hat{y}_i)}{\partial w_1} \neq 0$

$$w_0 = w_0^{(t+1)}$$

$$w_1 = w_1^{(t+1)}$$

Q3:

(1)



The data was plotted above, obviously, the +1 and -1 cannot be classified by linear regression.

Suppose we set $m=0, d=2$,

$$\text{Therefore } k(x,y) = (x^T y)^2$$

I am not sure m and d right or not, just to calculate the question below

(2)

By selecting the vector (w, x) , the transformed data is linearly separable.

(3)

The code shows below, this still have some bugs, the code can not converge when restart iterations. The main idea of my model is to multi vector to make data linearly separate and compute wTx and $ywTx$, use counter to check if all the cases have been converged, if counter == 8 which means all data converged.

```
np_x0 = np.asarray(x0)
np_x1 = np.asarray(x1)
np_x2 = np.asarray(x2)
np_y = np.asarray(y)
w_start = [1.00, 1.00, 1.00]
learning_rate = 0.2
counter = 0

def check_positive(x0, x1, x2, y):
    print("-----")
    global w_start, counter
    wTx = (w_start[0] * x0) + (w_start[1] * x1) + (w_start[2] * x2)
    ywTx = wTx * y

    if ywTx >= 0:
        counter += 1
    else:
        w_start = [w_start[0] - y * (learning_rate * x0), w_start[1] - y * (learning_rate * x1), w_start[2] - y * (learning_rate * x2)]
```

print("y(x_i) ~~~~~~ ywTx(x_i)w")

```
counter2 = 0
while True:
    counter = 0
    for i in range(8):
        check_positive(np_x0[i], np_x1[i], np_x2[i], np_y[i])
        list_w = [w_start[0], w_start[1], w_start[2]]
```

if list_w:

```
    b = "r_i > 0"
    print(list_w, b)
```

if counter == 8:

```
    break
```

```

 $\phi(x_i) \quad y_i \phi^T(x_i) w^*$ 
-----
[1.0, 1.0, 1.0]  $r_i > 0$ 
-----
[1.2, 1.78, 1.08]  $r_i > 0$ 
-----
[1.2, 1.78, 1.08]  $r_i > 0$ 
-----
[1.2, 1.78, 1.08]  $r_i > 0$ 
-----
[1.4, 2.02, 1.62]  $r_i > 0$ 
-----
[1.4, 2.02, 1.62]  $r_i > 0$ 
-----
[1.2, 2.3200000000000003, 1.7200000000000002]  $r_i > 0$ 
-----
[1.2, 2.3200000000000003, 1.7200000000000002]  $r_i > 0$ 

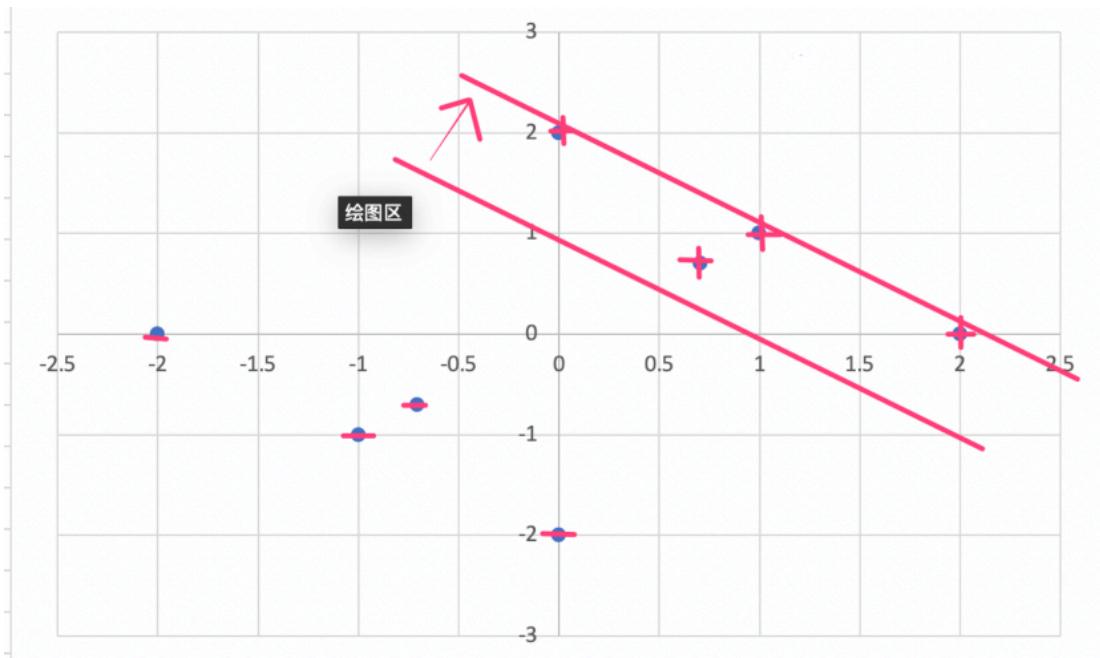
```

