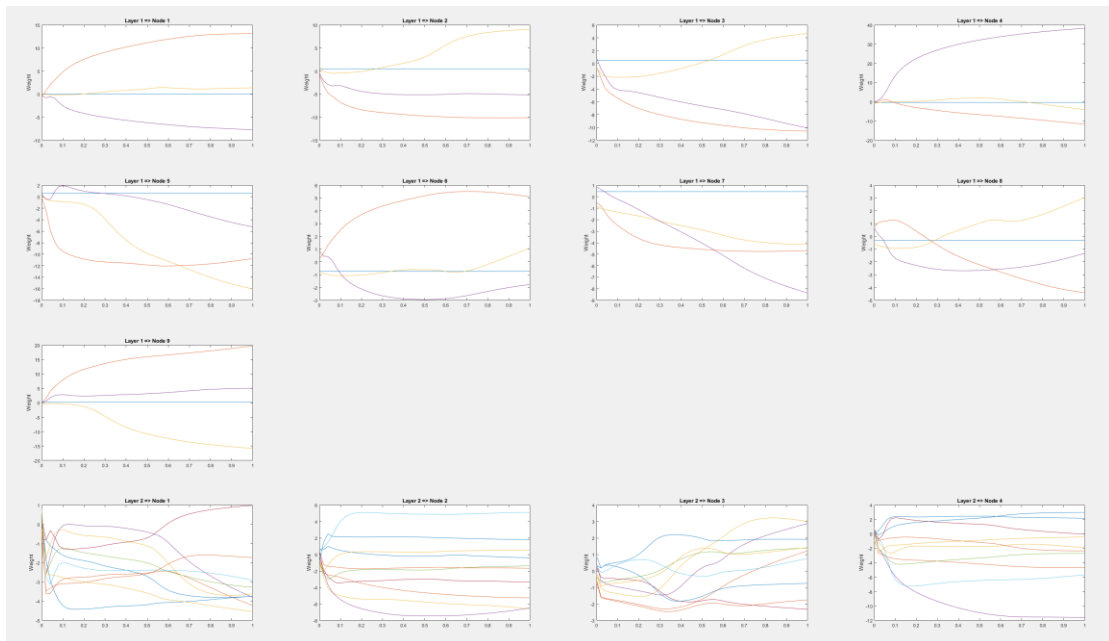
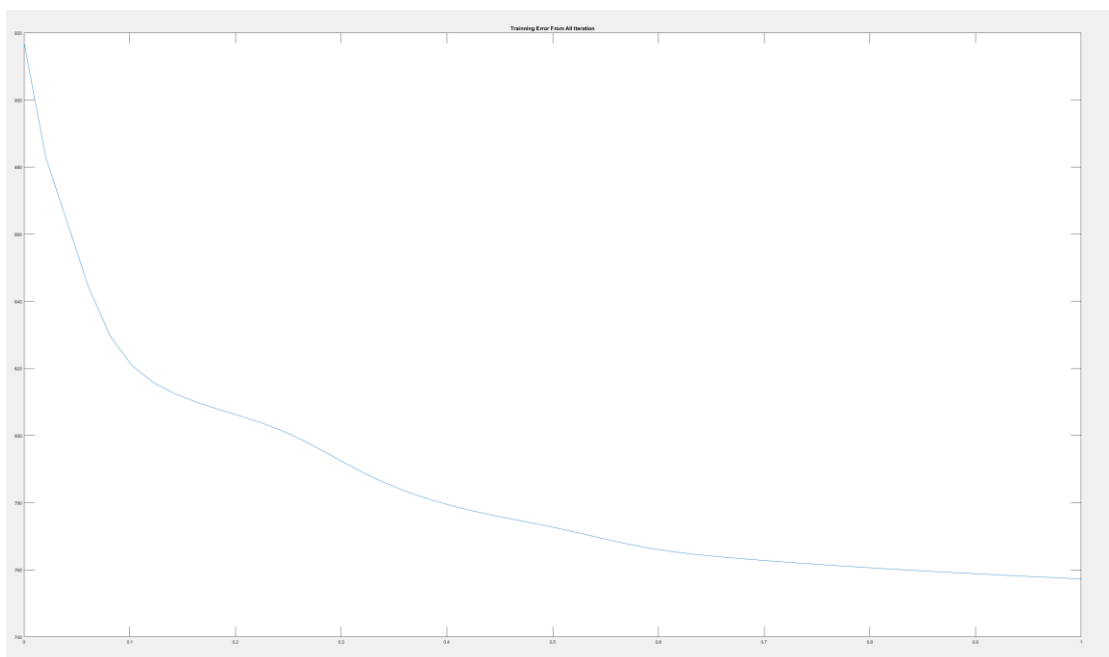


