# Convolutional Neural Networks

**Akshat Dave,   Muhammad Ahmed Riaz,   Nicholas Kinkade**
Computer Science and Engineering
University of California, San Diego
{akdave,mriaz,nkinkade}@ucsd.edu

(Programming Assignment 4, Submitted : 5-June-2015)

## Abstract

Convolutional neural networks have been used successfully in recent literature for learning applications, and particularly for classification of images. In this report, we explore the application of convolutional neural networks to multiclass classification, specifically, in the classification of handwritten digits in the *MNIST* dataset (as directed in the fourth programing assignment of **CSE291: Neural Networks**). We present a simple to use reconfigurable architecture for developing a convolutional neural network model. We evaluate the performance of some of these architectures/models on the *MNIST* dataset. Interestingly, we discover that the max pooling marginally improves accuracy and that the network architecture has a significant impact on the iteration profile. We report an accuracy of greater than $91\%$ on the test data classification of the *MNIST* images.

## 1   Introduction

This report discusses the problems posed in the fourth programming assignment[1] for **CSE291: Neural Networks**. We use convolutional neural networks for the purpose of classification, namely for digit recognition from the *MNIST* dataset. Our implementational contributions include:

1. Support for a reconfigurable user-defined architecture
2. Support for several activation functions including rectified linear units, logistic sigmoid and modified tanh
3. Support for both mean-pooling and max-pooling
4. Support for arbitrary convolutional kernel size, arbitrary pooling window and arbitrary feature map count at each convolutional layer
5. All parameter gradient checking utility

As we have designed a generalizable architecture, we claim and subsequently demonstrate in the experiments section that our system handles each of the specified requirements of PA4. We cite sources [1] and [2] as an inspiration for datastructures and optimizations for batch processing.

The remaining paper is organized as follows: Section 2 describes the method, Section 3 discusses the experiments, while Section 4 concludes the paper.

## 2   Method

Given an input space $X$, we wish to map this to an output space $y$. We refer to $X$ as the features while $y$ as the labels. In the case that we have the input space as the form of images, we model this mapping
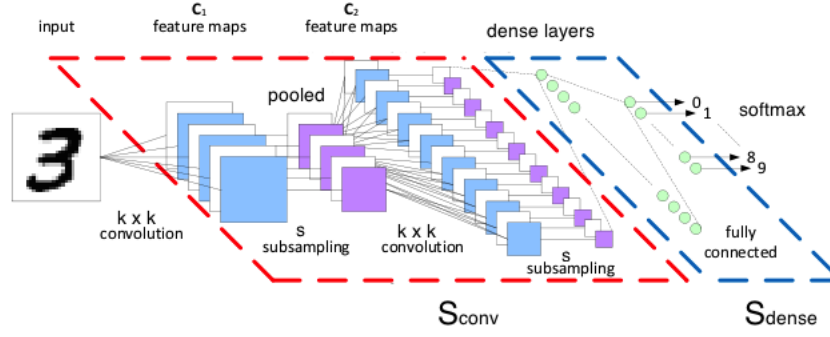
---

[1] http://tinyurl.com/pbcyp8c

Figure 1: Schematic of a generic convolutional neural network. The dense section $S_{dense}$ is reconfigurable to allow $n_d \geqslant 1$ layers where the last layer is a softmax layer. Similarly, the convolutional layer is reconfigurable to allow $n_c \geqslant 1$ layers where the first layer is an input layer. This image is an edited graphic from the TUE-course page

$h : X \mapsto y$ as a convolutional neural network parameterized by weights at each layer. The design of the neural network is similar to that presented in Fig. 1, where $S_{conv}$ is the convolutional section and $S_{dense}$ densely connected section. The convolutional section ($S_{conv}$) accomodates the convolution and the pooling layers, while the densely connected section ($S_{dense}$) accomodates a multilayered densely connected neural network. We constrain the design to allow $S_{conv}$ only before the $S_{dense}$ (as per requirements) however the order of layers within the network is free to configure. Our convolutional section supports both mean-pooling and max-pooling layers. The advantage of such a feed forward multilayered architecture is that it is able to learn internal representations. The mapping $h$ is learnt by the backpropagation algorithm. Since in our application we use these networks for multi-class classification, the final layer is a softmax layer. We demonstrate that training this network is efficient in that the gradients can be analytically computed for training. We proceed to discuss backpropagation (as described in [2]):

Suppose we have a fixed training set $(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})$ of $m$ training examples. We are able to train our neural network using batch gradient descent. Specifically, for a mini-batch of size $m$ with training examples $(x^{(i)}, y^{(i)})$, we have the cross entropy loss function to be:

$$L(W, b, \theta) = -\sum_{i=1}^{m} \sum_{k=1}^{K} \mathbf{1}\{y^{(i)} = k\} \log \frac{\exp \theta^{(k)T} z^{(n_d,i)}}{\sum_{j=1}^{K} \exp \theta^{(k)T} z^{(n_d,i)}} \text{ where } z^{(n_d,i)} = h_{W,b}(x^{(i)})$$

Note above that $W, b$ are the weights and biases of all the layers in the network except the final softmax layer which is parameteried by $\theta$ for explanatory purposes. More details can be found in the UFLDL lecture notes (notes-link). Our aim is to find the parameters that minimize the loss function, in other words:

$$\{W^*, b^*, \theta^*\} = \arg \min_{W,b,\theta} \left( -\sum_{i=1}^{m} \sum_{k=1}^{K} \mathbf{1}\{y^{(i)} = k\} \log \frac{\exp \theta^{(k)T} h_{W,b}(x^{(i)})}{\sum_{j=1}^{K} \exp \theta^{(k)T} h_{W,b}(x^{(i)})} \right)$$

We generally regularize this loss function by adding a penalty to the norm of the parameters. Thus, given a training set of $m$ examples, we can define the regularized loss function to be:

$$L(W, b, \theta) = \frac{1}{m} \sum_{i=1}^{m} L(W, b, \theta; x^{(i)}, y^{(i)}) + \frac{\lambda}{2} \left\| \mathbf{vec}\left( \begin{bmatrix} W \\ b \\ \theta \end{bmatrix} \right) \right\|^2 \tag{1}$$

---

[2]http://ufldl.stanford.edu/

2

In order to minimize the loss function, we proceed to take the derivative (since the loss function is differentiable) in order to design the update rule for gradient descent. We use the chain rule to minimize the loss function with respect to the weights at each layer. At this point we augment the weights $W$ and $\theta$ such that $W = (W, \theta)$. It can be shown that the gradient of the loss function with respect to the parameters of the densely connected layers is given by:

$$\nabla_{W^{(l)}} L(W, b; x, y) = \delta^{(l+1)} (a^{(l)})^T, \text{ and, } \nabla_{b^{(l)}} L(W, b; x, y) = \delta^{(l+1)} \ \forall n_c + 1 \leqslant l \leqslant n_d$$

while for the convolutional layers, the gradient with respect to the $k^t h$ filter is given by:

$$\nabla_{W_k^{(l)}} L(W, b; x, y) = \sum_{i=1}^{m} (a_i^{(l)})^T * \Phi(\delta_k^{(l+1)}), \text{ and, } \nabla_{b_k^{(l)}} L(W, b; x, y) = \sum_{a,b} (\delta_k^{(l+1)})_{a,b} \ \forall 1 \leqslant l \leqslant n_c$$

Where $a^{(l)}$ is the activation vector of the $l^{th}$ layer and $\Phi$ is a horizontal matrix flip/reflection operator and $(*)$ is the convolution operator. We recall that for a network with $n_c$ convolution/pooling layers and $n_d$ densely connected layers we use a backpropagation rule similar to multilayered neural networks. Here we use the shorthand $\delta$ to indicate the propagated error terms. For the densely connected layers we have the update rule:

$$\delta^{(l)} = \begin{cases} -(y - a^{(n_l)}) \cdot f'(z^{(n_l)}), \text{ if } l = n_d \\ (W^{(l)})^T \delta^{(l+1)}) \cdot f'(z^{(l)}) \text{ if } n_c + 1 \leqslant l \leqslant n_d - 1 \end{cases}$$

While for the convolutional/pooling layers we have:

$$\delta_k^{(l)} = \Psi_s((W_k^{(l)})^T \delta_k^{(l+1)})) \cdot f'(z_k^{(l)}), \text{ if } 1 \leqslant l \leqslant n_c$$

Where $f'(z^{(l)})$ is the derivative of the activation function at the layer $l$ with input $z$ and $\Psi_s$ is an upsampling function with sampling factor $s$. In the absence of a pooling layer, $s = 1$. We define $(\cdot)$ as the hadamard product.

