

Yu, Zhang (25238743)

Q: new document  $x^*$   $A=1$   $B=0$   $C=0$   $D=2$

(a) Multinomial:

$$E^+: A:5 \quad B:3 \quad C:9 \quad D:6 \quad \text{All: } 23$$

$$E^-: A:11 \quad B:3 \quad C:0 \quad D:2 \quad \text{All: } 16$$

$$\therefore P\theta^+ = \left( \frac{5}{23}, \frac{3}{23}, \frac{9}{23}, \frac{6}{23} \right)$$

$$P\theta^- = \left( \frac{11}{16}, \frac{3}{16}, \frac{0}{16}, \frac{2}{16} \right)$$

$$\therefore p(x_*^+|\theta^+) = \frac{\left(\frac{5}{23}\right)^1}{1!} \cdot \frac{\left(\frac{3}{23}\right)^0}{0!} \cdot \frac{\left(\frac{9}{23}\right)^0}{0!} \cdot \frac{\left(\frac{6}{23}\right)^2}{2!} \times 3! = 0.0443$$

$$\therefore p(x_*^-|\theta^-) = \frac{\left(\frac{11}{16}\right)^1}{1!} \cdot \frac{\left(\frac{3}{16}\right)^0}{0!} \cdot \frac{\left(\frac{0}{16}\right)^0}{0!} \cdot \frac{\left(\frac{2}{16}\right)^2}{2!} \cdot 3! = 0.0322$$

$$\therefore p(x_*^+|\theta^+) > p(x_*^-|\theta^-)$$

$\therefore X_*$  should be classified as positive

$$\begin{aligned} p(x_*) &= p(x_*^+|\theta^+) \cdot p(\theta^+) + p(x_*^-|\theta^-) \cdot p(\theta^-) \\ &= 0.0443 \times \frac{4}{8} + 0.0322 \times \frac{4}{8} = 0.03825 \end{aligned}$$

$$\therefore p(+|x_*) = \frac{p(x_*^+|\theta^+) \cdot p(\theta^+)}{p(x_*)} = \frac{0.0443 \times \frac{4}{8}}{0.03825} = 0.5795$$

(b). (add-1) use smooth probabilities.

$$P\theta^+ = \left( \frac{6}{27}, \frac{4}{27}, \frac{10}{27}, \frac{7}{27} \right)$$

$$P\theta^- = \left( \frac{12}{20}, \frac{4}{20}, \frac{1}{20}, \frac{3}{20} \right)$$

$$\therefore P(x_*^+ | \theta^+) = \frac{\left(\frac{6}{27}\right)^1 \left(\frac{4}{27}\right)^0 \left(\frac{10}{27}\right)^0 \left(\frac{7}{27}\right)^2 \times 3!}{1! \cdot 0! \cdot 0! \cdot 2!} = 0.0448$$

$$\therefore P(x_*^- | \theta) = \frac{\left(\frac{12}{20}\right)^1 \left(\frac{4}{20}\right)^0 \left(\frac{1}{20}\right)^0 \left(\frac{3}{20}\right)^2 \cdot 3!}{1! \cdot 0! \cdot 0! \cdot 2!} = 0.0405$$

$$\begin{aligned} P(x_*) &= P(x_*^+ | \theta^+) \cdot P(\theta^+) + P(x_*^- | \theta) \cdot P(\theta) \\ &= 0.0448 \times \frac{4}{8} + 0.0405 \times \frac{4}{8} = 0.04265 \end{aligned}$$

$$\therefore P(- | x_*) = \frac{P(x_*^- | \theta) \cdot P(\theta)}{P(x_*)} = \frac{0.0405 \times \frac{4}{8}}{0.04265} = 0.4748$$

(c) Multivariate Bernoulli:

Document	A	B	C	D	class
1	1	0	1	1	+
2	0	1	1	0	+
3	1	0	0	1	+
4	0	0	1	0	+
5	0	0	0	1	-
6	1	0	0	0	-
7	1	1	0	0	-
8	1	0	0	1	-

$$E^+: A:2 \quad B:1 \quad C:3 \quad D:2 \quad P(\theta^+ : (\frac{2}{4}, \frac{1}{4}, \frac{3}{4}, \frac{2}{4}))$$

$$E^-: A:3 \quad B:1 \quad C:0 \quad D:2 \quad P(\theta^- : (\frac{3}{4}, \frac{1}{4}, \frac{0}{4}, \frac{2}{4}))$$

new document  $x^*$   $A:1 \quad B:0 \quad C:0 \quad D:2 \Rightarrow (1, 0, 0, 1)$

$$\therefore p(x^*|\theta^+) : \frac{2}{4} \times (1 - \frac{1}{4}) \times (1 - \frac{3}{4}) \times \frac{2}{4} = 0.046875$$

$$P(x^*|\theta^-) : \frac{3}{4} \times (1 - \frac{1}{4}) \times (1 - \frac{0}{4}) \times \frac{2}{4} = 0.28125$$

$$\begin{aligned} P(x^*) &= P(x^*|\theta^+) \cdot P(\theta^+) + P(x^*|\theta^-) \cdot P(\theta^-) \\ &= 0.046875 \times \frac{4}{8} + 0.28125 \times \frac{4}{8} = 0.1640625 \end{aligned}$$

$$\therefore P(+|x^*) = \frac{P(x^*|\theta^+) P(\theta^+)}{P(x^*)} = \frac{0.046875 \times \frac{4}{8}}{0.1640625} = 0.1429$$

(d).

Document	A	B	C	D	class
1	1	0	1	1	+
2	0	1	1	0	+
3	1	0	0	1	+
4	0	0	1	0	+
5	0	0	0	1	-
6	1	0	0	0	-
7	1	1	0	0	-
8	1	0	0	1	-
9	1	1	1	1	+
10	0	0	0	0	+
11	1	1	1	1	-
12	0	0	0	0	-

$$E^+ : A:2 \quad B:1 \quad C:3 \quad D:2 \quad P(\theta^+) = \left( \frac{3}{8}, \frac{2}{8}, \frac{4}{8}, \frac{3}{8} \right)$$

$$E^- : A:3 \quad B:1 \quad C:0 \quad D:2 \quad P(\theta^-) = \left( \frac{4}{8}, \frac{2}{8}, \frac{1}{8}, \frac{3}{8} \right)$$

new document  $x_*$   $A:1 \quad B:0 \quad C:0 \quad D:2 \Rightarrow (1, 0, 0, 1)$

$$\therefore p(x_* | \theta^+) : \frac{3}{8} \times (1 - \frac{2}{8}) \times (1 - \frac{4}{8}) \times \frac{3}{8} = 0.052$$

$$P(x_* | \theta^-) : \frac{4}{8} \times (1 - \frac{2}{8}) \times (1 - \frac{1}{8}) \times \frac{3}{8} = 0.1230$$

$$\begin{aligned} p(x_*) &= p(x_* | \theta^+) \cdot p(\theta^+) + p(x_* | \theta^-) \cdot p(\theta^-) \\ &= 0.052 \times \frac{4}{8} + 0.1230 \times \frac{4}{8} = 0.0879 \end{aligned}$$

$$\therefore p(-|x_*) = \frac{p(x_* | \theta^-) \cdot p(\theta^-)}{p(x_*)} = \frac{0.1230 \times \frac{4}{8}}{0.0879} = 0.7692$$

Q2:

$$(a) \mathcal{L}_C(y, \hat{y}) = \sum_{i=1}^n \mathcal{L}_C(y_i, \hat{y}_i) = \sum_{i=1}^n \left[ \sqrt{\frac{1}{C^2}(y_i - w^\top x_i)^2 + 1} - 1 \right]$$

$\because \hat{y}_i = w^\top x_i = w_0 + w_1 x_i$  for  $i = 1, \dots, n$ .