(d)

by while $\frac{\partial L(y_i, \hat{y}_i)}{\partial w_0}$ and $\frac{\partial L(y_i, \hat{y}_i)}{\partial w_1} \neq 0$

 $w_0 = w_0^{(t+1)}$
 $w_1 = w_1^{(t+1)}$

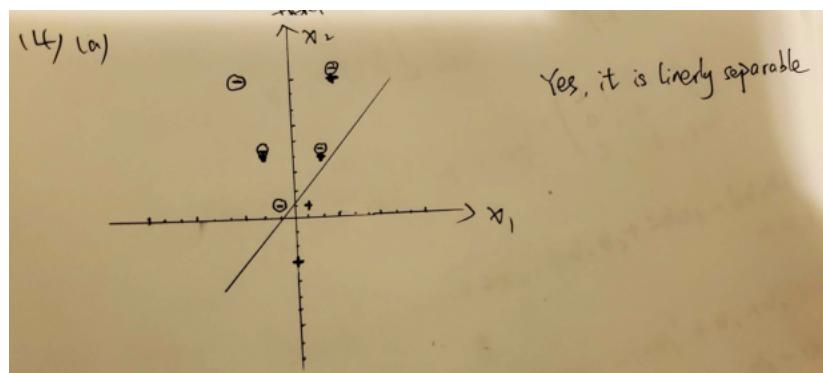
(e)



I agree with the claim, by plotting the data, it is obviously that when we move the line by applying the learning rate = 1, the most mistakes it can made is show as figure is 4. Which means make 5 mistakes during training will never happen.

Q4:

(a)



(b)

Q4 (b) According to the graph, we plot last problem, we can calculate the (d_1, \dots, d_3) by only selecting $(-1, 1) (1, 1) (2, 4)$, the ~~line~~ below.

x_1	x_2	y
-1	1	-1
1	1	1
2	4	1

therefore we can get $\mathbf{x} = \begin{pmatrix} -1 \\ 1 \\ 2 \end{pmatrix}, \mathbf{y} = \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}$

$$\mathbf{x}' = \mathbf{x}\mathbf{y} = \begin{pmatrix} 1 & -1 \\ 1 & 1 \\ 2 & -4 \end{pmatrix} \quad \mathbf{x}'^T = \begin{pmatrix} 1 & 1 & -2 \\ -1 & 1 & -4 \end{pmatrix}$$

therefore $\mathbf{x}' \cdot \mathbf{x}'^T = \begin{pmatrix} 1 & -1 \\ 1 & 1 \\ 2 & -4 \end{pmatrix} \times \begin{pmatrix} 1 & 1 & -2 \\ -1 & 1 & -4 \end{pmatrix} = \begin{pmatrix} 2 & 0 & 2 \\ 0 & 2 & -6 \\ 2 & -6 & 20 \end{pmatrix}$

(therefore,

$$= \underset{\mathbf{d}}{\operatorname{argmax}} -\frac{1}{2} (\sum d_i^2 + 2d_1d_3 + 2d_2^2 - 2d_2d_3 + 2d_1d_3 - 2d_2d_3 + 2d_3^2) + d_1 + d_2 + d_3$$

$$\leq \underset{\mathbf{d}}{\operatorname{argmax}} -d_1^2 - 2d_1d_3 - d_2^2 + 2d_2d_3 - (d_3^2 + d_1 + d_2 + d_3)$$

Because $\sum d_i \cdot y_i = 0$ therefore $-d_1 + d_2 - d_3 = 0$
 $d_2 = d_1 + d_3$

~~$\frac{\partial L}{\partial d_1}$~~

$$L = \underset{\mathbf{d}}{\operatorname{argmax}} -d_1^2 - 2d_1d_3 - (d_1 + d_3)^2 + (d_1 + d_3) \cdot d_3 - (d_3^2 + 2d_1 + 2d_3)$$

$$\frac{\partial L}{\partial d_1} = 2d_3 - 4d_1 + 2 \quad \frac{\partial L}{\partial d_2} = 2d_1 - 10d_3 + 2$$

$$\frac{\partial L}{\partial d_3} = 0 \quad \frac{\partial L}{\partial d_2} = 0$$

$$\begin{cases} 2d_3 - 4d_1 + 2 = 0 \\ 2d_1 - 10d_3 + 2 = 0 \\ d_2 = d_1 + d_3 \end{cases}$$

$$\begin{cases} d_1 = \frac{2}{3} \\ d_2 = 1 \\ d_3 = \frac{1}{3} \end{cases}$$

therefore
 $d_1 \text{ and } d_3$
 $= (0, 0, \frac{2}{3}, 0, 1, \frac{1}{3}, 0)$

(c)