We allow our network to work with various activation functions including the logistic sigmoid, the modified hyperbolic tangent and the rectified linear function. The derivative $f'(z^{(l)})$ for the logistic sigmoid function is simply given by $f'(z^{(l)}) = z^{(l)} \cdot (1 - z^{(l)})$. Here, $z^{(l)}$ is computed as (for all densely connected layers):

$$z^{(l)} = W^{(l)} f(z^{(l-1)}), \ \forall n_c + 1 \leqslant l \leqslant n_d$$

While for the convolutional layers it is given by:

$$z_k^{(l)} = \sum_{i=1}^{|l-1|} W_{ki}^{(l)} * f(z_i^{(l-1)}), \ \forall 1 \leqslant l \leqslant n_c$$

Where $|l - 1|$ is the number of feature maps output from layer $l - 1$. Finally, we have the update equations:

$$W^{(l)} = W^{(l)} - \alpha(\frac{\Delta W^{(l)}}{m} + \lambda W^{(l)}), \text{ and,} \tag{2a}$$

$$b^{(l)} = b^{(l)} - \alpha(\frac{\Delta b^{(l)}}{m}) \tag{2b}$$

Where the update for $\Delta W^{(l)}$ and $\Delta b^{(l)}$ are given by

$$\Delta W^{(l)} := \Delta W^{(l)} + \nabla_{W^{(l)}} L(W, b; x, y) \text{ and } \Delta b^{(l)} := \Delta b^{(l)} + \nabla_{b^{(l)}} L(W, b; x, y)$$

We optimize this function by minibatch gradient descent. The algorithm for gradient descent is generic, in-that if the batch size is $m = 1$, we have stochastic gradient descent, is $1 < m < |\mathbf{X}|$, we have minibatch gradient descent, and if $m = |\mathbf{X}|$ we have batch gradient descent. The algorithm is given in Alg. 2. For gradient descent we use momentum for faster convergence. Our stopping criteria is based upon maximum iterations.

For verifying correctness of our gradients, we perform a gradient check using the following test:

$$\text{verify if } (y^{(i)} - h_\theta(x^{(i)}))x_j^{(i)} \approx \lim_{\epsilon \to 0} \frac{L(\theta_j + \epsilon) - L(\theta_j - \epsilon)}{2\epsilon}$$

The algorithm is provided in Alg. 1.

---

**Algorithm 1** Gradient check algorithm

---

1: **procedure** DOGRADIENTCHECK($\theta, \epsilon$)
2:     $t \leftarrow \phi$ (equality test flag)
3:     **for each** $j$ **do**
4:         $\hat{g}_j \leftarrow \big(L(\theta_j + \epsilon) - L(\theta_j - \epsilon)\big)/2\epsilon$
5:         $g_j \leftarrow \nabla_{\theta_j} L(\theta)$
6:         **if** $g_j \approx \hat{g}_j$ **then**
7:             $t_j \leftarrow 1$ (success)
8:         **else**
9:             $t_j \leftarrow 0$ (fail)
10:        **end if**
11:    **end for**
12: **return** $t$
13: **end procedure**

---

**Algorithm 2** Gradient descent for error backpropagation. This algorithm is used for stochastic gradient descent when $m = 1$, for mini-batch gradient descent when $1 < m < |\mathbf{X}|$, and, for batch gradient descent when $m = |\mathbf{X}|$, where $|\mathbf{X}|$ is the size of the training dataset). Note that this algorithm was modified to include momentum (found in the appendix) to improve the results

---

1: **procedure** DOBACKPROP($\alpha, m, N$)
2: Inputs: (1) $\alpha$ : learning rate (2) $m$ : batch size (3) $N$ : maximum iterations
3:     $i \leftarrow 1$, randomly initialize weights
4:     **while** $i \leqslant N$ **do**
5:         Perform a feedforward pass, computing the activations for each layer
6:         For the output layer (layer $n_l$), $\delta^{(n_l)} \leftarrow -(y - a^{(n_l)}) \cdot f'(z^{(n_l)})$
7:         $\forall l \in \{n_l - 1, n_l - 2, n_l - 3, \ldots, n_c + 1\}$ set $\delta^{(l)} \leftarrow ((W^{(l)})^T \delta^{(l+1)}) \cdot f'(z^{(l)})$
8:         $\forall l \in \{n_c, n_l - 1, n_l - 2, \ldots, 1\}$ set $\delta_k^{(l)} \leftarrow \Psi_s(((W_k^{(l)})^T \delta_k^{(l+1)})) \cdot f'(z_k^{(l)})$
9:         Compute : $\nabla_{W^{(l)}} L(W, b; x, y) = \delta^{(l+1)}(a^{(l)})^T$ and $\nabla_{b^{(l)}} L(W, b; x, y) = \delta^{(l+1)}$
10:        $\forall l, \Delta W^{(l)} \leftarrow 0, \Delta b^{(l)} \leftarrow 0$
11:        **for** $1 \leqslant i \leqslant m$ **do**
12:            Use backprop from sec. 2 to compute $\nabla_{W^{(l)}} L(W, b; x, y)$ and $\nabla_{b^{(l)}} L(W, b; x, y)$
13:            $\Delta W^{(l)} \leftarrow \Delta W^{(l)} + \nabla_{W^{(l)}} L(W, b; x, y)$ and $\Delta b^{(l)} \leftarrow \Delta b^{(l)} + \nabla_{b^{(l)}} L(W, b; x, y)$
14:        **end for**
15:        Update parameters
16:        $W^{(l)} \leftarrow W^{(l)} - \alpha(\frac{\Delta W^{(l)}}{m} + \lambda W^{(l)})$ and $b^{(l)} \leftarrow b^{(l)} - \alpha(\frac{\Delta b^{(l)}}{m})$
17:        $i \leftarrow i + 1$
18:    **end while**
19: **return** $W, b$
20: **end procedure**

---

| model architecture name | architecture description | layer details |
|---|---|---|
| **MA1** | 1 input layer<br>2 convolutional layers<br>1 pooling layer<br>1 densely connected layer<br>1 softmax output layer | 1 channel input<br>2 conv. feat maps from $9 \times 9$ kernels<br>4 conv. feat maps from $5 \times 5$ kernels<br>$4 \times 4$ max pooling window<br>30 fully connected units layer<br>10 output units |
| **MA2** | 1 input layer<br>2 convolutional layers<br>2 pooling layers<br>1 densely connected layer<br>1 softmax output layer | 1 channel input<br>2 conv. feat maps from $5 \times 5$ kernels<br>$2 \times 2$ pooling window<br>2 conv. feat maps from $5 \times 5$ kernels<br>$4 \times 4$ pooling window<br>10 fully connected units layer<br>10 output units |
| **MA3** | 1 input layer<br>1 convolutional layer<br>1 pooling layer<br>1 softmax output layer | 1 channel input<br>20 conv. feat maps from $9 \times 9$ kernels<br>$2 \times 2$ pooling window<br>10 output units |

Table 1: Detailed specification of the models developed for experiments

# 3 Experiments

In this section we discuss the datasets used for each task, the experimental results and the parameters used by our models.

## 3.1 Dataset description

We use the *MNIST* dataset for our experiments. The dataset comprises $60,000$ labeled training examples and $10,000$ labeled test images. The task is to identify the digit class $1, 2, \ldots 10$ for the test data. We train a few convolutional neural network models and discuss these here.

