

## Homework #5

### 5.1 Powerfulness of Neural Networks

- (1) We implement  $f_A(x) = g_A(x|\theta, w_1, w_2, \dots, w_d)$ , with  $W = (\theta, w_1, w_2, \dots, w_d)$ .

$$W_{i1}^{(1)} = \begin{cases} d-1, & \text{if } i = 0 \\ 1, & \text{otherwise} \end{cases}$$

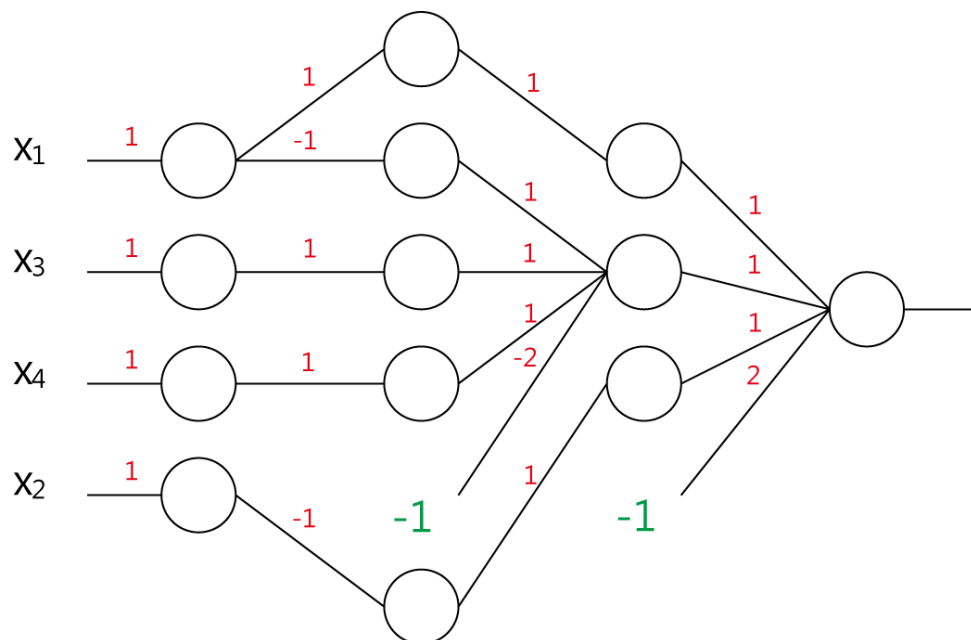
To prove it can work, we show the function  $g_A(x) = 1$  if and only if all inputs  $x$  are 1. The statement is satisfied because the truth that we only have  $d$ 's inputs argument, each of them either 1 or -1. All of inputs should be 1 and sum up to greater than  $d-1$ . Otherwise, the function  $g_A(x) = -1$ .

- (2) We implement  $f_O(x) = g_O(x|\theta, w_1, w_2, \dots, w_d)$ , where

$$W_{i1}^{(1)} = \begin{cases} 1-d, & \text{if } i = 0 \\ 1, & \text{otherwise} \end{cases}$$

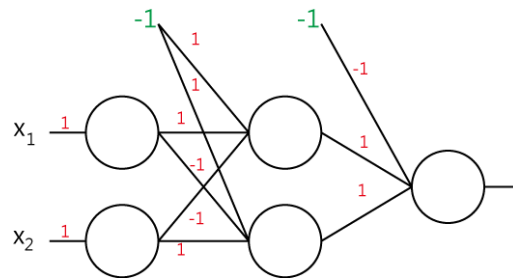
To prove it works, we show the function  $g_O(x) = -1$  if and only if all inputs  $x$  are -1. The statement is satisfied with the similar reason. All of inputs should be -1 and sum up to smaller than  $1-d$ . Otherwise, the function  $g_O(x) = 1$ .

- (3)  $f_N(x) = g_N(x|\theta, w_1)$ , where  $(\theta, w_1) = (0, -1)$ . The result is trivial.
- (4) Illustrated the result by the combination of  $f_A$ ,  $f_O$  and  $f_N$  perceptron we define above as to a Neural Networks (NN) (5.1.1).



