

Problem Set 1

Bangqi Wang

Handed In: March 11, 2017

1. Answer to problem 1

- a. The main idea of the algorithm is to construct an initial hypothesis and then update the hypothesis when iterating the training set. For instance, the algorithm constructs an initial complete hypothesis that is the conjunction of all variables and their negations. i.e. $\mathbf{H} = (x_1 \wedge x_2 \wedge \dots \wedge x_n \wedge \neg x_1 \wedge \neg x_2 \wedge \dots \wedge \neg x_n)$. Then, for each positive sequences in training set, the algorithm updates the hypothesis by removing the contradictions. i.e. $x_i = a$ in hypothesis, but $x_i = b$ in sequence. After finishing iterating positive sequences, using negative sequences to check the hypothesis. If there is any contradiction between negative sequence and hypothesis, the algorithm will indicate that there is no consistent hypothesis.

Algorithm 1 Pseudocode:

```

1: procedure MYPROCEDURE
2:    $result \leftarrow (x_1 = 1) \wedge (x_2 = 1) \wedge \dots \wedge (x_n = 1) \wedge (x_1 = 0) \wedge (x_2 = 0) \wedge \dots \wedge (x_n = 0)$ 
3:   for each '+' sequence  $S$  in training set do
4:     if  $(x_i = v) \neq S_i$  then
5:       remove  $(x_i = v)$  from  $result$ 
6:     end if
7:   end for
8:   for each '-' sequence  $S$  in training set do
9:     if  $S$  contradict with  $result$  then
10:      Indicate: no consistent hypothesis
11:      halt
12:    end if
13:   end for
14:   return  $result$ 
15: end procedure

```

- b. Inorder to prove the correctness of this algorithm, I will prove that there are neither **false positive** nor **false negative**.
- b.1 If there is a **false positive**, there are at least one condition absent from $result$, say $(x_i = v)$. However, in the first loop, the algorithm eliminated $(x_i = v)$ from complete conjunction, because it contradicts with other positive examples. Those two assumptions are mutually conflicting.
- b.2 If there is a **false negative**, there are at least one incorrect condition in $result$, say $(x_i = v)$. However, in the first loop, the algorithm eliminated all conditions that contradict with other positive examples. The remaining

$(x_i = v)$ should not contradict with any positive sequence. Therefore, those two assumptions are mutually conflicting.

Therefore, the *result* contains neither **false positive** nor **false negative**. The algorithm is correct.

- c. The algorithm traverses m training examples and checks each example with hypothesis *result* that has $2n$ conditions at worst. The runtime should be $\mathbf{O}(mn)$.
- d. There might be **false negative** in this situation. The example labeled as positive must be positive. However, the example labeled as negative might be positive, because the algorithm updates the hypothesis according to positive examples in training examples. The hypothesis changed from general to specific. The hypothesis *result* is the superset of real hypothesis \mathbf{H} .

2. Answer to problem 2

- a. Assume a vector \vec{w} that is perpendicular to the hyperplane. Point $x^* = (x_1^*, x_2^*, \dots, x_n^*)$ stands for any point in hyperplane. Define \vec{v} as vector from point $x^0 = (x_0^0, x_1^0, \dots, x_n^0)$ to point x^* . (I use x^* and x^0 to express the two points. $*$ and 0 are just labels.)

$$\vec{v} = [x_0^* - x_0^0, \quad x_1^* - x_1^0, \quad \dots, \quad x_n^* - x_n^0]^T$$

Then, the length D of the projection of \vec{v} onto \vec{w} is the distance from the point to the plane.

