

Project-3 Neural Network

Tianran Zhang 15300180104

November 10, 2019

Contents

1	Neural Network	1
1.1	Change the network structure	2
1.2	Change the step-size	2
1.3	Vectorize evaluating the function	4
1.4	Add l2 regrlarization	5
1.5	Softmax layer	7
1.6	Add bias in each layer	7
1.7	Implement 'dropout'	9
1.8	'fine-tuning'	10
1.9	Create more training examples	11
1.10	2D convolutional layer	12
2	Final script	13
3	Matlab Code	13
3.1	LossSoftmax	13
3.2	PredictSoftmax	15
3.3	ConvPredict	16
3.4	final_MLP	17
3.5	Loss_final	19
3.6	Predict_final	21

1 Neural Network

In this problem we aimed to investigate handwritten digit classification. The inputs are handwritten digits, and our goal is to predict the number of the given image. Now we choose Neural Network to train the model. The basic training procedure have already given, so my work is to modify the given training procedure and to optimize the performance.

To achieve the best test error, I do some modifications to the procedure step by step and integrate them together in the end. In the following section, I will report the modifications I made together with the best test error I have achieved.

1.1 Change the network structure

In a neural network procedure, we should first set the number of layers and the hidden units in each layer. To work out the best solution, I looked up much information and find that most classification using neural network have one layer and the experienced best hidden units is $(\text{inputs number} + \text{output number}) * 2/3$, which is equal to 178. To get the right number of units, I wrote a script named `find_nHidden.m`(which is shown below in section 2.1) and draw a plot of validation error - units. The plots are below:

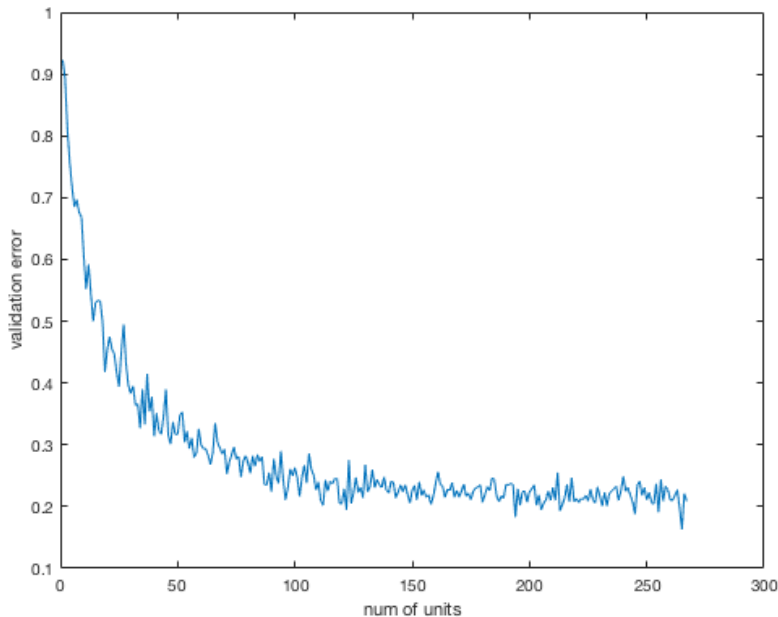


Figure 1: The plot of validation error - units

From the plot we could see that the validation error gets lowest when the number of units is within 100 to 250, which also contains the experienced number 178(the experienced best hidden units I have calculated before). So I choose 100 units which may save some time of the procedure.

After setting the little change of the network structure as:

```
1 nHidden = [100];
```

I trained the model for several times and got the Test error is approximately 0.24, which is half lower then the initial Test error.

1.2 Change the step-size

In the network training procedure, the way we modify w with different step-size is very important. Choosing small step-size, we may fall into local minimum and waste much time, On

the other side, choosing big step-size may be difficult for our model to get to the minimum of training error. Either big or small step-size will lead to a bad result. Now we try to deal with the problem by two ways:

a. Modify the sequence of step-size.

As the sequence of step-size is : $w^{t+1} = w^t - \alpha_t \nabla f(w^t)$, where α_t is the learning rate (step size). we add an item : $\beta_t(w^t - w^{t-1})$ (where $\beta_t = 0.9$) to reduce the influence of each training on weights. So we replace the sequence of step-size by : $w^{t+1} = w^t - \alpha_t \nabla f(w^t) + \beta_t(w^t - w^{t-1})$. Now we could avoid falling into local minimum, and thus we could make the initial step-size a little bigger, and to improve the convergence.

The modified part :

```

1      i = ceil(rand*n);
2      [errors(iter), g] = funObj(w, i);
3
4      w = w - stepSize*g + 0.9 * (w - w0);
5      w0 = w;
```

b. Use different step-size during the training.