(5.1.1)

- (5) By (1), (2) and (3). We know that the 3 element gates can be generated as a NN form. For any Boolean system, we can change it to conjunctive normal form. Thus, we can systematically implement any Boolean function by back propagate each gates to a NN by Boolean operation.
- (6) Try the simpler case  $XOR((x)_1, (x)_2) = OR(AND(NOT((x)_1), (x)_2), AND((x)_1, NOT((x)_2)))$  and get the figure (5.1.2)



(5.1.2)

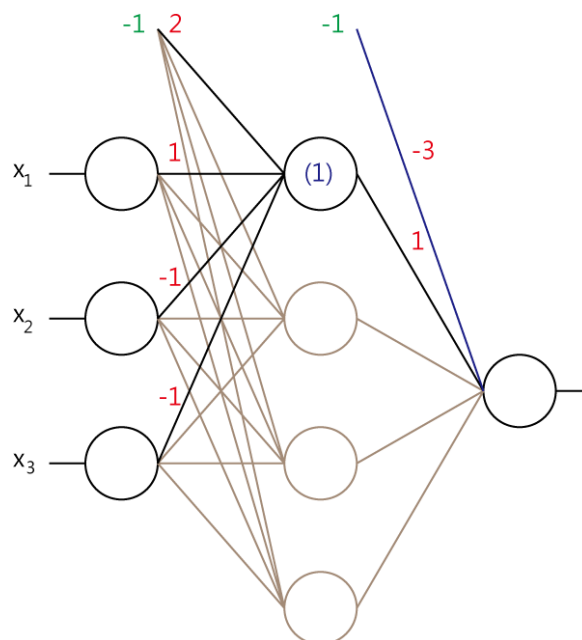
Use this form we can derive the Boolean equation (5.1.3) that implement a 3 inputs  $XOR$  with a hidden layer for 4 perceptrons.

$$XOR((x)_1, (x)_2, (x)_3) =$$

$$OR(AND((x)_1, \overline{(x)_2}, \overline{(x)_3}), AND(\overline{(x)_1}, (x)_2, \overline{(x)_3}), AND(\overline{(x)_1}, \overline{(x)_2}, (x)_3), AND((x)_1, (x)_2, (x)_3)) \quad (1)$$

(5.1.3)

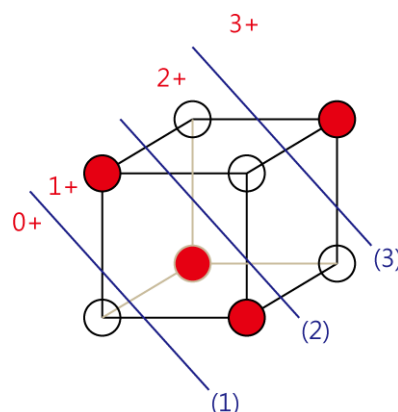
We can get the figure (5.1.4) by the equation (5.1.3)



(5.1.4)

To simplify figure, we use black edges to represent the  $AND$  gate structures, and the gray edges are similar to black edges but notice the position of  $NOT$  operation.

- (7) Figure (5.1.5) show that we only need 3 perceptrons to separate the hypercube in three dimensions.



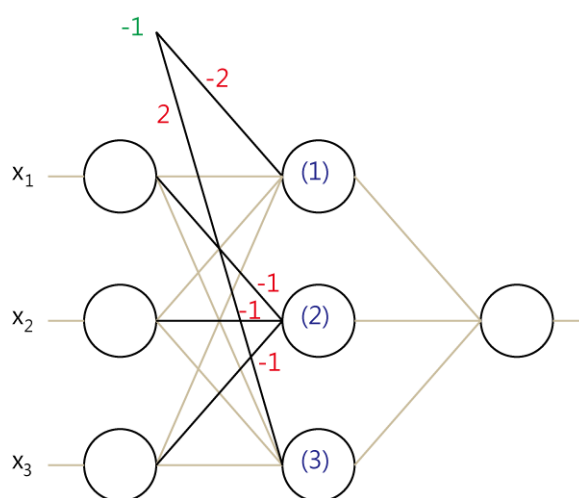
(5.1.5)

The formal way to represent our finding is using the table (5.1.6)

Patterns	Perceptron (1)	Perceptron (2)	Perceptron (3)	SUM	
+++	+1	−1	+1	+1	
++−			−1	−1	
+−+					+1
−++					
+−−		−1			
−+−				+1	
−−+					
−−−	−1				−1

Table (5.1.6)

Our implement as picture (5.1.7), note that all edge without notification weighted 1



(5.1.7)

(8) To extend the result we found at table (5.1.6). We try the case in four dimensions

