

Notes on Deep Feedforward Network (MLP)

Dr. Hao Wang

Late May-Early June, 2016

The task is to construct a proper computing method to obtain the gradient needed for very much a graph that describes the interactions of the perceptrons, also widely known as multilayer perceptron. We visualise this type of deep feedforward network as shown in Figure 1. Given m training sample data, a final cost/loss function without regularization term can be expressed as

$$J_A = \frac{1}{m} \sum_{i=1}^m J(\mathbf{x}, y; \mathbf{W}, \mathbf{b}).$$

A gradient descendant method has one key step in that parameters like W is updated each iteration by

$$W^{\text{new}} = W^{\text{old}} - \alpha \frac{\partial J_A}{\partial W},$$

where α is usually termed as *learning rate*. Here the cost function takes a 2-norm squared form i.e.

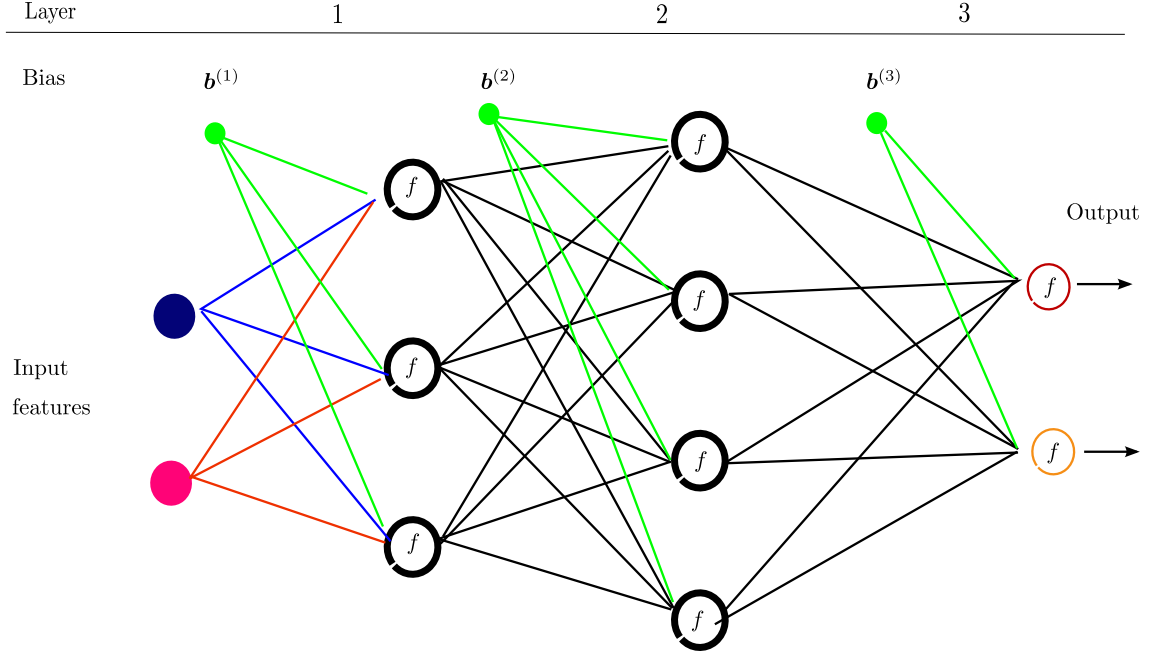
$$J = \frac{1}{2} \|\hat{y}(\mathbf{x}, \mathbf{W}, \mathbf{b}) - y\|^2.$$

1 The Flow of Back Propagation

Before we compute the gradient using BP method, we need do a feedforward step by some randomly initialised weights and biases so that we have information about $\mathbf{a}^{(i)}$ and $\mathbf{h}^{(i)}$. The full flow of the backward propagation scheme is shown in Figure 2. Here, theoretically the chain rule can be applied without much pain provided we stick to the *numerator layout* convention for vector/matrix derivatives. (For general purpose, we encounter a 3rd order tensor in the derivatives, which can be mitigated using index computing. We shall see this in the next section.) Also here the case is for a single training sample. Therefore,