One way to change the step-size is to set a minStep-size and a maxStep-size. In each iteration, the step-size is : $\text{step-size} = \text{maxStep-size} - \frac{\text{iter}}{\text{maxIter}} (\text{maxStep-size} - \text{minStep-size})$

The modified part :

```

1      %initialize
2      stepSize = 1e-3 ;
3      maxStepsize = 1e-3 * 5;
4      minStepsize = maxStepsize/10;
```

For each item in the loop:

```

1      i = ceil(rand * n);
2      [errors(iter), g] = funObj(w, i);
3
4      stepSize = maxStepsize - (maxStepsize - minStepsize)* iter / maxIter;
5      w = w - stepSize*g;
```

I run the model for several times and get the result that:

If we do nothing to the step-size, the validation error convergence to 50%;

If we do a (modify the sequence of step-size), the validation error convergence to 32%;

If we do b (using different step-size), the validation error convergence to 30%.

So apparently, our modification does improve the convergence.

1.3 Vectorize evaluating the function

Notice that the training procedure is quite slowly, I managed to vectorize the loss function together with the predict function by:

a. Try to initialize the cell, vector, and matrix, and compute the values before we use them for a lot of times.

In the predict function:

```
1 nH = length(nHidden);
2 ip = cell(1, nH);
3 fp = cell(1, nH);
```

And we do the same to the loss function as well.

b. Try to express as much as possible in terms of matrix operations.

In the loss function:

When calculating the output gradients, instead of using the loop:

```
1 for c = 1:nLabels
2     gOutput(:,c) = gOutput(:,c) + err(c)*fp{end}';
3 end
```

We can save a bunch of time by using:

```
1 gOutput = gOutput + fp{end}' * err;
```

When calculating the input gradients, instead of using the loop:

```
1 for c = 1:nLabels
2     gInput = gInput + err(c)*X(i,:)'.*(sech(ip{end}
3     outputWeights(:,c)'));
4 end
```

we could use:

```
1 gInput = gInput + X(i,:)'.*(sech(ip{end}.* ...
2     (outputWeights * err')'));;
```

And if we have more than one hidden layers, instead of using the loop:

```
1 for c = 1:nLabels
2     backprop(c,:) = err(c)*(sech(ip{end}outputWeights(:,c)'));
3     gHidden{end} = gHidden{end} + fp{end-1}'*backprop(c,:);
4 end
5 backprop = sum(backprop, 1);
```

We could vectorize it by :

```
1 backprop = err' * sech(ip{end}.* outputWeights');
2 backprop = sum(backprop,1);
3 gHidden{end} = gHidden{end} + fp{end-1}' * backprop;
```

After the speed up modification, I run the code and find out that the time to finish the code drops from 21s to 7s ! Which proves that the modification is effectiveness.

1.4 Add l2 regrlarization

One important modification to model training is to add a l2 regularization of the weights to the loss function, which is called *weight decay*. After adding the l2 regularization, the gradients change,

$$E_{l2} = E + \frac{1}{2}|weights|^2$$

$$gradient_{l2} = gradient + weights$$

so that we make a little change to the loss function:

$$gOutput_{l2} = gOutput + outputweights$$

$$gHidden_{l2} = gHidden + Hiddenweights$$

$$gInput_{l2} = gInput + inputweights$$

Output gradients:

```

1 gOutput = gOutput + fp{end}' * err + lambda * outputWeights .* ...
2   [zeros(1, size(outputWeights, 2)); ...
3   ones(size(outputWeights, 1)-1, ...
4   size(outputWeights, 2))];

```

Input gradients:

```

1 gInput = gInput + X(i,:) * (sech(ip{end} .* ...
2   (outputWeights * err')') + lambda * inputWeights .* ...
3   [zeros(1, size(inputWeights, 2)); ...
4   ones(size(inputWeights, 1)-1, ...
5   size(inputWeights, 2))];

```

Hidden gradients (if there are more than one hidden layers):

```

1 backprop = err' * sech(ip{end} .* outputWeights');
2   backprop = sum(backprop, 1);
3   gHidden{end} = gHidden{end} + fp{end-1}' * backprop + ...
4   lambda * hiddenWeights{end} .* ...
5   [zeros(1, size(hiddenWeights{end}, 2)); ...
6   ones(size(hiddenWeights{end}, 1)-1, ...
7   size(hiddenWeights{end}, 2))];
8
9   % Other Hidden Layers
10  for h = nH-2:-1:1
11      backprop = (backprop * hiddenWeights{h+1}') .* ...
12      sech(ip{h}
13      gHidden{h} = gHidden{h} + fp{h}' * backprop + ...
14      lambda * hiddenWeights{h} .* ...
15      [zeros(1, size(hiddenWeights{h}, 2)); ...

```

```

16         ones(size(hiddenWeights{h}, 1)-1, ...
17             size(hiddenWeights{h}, 2)));
18
19     end

```

Now, is quite important to find out the best λ , I wrote a loop that make λ ranged from 0 to 1 and find that it works better in range 0.01 to 0.06, so I adjust the loop make λ ranged from 0.01 to 0.05. The plot is shown below:

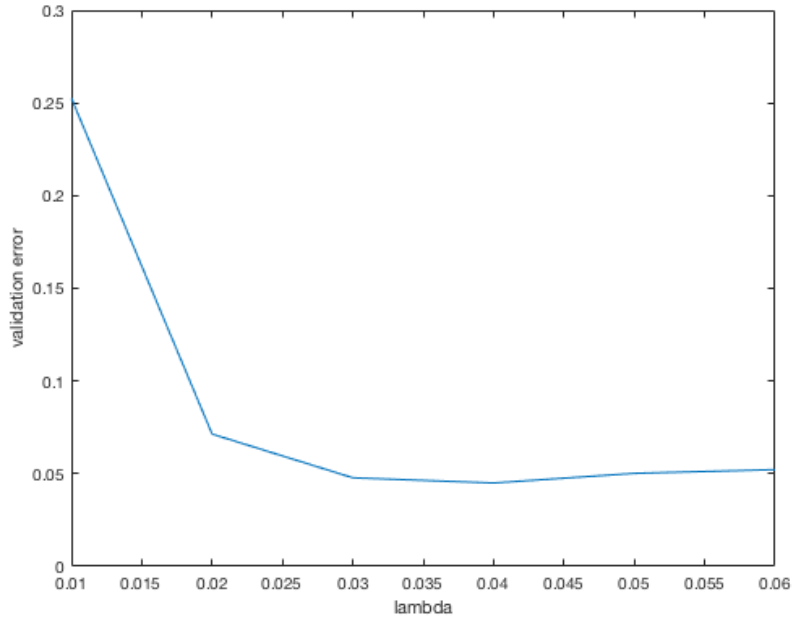


Figure 2: The plot of validation error - lambda

From the plot we could see that λ works best when ranged in 0.03 to 0.05, so I choose $\lambda = 0.03$ and run the script for several times and got the Test error with final model = 0.041000.

Except for adding a l2 regularization, I tried early stopping in the main script, which is used to stop training when the error on the validation set stops decreasing consistently for 3 times. That is : if $validation_error_t > validation_error_{t-1} > validation_error_{t-2}$ then stop training. The modification in the loop is :

```

1     a2 = sum(yhatyvalid)/t;
2     if (a2 > a1 )&(a1 > a0)
3         break;
4     else
5         a0 = a1;
6         a1 = a2;
7     end

```

After running the script, I found out that the test error and time caused are both reduced largely, so I will add this method to the final script.

1.5 Softmax layer

Instead of using the squared error, I use a softmax layer at the end of the network. So the 10 outputs can be interpreted as the probabilities of each class. To finish the softmax training procedure, I make some modifications to the loss function:

a. In the last layer, $z_i = fp\{end\}outputWeights$, where $i = 1 \dots nLabels$, we put them in the softmax function : $p(y_i) = \frac{\exp(z_i)}{\sum_{j=1}^J \exp(z_j)}$, so we calculate the probabilities in each class, and the error is $E = -\log p(y_{trueLabel})$:

```

1      yhat1 = exp(fp{end}*outputWeights);
2      yhat = yhat1/sum(yhat1);
3      y_true] = max(y(i, :));
4      yhat_true = yhat(y_true);
5      err = -log(yhat_true);
6      f=f + err;

```

b. Since we changed the output numbers, the gradients changed, too. So we made the calculation by :

$$\begin{aligned}
 \frac{\partial E}{\partial outputWeights} &= \frac{\partial E}{\partial P(y_{trueLabel})} \frac{\partial P(y_{trueLabel})}{\partial z} \frac{\partial z}{\partial outputWeights} \\
 &= -\frac{1}{p(y_{trueLabel})} \frac{\partial P(y_{trueLabel})}{\partial z} fp\{end\} \\
 \\
 \frac{\partial E}{\partial inputWeights} &= \frac{\partial E}{\partial P(y_{trueLabel})} \frac{\partial P(y_{trueLabel})}{\partial z} \frac{\partial z}{\partial fp\{1\}} \frac{\partial fp\{1\}}{\partial inputWeights} \\
 &= -\frac{1}{p(y_{trueLabel})} \frac{\partial P(y_{trueLabel})}{\partial z} outputWeights(:)X(i,:)
 \end{aligned}$$

I applied this on matlab:

```

1  gOutput = gOutput - fp{end}' * (1 - yhat_true) * (y(i,:)==1);
2
3  gInput = gInput - (1 - yhat_true) * X(i,:) ' * ...
4  ( sech(ip{end}.* outputWeights(:, y(i,:)==1)) ' );

```

Then I run the script (which is shown below in Sec.3) for several times and got the result that the validation error is convergence to 18%, so I will add this method to the final script.

1.6 Add bias in each layer

Instead of just having a bias variable at the beginning, I made the last units of each layer a constant, so that each layer has a bias.

