

Problem Set 1

Haoran Wang

Handed In: February 3, 2017

1. Answer to problem 1

- (a) The algorithm is designed to learn conjunctions in the hypothesis space, so we should begin with the complete conjunctions of all elements:

$$G = x_1 \wedge (\neg x_1) \wedge x_2 \wedge (\neg x_2) \wedge \dots \wedge x_n \wedge (\neg x_n) \quad (1)$$

So the result of initial conjunction G is always negative. Then for every positive training example, we need to eliminate x_i or $(\neg x_i)$ based on the value of x_i in the training example. If $x_i = 1$, then we eliminate $(\neg x_i)$ in G ; If $x_i = 0$, then we eliminate (x_i) in G . After going through all positive training examples, we could get a final conjunction function G , if G is empty, then we halt. If it is not empty, we could verify the modified conjunction function G by passing all training examples, comparing the labels it will produce with the actual label of the example. If there is any inconsistent example, we halt. If we halt at any time in the mentioned process, we can say that there is no consistent hypotheses.

Algorithm 1 LearnConjunctions($T[1,2,3,\dots,m]$)

```

1:  $G \leftarrow x_1 \wedge (\neg x_1) \wedge x_2 \wedge (\neg x_2) \wedge \dots \wedge x_n \wedge (\neg x_n)$ 
2: for  $i \leftarrow 1$  to  $m$  do
3:   if  $T[i]$  is labeled as '+' then
4:     if  $T[i]_j = 1$  then
5:       Remove  $(\neg x_j)$  from  $G$ 
6:     end if
7:     if  $T[i]_j = 0$  then
8:       Remove  $x_j$  from  $G$ 
9:     end if
10:  end if
11: end for
12: if  $G$  is empty then
13:   halt, and return "No Consistent Hypothesis"
14: end if
15: for  $i \leftarrow 1$  to  $m$  do
16:   if Label of  $T[i]$  is not the same with output of  $G$  then
17:     halt, and return "No Consistent Hypothesis"
18:   end if
19: end for
20: return  $G$ 

```

(b) Proof:

There are two possible results from the algorithm: first, we get a conjunction function which is consistent with the hypothesis. In this circumstance, there is one step in the algorithm that it will verify the found conjunction function G by going through all training examples. So if we could get a conjunction function from the algorithm, we could say that it is consistent with all the examples we have seen. So this part is verified.

Second, we may get "No consistent hypothesis" from the algorithm. 1) If we halt at the first loop after going through all positive examples, we could know that G is empty now, which means we do not find a consistent conjunction function with all positive examples, so it is impossible to find another function consistent with all positive training examples because the initial form of G is already the most general case. 2) If we halt at the second loop when going through all training examples (Both positive and negative), the function is inconsistent with at least one negative training example. Because we already eliminate the x_i or $(\neg x_i)$ or both due to it is inconsistent with the positive examples, so we cannot add back to the conjunction function G . So it is a contradiction to the previous loop, which means there is no consistent hypothesis existed.

Finish

- (c) As we can see from the algorithm pseudo-code, the first loop will compare the value of elements contained in the positive training example. So in worst case, we need to compare n elements if all of them are in training example. The time is $O(n)$ if we ignore the time of removing elements from G . As we have two loops here, so we need go through all m training examples, and we could infer that the time complexity is $O(mn)$.
- (d) If the conjunction function G is consistent with all training examples we have seen, we could say that we are confident it could correctly predict positive examples. In other words, if the output of G is positive, the label of new example must be positive. If the output is negative, we cannot say that our output is one hundred percent accurate. It could be wrong. Because we have eliminated the inconsistent elements in the process to find consistent function G as initialized in equation (1). The initial conjunction will output negative by default. When going through positive examples, we remove some elements in G to make it consistent with positive cases, but it is always the most general case of all positive cases. So based on this process, we are certain that if we predict a positive label, it must be positive. But when we predict a negative label, maybe it is wrong because we are not guaranteed to have seen all positive examples, so maybe we will misclassify some extreme positive cases.

In a nutshell, when the function output is positive, the label must be positive. If the output is negative, the label maybe negative or positive.