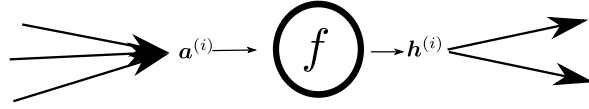
$$\begin{aligned} \frac{\partial J}{\partial \mathbf{W}^{(n_L)}} &= \frac{\partial J}{\partial \mathbf{h}^{(n_L)}} \frac{\partial \mathbf{h}^{(n_L)}}{\partial \mathbf{a}^{(n_L)}} \frac{\partial \mathbf{a}^{(n_L)}}{\partial \mathbf{W}^{(n_L)}} \\ \frac{\partial J}{\partial \mathbf{b}^{(n_L)}} &= \frac{\partial J}{\partial \mathbf{h}^{(n_L)}} \frac{\partial \mathbf{h}^{(n_L)}}{\partial \mathbf{a}^{(n_L)}} \frac{\partial \mathbf{a}^{(n_L)}}{\partial \mathbf{b}^{(n_L)}} \end{aligned} \quad (1)$$

and similar results can be obtained for other terms.



From i^{th} layer, the weight from node $j \rightarrow k$: $W_{jk}^{(i)}$

$$\mathbf{a}^{(i)} = \mathbf{b}^{(i)} + \mathbf{W}^T \mathbf{h}^{(i-1)}$$



$$\mathbf{h}^{(i)} = f(\mathbf{a}^{(i)})$$

Figure 1 Illustration of deep feedforward networks

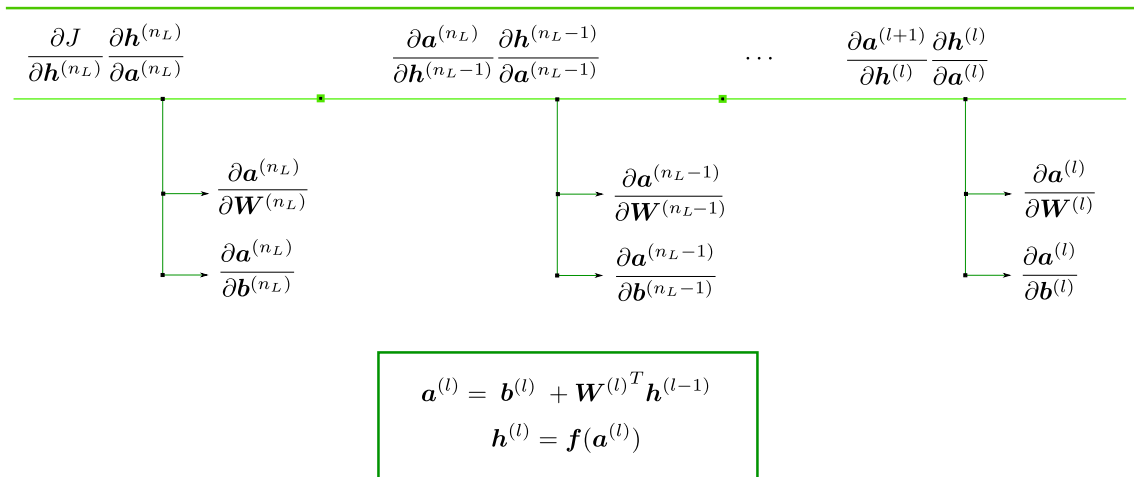


Figure 2 Illustration on how to compute the relevant gradients backwards, the squared nodes indicate multiplication action.

2 The Proposed BP Algorithm

First we compute a few terms occur in the above theoretic framework in Figure 2. Assume all the input and output are column vectors, and we adopt numerator layout convention for derivative computing and the indices in parentheses are indicator of the layer number. For i^{th} sample at the output layer, we have

$$\frac{\partial J}{\partial \mathbf{h}^{(n_L)}} = \left(\mathbf{h}^{(n_L)} - y_i \right)^T, \quad (2)$$

which by the way is a *row* vector.