## 3.2 Training and results

For training the parameters we use backpropagation with loss minimization using mini-batch gradient descent (MGD) as discussed in section 2. Gradient descent is a first order optimization method which uses a first order Taylor approximation to update the parameters in the direction of the negative gradient of the loss function.

Table 1 enlists the model architectures that we evaluated with gradient checking. We verify whether the analytical solution agrees with the numerical approximation as described in section 2 for each of the architectures (**MA1**, **MA2**, **MA3**) using each of the three activation functions viz. logistic sigmoid ($\sigma$), modified tanh ($\tau$) and the rectified linear activation ($\rho$). The gradient check passed for each of the activation functions for each of the architectures and some of the results of the gradient check are shown in Table 4

Table 2 enlists the models developed for evaluation on the *MNSIT* dataset. The parameters of the model are provided in the same table along with the architecture used.

For each of the models (**CNN1** and **CNN2**) discussed, the results on *MNIST* are provided in Table 3. Finally, Fig. 2 provides all the model training iteration profiles.

| Model | parameter | description | value |
|---|---|---|---|
| **CNN1** | $\alpha$ | learning rate | $1e-1$ |
| | $N_m$ | max iterations for mini-batch gradient descent | $2e2$ |
| | $m$ | mini-batch size | 256 |
| | $\gamma$ | momentum | 0.95 |
| | $\epsilon$ | epochs | 3 |
| | - | pooling type | max |
| | $f()$ | activation function | tanh |
| | - | architecture used | **MA2** |
| **CNN2** | $\alpha$ | learning rate | $1e-1$ |
| | $N_m$ | max iterations for mini-batch gradient descent | $1e2$ |
| | $m$ | mini-batch size | 256 |
| | $\gamma$ | momentum | 0.95 |
| | $\epsilon$ | epochs | 3 |
| | - | pooling type | max |
| | $f()$ | activation function | rectified linear |
| | - | architecture used | **MA3** |
| **CNN3** | $\alpha$ | learning rate | $1e-1$ |
| | $N_m$ | max iterations for mini-batch gradient descent | $1e2$ |
| | $m$ | mini-batch size | 256 |
| | $\gamma$ | momentum | 0.95 |
| | $\epsilon$ | epochs | 3 |
| | - | pooling type | max |
| | $f()$ | activation function | modified tanh |
| | - | architecture used | **MA3** |
| **CNN4** | $\alpha$ | learning rate | $1e-1$ |
| | $N_m$ | max iterations for mini-batch gradient descent | $1e2$ |
| | $m$ | mini-batch size | 256 |
| | $\gamma$ | momentum | 0.95 |
| | $\epsilon$ | epochs | 3 |
| | - | pooling type | mean |
| | $f()$ | activation function | modified tanh |
| | - | architecture used | **MA3** |

Table 2: Shows the parameters for each of the model architectures used for *MNIST* classification

| model | training accuracy | test accuracy | training time |
|---|---|---|---|
| **CNN1** | 90.1% | 90.7% | 163s |
| **CNN2** | 91.2% | 92.1% | 625s |
| **CNN3** | 91.0% | 91.68% | 632s |
| **CNN4** | 89.9% | 90.50% | 467s |

Table 3: Results of experiments for each of the models

| Model Architecture | True gradient $\nabla L(\theta)$ | Numerical approximation $\frac{(L(\theta+\epsilon)-L(\theta+\epsilon))}{2\epsilon}$ | Relative error | $\epsilon$ |
|---|---|---|---|---|
| **MA1** | $9.506783e-7$ | $9.506751e-7$ | $3.20e-12$ | $1e-4$ |
|  | $9.540871e-7$ | $9.540813e-7$ | $5.81e-12$ | $1e-4$ |
|  | $2.344425e-7$ | $2.344502e-7$ | $7.74e-12$ | $1e-4$ |
|  | $3.84988e-7$ | $3.849920e-7$ | $3.79e-12$ | $1e-4$ |
| **MA2** | $6.4193e-4$ | $6.4193e-4$ | $5.81e-13$ | $1e-4$ |
|  | $6.539991e-7$ | $6.539902e-7$ | $8.93e-12$ | $1e-4$ |
|  | $7.383305e-7$ | $7.383316e-7$ | $1.10e-12$ | $1e-4$ |
|  | $5.741183e-7$ | $5.741252e-7$ | $7.96e-13$ | $1e-4$ |
|  | $2.214921e-7$ | $2.214905e-7$ | $1.65e-12$ | $1e-4$ |
| **MA3** | $2.184758e-4$ | $2.184758e-4$ | $2.43e-12$ | $1e-4$ |
|  | $1.547129e-4$ | $1.547129e-4$ | $9.33e-13$ | $1e-4$ |
|  | $1.489933e-5$ | $1.489933e-5$ | $2.23e-12$ | $1e-4$ |
|  | $-1.406128e-4$ | $-1.406128e-4$ | $8.41e-12$ | $1e-4$ |

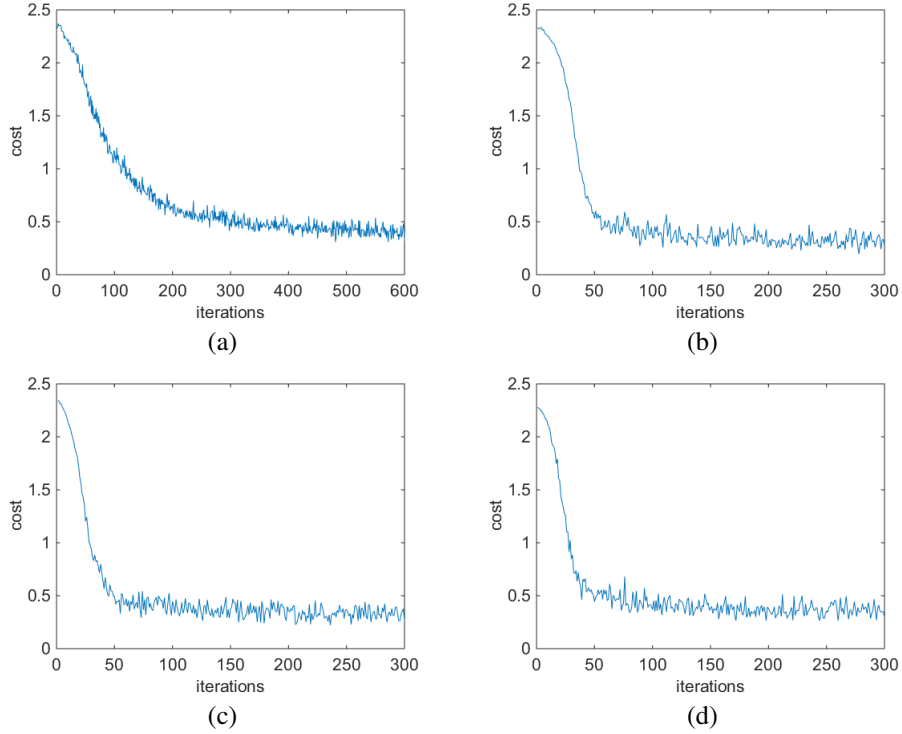Table 4: Shows numerical approximation and its comparison to the true gradient



Figure 2: Figures show the iteration profiles for the models developed (a) **CNN1** (b) **CNN2** (c) **CNN3** (d) **CNN4**

## 3.3 Discussion

We find that classification between the *MNIST* digits yields over $90\%$ accuracy. The four models we train perform relatively well with **CNN1** reporting an accuracy of $90.7\%$ ,**CNN2** with an accuracy of $92.1\%$, **CNN3** with an accuracy 91.68, **CNN4** with an accuracy of $90.50\%$.

We observe max pooling increases the accuracy as compared to mean pooling as seen in the $1\%$ increase in accuracy between **CNN4** and **CNN3**. Additionally, we observe that the architecture has an effect on the convergence trend. We see that architectures **CNN2**, **CNN3** and **CNN4** have similar convergence trends as compared to **CNN1** which has a different architecture. As **CNN1** is deeper, it drops slower in cost as compared to **CNN2**, **CNN3** and **CNN4**.