19) as we know:

$$w = \frac{1}{N} \sum_{i=1}^N \alpha_i \cdot x_i \cdot y_i = \frac{1}{3} \times (-1) \times \begin{pmatrix} -1 \\ 1 \end{pmatrix} + 1 \times 1 \times \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \frac{1}{3} \times (-1) \times \begin{pmatrix} 2 \\ -4 \end{pmatrix}$$
$$= \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

Because $y_1(wx_1 - t) = 1$

$$-1 \times \left[\begin{pmatrix} 1 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} -1 & 1 \end{pmatrix} - t \right] = 1$$

$$t = -1$$

$$\text{Margin}_w = \frac{1}{\|w\|} = \frac{1}{\sqrt{1+1}} = \frac{\sqrt{2}}{2}$$

(d)

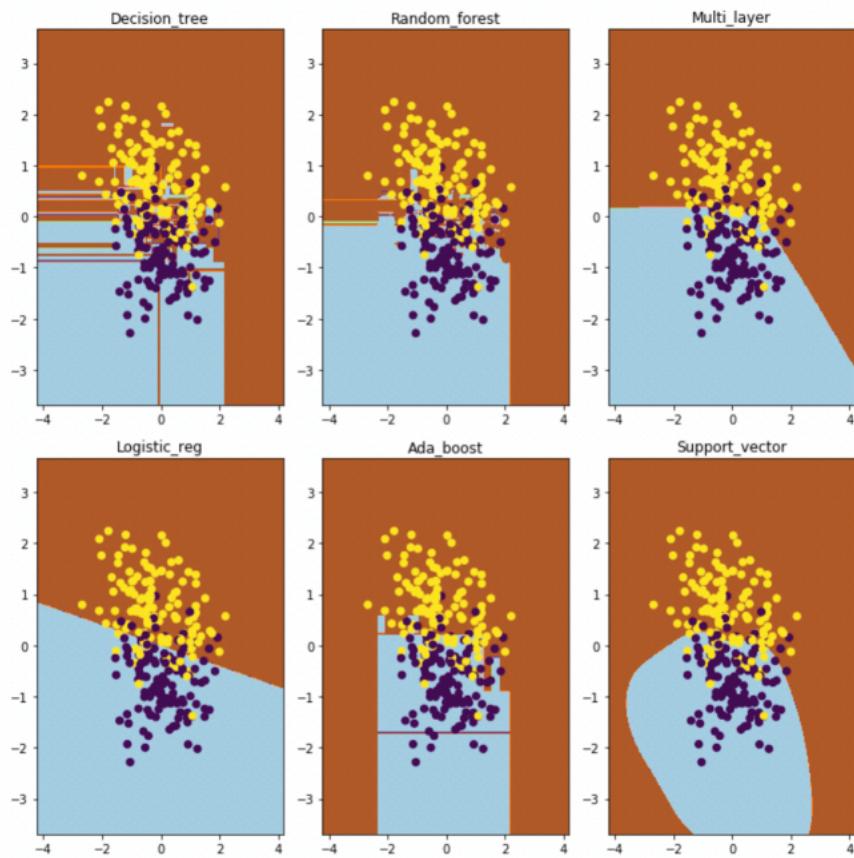
As we know, linear classifier uses a linear combination of features to make classification decisions, it can only distinguish between positive and negative samples. Most of non-linear function, the positive and negative samples mix together, therefore Linear classifiers are unable to represent non-linear functions.

(e)

Because linearity is inseparable in low-dimensional space, we can improve the probability of linearly separable by mapping it to high-dimensional space. Especially, it can solve the possibility of each result very close to each other, Therefore, we need to map indivisible data from low-dimensional space to high-latitude space to increase accuracy.