Since I made the last units of each layer a constant, the weights for each layers should be redistributed as:

```

1 inputWeights = reshape(w(1:nVars*(nHidden(1)-1)), nVars, nHidden(1)-1);
2 offset = nVars * (nHidden(1)-1);
3 for h = 2:nH
4     hiddenWeights{h-1} = reshape(...
5         w(offset+1 : offset+nHidden(h-1)*(nHidden(h)-1)),...
6         nHidden(h-1) , nHidden(h)-1);
7     offset = offset + nHidden(h-1)*(nHidden(h)-1);
8 end
9 outputWeights = w(offset+1:offset+nHidden(end)*nLabels);
10 outputWeights = reshape(outputWeights, nHidden(end), nLabels);

```

And since each layers have one units remains const, we should calculate the gradients and backprop without the consts' influence, that is:

```

1 % Output Weights
2 gOutput = gOutput + fp{end}' * err;
3
4 if nH > 1
5     % Last Layer of Hidden Weights
6     backprop = err' * sech(ip{end}.* outputWeights');
7     backprop = sum(backprop,1);
8     g1 = fp{end-1}' * backprop;
9     gHidden{end} = gHidden{end} + g1(:, 1: nHidden(end)-1);
10
11 % Other Hidden Layers
12 for h = nH-2 : -1 : 1
13     g1 = (backprop * hiddenWeights{h+1}') .* ...
14     sech(ip{h
15     backprop = g1(:, 1: nHidden(h)-1);
16     gHidden{h} = gHidden{h} + fp{h}'*backprop;
17 end
18
19 % Input Weights
20 g1 = (backprop(:, 1:nHidden(2)-1)*hiddenWeights{1}') .* ...
21     sech(ip
22     backprop = g1(:, 1:nHidden(1)-1);
23     gInput = gInput + X(i,:)'*backprop;
24 else
25     % Input Weights
26     gx = X(i,:) ' * (sech(ip{end}.* ...
27         (outputWeights * err')' );
28     gInput = gInput + gx(:, 1: nHidden(1)-1);
29 end

```

When we do the prediction of valid data or test data, we shall also make a little change:

```

1 for i = 1:nInstances

```



```

2      ip{1} = [X(i,:) * inputWeights, 1];
3      fp{1} = tanh(ip{1});
4      for h = 2:nH
5          ip{h} = [fp{h-1} * hiddenWeights{h-1}, 1];
6          fp{h} = tanh(ip{h});
7      end
8      y(i,:) = fp{end} * outputWeights;
9  end
10  y] = max(y, [], 2);

```

1.7 Implement 'dropout'

In this section, we randomly dropped out some hidden units with probability p during the training. This is the dropout method. When we choose $p = 0.5$ as a common choice, we randomly dropped half units, which is a good way to prevent over-fitting.

We used to train the neural network with:

$$ip(l+1) = w^{l+1}y^l$$

$$fp(l+1) = \tanh(ip(l+1))$$

But after we adapt dropout method, the calculation formula comes to :

$$r_j^l \sim \text{Bernoulli}(p)$$

$$\tilde{y}^l = r^l * y^l$$

$$ip(l+1) = w^{l+1}\tilde{y}^l$$

$$fp(l+1) = \tanh(ip(l+1))$$

I applied it on the loss function:

```

1      p{1} = [rand(1, nHidden(1)) > 0.5];
2
3      ip{1} = X(i,:) * inputWeights .* p{1};
4      fp{1} = tanh(ip{1});
5      for h = 2:nH
6          p{h} = [rand(1, nHidden(h)) > 0.5];
7          ip{h} = fp{h-1} * hiddenWeights{h-1} .* p{h};
8          fp{h} = tanh(ip{h});
9      end

```

When we calculate the gradients, we should still drop out the chosen units:

```

1  gOutput = (gOutput + fp{end}' * err) .* repmat(p{end}', 1, nLabels);
2
3  if nH > 1
4      % Last Layer of Hidden Weights
5      backprop = err' * sech(ip{end} .* outputWeights');

```

```

6      backprop = sum(backprop,1);
7      gHidden{end} = (gHidden{end} + fp{end-1}' * backprop) .* ...
8      repmat(p{end-1}', 1, length(p{end})) .* ...
9      repmat(p{end}, length(p{end-1}), 1);
10
11      % Other Hidden Layers
12      for h = nH-2:-1:1
13          backprop = (backprop * hiddenWeights{h+1}') .* ...
14                      sech(ip{h
15          gHidden{h} = (gHidden{h} + fp{h}'*backprop) .* ...
16                      (p{h}'*p{h+1}));
17      end
18
19      % Input Weights
20      backprop = (backprop*hiddenWeights{1}') .* sech(ip
21      gInput = (gInput + X(i,:) '*backprop) .* repmat(p{1}, nVars, 1);
22  else
23      % Input Weights
24      gInput = gInput + X(i,:) ' * (sech(ip{end} .* ...
25                      (outputWeights * err ') ');
26
27      gInput = gInput .* repmat(p{1}, nVars, 1);

```

Now I run the main script and find out that the validation convergence to 22% with the convergence speed improved largely. So I will add this method into the final script.

1.8 'fine-tuning'

We do 'fine-tuning' to the last layer : Fix the parameters of all the layers except the last one, and solve for the outputWeights exactly as a convex optimization problem :

$$\begin{aligned}
 fp\{end\}outputWeights &= y_i \\
 fp\{end\}'fp\{end\}outputWeights &= fp\{end\}'y_i \\
 outputWeights &= (fp\{end\}'fp\{end\})^{-1}fp\{end\}'y_i
 \end{aligned}$$

I applied it on Matlab:

```

1 outputWeights = pinv(fp{end}' * fp{end}) * fp{end}' * y(i,:);

```

After working out the exact outputWeights, we use it to form the new weights:

```

1 w(offset+1:offset+nHidden(end)*nLabels) = outputWeights(:);

```

Through the new weights to calculate the gradients of weights as we have done before.