Besides we know that the activation function \mathbf{f} is acting element-wise, so

$$\frac{\partial \mathbf{h}^{(l)}}{\partial \mathbf{a}^{(l)}} = \text{Diag}(\mathbf{f}'(\mathbf{a}^{(l)})) \quad (3)$$

Let us denote \mathbf{g} as the product of (2) and (3) (let $l = n_L$), then the computing can be simplified as

$$\mathbf{g} = \frac{\partial J}{\partial \mathbf{h}^{(n_L)}} \frac{\partial \mathbf{h}^{(n_L)}}{\partial \mathbf{a}^{(n_L)}} = \left(\mathbf{h}^{(n_L)} - y_i \right)^T \odot \mathbf{f}'(\mathbf{a}^{(n_L)})^T, \quad (4)$$

where operator \odot represents element-wise multiplication, so \mathbf{g} is also a row vector.

Using numerator convention, we can easily verify that

$$\frac{\partial \mathbf{a}^{(l+1)}}{\partial \mathbf{h}^{(l)}} = W^{(l+1)T} \quad (5)$$

and

$$\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{b}^{(l)}} = \mathbf{I}, \quad (6)$$

where \mathbf{I} is the identity matrix. The slightly tricky bit is $\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{W}^{(l)}}$ since according to definition we end up with a 3^{rd} order tensor. For notational simplicity, we omit the upper script (l) and work on a general gradient term $\frac{\partial J}{\partial \mathbf{W}}$. With a proper row vector \mathbf{g} , we can have

$$\frac{\partial J}{\partial \mathbf{W}} = \mathbf{g} \frac{\partial \mathbf{a}}{\partial \mathbf{W}}.$$

In the following, we use the summation convention, where the identical indices imply summation. So

$$\frac{\partial J}{\partial W_{jk}} = g_i \frac{\partial W_{qi} h_q}{\partial W_{jk}} = g_i h_q \delta_{qj} \delta_{ik} = g_k h_j, \quad (7)$$

transforming into a matrix form (also in numerator convention) we have

$$\frac{\partial J}{\partial \mathbf{W}} = (\mathbf{h} \mathbf{g})^T.$$

(\mathbf{h} is a column vector, and for actual updating W , no need to transpose as shown in the algorithm.)

The other identity:

$$\frac{\partial J}{\partial \mathbf{b}} = \mathbf{g} \mathbf{I} = \mathbf{g}.$$

Putting in column vector form, we can transpose \mathbf{g} . The algorithm can be constructed now as in Algorithm 1. (Note that all algorithm index starts with zero.)

Algorithm 1 Back propagation algorithm to compute gradients (single training sample)

Require: One training sample $(\mathbf{h}_i^{(0)}, y^i)$; Initialised \mathbf{W}, \mathbf{b} ; Activation function \mathbf{f}

Ensure: Gradients $\partial J / \partial \mathbf{W}$ and $\partial J / \partial \mathbf{b}$

```

1: procedure BP PROCEDURE
2:   First run a feedforward pass, to get  $\mathbf{a}$  and  $\mathbf{h}$ 
3:    $\mathbf{g} := \left( \mathbf{h}^{(n_L-1)} - y^i \right)^T \odot \mathbf{f}'(\mathbf{a}^{(n_L-1)})^T$   $\triangleright$  Starts from the output layer
4:    $\partial J / \partial \mathbf{W}^{(n_L)} := \mathbf{h}^{(n_L-1)} \mathbf{g}$ 
5:    $\partial J / \partial \mathbf{b}^{(n_L)} := \mathbf{g}^T$ 
6:   for  $l = n_L - 2, n_L - 2, \dots, 0$  do
7:      $\mathbf{g} := \left( \mathbf{g} \left[ \mathbf{W}^{(l+1)} \right]^T \right) \odot \mathbf{f}'(\mathbf{a}^{(l)})^T$ 
8:      $\partial J / \partial \mathbf{W}^{(l)} := \mathbf{h}^{(l-1)} \mathbf{g}$ 
9:      $\partial J / \partial \mathbf{b}^{(l)} := \mathbf{g}^T$ 
10:  end for
11: end procedure