$$D = |\text{proj}_{\vec{w}} \vec{v}| \tag{1}$$

$$= \frac{|\vec{w}^T \cdot \vec{v}|}{\|\vec{w}\|} \tag{2}$$

$$= \frac{|w_0 \cdot (x_0^* - x_0^0) + w_1 \cdot (x_1^* - x_1^0) + \dots + w_n \cdot (x_n^* - x_n^0)|}{\sqrt{w_0^2 + w_1^2 + \dots + w_n^2}} \tag{3}$$

$$= \frac{|w_0 \cdot x_0^* - w_0 \cdot x_0^0 + w_1 \cdot x_1^* - w_1 \cdot x_1^0 + \dots + w_n \cdot x_n^* - w_n \cdot x_n^0|}{\|\vec{w}\|} \tag{4}$$

$$= \frac{|\theta + w_0 \cdot x_0^* + w_1 \cdot x_1^* + \dots + w_n \cdot x_n^*|}{\|\vec{w}\|} \tag{5}$$

$$= \frac{|\vec{w}^T \cdot \vec{x}^* + \theta|}{\|\vec{w}\|} \tag{6}$$

- b. The distance D between the hyperplanes:

$$\vec{w}^T \cdot \vec{x} + \theta_2 = 0 \tag{7}$$

$$\vec{w}^T \cdot \vec{x} = -\theta_2 \tag{8}$$

$$D = \frac{|\vec{w}^T \cdot \vec{x} + \theta_1|}{\|\vec{w}\|} \tag{9}$$

$$= \frac{|\theta_1 - \theta_2|}{\|\vec{w}\|} \tag{10}$$

3. Answer to problem 3

- a. a.1 For this problem, I will prove that the data set D is linearly separable if (a.1.1) *there exists a hyperplane that satisfies condition (3) with $\delta = 0$* , and (a.1.2) *only if there exists the hyperplane*. Then I will consider the case (a.1.3) *that $\delta > 0$* .

(a.1.1) if hyperplane with $\delta = 0$, then dataset D is linear separable:

- For case $y_i = -1$, $\delta = 0$:

$$y_i(\vec{w}\vec{x}_i + \theta) \geq 1 - \delta \quad (11)$$

$$-1 \cdot (\vec{w}\vec{x}_i + \theta) \geq 1 - 0 \quad (12)$$

$$(\vec{w}\vec{x}_i + \theta) \leq -1 \quad (13)$$

$$(\vec{w}\vec{x}_i + \theta) \leq 0 \quad (14)$$

$$(15)$$

- For case $y_i = 1$, $\delta = 0$:

$$y_i(\vec{w}\vec{x}_i + \theta) \geq 1 - \delta \quad (16)$$

$$1 \cdot (\vec{w}\vec{x}_i + \theta) \geq 1 - 0 \quad (17)$$

$$(\vec{w}\vec{x}_i + \theta) \geq 1 \quad (18)$$

$$(\vec{w}\vec{x}_i + \theta) \geq 0 \quad (19)$$

$$(20)$$

The result is consistent with condition (1). Therefore, there exists linear separable dataset D if there exists a hyperplane that satisfies condition (3) with $\delta = 0$.

- (a.1.2) if dataset D is linear separable, then there exists hyperplane with $\delta = 0$:
We need to prove that:

$$(\vec{w}\vec{x}_i + \theta) \geq 0 \Rightarrow y_i = 1 \quad (21)$$

$$(\vec{w}\vec{x}_i + \theta) \leq 0 \Rightarrow y_i = -1 \quad (22)$$

We assume that there is a hyperplane that separates the dataset and touches a positive point. We define D as the distance between this hyperplane and the closest negative point and \vec{v} as vector.

$$D = \frac{|\vec{w}^T \cdot \vec{x}_i + \theta|}{\|\vec{w}\|} \quad (23)$$

$$(\vec{w}\vec{x}_i + \theta) = 0 \quad (24)$$

We shift the hyperplane by $D/2$ along \vec{v} . The hyperplane is still linear separable. After shifting $D/2$, $(\vec{w}\vec{x}_i + \theta) > 0$. Scaling the hyperplane by multiplying a factor in both sides. We get equation: $(\vec{w}\vec{x}_i + \theta) \geq 1$ and $(\vec{w}\vec{x}_i + \theta) \leq 1$. This is consistent with condition (3) with $\delta = 0$.