Patterns	1+↑	1+↓	3+↑	3+↓	SUM − 1			
++++	+1	−1	+1	−1	−1			
+++−				+1	+1			
++−+								
+−++								
−+++								
++−−		+1	−1		−1			
+−−+								
−−++								
−++−								
−+−+								
+−+−								
−−−+		−1			+1			
−−+−								
−+−−								
+−−−								
−−−−	−1				−1			

Table (5.1.8)

We can add -1 on the SUM column, and can get the desire output we want. To formalize the table, we combine the finding at table (5.1.6) and (5.1.8):

(a) If the  $d$  is odd number:

The table is constructed by some rules:

- 1 The perceptrons given with that more than odd number +’s fire a positive one, and less than even number +’s fire a negative one. Where odd numbers are belong to  $\{1, 3 \dots d\}$ .
- 2 Sum up the table values as output.

(b) If the  $d$  is even number:

The table is constructed by some rules:

- 1 The perceptrons given with that more than odd number +’s fire a positive one, and less than even number +’s fire a negative one. Where even numbers are belong to  $\{2, 4 \dots d\}$ .
- 2 Sum up the table values with a  $-1$  as output.

The way we transform the table to neural networks is by a simple observation on linear combination of inputs value. With proper weight, the outputs value will be monotonic function. So the neural network is existed.

Consider about  $k$ 's perceptron in  $d$  dimensions. If  $k$  is an odd number, than we fire the input more than and equal to  $k$ 's + in +1. Give the weight  $(w_1, w_2, \dots, w_d) = (1, 1, \dots, 1)$ . We have the inequality (5.1.9)

$$\begin{aligned} k - (d - k) - C_{odd} &\geq +1 \\ C_{odd} &\leq 2k - d - 1 \end{aligned} \quad (5.1.9)$$

If  $k$  is an even number, than we fire the input less than  $k$ 's + in +1. Give the weight  $(w_1, w_2, \dots, w_d) = (-1, -1, \dots, -1)$ . We have the inequality (5.1.10)

$$\begin{aligned} -(k - 1) + (d - (k - 1)) - C_{even} &\geq +1 \\ C_{even} &\leq d - 2k + 1 \end{aligned} \quad (5.1.10)$$

Therefore we have perceptrons:

- Hidden Layer:
  - 1 Perceptron( $k = 0 \bmod 2$ ) =  $g_{even}(x|d - 2k + 1, -1, -1, \dots, -1)$
  - 2 Perceptron( $k = 1 \bmod 2$ ) =  $g_{odd}(x|2k - d - 1, 1, 1, \dots, 1)$
- Outputs Layer:
  - 1 If  $d = 0 \bmod 2$ :  
Perceptron(1) =  $g(x|1, 1, 1, \dots, 1)$
  - 2 Else  $d = 1 \bmod 2$ :  
Perceptron(1) =  $g(x|0, 1, 1, \dots, 1)$

For a conclusion, we get the  $d$ - $d$ -1 neural networks.

$$W_{ij}^{(1)} = \begin{cases} (2j - d - 1)(-1)^{j+1}, & \text{if } i = 0 \\ (-1)^{j+1}, & \text{otherwise} \end{cases}$$

$$W_{ij}^{(2)} = \begin{cases} d + 1 \bmod 2, & \text{if } i = 0 \\ 1, & \text{otherwise} \end{cases}$$

## 5.2 Limitations of Neural Networks

- (1) Express the formula at left hand side inequality

$$(N^D + 1)^3 = N^{3D} + N^{2D} + N^D + 1 \leq 3N^{3D} + 1 \leq N^{3D+1} + 1 \quad (5.2.1)$$

The inequality is satisfy for  $D \geq 1$  and  $N \geq 3$ .

- (2) Proof  $2^N > N^k + 1$ , if  $N \geq 3k \log_2 k$ , assume  $k = \Delta \geq 1$  and  $N > 1$ .