1. Only Training have iterations, so Training Weight, Error are based on iteration. Training Output is based on the the single minimized error round.
For the testing case, there's no iteration need. Also, weight does not change. So there are only plot for single minimized error round Error and Output.

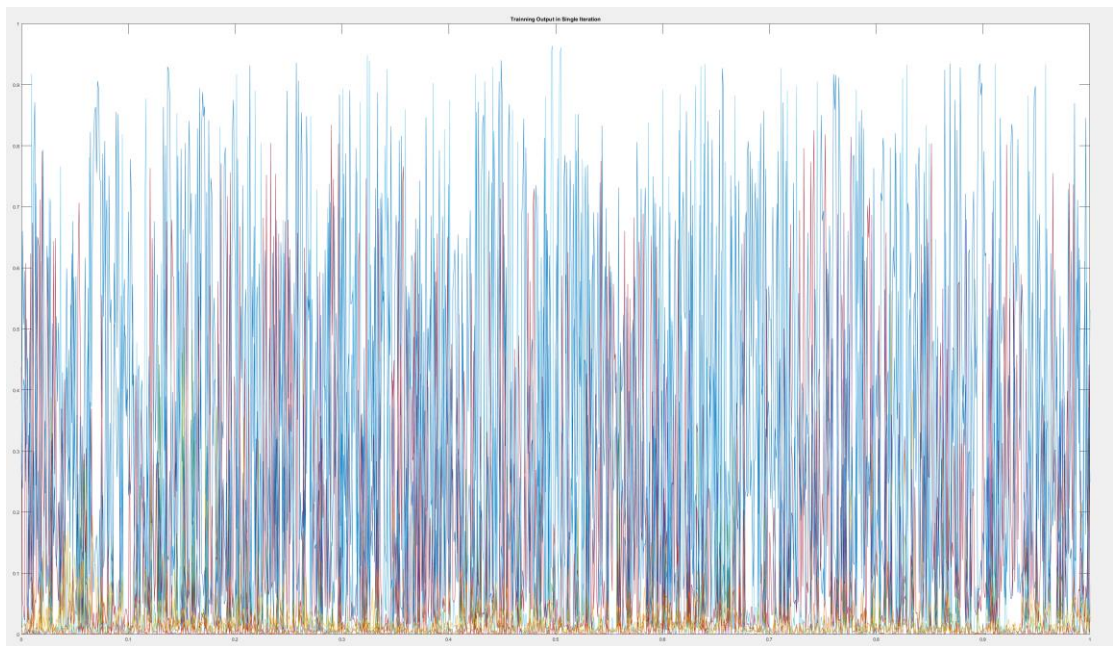
Training Weight Graph



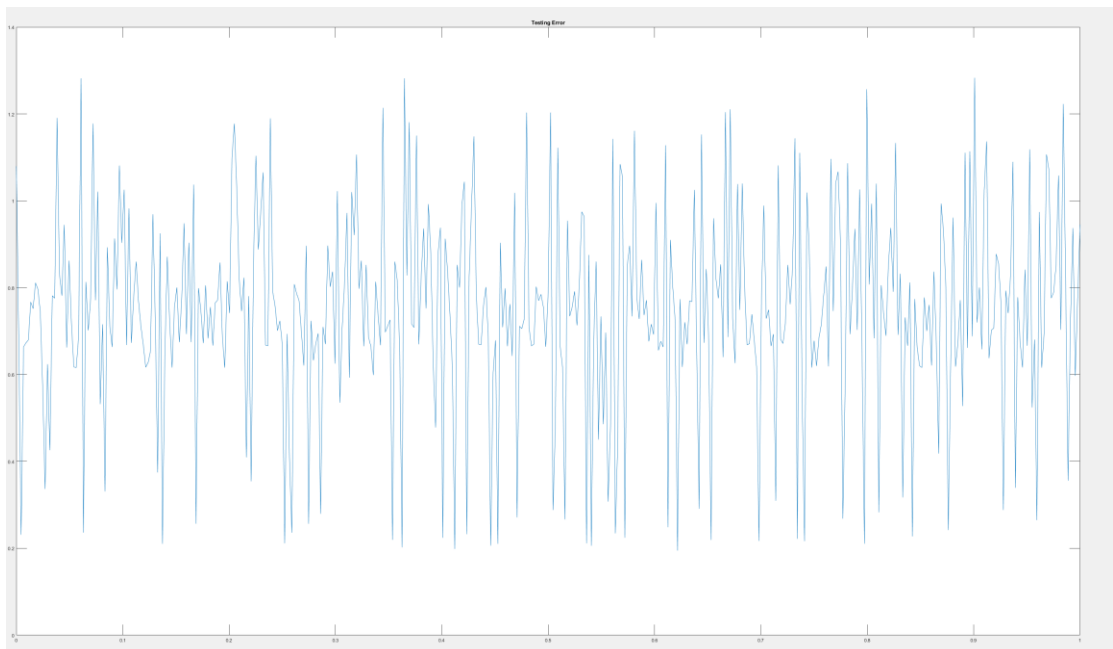
Training Error Graph



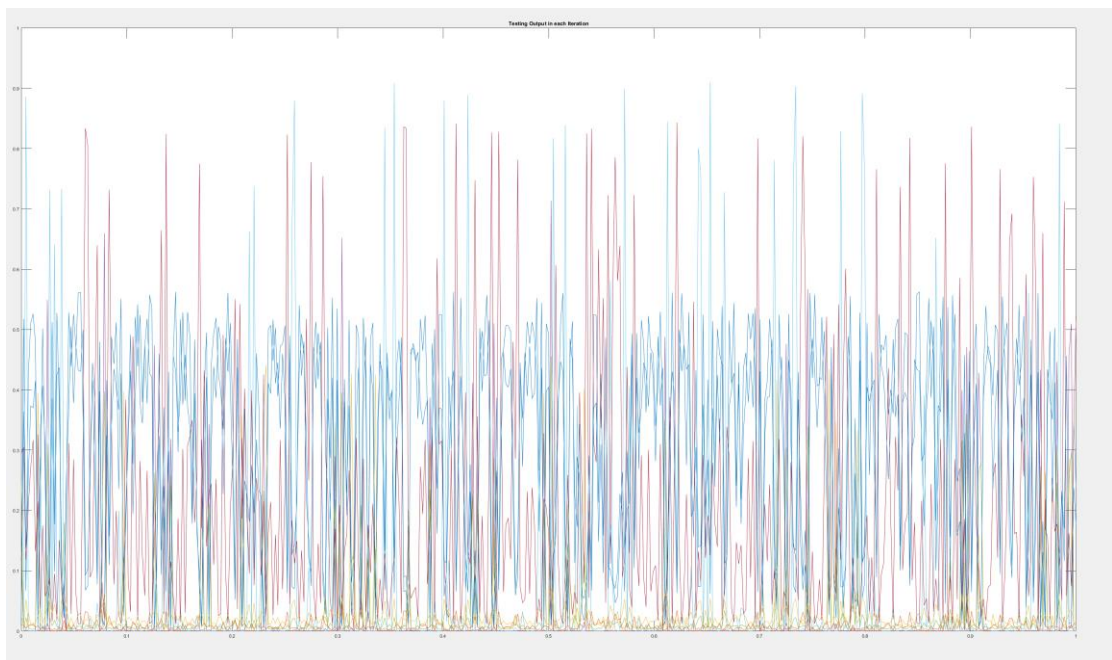
Training Output Graph



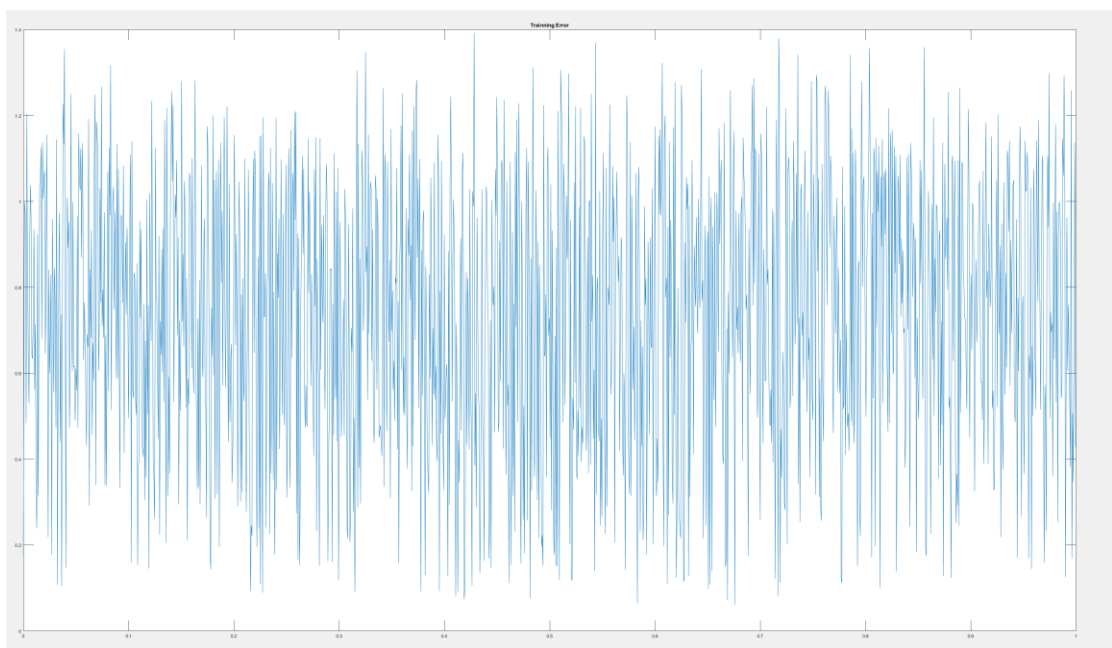
Testing Error Graph



Testing Output Graph



2. Training Error