```

3 Gradient Descendant Algorithm

This section is try to incorporate the former BP Procedure to obtain the full version gradient descendant (optimization) algorithm. Everything derived so far is for **single training data** only. We might add a regularization term for the weights $W_{ij}^{(l)}$ to the loss/cost function J_A , so it becomes

$$J_A = \frac{1}{m} \sum_{i=1}^m J_i(\mathbf{x}, y; \mathbf{W}, \mathbf{b}) + \frac{\lambda}{2} \sum_{i,j,l} (W_{ij}^{(l)})^2. \quad (8)$$

The overall gradients $\partial J_A / \partial \mathbf{W}$ is known by (the form breaks the numerator convention, no transpose, to simplify computation)

$$\frac{\partial J_A}{\partial \mathbf{W}} = \frac{1}{m} \sum_i^m \nabla_{\mathbf{W}} J_i + \lambda \mathbf{W}. \quad (9)$$

The gradient descendant algorithm is shown in Algorithm 2.

4 Implementation Tricks: Not Reinventing the Wheels

The scikit-learn machine learning library suggests to call the standard scipy optimization library **lbfgs** (Limited memory BFGS) algorithm, which ought to perform better than purely gradient descendant, since Hessian information is employed. All left to do is to transform the deep learning (network) optimization problem to the conventional one.

In BFGS, we need two piece of information:

1. the underlying objective function $f(x)$,
2. the gradient to the objective function $\nabla f(x)$,

which are cost function J_A and derivatives (gradient) $\frac{\partial J_A}{\partial \mathbf{W}}$ and $\frac{\partial J_A}{\partial \mathbf{b}}$ under our notation. Next we deal with 1-feature(classifier) output problem, and the row vector \mathbf{g} To account for all m training

Algorithm 2 Full gradients descendant algorithm

Require: All training samples (\mathbf{x}, \mathbf{y}) ; Initialised \mathbf{W}, \mathbf{b} ; Activation function \mathbf{f} ; Learning rate α , regularization parameter λ

Ensure: Trained weights and biases: \mathbf{W}, \mathbf{b}

```
1: procedure GRADIENT DESCENDANT PROCEDURE
2:   while Cost  $J_A$  is not small enough do
3:     Set  $\Delta \mathbf{b}^{(l)} := \mathbf{0}$ ,  $\Delta \mathbf{W}^{(l)} := \mathbf{0}$  for all layers  $l$ 
4:     for  $l = 0, 2, \dots, m-1$  do
5:       Call BP PROCEDURE to get  $\partial J / \partial \mathbf{W}^{(l)}$  and  $\partial J / \partial \mathbf{b}^{(l)}$ 
6:        $\Delta \mathbf{b}^{(l)} := \Delta \mathbf{b}^{(l)} + \partial J / \partial \mathbf{b}^{(l)}$ 
7:        $\Delta \mathbf{W}^{(l)} := \Delta \mathbf{W}^{(l)} + \partial J / \partial \mathbf{W}^{(l)}$ 
8:     end for
9:      $\mathbf{b}^{(l)} := \mathbf{b}^{(l)} - \frac{\alpha}{m} \Delta \mathbf{b}^{(l)}$ 
10:     $\mathbf{W}^{(l)} := \mathbf{W}^{(l)} - \alpha \left[ \frac{1}{m} \Delta \mathbf{W}^{(l)} + \lambda \mathbf{W}^{(l)} \right]$ 
11:    Compute the current cost  $J_A(\mathbf{W}^{(i)}, \mathbf{b}^{(l)})$ 
12:  end while
13: end procedure
```

data, we can arrange the input/output at each node as a $p \times m$ and $q \times m$ matrix.