(a) If  $k = 1$ ,  $2^N > N + 1, \forall N \in \mathbb{N}$

(b) If  $k \geq 2$ , Consider a function  $f$  such that

$$f(n) = n - k \log_2 n - \frac{1}{2} \quad (5.2.2)$$

-1 If  $n > k \log_2 e$ ,

$$f'(n) = 1 - \frac{k}{n} \log_2 e > 0 \quad (5.2.3)$$

-2 If  $n = 3k \log_2 k$ ,

$$\begin{aligned} f(n) &= 3k \log_2 k - k(\log_2 3 + \log_2 k + \log_2 \log_2 k) - \frac{1}{2} \\ &= k(2 \log_2 k - \log_2 3 - \log_2 \log_2 k) - \frac{1}{2} \\ &= k \left( \log_2 \frac{k^2}{3 \log_2 k} \right) - \frac{1}{2} > 0 \end{aligned} \quad (5.2.4)$$

By inequality (5.2.3) and (5.2.4), and  $3k \log_2 k \geq k \log_2 e$ . We get the inequality

$$f(n) > 0, \text{ for } n \geq 3k \log_2 k, k \geq 2 \quad (5.2.5)$$

- (c) Use the inequality (5.2.5), we have

$$\begin{aligned} f(n) &> 0 \\ &\rightarrow n > k \log_2 n + \frac{1}{2} \\ &\rightarrow 2^n > \sqrt{2} n^k > n^k + 1 \end{aligned} \quad (5.2.6)$$

Proof is done.

(3) Suppose we have  $N$  nodes, assume

$$N \geq 3(3(d+1) + 1) \log_2(3(d+1) + 1) \geq 3 \quad (5.2.7)$$

By proof of (2), with  $\Delta = 3(d+1) + 1 > 1$ , we get

$$N^\Delta + 1 < 2^N \quad (5.2.8)$$

By proof of (1), with  $D = d + 1 > 1$  and  $N \geq 3$ , we get

$$(N^D + 1)^3 \leq N^\Delta + 1 < 2^N \quad (5.2.9)$$

We know that several edges of our  $d$ -3-1 neural networks has fixed. We can only change three perceptrons of hidden layer. Each of them has  $D$  inputs.

By problem 3.2, we know that the

$$\text{VC dimension of linear model with 'sign gate' with } d\text{'s input is } D. \quad (5.2.10)$$

Suppose each perceptron can generate  $L_{GA}(N)$  patterns. We have the inequality

$$L_{GA}(N) \leq B(N, D + 1) \leq N^D + 1 \quad (5.2.11)$$

Therefore, three independent perceptrons can generate  $L_{3GA}(N)$  patterns but less than  $2^N$

$$L_{G3A}(N) \leq (N^D + 1)^3 \leq 2^N \quad (5.2.12)$$

It means we cannot generate all patterns. So we have

$$VCD(G_{3A}) < 3(3(d+1) + 1) \log_2(3(d+1) + 1) \geq 3 \quad (5.2.13)$$

(4) Neural networks can implement any Boolean function and  $VCD(G_M)$  is bounded, this statement is true. For a given dimensions  $d$ , our Boolean functions only can produce  $2^d$ 's points at most, which mean our maximum  $VCD(G_M)$  is bounded by  $d$ .

- (5) Let  $\mathcal{S}$  be some set of  $N$  input vectors for the neural networks  $NN$ . By (5.2.10), we know that a perceptron  $G$  can only compute at most  $|\mathcal{X}|^D + 1$  different functions from any finite set  $\mathcal{X} \subset \mathbb{R}^d \rightarrow \{0,1\}$ . Denote the total number of weight  $w = M * d$ .

Consider the relation as a mapping, so our  $NN$  generate at most have

$$\prod_G (N^D + 1) \leq N^{2w} \quad (5.2.14)$$

different functions from  $\mathcal{S}$  into  $\{0,1\}$ .

Suppose we yield all possible patterns, which means  $2^N$  functions. By (5.2.14) we get