I run the script for several times, the results doesn't seem well:

```

1 >> MLP8
2 Training iteration = 0, validation error = 0.901200
3 Training iteration = 5000, validation error = 0.856600

```

```

4 Training iteration = 10000, validation error = 0.851400
5 Training iteration = 15000, validation error = 0.912400
6 Training iteration = 20000, validation error = 0.919000
7 Training iteration = 25000, validation error = 0.865400
8 Training iteration = 30000, validation error = 0.864000
9 Training iteration = 35000, validation error = 0.892600
10 Training iteration = 40000, validation error = 0.907600
11 Training iteration = 45000, validation error = 0.878200
12 Training iteration = 50000, validation error = 0.863600
13 Training iteration = 55000, validation error = 0.890200
14 Training iteration = 60000, validation error = 0.878200
15 Training iteration = 65000, validation error = 0.905600
16 Training iteration = 70000, validation error = 0.849800
17 Training iteration = 75000, validation error = 0.853800
18 Training iteration = 80000, validation error = 0.879800
19 Training iteration = 85000, validation error = 0.866000
20 Training iteration = 90000, validation error = 0.894800
21 Training iteration = 95000, validation error = 0.882200
22 Elapsed time is 16.472326 seconds.
23 Test error with final model = 0.904000

```

It seems that the validation error doesn't change much. This is because when the `outputWeights` being exactly the matrix which $fp\{end\}outputWeights = y_i$, then the error : $E = y_i - fp\{end\}outputWeights$, which is approximately 0. So when we calculate the gradients, they are approximately equal to 0, too. So, during the training procedure, we don't change the weights and thus lead to its bad performance. And I won't add this method to my final script.

1.9 Create more training examples

Since the handwriting could be less precise (too big or too small, rotate clockwise or anti-clockwise, etc.) So we artificially creat more training examples by applying small transforations (like rotations and resizing, etc.) to the original images :

```

1 for i = 1:size(X,1)
2     X1 = imrotate(reshape(X(i,:),16,16), 5, 'crop');
3     X_clock(i,:) = X1(:);
4     X1 = imrotate(reshape(X(i,:),16,16), -5, 'crop');
5     X_anticlock(i,:) = X1(:);
6     X1 = imresize(reshape(X(i,:),16,16), 1.1, 'OutputSize', [16,16]);
7     X_big(i,:) = X1(:);
8     X1 = imresize(reshape(X(i,:),16,16), 0.9, 'OutputSize', [16,16]);
9     X_small(i,:) = X1(:);
10 end
11
12 X = [X; X_clock; X_anticlock; X_big; X_small];

```

```
13 y = repmat(y,5,1);
```

After create the new training set (X, y), we train our model as before with the new training set and find that the validation error convergent to 20%. The method works better when we got less dataset, but I think that is because our training model is not big enough. This method will work better when the training data is more large and more chaos, so I will not use it in the final script.

1.10 2D convolutional layer

In this section, we replace the first layer of the network with a 2D convolutional layer and train the model by the following step:

a. Reshape the USPS images back to their original 16 by 16 format:

```
1 [X,mu,sigma] = standardizeCols(X);
2 X = reshape(X',16,16,n);
3
4 Xvalid = standardizeCols(Xvalid,mu,sigma);
5 Xvalid = reshape(Xvalid',16,16,t);
6
7 Xtest = standardizeCols(Xtest,mu,sigma);
8 Xtest = reshape(Xtest',16,16,t2);
```

b. For each iteration, we randomly choose i from 1 to n(number of training data), and then finish the convolution layer c1 use the *conv2* function:

```
1 i = ceil(rand*n);
2 X1 = X(:, :, i);
3
4 c1 = zeros(12,12,nConv);
5 for k = 1:nConv
6     c1(:, :, k) = conv2(X1, rot90(cWeights(:, :, k), 2), 'valid');
7     c1(:, :, k) = tanh(c1(:, :, k) + cBias(1,k));
8 end
```

c. Now we finish the full-connected layer: put the convolution layer c1 into norml 16 x 16 format and make the connection to the hidden layer:

```
1 for n = 1:nHidden
2     count = 0;
3     for m = 1:nConv
4         count = count + c1(:, :, m) * hiddenWeights(m,n);
5     end
6     e(:, :, n) = count;
7     f1(:, :, n) = conv2(e(:, :, n), ...
8         rot90(fWeights(:, :, n), 2), 'valid');
9 end
10
```

```

11 for h = 1:nLabels
12     output(1,h) = exp( f0*outputWeights(:,h) ) / ...
13     sum( exp(f0*outputWeights) );
14 end

```

c. Using softmax method to generate the output layer because we don't have the exact prediction to each label :

```

1 output = zeros(1, nLabels);
2 for h = 1:nLabels
3     output(1,h) = exp( f0*outputWeights(:,h) ) / ...
4     sum( exp(f0*outputWeights) );
5 end

```

d. To updating the weights and bias, I wrote a function named *ConvUpdate* and the predict function named *ConvPredict* is shown in Sec.3:

```

1 % Update weights, kernels and bias.
2 [cWeights, fWeights, hiddenWeights, outputWeights, cBias, fBias] = ...
3 WechatCNN_update(stepSize, y(i), X1, c1, f0, ...
4 e, output, cWeights, fWeights, hiddenWeights, ...
5 outputWeights, cBias, fBias);

```

After running the script, I got the result that validation error convergence to 25%. For the training procedure is too slowly, I just set the maxIter = 10000, and The result seems not bad, but it still causes a lot of time, so I choose not use this method in my final script.