$$\therefore \mathcal{L}_C(y, \hat{y}) = \sum_{i=1}^n \mathcal{L}_C(y_i, \hat{y}_i) = \sum_{i=1}^n \left\{ \sqrt{\frac{1}{C^2}(y_i - (w_0 + w_1 x_i))^2 + 1} - 1 \right\}$$

$$\frac{\partial \mathcal{L}_C(y, \hat{y})}{\partial w_0} = \sum_{i=1}^n \frac{1}{2} \left( \frac{1}{C^2} (y_i - (w_0 + w_1 x_i))^2 + 1 \right)^{-\frac{1}{2}} \cdot 2(y_i - (w_0 + w_1 x_i)) \cdot (-1)$$

$$= \sum_{i=1}^n \frac{1}{C^2} (y_i - w_0 - w_1 x_i)^2 + 1)^{-\frac{1}{2}} (y_i - w_0 - w_1 x_i) \cdot (-1)$$

$$= \sum_{i=1}^n \frac{1}{C^2} \cdot \frac{-(y_i - w_0 - w_1 x_i)}{\sqrt{(y_i - w_0 - w_1 x_i)^2 + 1}} .$$

$$\frac{\partial \mathcal{L}_C(y, \hat{y})}{\partial w_1} = \sum_{i=1}^n \frac{1}{2} \left( \frac{1}{C^2} (y_i - (w_0 + w_1 x_i))^2 + 1 \right)^{-\frac{1}{2}} \cdot 2(y_i - (w_0 + w_1 x_i)) \cdot (-x_i)$$

$$= \sum_{i=1}^n \frac{1}{C^2} (y_i - w_0 - w_1 x_i)^2 + 1)^{-\frac{1}{2}} \cdot (y_i - w_0 - w_1 x_i) \cdot (-x_i)$$

$$= \sum_{i=1}^n \frac{1}{C^2} \cdot \frac{(y_i - w_0 - w_1 x_i) \cdot (-x_i)}{\sqrt{(y_i - w_0 - w_1 x_i)^2 + 1}}$$

(b) Initialise  $t = 0$

$$w = (w_0^t, w_1^t)$$

$t = 1, 2, 3, 4, \dots, n$ ,

while True:

$$w_0^{(t+1)} = w_0^{(t)} - \alpha \cdot \frac{\partial c(y, \hat{y})}{\partial w_0}$$

$$w_1^{(t+1)} = w_1^{(t)} - \alpha \cdot \frac{\partial c(y, \hat{y})}{\partial w_1}$$

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

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

if change in  $L(w^{(t)})$  is negligible,

break

(c). The result is:

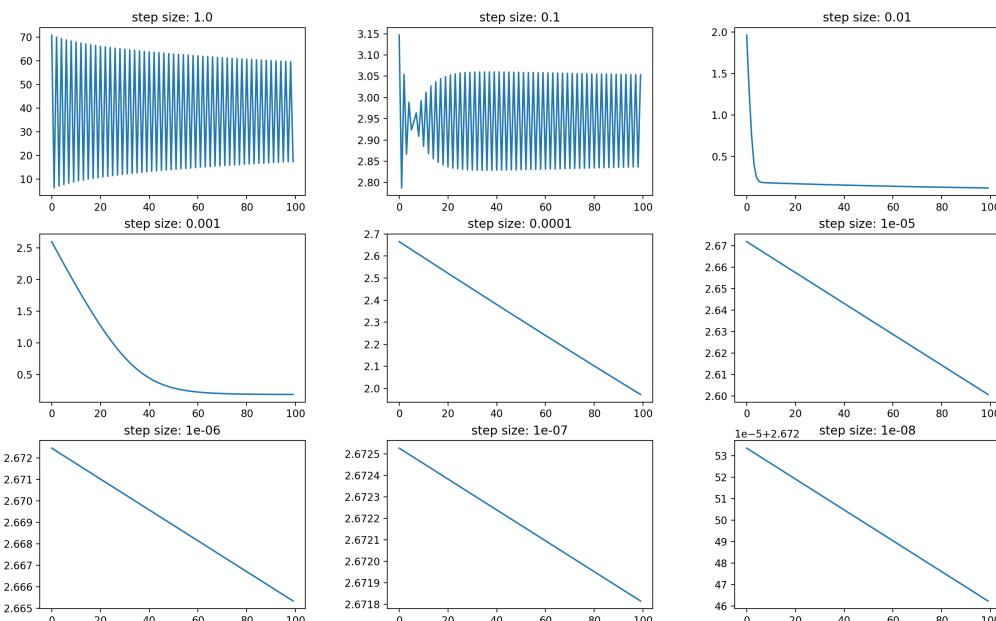


figure 2.1

```

6   ## Question 2
7
8   import numpy as np
9   import matplotlib.pyplot as plt
10
11  np.random.seed(42)      # make sure you run this line for consistency
12  x = np.random.uniform(1, 2, 100)
13  y = 1.2 + 2.9 * x + 1.8 * x*x + np.random.normal(0, 0.9, 100)
14  plt.scatter(x,y)
15  plt.show()
16
17  ##_(c)
18  c = 2
19  interation = 100
20  import math
21  def loss_function(c,x,y):
22      result = math.sqrt((pow((x-y),2)/c+c+1))-1
23      return result
24
25  def w0_cost_function(c,x,y,w0,w1):
26      result =0
27      i= 0
28      while i < interation:
29          cy = x[i]*w1 +w0
30          ty = y[i]
31          result = result + (- (1/pow(c,2))* (ty-cy)/(math.sqrt((1/pow(c,2))*pow((ty-cy),2)+1)))
32          i = i +1
33      return result
34
35  def w1_cost_function(c,x,y,w0,w1):
36      result =0
37      i = 0
38      while i < interation:
39          cy = x[i]*w1 + w0
40          ty = y[i]
41          result =result + (- (1/pow(c,2))*x[i]*(ty -cy)/(math.sqrt((1/pow(c,2))*pow((ty-cy),2)+1)))
42          i = i +1
43      return result
44
45  def update_weight_function(w0,w1,c,x,y,alphy):
46      nw0 = w0-alphy *w0_cost_function(c,x,y,w0,w1)
47      nw1 = w1-alphy *w1_cost_function(c,x,y,w0,w1)
48      return nw0,nw1
49
50  def run_function():
51      lista =[]
52      i = 0
53      while i < 9:
54          lista.append(1*pow(10,-i))
55          i = i+1
56
57      temploss =[]
58      losses =[]
59      for i in lista:
60          temploss =[]
61          w0 =1
62          w1 =1
63          j = 0
64          while(j<interation):
65              w0,w1 = update_weight_function(w0,w1,c,x,y,i)
66              two_temploss = []
67              k =0
68              while(k<interation):
69                  cy = x[k]*w1 +w0
70                  two_temploss.append(loss_function(c, cy, y[k]))
71                  k = k +1
72              temploss.append(np.mean(two_temploss))
73              j = j +1
74          losses.append(temploss)
75      return losses
76
77
78  losses = run_function()
79
80
81  ## plotting help
82  fig, ax = plt.subplots(3,3, figsize=(10,10))
83  alphas = [10e-1, 10e-2, 10e-3,10e-4,10e-5,10e-6,10e-7, 10e-8, 10e-9]
84  for i, ax in enumerate(ax.flat):
85      # losses is a list of 9 elements. Each element is an array of length 100 storing the loss at each iteration for
86      # that alpha value.
87      ax.plot(losses[i])
88      ax.set_title(f"step size: {alphas[i]}") # plot titles
89      plt.tight_layout() # plot formatting
90      plt.show()

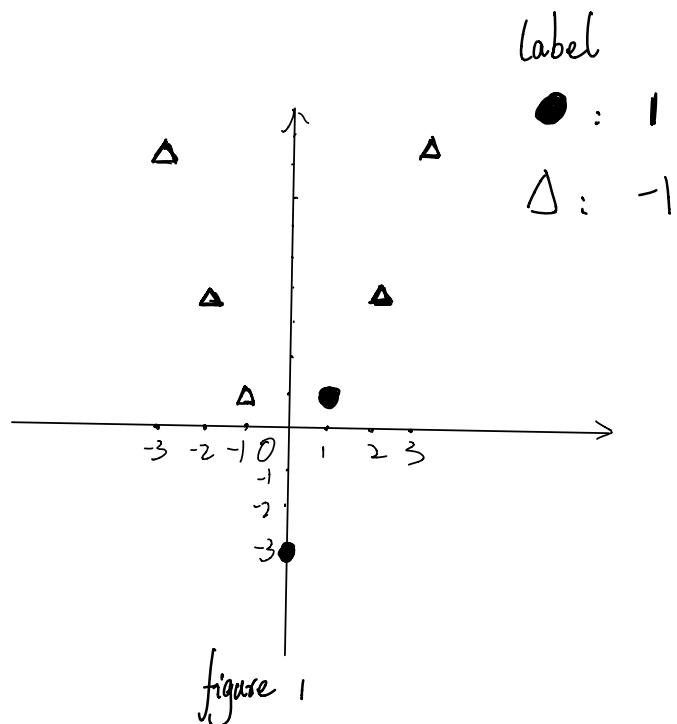
```