2. Answer to problem 2

- (a) Analytical solution: suppose x_1 is a point on the hyperplane $\vec{w}^T \vec{x} + \theta = 0$, so we have the following equation:

$$\vec{w}^T \vec{x}_1 + \theta = 0 \quad (2)$$

The distance between x_0 and x_1 is:

$$d = \|\vec{x}_0 - \vec{x}_1\| \quad (3)$$

We know that the normal vector of the hyperplane is \vec{w} , so in order to find the distance between the point and the hyperplane, we only need to calculate the projection of d on the direction of \vec{w} , which is calculated as:

$$h = d \cdot \cos\theta, \text{ where } \theta \text{ is the angle between } \vec{w} \text{ and } \vec{x}_0 - \vec{x}_1 \quad (4)$$

Thus, the distance h could be calculated as:

$$\begin{aligned} h = d \cdot \cos\theta &= \frac{\vec{w} \cdot (\vec{x}_0 - \vec{x}_1)}{\|\vec{w}\| \cdot \|\vec{x}_0 - \vec{x}_1\|} \cdot \|\vec{x}_0 - \vec{x}_1\| \\ &= \frac{\|\vec{w} \cdot \vec{x}_0 - \vec{w} \cdot \vec{x}_1\|}{\|\vec{w}\|} \\ &= \frac{|\vec{w} \cdot \vec{x}_0 + \theta|}{\|\vec{w}\|} \end{aligned} \quad (5)$$

- (b) Suppose \vec{x}_1 is a point on the hyperplane $\vec{w}^T \vec{x} + \theta_1 = 0$, similarly, \vec{x}_2 is a point on the other hyperplane $\vec{w}^T \vec{x} + \theta_2 = 0$. Suppose the distance between \vec{x}_1 and \vec{x}_2 could be represented as the distance between two hyperplanes.

We know the normal vector \vec{w} is perpendicular to the two hyperplanes, so the unit vector of \vec{w} could be used as:

$$\vec{u} = \frac{\vec{w}}{\|\vec{w}\|}$$

So the vector between \vec{x}_1 and \vec{x}_2 could be represented as:

$$\vec{x}_1 - \vec{x}_2 = d \cdot \vec{u} \quad (6)$$

where d is the distance (scalar) between two hyperplanes. So we can get the representation of d by dot product \vec{w} on both sides in order to use the equations $\vec{w}^T \vec{x} + \theta_1 = 0$ and $\vec{w}^T \vec{x} + \theta_2 = 0$, the equation is:

$$\begin{aligned} \vec{w} \cdot \vec{x}_1 - \vec{w} \cdot \vec{x}_2 &= d \cdot \frac{\vec{w}}{\|\vec{w}\|} \cdot \vec{w} \\ d &= \frac{|\vec{w} \cdot \vec{x}_1 - \vec{w} \cdot \vec{x}_2|}{\|\vec{w}\|} \\ d &= \frac{|\theta_1 - \theta_2|}{\|\vec{w}\|} \end{aligned} \quad (7)$$

So the distance of two hyperplanes is $d = \frac{|\theta_1 - \theta_2|}{\|\vec{w}\|}$

3. Answer to problem 3

- (a) i. First we need to prove from LHS to RHS, suppose we have data set D is linearly separable, we need to prove that $\delta = 0$ in this situation.

We know that there exists a hyperplane $\vec{w}^T \vec{x} + \theta = 0$ that could separate the data set, we choose the point that is closet to the hyperplane from the positive side. And it is the same that we need select a point from negative side. We suppose V_1 is the point from positive side, we have:

$$V_1 = \vec{w}^T \vec{x}_1 + \theta \geq 0$$

Suppose V_2 is the point from negative side, we have:

$$V_2 = \vec{w}^T \vec{x}_2 + \theta < 0$$

Suppose there exists a small α that we could get a new hyperplane that could also separate the data set if we move the current hyperplane by α . We have $V_1 - \alpha \geq 0 > V_2 - \alpha$, we could get the value of α analytically. Since we move the hyperplane on both sides, which means two hyperplanes are symmetric along with a central hyperplane that has the biggest margin between the two points. And the distance between the central hyperplane and either of the two hyperplanes is the same, we can get the equation of this distance by using α :