2 Final script

After finishing all the 10 questions, I did all the following adjustment to my final script:

- I choose nHidden=[100];
- Modify the sequence of step-size by $w^{t+1} = w^t - \alpha_t \nabla f(w^t) + \beta_t (w^t - w^{t-1})$, where $\beta_t = 0.9$
- Vectorize the evaluating the function to speed up the training;
- Add a l2 regularization $\lambda = 0.03$;
- Use a softmax layer at the end of the network;. Run the final script for several times, The test error is approximately 6%. The final script is shown in Sec.3.

3 Matlab Code

3.1 LossSoftmax

```

1 function [f,g] = LossSoftmax(w,X,y,nHidden,nLabels)
2 % g : the new weights
3 % f : the gradient
4 %
5 % Tianran Zhang Dec.6th

```

```

6
7 [nInstances,nVars] = size(X);
8 nH = length(nHidden);
9
10 % Form Weights
11 inputWeights = reshape(w(1:nVars*nHidden(1)),nVars,nHidden(1));
12 offset = nVars*nHidden(1);
13 for h = 2:length(nHidden)
14     hiddenWeights{h-1} = reshape(w(offset+1 : offset+nHidden(h-1) * ...
15         nHidden(h)), nHidden(h-1), nHidden(h));
16     offset = offset + nHidden(h-1) * nHidden(h);
17 end
18 outputWeights = w(offset+1 : offset+nHidden(end)*nLabels);
19 outputWeights = reshape(outputWeights, nHidden(end), nLabels);
20
21 f = 0;
22 ip = cell(1, nH);
23 fp = cell(1, nH);
24 if nargin > 1
25     gInput = zeros(size(inputWeights));
26     gHidden = cell(1, nH-1);
27     for h = 2:nH
28         gHidden{h-1} = zeros(size(hiddenWeights{h-1}));
29     end
30     gOutput = zeros(size(outputWeights));
31 end
32
33 % Compute Output
34 for i = 1:nInstances
35     ip{1} = X(i,:) * inputWeights;
36     fp{1} = tanh(ip{1});
37     for h = 2:length(nHidden)
38         ip{h} = fp{h-1} * hiddenWeights{h-1};
39         fp{h} = tanh(ip{h});
40     end
41
42     yhat1 = exp(fp{end}*outputWeights);
43     yhat = yhat1/sum(yhat1);
44     y_true] = max(y(i, :));
45     yhat_true = yhat(y_true);
46     err = -log(yhat_true);
47     f=f + err;
48
49     if nargin > 1
50

```

```

51      % Output Weights
52      gOutput = gOutput - fp{end}' * (1 - yhat_true) * (y(i,:)==1);%
53      if nH > 1
54          % Last Layer of Hidden Weights
55          backprop = err' * sech(ip{end} .* outputWeights');
56          gHidden{end} = gHidden{end} + fp{end-1}' * sum(backprop,1);
57          backprop = sum(backprop,1);
58
59
60          % Other Hidden Layers
61          for h = nH-2:-1:1
62              backprop = (backprop*hiddenWeights{h+1}') .* ...
63                  sech(ip{h}
64                      gHidden{h} = gHidden{h} + fp{h}' * backprop;
65          end
66
67          % Input Weights
68          backprop = (backprop*hiddenWeights{1}') .* sech(ip
69          gInput = gInput + X(i,:) ' * backprop;
70      else
71          % Input Weights
72          gInput = gInput - (1 - yhat_true) * X(i,:) ' * ...
73              ( sech(ip{end} .* outputWeights(:, y(i,:)==1))' );
74
75      end
76  end
77 end
78
79 % Put Gradient into vector
80 if nargout > 1
81     g = zeros(size(w));
82     g(1:nVars*nHidden(1)) = gInput(:);
83     offset = nVars*nHidden(1);
84     for h = 2:nH
85         g(offset+1:offset+nHidden(h-1)*nHidden(h)) = gHidden{h-1};
86         offset = offset+nHidden(h-1)*nHidden(h);
87     end
88     g(offset+1:offset+nHidden(end)*nLabels) = gOutput(:);
89 end