(d). the first two step size can not find the correct weights until the 5 step size that can be train. If step size is large, the classifier may miss the optimal solution. If step size is too small, the classifier trains will becomes slow and require more iteration

(e).

Q4:

(a) According to figure 1,  
the data is linearly separable.



(b) according to (a), figure 1. the max margin is close to (-1, 1) (1, 1) (2, 4)  
the classifier only near to (-1, 1), (1, 1), (2, 4)

$$\therefore X = \begin{pmatrix} -1 & 1 \\ 1 & 1 \\ 2 & 4 \end{pmatrix} \quad Y = \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix} \quad X^T = X^T = \begin{pmatrix} 1 & -1 \\ 1 & 1 \\ -2 & -4 \end{pmatrix}$$

$$X^T \cdot X^T = \begin{pmatrix} 2 & 0 & 2 \\ 0 & 2 & -6 \\ 2 & -6 & 20 \end{pmatrix}$$

$$\begin{aligned} & \arg \max_{a_1, \dots, a_n} \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_i a_j y_i y_j x_i \cdot x_j \quad \text{subject to } a_i \geq 0, \text{ for all } i, \sum_{i=1}^n a_i y_i \leq 0 \\ & = \arg \max_{a_1, \dots, a_3} (a_1 + a_2 + a_3) - \frac{1}{2} (2a_1^2 + 2a_2^2 + 20a_3^2 + 2a_1 a_2 + 2a_1 a_3 + 2a_2 a_3 - 6a_2 a_3 + 2a_3 a_1 - 6a_3 a_2) \end{aligned}$$

$$= a_1 + a_2 + a_3 - a_1^2 - a_2^2 - 10a_3^2 - 2a_1 a_3 + 6a_2 a_3$$

$$\therefore \sum_{i=1}^n a_i y_i = 0$$

$$\therefore -a_1 + a_2 - a_3 = 0$$

$$a_2 - a_1 = a_3$$

$$\therefore a_1 + a_2 + (a_2 - a_1) - a_1^2 - a_2^2 - 10(a_2 - a_1)^2 - 2a_1 \cdot (a_2 - a_1) + 6a_2 \cdot (a_2 - a_1)$$

$$= 2a_2 - a_1^2 - a_2^2 - 10a_2^2 + 20a_2a_1 - 10a_1^2 - 2a_1a_2 + 2a_1^2 + 6a_2^2 - 6a_2a_1$$

$$\text{Assume: } \frac{\partial L}{\partial a_1} = \frac{\partial L}{\partial a_2} = 0$$

$$\frac{\partial L}{\partial a_1} = -2a_1 + 20a_2 - 20a_1 - 2a_2 + 4a_1 - 6a_2$$

$$\frac{\partial L}{\partial a_2} = 2 - 2a_2 - 20a_2 + 20a_1 - 2a_1 + 12a_2 - 6a_2$$

$$\begin{cases} 12a_2 - 18a_1 = 0 \\ 2 - 4a_2 + 18a_1 = 0 \end{cases} \Rightarrow \begin{cases} a_1 = \frac{2}{3} \\ a_2 = 1 \\ a_3 = \frac{1}{3} \end{cases}$$

$$(c) w = \sum_{i=1}^m a_i y_i x_i$$

$$= a_1 y_1 x_1 + a_2 y_2 x_2 + a_3 y_3 x_3$$

$$= -\frac{2}{3} \times \begin{pmatrix} -1 \\ 1 \end{pmatrix} + 1 \times \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \frac{1}{3} \times \begin{pmatrix} 2 \\ 4 \end{pmatrix}$$