$$|V_1 - \alpha| = |V_2 - \alpha|$$

$$V_1 - \alpha = -V_2 + \alpha$$

$$\alpha = \frac{V_1 + V_2}{2}$$

To get the format of $y_i(\vec{w}^T \vec{x}_i + \theta) \geq 1 - \delta$, $\forall (\vec{x}_i, y_i) \in D$, we have:

$$y_i(\vec{w}^T \vec{x}_i + \theta - \alpha) = \frac{V_1 - V_2}{2}$$

Since $V_1 > V_2$, we have :

$$y_i\left(\frac{2}{V_1 - V_2} \cdot \vec{w}^T + \frac{2}{V_1 - V_2} \cdot (\theta - \alpha)\right) = 1$$

$$y_i(\vec{w}'^T + \theta') \geq 1 - \delta$$

So we could see that $\vec{w}' = \frac{2}{V_1 - V_2} \cdot \vec{w}$, $\theta' = \frac{2}{V_1 - V_2} \cdot (\theta - \alpha)$, and $1 = 1 - \delta$, which implies $\delta = 0$

Second, we also need to prove from RHS to LHS, when we have $\delta = 0$, we need to prove that

$$y_i = \begin{cases} 1 & \text{if } \vec{w}^T \vec{x}_i + \theta \geq 0 \\ -1 & \text{if } \vec{w}^T \vec{x}_i + \theta < 0. \end{cases}$$

So, $\delta = 0$, we have

$$\begin{aligned} y_i(\vec{w}^T \vec{x} + \theta) &\geq 1 - \delta \\ y_i(\vec{w}^T \vec{x} + \theta) &= 1, \text{ where } |y_i| = 1 \\ \vec{w}^T \vec{x} + \theta &= 1 > 0 \\ \vec{w}^T \vec{x} + \theta &= -1 < 0 \end{aligned}$$

Therefore, it is obvious that the equations satisfy the condition in (1). So we finish the proof.

If we have $\delta > 0$, we have several conditions to consider. First, if $1 - \delta > 0$, we can prove that data set could be separated in the same way as mentioned for $\delta = 0$. But if $\delta \geq 1$, y_i and $\vec{w}^T \vec{x}_i + \theta$ have different signs, so the data set is not separable.

ii. The trivial solution:

it is easy to see that if we set $\delta = 0$, then \vec{w}, θ are also 0. This solution could satisfy the conditions we get from [2] to [4]. And $\delta = 0$ could satisfy the condition [4].

We don't use the formulation because the trivial optimal solution is not a hyperplane to separate the data set.

iii. There are only two data points in the data set. It is obvious that we could separate the data set.

$$w_1 + w_2 + w_3 + \dots + w_n + \theta \geq 1 \quad (8)$$

$$-(w_1 + w_2 + w_3 + \dots + w_n - \theta) \leq -1 \quad (9)$$

So from (8) and (9), we have $w_1 + w_2 + w_3 + \dots + w_n \geq 1 + |\theta|$, and $\delta = 0$. This is the optimal solution.

(b) i. The function findLinearDiscriminant is written as:

```

1 % This function finds a linear discriminant using LP
2 % The linear discriminant is represented by
3 % the weight vector w and the threshold theta.
4 % YOU NEED TO FINISH IMPLEMENTATION OF THIS FUNCTION.
5
6 function [w, theta, delta] = findLinearDiscriminant(data)
7 %% setup linear program
8 [m, np1] = size(data);
9 n = np1-1;
10
11 % write your code here
12 % representation of A, b, c

```

```

13 A = zeros(m+1, n+2);
14 b = zeros(m+1,1);
15 c = zeros(n+2,1);
16 c(n+2,1) = 1;
17 b(1:m,1) = 1;
18 % label of the training examples
19 y_i = data(:,np1);
20 for i = 1:m
21     A(i,1:n) = y_i(i) * data(i,1:n);
22 end
23 A(1:m,n+1) = y_i(1:m);
24 A(1:m,n+2) = 1;
25 A(m+1,n+1) = 0;
26 A(m+1,n+2) = 1;
27 %% solve the linear program
28 %adjust for matlab input: A*x <= b
29 [t, z] = linprog(c, -A, -b);
30
31 %% obtain w,theta,delta from t vector
32 w = t(1:n);
33 theta = t(n+1);
34 delta = t(n+2);
35
36 end