$$\prod_G (N^D + 1) = 2^N \leq N^{2w} \quad (5.2.15)$$

It imply

$$N \leq 2w * \log_2 N \quad (5.2.16)$$

Consider about big-O notation, we claim that

$$\log_2 N = O(\log w) \quad (5.2.17)$$

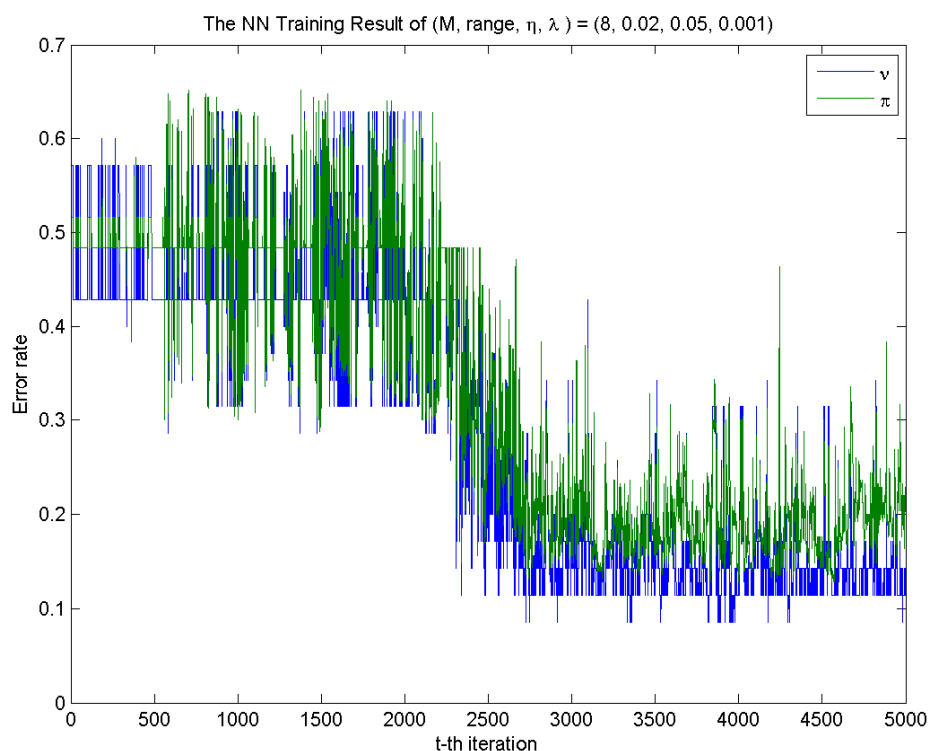
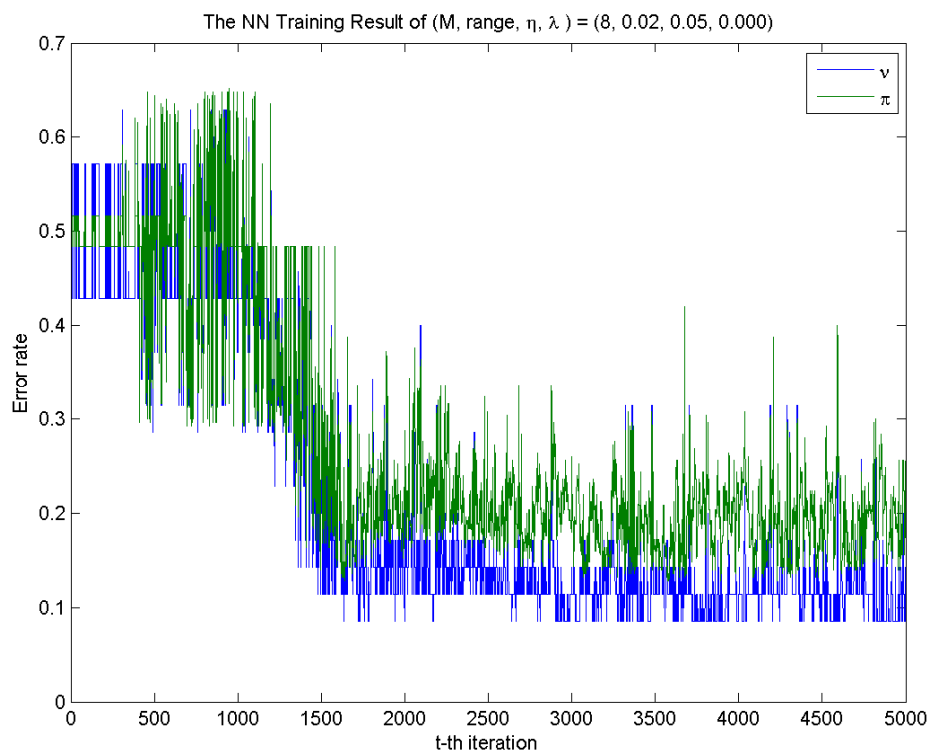
Combine (5.2.16) and (5.2.17), we have the result