```

3.2 PredictSoftmax

```

1  %Tianran Zhang Dec.6th
2
3  function [y] = MLPclassificationPredict(w,X,nHidden,nLabels)

```

```

4 [nInstances , nVars] = size(X);
5
6 % Form Weights
7 inputWeights = reshape(w(1:nVars*nHidden(1)), nVars, nHidden(1));
8 offset = nVars*nHidden(1);
9
10 for h = 2:length(nHidden)
11     hiddenWeights{h-1} = reshape(w(offset+1 : offset + ...
12         nHidden(h-1)*nHidden(h)), nHidden(h-1), nHidden(h));
13     offset = offset+nHidden(h-1)*nHidden(h);
14 end
15
16 outputWeights = w(offset+1:offset+nHidden(end)*nLabels);
17 outputWeights = reshape(outputWeights, nHidden(end), nLabels);
18
19 % Compute Output
20 y = zeros(nInstances, 1);
21 for i = 1:nInstances
22     ip{1} = X(i,:) * inputWeights;
23     fp{1} = tanh(ip{1});
24     for h = 2:length(nHidden)
25         ip{h} = fp{h-1} * hiddenWeights{h-1};
26         fp{h} = tanh(ip{h});
27     end
28     yhat1 = exp(fp{end} * outputWeights);
29     yhat(i,:) = yhat1 / sum(yhat1);
30 end
31 [v, y] = max(yhat, [], 2);

```

3.3 ConvPredict

```

1 function [yhat] = ConvPredict(X, cWeights, fWeights, hiddenWeights, ...
2     outputWeights, cBias, fBias)
3
4 nInstances = size(X, 3);
5 yhat = zeros(nInstances, 1);
6 nConv = size(cWeights, 3);
7 [nHidden, nLabels] = size(outputWeights);
8 c1 = size(X, 1) - size(cWeights, 1) + 1;
9 c2 = size(X, 2) - size(cWeights, 2) + 1;
10
11 for i = 1:nInstances
12     train_data = X(:, :, i);
13
14     c0 = zeros(c1, c2, nConv);

```



```

15     for k = 1:nConv
16         c0(:, :, k) = conv2(train_data, rot90(cWeights(:, :, k), 2), 'valid');
17         c0(:, :, k) = tanh(c0(:, :, k) + cBias(1, k));
18     end
19
20     [f1 = conv(c0, fWeights, hiddenWeights);
21     f = zeros(1, nHidden);
22     for h = 1:nHidden
23         f(1, h) = tanh(f1(:, :, h) + fBias(1, h));
24     end
25
26     output = zeros(1, nLabels);
27     for h = 1:nLabels
28         output(1, h) = exp( f * outputWeights(:, h) ) / ...
29             sum( exp(f * outputWeights) );
30     end
31     yhat(i)] = max(output);
32 end
33
34 end

```

3.4 final_MLP

```

1  %Final script
2  %
3  % Tianran Zhang, Dec. 5, 2017.
4
5  load digits.mat
6  [n, d] = size(X);
7  nLabels = max(y);
8  yExpanded = linearInd2Binary(y, nLabels);
9  t = size(Xvalid, 1);
10 t2 = size(Xtest, 1);
11
12 % Standardize columns and add bias
13 [X, mu, sigma] = standardizeCols(X);
14 X = [ones(n, 1) X];
15 d = d + 1;
16
17 % Make sure to apply the same transformation to the validation/test data
18 Xvalid = standardizeCols(Xvalid, mu, sigma);
19 Xvalid = [ones(t, 1) Xvalid];
20 Xtest = standardizeCols(Xtest, mu, sigma);
21 Xtest = [ones(t2, 1) Xtest];
22

```

```

23 % Choose network structure
24 nHidden = [100];
25
26 % Count number of parameters and initialize weights 'w'
27 nParams = d*nHidden(1);
28 for h = 2:length(nHidden)
29     nParams = nParams+nHidden(h-1)*nHidden(h);
30 end
31 nParams = nParams+nHidden(end)*nLabels;
32 maxIter = 100000;
33 w = randn(nParams,1);
34
35 % Train with stochastic gradient
36 stepSize = 1e-3;
37
38 funObj = @(w,i) Loss_final(w,X(i,:),yExpanded(i,:),nHidden,nLabels, 0.03);
39
40 tic;
41 w0 = 0;
42 a0 = 1;
43 a1 = 1;
44 for iter = 1:maxIter
45     if mod(iter-1,round(maxIter/20)) == 0
46         yhat = Predict5(w,Xvalid,nHidden,nLabels);
47         a2 = sum(yhatyvalid)/t;
48         fprintf('Training iteration = %d, validation error = %f\n',iter-1, a2);
49
50         %%early stop
51         if (a2 > a1 ) & (a1 > a0)
52             break;
53         else
54             a0 = a1;
55             a1 = a2;
56         end
57     end
58
59     i = ceil(rand*n);
60     g] = funObj(w,i);
61     w = w - stepSize*g + 0.9 * (w - w0);
62     w0 = w;
63 end
64 toc;
65
66 % Evaluate test error
67 yhat = Predict5(w, Xtest, nHidden, nLabels);

```

```
68 fprintf('Test error with final model = %f\n',sum(yhat ytest)/t2);
```