```

ii. 1)The monotone conjunction data set is:

```

1 1 1 1
2 0 0 -1
3 0 1 -1
4 1 0 -1

```

2)The plot2dSeparator function is:

```

1 % This function plots the linear discriminant.
2 % YOU NEED TO IMPLEMENT THIS FUNCTION
3
4 function plot2dSeparator(w, theta)
5 x = [-2:0.01:2];
6 y = -(w(1)*x+theta)/w(2);
7 plot(x, y);
8 end

```

3)The picture of the line separator with data is shown as Figure 1:

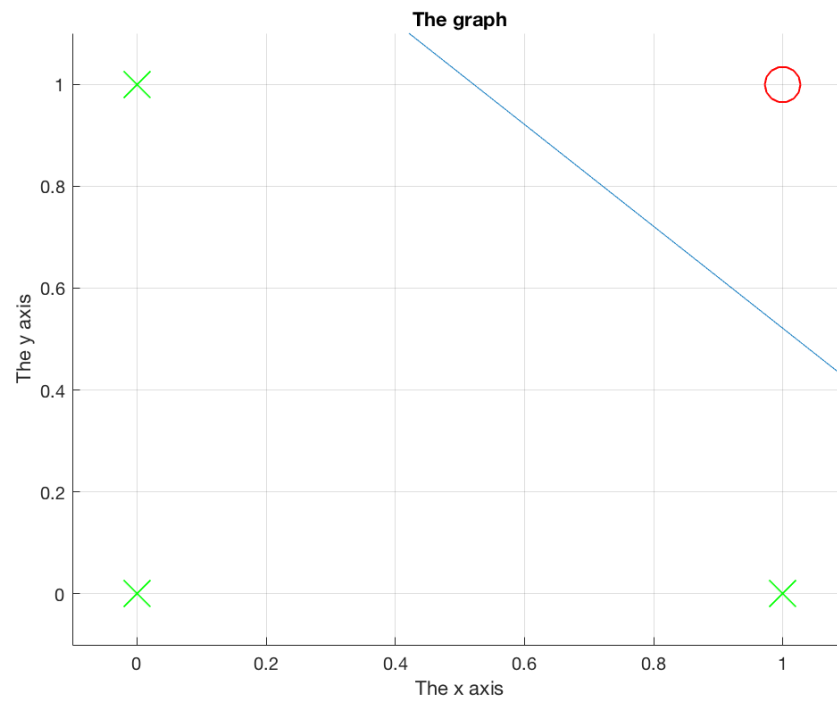


Figure 1: Separator with data points

After using the `findLinearDiscriminant` function, we get the hyperplane of 2d monotone conjunction function as:

$$156.2452x_1 + 156.2452x_2 - 237.6709 = 0$$

4) For learning `hw1conjunction.txt`, we get the output as shown at Figure 2:

```

Command Window

w =

    2.9104
   -2.0498
    0.1775
   190.5196
    0.1399
   -3.1010
   -2.9534
  -193.2778
    1.1679
   -8.8942

theta =

   -90.2115

delta =

   -2.4158e-13

```

Figure 2: Output for hw1conjunction.txt

According to the result from Figure 6, the answer is: $\vec{w}^T x - \theta = 0$. $\theta = -90.2115$, $\delta = -2.4158 \times 10^{-13}$ and $\delta \approx 0$, \vec{w} is listed in the Figure 6. Therefore, from this function, we could conclude that the data set is separable. δ is the target we want to minimize in this function. θ is the threshold for the function $\vec{w}^T \vec{x}$, which means it could be a range of values that could fit the hyperplane within the gap between the closest positive point and negative point. As for \vec{w} , we could observe from Figure 6 that the 4th and 8th element of \vec{w} have larger weights compared with others, although the value of 8th weight is negative. To represent the weights in a conjunction function, we could denote it as:

$$G = x_4 \wedge \neg x_8 \quad (10)$$

We have $\neg x_8$ because the value is negative in the actual weights, so we use negation in this conjunction.