$$\begin{aligned} \mathbf{a}^{(l)}_{q \times m} &= \mathbf{b}^{(l)}_{q \times 1} \mathbf{I}_{1 \times m} + \mathbf{W}^{(i)T}_{q \times p} \mathbf{h}^{(i-1)}_{p \times m} \\ \mathbf{h}^{(l)}_{q \times m} &= \mathbf{f}(\mathbf{a}^{(l)}_{q \times m}) \end{aligned} \quad (10)$$

where $\mathbf{I}_{1 \times m}$ is all-ones row vector. Cost function J_A is calculated by (8). After some careful manipulations we can have the following algorithm to compute $\frac{\partial J_A}{\partial \mathbf{W}}$ and $\frac{\partial J_A}{\partial \mathbf{b}}$:

Algorithm 3 Algorithm to compute gradients (all training sample)

Require: Training sample (\mathbf{x}, \mathbf{y}) ; Initialised \mathbf{W}, \mathbf{b} ; Activation function \mathbf{f}

Ensure: Gradients $\partial J_A / \partial \mathbf{W}$ and $\partial J_A / \partial \mathbf{b}$

```
1: procedure REVISED BP PROCEDURE
2:   First run a feedforward pass based on (10)
3:   for  $l = n_L - 1, n_L - 1, \dots, 0$  do
4:     if  $l == n_L - 1$  then
5:        $\mathbf{g} := (\mathbf{h}^{(l)} - \mathbf{y})^T \odot \mathbf{f}'(\mathbf{a}^{(l)})^T$   $\triangleright$  Starts from the output layer,  $\mathbf{g}: m \times q$ 
6:     else
7:        $\mathbf{g} := \left( \mathbf{g} \left[ \mathbf{W}^{(l+1)} \right]^T \right) \odot \mathbf{f}'(\mathbf{a}^{(l)})^T$   $\triangleright$  for other layers,  $\mathbf{g}: q_l \times m$ 
8:     end if
9:      $\partial J_A / \partial \mathbf{W}^{(l)} := \frac{1}{m} \mathbf{h}^{(l-1)} \mathbf{g} + \lambda \mathbf{W}^{(l)}$ 
10:     $\partial J_A / \partial \mathbf{b}^{(l)} := \frac{1}{m} (\mathbf{I} \mathbf{g})^T$   $\triangleright \mathbf{I}: 1 \times m$ 
11:  end for
12: end procedure
```

5 Experimental Run

The training dataset and test set are shown in Figure 3. The classified results is shown in Figure 4–Figure 8.