Testing Training Error is : 1.094755e+03

Final Activation Function is : 0.1250 0.0000 0.0000 0.0001 0.0058 0.1363 0.0176 0.4829 0.0006 0.0052

3. Output From Code (First is Weight for first Layer, Second is Weight for Second Layer)

weight{1, 1}				
	1	2	3	4
1	0.6018	-0.1310	-1.0360	0.1442
2	0.8246	-0.7885	-0.2397	0.5892
3	-0.8556	0.6000	0.0325	0.4704
4	0.1665	0.4710	0.9629	0.5869
5	0.8063	-0.6864	-0.6578	-0.3123
6	-0.9170	0.5934	-0.2250	0.5619
7	-0.4790	-0.9973	0.6915	-0.3575
8	-0.1676	0.0285	0.8255	-0.0788
9	0.8780	0.7259	0.4954	0.3929

weight{1, 2}										
	1	2	3	4	5	6	7	8	9	10
1	0.2467	0.6043	-0.2333	0.5955	-0.5609	-1.1697	-0.7772	-1.0516	-0.0513	-0.1864
2	0.6864	-0.8291	0.0572	0.0206	-0.1564	-0.4832	-0.6069	0.2620	-0.4375	-0.3079
3	-0.2821	-0.5880	0.3947	-0.5448	-0.5691	0.1893	-0.9394	0.7024	0.5814	0.6418
4	0.5240	-0.9158	-0.2852	-0.4176	-0.8320	0.7769	0.2309	-0.4284	-0.7523	0.3512

Ecs 171 HW2 Question 3

3 Layers = $L_1 \Rightarrow X_0^1$ $L_2 \Rightarrow a_0^2$ $L_3 \Rightarrow a_0^3$

X_0^1
 X_1^1
 X_2^1
 X_3^1
 X_4^1
 X_5^1
 X_6^1
 X_7^1
 X_8^1

a_0^2
 a_1^2
 a_2^2
 a_3^2
 a_4^2
 a_5^2
 a_6^2
 a_7^2
 a_8^2
 a_9^2

a_0^3
 a_1^3
 a_2^3
 a_3^3
 a_4^3
 a_5^3
 a_6^3
 a_7^3
 a_8^3
 a_9^3

Where in Layer 1 and 2, We add a extra node for error (W_{i0})