iii. 1) The computeLabel function is implemented as follows:

```

1 % This function computes the label for the given
2 % feature vector x using the linear separator represented by
3 % the weight vector w and the threshold theta.
4 % YOU NEED TO WRITE THIS FUNCTION.
5
6 function y = computeLabel(x, w, theta)
7 temp = w.' * x + theta;
8 % w_dim = n x 1; where x_dim = n x 1;
9 disp(w)
10 if temp < 0
11     y = -1;

```



```

12 else
13     y = 1;
14 end
15 end

```

2) The result of the script is shown at Figure 3:

```

delta =

    1.2619e-09

accuracyInTrain =

    1

accuracyInTest =

    1
fx

```

Figure 3: Result of learnBadges

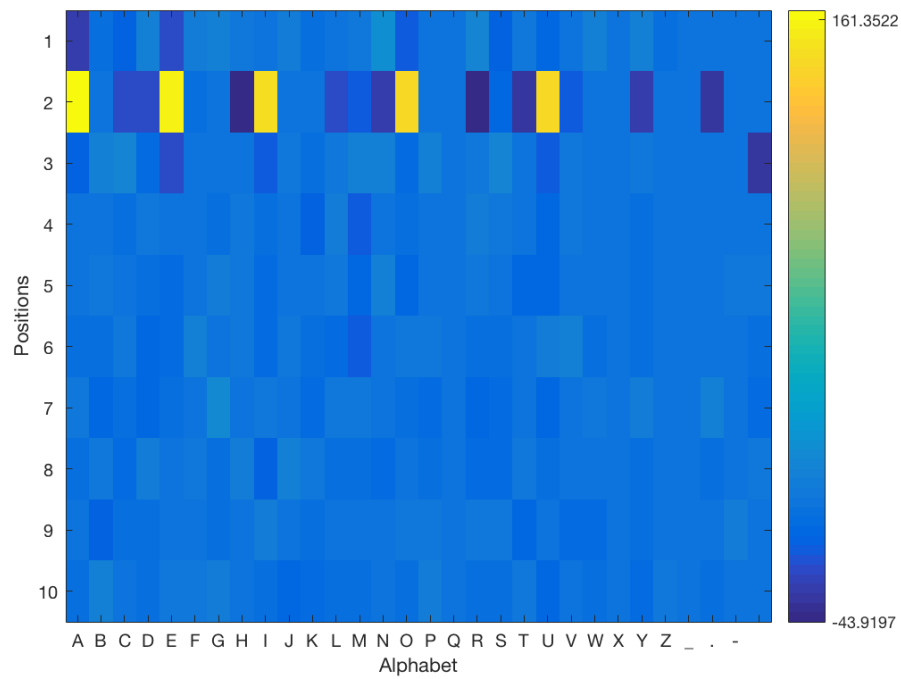


Figure 4: Original Alphabets and Positions

From the result, we could observe that $\delta \approx 0$, which means we found a linear separator for the data set. The train accuracy is 1, which also indicates that we successfully separate the data set. Besides, using this function also separates the test data set with no error, the accuracy is also 1. Figure 9 shows the weights for our features in the Alphabet vector. We could see that the letter 'A', 'E', 'I', 'O', 'U' have the largest weights compared with others at the 2nd position. Based on this feature, we could propose a hypothesis that the name with vowel at 2nd position will be classified as positive, otherwise, it will be negative.

4) To achieve 100% accuracy in training examples, but less than 100% on test data, I found two ways to achieve that.

First, we could remove letter 'U' at the Alphabets, then we could get result shown at Figure 5.

```
delta =  
    1.6200e-12  
  
accuracyInTrain =  
    1  
  
accuracyInTest =  
    0.9574
```