(a.1.3) For case $\delta > 0$:

- if $\delta < 1$: we can use similar trick to prove that the dataset is linear separable.
- if $\delta \geq 1$: the dataset is not separable.

a.2 The trivial solution is:

$$\vec{w} = [0_0, 0_1, \dots, 0_n], \theta = 0, \delta = 0$$

This is the optimal solution. $\theta = 0$ is the smallest θ that minimizes the δ with condition $1 - \delta \geq 0$. To avoid this trivial solution, we add a positive constant to make sure there is a hyperplane.

a.3 By plugging in \vec{x}_1 and \vec{x}_2 :

$$\sum_{i=1}^n \vec{w}_i + \theta \geq 1 \quad (25)$$

$$\sum_{i=1}^n -\vec{w}_i + \theta \leq 1 \quad (26)$$

b. b.1 **findLinearDiscriminant.m**

```

findLinearDiscriminant.m x
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 y = data(:, np1);
13 features = data(:, 1:n);
14 sign = zeros(m, n);
15 for i = 1:m
16     sign(i,:) = y(i) * features(i,:);
17 end
18 c = [zeros(np1,1); 1];
19 A = [sign, y, ones(m,1); zeros(1,np1), 1];
20 b = [ones(m,1); 0];
21 %% solve the linear program
22 %adjust for matlab input: A*x <= b
23 [t, z] = linprog(c, -A, -b);
24
25 %% obtain w,theta,delta from t vector
26 w = t(1:n);
27 theta = t(n+1);
28 delta = t(n+2);
29
30 end

```

b.2 hw1sample2d.txt

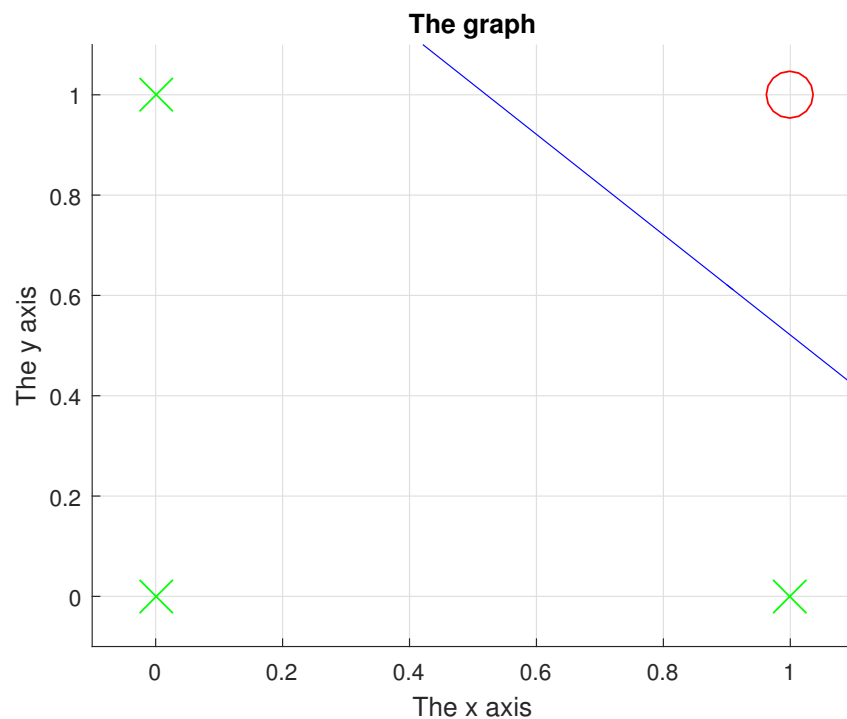
```
hw1sample2d.txt
1  0 0 -1
2  0 1 -1
3  1 0 -1
4  1 1 1
5
```