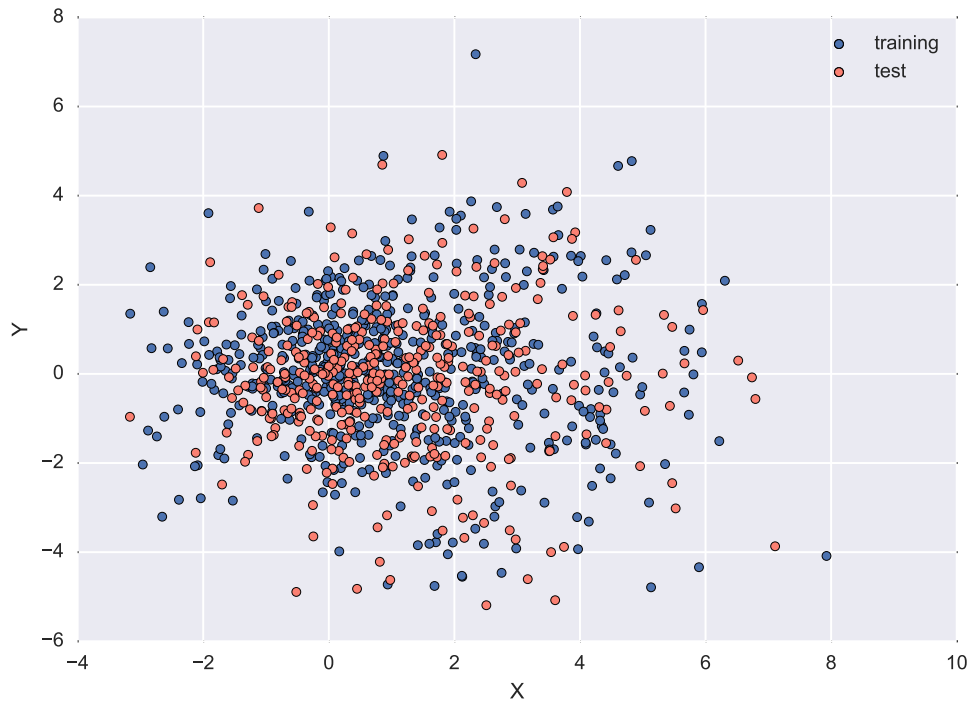


Figure 3 A Gaussian distributed data sample, divided by training and test set.

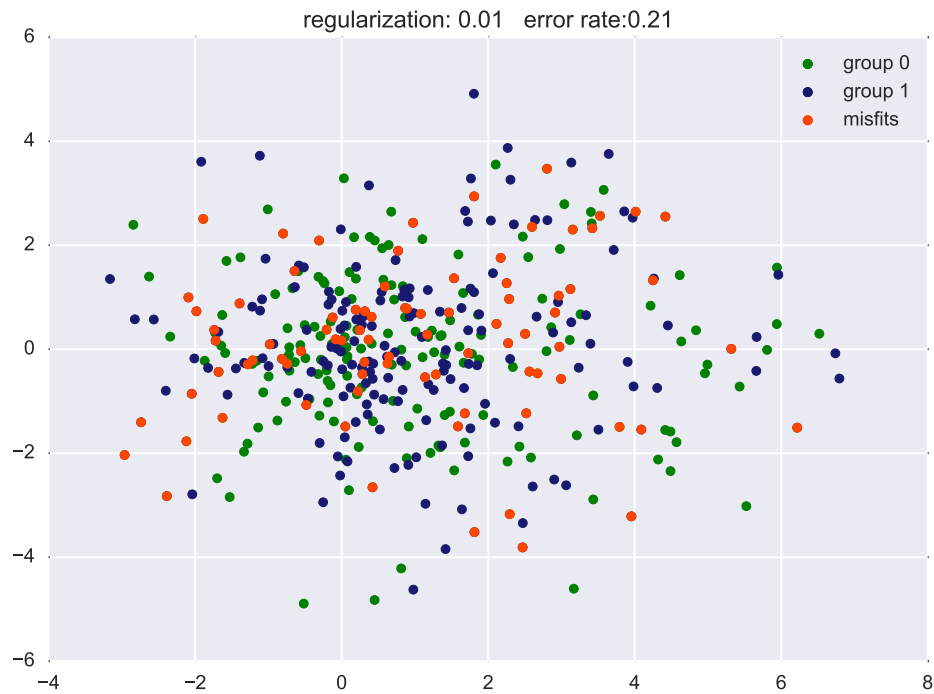


Figure 4 A trained network with two hidden layers (100, 4), regularization parameter $\lambda = 0.01$