We observed that the logistic sigmoid (iteration profiles not included) gives slightly lower performance ($85.5\%$ training accuracy and $86.5\%$ test accuracy) compared to using the modified tanh and the rectified linear units. We also observe that the training time increases with the size of the batch.

The gradient check algorithm ensures that the relative error in between the approximation and the true gradient is always less than $O(\epsilon)$.

# 4 Conclusion

We have demonstrated that our proposed system allows the design of reconfigurable convolutional neural network architectures with several convolutional/pooling layers and several densely connected layers. We have presented various exemplar model architectures, viz. **MA1**, **MA2** and **MA3**. Finally, we have demonstrated that the models created by our system, namely, **CNN1**, **CNN2**, **CNN3** and **CNN4** all perform classifaction on the *MNIST* dataset with over $90\%$ accuracy.

# References

[1] A. Vedaldi and K. Lenc, "MatConvNet – Convolutional Neural Networks for MATLAB", *CoRR* (2014).

[2] "Prediction as a candidate for learning deep hierarchical models of data", *Palm* (2012)

# Appendix : Source Code

Listing 1: calculate_cost.m

```matlab
function cost = calculate_cost(netOut, labels, model, lambda)
%regularization term
reg = 0;
for l=1:numel(model.conv.layers)
    if(strcmp(model.conv.layers{1}.type,'conv'))
        for i = 1:numel(model.conv.layers{1}.W)
            for j = 1:numel(model.conv.layers{1}.W{i})
                reg = reg + sum(sum((model.conv.layers{1}.W{i}{j
                    }.^2)));
            end
        end
    end
end
for i=1:numel(model.dense.layers)
    reg = reg + sum(sum((model.dense.layers{i}.W.^2)));
end
cost = reg * lambda/2;

%cost function
m = size(netOut,2);
for i = 1:m
    cost = cost - log(netOut(labels(i),i))/m;
end
end
```

Listing 2: cnnConvolve.m

```matlab
 1 function convolvedFeatures = cnnConvolve(filterDim, numFilters, ...
 2                                           images, W, b, activationType)
 3 %cnnConvolve Returns the convolution of the features given by W and b with
 4 %the given images
 5 %
 6 % Parameters:
 7 %  filterDim - filter (feature) dimension
 8 %  numFilters - number of feature maps
 9 %  images - large images to convolve with, matrix in the form
10 %           images(r, c, image number)
11 %  W, b - W, b for features from the sparse autoencoder
12 %         W is of shape (filterDim, filterDim, numFilters)
13 %         b is of shape (numFilters, 1)
14 %  activationType - 1 : logistic sigmoid function
15 %                   2 : rectified linear unit function
16 %                   3 : funny tanh function
17 %
18 % Returns:
19 %  convolvedFeatures - matrix of convolved features in the form
20 %                      convolvedFeatures(imageRow, imageCol, featureNum, imageNum)
21
22 numImages = size(images, 3);
23 imageDim = size(images, 1);
24 convDim = imageDim - filterDim + 1;
25
26 convolvedFeatures = zeros(convDim, convDim, numFilters, numImages);
27
28 [activationFunc, ~] = select_act(activationType);
29
30 % Instructions:
31 %    Convolve every filter with every image here to produce the
32 %    (imageDim - filterDim + 1) x (imageDim - filterDim + 1) x numFeatures x numImages
33 %    matrix convolvedFeatures, such that
34 %    convolvedFeatures(imageRow, imageCol, featureNum, imageNum) is the
35 %    value of the convolved featureNum feature for the imageNum image over
36 %    the region (imageRow, imageCol) to (imageRow + filterDim - 1, imageCol + filterDim - 1)
37 %
38 % Expected running times:
39 %    Convolving with 100 images should take less than 30 seconds
40 %    Convolving with 5000 images should take around 2 minutes
41 %    (So to save time when testing, you should convolve with less images, as
42 %    described earlier)
43
44
45 for imageNum = 1:numImages
46   for filterNum = 1:numFilters
47
48     % convolution of image with feature matrix
```

```
49        convolvedImage = zeros(convDim, convDim);
50
51        % Obtain the feature (filterDim x filterDim) needed during the
               convolution
52        %OUR CODE
53        filter = W(:,:,filterNum);
54
55        % Flip the feature matrix because of the definition of
               convolution, as explained later
56        filter = rot90(squeeze(filter),2);
57
58        % Obtain the image
59        im = squeeze(images(:, :, imageNum));
60
61        % Convolve "filter" with "im", adding the result to
               convolvedImage
62        % be sure to do a 'valid' convolution
63        %OUR CODE
64        convolvedImage = convolvedImage + conv2(im, filter, 'valid');
65
66        % Add the bias unit
67        % Then, apply the activation function to get the hidden
               activation
68        %OUR CODE
69        bias = b(filterNum);
70        convolvedImage = convolvedImage + bias;
71        convolvedImage = activationFunc(convolvedImage);
72
73        convolvedFeatures(:, :, filterNum, imageNum) = convolvedImage;
74    end
75 end
76 end
```

Listing 3: cnnCost.m

```
 1 function [cost, grad, predDist] = cnnCost(theta,images,labels,
     numClasses ,...
 2                                  model, lambda,   ...
 3                                  activationType, ...
 4                                  pred_only)
 5 %default to no prediction
 6 if ~exist('pred_only','var')
 7     pred_only = false;
 8 end;
 9
10 %reshape parameters
11 [model] = cnnParamsToStack(theta, model);
12
13 %deciding activation function and gradient of activation function
14 [act_fun, grad_fun, ~] = select_act(activationType);
15
16 %forward propogate
17 [model] = conv_forward_pass(images, model, act_fun, grad_fun);
18 [convOut] = vectorize_conv_out(model);
19 [model, netOut] = hidden_forward_pass(convOut, model, act_fun,
     grad_fun);
20
21 %if we are only performing prediction
22 if pred_only
23     predDist = netOut; % prediction distributions
```

```
24      grad = [];
25      cost = -1;
26      return ;
27 end
28
29 %compute cost
30 cost = calculate_cost(netOut, labels, model, lambda);
31
32 %backprop
33 [model, delta] = hidden_backprop(labels, model, netOut, numClasses
       , lambda);
34 [deltaOut] = vectorize_delta_out(delta);
35 [model] = conv_backprop(model, deltaOut, lambda);
36
37 %convert deltas
38 gradModel = makeGradModel(model);
39
40 %convert gradients to vector
41 grad = cnnStackToParams(gradModel);
42 end
```

Listing 4: cnnExercise.m

```
 1 %% Convolution and Pooling Exercise
 2
 3 %   Instructions
 4 %   ------------
 5 %
 6 %   This file contains code that helps you get started on the
 7 %   convolution and pooling exercise. In this exercise, you will
       only
 8 %   need to modify cnnConvolve.m and cnnPool.m. You will not need
       to modify
 9 %   this file.
10
11 %
       %===================================================================
12 %% STEP 0: Initialization and Load Data
13 %   Here we initialize some parameters used for the exercise.
14
15 imageDim = 28;            % image dimension
16
17 filterDim = 8;            % filter dimension
18 numFilters = 100;         % number of feature maps
19
20 numImages = 60000;      % number of images
21
22 poolDim = 3;              % dimension of pooling region
23
24 % Here we load MNIST training images
25 addpath ../common/;
26 images = loadMNISTImages('../common/train-images-idx3-ubyte');
27 images = reshape(images,imageDim,imageDim,numImages);
28
29 W = randn(filterDim, filterDim, numFilters);
30 b = rand(numFilters, 1);
31
```