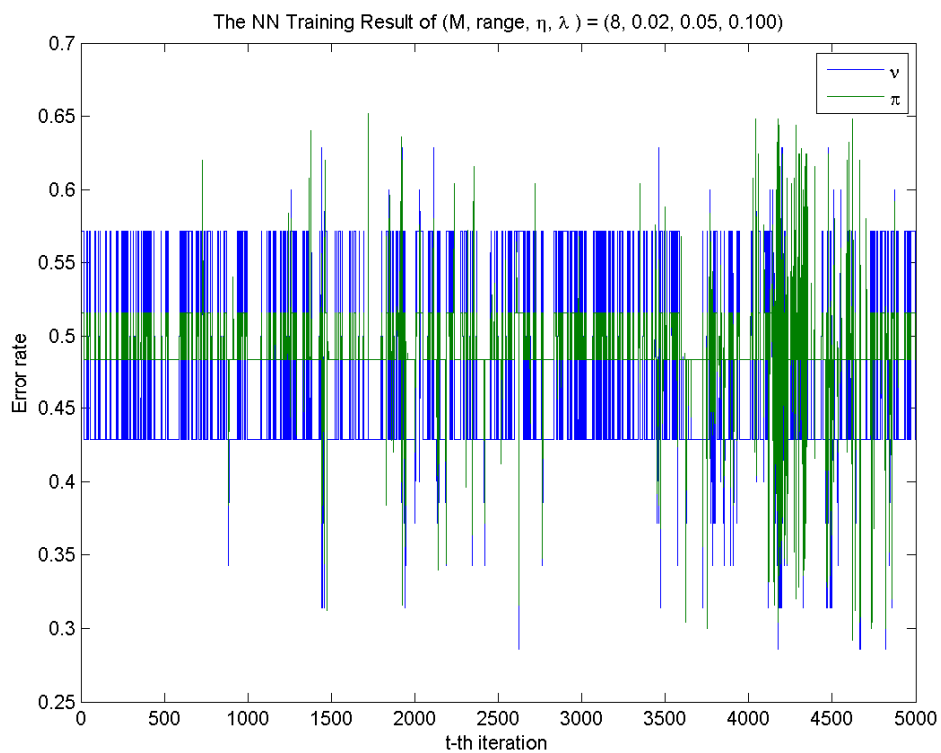
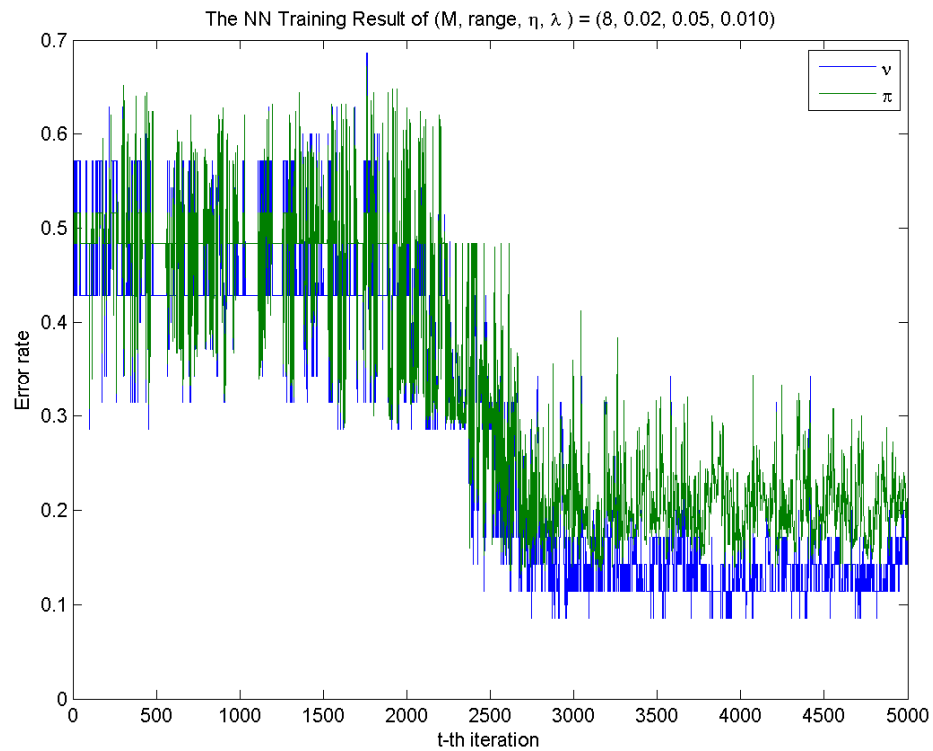
$$VCD(G_M) = N = O(w * \log w) = O(M * d * \log(M * d)) \quad (5.2.18)$$