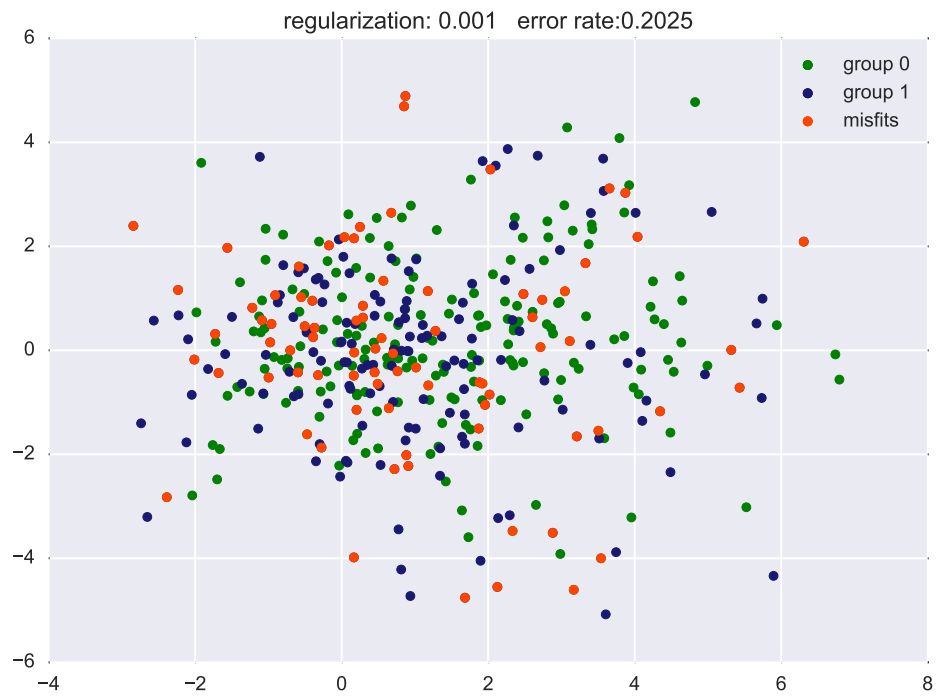


Figure 5 A trained network with two hidden layers (100, 4), regularization parameter $\lambda = 0.001$

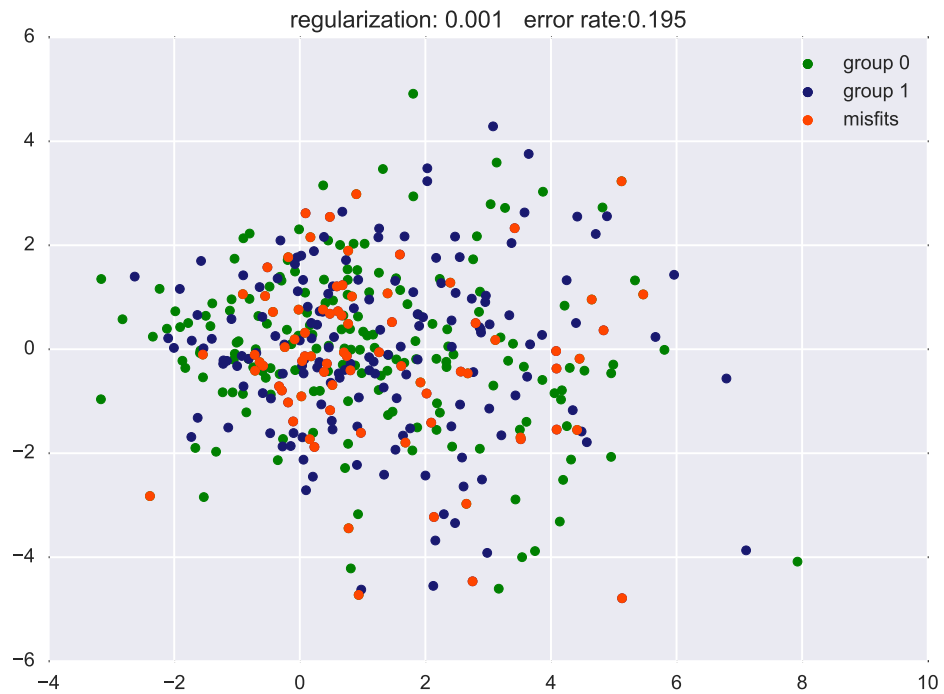


Figure 6 A trained network with two hidden layers (10, 4), regularization parameter $\lambda = 0.001$