Then, the Weight for first layer is W_{ji} with Dim $(p \times n)$ =

$$\begin{bmatrix} W_{00} & W_{01} & \dots & W_{0n} \\ W_{10} & W_{11} & \dots & W_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ W_{p0} & W_{p1} & \dots & W_{pn} \end{bmatrix}$$

i is index for Layer 2 node
 j is index for Layer 1 node

So Weight Matrix for first two Layer is

$W^{(1)} = \begin{bmatrix} 0.6018 & -0.1310 & -1.0360 & 0.1442 \\ 0.8246 & -0.7885 & -0.2397 & 0.5892 \\ -0.8556 & 0.6000 & 0.0325 & 0.4704 \\ 0.1665 & 0.4710 & 0.9629 & 0.5869 \\ 0.8063 & -0.6864 & -0.6578 & -0.3123 \\ -0.9170 & 0.5934 & -0.2250 & 0.5619 \\ -0.4790 & -0.9973 & 0.6915 & -0.3575 \\ -0.1676 & 0.0285 & 0.8255 & -0.0788 \\ 0.8780 & 0.7259 & 0.4954 & 0.3929 \end{bmatrix}$

$W^{(2)} = \begin{bmatrix} 0.2467 & 0.6043 & -0.2333 & 0.5955 & -0.5609 & -1.1697 & -0.7772 & -1.0516 & -0.0513 & -0.1864 \\ 0.6864 & -0.8291 & 0.0572 & 0.0206 & -0.1564 & -0.4832 & -0.6069 & 0.2620 & -0.4375 & -0.3079 \\ -0.2821 & -0.5880 & 0.3947 & -0.5448 & -0.5691 & 0.1893 & -0.9394 & 0.7024 & 0.5814 & 0.6418 \\ 0.5240 & -0.9158 & -0.2852 & -0.4176 & -0.8320 & 0.7769 & 0.2309 & -0.4284 & -0.7523 & 0.3512 \end{bmatrix}$

I randomly Initialized W to values $\in (-1, 1)$

$a_0^{(1)} = g(\sum_{i=0}^n W_{i0}^{(1)} x_i) = a_0^{(1)} = \text{Input}$, $a_1^{(1)} = 1$ for error first Soplein
 $[1, 0.34, 0.46, 0.17, 0.09, 0.05, 0.08, 0.27]$

$a^1 = [1, 0.34, 0.37, 0.46, 0.15, 0.5, 0, 0.5, 0.27]$

$a^2 = [0.6261, 0.5726, 0.4641, 0.7610]$

$a^3 = [0.6374, 0.0742, 0.3632, 0.3508, 0.2460, 0.4437, 0.2882, 0.4425, 0.4282, 0.6787]$

Then I set $a_i^{(i)}$ to 1 for error if Layer \neq final Layer

Then $\delta_j^{(L)} = a_i^{(L)} - y^{(L)}$ for final Layer

$\therefore \delta^{(L)} = [-0.3606, 0.2782, 0.5632, 0.5408, 0.2360, 0.3547, 0.2882, 0.3488, 0.4282, 0.6087]$

Then $\delta^{(i)}$, $i \neq \text{num-Layers}$, $\delta^{(i)} = (W^{(i+1)})^T \delta^{(i+1)} \cdot a^{(i)} \cdot (1 - a^{(i)})$

$\delta^{(2)} = [0, -0.1659, 0.2328, -0.0421]$

$\delta^{(1)} = [0, 0.0140, 0.0298, -0.0033, 0, 0.0176, -0.0043]$

$W_{ij}^{(L)} = W_{ij}^{(L)} - \alpha \left(\sum_{a=1}^n a_i^{(L)} \delta_i^{(L+1)} \right)$