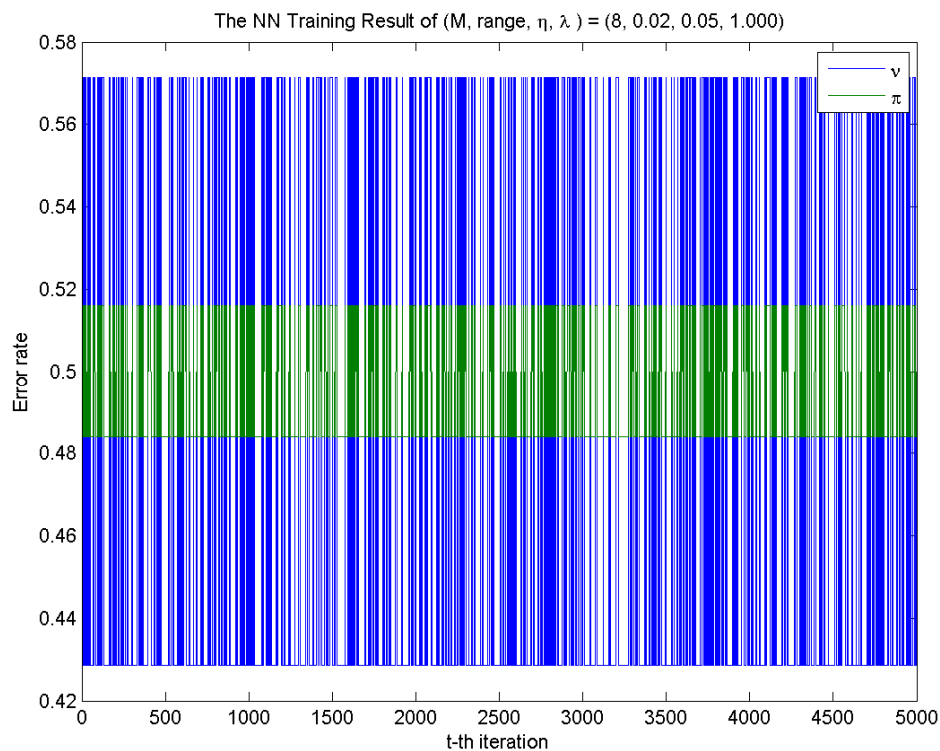


## 5.3 Experiment with Weight Decay in Neural Network

- (1) The implement of the algorithm isn't very complicated, but should be careful about one thing, that we have to change the  $\delta^2$  term of back propagation. Plot the figure as fellow,







Conclusion:

- 1 Error rate doesn't change much at  $K = \{0.001, 0.01\}$ . But we can see the curve kind of border when  $K$  got smaller.
- 2 When  $K = 0.1$  or higher, the curve cannot converge at  $T = 5000$ .

## 5.4 Experiment with Naïve Bayes

- (1) The results show as Table (5.4.1), (5.4.2) and (5.4.3)

	$K = 2$	$K = 5$	$K = 10$	$K = 20$	$K = 50$
$v(g_K)$	0.22	0.17	0.12	0.11	0.02
$\hat{\pi}(g_K)$	0.211	0.204	0.211	0.234	0.301

Table (5.4.1)

	$K = 2$	$K = 5$	$K = 10$	$K = 20$	$K = 50$
$v(g_K)$	0.22	0.17	0.12	0.11	0.02
$\hat{\pi}(g_K)$	0.211	0.204	0.215	0.267	0.384

Table (5.4.2)

	$K = 2$	$K = 5$	$K = 10$	$K = 20$	$K = 50$
$v(g_K)$	0.22	0.17	0.12	0.11	0.02
$\hat{\pi}(g_K)$	0.211	0.204	0.211	0.238	0.299

Conclusion:

- 1 It may not good for a big  $K$ , looks like an over fitting result.
- 2 For a testing data in which never showed up our training bin, it may raise a problem.  
We don't know how to classify such case, so we sign it as a guess.
  1. Guess it is  $+1$ , we get Table (5.4.1).
  2. Guess it is  $-1$ , we get Table (5.4.2).
  3. Guess it by random guess with function,  $\text{sign}(\text{random} \in (-0.1, 0.9))$ , we get Table (5.4.3).