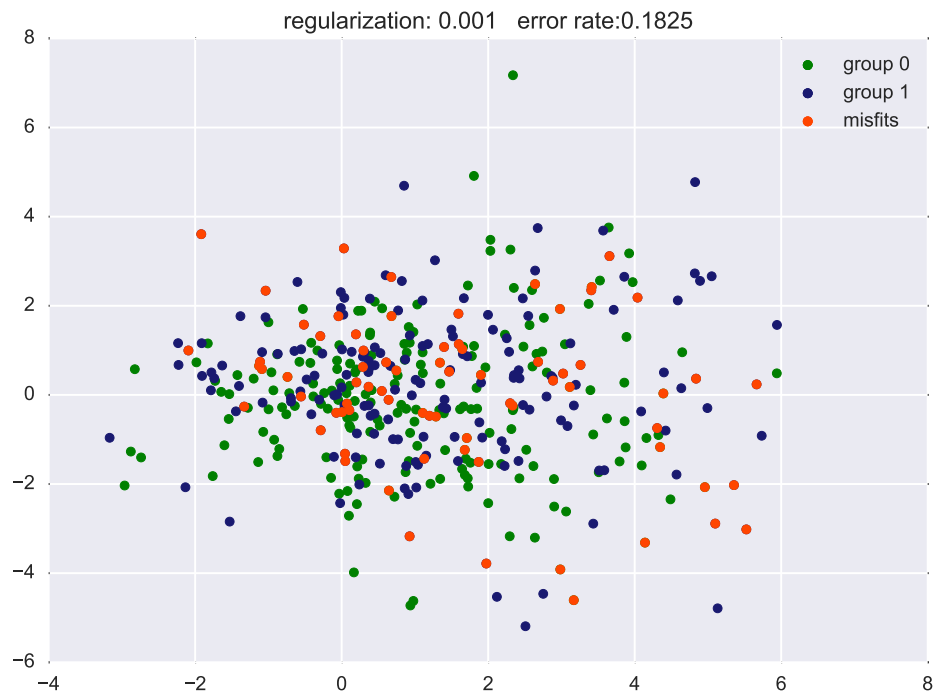


Figure 7 A trained network with two hidden layers (10,400), regularization parameter $\lambda = 0.001$

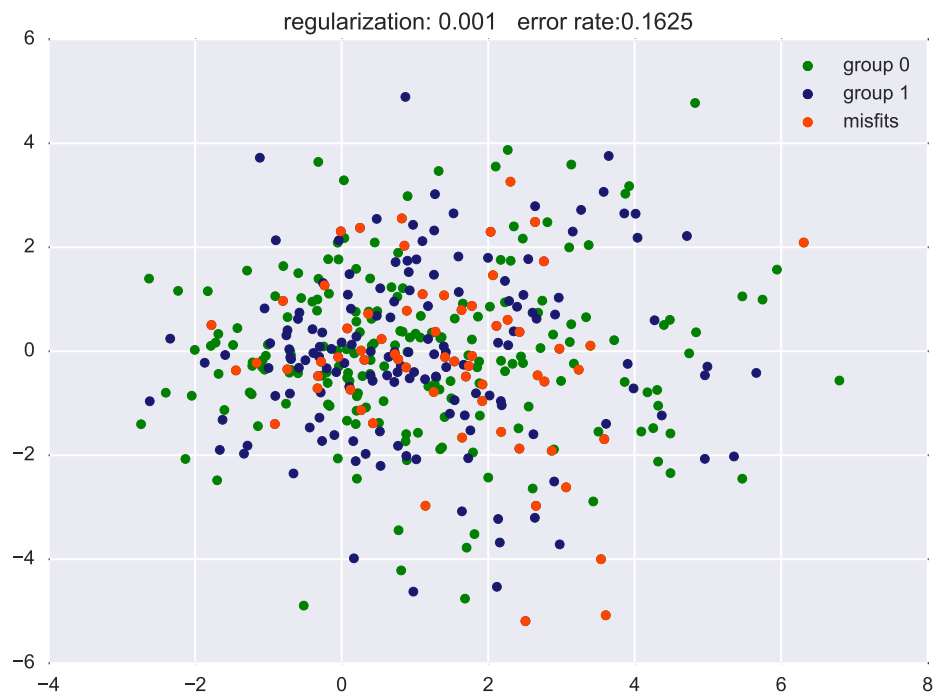


Figure 8 A trained network with three hidden layers (10,400,2), regularization parameter $\lambda = 0.001$