3.5 Loss_final

```
1 function [f,g] = Loss_final(w,X,y,nHidden,nLabels,lambda)
2 % Tianran Zhang, Dec. 7th
3 % g : the new weights
4 % f : the gradient
5
6 [nInstances,nVars] = size(X);
7 nH = length(nHidden);
8
9 % Form Weights
10 inputWeights = reshape(w(1:nVars*nHidden(1)),nVars,nHidden(1));
11 offset = nVars*nHidden(1);
12 for h = 2:length(nHidden)
13     hiddenWeights{h-1} = reshape(w(offset+1 : offset+nHidden(h-1) * ...
14         nHidden(h)), nHidden(h-1), nHidden(h));
15     offset = offset + nHidden(h-1) * nHidden(h);
16 end
17 outputWeights = w(offset+1 : offset+nHidden(end)*nLabels);
18 outputWeights = reshape(outputWeights, nHidden(end), nLabels);
19
20 f = 0;
21 ip = cell(1, nH);
22 fp = cell(1, nH);
23 if nargout > 1
24     gInput = zeros(size(inputWeights));
25     gHidden = cell(1, nH-1);
26     for h = 2:nH
27         gHidden{h-1} = zeros(size(hiddenWeights{h-1}));
28     end
29     gOutput = zeros(size(outputWeights));
30 end
31
32 % Compute Output
33 for i = 1:nInstances
34     ip{1} = X(i,:) * inputWeights;
35     fp{1} = tanh(ip{1});
36     for h = 2:length(nHidden)
37         ip{h} = fp{h-1} * hiddenWeights{h-1};
38         fp{h} = tanh(ip{h});
39     end
40
41     yhat1 = exp(fp{end}*outputWeights);
```

```

42     yhat = yhat1/sum(yhat1);
43     y_true] = max(y(i, :));
44     yhat_true = yhat(y_true);
45     err = -log(yhat_true);
46     f=f + err;
47
48     if nargout > 1
49
50         % Output Weights
51         gOutput = gOutput - fp{end}' * (1 - yhat_true) * (y(i,:)==1) + ...
52             lambda * outputWeights .* ...
53             [zeros(1, size(outputWeights, 2)); ...
54             ones(size(outputWeights, 1)-1, size(outputWeights, 2))]];
55
56         if nH > 1
57             % Last Layer of Hidden Weights
58             backprop = err' * sech(ip{end} .* outputWeights');
59             gHidden{end} = gHidden{end} + fp{end-1}' * sum(backprop,1);
60
61             backprop = sum(backprop,1);
62             % gHidden{end} = gHidden{end} + fp{end-1}' * backprop;
63
64             % Other Hidden Layers
65             for h = nH-2:-1:1
66                 backprop = (backprop*hiddenWeights{h+1}') .* ...
67                     sech(ip{h}
68                     gHidden{h} = gHidden{h} + fp{h}' * backprop;
69             end
70
71             % Input Weights
72             backprop = (backprop*hiddenWeights{1}') .* sech(ip
73             gInput = gInput + X(i,:) ' * backprop;
74         else
75             % Input Weights
76             gInput = gInput - (1 - yhat_true) * X(i,:) ' * ...
77                 ( sech(ip{end} .* outputWeights(:, y(i,:)==1)')+...
78                 lambda * inputWeights .* ...
79                 [zeros(1, size(inputWeights, 2)); ...
80                 ones(size(inputWeights, 1)-1, size(inputWeights, 2))]];
81
82         end
83     end
84 end
85
86 % Put Gradient into vector

```

```

87 if nargout > 1
88     g = zeros(size(w));
89     g(1:nVars*nHidden(1)) = gInput(:);
90     offset = nVars*nHidden(1);
91     for h = 2:nH
92         g(offset+1:offset+nHidden(h-1)*nHidden(h)) = gHidden{h-1};
93         offset = offset+nHidden(h-1)*nHidden(h);
94     end
95     g(offset+1:offset+nHidden(end)*nLabels) = gOutput(:);
96 end

```

3.6 Predict_final

```

1 function [y] = Predict_final(w,X,nHidden,nLabels)
2 % Tianran Zhang Dec. 8th
3 [nInstances,nVars] = size(X);
4
5 % Form Weights
6 inputWeights = reshape(w(1:nVars*nHidden(1)),nVars,nHidden(1));
7 offset = nVars*nHidden(1);
8
9 for h = 2:length(nHidden)
10     hiddenWeights{h-1} = reshape(w(offset+1 : offset + ...
11         nHidden(h-1)*nHidden(h)),nHidden(h-1),nHidden(h));
12     offset = offset+nHidden(h-1)*nHidden(h);
13 end
14
15 outputWeights = w(offset+1:offset+nHidden(end)*nLabels);
16 outputWeights = reshape(outputWeights,nHidden(end),nLabels);
17
18 % Compute Output
19 y = zeros(nInstances, 1);
20 for i = 1:nInstances
21     ip{1} = X(i,:)*inputWeights;
22     fp{1} = tanh(ip{1});
23     for h = 2:length(nHidden)
24         ip{h} = fp{h-1}*hiddenWeights{h-1};
25         fp{h} = tanh(ip{h});
26     end
27     yhat1 = exp(fp{end}*outputWeights);
28     yhat(i,:) = yhat1/sum(yhat1);
29     % y(i) = max(yhat(i,:));
30 end
31 [v,y] = max(yhat,[],2);

```