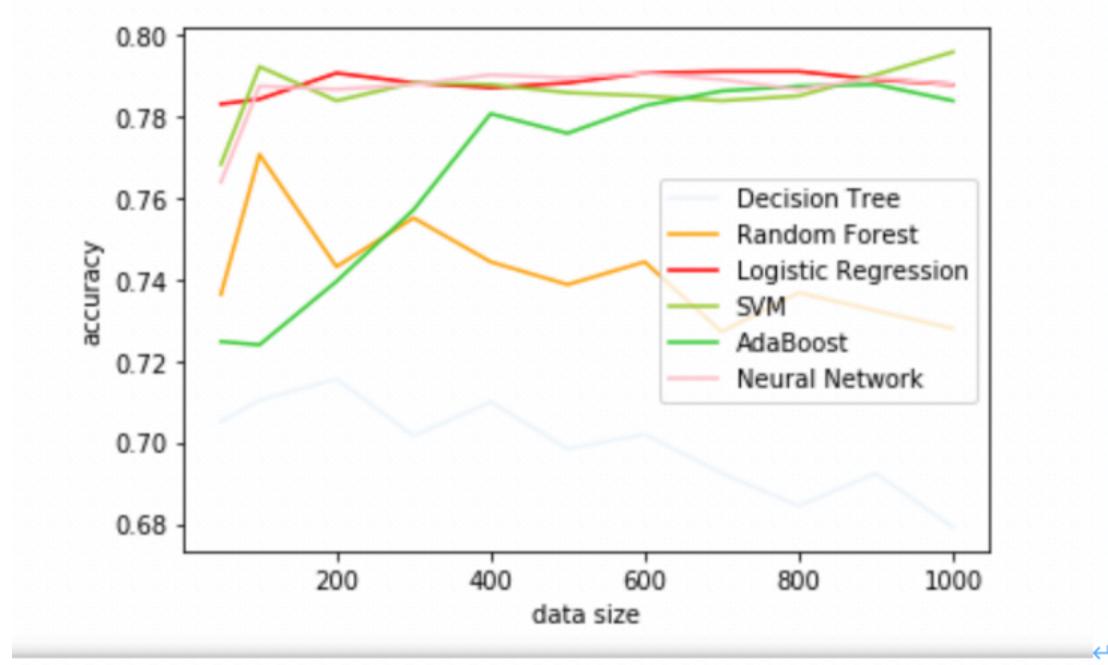
Q5:

(a)



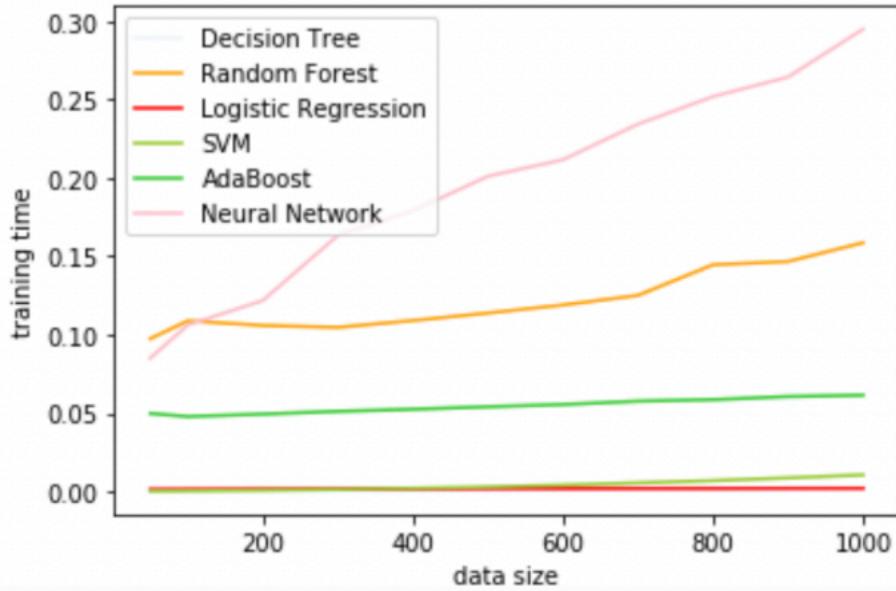
The dataset and decision boundaries show above.

(b)



The accuracy of each model show above, Obviously, most of model except decision tree's accuracy is increasing with the larger data size. After 400, random forest and Decision tree keep decreasing. At the size of 600+, logistic regression, SVM, Neural network, and Adaboost keep stability at rate of 0.75, much better than Decision tree and random forest. (the color shown above may not right, because I am a little bit color blindness, sorry for that==)

(c)



According the figure shown above, Random tree and neural network increasing significantly, in contract, Adaboost, logistic regression and SVM keep at a lower time. According to the accuracy mentioned in the last problem, the SVM and Adaboost is much better to use binary classification at some certain situation.