Then $W^2 = \begin{bmatrix} 0.2867 & 0.6053 & -0.2093 & 0.4755 & -0.1607 & -1.1677 & -0.7772 & -1.0216 & -0.6413 & -0.1864 \\ 0.6867 & -0.2071 & 0.4472 & 0.4206 & -0.6444 & -0.4532 & -0.6467 & 0.2620 & -0.4376 & -0.3077 \\ -0.2821 & -0.0880 & 0.7747 & -0.0898 & -0.4571 & 0.1873 & -0.7317 & 0.7024 & 0.1814 & 0.6418 \\ 0.4240 & -0.708 & -0.2242 & -0.4476 & -0.2320 & 0.7867 & 0.2207 & -0.4287 & -0.7422 & 0.312 \end{bmatrix}$

$W^1 = \begin{bmatrix} 0.6018 & -0.1310 & -1.0340 & 0.1942 \\ 0.2146 & -0.7888 & -0.2777 & 0.1872 \\ -0.2816 & 0.6 & 0.0345 & 0.4204 \\ 0.1661 & 0.4710 & 0.1602 & 0.1887 \\ 0.8063 & -0.6864 & -0.6178 & -0.3033 \\ -0.9710 & 0.4734 & -0.2800 & 0.4417 \\ -0.4770 & -0.7773 & 0.675 & -0.371 \\ -0.1674 & 0.0281 & 0.481 & -0.0788 \\ 0.2780 & 0.7217 & 0.4707 & 0.3727 \end{bmatrix}$

Both are in agreement.

4. Rows are Number of Layers {1,2,3}, Columns are Number of Nodes per Layer {3,6,9,12}
The Second Graph is the corresponding Testing Classification Rate

all_test_error(1,1)

plot Plot as mult... PI

SELECTION

3x4 double

	1	2	3	4
1	333.0578	304.4382	315.0384	304.0126
2	360.9271	317.4250	321.5749	321.4009
3	393.3198	349.6258	332.0484	317.6584

all_test_percent(1,1)					
SELECTION					
3x4 double					
	1	2	3	4	
1	0.5281	0.5708	0.5191	0.5820	
2	0.3888	0.5191	0.5326	0.4944	
3	0.3011	0.4292	0.5056	0.5236	
4					

From the two matrix, the Optimal Configuration is 1 Hidden Layer with 12 Nodes per Hidden Layer.

The Relationship between these attributes can be summarized as: Generally, increase in Nodes per Hidden Layer will reduce the Error, while increase in Number of Layers might cause over-fitting which actually increase the Error in the test.

5.

Trainning percentage:

ans =

0.5512

Predicted Category is: NUC

6.

Exc 171 HW2 Q6

The uncertainty comes from the unknown sample data $X = \{X_1, \dots, X_m\}$
 Where m is the size of Sample, $X_i = \{X_i^1, \dots, X_i^n\}$, $n = \text{Size of Features or Input Nodes}$
 (rows)

Assume¹⁾ the error of X comes from a Gaussian white noise of each $X_i^j \sim N(\mu(X^j), \sigma^2(X^j))$

Then $y_i = f(X_i; \tilde{w}) + \epsilon(X_i) = f(X_i; \tilde{w}) + \sum_{j=1}^n \epsilon(X_i^j)$, $f(X_i; \tilde{w})$ is net-work procedure
 function to map X_i to y_i (X_i is Set of Single row data, y_i is Set of Output Node Value)
 Here, we assume²⁾ the function $f(X_i; \tilde{w})$ is a perfect deterministic mapping function without uncertainty.

Then the uncertainty is purely $\sum_{j=1}^n \epsilon(X_i^j)$ for a set of unknown Sample (single row)

If we are not to assume 1) condition, we need to add the uncertainty inside $f(X_i; \tilde{w})$,
 which is $D_f^2 = E\{[f(X_i; \tilde{w}) - E(f(X_i; \tilde{w}))]^2\}$ if f can be a known distribution (Assume³⁾)
 or

$\text{Var}(f(X_i; \tilde{w}))$ can be approximated using Taylor expansion

$$\text{Var}(f(X_i; \tilde{w})) \approx \left(f'(E(X_i; \tilde{w})) \right)^2 \text{Var}(X_i) + T_3$$

Second Moment other moments

Which I will not elaborate here