11

```matlab
32 %
    %==========================================================================

33 %% STEP 1: Implement and test convolution
34 %  In this step, you will implement the convolution and test it on
35 %  on a small part of the data set to ensure that you have
       implemented
36 %  this step correctly.
37
38 %% STEP 1a: Implement convolution
39 %  Implement convolution in the function cnnConvolve in
       cnnConvolve.m
40
41 %% Use only the first 8 images for testing
42 convImages = images(:, :, 1:8);
43
44 convolvedFeatures = cnnConvolve(filterDim, numFilters, convImages,
      W, b, 1);
45
46 %% STEP 1b: Checking your convolution
47 %  To ensure that you have convolved the features correctly, we
       have
48 %  provided some code to compare the results of your convolution
       with
49 %  activations from the sparse autoencoder
50
51 % For 1000 random points
52 for i = 1:1000
53     filterNum = randi([1, numFilters]);
54     imageNum = randi([1, 8]);
55     imageRow = randi([1, imageDim − filterDim + 1]);
56     imageCol = randi([1, imageDim − filterDim + 1]);
57
58     patch = convImages(imageRow:imageRow + filterDim − 1, imageCol
         :imageCol + filterDim − 1, imageNum);
59
60     feature = sum(sum(patch.*W(:,:,filterNum)))+b(filterNum);
61     feature = 1./(1+exp(−feature));
62
63     if abs(feature − convolvedFeatures(imageRow, imageCol,
         filterNum, imageNum)) > 1e−9
64         fprintf('Convolved feature does not match test feature\n')
             ;
65         fprintf('Filter Number    : %d\n', filterNum);
66         fprintf('Image Number     : %d\n', imageNum);
67         fprintf('Image Row        : %d\n', imageRow);
68         fprintf('Image Column     : %d\n', imageCol);
69         fprintf('Convolved feature : %0.5f\n', convolvedFeatures(
              imageRow, imageCol, filterNum, imageNum));
70         fprintf('Test feature : %0.5f\n', feature);
71         error('Convolved feature does not match test feature');
72     end
73 end
74
75 disp('Congratulations! Your convolution code passed the test.');
76
77 %
    %==========================================================================
```

```
78 %% STEP 2: Implement and test pooling
79 %   Implement pooling in the function cnnPool in cnnPool.m
80
81 %% STEP 2a: Implement pooling
82 % NOTE: Implement cnnPool in cnnPool.m first!
83 pooledFeatures = cnnPool(poolDim, convolvedFeatures, 1);
84
85 %% STEP 2b: Checking your pooling
86 %   To ensure that you have implemented pooling, we will use your
        pooling
87 %   function to pool over a test matrix and check the results.
88
89 testMatrix = reshape(1:64, 8, 8);
90 expectedMatrix = [mean(mean(testMatrix(1:4, 1:4))) mean(mean(
        testMatrix(1:4, 5:8))); ...
91                   mean(mean(testMatrix(5:8, 1:4))) mean(mean(
                        testMatrix(5:8, 5:8))); ];
92
93 testMatrix = reshape(testMatrix, 8, 8, 1, 1);
94
95 pooledFeatures = squeeze(cnnPool(4, testMatrix, 1));
96
97 if ~isequal(pooledFeatures, expectedMatrix)
98     disp('Pooling incorrect');
99     disp('Expected');
100    disp(expectedMatrix);
101    disp('Got');
102    disp(pooledFeatures);
103 else
104    disp('Congratulations! Your pooling code passed the test.');
105 end
```

Listing 5: cnnInitParams.m

```
 1 function [theta, model] = cnnInitParams(imageDim, model, init_zero)
 2 %initialize convolutional layers
 3 inCount = 1;
 4 mapDim = [imageDim imageDim];
 5 for l = 1 : numel(model.conv.layers)
 6         if strcmp(model.conv.layers{l}.type, 'conv')
 7         mapDim = mapDim - model.conv.layers{l}.kernelDim + 1;
 8         fanOut = model.conv.layers{l}.outCount * model.conv.layers
                {l}.kernelDim * model.conv.layers{l}.kernelDim ;
 9         for j = 1 : model.conv.layers{l}.outCount
10             fanIn = inCount * model.conv.layers{l}.kernelDim *
                    model.conv.layers{l}.kernelDim;
11             for i = 1 : inCount
12                 s = sqrt(6 / (fanIn + fanOut));
13                 model.conv.layers{l}.W{i}{j} = rand(model.conv.
                        layers{l}.kernelDim)*2*s - s;
14                 if (init_zero)
15                     model.conv.layers{l}.W{i}{j}  = model.conv.
                            layers{l}.W{i}{j}*0;
16                 end
17             end
18             model.conv.layers{l}.b{j} = 0;
19         end
20         inCount = model.conv.layers{l}.outCount;
21         end
22
```

```matlab
23          if strcmp(model.conv.layers{l}.type, 'pool')
24              mapDim = mapDim / model.conv.layers{l}.poolDim;
25          end
26      end
27
28      model.outDim = prod(mapDim) * inCount;
29
30      %initialize dense layers
31      for l = 1 : numel(model.dense.layers)
32          if l > 1
33              fanIn = model.dense.layers{l-1}.outCount;
34          else
35              fanIn = model.outDim;
36          end;
37          fanOut = model.dense.layers{l}.outCount;
38          s = sqrt(6) / sqrt(fanIn + fanOut);
39
40          model.dense.layers{l}.W = rand(fanOut, fanIn)*2*s - s;
41          if(init_zero)
42              model.dense.layers{l}.W = model.dense.layers{l}.W*0;
43          end
44          model.dense.layers{l}.b = zeros(fanOut, 1);
45      end
46
47      %convert to vector form
48      theta = cnnStackToParams(model);
49  end
```

Listing 6: cnnParamsToStack.m

```matlab
1  function [model] = cnnParamsToStack(theta, model)
2  cur_pos = 1;
3
4  %convolutional layers
5  inCount = 1;
6  for l = 1 : numel(model.conv.layers)
7      if strcmp(model.conv.layers{l}.type, 'conv')
8          for j = 1 : model.conv.layers{l}.outCount
9              for i = 1 : inCount
10                 wlen = model.conv.layers{l}.kernelDim ^ 2;
11                 model.conv.layers{l}.W{i}{j} = reshape(theta(
                       cur_pos:cur_pos+wlen-1), model.conv.layers{l}.
                       kernelDim, model.conv.layers{l}.kernelDim);
12                 cur_pos = cur_pos + wlen;
13             end
14             model.conv.layers{l}.b{j} = theta(cur_pos);
15             cur_pos = cur_pos + 1;
16         end
17         inCount = model.conv.layers{l}.outCount;
18     end
19  end
20
21  %dense layers
22  hiddenDepth = numel(model.dense.layers);
23  inCount = model.outDim;
24  for d = 1:hiddenDepth
25      outCount = model.dense.layers{d}.outCount;
26
27      wlen = double(outCount * inCount);
```

14

```
28        model.dense.layers{d}.W = reshape(theta(cur_pos:cur_pos+wlen
              -1), outCount, inCount);
29        cur_pos = cur_pos+wlen;
30
31        blen = outCount;
32        model.dense.layers{d}.b = reshape(theta(cur_pos:cur_pos+blen
              -1), outCount, 1);
33        cur_pos = cur_pos+blen;
34
35        inCount = outCount;
36 end
37 end
```