The dataset represents a monotone conjunction over two variables.

plot2dSeparator.m

```
plot2dSeparator.m
1 % This function plots the linear discriminant.
2 % YOU NEED TO IMPLEMENT THIS FUNCTION
3
4 function plot2dSeparator(w, theta)
5 xdata = -10:0.01:10;
6 y = -w(1)/w(2) * xdata - theta/w(2);
7 plot(xdata, y, 'b');
8 end
9
```

two dimension separator:



$$\theta = -90.2115, \delta = -1.5632e - 13$$

$$\vec{w} = [2.910, -2.050, 0.178, 190.520, 0.140, -3.101, -2.953, -193.278, 1.168, -8.894]$$

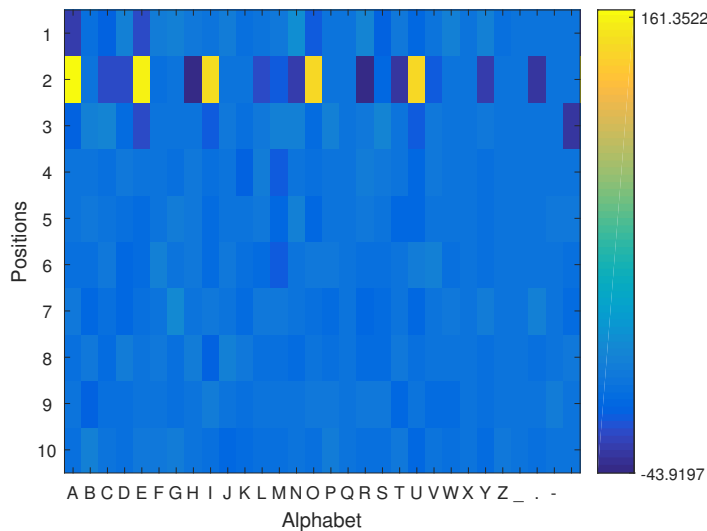
In this solution, $\delta \approx 0$ shows that the dataset is linearly separable. Since x_4 and x_8 are significantly larger than others, the label is mostly depend on x_4 and $\neg x_8$. The hyperplane should be $x_4 \wedge \neg x_8$. In this solution, $\theta \ll 0$ guarantees a negative output when $\vec{w}\vec{x}$ is small.

b.3 computeLabel.m

```

computeLabel.m
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  if w'*x + theta >= 0
8      y = 1;
9  else
10     y = -1;
11 end
12 end
13

```

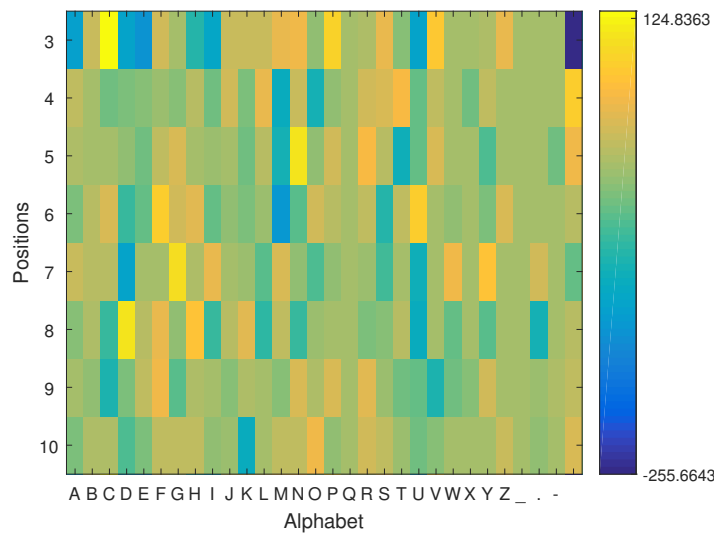


($\delta = 1.2619e - 09$, $accuracyInTrain = 1$, $accuracyInTest = 1$)

Since $\delta \approx 0$, there exists a hyperplane that separates dataset. This set of features shows a perfect accuracy in both training dataset and test dataset.

change the position from 1:10 to 3:10:

This change will ignore the first two positions, so the obvious yellow cube in image above will be ignored.