$$= \begin{pmatrix} \frac{2}{3} + 1 - \frac{2}{3} \\ -\frac{2}{3} + 1 - \frac{4}{3} \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$\text{Thus, resulting in a margin is } \frac{1}{\|w\|} = \frac{1}{\sqrt{1^2 + (-1)^2}} = \frac{1}{\sqrt{2}} = \frac{\sqrt{2}}{2}$$

$$\therefore y_1 (w^T x_1 - t) = 1$$

$$\therefore y_1 = -1 \quad x_1 = (-1, 1)$$

$$\therefore -1(-2 - t) = 1$$

$$t = 1$$

Thus, bias  $t = 1$ .

(d) Linear classifiers are unable to represent non-linear functions.

Linear classifier is a line that discriminate two class, and the two class are called linearly-separable. And in some condition, It is clear that drawing one straight line cannot identify all the points of the class correctly.

For example:

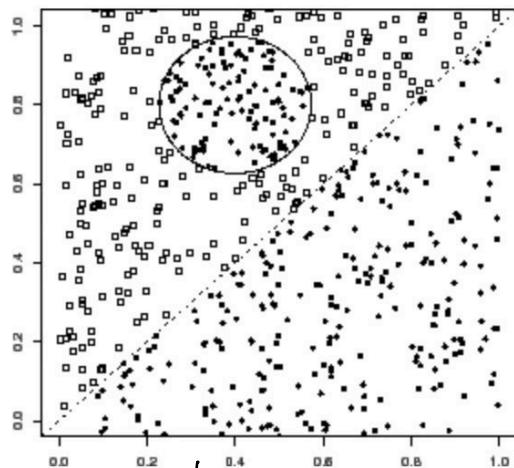


figure 1

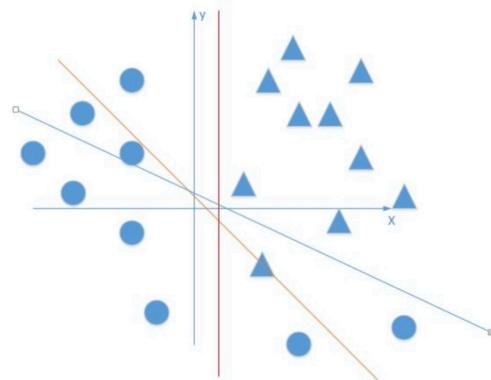


figure 2

Although 'kernel trick' will decrease dimension, it will not decrease to 1 dimension.

(e). the kernel trick is general task of pattern analysis is to find and study types of relations in dataset. for many algorithm that solve these tasks, the data in raw representation have to be explicitly transformed into feature vector via a user-specify feature map. In contrast, kernel methods require only a user-specified kernel and can operate in a high-dimensional

Thus, this is such an important technique for machine learning.

QF :

(a).