Listing 7: cnnPool.m

```
1 function pooledFeatures = cnnPool(poolDim, convolvedFeatures,
      poolType)
2 %cnnPool Pools the given convolved features
3 %
4 % Parameters:
5 %   poolDim - dimension of pooling region
6 %   convolvedFeatures - convolved features to pool (as given by
      cnnConvolve)
7 %                       convolvedFeatures(imageRow, imageCol,
      featureNum, imageNum)
8 %   poolType - 1 : MEAN
9 %              2 : MAX
10 %
11 % Returns:
12 %   pooledFeatures - matrix of pooled features in the form
13 %                    pooledFeatures(poolRow, poolCol, featureNum,
      imageNum)
14 %
15
16 numImages = size(convolvedFeatures, 4);
17 numFilters = size(convolvedFeatures, 3);
18 convolvedDim = size(convolvedFeatures, 1);
19
20 pooledFeatures = zeros(convolvedDim / poolDim, ...
21         convolvedDim / poolDim, numFilters, numImages);
22
23 % Instructions:
24 %   Now pool the convolved features in regions of poolDim x
      poolDim,
25 %   to obtain the
26 %   (convolvedDim/poolDim) x (convolvedDim/poolDim) x numFeatures
      x numImages
27 %   matrix pooledFeatures, such that
28 %   pooledFeatures(poolRow, poolCol, featureNum, imageNum) is the
29 %   value of the featureNum feature for the imageNum image pooled
      over the
30 %   corresponding (poolRow, poolCol) pooling region.
31 %
32 %   Use mean pooling here.
33 for imageNum = 1:numImages
34   for filterNum = 1:numFilters
35        convolvedFeature = convolvedFeatures(:, :, filterNum,
            imageNum);
36        if(poolType == 1)
37            %perform mean filtering
```

```
38            feature = conv2(convolvedFeature, ones(poolDim) / (
                 poolDim*poolDim) , 'same');
39        else
40            %perform max filtering
41            feature = ordfilt2(convolvedFeature, poolDim*poolDim,
                 true(poolDim));
42        end
43        %subsample
44        pooledFeatures(:, :, filterNum, imageNum) = feature(ceil(
             poolDim / 2) : poolDim : end, ceil(poolDim / 2) :
             poolDim : end);
45    end
46 end
47
48 end
```

<div align="center">Listing 8: cnnStackToParams.m</div>

```
1 function theta = cnnStackToParams(model)
2 theta = [];
3 %convolutional layers
4 inCount = 1;
5 for l = 1 : numel(model.conv.layers)
6     if strcmp(model.conv.layers{l}.type, 'conv')
7         for j = 1 : model.conv.layers{l}.outCount
8             for i = 1 : inCount
9                 theta = [theta; model.conv.layers{l}.W{i}{j}(:)];
10            end
11            theta = [theta; model.conv.layers{l}.b{j}];
12        end
13        inCount = model.conv.layers{l}.outCount;
14    end
15 end
16
17 %dense layers
18 for d = 1:numel(model.dense.layers)
19    theta = [theta ; model.dense.layers{d}.W(:) ; model.dense.
          layers{d}.b(:)];
20 end
21 end
```

<div align="center">Listing 9: cnnTrain.m</div>

```
1 clear;
2 %convolutional neural network
3 addpath(genpath('../../'));
4
5 %network configuration
6 activationType = 3;
7 imageDim = 28;
8 numClasses = 10;
9 grad_check=true;
10 init_zero=false;
11 model.conv.layers = {
12     struct('type', 'input')
13 %        struct('type', 'conv', 'outCount', 2, 'kernelDim', 5)
14 %        struct('type', 'pool', 'poolType', 'max', 'poolDim', 2)
15 %        struct('type', 'conv', 'outCount', 3, 'kernelDim', 5)
16 %        struct('type', 'pool', 'poolType', 'max', 'poolDim', 2)
17
```

```matlab
18        struct('type', 'conv', 'outCount', 20, 'kernelDim', 9)
19        struct('type', 'pool', 'poolType', 'mean', 'poolDim', 2)
20
21 %       struct('type', 'conv', 'outCount', 2, 'kernelDim', 5)
22 %       struct('type', 'pool', 'poolType', 'max', 'poolDim', 2)
23 %       struct('type', 'conv', 'outCount', 3, 'kernelDim', 5)
24 %       struct('type', 'pool', 'poolType', 'max', 'poolDim', 4)
25 };
26 model.dense.layers = {
27 %       struct('outCount', 10)
28 %       struct('outCount', 10)
29       struct('outCount', numClasses)
30 };
31
32 %load MNIST training images
33 addpath  ../common/;
34 images = loadMNISTImages('../common/train-images-idx3-ubyte');
35 images = reshape(images,imageDim,imageDim,[]);
36 labels = loadMNISTLabels('../common/train-labels-idx1-ubyte');
37 labels(labels==0) = 10; % Remap 0 to 10
38
39 %initialize parameters
40 [theta, model] = cnnInitParams(imageDim,model,init_zero);
41
42 %gradient check
43 if(grad_check)
44     ran = randperm(60000,100);
45     gradCheckImages = images(:,:,ran);
46     gradCheckLabels = labels(ran);
47     gradient_check(numel(theta), 0.1, @(theta)cnnCost(theta,
          gradCheckImages, gradCheckLabels, numClasses, model, ...
48         0.001, activationType, false), 3);
49 end
50
51 %training configuration
52 options.epochs = 3;
53 options.numIterations = 200;
54 options.minibatch = 256; %100
55 options.alpha = 0.1;
56 options.momentum = 0.95;
57
58 %train
59 tic;
60 opttheta = minFuncSGD(@(theta,images,labels)cnnCost(theta, images,
       labels, numClasses, model,...
61                         0.001, activationType, false),theta,images,
                           labels,options);
62 fprintf('Optimization took %f seconds.\n', toc);
63
64 %load MNIST testing images
65 testImages = loadMNISTImages('../common/t10k-images-idx3-ubyte');
66 testImages = reshape(testImages,imageDim,imageDim,[]);
67 testLabels = loadMNISTLabels('../common/t10k-labels-idx1-ubyte');
68 testLabels(testLabels==0) = 10; % Remap 0 to 10
69
70 %calculate training accuracy
71 [~,cost,preds] =  cnnCost(opttheta, images, labels, numClasses,
      model,...
72                         0.001, activationType, true);
```

```matlab
73 [~, pred ] = max( preds );
74 train_acc = mean( pred==reshape( labels , size ( pred ) ) );
75
76 %calculate test accuracy
77 [~, cost , preds ] = cnnCost( opttheta , testImages , testLabels ,
      numClasses , model ,...
78                        0.001 , activationType , true );
79 [~, pred ] = max( preds );
80 test_acc = mean( pred==reshape ( testLabels , size ( pred ) ) );
81
82 %display accuracies
83 fprintf ( 'Training Accuracy is %f\n' , train_acc );
84 fprintf ( 'Test Accuracy is %f\n' , test_acc );
```

Listing 10: conv_backprop.m

```matlab
1 function [model] = conv_backprop( model , deltaIn , lambda )
2 n = numel( model.conv.layers ) ;
3
4 %reshape feature vector deltas into output map style
5 outActSize = size ( model.conv.layers{n}.act{1} );
6 outActDim = outActSize(1) * outActSize(2) ;
7 for j = 1 : numel( model.conv.layers{n}.act )
8     if ( length ( outActSize )<3 )
9         outActSize(3) = 1;
10    end
11    model.conv.layers{n}.d{j} = reshape( deltaIn ((( j − 1) *
          outActDim + 1) : j * outActDim , :) , outActSize(1) ,
          outActSize(2) , outActSize(3) );
12 end
13
14 %calculate deltas
15 for l = (n − 1) : −1 : 1
16     %backprop previous layer
17     backDeltas = cell( numel( model.conv.layers{l}.act ) , 1 ) ;
18     if strcmp ( model.conv.layers{l+1}.type , 'conv' )
19         for i = 1 : numel( model.conv.layers{l}.act )
20             backDeltas{i} = zeros ( size ( model.conv.layers{l}.act
                  {1} ) ) ;
21             for j = 1 : numel( model.conv.layers{l + 1}.act )
22                 backDeltas{i} = backDeltas{i} + convn( model.conv.
                      layers{l + 1}.d{j} , flip ( flip ( model.conv.
                      layers{l + 1}.W{i}{j} , 1) , 2) , 'full' );
23             end
24         end
25     elseif strcmp ( model.conv.layers{l+1}.type , 'pool' )
26         for i = 1 : numel( model.conv.layers{l}.act )
27             if ( strcmp ( model.conv.layers{l+1}.poolType , 'mean' ) )
28                 backDeltas{i} = expand( model.conv.layers{l + 1}.d{
                      i} , [ model.conv.layers{l + 1}.poolDim model.
                      conv.layers{l + 1}.poolDim 1] ) / model.conv.
                      layers{l + 1}.poolDim ^ 2;
29             end
30             if ( strcmp ( model.conv.layers{l+1}.poolType , 'max' ) )
31                 backDeltas{i} = zeros ( size ( model.conv.layers{l +
                      1}.d{i},1)*model.conv.layers{l + 1}.poolDim ,
                      size ( model.conv.layers{l + 1}.d{i},2)*model.
                      conv.layers{l + 1}.poolDim , size ( model.conv.
                      layers{l + 1}.d{i},3));
32                 for j = 1 : size ( model.conv.layers{l+1}.mask{i},3)
```