($\delta = 1.2179e - 11$, $accuracyInTrain = 1$, $accuracyInTest = 0.7234$)

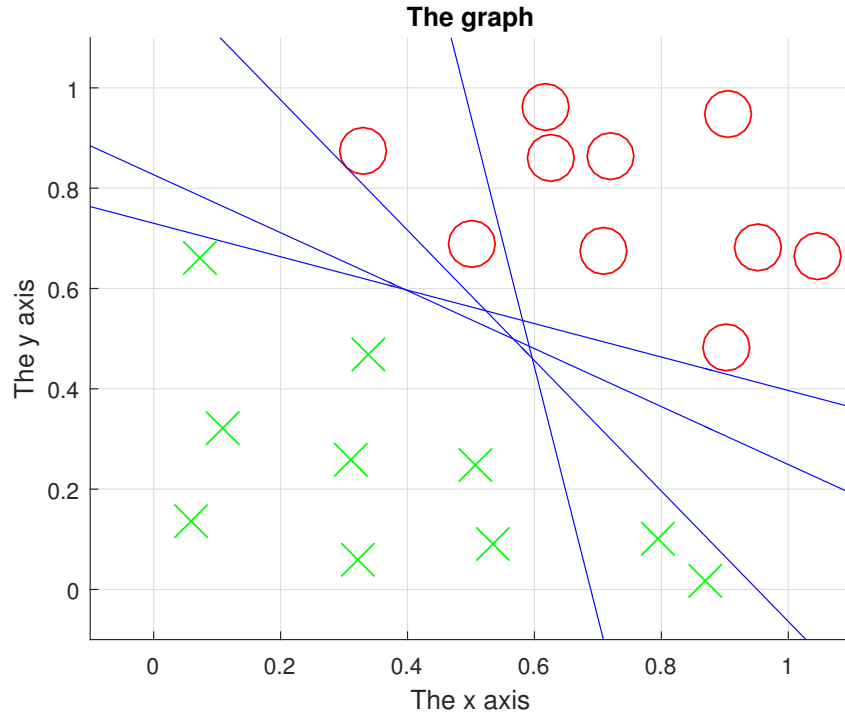
Since $\delta \approx 0$, the new hyperplane separates the dataset. The accuracy for training dataset is 1, but the accuracy for test is only 0.7. There is no obvious feature under this alphabet and position.

b.4 findLinearThreshold.m

```

findLinearThreshold.m
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 y = data(:, np1);
12 features = data(:, 1:n);
13 sign = zeros(m, n);
14 for i = 1:m
15     sign(i,:) = y(i) * features(i,:);
16 end
17 c = [zeros(np1,1); 1];
18 A = [sign, y, ones(m,1); zeros(1,np1), 1];
19 b = [ones(m,1); 0];
20 %% solve the linear program
21 %adjust for matlab input: A*x <= b
22 [t, z] = linprog(c, -A, -b, [], [], [w' -inf -inf], [w' inf inf]);
23
24
25 %% obtain w,theta,delta from t vector
26 w = t(1:n);
27 theta = t(n+1);
28 delta = t(n+2);
29
30 end

```



($\delta_1 = 1.0658e - 12$, $\delta_2 = -2.2737e - 13$, $\delta_3 = -2.0009e - 11$, $\delta_4 = 92.9650$)

Among those four hypothesis, only three hypothesis successfully separates the dataset, ($\delta_4 \gg 0$). The lines which left end located between 0.8 and 1 separates the dataset exactly. Since $\delta_4 \gg 0$, the most vertical line fails to separate the dataset. Other two lines are closer to positive or negative examples. We can infer that there are multiply lines that can separate the dataset.