Figure 5: Remove 'U' at Alphabets

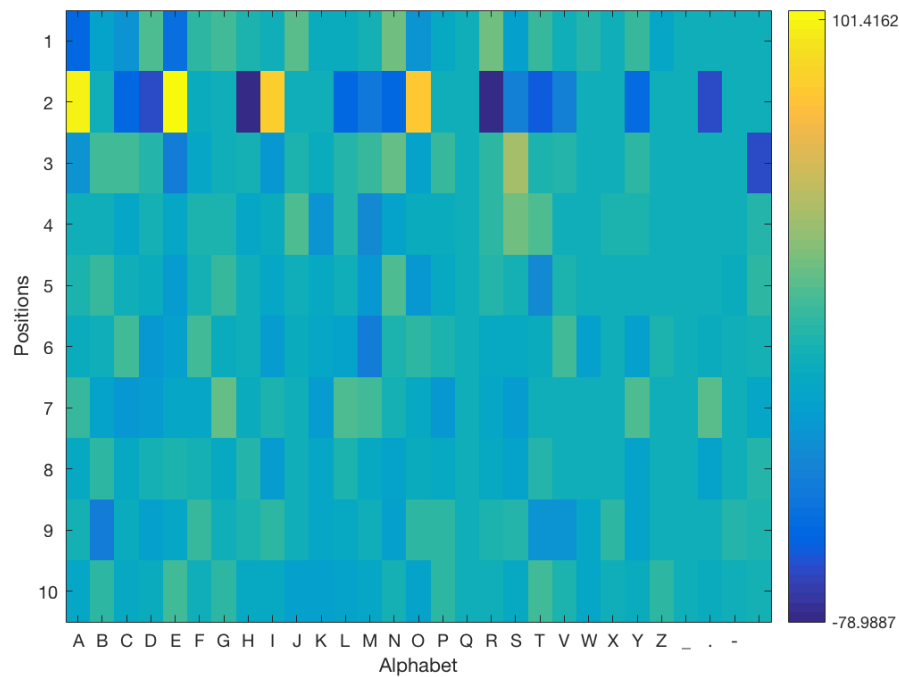


Figure 6: Weights graph when remove 'U' at Alphabets

We can get this result because we may have never seen the 'U' at 2nd position in our training data set, but we may have several cases at test data set, which explains why we cannot get 100% accuracy in the test data set.

Second, I found changing the positions could also generate the required result. Suppose I change the position from 1:10 to 3:10, then we can get:

```

delta =

    1.0232e-12

accuracyInTrain =

    1

accuracyInTest =

    0.7234

```

Figure 7: Change position to 3:10

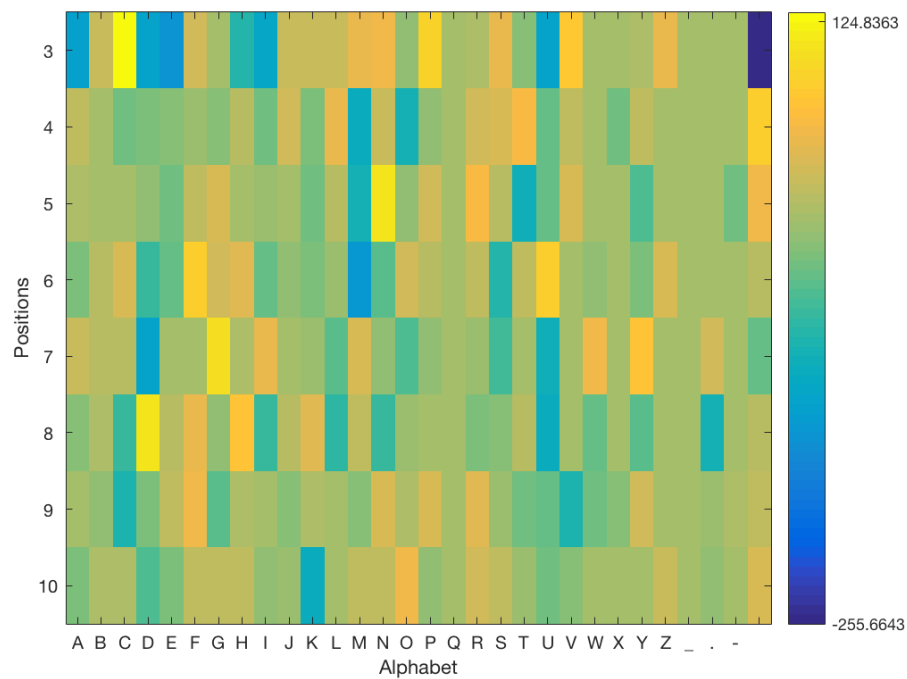


Figure 8: Weights graph when change position to 3:10

Therefore, from the Figure 8, we could conclude that the function finds a more complicated combination of the features to represent the desired linear

separator. We also have $\delta \approx 0$, which indicates it could separate the train data set. But the accuracy is lower at test data set, which indicates the separator could not completely separate the test data set, we may not choose the best separator.