18

```matlab
33                      mask = model.conv.layers{l+1}.mask{i}(:,:,j);
34                      deltas = model.conv.layers{l+1}.d{i}(:,:,j);
35                      mask = mask(:);
36                      deltas = deltas(:);
37                      curDelta = backDeltas{i}(:,:,j);
38                      curDelta(mask) = deltas;
39                      backDeltas{i}(:,:,j) = curDelta;
40                  end
41              end
42
43          end
44      end
45
46      %multiply by gradient of activation function
47      if ~strcmp(model.conv.layers{l}.type, 'input')
48          for i = 1 : numel(model.conv.layers{l}.act)
49              model.conv.layers{l}.d{i} = model.conv.layers{l}.
                  actGrad{i} .* backDeltas{i};
50          end
51      end
52 end
53
54 %calculate gradients
55 for l = 2 : n
56      if strcmp(model.conv.layers{l}.type, 'conv')
57          for j = 1 : numel(model.conv.layers{l}.act)
58              for i = 1 : numel(model.conv.layers{l-1}.act)
59                  model.conv.layers{l}.dW{i}{j} = convn(flipall(
                      model.conv.layers{l-1}.act{i}), model.conv.
                      layers{l}.d{j}, 'valid') / size(model.conv.
                      layers{l}.d{j}, 3);
60              end
61              model.conv.layers{l}.db{j} = sum(model.conv.layers{l}.
                  d{j}(:)) / size(model.conv.layers{l}.d{j}, 3);
62          end
63      end
64 end
65
66 %regularization
67 for l = 2 : n
68      if (strcmp(model.conv.layers{l}.type,'conv'))
69          for i = 1:numel(model.conv.layers{l}.W)
70              for j = 1:numel(model.conv.layers{l}.W{i})
71                  model.conv.layers{l}.dW{i}{j} = model.conv.layers{
                      l}.dW{i}{j} + model.conv.layers{l}.W{i}{j} *
                      lambda;
72              end
73          end
74      end
75 end
76 end
```

Listing 11: conv_forward_pass.m

```matlab
1 function [ model ] = conv_forward_pass( images, model, act_fun,
      grad_fun )
2 n = numel(model.conv.layers);
3 model.conv.layers{1}.act{1} = images;
4 inCount = 1;
5 %for every layer
```

```matlab
 6  for l = 2 : n
 7      %convolution layer
 8      if strcmp(model.conv.layers{l}.type, 'conv')
 9          for j = 1 : model.conv.layers{l}.outCount
10              sz = [model.conv.layers{l}.kernelDim - 1 model.conv.
                     layers{l}.kernelDim - 1 0];
11              z = zeros(size(model.conv.layers{l - 1}.act{1}) - sz);
12              for i = 1 : inCount
13                  z = z + convn(model.conv.layers{l - 1}.act{i},
                         model.conv.layers{l}.W{i}{j}, 'valid');
14              end
15              z = z + model.conv.layers{l}.b{j};
16              model.conv.layers{l}.act{j} = act_fun(z);
17              model.conv.layers{l}.actGrad{j} = grad_fun(z);
18          end
19          inCount = model.conv.layers{l}.outCount;
20      %pooling layer
21      elseif strcmp(model.conv.layers{l}.type, 'pool')
22          if strcmp(model.conv.layers{l}.poolType, 'mean')
23              for j = 1 : inCount
24                  z = convn(model.conv.layers{l - 1}.act{j}, ones(
                         model.conv.layers{l}.poolDim) / (model.conv.
                         layers{l}.poolDim ^ 2), 'valid');   % !!
                         replace with variable
25                  model.conv.layers{l}.act{j} = z(1 : model.conv.
                         layers{l}.poolDim : end, 1 : model.conv.layers
                         {l}.poolDim : end, :);
26                  model.conv.layers{l}.actGrad{j} = 1;
27              end
28          else
29              for j = 1 : inCount
30                  z = zeros(size(model.conv.layers{l-1}.act{j}));
31                  mask = zeros(size(model.conv.layers{l-1}.act{j}));
32                  for i = 1 : size(model.conv.layers{l-1}.act{j}, 3)
33                      [z(:,:,i), mask(:,:,i)] = slideMax(model.conv.
                             layers{l-1}.act{j}(:,:,i), model.conv.
                             layers{l}.poolDim);
34                  end
35                  model.conv.layers{l}.act{j} = z(1 : model.conv.
                         layers{l}.poolDim : end, 1 : model.conv.layers
                         {l}.poolDim : end, :);
36                  model.conv.layers{l}.actGrad{j} = 1;
37                  model.conv.layers{l}.mask{j} = mask(1 : model.conv
                         .layers{l}.poolDim : end, 1 : model.conv.
                         layers{l}.poolDim : end, :);
38              end
39          end
40      end
41  end
42  end
43
44  function [xMax, mask] = slideMax(x, poolDim)
45  %sliding max pooling, also gives mask with indications where
46  %max pooling came from
47  %can certainly be more optimized
48  poolDim = poolDim-1;
49  [width, height] = size(x);
50  xMax = zeros(width, height);
51  mask = zeros(width, height);
```

```
52  for i = 1 : width
53      for j = 1 : height
54          maxVal = -Inf;
55          maxK = -1;
56          maxL = -1;
57          for k = max(1,i) : min(i+poolDim,width)
58              for l = max(1,j) : min(j+poolDim,height)
59                  if(x(k,l) > maxVal)
60                      maxVal = x(k,l);
61                      maxK = k;
62                      maxL = l;
63                  end
64              end
65          end
66          xMax(i,j) = maxVal;
67          mask(i,j) = maxK + (maxL-1)*width;
68      end
69  end
70  end
```

Listing 12: funnyTan.m

```
1  function f = funnyTan(X)
2  f = 1.7159 * tanh(2/3 * X);
3  end
```

Listing 13: funnyTanGradient.m

```
1  function g = funnyTanGradient(X)
2  g = 1.7159*2/3 * sech(2/3.*X).^2;
3  end
```

Listing 14: gradient_check.m

```
1  function gradient_check(thetaSize, checkFrac, costGradientFunc,
       checkCount)
2  %inputs:
3  % thetaSize - number of elements in theta
4  % checkFrac - fraction of elements in theta to check each time
5  % costGradientFunc - function to calculate cost (with input theta)
6  %                    and gradient (with input theta)
7  % checkCount - number of checks to do
8
9  %checks checkCount random theta value gradients
10 for i = 1 : checkCount
11     %get random theta
12     theta = rand(thetaSize,1)*0.001;
13
14     [~, g] = costGradientFunc(theta);
15
16     %check gradient for this theta
17     single_gradient_check(theta, g, costGradientFunc, checkFrac);
18 end
19 end
20
21 function single_gradient_check(theta, g, costGradientFunc,
       checkFrac)
22 EPSILON = 0.0001;
23 n = size(theta, 1);
24 thetaIndices = 1:n;%randperm(n, uint16(checkFrac*n));
```

```matlab
25 k = 1;
26 for i = thetaIndices
27     theta_i_plus = theta;
28     theta_i_minus = theta;
29     theta_i_plus(i) = theta_i_plus(i) + EPSILON;
30     theta_i_minus(i) = theta_i_minus(i) - EPSILON;
31     [cost_plus, ~] = costGradientFunc(theta_i_plus);
32     [cost_minus, ~] = costGradientFunc(theta_i_minus);
33     g_i_approximation = (cost_plus - cost_minus) / (2 * EPSILON);
34     delta = abs(g(i) - g_i_approximation);
35     fprintf('%d\n',k);
36     if(g(i) == 0)
37         assert(delta < 1e-10);
38     else
39         delta_ratio = delta / g(i);
40         %assert(delta_ratio < EPSILON * 10);
41         if(delta_ratio >= EPSILON * 10)
42             fprintf('%e %e %0.2e  %0.2e\n', g(i),
                   g_i_approximation, EPSILON, delta);
43         end
44     end
45     if k<10
46         %print first 10 approximations
47         fprintf('%e %e %0.2e  %0.2e\n', g(i), g_i_approximation,
               EPSILON, delta);
48     end
49     k = k + 1;
50 end
51 end
```