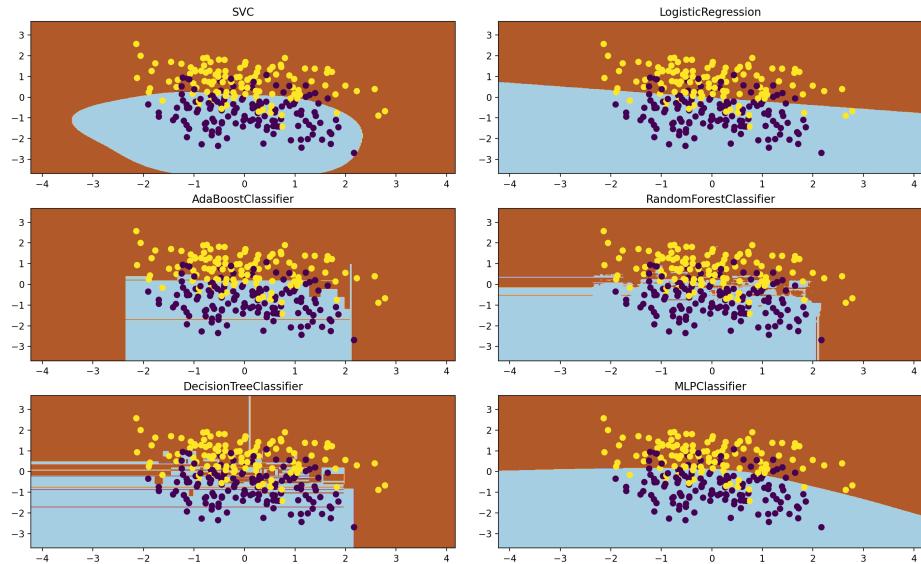
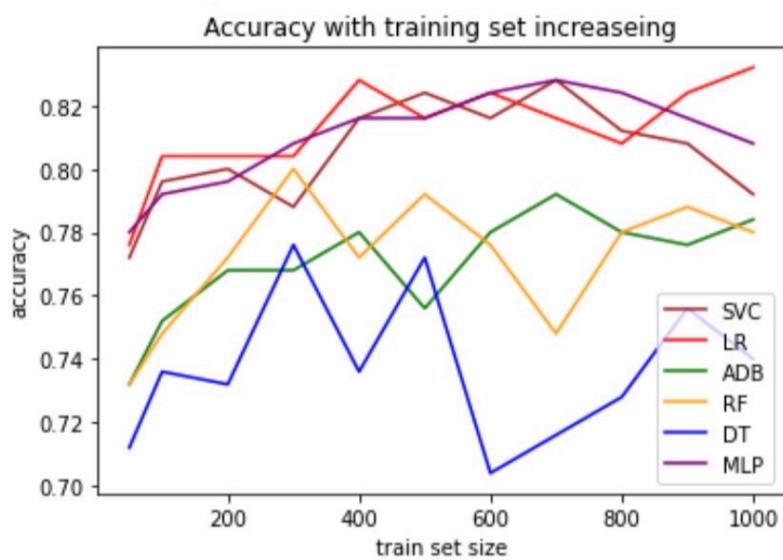


figure 1

(b).



with the increase of train set size. The accuracy of DT are not stability especially in 600 train set size, it have the largest accuracy. In contrast, with the increase of train set size, ADB has the increasing accuracy, SVC

and LR and MLP have high accuracy all the time. ✓ ✓

(c)

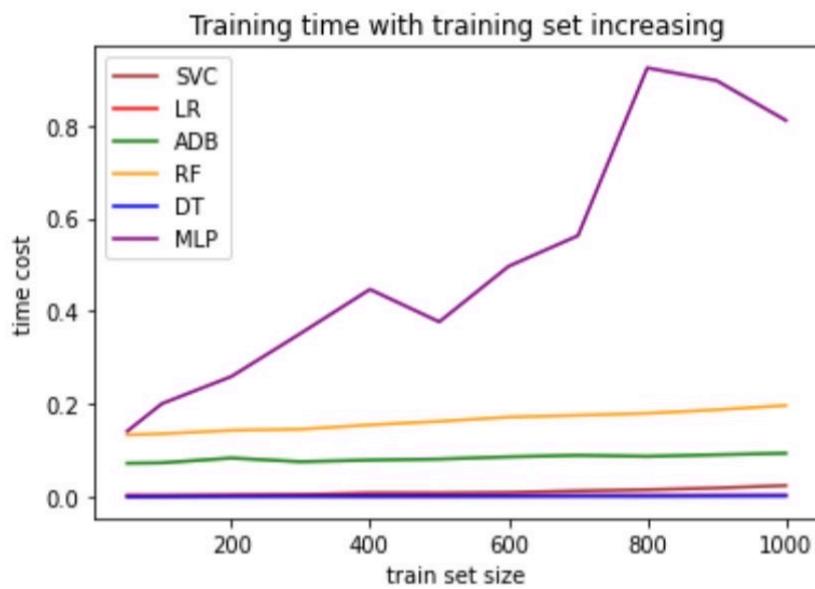


figure 3

with the increasing train set size, the time cost of MLP one still increase.  
Because MLP is dealing with data, it have a lot of layers so  
it will cost too much time