iv. 1) The findLinearThreshold function is implemented as:

```

1 % This function solves the LP problem for a given weight vector
2 % to find the threshold theta.
3 % YOU NEED TO FINISH IMPLEMENTATION OF THIS FUNCTION.
4
5 function [theta,delta] = findLinearThreshold(data,w)
6 %% setup linear program
7 [m, np1] = size(data);
8 n = np1-1;
9
10 % write your code here
11 A = zeros(m+1, n+2);
12 b = zeros(m+1,1);
13 c = zeros(n+2,1);
14 c(n+2,1) = 1;
15 b(1:m,1) = 1;
16 % label of the training examples
17 y_i = data(:,np1);
18 for i = 1:m
19     A(i,1:n) = y_i(i) * data(i,1:n);
20 end
21 A(1:m,n+1) = y_i(1:m);
22 A(1:m,n+2) = 1;
23 A(m+1,n+1) = 0;
24 A(m+1,n+2) = 1;
25 %% solve the linear program
26 %adjust for matlab input: A*x <= b
27 [t, z] = linprog(c, -A, -b, [], [], [w' -inf -inf], [w' inf inf])
28     ;
29
30 %% obtain w,theta,delta from t vector
31
32 theta = t(n+1);
33 delta = t(n+2);
34
35 end

```

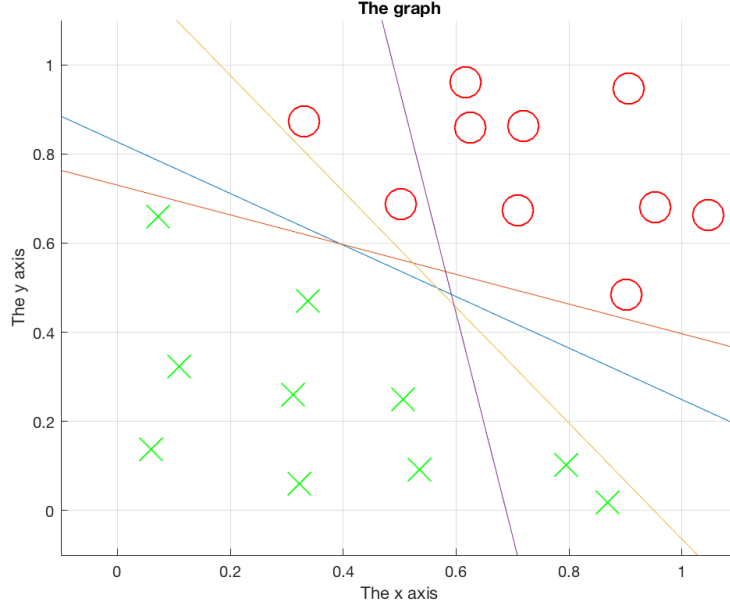


Figure 9: Line separators

2) From the output of learnMultipleHypothesis function, we get the Figure 9. We could see the blue one is the separator generated from the findLinearDiscriminant function, which is not given \vec{w} , so it has larger margin than others. So the other three lines (yellow, orange, purple) are generated by the given weights. From the graph, we could see the yellow, blue, orange lines could separate the data without errors. In other words, it could separate the data set exactly. But the purple line could not separate the data set. In addition, we could see that the yellow and orange lines are very close to positive or negative data points.

For yellow, blue and orange lines, the value of $\delta \approx 0$, which indicates it could separate the data set exactly.

As for which separator is better, basically speaking, the blue one maybe better because it has larger margin to the closest points on both sides. But we cannot say the other two separators (yellow and orange) are not good because we only have these data points, maybe they can do better job when it comes to lots of new training examples.

From the Figure 9, we found that the yellow and orange separators are also valid for separating the entire data set. So they might also be valid output from the process [2] to [4] in linear program. So we could infer that the solution is not unique.