Listing 15: hidden_backprop.m

```matlab
1 function [model, delta] = hidden_backprop(labels, model, netOut,
      numClasses, lambda)
2 numHidden = numel(model.dense.layers) - 1;
3 m = size(netOut,2);
4
5 %stack of gradients (initialize to zeros)
6 for i = 1:numel(model.dense.layers)
7     model.dense.layers{i}.dW = zeros(size(model.dense.layers{i}.W)
          );
8     model.dense.layers{i}.db = zeros(size(model.dense.layers{i}.b)
          );
9 end
10
11 delta = cell(m,1);
12 %backprop for every data point
13 for i_data = 1:m
14     %label indicator vector
15     labelIndicator = zeros(numClasses,1);
16     labelIndicator(labels(i_data)) = 1;
17     %compute delta for output layer
18     delta{i_data}{numHidden+2} = -(labelIndicator - netOut(:,
          i_data)); %HERE
19     %backprop delta
20     for i=(numHidden+1):-1:1
21         %activation
22         a = model.dense.layers{i}.act{i_data};
23         %gradient of activation
24         g = model.dense.layers{i}.actGrad{i_data};
```

```matlab
25          %compute delta for hidden layer
26          delta{i_data}{i} = (model.dense.layers{i}.W' *delta{i_data
                }{i+1}).*g;
27          %update gradient stack
28          model.dense.layers{i}.dW = model.dense.layers{i}.dW +
                delta{i_data}{i+1}*a'/m;
29          model.dense.layers{i}.db = model.dense.layers{i}.db +
                delta{i_data}{i+1}/m;
30      end
31 end
32
33 %regularization
34 for i = 1:numel(model.dense.layers)
35      model.dense.layers{i}.dW = model.dense.layers{i}.dW + model.
            dense.layers{i}.W*lambda;
36 end
37 end
```

Listing 16: hidden_forward_pass.m

```matlab
1 function [ model, netOut ] = hidden_forward_pass( data , model,
      act_fun , grad_fun )
2 %dense forward pass
3 numHidden = numel(model.dense.layers) − 1;
4 netOut = zeros(numel(model.dense.layers{end}.b),size(data,2));
5 for i_data = 1:size(data,2)
6      %first layer
7      model.dense.layers{1}.act{i_data} = data(:,i_data);
8      model.dense.layers{1}.actGrad{i_data} = 1;
9      %other layers
10     for i=1:(numHidden+1)
11         z = model.dense.layers{i}.W * model.dense.layers{i}.act{
              i_data} + model.dense.layers{i}.b;
12         if(i < numHidden+1)
13             model.dense.layers{i+1}.act{i_data} = act_fun(z);
14             model.dense.layers{i+1}.actGrad{i_data} = grad_fun(z);
15         end
16     end
17     netOut(:,i_data) = exp(z)/sum(exp(z));
18 end
19 end
```

Listing 17: invFunnyTanGradient.m

```matlab
1 function g = invFunnyTanGradient(act)
2      g = (1.7159*2/3) *(1−(act./(1.7159)).^2);
3 end
```

Listing 18: invLogisticSigmoidGradient.m

```matlab
1 function g = invLogisticSigmoidGradient(act)
2      g = act.*(1−act);
3 end
```

Listing 19: invReluGradient.m

```matlab
1 function g = invReluGradient(act)
2      g = double(act>0);
3 end
```

## Listing 20: logisticSigmoid.m

```
1 function f = logisticSigmoid(X)
2 f = 1 ./ (1 + exp(-X));
3 end
```

## Listing 21: logisticSigmoidGradient.m

```
1 function g = logisticSigmoidGradient(X)
2 g = logisticSigmoid(X) .* (1 - logisticSigmoid(X));
3 end
```

## Listing 22: makeGradModel.m

```
1 function gradModel = makegradModel(model)
2 %create a model to convert the cell structure to a param vector
3 for l = 1:numel(model.conv.layers)
4     if(strcmp(model.conv.layers{l}.type,'conv'))
5         gradModel.conv.layers{l}.outCount = model.conv.layers{l}.
            outCount;
6         for i = 1:numel(model.conv.layers{l}.W)
7             for j = 1:numel(model.conv.layers{l}.W{i})
8                 gradModel.conv.layers{l}.W{i}{j} = model.conv.
                    layers{l}.dW{i}{j};
9             end
10        end
11        for j = 1 : numel(model.conv.layers{l}.b)
12            gradModel.conv.layers{l}.b{j} = model.conv.layers{l}.
                db{j};
13        end
14        gradModel.conv.layers{l}.type = 'conv';
15    else
16        gradModel.conv.layers{l}.type = 'unused';
17    end
18 end
19
20 for l = 1 : numel(model.dense.layers)
21     gradModel.dense.layers{l}.W = model.dense.layers{l}.dW;
22     gradModel.dense.layers{l}.b = model.dense.layers{l}.db;
23 end
24 end
```

## Listing 23: minFuncSGD.m

```
1 function [opttheta] = minFuncSGD(funObj,theta,data,labels,options)
2 epochs = options.epochs; %epochs
3 numIterations = options.numIterations; %number of iterations
4 alpha = options.alpha; %learning rate
5 minibatch = options.minibatch; %minibatch size
6 m = length(labels); % training set size
7 mom = options.momentum; %momentum
8 velocity = zeros(size(theta));
9
10 all_cost =[];
11 for e = 1 : epochs
12     for i = 1 : numIterations
13         ind = randperm(m,minibatch);
14         mb_data = data(:,:,ind);
15         mb_labels = labels(ind);
16         [cost, grad] = funObj(theta,mb_data,mb_labels);
17         velocity = (1-mom) * alpha * grad + mom * velocity;
```

```
18            theta = theta − velocity;
19
20            all_cost = [all_cost, cost]; plot(all_cost);
21         drawnow update
22         disp(cost)
23      end
24      velocity = velocity * 0;
25      alpha = alpha / 2;
26 end
27
28 %save('costs_info', all_cost, theta);
29
30 opttheta = theta;
31 end
```

Listing 24: relu.m

```
1 function f = relu(X)
2 f = max(0, X);
3 end
```

Listing 25: reluGradient.m

```
1 function g = reluGradient(X)
2 g = double(relu(X) > 0);
3 end
```

Listing 26: vectorize_conv_out.m

```
1 function [convOut] = vectorize_conv_out(model)
2     %concatenate all end layer feature maps into vector
3     convOut = [];
4     for j = 1 : numel(model.conv.layers{end}.act)
5         actSize = size(model.conv.layers{end}.act{j});
6         if( length(actSize)<3 )
7             actSize(3) = 1;
8         end
9         convOut = [convOut; reshape(model.conv.layers{end}.act{j},
               actSize(1) * actSize(2), actSize(3))];
10     end
11 end
```

Listing 27: vectorize_delta_out.m

```
1 function deltaOut = vectorize_delta_out(delta)
2 %get deltaOut from delta
3 deltaOut = [];
4 for i = 1 : numel(delta)
5     deltaOut = [deltaOut delta{i}{1}];
6 end
7 end
```